

Also Helpful for GAIL, BARC, HPCL, BHEL, ONGC, SAIL, DRDO & Other PSU's

TRISHNA'S

SERIES

crack the

GATEWAY

GRADUATE APTITUDE TEST IN ENGINEERING

COMPUTER SCIENCE AND INFORMATION TECHNOLOGY 2020

HIGHLIGHTS

- ✓ Maximum Coverage/Explanations/Illustrations as per Latest Syllabus
- ✓ 770+ Solved Problems and 2550+ Practice Questions
- ✓ Elaborated Question Bank Covering Previous 12 Years' GATE Question Papers
- ✓ Unit - wise Time - bound Tests
- ✓ 2019 GATE Online Papers with Topic - wise Analysis
- ✓ Solution Manual available in online resources



AVAILABLE
AS e-book

About Pearson

Pearson is the world's learning company, with presence across 70 countries worldwide. Our unique insights and world-class expertise comes from a long history of working closely with renowned teachers, authors and thought leaders, as a result of which, we have emerged as the preferred choice for millions of teachers and learners across the world.

We believe learning opens up opportunities, creates fulfilling careers and hence better lives. We hence collaborate with the best of minds to deliver you class-leading products, spread across the Higher Education and K12 spectrum.

Superior learning experience and improved outcomes are at the heart of everything we do. This product is the result of one such effort.

Your feedback plays a critical role in the evolution of our products and you can contact us at reachus@pearson.com. We look forward to it.

This page is intentionally left blank

GATE

(Graduate Aptitude Test in Engineering)

Computer Science and Information Technology

TRISHNA KNOWLEDGE SYSTEMS



Copyright © 2019 Pearson India Education Services Pvt. Ltd

Published by Pearson India Education Services Pvt. Ltd, CIN: U72200TN2005PTC057128.

No part of this eBook may be used or reproduced in any manner whatsoever without the publisher's prior written consent.

This eBook may or may not include all assets that were part of the print version. The publisher reserves the right to remove any material in this eBook at any time.

ISBN 978-93-534-3394-9

eISBN:

Head Office: A-8(A), 7th Floor, Knowledge Boulevard, Sector 62, Noida 201 309, Uttar Pradesh, India.

Registered Office: The HIVE, 3rd Floor, Metro Zone, No.44, Pillayar Koil Street, Jawaharlal Nehru Road, Anna Nagar, Chennai 600 040, Tamil Nadu, India.

Phone: 044-66540100

Website: in.pearson.com, Email: companysecretary.india@pearson.com

Contents

PREFACE	<i>viii</i>
KEY PEDAGOGICAL FEATURES	<i>ix</i>
SYLLABUS: COMPUTER SCIENCE AND INFORMATION TECHNOLOGY	<i>xi</i>
CHAPTER-WISE ANALYSIS OF GATE PREVIOUS YEARS' PAPERS	<i>xii</i>
GENERAL INFORMATION ABOUT GATE	<i>xiii</i>
SOLVED PAPERS 2019	<i>xv</i>

UNIT I DIGITAL LOGIC **1.1**

CHAPTER 1	NUMBER SYSTEMS	1.3
CHAPTER 2	BOOLEAN ALGEBRA AND MINIMIZATION OF FUNCTIONS	1.14
CHAPTER 3	COMBINATIONAL CIRCUITS	1.35
CHAPTER 4	SEQUENTIAL CIRCUITS	1.56

UNIT II COMPUTER ORGANIZATION AND ARCHITECTURE **2.1**

CHAPTER 1	MACHINE INSTRUCTIONS, ADDRESSING MODES	2.3
CHAPTER 2	ALU AND DATA PATH, CPU CONTROL DESIGN	2.16
CHAPTER 3	MEMORY INTERFACE, I/O INTERFACE	2.33
CHAPTER 4	INSTRUCTION PIPELINING	2.47
CHAPTER 5	CACHE AND MAIN MEMORY, SECONDARY STORAGE	2.60

UNIT III PROGRAMMING AND DATA STRUCTURES **3.1**

PART A PROGRAMMING AND DATA STRUCTURES

CHAPTER 1	PROGRAMMING IN C	3.3
CHAPTER 2	FUNCTIONS	3.14
CHAPTER 3	ARRAYS, POINTERS AND STRUCTURES	3.30
CHAPTER 4	LINKED LISTS, STACKS AND QUEUES	3.47
CHAPTER 5	TREES	3.60

PART B ALGORITHMS

CHAPTER 1	ASYMPTOTIC ANALYSIS	3.81
CHAPTER 2	SORTING ALGORITHMS	3.98
CHAPTER 3	DIVIDE-AND-CONQUER	3.107
CHAPTER 4	GREEDY APPROACH	3.116
CHAPTER 5	DYNAMIC PROGRAMMING	3.135

UNIT IV	DATABASES	4.1
CHAPTER 1	ER MODEL AND RELATIONAL MODEL	4.3
CHAPTER 2	STRUCTURED QUERY LANGUAGE	4.21
CHAPTER 3	NORMALIZATION	4.49
CHAPTER 4	TRANSACTION AND CONCURRENCY	4.65
CHAPTER 5	FILE MANAGEMENT	4.82
UNIT V	THEORY OF COMPUTATION	5.1
CHAPTER 1	FINITE AUTOMATA AND REGULAR LANGUAGES	5.3
CHAPTER 2	CONTEXT FREE LANGUAGES AND PUSH DOWN AUTOMATA	5.24
CHAPTER 3	RECURSIVELY ENUMERABLE SETS AND TURING MACHINES, DECIDABILITY	5.37
UNIT VI	COMPILER DESIGN	6.1
CHAPTER 1	LEXICAL ANALYSIS AND PARSING	6.3
CHAPTER 2	SYNTAX DIRECTED TRANSLATION	6.27
CHAPTER 3	INTERMEDIATE CODE GENERATION	6.36
CHAPTER 4	CODE OPTIMIZATION	6.56
UNIT VII	OPERATING SYSTEM	7.1
CHAPTER 1	PROCESSES AND THREADS	7.3
CHAPTER 2	INTERPROCESS COMMUNICATION, CONCURRENCY AND SYNCHRONIZATION	7.17
CHAPTER 3	DEADLOCK AND CPU SCHEDULING	7.35
CHAPTER 4	MEMORY MANAGEMENT AND VIRTUAL MEMORY	7.54
CHAPTER 5	FILE SYSTEMS, I/O SYSTEMS, PROTECTION AND SECURITY	7.74
UNIT VIII	NETWORKS, INFORMATION SYSTEMS, SOFTWARE ENGINEERING AND WEB TECHNOLOGY	8.1
	PART A	NETWORK
CHAPTER 1	OSI LAYERS	8.3
CHAPTER 2	ROUTING ALGORITHMS	8.24
CHAPTER 3	TCP/UDP	8.36
CHAPTER 4	IP(v4)	8.52
CHAPTER 5	NETWORK SECURITY	8.66

PART B INFORMATION SYSTEMS

CHAPTER 1	PROCESS LIFE CYCLE	8.79
CHAPTER 2	PROJECT MANAGEMENT AND MAINTENANCE	8.89

PART C SOFTWARE ENGINEERING AND WEB TECHNOLOGY

CHAPTER 1	MARKUP LANGUAGES	8.111
------------------	------------------	-------

Preface

Graduate Aptitude Test in Engineering (GATE) is one of the preliminary tests for undergraduate subjects in Engineering/Technology/Architecture and postgraduate subjects in Science stream only.

Apart from giving the aspirant a chance to pursue an M.Tech. from institutions like the IITs /NITs, a good GATE score can be highly instrumental in landing the candidate a plush public sector job, as many PSUs are recruiting graduate engineers on the basis of their performance in GATE. The GATE examination pattern has undergone several changes over the years—sometimes apparent and sometimes subtle. It is bound to continue to do so with changing technological environment.

GATE Computer Science and Information Technology, as a complete resource helps the aspirants be ready with conceptual understanding, and enables them to apply these concepts in various applications, rather than just proficiency with questions type. Topics are handled in a comprehensive manner, beginning with the basics and progressing in a step-by-step manner along with a bottom-up approach. This allows the student to better understand the concept and to practice applicative techniques in a focused manner. The content has been systematically organized to facilitate easy understanding of all topics. The given examples will not only help the students to understand the concepts involved in the problems but also help to get a good idea about the different models of problems on that particular topic. Due care has also been taken to cover a very wide range of problems including questions that have been appearing over the last few years in GATE examination.

The practice exercises in every chapter, contain questions ranging simple to moderate to difficult level. These exercises are meant to hone the examination readiness over a period of time. At the end of each unit, practice tests have been placed. These tests will help the student assess their level of learning on a regular interval.

This book has been prepared by a group of faculty who are highly experienced in training GATE preparations and are also subject matter experts. As a result, this book would serve as an effective tool for GATE aspirant to crack the examination.

Salient Features

1. Elaborate question bank covering previous 16 years' GATE question papers
2. 5 free online mock tests for practice
3. Detailed coverage of key topics
4. Complete set of solved 2019 GATE online papers with chapter-wise analysis
5. Exhaustive pedagogy:
 - (a) More than 770+ Solved Examples
 - (b) More than 2550+ Practice Questions
 - (c) Unit-wise time-bound tests
 - (d) Modular approach for easy understanding

We would like to thank the below mentioned reviewers for their valuable feedback and suggestions which has helped in shaping this book.

R. Marudhachalam	Professor (Sr. Grade), Kumaraguru College of Technology Coimbatore, Tamil Nadu
Daya Gupta	Professor, Delhi Technological University, Main Bawana Road, Delhi
Manoj Kumar Gupta	Associate Professor, Delhi Technological University Main Bawana Road, Delhi
Gurpreet Kour	Lecturer, Lovely Professional University Phagwara, Punjab
Pinaki Chakraborty	Assistant Professor, Netaji Subhas Institute of Technology Dwarka, Delhi
Gunit Kaur	Lecturer, Lovely Professional University Phagwara, Punjab

Despite of our best efforts, some errors may have inadvertently crept into the book. Constructive comments and suggestions to further improve the book are welcome and shall be acknowledged gratefully.

Wishing you all the very best..!!!

—Trishna Knowledge Systems

Key Pedagogical Features

Learning Objectives
List of important topics which are covered in chapter.

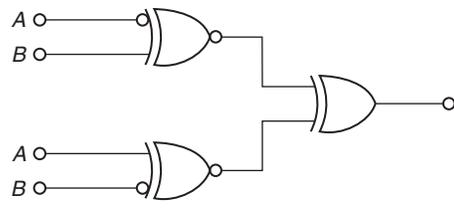
- LEARNING OBJECTIVES**
- ☞ Digital circuits
 - ☞ Number system with different base
 - ☞ Conversion of number systems
 - ☞ Complements
 - ☞ Subtraction with complement
 - ☞ Numeric codes
 - ☞ Weighted and non-weighted codes
 - ☞ Error detection and correction code
 - ☞ Sequential, reflective and cyclic codes
 - ☞ Self complementing code

Solved Examples

Example 1: Simplify the Boolean function, $xy + x'z + yz$

Solution: $xy + x'z + yz$
By using consensus property
 $xy + x'z + yz = xy + x'z$
 $Y = xy + x'z$

Example 2: The output of the given circuit is equal to



Solved Example
Solved problems given topic-wise to learn to apply the concepts learned in a particular section as per exam pattern.

Exercises
Practice problems for students to master the concepts studied in chapter. Exercises consist of two levels of problems “Practice Problem I” and “Practice Problem II” based on increasing difficulty level.

Exercises

Practice Problems I
Directions for questions 1 to 15: Select the correct alternative from the given choices.

- Assuming all the numbers are in 2's complement representation, which of the following is divisible by 11110110?
(A) 11101010 (B) 11100010
(C) 11111010 (D) 11100111
- If $(84)_x$ (in base x number system) is equal to $(64)_y$ (in base y number system), then possible values of x and y are
(A) 12, 9 (B) 6, 8
(C) 9, 12 (D) 12, 18
- Let $A = 1111\ 1011$ and $B = 0000\ 1011$ be two 8-bit signed 2's complement numbers. Their product in 2's complement representation is

Practice Problems 2
Directions for questions 1 to 20: Select the correct alternative from the given choices.

- The hexadecimal representation of $(567)_8$ is
(A) 1AF (B) D77
(C) 177 (D) 133
- $(2326)_8$ is equivalent to
(A) $(14D6)_{16}$ (B) $(103112)_4$
(C) $(1283)_{10}$ (D) $(09AC)_{16}$
- Signed 2's complement representation of $(-15)_{10}$ is
(A) 11111 (B) 10001
(C) 01111 (D) 10000
- $(0.25)_{10}$ in binary number system is?
(A) (0.01) (B) (0.11)
(C) 0.001 (D) 0.101
- The equivalent of $(25)_6$ in number system with base 7 is?
(A) 22 (B) 23

PREVIOUS YEARS' QUESTIONS

- (1217)₈ is equivalent to [2009]
 (A) (1217)₁₆ (B) (028F)₁₆
 (C) (2297)₁₀ (D) (0B17)₁₆
- P is a 16-bit signed integer. The 2's complement representation of P is (F87B)₁₆. The 2's complement representation of 8*P is [2010]
 (A) (C3D8)₁₆ (B) (187B)₁₆
 (C) (F878)₁₆ (D) (987B)₁₆
- The smallest integer that can be represented by an 8-bit number in 2's complement form is [2013]
 (A) -256 (B) -128
 (C) -127 (D) 0
- The base (or radix) of the number system such that the following equation holds is $\frac{312}{20} = 13.1$ [2014]
- Consider the equation (123)₅ = (x8)₉ with x and y as unknown. The number of possible solutions is [2014]
- Consider the equation (43)_x = (y3)₈ where x and y are unknown. The number of possible solutions is [2015]
- Define another random variable $Y = X_1 X_2 \oplus X_3$, where \oplus denotes XOR. Then $Pr[Y = 0 | X_3 = 0] =$ [2015]
- The 16-bit 2's complement representation of an integer is 1111 1111 1111 0101; its decimal representation is [2016]
- Consider an eight-bit ripple-carry adder for computing the sum of A and B, where A and B are integers represented in 2's complement form. If the decimal value of A is one, the decimal value of B that leads to the longest latency for the sum to stabilize is [2016]
- Let $x_1 \oplus x_2 \oplus x_3 \oplus x_4 = 0$ where x_1, x_2, x_3, x_4 are Boolean variables, and \oplus is the XOR operator. Which one of the following must always be TRUE? [2016]
 (A) $x_1 x_2 x_3 x_4 = 0$
 (B) $x_1 x_3 + x_2 = 0$
 (C) $\bar{x}_1 \oplus \bar{x}_3 = \bar{x}_2 \oplus \bar{x}_4$
 (D) $x_1 + x_2 + x_3 + x_4 = 0$
- Consider a quadratic equation $x^2 - 13x + 36 = 0$ with

Previous Years' Questions

Contains previous 10 years GATE Questions at end of every chapter which help students to get an idea about the type of problems asked in GATE and prepare accordingly.

HINTS/SOLUTIONS

Practice Problems I

- X-OR of two equal inputs will give you result as zero. Hence, the correct option is (B).
- Positive level OR means negative level AND vice versa. Hence, the correct option is (D).
- $\bar{A} \cdot \bar{C} \cdot D \cdot E \cdot F \cdot G \cdot H$
 (De Morgan's law)
 $= (\bar{A} + \bar{B})(\bar{C} + \bar{D})(\bar{E} + \bar{F})(\bar{G} + \bar{H})$

Then $Z = \bar{A} + B$
 Hence, the correct option is (B).

- Error → transmits odd number of one's, for both cases. Hence, the correct option is (A).
- $\sum(0, 1, 2, 4, 6)$ P should contain minterms of each function of x as well as y. Hence, the correct option is (B).
- $AB + ACD + \bar{A}C$
 $= AB(C + \bar{C})(D + \bar{D}) + A(B + \bar{B})CD + \bar{A}(B + \bar{B})C(D + \bar{D})$
 $= AB(CD + \bar{C}\bar{D} + \bar{C}D + C\bar{D}) + ABCD + \bar{A}BCD + \bar{A}C(BD + \bar{B}\bar{D} + B\bar{D} + \bar{B}\bar{D})$
 $ABCD + ABC\bar{D} + ABCD + ABC\bar{D} + \bar{A}BCD +$

Hints/Solutions

This section gives complete solutions of all the unsolved questions given in the chapter. The Hints/Solutions are included in the CD accompanying the book.

Practice Tests

Time-bound Test provided at end of each unit for assessment of topics learned in the unit.

TEST

Time: 60 min.

DIGITAL LOGIC

Directions for questions 1 to 30: Select the correct alternative from the given choices.

- What is the range of signed decimal numbers that can be represented by 4-bit 1's complement notation?
 (A) -7 to +7 (B) -16 to +16
 (C) -7 to +8 (D) -15 to +16
- Which of the following signed representation have a unique representation of 0?
 (A) Sign-magnitude (B) 1's complement
 (C) 0's complement (D) 2's complement
- Find the odd one out among the following
 (A) EBCDIC (B) GRAY
 (C) Hamming (D) ASCII

9. The number of product terms in the minimized SOP from is

1	0	0	1
0	D	0	0
0	0	D	1
1	0	0	1

(A) 2 (B) 4
 (C) 5 (D) 3

10. The minimum number of 2 input NAND gates needed to implement $Z = XY + VW$ is
 (A) 2 (B) 3
 (C) 4 (D) 5

Syllabus: Computer Science and Information Technology

Computer Science and Information Technology

Digital Logic: Boolean algebra. Combinational and sequential circuits. Minimization. Number representations and computer arithmetic (fixed and floating point).

Computer Organization and Architecture: Machine instructions and addressing modes. ALU, data-path and control unit. Instruction pipelining. Memory hierarchy: cache, main memory and secondary storage; I/O interface (interrupt and DMA mode).

Programming and Data Structures: Programming in C. Recursion. Arrays, stacks, queues, linked lists, trees, binary search trees, binary heaps, graphs.

Algorithms: Searching, sorting, hashing. Asymptotic worst case time and space complexity. Algorithm design techniques: greedy, dynamic programming and divide-and-conquer. Graph search, minimum spanning trees, shortest paths.

Theory of Computation: Regular expressions and finite automata. Context-free grammars and push-down automata. Regular and context-free languages, pumping lemma. Turing machines and undecidability.

Compiler Design: Lexical analysis, parsing, syntax-directed translation. Runtime environments. Intermediate code generation.

Operating System: Processes, threads, interprocess communication, concurrency and synchronization. Deadlock. CPU scheduling. Memory management and virtual memory. File systems.

Databases: ER model. Relational model: relational algebra, tuple calculus, SQL. Integrity constraints, normal forms. File organization, indexing (e.g., B and B+ trees). Transactions and concurrency control.

Computer Networks: Concept of layering. LAN technologies (Ethernet). Flow and error control techniques, switching. IPv4/IPv6, routers and routing algorithms (distance vector, link state). TCP/UDP and sockets, congestion control. Application layer protocols (DNS, SMTP, POP, FTP, HTTP). Basics of Wi-Fi. Network security: authentication, basics of public key and private key cryptography, digital signatures and certificates, firewalls.

Chapter-wise Analysis of GATE Previous Years' Papers

Subject	2005	2006	2007	2008	2009	2010	2011	2012	2013	2014	2015	2016	2017	2018	2019
Theory of Computation															
1 Mark	0	2	2	3	4	1	3	4	1	5	1	3	2	2	2
2 Marks	7	6	5	6	3	3	3	1	2	6	3	3	4	3	3
Total Marks	14	14	12	15	10	7	9	6	5	17	7	9	10	8	8
Digital Logic															
1 Mark	4	1	3	4	2	3	3	2	3	3	1	2	2	2	3
2 Marks	5	5	5	1	0	2	3	2	1	5	2	1	2	1	2
Total Marks	14	11	13	6	2	7	9	6	5	13	5	4	6	4	7
Computer Organization and Architecture															
1 Mark	4	0	2	0	2	1	3	2	1	2	1	2	3	3	3
2 Marks	9	7	6	12	4	4	2	4	7	2	2	2	4	4	1
Total Marks	22	14	14	24	10	9	7	10	15	6	5	6	11	11	5
Programming and Data Structures															
1 Mark	5	0	1	1	1	3	4	2	2	0	5	2	4	3	2
2 Marks	3	6	3	3	3	5	7	6	5	2	3	5	4	3	5
Total Marks	11	12	7	7	7	13	18	14	12	4	11	12	12	9	12
Algorithm															
1 Mark	2	8	3	2	3	1	1	4	5	1	4	4	2	1	2
2 Marks	10	7	12	15	6	3	0	2	3	2	4	5	2	5	3
Total Marks	22	22	27	32	15	7	1	8	11	5	12	14	6	11	8
Compiler Design															
1 Mark	1	1	1	2	1	2	1	1	2	1	2	1	2	1	2
2 Marks	5	5	5	2	0	1	0	3	2	2	1	1	1	2	2
Total Marks	11	11	11	6	1	4	1	7	6	5	4	3	4	5	6
Operating System															
1 Mark	0	1	2	2	2	3	3	1	1	0	2	1	2	2	2
2 Marks	2	8	6	5	5	2	2	3	1	2	2	4	2	3	4
Total Marks	4	17	14	12	12	7	7	7	3	4	6	9	6	8	10
Database															
1 Mark	3	1	0	1	0	3	0	3	1	3	1	3	2	2	2
2 Marks	4	4	6	5	5	2	3	3	4	2	2	1	3	2	3
Total Marks	11	9	12	11	5	7	6	9	9	7	5	5	8	6	8
Computer Networks															
1 Mark	5	1	2	1	0	2	5	3	4	4	2	2	2	3	1
2 Marks	2	5	6	4	5	3	2	3	2	2	3	4	3	2	4
Total Marks	9	11	14	9	5	8	9	9	8	8	8	10	8	7	9
Software Engineering															
1 Mark					1	0	1	0	0	1	1				
2 Marks					0	0	0	0	1	0	1				
Total Marks					1	0	1	0	2	1	3				
Web Technology															
1 Mark					1	0	1	0	0	0	1				
2 Marks					0	0	0	0	0	0	1				
Total Marks					1	0	1	0	0	0	3				

General Information about GATE

STRUCTURE OF GATE

The GATE examination consists of a single online paper of 3-hour duration, in which there will be a total of 65 questions carrying 100 marks out of which 10 questions carrying a total of 15 marks are in General Aptitude (GA).

Section Weightage and Marks

70% of the total marks is given to the technical section while 15% weightage is given to General Aptitude and Engineering Mathematics each.

Weightage	Questions (Total 65)	
Respective Engineering Branch	70 Marks	25—1 mark questions 30—2 mark questions
Engineering Maths	15 Marks	
General Aptitude	15 Marks	5—1 mark questions 5—2 mark questions

Particulars

For 1-mark multiple-choice questions, 1/3 marks will be deducted for a wrong answer. Likewise, for 2-mark multiple-choice questions, 2/3 marks will be deducted for a wrong answer. There is no negative marking for numerical answer type questions.

Question Types

- 1. Multiple Choice Questions (MCQ)** carrying 1 or 2 marks each in all papers and sections. These questions are objective in nature, and each will have a choice of four answers, out of which the candidate has to mark the correct answer.
- 2. Numerical Answer** carrying 1 or 2 marks each in all papers and sections. For numerical answer questions, choices will not be given. For these questions the answer is a real number, to be entered by the candidate using the virtual keypad. No choices will be shown for this type of questions.

Design of Questions

The fill in the blank questions usually consist of 35%– 40% of the total weightage.

The questions in a paper may be designed to test the following abilities:

- 1. Recall:** These are based on facts, principles, formulae or laws of the discipline of the paper. The candidate is expected to be able to obtain the answer either from his/her memory of the subject or at most from a one-line computation.
- 2. Comprehension:** These questions will test the candidate's understanding of the basics of his/her field, by requiring him/her to draw simple conclusions from fundamental ideas.
- 3. Application:** In these questions, the candidate is expected to apply his/her knowledge either through computation or by logical reasoning.
- 4. Analysis and Synthesis:** In these questions, the candidate is presented with data, diagrams, images etc. that require analysis before a question can be answered. A Synthesis question might require the candidate to compare two or more pieces of information. Questions in this category could, for example, involve candidates in recognising unstated assumptions, or separating useful information from irrelevant information.

About Online Pattern

The examination for all the papers will be carried out in an ONLINE Computer Based Test (CBT) mode where the candidates will be shown the questions in a random sequence on a computer screen. The candidates are required to either select the answer (for MCQ type) or enter the answer for numerical answer-type question using a mouse on a virtual keyboard (keyboard of the computer will be disabled). The candidates will also be allowed to use a calculator with which the online portal is equipped with.

IMPORTANT TIPS FOR GATE

The followings are some important tips which would be helpful for students to prepare for GATE exam

1. Go through the pattern (using previous year GATE paper) and syllabus of the exam and start preparing accordingly.
2. Preparation time for GATE depends on many factors, such as, individual's aptitude, attitude, fundamentals, concentration level etc., Generally rigorous preparation for 4 to 6 months is considered good but it may vary from student to student.
3. Make a list of books which cover complete syllabus, contains solved previous year questions and mock tests for practice based on latest GATE pattern. Purchase these books and start your preparation.
4. Make a list of topics which needs to be studied and make priority list for study of every topic based upon the marks for which that particular topic is asked in GATE exam. Find out the topics which fetch more marks and give more importance to those topics. Make a timetable for study of topics and follow the timetable strictly.
5. An effective way to brush up your knowledge about technical topics is group study with your friends. During group study you can explore new techniques and procedures.
6. While preparing any subject highlight important points (key definitions, equations, derivations, theorems and laws) which can be revised during last minute preparation.
7. Pay equal attention to both theory and numerical problems. Solve questions (numerical) based on latest exam pattern as much as possible, keeping weightage of that topic in mind. Whatever topics you decide to study, make sure that you know everything about it.
8. Try to use short-cut methods to solve problems instead of traditional lengthy and time consuming methods.
9. Go through previous year papers (say last ten years), to check your knowledge and note the distribution of different topics. Also analyze the topics in which you are weak and concentrate more on those topics. Always try to solve papers in given time, to obtain an idea how many questions you are able to solve in the given time limit.
10. Finish the detail study of topics one and a half month before your exam. During last month revise all the topics once again and clear leftover doubts.

GATE 2019 SOLVED PAPER

CS: COMPUTER SCIENCE AND INFORMATION TECHNOLOGY

Number of Questions: 46

Total Marks: 67

Wrong answer for MCQ will result in negative marks, (-1/3) for 1 Mark Questions and (-2/3) for 2 Marks Question.

Q.1–Q.25 carry one mark each.

Question Number: 1 **Question Type: MCQ**

Consider $Z = X - Y$, where X, Y and Z are all in sign-magnitude form. X and Y are each represented in n bits. To avoid overflow, the representation of Z would require a minimum of:

- (A) n bits (B) $n + 2$ bits
(C) $n + 1$ bits (D) $n - 1$ bits

Solution: $Z = X - Y$

X is n -bit sign magnitude number

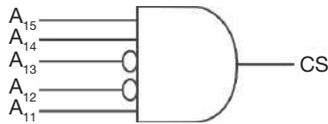
Y is n -bit sign magnitude number

To avoid overflow, the representation of Z would require a minimum of $n + 1$ bits

Hence, the correct option is (C).

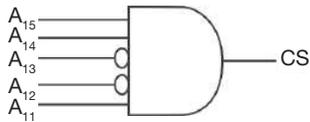
Question Number: 2 **Question Type: MCQ**

The chip select logic for a certain DRAM chip in a memory system design is shown below. Assume that the memory system has 16 address lines denoted by A_{15} to A_0 . What is the range of addresses (in hexadecimal) of the memory system that can get enabled by the chip select (CS) signal?



- (A) CA00 to CAFF (B) C800 to CFFF
(C) C800 to C8FF (D) DA00 to DFFF

Solution:



A_{15}	A_{14}	A_{13}	A_{12}	A_{11}	-	A_4	A_3	A_2	A_1	A_0	
1	1	0	0	1	-	0	0	0	0	0	C800H
					.						
					.						
1	1	0	0	1		1	1	1	1	1	CFFFH

Hence, the correct option is (B).

Question Number: 3 **Question Type: MCQ**

For $\Sigma = \{a, b\}$, let us consider the regular language $L = \{x \mid x = a^{2+3k} \text{ or } x = b^{10+12k}, k \geq 0\}$. Which one of the following can be a pumping length (the constant guaranteed by the pumping lemma) for L ?

- (A) 5 (B) 3
(C) 24 (D) 9

Solution: If L is a regular language, then there is a number P (the pumping length) such that S is any string in L of length P or more can be written as $S = xyz$, satisfying the following conditions.

for each $i \geq 0, xy^i z \in L,$

$|y| > 0$ and $|xy| \leq P$

So we need to find the minimum length string $s = xyz \in L$ such that $xy^i z$ should also be L .

$L = \{a^2, a^5, a^8, a^{11}, \dots\}$

Here we can take pumping length as 3.

OR

$L = \{b^{10}, b^{22}, b^{34}, \dots\}$

Here, the pumping length can be 12.

If we take pumping length as 24, in every repetition we will get multiple of 3 and 12.

Hence, the correct option is (C).

Question Number: 4 **Question Type: MCQ**

Which one of the following is NOT a valid identity?

- (A) $x \oplus y = x + y$, if $xy = 0$
(B) $(x \oplus y) \oplus z = x \oplus (y \oplus z)$
(C) $(x + y) \oplus z = x \oplus (y + z)$
(D) $x \oplus y = (xy + x'y)'$

Solution: $(x + y) \oplus z = x \oplus (y + z)$ is not a valid statement. Hence, the correct option is (C).

Question Number: 5 **Question Type: MCQ**

Which one of the following kinds of derivation is used by LR parsers?

- (A) Leftmost in reverse
(B) Leftmost
(C) Rightmost in reverse
(D) Rightmost

Solution: LR parser is a bottom up parser. Bottom up parser uses reverse of right most derivation. Hence, the correct option is (C).

Question Number: 6 **Question Type: NAT**

Consider a sequence of 14 elements; $A = [-5, -10, 6, 3, -1, -2, 13, 4, -9, -1, 4, 12, -3, 0]$. The subsequence sum $s(i, j)$, where $0 \leq i \leq j < 14$. (Divide and conquer approach may be used.)

Solution: $A[-5, -10, 6, 3, -1, -2, 13, 4, -9, -1, 4, 12, -3, 0]$

$$\text{Max}(S(i, j)) = S(2, 11) = \sum_{k=2}^{11} A[k] = 29$$

Hence, the correct answer is (29).

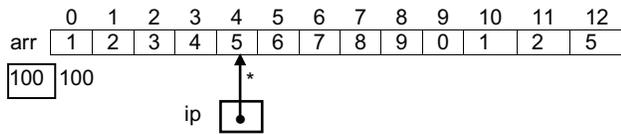
Question Number: 7 **Question Type: NAT**

Consider the following C program:

```
#include <stdio.h>
int main ( ) {
    int arr [ ] = {1, 2, 3, 4, 5, 6, 7,
        8, 9, 0, 1, 2, 5}, *ip = arr + 4;
    printf("%d\n", ip[1]);
    return 0;
}
```

The number that will be displayed on execution of the program is _____.

Solution:



Hence, the correct answer is (6).

Question Number: 8 **Question Type: MCQ**

Consider the following two statements about database transaction schedules:

- I. Strict two-phase locking protocol generates conflict serializable schedules that are also recoverable.
- II. Timestamp-ordering concurrency control protocol with Thomas' Write Rule can generate view serializable schedules that are not conflict serializable.

- (A) II only
- (B) I only
- (C) Both I and II
- (D) Neither I nor II

Solution: The strict two-phase locking protocol guarantees strict schedules (strict schedules are conflict serializable.)

The Timestamp-ordering concurrency control with Thomas, Write Rule does not enforce conflict serializability.

The given both statements are true.

Hence, the correct option is (C).

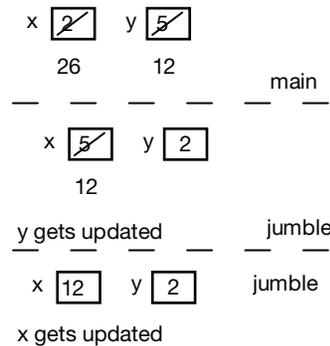
Question Number: 9 **Question Type: NAT**

Consider the following C program:

```
#include <stdio.h>
int jumble (int x, int y) {
    x = 2 *x + y;
    return x;
}
int main ( ) {
    int x = 2, y = 5;
    y = jumble (y, x);
    x = jumble (y, x);
    printf("%d \n", x);
    return 0;
}
```

The value printed by the program is _____.

Solution:



Hence, the correct answer is (26).

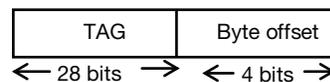
Question Number: 10 **Question Type: MCQ**

A certain processor uses a fully associative cache of size 16 kB. The cache block size 16 bytes. Assume that the main memory is byte addressable and uses a 32-bit address. How many bits are required for the Tag and the Index fields respectively in the addresses generated by the processor?

- (A) 28 bits and 4 bits
- (B) 28 bits and 0 bits
- (C) 24 bits and 4 bits
- (D) 24 bits and 0 bits

Solution: Cache block size = 16 Bytes
 $= 2^4 \text{ B} \Rightarrow \text{Block offset} = 4 \text{ bits}$

Fully associative cache size = 16 kB



In fully associative cache, there is no any index.

Indexing requires 0 bits and Tag bits = 28.

Hence, the correct option is (B).

Question Number: 11

Question Type: NAT

The value of $3^{51} \bmod 5$ is _____.

Solution: We have to find $3^{51} \bmod 5$

$$3^1 \bmod 5 = 3$$

$$3^2 \bmod 5 = 4$$

$$3^3 \bmod 5 = 2$$

$$3^4 \bmod 5 = 1$$

$$3^5 \bmod 5 = 3$$

$$3^6 \bmod 5 = 4$$

$$3^7 \bmod 5 = 2$$

$$3^8 \bmod 5 = 1$$

In general $3^{4m} \bmod 5 = 1$ for $m \in \mathbb{Z}^+$

$$\begin{aligned} \therefore 3^{51} \bmod 5 &= 3^{4 \times 12 + 3} \bmod 5 \\ &= 3^3 \bmod 5 = 2 \end{aligned}$$

Hence, the correct answer is (2).

Question Number: 12

Question Type: MCQ

Which one of the following statements is NOT correct about the B+ tree data structure used for creating an index of a relational database table?

- (A) Non-leaf nodes have pointers to data records
- (B) Key values in each node are kept in sorted order
- (C) B+ Tree is a height-balanced tree
- (D) Each leaf node has a pointer to the next leaf node

Solution: In a B+ tree, data pointers are stored only at the leaf nodes of the tree.

Hence, the correct option is (A).

Question Number: 13

Question Type: MCQ

If L is a regular language over $\Sigma = \{a, b\}$, which one of the following language is NOT regular?

- (A) Suffix (L) = $\{y \in \Sigma^* \mid \exists x \in \Sigma^* \text{ such that } xy \in L\}$
- (B) $\{\omega\omega^R \mid \omega \in L\}$
- (C) $L \cdot L^R = \{xy \mid x \in L, y^R \in L\}$
- (D) Prefix (L) = $\{x \in \Sigma^* \mid \exists y \in \Sigma^* \text{ such that } xy \in L\}$

Solution: The regular languages are closed under reversal and concatenation.

So, $L \cdot L^R$ is regular WW^R needs a memory (stack), therefore the language is not regular.

prefix (x) and suffix (x) is also regular.

Hence, the correct option is (B).

Question Number: 14

Question Type: NAT

Consider three concurrent processes P1, P2 and P3 as shown below, which access a shared variable D that has been initialized to 100.

P1	P2	P3
:	:	:
:	:	:
$D = D + 20$	$D = D - 50$	$D = D + 10$
:	:	:
:	:	:

The processes are executed on a uniprocessor system running a time-shared operating system. If the minimum and maximum possible values of D after the three processes have completed execution are X and Y respectively, then the value of $Y - X$ is _____.

Solution: Assembly code of the three process are

Process P₁:

- (1) Load R₁, M[D]
- (2) ADD R₁, #20
- (3) STORE M[D], R₁

Process P₂:

- (A) Load R₂, M[D]
- (B) SUB R₂, #50
- (C) STORE M[D], R₂

Process P₃:

- (X) Load R₃, M[D]
- (Y) ADD R₃, #10
- (Z) STORE M[D], R₃

Minimum value of D

Process P₂ executes (A), (B) instructions and got preempted, i.e., it did not store the value 50 to 'D'. Now, process P₁ and P₂ executes their instructions, i.e., (1), (2), (3) and (X), (Y), (Z). The value of D will be 130 and now, the instruction (C) is executed, it stores the value 50 to D.

Maximum value of D

Process P₁ first completed its execution, its value will be 120.

Process P₃ executes (X), (Y) and get preempted, P₂ completes its execution. Now, process P₃ completes (Z) instruction the value will be 130.

Therefore, value of $(Y - X)$ will be 80.

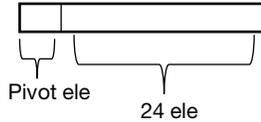
Hence, the correct answer is (80).

Question Number: 15

Question Type: NAT

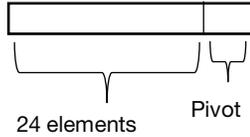
An array of 25 distinct elements is to be sorted using quick-sort. Assume that the pivot element is chosen uniformly at random. The probability that the pivot element gets placed in the worst possible location in the first round of partitioning (rounded off to 2 decimal place) is _____.

Solution: For a quick sort, worst case is when array is sorted. Let us consider the first element is considered as a pivot element



Or

The last element is considered as a pivot



Except first and last position all are either average case or best case, as list will get divided in to 2 sublists.

$$\text{Probability} = \frac{2}{25} = 0.08.$$

Hence, the correct answer is (0.08).

Question Number: 16 **Question Type: MCQ**

Which of the following protocol pairs can be used to send and retrieve e- mails (in that order)?

- (A) IMAP, SMTP
- (B) SMTP, MIME
- (C) SMTP, POP3
- (D) IMAM, POP3

Solution: Mails can be sent using SMTP and to retrieve emails POP3 is used.

Hence, the correct option is (C).

Question Number: 17 **Question Type: NAT**

The following C program is executed on a Unix/Linux system

```
#include <unistd.h>
int main( )
{
    int i;
    for (i = 0 ; i + +)
        if (i % 2 == 0 ) fork( ) :
    return 0;
}
```

The total number of child processes created is _____.

Solution: Number of times fork() get executed is 5 times i.e for i = 0, 2, 4, 6, 8.

Number of child process created for 'n' fork() calls are $2^n - 1$

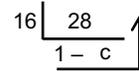
Hence, the correct answer is (31).

Question Number: 18 **Question Type: MCQ**

In 16-bit 2's complement representation, the decimal number - 28 is:

- (A) 1000 0000 1110 0100
- (B) 1111 1111 0001 1100
- (C) 0000 0000 1110 0100
- (D) 1111 1111 1110 0100

Solution:



$$(28)_{10} = (1C)_{16} = (00011100)_2$$

$$+ 28_{10} : 0000\ 0000\ 0001\ 1100$$

$$- 28_{10} : 1111\ 1111\ 1110\ 0100$$

Hence, the correct option is (D).

Question Number: 19 **Question Type: NAT**

Consider the grammar given below:

$$S \rightarrow Aa$$

$$A \rightarrow BD$$

$$B \rightarrow b| \in$$

$$D \rightarrow d| \in$$

Let a, b, d, and \$ be indexed as follows:

a	b	d	\$
3	2	1	0

Compute the FOLLOW set of the non-terminal B and write the index values for the symbols in the FOLLOW set in the descending order. (For example, if the FOLLOW set is {a, b, d, \$}, then the answer should be 3210)

Solution: Follow (B) = First (D) U

Follow (A) - { ∈ }

= {a, d}

It is indexed as 31.

Hence, the correct answer is (31).

Q.26–Q.55 carry two marks each.

Question Number: 20 **Question Type: NAT**

Let Σ be the set of all bijections from $\{1, \dots, 5\}$ to $\{1, \dots, 5\}$, where id denotes the identity function, i.e. $id(j) = j, \forall j$. Let \circ denotes composition on functions. For a string $x = x_1, x_2, \dots, x_n \in \Sigma^n, n \geq 0$, let $\pi(x) = x_1 \circ x_2 \circ \dots \circ x_n$. Consider the language $L = \{x \in \Sigma^* | \pi(x) = id\}$. The minimum number of states in any DFA accepting L is _____.

Solution: For a set of 5 elements number of bijections possible with itself are $5! = 120$.

Minimum number of states required for the language given in the question are 120.

Hence, the correct answer is (120).

Question Number: 21 **Question Type: NAT**

Consider the following relations P(X,Y,Z) Q(X,Y,T) and R(Y,V).

P		
X	Y	Z
X1	Y1	Z1
X1	Y1	Z2

Q		
X	Y	T
X2	Y1	2
X1	Y2	5

R	
Y	V
Y1	V1
Y3	V2

X2	Y2	Z2
X2	Y4	Z4

X1	Y1	6
X3	Y3	1

Y2	V3
Y2	V2

How many tuples will be returned by the following relational algebra query?

$$\Pi_x \left((\sigma_{(P,Y=R,Y \wedge R,V=V_2)}(P \times R)) - \Pi_x \left((\sigma_{(Q,Y=R,Y \wedge Q,T>2)}(Q \times R)) \right) \right)$$

Solution:

P × R				
X	Y	Z	Y	V
X ₁	Y ₁	Z ₁	Y ₁	✓ ₁
X ₁	Y ₁	Z ₁	Y ₃	V ₂
X ₁	Y ₁	Z ₁	V ₂	✓ ₃
X ₁	Y ₁	Z ₁	Y ₂	V ₂
X ₁	Y ₁	Z ₂	Y ₁	✓ ₁
X ₁	Y ₁	Y ₂	Y ₃	V ₂
X ₁	Y ₁	Z ₂	Y ₂	✓ ₃
X ₁	Y ₁	Z ₂	Y ₂	V ₂
X ₂	Y ₂	Z ₂	Y ₁	✓ ₁
X ₂	Y ₂	Z ₂	Y ₃	V ₂
X ₂	Y ₂	Z ₂	✓ ₂	✓ ₃
X ₂	Y ₂	Z ₂	Y ₂	V ₂
X ₂	Y ₄	Z ₄	Y ₁	✓ ₁
X ₂	Y ₄	Z ₄	Y ₃	V ₂
X ₂	Y ₄	Z ₄	Y ₂	✓ ₃
X ₂	Y ₄	Z ₄	Y ₂	V ₂

$$\pi_x (\sigma_{(P,Y=R,Y \wedge R,V=V_2)}(P \times R))$$

X
X ₂

$$\pi_x (\sigma_{(Q,Y=R,Y \wedge Q,T>2)}(Q \times R))$$

X
X ₁

$$X_2 - X_1 = X_2$$

Hence, the correct answer is (1).

Question Number: 22

Question Type: MCQ

Let the set of functional dependencies $F = \{QR \rightarrow S, R \rightarrow P, S \rightarrow Q\}$ hold on a relation schema $X = (PQRS)$. X is not in BCNF. Suppose X is decomposed into two schemas Y and Z , where $Y = (PR)$ and $Z = (QRS)$.

Consider the two statements given below.

- I. Both Y and Z are in BCNF.
- II. Decomposition of X and Y and Z is dependency preserving and lossless.

Which of the above statements is/are correct?

- (A) II only (B) Both I and II
(C) I only (D) Neither I nor II

Solution: $Y = PR$ and $Z = QRS$

$Y \cap Z = R$ (Key)

As FD $R \rightarrow P$ is present, R is key.

The given relation is lossless.

The relation Y is in BCNF but relation Z is not in BCNF as in $S \rightarrow Q$ S is not a superkey.

Hence, the correct option is (A).

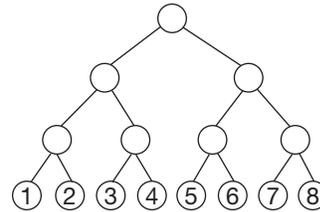
Question Number: 23

Question Type: NAR

Let T be a full binary tree with 8 leaves. (A full binary tree has every level full.) Suppose two leaves a and b of T are chosen uniformly and independently at random. The expected value of the distance between a and b in T (i.e., the number of edges in the unique path between a and b) is (rounded off to 2 decimal places) _____.

Solution:

Full binary tree with 8 leaf nodes



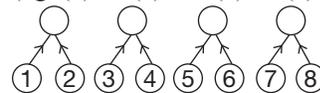
Expected value of the distance between a and $b = 0$

$$0 \times \frac{8}{64} + 2 \times \frac{8}{64} + 4 \times \frac{16}{64} + 6 \times \frac{32}{64} = \frac{272}{64} = 4.25$$

Out of 8 leaf nodes any 2 nodes can be chosen ${}^8P_2 + 8$ (selection of nodes with distance '0') = 64.

For length '2' between a and b will have '8' possibilities

(i.g. (1) \rightarrow (2) and (2) \rightarrow (1) are two different paths)



Hence, the correct answer is (4.25 to 4.25).

Question Number: 24

Question Type: NAT

In a RSA cryptosystem, the value of the public modulus parameter n is 3007. If it is also known that $\phi(n) = 2880$, where $\phi()$ denotes Euler's Totient function, then the prime factor of n which is greater than 50 is _____.

Solution: Let p, q be prime numbers.

Given

$$p * q = 3007$$

and

$$(p - 1)(q - 1) = 2880$$

$$pq - p - q + 1 = 2880$$

$$3007 - p - q + 1 = 2880$$

$$(p + q) = 128$$

$$p + \frac{3007}{p} = 128$$

$$p^2 + 3007 = 128p$$

$$p^2 - 128p + 3007 = 0$$

On solving

$$p = 97, 31$$

Hence, the correct answer is (97).

Question Number: 25

Question Type: MCQ

Assume that in a certain computer, the virtual addresses are 64 bits long and the physical addresses are 48 bits long. The memory is word addressable. The page size is 8kB and the word size is 4 bytes. The translation Look-aside Buffer (TLB) in the address translation path has 128 valid entries. At most how many distinct virtual addresses can be translated without any TLB miss?

- (A) 8×2^{20}
- (B) 16×2^{10}
- (C) 4×2^{20}
- (D) 256×2^{10}

Solution: Logical address = 64 bits

Physical address = 48 bits

Page size = 8KB

Word size = 4B

TLB entries = 128

$$\text{Number of entries in page} = \frac{8\text{KB}}{4\text{B}} = 2k$$

As the number of entries in TLB are 128. 128 entries translates 128 page numbers into frame numbers.

Distinct virtual addresses can be translated with TLB miss is

$$128 \times 2^{11}$$

$$256 \times 2^{10}$$

Hence, the correct option is (D).

Question Number: 26

Question Type: MCQ

Consider the first order predicate formula ϕ :

$\forall x [(\forall z z|x \Rightarrow ((z = x) \vee (z = 1))) \Rightarrow \exists w (w > x) \wedge (\forall z z|w \Rightarrow ((w = z) \vee (z = 1)))]$ Here ' $a|b$ ' denotes that ' a divides b ', where a and b are integers. Consider the following sets:

S1. $\{1, 2, 3, \dots, 100\}$

S2. Set of all positive integers

S3. Set of all integers

Which of the above sets satisfy ϕ ?

- (A) S1, S2 and S3
- (B) S1 and S3
- (C) S1 and S2
- (D) S2 and S3

Solution: Given

$$\phi : \forall x [(\forall z z|x \Rightarrow ((z = x) \vee (z = 1))) \Rightarrow \exists w (w > x) \wedge (\forall z z|w \Rightarrow ((w = z) \vee (z = 1)))]$$

Means, for every prime number, x , we can find another prime number w such that $w > x$.

Consider S1: $\{1, 2, 3, \dots, 100\}$

Take $x = 97$

Then there is no prime number w in S1 such that $w > x$

\therefore S1 does not satisfy ϕ

Clearly, S2 and S3 satisfy ϕ .

Hence, the correct option is (D).

Question Number: 27

Question Type: NAT

A relational database contains two tables Students and Performance as shown below:

Student		Performance		
Roll_No	Student name	Roll_No	Subject_code	Marks
1	Amit	1	A	86
2	Priya	1	B	95
3	Vinit	1	C	90
4	Rohan	2	A	89
5	Smita	2	C	92
		3	C	80

The primary key of the Student table is Roll_no. For the Performance table, the columns Roll_no and Subject_code together form the primary key. Consider the SQL query given below:

SELECT S.Student name, sum(P.Marks)

FROM Student S, Performance P

WHERE P.Marks > 84

GROUP BY S. Student_name;

The number of rows returned by the above SQL query is _____.

Solution:

Student_name	Sum (P.marks)
Amit	452
Priya	452
Rohan	452
Smita	452
Vinit	452

Hence, the correct answer is (5).

Question Number: 28

Question Type: MCQ

There are n unsorted arrays: A_1, A_2, \dots, A_n . Assume that n is odd. Each of A_1, A_2, \dots, A_n contains n distinct elements. There are no common elements between any two arrays. The worst-case time complexity of computing the median of A_1, A_2, \dots, A_n is

- (A) $O(n)$
- (B) $O(n \log n)$
- (C) $\Omega(n^2 \log n)$
- (D) $O(n^2)$

Solution: To compute the median of unsorted array of n elements, it takes $O(n)$ time. To find the medians of median, it will take $O(n^2)$ time.

Hence, the correct option is (D).

Question Number: 29 **Question Type: NAT**

Consider the following four processes with arrival times (in milliseconds) and their length of CPU bursts (in milliseconds) as shown below:

Process	P1	P2	P3	P4
Arrival time	0	1	3	4
CPU burst time	3	1	3	Z

These processes are run on a single processor using preemptive Shortest Remaining Time First scheduling algorithm. If the average waiting time of the processes is 1 millisecond, then the value of Z is _____.

Solution:

Process	P ₁	P ₂	P ₃	P ₄
Arrival time	0	1	3	4
Burse time	3	1	3	Z

Gantt chart

P ₁	P ₂	P ₃	P ₄
0	1	2	3
			4

at t = 4

Process P₁ and P₂ completes its execution. Waiting time of process P₁ and P₂ will be 1 and 0. Process P₃ waiting time is 1. If Process P₄ burst time can't be 3, then the average waiting time (AWS) exceeds 1 ms. It should be either 1 or 2. If it is 1, then AWS will be less than 1. So, process P₂ burst time will be 2. Then AWS will be 1.

Hence, the correct answer is (2).

Question Number: 30 **Question Type: MCQ**

Consider three machines M, N and P with IP addresses 100.10.5.2, 100.10.5.5 and 100.10.5.6 respectively. The subnet mask is set to 255.255.255.252 for all the three machines. Which one of the following is true?

- (A) Only M and N belong to the same subnet
- (B) Only N and P belong to the same subnet
- (C) M, N and P belong to three different subnets
- (D) M, N and P all belong to the same subnet

Solution: Subnet mask = 255.255.255.252

Machine M

100.10. 5. 0 0 0 0 0 0 1 0
 255.255. 255. 1 1 1 1 1 1 0 0

100.10. 5. 0

Machine N

100.10. 5. 0 0 0 0 0 1 0 1
 255.255. 255. 1 1 1 1 1 1 0 0

100.10. 5. 4

Machine P

100.10. 5. 0 0 0 0 0 1 1 0
 255.255. 255. 1 1 1 1 1 1 0 0

100.10. 5. 4

Machine N and P belong to the same subnet.

Hence, the correct option is (B).

Question Number: 31 **Question Type: MCQ**

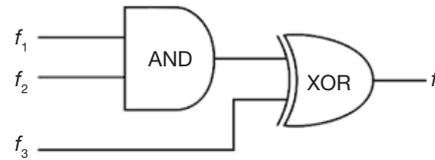
Consider three 4-variable functions $f_1, f_2,$ and $f_3,$ which are expressed in sum-of-minterms as

$$f_1 = \sum(0, 2, 5, 8, 14),$$

$$f_2 = \sum(2, 3, 6, 8, 14, 15),$$

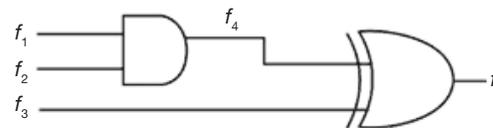
$$f_3 = \sum(2, 7, 11, 14)$$

For the following circuit with one AND gate and one XOR gate, the output function f can be expressed as:



- (A) $\sum(0, 2, 3, 5, 6, 7, 8, 11, 14, 15)$
- (B) $\sum(7, 8, 11)$
- (C) $\sum(2, 7, 8, 11, 14)$
- (D) $\sum(2, 14)$

Solution:



$$f_1 = \sum m(0, 2, 5, 8, 14)$$

$$f_2 = \sum m(2, 3, 6, 8, 14, 15)$$

$$f_4 = \sum m(2, 8, 14)$$

$$f_3 = \sum m(2, 7, 11, 14)$$

$$f = \sum m(7, 8, 11)$$

Hence, the correct option is (B).

Question Number: 32 **Question Type: NAT**

Consider the augmented grammar given below:

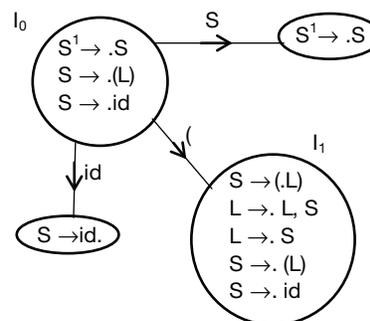
$$S' \rightarrow S$$

$$S \rightarrow (L) \mid id$$

$$L \rightarrow L, S \mid S$$

Let $I_0 = \text{CLOSURE}(\{[S' \rightarrow \cdot S]\})$. The number of items in the set GOTO($I_0, ()$) is: _____.

Solution:



Number of items in I_1 are 5.

Hence, the correct answer is (5).

Question Number: 33 **Question Type: MCQ**

Suppose that in a IP-over-Ethernet network, a machine X wishes to find the MAC address of another machine Y in its subnet. Which one of the following techniques can be used for this?

- (A) X send an ARP request packet to the local gateway's IP address which then finds the MAC address of Y and sends to X .
- (B) X send an ARP request packet to the local gateway's MAC address which then finds the MAC address of Y and sends to X .
- (C) X sends an ARP request packet with broadcast IP address in its local subnet.
- (D) X sends an ARP request packet with broadcast MAC address in its local subnet.

Solution: In IP over Ethernet network, If machine X wants to find MAC address of another machine Y in its subnet, then X sends an ARP request packet with broadcast MAC address in its local subnet.

Hence, the correct option is (D).

Question Number: 34 **Question Type: NAT**

The index node (inode) of Unix-like file system has 12 direct, one single – indirect and one double – indirect pointers. The disk block size is 4 kB, and the disk block address is 32 – bits long. The maximum possible file size is (rounded off to 1 decimal place) _____GB.

Solution: Maximum file size

$$= \left[\text{Number of direct pointers} + \left(\frac{\text{Disk block size}}{\text{Disk block address}} \right) + \left(\frac{\text{Disk block size}}{\text{disk block address}} \right) \right]$$

* (Disk Block size)

$$= \left[12 + \frac{4\text{KB}}{4\text{B}} + \left(\frac{4\text{KB}}{4\text{B}} \right)^2 \right] 4\text{KB}$$

$$= [12 + 4\text{K} + 1\text{M}] * 4\text{KB}$$

$$= 48\text{KB} + 16\text{MB} + 4\text{GB}$$

$$\cong 4\text{GB}$$

Hence, the correct answer is (3.7 to 3.8 or 4.0 to 4.1).

Question Number: 35 **Question Type: MCQ**

Consider the following snapshot of a system running n concurrent processes. Processes i is holding X_i instances of a resource R , $1 \leq i \leq n$. Assume that all instances of R are currently in use. Further, for all i , process i can place a request for at most Y_i additional instances of R while holding the X_i instances it already has. Of the n processes, there are exactly

two processes p and q such that $Y_p = Y_q = 0$. Which one of the following conditions guarantees that no other process apart from p and q can complete execution?

- (A) $\text{Min}(X_p, X_q) \leq \text{Max}\{Y_k \mid 1 \leq k \leq n, k \neq p, k \neq q\}$
- (B) $X_p + X_q < \text{Max}\{Y_k \mid 1 \leq k \leq n, k \neq p, k \neq q\}$
- (C) $X_p + X_q < \text{Min}\{Y_k \mid 1 \leq k \leq n, k \neq p, k \neq q\}$
- (D) $\text{Min}(X_p, X_q) \geq \text{Min}\{Y_k \mid 1 \leq k \leq n, k \neq p, k \neq q\}$

Solution: Given

$X_i \rightarrow$ Holding resources for process p_i

$Y_i \rightarrow$ Additional resources for process p_i

As process p and q doesn't require any additional resources, it completes its execution and available resources are $(X_p + X_q)$

There are $(n - 2)$ process p_i ($1 \leq i \leq n, i \neq p, q$) with their requirements as Y_i ($1 \leq i \leq n, i \neq p, q$). In order to not execute process p_i , no instance of Y_i should be satisfied with $(X_p + X_q)$ resources, i.e., minimum of Y_i instances should be greater than $(X_p + X_q)$.

Hence, the correct option is (C).

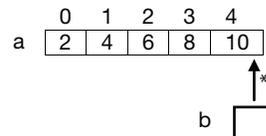
Question Number: 36 **Question Type: NAT**

Consider the following C program:

```
# include < stdio. h>
int main ( )
{
    Int a [ ] = {2, 4, 6, 8, 10}
    int i, sum = 0, *b = a + 4;
    for (i = 0; i < 5; i ++ )
        sum = sum + (*b - i) - *(b - i);
    printf("%d\n", sum);
    return 0;
}
```

The output of the above C program is ____

Solution:



$$i = 0$$

$$\text{sum} = 0 + 10 - 10 = 0$$

$$i = 1$$

$$\text{sum} = 0 + 9 - 8 = 1$$

$$i = 2$$

$$\text{sum} = 1 + 8 - 6 = 3$$

$$i = 3$$

$$\text{sum} = 3 + 7 - 4 = 6$$

$$i = 4$$

$$\text{sum} = 6 + 6 - 2 = 10$$

Hence, the correct answer is (10).

Question Number: 37 **Question Type: MCQ**

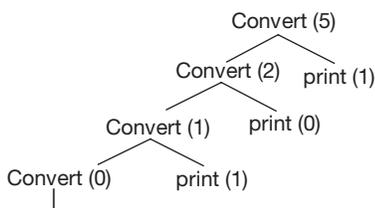
Consider the following C function.

```
void convert (int n) {
    if (n < 0)
        printf ("%d", n);
    else {
        convert (n/2);
        printf ("%d", n %2);
    }
}
```

Which one of the following will happen when the function convert is called with any positive integer n as argument?

- (A) It will print the binary representation of n and terminate.
- (B) It will print the binary representation of n but will not terminate.
- (C) It will not print anything and will not terminate.
- (D) It will print the binary representation of n in the reverse order and terminate.

Solution: Let, $n = 5$



As the program will terminate for negative value and if we divide zero by 2, it will give infinite value. The program will not terminate and doesn't print anything.

Hence, the correct option is (C).

Question Number: 38 **Question Type: NAT**

Consider that 15 machines need to be connected in a LAN using 8 – port Ethernet switches, assume that these switches do not have any separate uplink ports. The minimum number of switches needed is _____

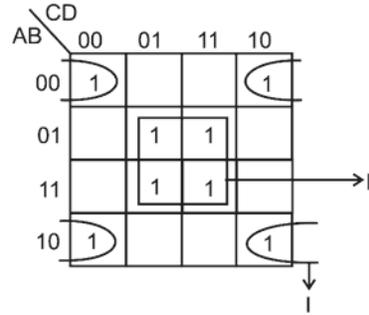
Solution: In an 8-port Ethernet switch, one port is used for networks connection and one port for the other switch, total 6 ports are used for connecting machines. For 15 machines, it requires 3 switches.

Hence, the correct answer is (3).

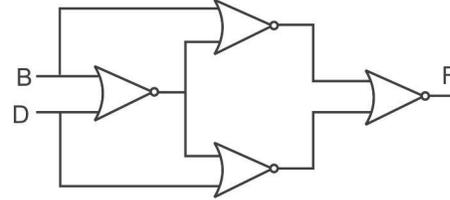
Question Number: 39 **Question Type: NAT**

What is the minimum number 2-input NOR gates required to implement a 4-variable function expressed in sum-of-minterms form as $F = \sum(0, 2, 5, 7, 8, 10, 13, 15)$?. Assume that all the inputs and their complements are available.

Solution:



$$F = I + II = \bar{B}\bar{D} + BD$$



$$F = \bar{B}\bar{D} + BD = B \odot D$$

Hence, to implement the given function in sum of minterms form, we need to take four 2-input NOR gates.

Hence, the correct answer is (3).

Question Number: 40 **Question Type: NAT**

A certain processor deploys a single-level cache. The cache block size is 8 words and the word size is 4 bytes. The memory system uses a 60-MHz clock. To service a cache miss, the memory controller first takes 1 cycle to accept the starting address of the block, it then takes 3 cycles to fetch all the eight words of the block, and finally transmits the words of the requested block at the rate of 1 word per cycle. The maximum bandwidth for the memory system when the program running on the processor issues a series of read operations is $___ \times 10^6$ bytes/sec.

Solution: Block is 8 words
words is 4 bytes

$$\text{Block size} = 8 \times 4 = 32 \text{ Bytes}$$

$$\text{To transfer a block from memory} = (1 + 3 + 8) = 12 \text{ clocks}$$

In 12 clocks, it transfers 32 bytes

$$12 \text{ clocks} \rightarrow 32 \text{ Bytes}$$

$$60 \times 10^6 \rightarrow ?$$

$$= \frac{60 \times 10^6 \times 32}{12}$$

$$= 160 \text{ bytes/sec}$$

Hence, the correct answer is (160).

Question Number: 41 **Question Type: MCQ**

Consider the following statements;

- I. The smallest element in a max-heap is always at a leaf node.

- II. The second largest element in a max-heap is always a child of the root node.
- III. A max-heap can be constructed from a binary search tree in $\theta(n)$ time.
- IV. A binary search tree can be constructed from a max-heap in $\theta(n)$ time.

Which of the above statements are TRUE?

- (A) I, II and III
- (B) I, III and IV
- (C) II, III and IV
- (D) I, II and IV

Solution: I. The smallest element in a max-heap is always present at a leaf node, it takes $\theta(n)$ time to find smallest element.

- II. The largest element in a max-heap is root element whereas second largest is either left child or right child.
- III. Construction of max-heap from binary search tree will take $\theta(n)$ time.
- IV. A binary search tree can not be constructed from a max-heap in $\theta(n)$ time.

Hence, the correct option is (A).

Question Number: 42 **Question Type: MCQ**

Consider the following C program:

```
#include <stdio.h>
int r ( ) {
    static int num = 7;
    return num --;
}
int main( ) {
    for ( r ( ); r ( ); r ( ))
        printf("%d", r ( ));
    return 0;
}
```

Which one of the following values will be displayed on execution of the programs?

- (A) 63
- (B) 52
- (C) 41
- (D) 630

Solution: r() // initialization

```
num  6
```

As post – decrement is present first the values ‘7’ is returned and then it will get decremented.

```
r()// condition
```

```
num  5
```

```
print (r()) num  4
```

Here ‘5’ will be printed first then it will be decremented.

```
r() // increment / decrement
```

```
num  3
r() // condition
```

```
num  2
```

As non-zero value is returned by function r(), print statement will get executed.

```
print (r()) num  1
```

Value ‘2’ will get printed first.

```
r()// increment/decrement
```

```
num  0
```

```
r() // condition
```

```
num  -1
```

As ‘zero’ is returned, the loop terminates,

The value printed is 52.

Hence, the correct option is (B).

Question Number: 43 **Question Type: NAT**

Consider the following C program:

```
#include <stdio.h>
int main( ) {
    float sum = 0.0, j = 1.0, i = 2.0;
    while (i/j >0.0625) {
        j = j + j;
        sum = sum + i/j;
        printf("%f\n", sum);
    }
    return 0;
}
```

The number of times the variable sum will be printed, when the above program is executed, is _____.

Solution:

```
Sum  0.0 j  1.0 i  2.0
```

```
2.0 > 0.0625
```

```
j = 2.0
```

```
Sum = 0 + 1.0
```

```
Print (1.0)//1
```

```
2.0/2.0 > 0.0625
```

```
j = 4.0
```

```
sum = 1.0 + 0.5 = 1.5
```

```
print (1.5)//2
```

```
4.0/2.0 > 0.0625
```

```
j = 8.0
```

```
sum = 1.5 + 0.25 = 1.75
```

```
print (1.75)//3
```

```
8.0/2.0 > 0.0625
```

```
j = 16.0
```

```
sum = 1.75 + 0.125 = 1.875
```

```
print (1.875)//4
16.0/2.0 > 0.0625
j = 32.0
sum = 1.875 + 0.0625 = 1.9375
print (1.9375)//5
```

The sum get printed 5 times.
Hence, the correct answer is (5).

Question Number: 44 **Question Type: MCQ**

Consider the following sets:

- S1. Set of all recursively enumerable languages over the alphabet $\{0, 1\}$
- S2. Set of all syntactically valid C programs
- S3. Set of all languages over the alphabet $\{0, 1\}$
- S4. Set of all non-regular languages over the alphabet $\{0, 1\}$

Which of the above sets are uncountable?

- (A) S2 and S3
- (B) S1 and S2
- (C) S3 and S4
- (D) S1 and S4

Solution: Recursively enumerable languages are countable.

Syntactically valid C program can be represented with CFG. CFG generates CFL, CFL is countable.

All languages over $\{0, 1\}$ may not be countable.

Set of regular languages are countable, non-regular languages may not be countable.

Hence, the correct option is (C).

Question Number: 45 **Question Type: MCQ**

Consider the following grammar and the semantic actions to support the inherited type declaration attributes. Let X_1, X_2, X_3, X_4, X_5 and X_6 be the placeholders for the non-terminals D, T, L or L_1 in the following table:

Production rule	Semantic action
$D \rightarrow TL$	$X_1.type = X_2.type$
$T \rightarrow int$	$T.type = int$

$T \rightarrow float$	$T.type = float$
$L \rightarrow L_1, id$	$X_3.type = X_4.type$ add Type (id entry, $X_5.type$)
$L \rightarrow id$	Add Type (id.entry, $X_6.type$)

Which one of the following are the appropriate choice for $X_1, X_2, X_3,$ and X_4 ?

- (A) $X_1 = T, X_2 = L, X_3 = L_1, X_4 = T$
- (B) $X_1 = L, X_2 = T, X_3 = L_1, X_4 = L$
- (C) $X_1 = T, X_2 = L, X_3 = T, X_4 = L_1$
- (D) $X_1 = L, X_2 = L_1, X_3 = L_1, X_4 = T$

Solution: Given, inherited attributes are evaluated by bottom up evaluation.

For production

$D \rightarrow TL$

The semantic action will be

$L.type = T.type$

As the L.type is decided by the T.type Similarly, $L_1.Type$ is decided by L.type, for the production

$L \rightarrow L_1, id$

\therefore the values of X_1, X_2, X_3 and X_4 are

L, T, L_1 and L

Hence, the correct option is (B).

Question Number: 46 **Question Type: MCQ**

Which one of the following languages over $\Sigma = \{a, b\}$ is NOT context-free?

- (A) $\{ww^R \mid w \in \{a, b\}^*\}$
- (B) $\{wa^n b^n w^R \mid w \in \{a, b\}^*, n \geq 0\}$
- (C) $\{wa^n w^R b^n \mid w \in \{a, b\}^*, n \geq 0\}$
- (D) $\{a^n b^i \mid i \in \{n, 3n, 5n\}, n \geq 0\}$

Solution: It is not possible to draw a PDA for language

$L = \{wa^n w^R b^n \mid w \in \{a, b\}^*, n \geq 0\}$

Hence, the correct option is (C).

This page is intentionally left blank

Digital Logic

Chapter 1: Number Systems	1.3
Chapter 2: Boolean Algebra and Minimization of Functions	1.14
Chapter 3: Combinational Circuits	1.35
Chapter 4: Sequential Circuits	1.56

U

N

I

T

1

This page is intentionally left blank

Chapter 1

Number Systems

LEARNING OBJECTIVES

- ☞ Digital circuits
- ☞ Number system with different base
- ☞ Conversion of number systems
- ☞ Complements
- ☞ Subtraction with complement
- ☞ Numeric codes
- ☞ Weighted and non-weighted codes
- ☞ Error detection and correction code
- ☞ Sequential, reflective and cyclic codes
- ☞ Self complementing code

DIGITAL CIRCUITS

Computers work with binary numbers, which use only the digits '0' and '1'. Since all the digital components are based on binary operations, it is convenient to use binary numbers when analyzing or designing digital circuits.

Number Systems with Different Base

Decimal number system

Decimal numbers are usual numbers which we use in our day-to-day life. The base of the decimal number system is 10. There are ten numbers 0 to 9.

The value of the n th digit of the number from the right side = n th digit \times (base) $^{n-1}$

Example 1: $(99)_{10} \rightarrow 9 \times 10^1 + 9 \times 10^0$
 $= 90 + 9 = 99$

Example 2: $(332)_{10} \rightarrow 3 \times 10^2 + 3 \times 10^1 + 2 \times 10^0$
 $= 300 + 30 + 2$

Example 3: $(1024)_{10} \rightarrow 1 \times 10^3 + 0 \times 10^2 + 2 \times 10^1 + 4 \times 10^0$
 $= 1000 + 0 + 20 + 4 = 1024$

Binary number system

In binary number system, there are only two digits '0' and '1'. Since there are only two numbers, its base is 2.

Example 4: $(1101)_2 = 1 \times 2^3 + 1 \times 2^2 + 0 \times 2^1 + 1 \times 2^0$
 $= 8 + 4 + 1 = (13)_{10}$

Octal number system

Octal number system has eight numbers 0 to 7. The base of the number system is 8. The number $(8)_{10}$ is represented by $(10)_8$.

Example 5: $(658)_8 = 6 \times 8^2 + 5 \times 8^1 + 8 \times 8^0$
 $= 384 + 40 + 8 = (432)_{10}$

Hexadecimal number system

In hexadecimal number system, there are 16 numbers 0 to 9, and digits from 10 to 15 are represented by A to F, respectively. The base of hexadecimal number system is 16.

Example 6: $(1A5C)_{16} = 1 \times 16^3 + A \times 16^2 + 5 \times 16^1 + C \times 16^0$
 $= 1 \times 4096 + 10 \times 256 + 5 \times 16 + 12 \times 1$
 $= 4096 + 2560 + 80 + 12 = (6748)_{10}$

Table 1 Different number systems

Decimal	Binary	Octal	Hexadecimal
0	000	0	0
1	001	1	1
2	010	2	2
3	011	3	3
4	100	4	4
5	101	5	5
6	110	6	6
7	111	7	7
8	1000	10	8
9	1001	11	9
10	1010	12	A
11	1011	13	B
12	1100	14	C
13	1101	15	D
14	1110	16	E
15	1111	17	F

(Continued)

Table 1 (Continued)

Decimal	Binary	Octal	Hexadecimal
16	10000	20	10
17	10001	21	11
18	10010	22	12
19	10011	23	13
20	10100	24	14

1. For a number system with base n , the number of different symbols in the number system will be n . Example: octal number system will have total of 8 numbers from 0 to 7.
2. The number ‘ n ’ in the number system with base ‘ n ’ is represented as $(10)_n$.
3. The equivalent of number $(a_3a_2a_1a_0 \cdot a_{-1}a_{-2})_n$ in decimal is $a_3 \times n^3 + a_2 \times n^2 + a_1 \times n^1 + a_0 \times n^0 + a_{-1} \times n^{-1} + a_{-2} \times n^{-2}$.

Conversion of Number Systems

The conversion of decimal to any other number system involves successive division by the radix until the dividend reaches 0. At each division, the remainder gives a digit of converted number; and the last one is most significant digit, the remainder of the first division is least significant digit.

The conversion of other number system to decimal involves multiplying each digit of number system with the weight of the position (in the power of radix) and sum the products calculated, the total is the equivalent value in decimal.

Decimal to binary conversion

Example 7: $(66)_{10}$

$$\begin{array}{r}
 2 \overline{) 66} \\
 \underline{2 \ 33} \ 0 \\
 2 \overline{) 16} \ 1 \\
 \underline{2 \ 8} \ 0 \\
 2 \overline{) 4} \ 0 \\
 \underline{2 \ 2} \ 0 \\
 \underline{1 \ 0} \
 \end{array}
 \begin{array}{l}
 \text{Reading remainders} \\
 \text{from bottom to top}
 \end{array}$$

$= (1000010)_2$

Example 8: $(928)_{10}$

$$\begin{array}{r}
 2 \overline{) 928} \\
 \underline{2 \ 464} \ 0 \\
 2 \overline{) 232} \ 0 \\
 \underline{2 \ 116} \ 0 \\
 2 \overline{) 58} \ 0 \\
 \underline{2 \ 29} \ 0 \\
 2 \overline{) 14} \ 1 \\
 \underline{2 \ 7} \ 0 \\
 2 \overline{) 3} \ 1 \\
 \underline{1 \ 1} \
 \end{array}$$

$= (1110100000)_2$

Example 9: $(105.75)_{10}$

$$\begin{array}{r}
 2 \overline{) 105} \\
 \underline{2 \ 52} \ 1 \\
 2 \overline{) 26} \ 0 \\
 \underline{2 \ 13} \ 0 \\
 2 \overline{) 6} \ 1 \\
 \underline{2 \ 3} \ 0 \\
 \underline{1 \ 1} \
 \end{array}$$

$(105)_{10} = (1101001)_2$

$(0.75)_{10}$

Multiply 0.75 by 2 = 1.50

Multiply 0.50 by 2 = 1.00

Reading integers from top to bottom $0.75 = (0.11)_2$

$\therefore (105.75)_{10} = (1101001.11)_2$

Binary to decimal conversion

Example 10: $(10100011)_2$

$$\begin{aligned}
 &= 1 \times 2^7 + 0 \times 2^6 + 1 \times 2^5 + 0 \times 2^4 + 0 \times 2^3 + 0 \times 2^2 + 1 \\
 &\quad \times 2^1 + 1 \times 2^0 \\
 &= 128 + 0 + 32 + 0 + 0 + 0 + 2 + 1 \\
 &= (163)_{10}
 \end{aligned}$$

Example 11: $(11010011.101)_2$

$$\begin{aligned}
 &= 1 \times 2^7 + 1 \times 2^6 + 0 \times 2^5 + 1 \times 2^4 + 0 \times 2^3 + 0 \times 2^2 + 1 \\
 &\quad \times 2^1 + 1 \times 2^0 + (1 \times 2^{-1}) + (0 \times 2^{-2}) + (1 \times 2^{-3}) \\
 &= 128 + 64 + 0 + 16 + 0 + 0 + 2 + 1 + 0.5 + 0 + 0.125 \\
 &= (211.625)_{10}
 \end{aligned}$$

Decimal to octal conversion

Example 12: $(16)_{10}$

$$\begin{array}{r}
 8 \overline{) 16} \\
 \underline{2} \ 0
 \end{array}$$

Remainder from bottom to top = $(20)_8$

Example 13: $(347.93)_{10}$

$$\begin{array}{r}
 (93)_{10} \\
 0.93 \times 8 = 7.44 \\
 0.44 \times 8 = 3.52 \\
 0.52 \times 8 = 4.16 \\
 0.16 \times 8 = 1.28 \\
 \dots\dots
 \end{array}$$

Read the integers of octal point from top to bottom.

$\therefore (0.93)_{10} = (0.7341)_8$

$(347)_{10}$

$$\begin{array}{r}
 8 \overline{) 347} \ 3 \\
 \underline{8 \ 43} \ 3 \\
 \ 5
 \end{array}$$

$\therefore (347)_{10} = (533)_8$

Ans: $(533.7341)_8$

Octal to decimal conversion

Example 14: $(33)_8$
 $3 \times 8^1 + 3 \times 8^0 = 24 + 3$
 $(27)_{10}$

Example 15: $(1023.06)_8$
 $1 \times 8^3 + 0 \times 8^2 + 2 \times 8^1 + 3 \times 8^0 + 0 \times 8^{-1} + 6 \times 8^{-2}$
 $= 512 + 0 + 16 + 3 + 0 + 0.0937 = (2095.0937)_{10}$

Octal to binary conversion

To convert octal to binary, replace each octal digit with their equivalent 3-bit binary representation.

Example 16: $(7777)_8$
 Convert each octal digit to binary
 $= \frac{7}{111} \frac{7}{111} \frac{7}{111} \frac{7}{111}$
 $= (111\ 111\ 111\ 111)_2$

Example 17: $(567.62)_8$
 $\begin{array}{ccc} 5 & 6 & 7 \\ 101 & 110 & 111 \end{array} . \begin{array}{cc} 6 & 2 \\ 110 & 010 \end{array}$
 $= (101110111.110010)_2$

Binary to octal conversion

To convert a binary number to an octal number, starting from the binary point, make groups of 3-bits each on either side of the binary point, and replace each 3-bit binary group by the equivalent octal digit.

Example 18: $(010011101)_2$
 $\frac{010}{2} \frac{011}{3} \frac{101}{5} = (235)_8$

Example 19: $(10010111011.1011)_2$
 $\frac{010}{2} \frac{010}{2} \frac{111}{7} \frac{011}{3} . \frac{101}{5} \frac{100}{4} = (2273.54)_8$

Decimal to hexadecimal conversion

Example 20: $(527)_{10}$

16		527	
16		32	15
		2	0

Decimal		Hexa
2	→	2
0	→	0
15	→	F

$= (20F)_{16}$

Example 21: $(18.675)_{10}$
 $(18)_{10}$

16		18	
		1	2

Decimal		Hexa
1	-	1
2	-	2

$(18)_{10} = (12)_{16}$
 $(0.675)_{10}$

0.675 × 16		10.8
0.800 × 16		12.8
0.800 × 16		12.8
0.800 × 16		12.8

Decimal	Hexa
10	A
12	C
12	C
12	C

$= (0.ACCC)_{16}$
 \therefore Hexadecimal equivalent is
 $= (12.AC\ CC)_{16}$

Hexadecimal to decimal conversion

Example 22: $(A3F)_{16}$

Decimal	Hexa
A	- 10
3	- 3
F	- 15

$\rightarrow 10 \times 16^2 + 3 \times 16^1 + 15 \times 16^0$
 $\rightarrow 2560 + 48 + 15 \rightarrow (2623)_{10}$

Example 23: $(1F63.0EB)_{16}$

1	1
F	15
6	6
3	3
0	0
E	14
B	11

$\rightarrow 1 \times 16^3 + 15 \times 16^2 + 6 \times 16^1 + 3 \times 16^0 + (0 \times 16^{-1})$
 $+ (14 \times 16^{-2}) + (11 \times 16^{-3})$
 $\rightarrow 4096 + 3840 + 96 + 3 + 0 + 0.0546 + 0.0026$
 $\rightarrow (8035.0572)_{10}$

Hexadecimal to binary number system

To represent hexadecimal in binary, represent each HEX number with its 4-bit binary equivalent.

Example 24: $(34F)_{16}$

Hexa	Decimal	Binary
3	3	0011
4	4	0100
F	15	1111

$= (001101001111)_2$

Example 25: $(AFBC \cdot BED)_{16}$

Hexa	Decimal	Binary
A	10	1010
F	15	1111
B	11	1011
C	12	1100
B	11	1011
E	14	1110
D	13	1101

$= (1010111110111100.101111101101)_2$

Binary to hexadecimal number system

To convert binary number to a hexadecimal number, starting from the binary point, make groups of 4-bits each on either side of the binary point and replace each 4-bit group by the equivalent hexadecimal digit.

Example 26: $(11001001)_2$
 $\rightarrow \frac{1100}{12} \frac{1001}{9}$
 $\rightarrow (C9)_{16}$

Example 27: $(1011011011.01111)_2$
 $\frac{0010}{2} \frac{1101}{D} \frac{1011}{B} \frac{0111}{7} \frac{1000}{8} = (2DB.78)_{16}$

Hexadecimal to octal number system

The simplest way to convert hexadecimal to octal is, first convert the given hexadecimal number to binary and then the Binary number to Octal.

Example 28: $(C3AF)_{16}$
 $\rightarrow 001100001110101111$
 $\rightarrow (141657)_8$

Example 29: $(C6.AE)_{16}$
 $\rightarrow 0011000110.10101110$
 $\rightarrow (306.534)_8$

Octal to hexadecimal number system

The simplest way to convert octal to hexadecimal is first convert the given octal number to binary and then the binary number to hexadecimal.

Example 30: $(775)_8$
 $\rightarrow (000111111101)_2$
 $\rightarrow (1FD)_{16}$

Example 31: $(34.7)_8$
 $\rightarrow (00011100.1110)_2$
 $\rightarrow (1C.E)_{16}$

COMPLEMENTS

Complements are used in digital computers to simplify the subtraction operation and for logical manipulation.

There are two types of complements for each base r -system.

1. Radix complement (or) r 's complement: the r 's complement of an m digit number N in base r is $r^m - N$ for $N \neq 0$.

For example, $N = 0$, r 's complement is 0.

2. Diminished radix complement: (or) $(r - 1)$'s complement: Given a number N in base r having m digits, then $(r - 1)$'s complement is $(r^m - 1) - N$.

For example, decimal number system will have 10's complement and 9's complement.

Similarly, binary number system will have 2's complement and 1's complement.

Example 32: 10's complement of $(2657)_{10}$ is $(10^4) - 2657$

$$\begin{array}{r} 10000 \\ - 2657 \\ \hline 7343 \end{array}$$

Example 33: 9's complement of $(2657)_{10}$ is $(10^4 - 1) - 2657$

$$\begin{array}{r} 10000 \\ - 1 \\ \hline 9999 \\ - 2657 \\ \hline 7342 \end{array}$$

- r 's complement can be obtained by adding 1 to $(r - 1)$'s complement.

$$r^m - N = \{(r^m - 1) - N\} + 1$$

Example 34: 2's complement of $(101101)_2$ is

$$\begin{array}{r} = (2^6) - 101101 \\ (2^6)_{10} = (100000)_2 \\ 2\text{'s complement is } 100000 \\ - 101101 \\ \hline 010011 \end{array}$$

Example 35: 1's complement of $(101101)_2$ is

$$2^6 - 1 = 1000000$$

$$\begin{array}{r} - 1 \\ 11111 \\ \hline 101101 \end{array}$$

1's complement -101010

The one's complement of a binary number is formed by changing 1's to 0's and 0's to 1's, The 2's complement can be formed by leaving all least significant 0's and the first 1 unchanged, and replacing 1's with zeros and zeros with 1's in all other bits.

If the number M contains radix point, the point should be removed temporarily in order to form r 's/ $(r - 1)$'s complement.

The radix point is then restored to the complemented number in the same relative position.

Example 36: What is 1's complement of $(1001.011)_2$?

\rightarrow Consider without radix point 1001011

Take 1's complement 0110100

Place radix point again $(0110.100)_2$

Example 37: What is 2's complement of $(1001.011)_2$?

Consider without radix point 1001011

Take 2's complement 0110101

Place radix point again $(0110.101)_2$

Complement of a complement is equal to the original number $r^m - (r^m - M) = M$

Subtraction with Complements

Subtraction of two n digit unsigned numbers $A - B$ in base r can be done as follows by r 's complement method.

Add A to the r 's complement of B . Mathematically $A + (r^n - B) = A - B + r^n$

If $A \geq B$ the sum will produce an end carry r^n ; which can be discarded. (Discarding carry is equivalent to subtracting r^n from result). What is left is the result $A - B$?

$$\begin{array}{r}
 A = 1100 \rightarrow 1100 \\
 B = 1010 \quad (2\text{'s complement}) + \underline{0110} \\
 \text{Sum: } 10010 \\
 \text{discard carry } (-r^n) \quad - \underline{10000} \\
 A - B: \underline{0010}
 \end{array}$$

If $A < B$, the sum does not produce an end carry and result is $r^n - (B - A)$. Then take r 's complement of the sum, and place a negative sign in front.

$$\begin{array}{r}
 \text{If } A = 1010 \\
 B = 1100 \\
 A - B \text{ can be done as} \\
 A \rightarrow 1010 \\
 B \rightarrow 2\text{'s complement} + \underline{0100} \\
 \text{Sum: } 1110
 \end{array}$$

Here, no carry generated, so result is a negative number. 2's complement of result $\rightarrow 0010 = 2$
result = -2

Subtraction of unsigned numbers by using $(r - 1)$'s complement can be done in similar way. However, $(r - 1)$'s complement is one less than the r 's complement. Because of this, the sum produced is one less than the correct difference when an end carry occurs. So end carry will be added to the sum. Removing the end carry and adding 1 to the sum is referred to as an end-around-carry.

$$\begin{array}{r}
 \text{Consider } A = 1100, B = 1010 \\
 \text{For } A - B \\
 A \rightarrow 1100 \\
 B \rightarrow (1\text{'s complement}) + \underline{0101} \\
 \text{Sum: } \underline{10001} \\
 \text{End around carry } + \xrightarrow{1} \\
 A - B = 0010
 \end{array}$$

$$\begin{array}{r}
 \text{For } B - A \\
 B \rightarrow 1010 \\
 A \rightarrow (1\text{'s complement}) + \underline{0011} \\
 \text{Sum: } \underline{1101}
 \end{array}$$

There is no end carry, for there result is $-(B - A) = -(1\text{'s complement of } 1101) = -0010 = -2$

Signed Binary Numbers

Positive integers can be represented as unsigned numbers; but to represent negative integer, we need a notation for negative values in binary.

It is customary to represent the sign with a bit placed in the left most position of the number. The convention is to make the sign bit 0 for positive and 1 for negative. This representation of signed numbers is referred to as sign-magnitude convention

S Magnitude

$$\begin{array}{l}
 +24 \text{ is } \underline{0} \underline{11000} \\
 \text{sign magnitude} \\
 -24 \text{ is } \underline{1} \underline{11000} \\
 \text{sign magnitude}
 \end{array}$$

Other notation for representation of signed numbers is signed complement system. This is convenient to use in a computer for arithmetic operations. In this system, a negative number is indicated by its complement (i.e., complement of corresponding positive number) whereas the sign-magnitude system negates a number by changing its sign bit, the signed-complement system negates a number by taking its complement. Positive numbers use same notation in sign-magnitude as well as sign-complement systems.

The signed-complement system can be used either as the 1's complement or the 2's complement.

But 2's complement is the most common.

+24 in 1's/2's complement representation is 011000

-24 in 1's complement representation 100111

-24 in 2's complement representation 101000

Table 2 Signed binary numbers - (4-bits)

Decimal	Signed-Magnitude	Signed 1's Complement	Signed 2's Complement
+7	0111	0111	0111
+6	0110	0110	0110
+5	0101	0101	0101
+4	0100	0100	0100
+3	0011	0011	0011
+2	0010	0010	0010
+1	0001	0001	0001
+0	0000	0000	0000
-0	1000	1111	-
-1	1001	1110	1111
-2	1010	1101	1110
-3	1011	1100	1101
-4	1100	1011	1100
-5	1101	1010	1011
-6	1110	1001	1010
-7	1111	1000	1001
-8	-	-	1000

The ranges of signed binary numbers with n-bits

Signed-magnitude: $-2^{n-1} + 1$ to $+2^{n-1} - 1$

1's complement representation: $-2^{n-1} + 1$ to $+2^{n-1} - 1$

2's complement representation: -2^{n-1} to $+2^{n-1} - 1$

Signed 2's complement representation can be directly used for arithmetic operations. The carryout of the sign bit position is discarded.

In order to obtain a correct answer, we must ensure that the result has a sufficient number of bits to accommodate the sum/product.

Example 38: $X = 00110$, $Y = 11100$ are represented in 5-bit signed 2's complement system

Then their sum $X + Y$ in 6-bit signed 2's complement representation is?

Solution: $X = 00110$
 $Y = 11100$

are 5-bit numbers

But result needs to be in 6-bit format.

Operands X and Y also should be in 6-bit format

$$X = \quad 000110$$

$$Y = \quad \underline{111100}$$

$$X + Y = (1) 000010$$

The carry out of sign bit position is discarded result is 000010.

Example 39: $(36x70)_{10}$ is 10's complement of $(yz0)_{10}$. Then values of x, y, z are?

- (A) 4, 5, 2 (B) 4, 6, 3
 (C) 3, 6, 3 (D) 3, 5, 4

Solution: (C)

$(36x70)_{10}$ is 10's complement of $(yz0)_{10}$.

10's complement of $(yz0)_{10}$ is

$$10^5 - yz0 = 36 \times 70$$

$$\text{So } 36x70 + yz0 = 100000$$

$$\begin{array}{r} 36x70 \\ +yz0 \\ \hline 100000 \end{array}$$

$$100000$$

$$\text{so } 7 + z = 10,$$

$$1 + x + y = 10 \quad z = 3$$

$$1 + 6 + z = 10 \quad y = 6$$

$$1 + 3 + y = 10,$$

$$\rightarrow x = 3$$

Example 40: The 10's complement of $(843)_{11}$ is?

- (A) $(157)_{11}$ (B) $(267)_{11}$
 (C) $(156)_{11}$ (D) $(268)_{11}$

Solution: (B)

Given $(843)_{11}$ is base 11 number system and the number in the number system range from 0 to 9 & A ($A = 10$)

10's complement for $(843)_{11}$ means $(r - 1)$'s complement.

$$\text{So } (r^n - 1) - N = [(11)^n - 1] - N$$

$$(11)^n - 1 \Rightarrow 1000$$

$$\begin{array}{r} -1 \\ AAA \\ \hline -843 \\ 267 \end{array}$$

$$AAA$$

$$-843$$

$$267$$

10's complement is $(267)_{11}$

Example 41: Consider the signed binary number to be 10111011. What is the decimal equivalent of this number if it is in Sign-Magnitude form, or 1's complement form, or 2's complement form?

Solution: Given binary number = 10111011. As sign bit is 1, it is a negative number. If it is in sign-magnitude format, then MSB is sign bit, and remaining bits represent the magnitude,

$(0111011)_2 = 32 + 16 + 8 + 2 + 1 = 59$. So if the given number is in sign-magnitude format, then the number is -59 .

If it is in 1's complement/2's complement form, then the magnitude of negative number can be obtained by taking 1's complement/2's complement for the number, respectively.

$$10111011 \Rightarrow 1\text{'s complement} \Rightarrow 01000100 = (68)_{10}$$

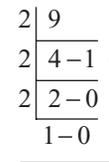
In 1's complement format, the number is -68 .

$$10111011 \Rightarrow 2\text{'s complement} \Rightarrow 01000101 = (69)_{10}$$

In 2's complement format, the number is -69 .

Example: Find $(-9.625)_{10}$ in signed 2's complement representation.

Signed binary fraction can be represented in the same way of signed integer.



$$0.625 \times 2 = 1.25$$

$$0.25 \times 2 = 0.5$$

$$0.5 \times 2 = 1.0$$

$$= 0.101$$

$$+(9.625) = 01001.101$$

$$-9.625 = 10110.011 \text{ (by taking 2's complement)}$$

Binary Multipliers

Multiplication of binary number is done in the same way as multiplication of decimal.

The multiplicand (m) is multiplied by each bit of the multiplier (N), starting from the LSB.

Let

$$M = B_3 B_2 B_1 B_0$$

$$N = A_3 A_2 A_1 A_0$$

$$\text{If } M \times N = P$$

			$A_0 B_3$	$A_0 B_2$	$A_0 B_1$	$A_0 B_0$	
		$A_1 B_3$	$A_1 B_2$	$A_1 B_1$	$A_1 B_0$		
	$A_2 B_3$	$A_2 B_2$	$A_2 B_1$	$A_2 B_0$			
$A_3 B_3$	$A_3 B_2$	$A_3 B_1$	$A_3 B_0$				
$P_7 P_6$	P_5	P_4	P_3	P_2	P_1	P_0	$= P$

Example: Let $M = 1011$

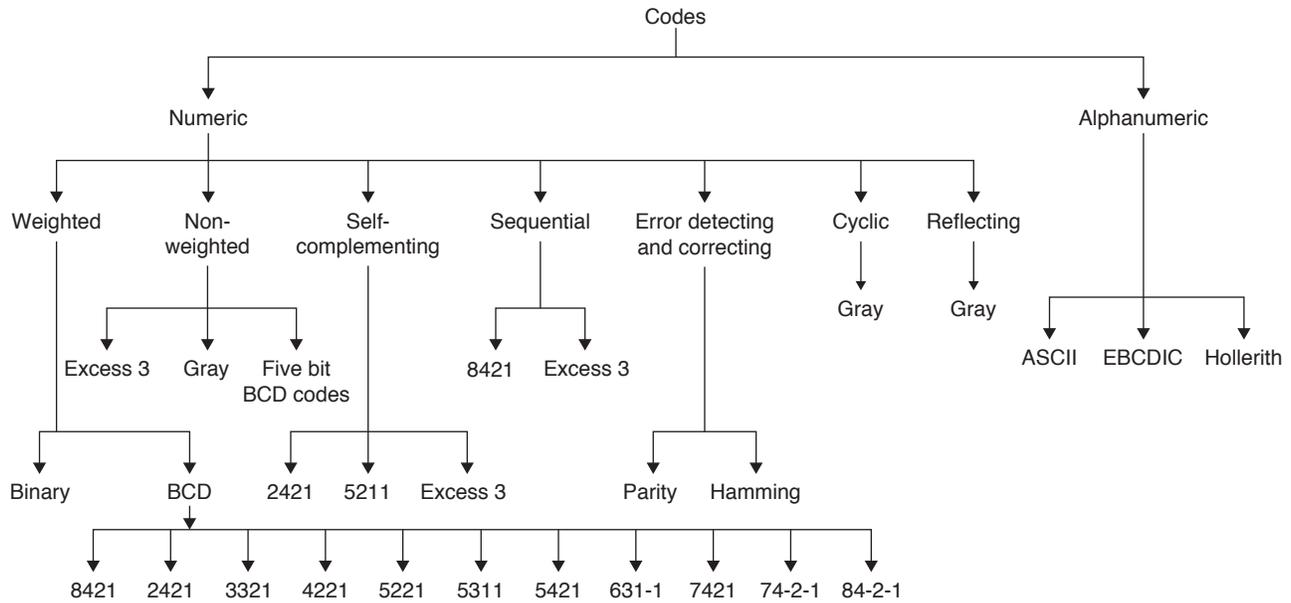
$$N = 1100$$

$$M \times N = P$$

$$\begin{array}{r} 1011 \\ \times 1100 \\ \hline 0000 \\ 0000 \\ 1011 \\ 1011 \\ \hline 1111 \\ \hline 10000100 = P \end{array}$$

Binary Codes

Binary codes can be classified as numeric codes and alpha-numeric codes. The figure below shows the classification of codes.



Numeric Codes

Numeric Codes are the codes which represent numerals in binary, i.e., only numbers as a series of 0s and 1s.

Weighted and non-weighted codes

- The weighted codes are those which obey the position-weighting principle. Each position of a number represents a specific weight.

Example: BCD, Binary, 84-2-1, 2421,

- Non-weighted codes are codes which are not assigned fixed values.

Example: Excess-3, Gray code

2421, 5211, 84-2-1 are examples of weighted codes, in which weight is assigned to each position in the number.

(27)₁₀ in 2421 code → 0010 1101

(45)₁₀ in 5211 code → 0111 1000

(36)₁₀ in 84-2-1 code → 0101 1010

Any digit in decimal will be represented by the weights represented by the code.

Error-detecting and correcting codes

Codes which allow only error detection are error-detecting codes.

Example: Parity

Codes which allow error detection as well as correction are called error correcting codes.

Example: Hamming codes

Sequential codes

A sequential code is one in which each succeeding code word is one binary number greater than the preceding code word.

Example: XS-3, BCD

Cyclic codes (unit distance codes)

Cyclic codes are those in which each successive code word differs from the preceding one in only one bit position.

Example: Gray code

Reflective codes

Binary code in which the n least significant bits for code words 2^n through $2^{n+1} - 1$ are the mirror images of than for 0 through $2^n - 1$ is called reflective codes.

Example: Gray Code

Self-complementing codes

A code is said to be self-complementing, if the code word of the 9's complement of number ' N ', i.e., of " $9-N$ " can be obtained from the code word of ' N ' by interchanging all the zeros and ones, i.e., by taking 1's complement. In other words, logical complement of number code is equivalent to representation of its arithmetic complement.

Example: 84-2-1, 2421, XS -3.

All weighted BCD codes are self-complementing codes.

Binary-coded decimal (BCD)

In BCD, each decimal digit 0 to 9 is coded by a 4-bit binary number. BCD codes are convenient to convert to/from decimal number system.

Decimal	BCD Digit
0	0000
1	0001
2	0010
3	0011
4	0100
5	0101
6	0110
7	0111
8	1000
9	1001

Example 42: $(628)_{10} = (0110\ 0010\ 1000)_{BCD}$

BCD addition

- BCD addition is performed by individually adding the corresponding digits of the decimal number expressed in 4-bit binary groups starting from the LSB.
- If there is no carry and the sum term is not an illegal code, no correction is needed.
- If there is a carry out of one group to the next group or if the sum term is an illegal code, the $(6)_{10}$ is added to the sum term of that group, and the resulting carry is added to the next group.

Example 43: $44 + 12$

$$\begin{array}{r} 0100 \quad 0100 \text{ (44 in BCD)} \\ \underline{0001 \quad 0010 \text{ (12 in BCD)}} \\ 0101 \quad 0110 \text{ (56 in BCD)} \end{array}$$

Example 44: $76.9 + 56.6$

$$\begin{array}{r} 0111 \quad 0110 \quad . \quad 1001 \\ \underline{0101 \quad 0110 \quad . \quad 0110} \\ 1100 \quad 1100 \quad . \quad 1111 \quad \text{(all are illegal codes)} \\ \underline{0110 \quad 0110 \quad . \quad 0110} \\ 0010 \quad 0010 \quad . \quad 0101 \\ \underline{\quad +1 \quad +1 \quad +1 \quad \text{(propagate carry)}} \\ 0001 \quad 0011 \quad 0011 \quad . \quad 0101 \\ 1 \quad 3 \quad 3 \quad . \quad 5 \end{array}$$

BCD subtraction BCD subtraction is performed by subtracting the digits of each 4-bit group of the subtrahend from the corresponding 4-bit group of the minuend in binary starting from the LSB.

Example 45: $42 - 12$

$$\begin{array}{r} 0100 \quad 0010 \text{ (42 in BCD)} \\ \underline{-0001 \quad 0010 \text{ (12 IN BCD)}} \\ 0011 \quad 0000 \text{ (No borrow, so this is the correct difference)} \end{array}$$

Example 46:

$$\begin{array}{r} 247.7 \quad 0010 \quad 0100 \quad 0111 \quad . \quad 0111 \quad \text{(Borrow are present, subtract 0110)} \\ \underline{-156.9 \quad 0001 \quad 0101 \quad 0110 \quad . \quad 1001} \\ 90.8 \quad 0000 \quad 0111 \quad 0000 \quad . \quad 1110 \\ \underline{\quad \quad \quad -01001 \quad -0110} \\ 1001 \quad 000 \quad . \quad 1000 \quad \text{Corrected difference (90.8)} \end{array}$$

Excess-3 (XS-3) code

Excess-3 code is a non-weighted BCD code, where each digit binary code word is the corresponding 8421 code word plus 0011.

Find the XS-3 code of

Example 47: $(3)_{10} \rightarrow (0011)_{BCD} = (0110)_{XS3}$

Example 48: $(16)_{10} \rightarrow (0001 \ 0110)_{BCD} \rightarrow (0100 \ 1001)_{XS3}$

Gray code

Each gray code number differs from the preceding number by a single bit.

Decimal	Gray Code
0	0000
1	0001
2	0011
3	0010
4	0110
5	0111

Binary to gray conversion

Step I: Shift the binary number one position to the right, LSB of the shifted number is discarded.

Step II: Exclusive or the bits of the binary number with those of the binary number shifted.

Example 49: Convert $(1001)_2$ to gray code

$$\begin{array}{l} \text{Binary} \quad \rightarrow 1010 \\ \text{Shifted Binary} \rightarrow \underline{101} \oplus \\ \text{Gray} \quad \rightarrow 1111 \end{array}$$

Gray to binary conversion

- Take the MSB of the binary number is same as MSB of gray code number.
- X-OR the MSB of the binary to the next significant bit of the gray code.
- X-OR the 2nd bit of binary to the 3rd bit of Gray code to get 3rd bit binary and so on.
- Continue this till all the gray bits are exhausted.

Example 50: Convert, gray code 1010 to binary

$$\begin{array}{r} \text{Gray} \quad \quad \quad 1 \quad 0 \quad 1 \quad 0 \\ 1010 \quad \quad \quad \downarrow \oplus \parallel \oplus \parallel \oplus \parallel \\ 1100 \quad \quad \quad 1 \quad 1 \quad 0 \quad 0 \\ = (1100)_2 \end{array}$$

Exercises

Practice Problems I

Directions for questions 1 to 15: Select the correct alternative from the given choices.

- Assuming all the numbers are in 2's complement representation, which of the following is divisible by 11110110?
 (A) 11101010 (B) 11100010
 (C) 11111010 (D) 11100111

- If $(84)_x$ (in base x number system) is equal to $(64)_y$ (in base y number system), then possible values of x and y are
 (A) 12, 9 (B) 6, 8
 (C) 9, 12 (D) 12, 18
- Let $A = 1111 \ 1011$ and $B = 0000 \ 1011$ be two 8-bit signed 2's complement numbers. Their product in 2's complement representation is

- (A) 11001001 (B) 10011100
(C) 11010101 (D) 10101101
4. Let r denotes number system's radix. The only value(s) of r that satisfy the equation $\sqrt[3]{(1331)_r} = (11)_r$ is/are
(A) 10 (B) 11
(C) 10 and 11 (D) any $r > 3$
5. X is 16-bit signed number. The 2's complement representation of X is $(F76A)_{16}$. The 2's complement representation of $8 \times X$ is
(A) $(1460)_{16}$ (B) $(D643)_{16}$
(C) $(4460)_{16}$ (D) $(BB50)_{16}$
6. The HEX number $(CD.EF)_{16}$ in octal number system is
(A) $(315.736)_8$ (B) $(513.637)_8$
(C) $(135.673)_8$ (D) $(531.367)_8$
7. 8-bit 2's complement representation a decimal number is 10000000. The number in decimal is
(A) +256 (B) 0
(C) -128 (D) -256
8. The range of signed decimal numbers that can be represented by 7-bit 1's complement representation is
(A) -64 to + 63 (B) -63 to + 63
(C) -127 to + 128 (D) -128 to +127
9. Decimal 54 in hexadecimal and BCD number system is respectively
(A) 63, 10000111 (B) 36, 01010100
(C) 66, 01010100 (D) 36, 00110110
10. A new binary-coded hexary (BCH) number system is proposed in which every digit of a base -6 number system is represented by its corresponding 3-bit binary code. For example, the base -6 number 35 will be represented by its BCH code 011101.
In this numbering system, the BCH code 001001101011 corresponds to the following number in base -6 system.
(A) 4651 (B) 4562
(C) 1153 (D) 1353
11. The signed 2's complement representation of $(-589)_{10}$ in Hexadecimal number system is
(A) $(F24D)_{16}$ (B) $(FDB3)_{16}$
(C) $(F42D)_{16}$ (D) $(F3BD)_{16}$
12. The base of the number system for which the following operation is to be correct $\frac{66}{5} = 13$
(A) 6 (B) 7
(C) 8 (D) 9
13. The solution to the quadratic equation $x^2 - 11x + 13 = 0$ (in number system with radix r) are $x = 2$ and $x = 4$. Then base of the number system is $(r) =$
(A) 7 (B) 6
(C) 5 (D) 4
14. The 16's complement of BADA is
(A) 4525 (B) 4526
(C) ADAB (D) 2141
15. $(11A1B)_8 = (12CD)_{16}$, in the above expression A and B represent positive digits in octal number system and C and D have their original meaning in Hexadecimal, the values of A and B are?
(A) 2, 5 (B) 2, 3
(C) 3, 2 (D) 3, 5

Practice Problems 2

Directions for questions 1 to 20: Select the correct alternative from the given choices.

1. The hexadecimal representation of $(567)_8$ is
(A) 1AF (B) D77
(C) 177 (D) 133
2. $(2326)_8$ is equivalent to
(A) $(14D6)_{16}$ (B) $(103112)_4$
(C) $(1283)_{10}$ (D) $(09AC)_{16}$
3. $(0.46)_8$ equivalent in decimal is?
(A) 0.59375 (B) 0.3534
(C) 0.57395 (D) 0.3435
4. The 15's complement of $(CAFA)_{16}$ is
(A) $(2051)_{16}$ (B) $(2050)_{16}$
(C) $(3506)_{16}$ (D) $(3505)_{16}$
5. 53 in 2's complement from is?
(A) 1001011 (B) 001010
(C) 0110101 (D) 001011
6. Signed 2's complement representation of $(-15)_{10}$ is
(A) 11111 (B) 10001
(C) 01111 (D) 10000
7. $(0.25)_{10}$ in binary number system is?
(A) (0.01) (B) (0.11)
(C) 0.001 (D) 0.101
8. The equivalent of $(25)_6$ in number system with base 7 is?
(A) 22 (B) 23
(C) 24 (D) 26
9. The operation $35 + 26 = 63$ is true in number system with radix
(A) 7 (B) 8
(C) 9 (D) 11
10. The hexadecimal equivalent of largest binary number with 14-bits is?
(A) 2FFF (B) 3FFFF
(C) FFFF (D) 1FFFF

1.12 | Unit 1 • Digital Logic

11. If x is radix of number system, $(211)_x = (152)_8$, then x is
 (A) 6 (B) 7
 (C) 9 (D) 5
12. The value of r for which $\sqrt{(224)_r} = (13)_r$ is valid expression, in number system with radix r is?
 (A) 5 (B) 6
 (C) 7 (D) 8
13. Which of the representation in binary arithmetic has a unique zero?
 (A) sign-magnitude (B) 1's compliment
 (C) 2's complement (D) All of these
14. For the binary number 101101111 the equivalent hexadecimal number is
 (A) 14E (B) 9E
 (C) B78 (D) 16F
15. Subtract 1001 from 1110
 (A) 0010 (B) 0101
 (C) 1011 (D) 1010
16. Which of the following is a positively weighted code?
 (A) 8421 (B) 84-2-1
 (C) EXS-3 (D) 74-2-1

17. Match the items correctly

Column 1	Column 2
(P) 8421	(1) Cyclic code
(Q) 2421	(2) self-complementing
(R) 5212	(3) sequential code
(S) Gray code	(4) non-sequential code

- (A) P-2, Q-4, R-3, S-1
 (B) P-1, Q-4, R-3, S-2
 (C) P-3, Q-2, R-4, S-1
 (D) P-2, Q-4, R-1, S-2
18. Perform the subtraction in XS-3 code $57.6 - 27.8$
 (A) 0101 1100.1011 (B) 0010 1001.1100
 (C) 00011101.1100 (D) 1010 1110.1011
19. The 2's complement representation of -17 is
 (A) 101110 (B) 111110
 (C) 101111 (D) 110001
20. The decimal 398 is represented in 2421 code by
 (A) 110000001000 (B) 001110011000
 (C) 001111111110 (D) 010110110010

PREVIOUS YEARS' QUESTIONS

1. $(1217)_8$ is equivalent to [2009]
 (A) $(1217)_{16}$ (B) $(028F)_{16}$
 (C) $(2297)_{10}$ (D) $(0B17)_{16}$
2. P is a 16-bit signed integer. The 2's complement representation of P is $(F87B)_{16}$. The 2's complement representation of $8*P$ is [2010]
 (A) $(C3D8)_{16}$ (B) $(187B)_{16}$
 (C) $(F878)_{16}$ (D) $(987B)_{16}$
3. The smallest integer that can be represented by an 8-bit number in 2's complement form is [2013]
 (A) -256 (B) -128
 (C) -127 (D) 0
4. The base (or radix) of the number system such that the following equation holds is $\frac{312}{20} = 13.1$ [2014]
5. Consider the equation $(123)_5 = (x8)_y$, with x and y as unknown. The number of possible solutions is [2014]
6. Consider the equation $(43)_x = (y3)_8$ where x and y are unknown. The number of possible solutions is [2015]
7. Suppose X_i for $i = 1, 2, 3$ are independent and identically distributed random variables whose probability mass functions are $Pr[X_i = 0] = Pr[X_i = 1] = \frac{1}{2}$ for $i =$

- 1, 2, 3. Define another random variable $Y = X_1 X_2 \oplus X_3$, where \oplus denotes XOR. Then
 $Pr[Y = 0 | X_3 = 0] =$ [2015]
8. The 16-bit 2's complement representation of an integer is 1111 1111 1111 0101; its decimal representation is [2016]
9. Consider an eight-bit ripple-carry adder for computing the sum of A and B , where A and B are integers represented in 2's complement form. If the decimal value of A is one, the decimal value of B that leads to the longest latency for the sum to stabilize is [2016]
10. Let $x_1 \oplus x_2 \oplus x_3 \oplus x_4 = 0$ where x_1, x_2, x_3, x_4 are Boolean variables, and \oplus is the XOR operator. Which one of the following must always be TRUE? [2016]
 (A) $x_1 x_2 x_3 x_4 = 0$
 (B) $x_1 x_3 + x_2 = 0$
 (C) $\bar{x}_1 \oplus \bar{x}_3 = \bar{x}_3 \oplus \bar{x}_4$
 (D) $x_1 + x_2 + x_3 + x_4 = 0$
11. Consider a quadratic equation $x^2 - 13x + 36 = 0$ with coefficients in a base b . The solutions of this equation in the same base b are $x = 5$ and $x = 6$. Then $b =$ [2017]

ANSWER KEYS**EXERCISES****Practice Problems 1**

1. B 2. C 3. A 4. D 5. D 6. A 7. C 8. B 9. B 10. C
11. B 12. D 13. C 14. B 15. D

Practice Problems 2

1. C 2. B 3. A 4. D 5. D 6. B 7. A 8. B 9. B 10. B
11. B 12. A 13. C 14. D 15. B 16. A 17. C 18. A 19. C 20. C

Previous Years' Questions

1. B 2. A 3. B 4. 5 5. 3 6. 5 7. 0.75 8. -11 9. -1.0 10. C
11. 8.0 to 8.0

Chapter 2

Boolean Algebra and Minimization of Functions

LEARNING OBJECTIVES

- Logic gates
- Boolean algebra
- AXIOMS and Laws of Boolean algebra
- Properties of Boolean algebra
- Conversion from Min term to Max term
- Minimization of Boolean function
- K-map method
- Prime implicant
- Implementation of function by using NAND-NOR Gates
- EX-OR, EX-NOR GATE

LOGIC GATES

- Inverter or NOT gate (7404 IC):** The inverter performs a basic logic operation called inversion or complementation. The purpose of the inverter is to change one logic level to the opposite level. In terms of digital circuits, it converts 1 to 0 and 0 to 1.

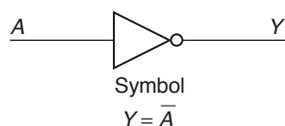


Table 1 Truth Table

Input		Output
A		Y
0		1
1		0

- AND gate (logical multiplier 7408 IC):** The AND gate performs logical multiplication more commonly known as AND function. The AND gate is composed of 2 or more inputs and a single output

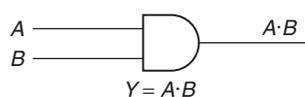


Figure 1 2 input AND gate

Table 2 Truth Table

Input		Output
A	B	Y
0	0	0
0	1	0
1	0	0
1	1	1

- OR gate (logical adder 7432 IC):** The OR gate performs logical addition commonly known as OR function.

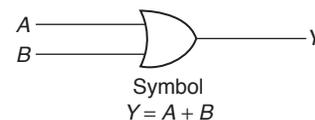


Figure 2 2 input OR gate

Table 3 Truth Table

Input		Output
A	B	Y
0	0	0
0	1	1
1	0	1
1	1	1

4. **NAND gate (7400 IC):** The NAND gate's function is basically AND + NOT function.

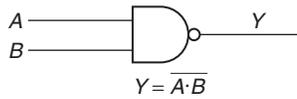


Figure 3 2 input NAND gate

Table 4 Truth Table

Input			Output
A	B	A · B	$\overline{A \cdot B}$ (Y)
0	0	0	1
0	1	0	1
1	0	0	1
1	1	1	0

5. **NOR gate (7402 IC):** The NOR gate is basically OR + NOT function.

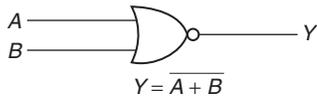


Figure 4 2 input NOR gate

Table 5 Truth Table

Input			Output
A	B	A + B	$\overline{A + B}$ (Y)
0	0	0	1
0	1	1	0
1	0	1	0
1	1	1	0

6. **Exclusive OR gate X-OR (7486 IC):** X-OR is a gate in which unequal inputs create a high logic level output and if both inputs are equal, the output will be low. Other name for EX-OR gate is unequivalent gate. 2 input X-OR Gate

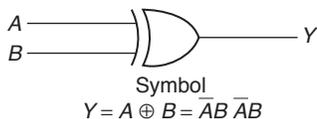


Figure 5 2 input X-OR Gate

Table 6 Truth Table

A	B	Y
0	0	0
0	1	1
1	0	1
1	1	0

7. **Exclusive NOR gate (X-NOR):** X-NOR is a gate in which equal inputs create a high logic level output; and

if both inputs are unequal, then the output will be low. Other name for X-NOR gate is equivalent gate.

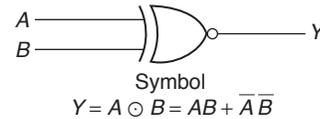


Figure 6 2 input X-NOR Gate

Table 7 Truth Table

A	B	Y
0	0	1
0	1	0
1	0	0
1	1	1

X-NOR Gate is complement of X-OR Gate.

BOOLEAN ALGEBRA

Boolean algebra is a system of mathematical logic. It is an algebraic system consisting of the set of elements (0, 1), two binary operators OR and AND and one unary operator NOT. The Boolean algebra is governed by certain well-developed rules and laws.

AXIOMS and Laws of Boolean Algebra

1. AXIOMS

(a) AND operation

- (1) $0 \cdot 0 = 0$
- (2) $0 \cdot 1 = 0$
- (3) $1 \cdot 0 = 0$
- (4) $1 \cdot 1 = 1$

(b) OR operation

- (5) $0 + 0 = 0$
- (6) $0 + 1 = 1$
- (7) $1 + 0 = 1$
- (8) $1 + 1 = 1$

(c) NOT operation

- (9) $\bar{1} = 0$
- (10) $\bar{0} = 1$

2. Laws

(a) Complementation law

- (1) $\bar{\bar{0}} = 1$
- (2) $\bar{\bar{1}} = 0$
- (3) If $A = 0$, then $\bar{A} = 1$
- (4) If $A = 1$, then $\bar{A} = 0$
- (5) $\overline{\bar{A}} = A$

(b) AND laws

- (1) $A \cdot 0 = 0$ (NULL Law)
- (4) $A \cdot 1 = A$ (Identity Law)
- (3) $A \cdot A = A$
- (4) $A \cdot \bar{A} = 0$

(c) OR laws

- (1) $A + 0 = A$ (NULL Law)
- (2) $A + 1 = 1$ (Identity Law)
- (3) $A + A = A$
- (4) $A + \bar{A} = 1$

(d) Commutative laws

- (1) $A + B = B + A$
- (2) $A \cdot B = B \cdot A$

(e) Associative laws

- (1) $(A + B) + C = A + (B + C)$
- (2) $(A \cdot B)C = A(B \cdot C)$

(f) Distributive laws

- (1) $A(B + C) = AB + AC$
- (2) $A + BC = (A + B)(A + C)$

(g) Redundant literal rule (RLR)

- (1) $A + \bar{A}B = A + B$
- (2) $A(\bar{A} + B) = AB$

(h) Idempotence laws

- (1) $A \cdot A = A$
- (2) $A + A = A$

(i) Absorption laws

- (1) $A + A \cdot B = A$
- (2) $A(A + B) = A$

3. Theorems

(a) Consensus theorem

Theorem 1:

$$AB + \bar{A}C + BC = AB + \bar{A}C$$

Proof:

$$\begin{aligned} \text{LHS} &= AB + \bar{A}C + BC \\ &= AB + \bar{A}C + BC(A + \bar{A}) \\ &= AB + \bar{A}C + BCA + BC\bar{A} \\ &= AB(1 + C) + \bar{A}C(1 + B) \\ &= AB(1) + \bar{A}C(1) \\ &= AB + \bar{A}C \\ &= \text{RHS.} \end{aligned}$$

Theorem 2:

$$(A + B)(\bar{A} + C)(B + C) = (A + B)(\bar{A} + C)$$

Proof:

$$\begin{aligned} \text{LHS} &= (A + B)(\bar{A} + C)(B + C) \\ &= (A\bar{A} + AC + B\bar{A} + BC)(B + C) \\ &= (AC + BC + \bar{A}B)(B + C) \\ &= ABC + BC + \bar{A}B + AC + BC + \bar{A}BC \\ &= AC + BC + \bar{A}B \\ \text{RHS} &= (A + B)(\bar{A} + C) \\ &= A\bar{A} + AC + BC + \bar{A}B \\ &= AC + BC + \bar{A}B \\ &= \text{LHS.} \end{aligned}$$

(b) Transposition theorem

$$AB + \bar{A}C = (A + C)(\bar{A} + B)$$

Proof:

$$\begin{aligned} \text{RHS} &= (A + C)(\bar{A} + B) \\ &= A\bar{A} + C\bar{A} + AB + CB \end{aligned}$$

$$\begin{aligned} &= 0 + \bar{A}C + AB + BC \\ &= \bar{A}C + AB + BC(A + \bar{A}) \\ &= AB + ABC + \bar{A}C + \bar{A}BC \\ &= AB + \bar{A}C \\ &= \text{LHS} \end{aligned}$$

(c) De Morgan's theorem

$$\text{Law 1: } \overline{A + B} = \bar{A} \cdot \bar{B}$$

This law states that the complement of a sum of variable is equal to the product of their individual complements.

$$\text{Law 2: } \overline{A \cdot B} = \bar{A} + \bar{B}$$

This law states that the complement of the product of variables is equal to the sum of their individual complements.

Example 1: Simplify the Boolean function $Y = A(A + \bar{B})$

$$Y = A \cdot A + A \cdot \bar{B}$$

$$\begin{aligned} \text{Solution: } Y &= A + A\bar{B} \\ &= A(1 + \bar{B}) \\ &= A \end{aligned}$$

Example 2: Simplify the Boolean function $Y = A + \bar{A}B$

$$\begin{aligned} \text{Solution: } Y &= A \cdot (B + 1) + \bar{A} \cdot B \\ &= A \cdot B + A + \bar{A}B \\ &= B(A + \bar{A}) + A \\ &= A + B \end{aligned}$$

Example 3: Simplify the Boolean function

$$Y = A(A + B) + B(\bar{A} + B)$$

$$\begin{aligned} \text{Solution: } Y &= A \cdot A + A \cdot B + B \cdot \bar{A} + B \cdot B \\ &= A + B(A + \bar{A}) + B \\ &= A + B \cdot 1 + B \\ &= A + B + B \\ &= A + B \end{aligned}$$

Example 4: Simplify the Boolean function

$$Y = \bar{A}\bar{B}\bar{C} + \bar{A}B\bar{C} + AB\bar{C} + A\bar{B}\bar{C}$$

$$\begin{aligned} \text{Solution: } Y &= \bar{A}\bar{C}(\bar{B} + B) + A\bar{C}(B + \bar{B}) \\ &= \bar{A}\bar{C} + A\bar{C} \\ &= \bar{C}(\bar{A} + A) \\ &= \bar{C} \end{aligned}$$

Example 5: Simplify the Boolean function

$$Y = \overline{ABC + \bar{A}B\bar{C} + \bar{A}\bar{B}C}$$

$$\begin{aligned} \text{Solution: } &= \overline{AC(B + \bar{B}) + \bar{A}B\bar{C}} \\ &= \overline{\bar{A}C + \bar{A}B\bar{C}} \\ &= \overline{\bar{A}(C + B\bar{C})} \\ &= \overline{\bar{A}(C + B)} \\ &= A + \bar{C}\bar{B} \end{aligned}$$

Example 6: Simplify the Boolean function

$$Y = AB + C\bar{B} + CA + ABD$$

Solution: $Y = AB(1 + D) + C\bar{B} + CA$

$$= AB + C\bar{B} + CA$$

$$= AB + C\bar{B}$$

PROPERTIES OF BOOLEAN ALGEBRA

With n variables, maximum possible distinct functions $= 2^{2^n}$.

Duality consider the distributive law

1. $x(y + z) = xy + xz$
2. $x + yz = (x + y)(x + z)$

Second one can be obtained from the first law if the binary operators and the identity elements are interchanged. This important property of Boolean algebra is called the duality principle.

The dual of an algebraic expression can be written by interchanging OR and AND operators, 1s by 0, and 0s by 1s.

Example 7: $x + x' = 1 \xrightarrow{\text{Dual}} x \cdot x' = 0$

Solution: $xy + xy' = x \xrightarrow{\text{Dual}} (x + y)(x + y') = x$

$$x + x'y = x + y \xrightarrow{\text{Dual}} x(x' + y) = xy$$

Example 8: The dual of $F = xy + xz + yz$ is?

Solution: Dual of $F = (x + y)(x + z)(y + z)$
 $= (x + xz + xy + yz)(y + z) = xy + yz + xz$

So dual of $xy + xz + yz$ is same as the function itself; For N variables maximum possible self-dual functions $= 2^{2^{n-1}} = 2^{(2^n/2)}$

Example 9: Which of the following statement/s is/are true

- S₁: The dual of NAND function is NOR
- S₂: The dual of X-OR function is X-NOR
- (A) S₁ and S₂ are true
- (B) S₁ is true
- (C) S₂ is true
- (D) None of these

Solution: (A)

$$\text{NAND} = (xy)' = x' + y'$$

$$\text{Dual of NAND} = (x + y)' = x'y'$$

$$\text{X-OR} = xy' + x'y$$

$$\text{Dual of X-OR} = (x + y')(x' + y) = xy + x'y' = \text{X-NOR}$$

Both S₁ and S₂ are true

Operator precedence The operator precedence for evaluating Boolean expression is

1. Parentheses
2. NOT
3. AND
4. OR

So the expression inside the parentheses must be evaluated before all the operations. The next operation to be performed is the complement and then follows AND and finally the OR.

Complement of function The complement of a function F is F' is obtained from an interchange of 0s for 1s and 1s for 0s in the value of F . The complements of a function may be derived algebraically through De Morgan's theorems.

$$(x_1 \cdot x_2 \cdot x_3 \dots x_n)' = x_1' + x_2' + x_3' + \dots + x_n'$$

$$(x_1 + x_2 + x_3 + \dots + x_n)' = x_1' \cdot x_2' \cdot x_3' \cdot x_4' \dots x_n'$$

Example 10: The complement of function $F = a(b'c + bc')$ is?

Solution: $(F)' = [a(b'c + bc')]'$
 $= a' + (b'c + bc)'$
 $= a' + (b'c)' \cdot (bc)'$
 $= a' + (b + c)'(b' + c)$
 $F' = a' + bc + b'c'$

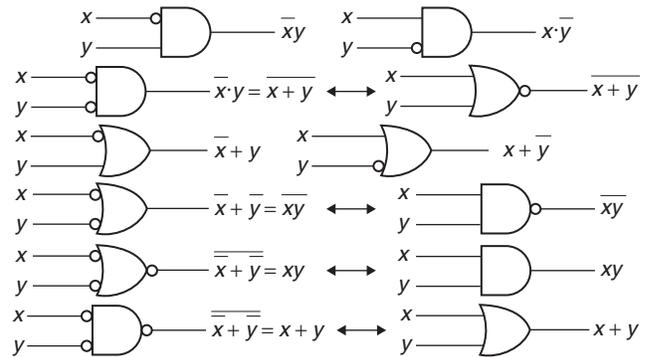


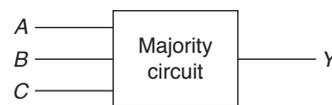
Figure 7 Gates with inverted inputs

BOOLEAN FUNCTIONS, MIN TERMS AND MAX TERMS

The starting point for designing most logic circuits is the truth table, which can be derived from the statement of problem. The truth table is then converted into a Boolean expression and finally create the assembly of logic gates accordingly.

Let us consider the example of majority circuit. This circuit takes three inputs (A, B, C) and have one output (Y) which will give the majority of the inputs, i.e., if A, B, C are having more number of zeros, $Y = 0$ else if A, B, C are having more number of 1s, $Y = 1$.

So from the statements we can derive the truth table as follows:



As we are using three Boolean variables A, B, C , total number of combinations in truth table are $2^3 = 8$.

Similarly for n variables, the truth table will have total of 2^n combinations, for a Boolean function.

Sl. no.	Input			Output
	A	B	C	Y
1	0	0	0	0 \rightarrow $Y=0$, If inputs are having more zeros.
2	0	0	1	
3	0	1	0	
4	0	1	1	1 \rightarrow $Y=1$, If inputs are having more 1's
5	1	0	0	
6	1	0	1	
7	1	1	0	
8	1	1	1	

For some combinations, output $Y = 1$, and for others $Y = 0$. The input combinations for which output $Y = 1$ are called as min terms.

Similarly the input combinations for which output $Y = 0$ are called as max terms.

Min terms are expressed as product terms, Similarly, max terms are expressed as sum terms.

The output $Y = 1$, only in rows 4, 6, 7, 8.

So the min terms combinations are 011, 101, 110, 111 in Boolean Algebra, 1 input will be written as A, B, C and 0 input will be written as $\bar{A}, \bar{B}, \bar{C}$ in complement form, we express these min terms as product terms, $\bar{A}BC, A\bar{B}C, ABC, \bar{A}BC$.

To express Y as Boolean expression, we can write it as sum of the min terms.

$$Y = \bar{A}BC + A\bar{B}C + ABC + \bar{A}BC$$

We know that AND operation is a product while OR is sum. So the above equation is a sum of the products (SOP), (or) min terms expression.

The other way of expressing Y is $Y = \sum m(3, 5, 6, 7)$.

$$Y = m_3 + m_5 + m_6 + m_7.$$

The min term numbers are the decimal equivalent of input binary combinations.

Similar to SOP we can have product of sums (POS) Boolean expression.

The output $Y = 0$ for the input combinations 000, 001, 010, 100. For max terms 1 input will be indicated as $\bar{A}, \bar{B}, \bar{C}$ in complement form, 0 input will be indicated as A, B, C and max terms are expressed as sum terms.

$$A + B + C, A + B + \bar{C}, A + \bar{B} + C, \bar{A} + B + C$$

Any function can be expressed as product of max terms.

$$\text{So } Y = (A + B + C)(A + B + \bar{C})(A + \bar{B} + C)(\bar{A} + B + C)$$

The above equation is a product of sum expression (POS) or max terms expression.

$$\begin{aligned} \text{In other way } Y &= \pi M(0, 1, 2, 4) \\ &= M_0 \cdot M_1 \cdot M_2 \cdot M_4 \end{aligned}$$

The max term numbers are decimal equivalents of corresponding input binary combinations.

Min Term and Max Term

All the Boolean expressions can be expressed in a standard sum of product (SOP) form or in a standard product of sum (POS) form.

- A standard SOP form is one in which a number of product terms, each contains all the variables of the function either in complement or non-complement form are summed together.
- A standard POS form is one in which a number of sum terms, each one of which contain all the variable of the function either in complemented or non-complement form are multiplied together.
- Each of the product term in standard SOP form is called a min term.
- Each of the sum term in the standard POS form is called a max term.

Conversion from min terms to max terms representation

$$Y = \bar{A}BC + A\bar{B}C + ABC + \bar{A}BC$$

$$\begin{aligned} Y' &= (\bar{A}BC + A\bar{B}C + ABC + \bar{A}BC)' \\ &= (\bar{A}BC)'(A\bar{B}C)'(ABC)'(\bar{A}BC)' \end{aligned}$$

$$\begin{aligned} (Y')' &= [(A + \bar{B} + \bar{C})(\bar{A} + B + \bar{C})(\bar{A} + \bar{B} + C)(\bar{A} + \bar{B} + \bar{C})]' \\ &= [\pi(3, 5, 6, 7)]' = \pi(0, 1, 2, 4) \end{aligned}$$

$$Y = (A + B + C)(A + B + \bar{C})(A + \bar{B} + C)(\bar{A} + B + C)$$

$$\text{or } Y = \Sigma(3, 5, 6, 7) = \pi(0, 1, 2, 4)$$

Conversion from normal SOP/POS form to canonical SOP/POS

Let us consider $f(A, B, C) = A + BC + \bar{A}C$

The above function is in normal (minimized) SOP form, to convert this function to standard SOP(or) canonical SOP form, include missing variable in each and every term, to make it complete. First term A , Missing literals are B, C . Consider $A X X$, so possible combinations are $\bar{A}BC, A\bar{B}C, ABC, \bar{A}BC$ or we can write

$$A = A = A(B + \bar{B})(C + \bar{C}) = ABC + A\bar{B}C + \bar{A}BC + \bar{A}BC$$

Second term BC -missing literal is A . Consider $XBC \Rightarrow$ So possible combinations are $ABC, \bar{A}BC$ or we can write

$$\begin{aligned} BC &= (A + \bar{A})BC \\ &= ABC + \bar{A}BC \end{aligned}$$

Third term $\bar{A}C$ = missing literal is B . Consider $\bar{A}XC \rightarrow$ so possible combinations are $\bar{A}BC, \bar{A}\bar{B}C$ or we can write

$$\begin{aligned} \bar{A}C &= \bar{A}(B + \bar{B})C \\ &= \bar{A}BC + \bar{A}\bar{B}C \end{aligned}$$

Now, $f(A, B, C) = ABC + ABC\bar{C} + A\bar{B}C + \bar{A}BC + \bar{A}\bar{B}C$, after removing the redundant terms.

Now consider

$$f(A, B, C) = (A + B)(\bar{A} + C)$$

To convert this expression to canonical form or standard POS form we can write

$$f(A, B, C) = (A + B + C \cdot \bar{C})(\bar{A} + B + \bar{B} + C)$$

Here the C variables is absent from first term and B from second term. So add $C \cdot \bar{C} = (0)$ to first, and $B \cdot \bar{B}$ to second, and using distributive law arrive at the result.

$$f(A, B, C) = (A + B + C)(A + B + \bar{C})(\bar{A} + B + C)(\bar{A} + \bar{B} + C)$$

MINIMIZATION OF BOOLEAN FUNCTIONS

Simplification Procedure

- Obtain truth table, and write output in canonical form or standard form
- Generate K-map!
- Determine Prime implicants.
- Find minimal set of prime implicants.

Karnaugh Map (K-map) Method

Karnaugh map method is a systematic method of simplifying the Boolean expression. K-map is a chart or a graph composed of an arrangement of adjacent cell, each representing a particular combination of variable in sum or product form. (i.e., min term or max term).

Two-variable K-map

x	y	F
0	0	m_0
0	1	m_1
1	0	m_2
1	1	m_3

m_0	m_1
m_2	m_3

x \ y	0	1
0	$x'y'$	$x'y$
1	xy'	xy

Three-variable K-map

A three-variable map will have eight min terms (for three variables $2^3 = 8$) represented by 8 squares

x	y	z	F
0	0	0	m_0
0	0	1	m_1
0	1	0	m_2
0	1	1	m_3
1	0	0	m_4
1	0	1	m_5
1	1	0	m_6
1	1	1	m_7

m_0	m_1	m_3	m_2
m_4	m_5	m_7	m_6

x \ yz	00	01	11	10
0	$x'y'z$	$x'y'z'$	$x'yz$	$x'yz'$
1	$xy'z'$	$xy'z$	xyz	xyz'

3-variable K-map

Four-variable K-maps

The K-map for four variables is shown here, 16 min terms are assigned to 16 squares.

wx \ yz	00	01	11	10
00	0	1	3	2
01	4	5	7	6
11	12	13	15	14
10	8	9	11	10

The map is considered to lie on a surface with the top and bottom edges as well as the right and left edge touching each other to form adjacent squares.

- One square → a min term of four literals
- Two adjacent square → a term of three literals
- Four adjacent square → a term of two literals
- Eight adjacent square → a term of one literal
- Sixteen adjacent square → The constant one

Don't-care Combinations

It can often occur that for certain input combinations, the value of the output is unspecified either because the input combination are invalid or because the precise value of the output is of no consequence. The combination for which the values of the expression are not specified are called don't-care combinations. During the process of design using an SOP, K-map, each don't-care is treated as a 1 if it is helpful in Map Reduction, otherwise it is treated as a 0 and left alone. During the process of design using a POS K-map, each Don't-care is treated as a 0 if it is useful in Map Reduction, otherwise it is treated as a 1 and left alone.

Example 11: Find the Minimal expression

$$\Sigma m(1, 5, 6, 12, 13, 14) + d(2, 4)$$

Solution:

AB \ CD	00	01	11	10
00		1		×
01	×	1		1
11	1	1		1
10				

$$\therefore F = B\bar{C} + B\bar{D} + \bar{A}\bar{C}\bar{D}$$

Pairs, Quads and Octets

		BC			
		00	01	11	10
A	0		1	1	
	1				

The map contains a pair of 1s those are horizontally adjacent. Those cells represent $\overline{A}\overline{B}C, \overline{A}BC$.

For these two min terms, there is change in the form of variable B . By combining these two cells we can form a pair, which is equal to $\overline{A}\overline{B}C + \overline{A}BC = \overline{A}C(\overline{B} + B) = \overline{A}C$.

If more than one pair exists on K-map, OR the simplified products to get the Boolean function.

		BC			
		00	01	11	10
A	0	1			1
	1		1	1	

$$F = \overline{A}C + AC$$

		CD			
		00	01	11	10
AB	00	1	1	1	
	01	1			
11				1	
10	1	1		1	

$$F = \overline{A}C\overline{D} + \overline{A}BD + \overline{A}BC + AC\overline{D}$$

So Pair eliminates one variable by minimization.

Quad

Quad is a group of four 1s those are horizontally or vertically adjacent.

		BC			
		00	01	11	10
A	0		1	1	
	1		1	1	

$$F = C$$

		BC			
		00	01	11	10
A	0		1	1	
	1		1	1	

$$F = \overline{A}C + AC = (\overline{A} + A)C = C$$

By considering two pairs also it will be simplified to C . Quad eliminates two variables from the function

		CD			
		00	01	11	10
AB	00		1	1	
	01	1			1
11	1			1	
10		1	1		

$$F = \overline{B}\overline{D} + \overline{B}D$$

Corner min terms can form a Quad

		RS			
		00	01	11	10
PQ	00	1			1
	01				
11					
10	1			1	

$$F = \overline{Q}\overline{S}$$

Octet

The group of eight cells will form one octet.

		ZW			
		00	01	11	10
XY	00				
	01	1	1	1	1
11	1	1	1	1	
10					

$$F = Y$$

Other variable X, Z, W changes their form in octet. Octet can eliminate three variables and their complements.

		CD			
		00	01	11	10
AB	00	1			1
	01	1			1
11	1			1	
10	1			1	

$$F = \overline{D}$$

Other variable A, B, C are vanished.

Eliminating Redundant Groups

		BC			
		00	01	11	10
A	0		1	1	
	1		1	1	

$$AB + \overline{A}C + BC$$

		BC			
		00	01	11	10
A	0		1	1	
	1		1	1	

$$AB + \overline{A}C$$

Here BC is redundant pair, which covers already covered min terms of $AB, \overline{A}C$.

		RS			
		00	01	11	10
PQ	00		1		
	01	1		1	1
11	1	1	1		
10			1		

This K-map gives four pairs and one quad.

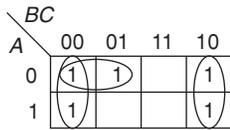
		RS			
		00	01	11	10
PQ	00		1		
	01		1	1	1
11	1	1	1		
10			1		

But only four pairs are enough to cover all the min times, Quad is not necessary.

$$\overline{P}\overline{R}S + \overline{P}QR + PQ\overline{R} + PRS \text{ is minimized function.}$$

Prime Implicant

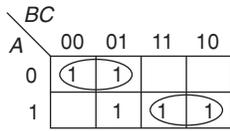
The group of adjacent min terms is called a Prime Implicant, i.e., all possible pairs, quads, octets, etc.



Prime implicants are $\overline{B}\overline{C}$, $\overline{B}C$, \overline{C} , $\overline{A}B$. Minimized function is $\overline{C} + \overline{A}B$

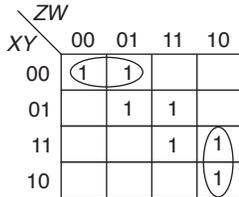
Essential Prime Implicant

The prime implicant which contains at least one min term which cannot be covered by any other prime implicant is called Essential prime implicant.



Min term 0, 6 can be grouped with only one pair each.

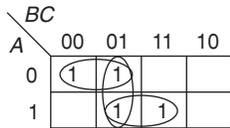
The total possible prime implicants are $\overline{A}B$, $\overline{B}C$, AC , AB but min term 0, 6 can be covered with $\overline{A}B$, AB . So we call them as essential prime implicants. Min term 5 can be paired with any of 1 or 7 min term.



The essential prime Implicants are $XZ\overline{W}$, $\overline{X}\overline{Y}\overline{Z}$

Redundant Prime Implicant

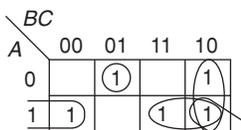
The prime implicant whose min terms are already covered by at least one min term is called redundant prime implicants.



Here prime implicants are $\overline{A}\overline{B}$, AC , $\overline{B}C$. But $\overline{B}C$ is already covered by other min terms So $\overline{B}C$ is redundant prime implicant.

Example 12: Find the minimal expression for $\Sigma m(1, 2, 4, 6, 7)$ and implement it using Basic gates.

Solution: K-map is



$$F = \overline{A}\overline{C} + \overline{A}B + \overline{B}\overline{C} + \overline{B}C + \overline{A}\overline{B}C$$

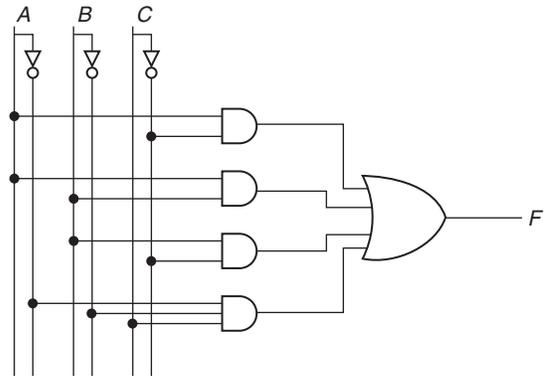
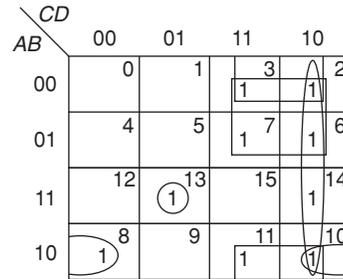


Figure 8 Logic diagram

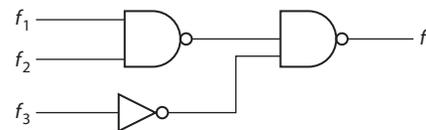
Example 13: Find the minimal expression for $\Sigma m(2, 3, 6, 7, 8, 10, 11, 13, 14)$

Solution: K-map is:



$$\therefore F(A, B, C, D) = ABC\overline{D} + A\overline{B}\overline{D} + \overline{A}C + \overline{B}C + C\overline{D}$$

Example 14: Three Boolean functions are defined as below $f_1 = \Sigma m(0, 1, 3, 5, 6)$, $f_2 = \Sigma m(4, 6, 7)$, $f_3 = \Sigma m(1, 4, 5, 7)$, then find f .



Solution: When two Boolean functions are ANDed, the resultant will contain the common min terms of both of the functions (like, intersection of min terms). If two Boolean functions are ORed, then resultant is the combination of all the min terms of the functions (like union of min terms)

$$\text{Here } f = \overline{\overline{f_1 f_2} \cdot \overline{f_3}} = f_1 f_2 + f_3$$

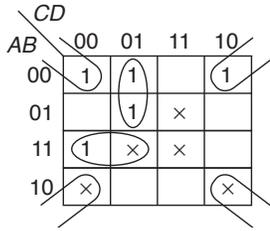
Here $f_1 \cdot f_2 =$ Common min terms in f_1 and $f_2 = \Sigma m(6)$

$$f_1 \cdot f_2 + f_3 = \text{Combination of min terms of } f_1 \cdot f_2 \text{ and } f_3 = \Sigma(1, 4, 5, 6, 7)$$

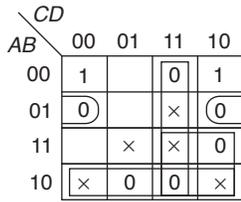
Example 15: What is the literal count for the minimized SOP, and minimized POS form for the following function?

$$f(A, B, C, D) = \Sigma m(0, 1, 2, 5, 12) + \phi d(7, 8, 10, 13, 15)$$

Solution: $f(A, B, C, D) = \Sigma m(0, 1, 2, 5, 12) + \phi(7, 8, 10, 13, 15)$



$f = 1$ quad + 2 pairs
 Literal count = $1 \times 2 + 2 \times 3 = 8$
 $f(A, B, C, D) = \pi M(3, 4, 6, 9, 11, 14) + \phi(7, 8, 10, 13, 15)$



f will consists of 3 quads + 1 pair
 $= 3 \times 2 + 1 \times 3 = 6 + 3 = 9$

IMPLEMENTATION OF FUNCTION BY USING NAND-NOR GATES

NAND or NOR gates are called as universal gates, because any function can be implemented by using only NAND gates or only using NOR gates.

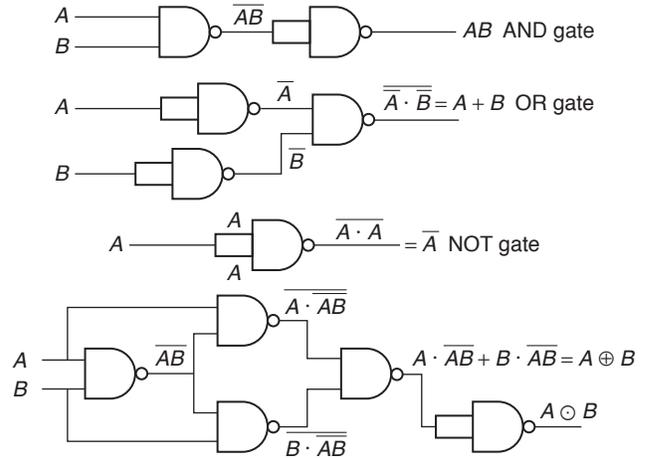


Figure 9 Implement of basics gates by using NAND gates

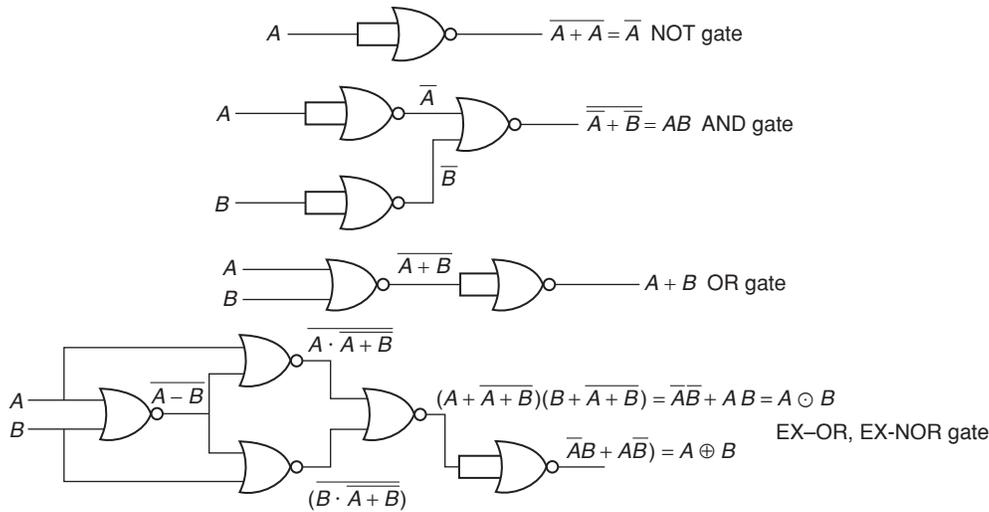
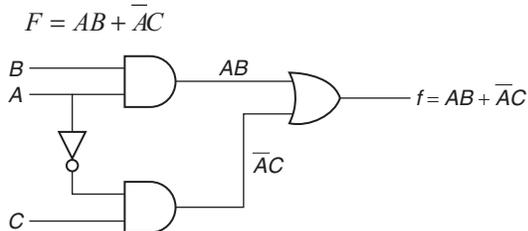
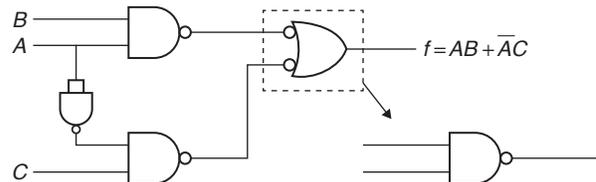


Figure 10 Implementation of basic gates by using NOR gates

Any function which is in the SOP form can be implemented by using AND-OR gates, which is also equivalent to NAND-NAND gates.



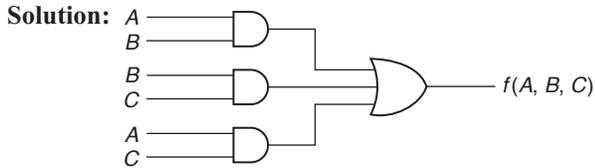
By considering bubble at AND gate output and OR gate input, and by changing NOT gates to NAND gates the circuit becomes as,



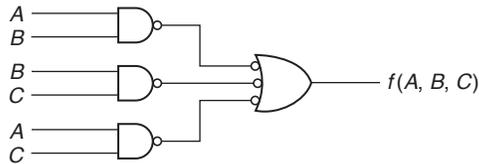
Now the circuit is in completely in NAND-NAND form. So the functions expressed in SOP form, can be implemented by using AND-OR, (or) NAND-NAND gates.

Any function in POS form, can be implemented by using OR-AND gates, which is similar to NOR-NOR gate.

Example 16: How many number of NAND gates are required to implement $f(A, B, C) = AB + BC + AC$
 (A) 3 (B) 4 (C) 5 (D) 6



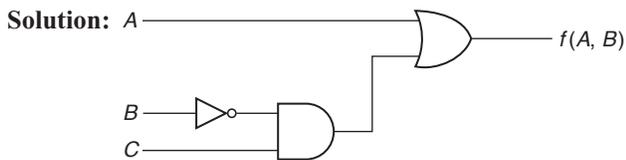
By considering bubbles at output of AND gate and input of OR gate.



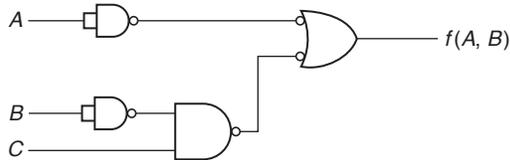
So four NAND gates are required.

Example 17: Number of NAND gates required for implementation of $f(A, B) = A + \overline{B}C$ is

- (A) 3 (B) 4 (C) 5 (D) 6



To convert the all gates into NAND gates, place bubble output of AND, and inputs of OR gates. Now, the circuit can be drawn as



Four NAND gates are required.

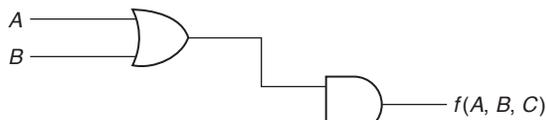
Example 18: $f = A + BC$, the number of NOR gates required to implement f , are?

- (A) 3 (B) 4 (C) 5 (D) 2

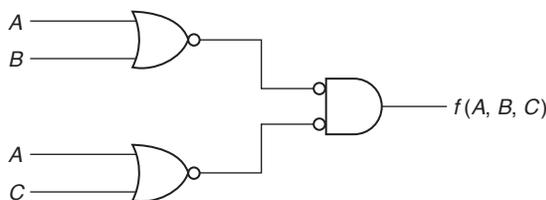
Solution: $A + BC$ is in SOP form.

To implement this function by using NOR gates, we can write $f(A, B, C) = A + BC = (A + B)(A + C)$

Which is in the form of POS?



By including bubbles at output of OR gate, and input of AND gate, the circuit becomes



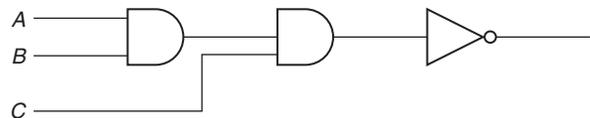
Now the circuit consists of all NOR gates. Three NOR Gates are required.

Example 19: How many number of two-input NAND-NOR gates are required to implement three-input NAND-NOR gates respectively?

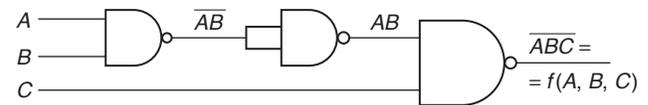
- (A) 2, 2 (B) 2, 3
(C) 3, 2 (D) 3, 3

Solution: $f(A, B, C) = \overline{ABC} = \overline{AB} + \overline{C}$

(1) Implement above function by using two-inputs gates

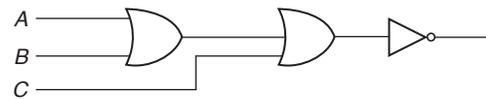


Now convert each gate to NAND gate

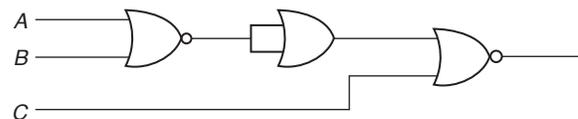


Three two-input NAND gates are required.

(2) $G(A, B, C) = \overline{A + B + C}$ Implement it by using two-input gates



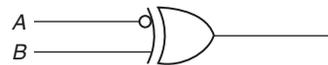
Now convert each gate to NOR gate



Three two-input, NOR gates are required.

EX-OR, EX-NOR GATES

Inverted inputs for EX OR, EX-NOR gates



$$\overline{A} \oplus B = \overline{\overline{A}B} + \overline{A\overline{B}} = AB + \overline{A}\overline{B} = A \odot B$$



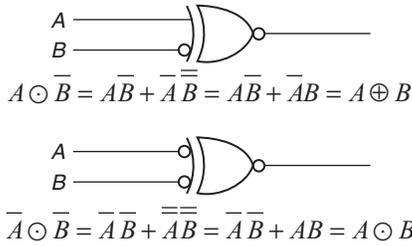
$$A \oplus \overline{B} = \overline{\overline{A}\overline{B}} + \overline{A\overline{B}} = AB + \overline{A}\overline{B} = A \odot B$$



$$\overline{A} \oplus B = \overline{\overline{A}\overline{B}} + \overline{A\overline{B}} = \overline{A}\overline{B} + A\overline{B} = A \oplus B$$



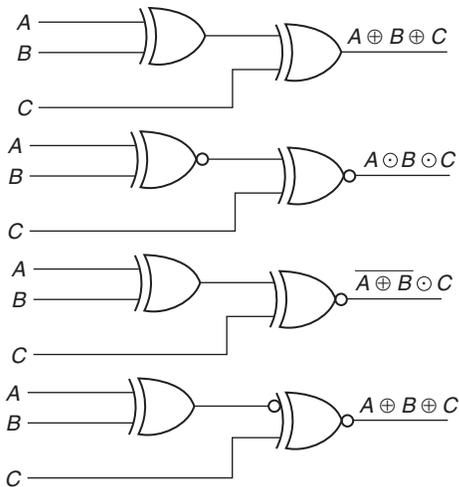
$$\overline{A} \odot B = \overline{\overline{A}\overline{B}} + \overline{A\overline{B}} = \overline{A}\overline{B} + A\overline{B} = A \oplus B$$



From the above discussions we can conclude that inverted input EXOR gate is EX-NOR gate.

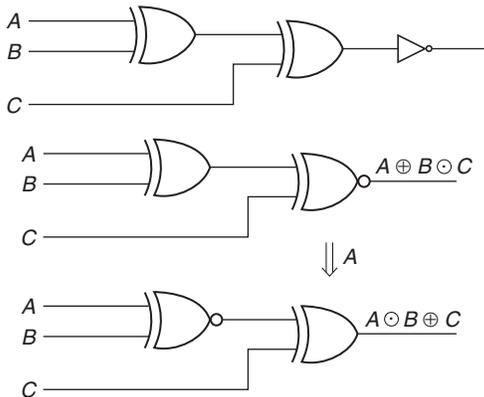
Similarly, inverted input EX-NOR gate is EX-OR gate. If both inputs are inverted the EX-OR / EX-NOR will remain as it is.

Consider a three-inputs X-OR gates by using two-input XOR gates.



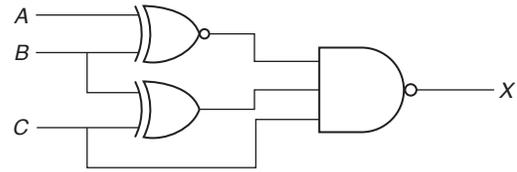
So we can conclude that $A \oplus B \oplus C = A \odot B \odot C$

$$\overline{A \oplus B \oplus C} = \overline{A \odot B \odot C}$$



$$\begin{aligned} \overline{A \oplus B \oplus C} &= \overline{A \odot B \odot C} = A \oplus B \odot C \\ &= A \odot B \oplus C \\ A \oplus B \oplus C \oplus D &= A \oplus B \odot C \odot D = A \odot B \odot C \oplus D = A \\ &\odot B \oplus C \odot D \\ A \odot B \odot C \odot D &= \overline{A \oplus B \oplus C \oplus D} \\ A \odot B \odot C \odot D &= A \oplus B \oplus C \odot D = A \oplus B \odot C \oplus D \end{aligned}$$

Example 20: For the logic circuit shown in figure, the required input condition (A, B, C) to make the output $X = 0$ is?



- (A) 1, 1, 1
- (B) 1, 0, 1
- (C) 0, 1, 1
- (D) 0, 0, 1

Solution: (D)

To get output $X = 0$, all inputs to the NAND gate should be 1, so $C = 1$.

When $C = 1$, the output of X-OR gate $B \oplus C = 1$ only when $B = 0$.

If $B = 0$ the output of X-NOR gate $A \odot B = 1$.

Only when $A = 0$

So $X = 1$, only when $(A, B, C) = (0, 0, 1)$.

Example 21: The minimized expression of

$$(A + \bar{B})(\bar{A}\bar{B} + AC)(\bar{A}\bar{C} + \bar{B})$$

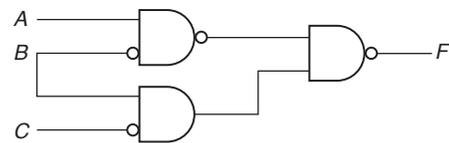
Solution: $(A + \bar{B})(\bar{A}\bar{B} + AC)(\bar{A}\bar{C} + \bar{B})$

$$\begin{aligned} &= (A + \bar{B})(\bar{A}\bar{B} \cdot \bar{A}\bar{C} + \bar{A}\bar{B} \cdot \bar{B} + AC \cdot \bar{A}\bar{C} + AC \cdot \bar{B}) \\ &= (A + \bar{B})(\bar{A}\bar{B} + \bar{A}\bar{B}C) = (A + \bar{B})\bar{A}\bar{B}(1 + C) \\ &= \bar{A}\bar{B} + \bar{A}\bar{B} = \bar{A}\bar{B} \end{aligned}$$

Example 22: The Boolean function f is independent of

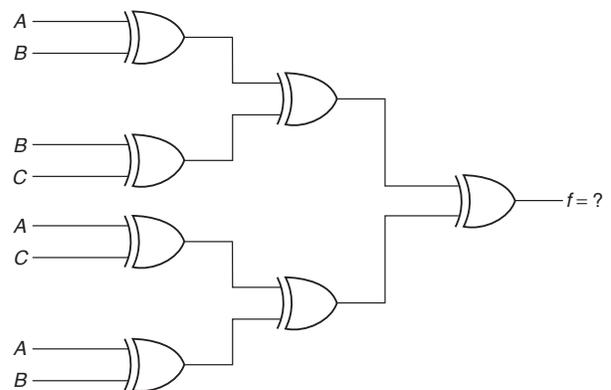
- (A) a
- (B) b
- (C) c
- (D) None of these

Solution: (A)



$$\begin{aligned} F &= \overline{\overline{ab} \cdot \overline{bc}} \\ &= ab + \overline{bc} = ab + b + \bar{c} \\ &= b + \bar{c} \text{ is independent of 'a'.} \end{aligned}$$

Example 23:



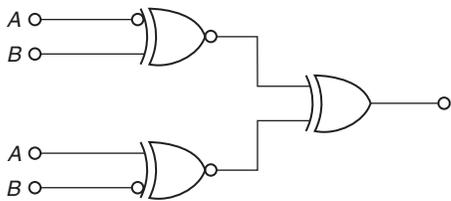
Solution: $f = \{A \oplus B \oplus B \oplus C\} \oplus \{A \oplus C \oplus B \oplus A\}$
 $= \{A \oplus 0 \oplus C\} \oplus \{0 \oplus C \oplus B\}$
 $= A \oplus C \oplus C \oplus B = A \oplus 0 \oplus B = A \oplus B$

Solved Examples

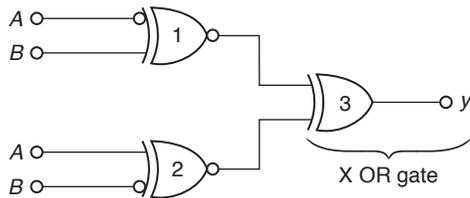
Example 1: Simplify the Boolean function, $xy + x'z + yz$

Solution: $xy + x'z + yz$
 By using consensus property
 $xy + x'z + yz = xy + x'z$
 $Y = xy + x'z$

Example 2: The output of the given circuit is equal to



Solution: $\bar{A} \odot B = \bar{A}B + A\bar{B}$



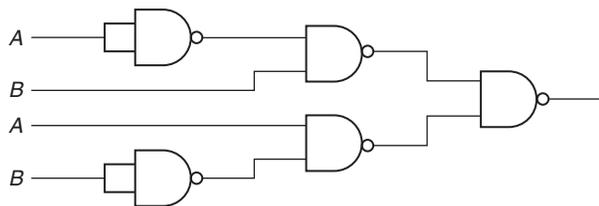
$A \odot \bar{B} = \bar{A}B + A\bar{B}$

So the output of above circuit is '0'. As two inputs are same at third gate.

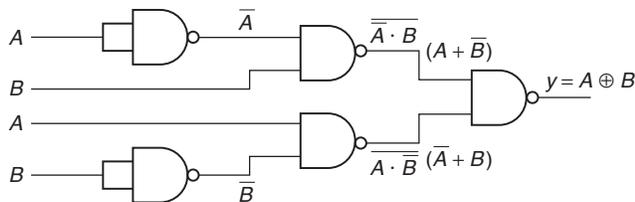
Output of XOR gate with two equal inputs is zero.

$\therefore y = 0$

Example 3: The circuit shown in the figure is functionally equivalent to



Solution:



$Y = \overline{(\bar{A} \cdot B \cdot \bar{A}B)} = \overline{(A + \bar{B})(\bar{A} + B)} \quad \because (\bar{A} \cdot B = \bar{A} + \bar{B})$
 $= \overline{(A + \bar{B}) + (\bar{A} + B)} = \bar{A} \cdot \bar{B} + \bar{A} \cdot \bar{B}$
 $= \bar{A} \cdot B + A \cdot \bar{B} = A \oplus B$

Example 4: Simplify the Boolean function $A \oplus \bar{A}B \oplus \bar{A}$

Solution: $A \oplus \bar{A}B \oplus \bar{A}$
 Associativity
 $= 1 \oplus \bar{A}B = \bar{A}B$
 $= A + \bar{B} \quad (\because \text{De Morgan's})$

Example 5:

		AB			
	CD	00	01	11	10
00		0	0	1	1
01		0	x	x	1
11		x	x	1	x
10		1	0	1	1

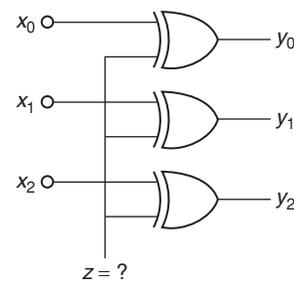
The minimized expression for the given K-map is

Solution:

		AB			
	CD	00	01	11	10
00		0	0	1	1
01		0	x	x	1
11		x	x	1	x
10		1	0	1	1

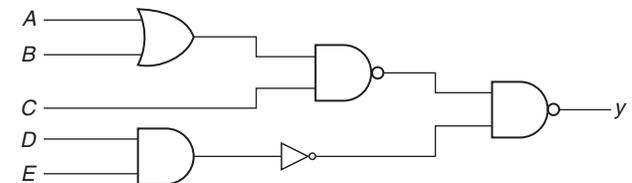
$= A + \bar{B}C$

Example 6: In the figure shown, y_2, y_1, y_0 will be 1s complement of x_2, x_1, x_0 if $z = ?$

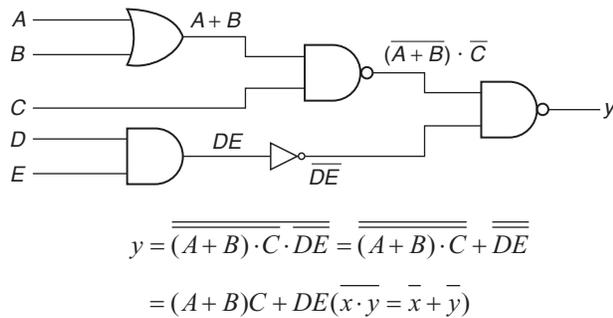


Solution: We are using X-OR gate
 \therefore XOR out-put is complement of input only when other input is high.
 $\therefore Z = 1$

Example 7: The output y of the circuit shown in the figure is



Solution:



Example 8: Simplify the following function

$$f = \overline{\overline{\overline{A(AB)}} \cdot \overline{\overline{\overline{B(AB)}}}}$$

Solution: $\overline{\overline{\overline{A(AB)}} \cdot \overline{\overline{\overline{B(AB)}}}}$

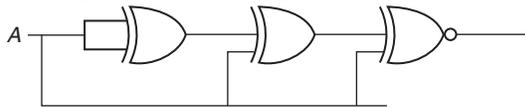
$$\begin{aligned} \overline{[A+(AB)]} \cdot \overline{[B+(AB)]} &= \overline{A+(AB)} + \overline{B+(AB)} \\ &= \overline{A} \cdot \overline{(AB)} + \overline{B} \cdot \overline{(AB)} = \overline{A} \cdot (\overline{A} + \overline{B}) + \overline{B} (\overline{A} + \overline{B}) \\ &= \overline{A} \cdot \overline{A} + \overline{A} \overline{B} + \overline{B} \overline{A} + \overline{B} \cdot \overline{B} = \overline{A} + \overline{B} = \overline{AB} \end{aligned}$$

EXERCISES

Practice Problems I

Directions for questions 1 to 25: Select the correct alternative from the given choices.

1. The output of the following circuit is



- (A) 0 (B) 1
(C) A (D) A'

2. The circuit which will work as OR gate in positive level will work as ___ gate in negative level logic

- (A) NOR gate
(B) NAND gate
(C) Both NAND and NOR gate
(D) AND gate

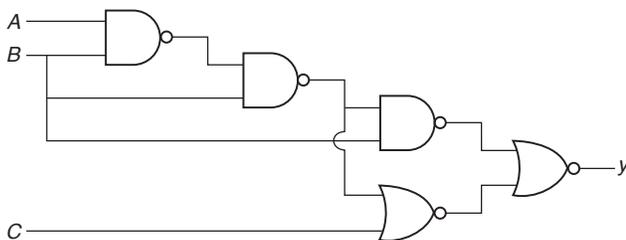
3. Four logical expressions are given below:

- (a) $\overline{A} \cdot \overline{B} \cdot \overline{C} \cdot \overline{D} \cdot \overline{E} \cdot \overline{F} \cdot \overline{G} \cdot \overline{H}$
(b) $\overline{AB} \cdot \overline{CD} \cdot \overline{EF} \cdot \overline{GH}$
(c) $\overline{A} + \overline{B} + \overline{C} + \overline{D} + \overline{E} + \overline{F} + \overline{G} + \overline{H}$
(d) $(\overline{A} + \overline{B})(\overline{C} + \overline{D})(\overline{E} + \overline{F})(\overline{G} + \overline{H})$

Two of these expression are equal. They are

- (A) c and d (B) b and d
(C) a and b (D) a and c

4. For the logic circuit shown in figure, the simplified Boolean expression for the output y is



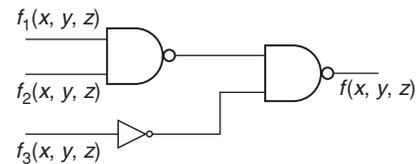
- (A) $A+B+C$ (B) $A \overline{B}$
(C) ABC (D) \overline{BC}

5. In a digital system, there are three inputs A, B and C. The output should be high when at least two inputs

are high. The minimized Boolean expression for the output is

- (A) $AB + BC + AC$
(B) $ABC + ABC + \overline{ABC} + \overline{ABC}$
(C) $ABC + \overline{ABC} + \overline{ABC}$
(D) $\overline{AB} + \overline{BC} + \overline{AC}$

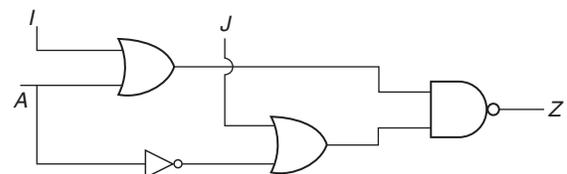
6. Consider the following logic circuit whose inputs are functions f_1, f_2, f_3 and output is f.



Given that $f_1(x, y, z) = \bullet(0, 1, 3, 5)$, $f_2(x, y, z) = \bullet(6, 7)$ and $f_3(x, y, z) = \bullet(1, 4, 5)$, then f_3 is

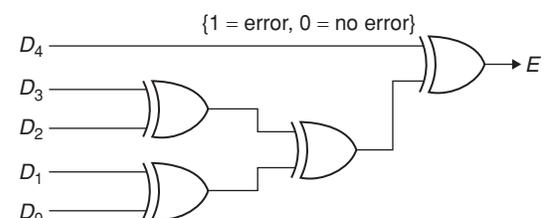
- (A) $\bullet(1, 4, 5)$ (B) $\bullet(6, 7)$
(C) $\bullet(0, 1, 3, 5)$ (D) None of these

7. The circuit shown above is to be used to implement the function $z = f(A, B) = \overline{A} + B$ what values are to be selected for I and J?



- (A) $I = 0, J = B$ (B) $I = 1, J = B$
(C) $I = B, J = 1$ (D) $I = B, J = 0$

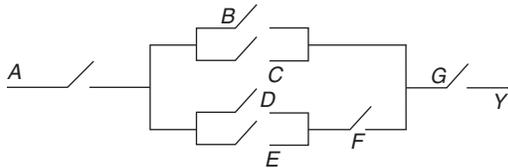
8. Parity checker output from the below figure, if input is 11111 ($D_4 D_3 D_2 D_1 D_0$) and 10000 ($D_4 D_3 D_2 D_1 D_0$).



Practice Problems 2

Directions for questions 1 to 25: Select the correct alternative from the given choices.

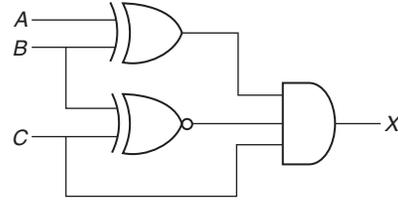
- An OR Gate has six inputs. How many input words are there in its truth table?
 (A) 6 (B) 36
 (C) 32 (D) 64
- Sum of product form can be implemented by using
 (A) AND–OR
 (B) NAND–NAND
 (C) NOR–NOR
 (D) Both A and B
- Which one of the following is equivalent to the Boolean expression?
 $Y = AB + BC + CA$
 (A) $\overline{AB + BC + CA}$
 (B) $(\overline{A+B})(\overline{B+C})(\overline{A+C})$
 (C) $\overline{(A+B)(B+C)(A+C)}$
 (D) $\overline{(\overline{A+B})(\overline{B+C})(\overline{C+A})}$
- What Boolean function does the following circuit represents?



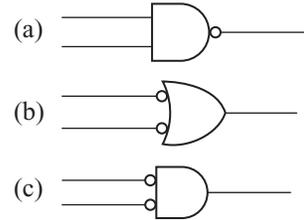
- $A [F + (B + C) \cdot (D + E)] G$
 - $A + BC + DEF + G$
 - $A [(B + C) + F (D + E)] G$
 - $ABG + ABC + F(D + E)$
- The minimum number of two input NOR gates are required to implement the simplified value of the following equation
 $f(w, x, y, z) = \sum m(0, 1, 2, 3, 8, 9, 10, 11)$
 (A) One (B) Two
 (C) Three (D) Four
- The out put of a logic gate is '1' when all inputs are at logic '0'. Then the gate is either
 (1) NAND or X-OR gate
 (2) NOR or X-OR gate
 (3) NOR or X-NOR gate
 (4) NAND or X-NOR gate
 (A) 1 and 2 (B) 2 and 3
 (C) 3 and 4 (D) 4 and 1
- If the functions $w, x, y,$ and z are as follows.
 $w = R + \overline{PQ} + \overline{RS}$
 $x = \overline{PQRS} + \overline{PQ\overline{RS}} + \overline{P\overline{Q}RS}$
 $y = \overline{RS} + \overline{PR} + \overline{PQ} + \overline{P} \cdot \overline{Q}$
 $z = \overline{R + S + PQ + P \cdot Q \cdot R + P \cdot Q \cdot S}$

- $w = z, x = y$ (B) $w = z, x = \overline{z}$
- $w = y$ (D) $w = y = z$

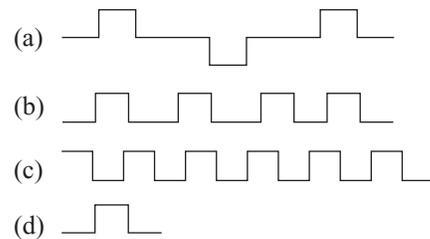
- For the logic circuit shown in the figure, the required input condition (A, B, C) to make the output (x) = 1 is



- 0, 0, 1 (B) 1, 0, 1
 - 1, 1, 1 (D) 0, 1, 1
- Which of the following is a basic gate?
 (A) AND (B) X-OR
 (C) X-NOR (D) NAND
 - Which of the following represent the NAND gate?



- a only (B) a, b, c
 - b, a (D) a, c
- The universal gates are
 (A) NAND and NOR (B) AND, OR, NOT
 (C) X-OR and X-NOR (D) All of these
 - In the circuit the value of input A goes from 0 to 1 and part of B goes from 1 to 0. Which of the following represent output under a static hazard condition?

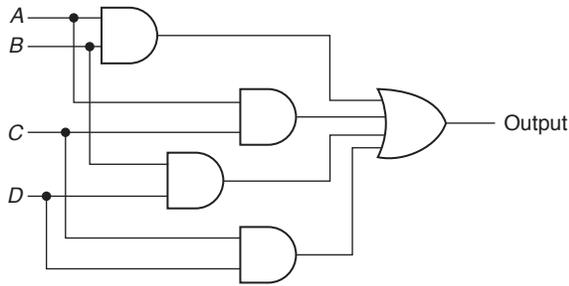


- Output a (B) Output b
 - Output c (D) Output d
- The consensus theorem states that
 (A) $A + \overline{A}B = A + B$
 (B) $A + AB = A$
 (C) $AB + \overline{A}C + BC = AB + \overline{A}C$
 (D) $(A + B) \cdot (A + \overline{B}) = A$
 - The dual form of expression
 $AB + \overline{A}C + BC = AB + \overline{A}C$ is

- (A) $(A+B)(\bar{A}+C)(B+C) = (A+B)(\bar{A}+C)$
 (B) $(A+B)(\bar{A}+C)(B+C) = (\bar{A}+\bar{B})(A+\bar{C})$
 (C) $(\bar{A}+\bar{B})(\bar{A}+\bar{C})(\bar{B}+\bar{C}) = (\bar{A}+\bar{B})(A+\bar{C})$
 (D) $\bar{A}\bar{B}+A\bar{C}+\bar{B}\bar{C} = \bar{A}\bar{B}+A\bar{C}$

15. The max term corresponding to decimal 12 is
 (A) $\bar{A}+\bar{B}+C+D$ (B) $A+B+\bar{C}+\bar{D}$
 (C) $\bar{A}\bar{B}CD$ (D) $ABC\bar{D}$

16. The given circuit is equivalent to



- (A) $(A+C)(B+D)$ (B) $AC+BD$
 (C) $(A+D)(B+C)$ (D) $(\bar{A}+\bar{B})(\bar{C}+\bar{D})$

17. Minimized expression for Karnaugh map is

AB \ C	00	01	11	10
0	1			1
1	1			1

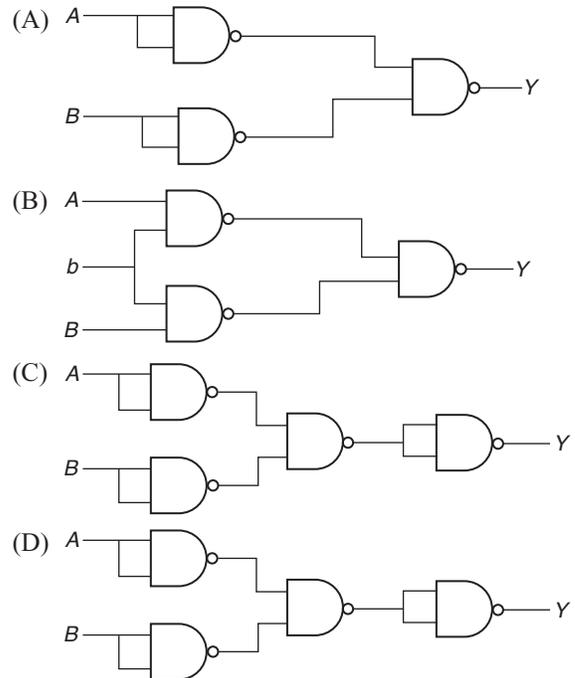
- (A) $AB+C$ (B) $\bar{A}B+C$
 (C) \bar{B} (D) $\bar{B}+C$

18. An XOR gate will act as _____ when one of its input is one and as _____ when one of its input is zero.
 (A) buffer, buffer (B) buffer, inverter
 (C) inverter, buffer (D) inverter, inverter
19. The minimum number of two input NAND gates required to implement $A \odot B$ if only A and B are available
 (A) 6 (B) 3
 (C) 5 (D) 4
20. Negative logic in a logic circuit is one in which
 (A) logic 0 and 1 are represented by GND and positive voltage respectively.
 (B) logic 0 and 1 are represented by negative and positive voltage.

- (C) logic 0 voltage level is lower than logic 1 voltage level.
 (D) logic 0 voltage level is higher than logic 1 voltage level.

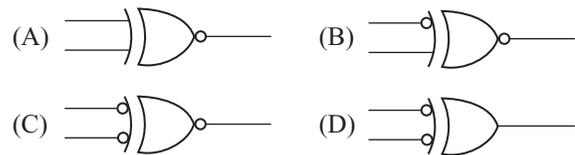
21. If the input to a gate is eight in number, then its truth table contains _____ input words.
 (A) 128 (B) 8
 (C) 64 (D) 256

22. The X-OR gate implementation using NAND gate is



23. The equivalent of AND-OR logic circuit is
 (A) NAND-NOR (B) NOR-AND
 (C) NAND-NAND (D) NAND-OR

24. The X-OR is equivalent to



25. Simplify $A\bar{B}C + B + B\bar{D} + AB\bar{D} + \bar{A}C$
 (A) B (B) $B+C$
 (C) $C+A$ (D) $\bar{A}+B$

PREVIOUS YEARS' QUESTIONS

- Let $f(w, x, y, z) = \sum(0, 4, 5, 7, 8, 9, 13, 15)$. Which of the following expressions are NOT equivalent to f ? [2007]
 - (P) $x'y'z' + w'xy' + wy'z + xz$
 - (Q) $w'y'z' + wx'y' + xz$
 - (R) $w'y'z' + wx'y' + xyz + xy'z$
 - (S) $x'y'z' + wx'y' + w'y$

(A) P only (B) Q and S
(C) R and S (D) S only
- Define the connective $*$ for the Boolean variables X and Y as: $X * Y = XY + X'Y'$. Let $Z = X * Y$. Consider the following expressions P, Q and R . [2007]

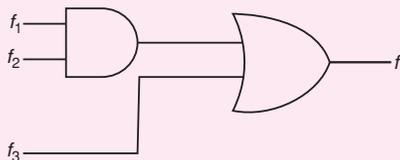
$P: X = Y * Z$ $Q: Y = X * Z$
 $R: X * Y * Z = 1$

Which of the following is TRUE?

(A) Only P and Q are valid.
(B) Only Q and R are valid.
(C) Only P and R are valid.
(D) All P, Q, R are valid.
- In the Karnaugh map shown below, \times denotes a don't-care term. What is the minimal form of the function represented by the Karnaugh map? [2008]

		ab			
		00	01	11	10
cd	00	1	1		1
	01	\times	1		
	11	\times			
	10	1	1		\times

- (A) $\bar{b} \cdot \bar{d} + \bar{a} \cdot \bar{d}$
 (B) $\bar{a} \cdot \bar{b} + \bar{b} \cdot \bar{d} + \bar{a} \cdot b \cdot \bar{d}$
 (C) $\bar{b} \cdot \bar{d} + \bar{a} \cdot b \cdot \bar{d}$
 (D) $\bar{a} \cdot \bar{b} + \bar{b} \cdot \bar{d} + \bar{a} \cdot \bar{d}$
- Given f_1, f_3 and f in canonical sum of products form (in decimal) for the circuit [2008]



$f_1 = \sum m(4, 5, 6, 7, 8)$

$f_3 = \sum m(1, 6, 15)$

$f = \sum m(1, 6, 8, 15)$

then f_2 is

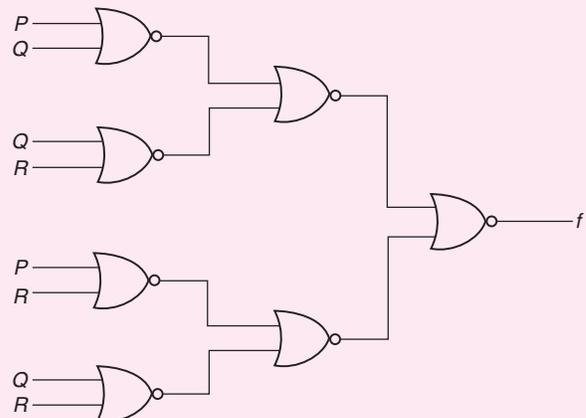
- (A) $\sum m(4, 6)$ (B) $\sum m(4, 8)$
(C) $\sum m(6, 8)$ (D) $\sum m(4, 6, 8)$

- If $\bar{P}, \bar{Q}, \bar{R}$ are Boolean variables, then $(P + \bar{Q})(\bar{P} \cdot \bar{Q} + P \cdot R)(\bar{P} \cdot \bar{R} + \bar{Q})$ simplifies to [2008]
 - (A) $P \cdot \bar{Q}$ (B) $P \cdot \bar{R}$
 - (C) $P \cdot \bar{Q} + R$ (D) $P \cdot \bar{R} + Q$
- What is the minimum number of gates required to implement the Boolean function $(AB + C)$, if we have to use only two-input NOR gates? [2009]
 - (A) 2 (B) 3
 - (C) 4 (D) 5
- The binary operation \square is defined as follows [2009]

P	Q	$P \square Q$
T	T	T
T	F	T
F	T	F
F	F	T

Which one of the following is equivalent to $P \square Q$?

- (A) $\neg Q \neg P$ (B) $P \square \neg Q$
 (C) $\neg P \square Q$ (D) $\neg P \square \neg Q$
- The min term expansion of $f(P, Q, R) = PQ + Q\bar{R} + P\bar{R}$ is [2010]
 - (A) $m_2 + m_4 + m_6 + m_7$
 - (B) $m_0 + m_1 + m_3 + m_5$
 - (C) $m_0 + m_1 + m_6 + m_7$
 - (D) $m_2 + m_3 + m_4 + m_5$
 - What is the Boolean expression for the output f of the combinational logic circuit of NOR gates given below? [2010]



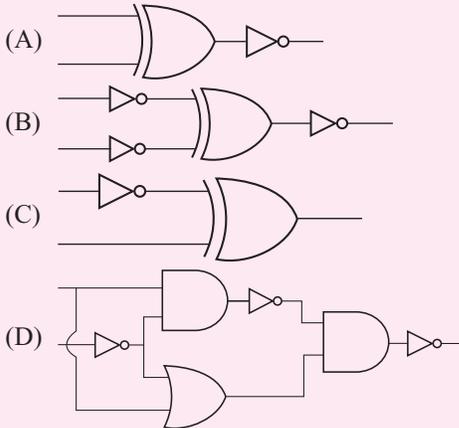
- (A) $\overline{Q + R}$ (B) $\overline{P + Q}$
(C) $\overline{P + R}$ (D) $\overline{P + Q + R}$

10. The simplified SOP (Sum of Product) form of the Boolean expression. [2011]

$(P + \bar{Q} + \bar{R})(P + \bar{Q} + R)(P + Q + \bar{R})$ is

(A) $(PQ + \bar{R})$ (B) $(P + \bar{Q}\bar{R})$
 (C) $(\bar{P}Q + R)$ (D) $(PQ + R)$

11. Which one of the following circuits is NOT equivalent to a two-input X-NOR (exclusive NOR) gate? [2011]



12. The truth table [2012]

X	Y	F(X, Y)
0	0	0
0	1	0
1	0	1
1	1	1

represents the Boolean function

- (A) X (B) $X + Y$
 (C) $X \oplus Y$ (D) Y
13. What is the minimal form of the Karnaugh map shown below? Assume that \times denotes a don't-care term. [2012]

ab \ cd	00	01	11	10
00	1	\times	\times	1
01	\times			1
11				
10	1			\times

- (A) $\bar{b}\bar{d}$ (B) $\bar{b}\bar{d} + \bar{b}c$
 (C) $\bar{b}\bar{d} + \bar{a}bcd$ (D) $\bar{b}\bar{d} + \bar{b}c + \bar{c}d$
14. Which one of the following expressions does not represent exclusive NOR of x and y ? [2013]
- (A) $xy + x'y'$ (B) $x \oplus y'$
 (C) $x' \oplus y$ (D) $x' \oplus y'$

15. Consider the following Boolean expression for F :
 $F(P, Q, R, S) = PQ + \bar{P}QR + \bar{P}Q\bar{R}S$

The minimal sum of products form of F is [2014]

(A) $PQ + QR + QS$ (B) $P + Q + R + S$
 (C) $\bar{P} + \bar{Q} + \bar{R} + \bar{S}$ (D) $\bar{P}R + \bar{P}\bar{R}S + P$

16. The dual of a Boolean function $f(x_1, x_2, \dots, x_n, +, \cdot, ', \bar{})$, written as F^D , is the same expression as that of F with $+$ and \cdot swapped. F is said to be self-dual if $F = F^D$. The number of self-dual functions with n Boolean variables is [2014]

(A) 2^n (B) 2^{n-1}
 (C) 2^{2^n} (D) $2^{2^{n-1}}$

17. Consider the following min term expression for F :

$F(P, Q, R, S) = \sum m(0, 2, 5, 7, 8, 10, 13, 15)$

The min terms 2, 7, 8 and 13 are 'don't-care terms'. The minimal sum of products form for F is [2014]

(A) $Q\bar{S} + \bar{Q}S$
 (B) $\bar{Q}\bar{S} + QS$
 (C) $\bar{Q}\bar{R}\bar{S} + \bar{Q}R\bar{S} + Q\bar{R}S + QRS$
 (D) $\bar{P}\bar{Q}\bar{S} + \bar{P}QS + PQS + P\bar{Q}\bar{S}$

18. The binary operator \neq is defined by the following truth table

p	q	p \neq q
0	0	0
0	1	1
1	0	1
1	1	0

Which one of the following is true about the binary operator \neq ? [2015]

- (A) Both commutative and associative
 (B) Commutative but not associative
 (C) Not commutative but associative
 (D) Neither commutative nor associative
19. Consider the operations [2015]

$f(X, Y, Z) = X^1 YZ + XY^1 + Y^1 Z^1$ and

$g(X, Y, Z) = X^1 YZ + X^1 YZ^1 + XY$.

Which one of the following is correct?

- (A) Both $\{f\}$ and $\{g\}$ are functionally complete
 (B) Only $\{f\}$ is functionally complete
 (C) Only $\{g\}$ is functionally complete
 (D) Neither $\{f\}$ nor $\{g\}$ is functionally complete
20. The number of min-terms after minimizing the following Boolean expression is [2015]

$[D^1 + AB^1 + A^1C + AC^1D + A^1C^1D]^1$

21. Let $\#$ be a binary operator defined as [2015]

$X\#Y = X^1 + Y^1$ where X and Y are Boolean variables.

Consider the following two statements.

(S1) $(P\#Q)\#R = P\#(Q\#R)$

(S2) $Q\#R = R\#Q$

Which of the following is/are true for the Boolean variables P, Q and R ?

- (A) Only S_1 is true
 (B) Only S_2 is true

- (C) Both S_1 and S_2 are true
 (D) Neither S_1 nor S_2 are true

22. Given the function $F = P^1 + QR$, where F is a function in three Boolean variables P , Q and R and $P^1 = !P$, consider the following statements. [2015]

- (S_1) $F = \Sigma(4, 5, 6)$
 (S_2) $F = \Sigma(0, 1, 2, 3, 7)$
 (S_3) $F = \pi(4, 5, 6)$
 (S_4) $F = \pi(0, 1, 2, 3, 7)$

Which of the following is true?

- (A) (S_1) – False, (S_2) – True, (S_3) – True, (S_4) – False
 (B) (S_1) – True, (S_2) – False, (S_3) – False, (S_4) – True
 (C) (S_1) – False, (S_2) – False, (S_3) – True, (S_4) – True
 (D) (S_1) – True, (S_2) – True, (S_3) – False, (S_4) – False

23. The total number of prime implicants of the function $f(w, x, y, z) = \Sigma(0, 2, 4, 5, 6, 10)$ is _____. [2015]

24. Consider the Boolean operator # with the following properties: [2016]

$x \# 0 = x, x \# 1 = \bar{x}, x \# x = 0$ and

$x \# \bar{x} = 1$. Then $x \# y$ is equivalent to

- (A) $x \bar{y} + \bar{x} y$ (B) $x \bar{y} + \bar{x} \bar{y}$
 (C) $\bar{x} y + x y$ (D) $x y + \bar{x} \bar{y}$

25. Consider the Karnaugh map given below, where X represents “don’t care” and blank represents 0.

		ba			
		00	01	11	10
dc					
00			X	X	
01		1			X
11		1			1
10			X	X	

Assume for all inputs (a, b, c, d) , the respective complements $(\bar{a}, \bar{b}, \bar{c}, \bar{d})$ are also available. The above logic is implemented using 2-input NOR gates only. The minimum number of gates required is _____.

[2017]

26. If w, x, y, z are Boolean variables, then which one of the following is INCORRECT? [2017]

- (A) $wx + w(x+y) + x(x+y) = x + wy$
 (B) $\overline{w\bar{x}(y+\bar{z})} + \bar{w}x = \bar{w} + x + \bar{y}z$
 (C) $(w\bar{x}(y+x\bar{z}) + \bar{w}\bar{x})y = x\bar{y}$
 (D) $(w+y)(wxy + wyz) = wxy + wyz$

27. Given $f(w, x, y, z) = \Sigma m(0,1,2,3,7,8,10) + \Sigma d(5,6,11,15)$, where d represents the don’t-care condition in Karnaugh maps. Which of the following is a minimum product-of-sums (POS) form of $f(w, x, y, z)$? [2017]

- (A) $f = (\bar{w} + \bar{z})(\bar{x} + z)$
 (B) $f = (\bar{w} + z)(x + z)$
 (C) $f = (w + z)(\bar{x} + z)$
 (D) $f = (w + \bar{z})(\bar{x} + z)$

28. Let \oplus and \odot denote the Exclusive OR and Exclusive NOR operations, respectively. Which one of the following is NOT CORRECT? [2018]

- (A) $\overline{P \oplus Q} = P \odot Q$
 (B) $\bar{P} \oplus Q = P \odot Q$
 (C) $\bar{P} \oplus \bar{Q} = P \oplus Q$
 (D) $(P \oplus \bar{P}) \oplus Q = (P \odot \bar{P}) \odot \bar{Q}$

29. Consider the minterm list form of a Boolean function F given below.

$$F(P, Q, R, S) = \sum m(0, 2, 5, 7, 9, 11)$$

$$+ d(3, 8, 10, 12, 14)$$

Here, m denotes a minterm and d denotes a don’t care term. The number of essential prime implicants of the function F is _____. [2018]

ANSWER KEYS**EXERCISES****Practice Problems 1**

- | | | | | | | | | | |
|-------|-------|-------|-------|-------|-------|-------|-------|-------|-------|
| 1. B | 2. D | 3. B | 4. C | 5. A | 6. A | 7. B | 8. A | 9. B | 10. A |
| 11. C | 12. D | 13. A | 14. A | 15. A | 16. A | 17. D | 18. B | 19. A | 20. D |
| 21. D | 22. C | 23. A | 24. C | 25. D | | | | | |

Practice Problems 2

- | | | | | | | | | | |
|-------|-------|-------|-------|-------|-------|-------|-------|-------|-------|
| 1. D | 2. D | 3. D | 4. C | 5. A | 6. C | 7. B | 8. D | 9. A | 10. C |
| 11. A | 12. D | 13. A | 14. A | 15. A | 16. C | 17. C | 18. C | 19. C | 20. D |
| 21. D | 22. C | 23. C | 24. B | 25. B | | | | | |

Previous Years' Questions

- | | | | | | | | | | |
|-------|-------|-------|-------|-------|-------|-------|-------|-------|-------|
| 1. D | 2. D | 3. A | 4. C | 5. A | 6. B | 7. B | 8. A | 9. A | 10. B |
| 11. D | 12. A | 13. B | 14. D | 15. A | 16. D | 17. B | 18. A | 19. B | 20. 1 |
| 21. B | 22. A | 23. 3 | 24. A | 25. 1 | 26. C | 27. A | 28. D | 29. 3 | |

Chapter 3

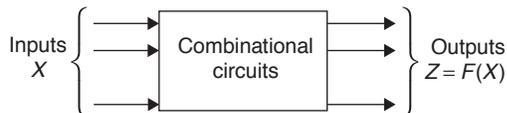
Combinational Circuits

LEARNING OBJECTIVES

- Combinational logic design
- Arithmetic circuit
- Half adder
- Full adder
- Half subtractor
- Full subtractor
- n -bit comparator
- Parity bit generator and parity bit checker
- Code converter
- Decoder
- Designing high order decoders from lower order decoder
- Combinational logic implementation
- Encoders
- Multiplexer
- Demultiplexer

INTRODUCTION

Combinational logic is a type of logic circuit whose output is a function of the present input only.



COMBINATIONAL LOGIC DESIGN

The design of combinational circuit starts from the problem, statement and ends with a gate level circuit diagram.

The design procedure involves the following steps:

- Determining the number of input variables and output variables required, from the specifications.
- Assigning the letter symbols for input and output.
- Deriving the truth table that defines the required relationship between input and output.
- Obtaining the simplified Boolean function for each output by using K-map or algebraic relations.
- Drawing the logic diagram for simplified expressions.

We will discuss combinational circuits under the following categories:

- Arithmetic circuits
- Code converters
- Data processing circuits

ARITHMETIC CIRCUITS

Arithmetic circuits are the circuits that perform arithmetic operation. The most basic arithmetic operation is addition.

Half Adder

Addition is an arithmetic operation, and here to implement addition in digital circuits we have to implement by logical gates. So the addition of binary numbers will be represented by the logical expressions. Half adder is an arithmetic circuit which performs the addition of two binary bits, and the result is viewed in two output—sum and carry.

The sum 'S' is the X-OR of 'A' and 'B' where A and B are inputs.

$$\therefore S = A\bar{B} + B\bar{A} = A \oplus B$$

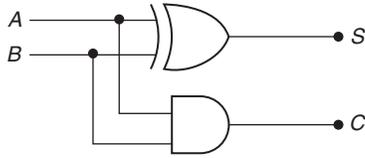
The carry 'C' is the AND of A and B.

$$\therefore C = AB$$

Table 1 Truth Table

Inputs		Outputs	
A	B	S	C
0	0	0	0
0	1	1	0
1	0	1	0
1	1	0	1

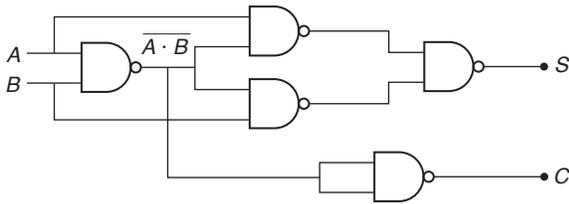
So, half adder can be realized by using one X-OR gate and one AND gate.



Half adder can also be realized by universal logic such as only NAND gate or only NOR gate as given below.

NAND logic

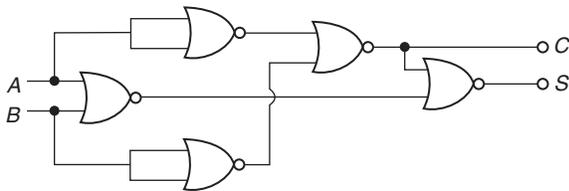
$$\begin{aligned}
 S &= A\bar{B} + \bar{A}B \\
 &= A\bar{B} + A\bar{A} + \bar{A}B + B\bar{B} \\
 &= A(\bar{A} + \bar{B}) + B(\bar{A} + \bar{B}) \\
 &= \overline{\overline{A\bar{B}} \cdot \overline{B\bar{A}}} \\
 C &= AB = \overline{\overline{A \cdot B}}
 \end{aligned}$$



Half adder using NAND logic

NOR logic

$$\begin{aligned}
 S &= A \cdot \bar{B} + \bar{A}B \\
 &= A\bar{B} + A\bar{A} + \bar{A}B + B\bar{B} \\
 &= A(\bar{A} + \bar{B}) + B(\bar{A} + \bar{B}) \\
 &= (A + B)(\bar{A} + \bar{B}) \\
 &= \overline{\overline{A + B} \cdot \overline{\bar{A} + \bar{B}}} \\
 C &= A \cdot B = \overline{\overline{A \cdot B}} = \overline{\bar{A} + \bar{B}}
 \end{aligned}$$



Half adder using NOR logic

Full Adder

Full adder is an arithmetic circuit that performs addition of two bits with carry input. The result of full adder is given by two outputs—sum and carry. The full adder circuit is used in parallel adder circuit as well as in serial adder circuit.

For full adder, if total number of 1's is odd at input lines, the sum output is equal to logic 1, and if total number of 1's at input lines are more than or equal to 2, then the carry output is logic 1.

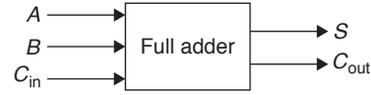


Figure 1 Block diagram

Table 2 Truth Table

A	B	C _{in}	S	C _{out}
0	0	0	0	0
0	0	1	1	0
0	1	0	1	0
0	1	1	0	1
1	0	0	1	0
1	0	1	0	1
1	1	0	0	1
1	1	1	1	1

$$\begin{aligned}
 S &= \bar{A}\bar{B}C_{in} + \bar{A}B\bar{C}_{in} + A\bar{B}\bar{C}_{in} + ABC_{in} \\
 &= A \oplus B \oplus C_{in} \\
 C_{out} &= \bar{A}BC_{in} + A\bar{B}C_{in} + AB\bar{C}_{in} + ABC_{in} \\
 &= AB + (A \oplus B)C_{in} \\
 &= AB + AC_{in} + BC_{in}
 \end{aligned}$$

Full adder can also be realized using universal logic gates, i.e., either only NAND gates or only NOR gates as explained below.

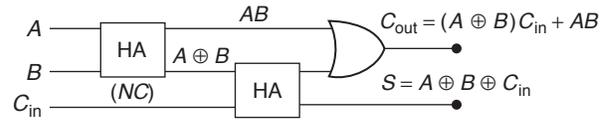


Figure 2 Block diagram of full adder by using Half adder

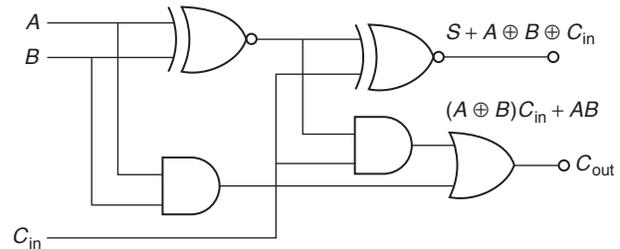


Figure 3 Logic diagram of full adder

NAND logic

$$A \oplus B = \overline{\overline{A \cdot B} \cdot \overline{B \cdot A}}$$

So $A \oplus B \oplus C_{in}$

$$\begin{aligned}
 \text{Let } A \oplus B = x \text{ then } s &= \overline{\overline{x \cdot C_{in}} \cdot \overline{C_{in} \cdot x}} \\
 &= x \oplus C_{in}
 \end{aligned}$$

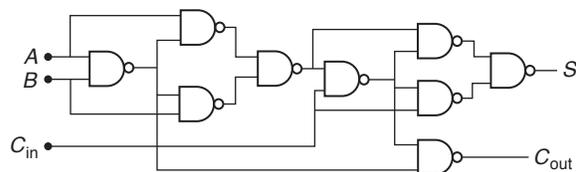


Figure 4 Logic diagram of a full adder using only 2-input NAND gates

NOR logic

Full adder outputs

Sum = $a \oplus b \oplus c$, carry = $ab + bc + ac$ are self dual functions

[∵ A function is called as self dual if its dual is same as the function itself $f^D = f$]

For self dual functions, the number of NAND gates are same as number of NOR gates.

By taking the dual for above NAND gate implementation, all gates will become NOR gates, and the output is dual of the sum and carry, but they are self dual ($f^D = f$).

So, output remain same, and only 9 NOR gates are required for full adder, structure similar to NAND gate circuit.

Half Subtractor

Half subtractor is an arithmetic circuit which performs subtraction of one bit (subtrahend) from other bit (minuend), and the result gives difference and borrow each of one bit. The borrow output is logic 1 only if there is any subtraction of 1 from 0.

When a bit 'B' is subtracted from another bit 'A', a difference bit (d) and a borrow bit (b) result according to the rule given below.

Table 3 Truth Table

A	B	d	b
0	0	0	0
1	0	1	0
1	1	0	0
0	1	1	1

$$d = A\bar{B} + B\bar{A}$$

$$= A \oplus B$$

$$b = \bar{A}B$$

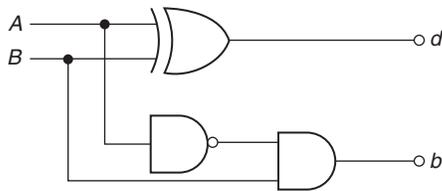


Figure 5 Logic diagram of a half subtractor

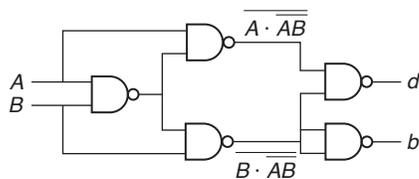
A half subtractor can also realized using universal logic either using only NAND gates or only NOR gates as explained below.

NAND logic

$$d = A \oplus B$$

$$= \overline{\overline{A \cdot B} \cdot \overline{B \cdot A}}$$

$$b = \bar{A}B = B(\bar{A} + \bar{B}) = \overline{\overline{B \cdot \bar{A}B}}$$



NOR logic

$$d = A \oplus B$$

$$= \overline{\overline{A\bar{B}} + \overline{\bar{A}B}}$$

$$= \overline{\overline{A\bar{B}} + \overline{\bar{A}B} + \overline{AA}}$$

$$= \overline{\overline{B(A+B)} + \overline{\bar{A}(A+B)}}$$

$$= \overline{\overline{B+A+B} + \overline{A+A+B}}$$

$$b = \bar{A}B$$

$$= \overline{\overline{\bar{A}(A+B)}}$$

$$= \overline{\overline{\bar{A}(A+B)}}$$

$$= A + \overline{(A+B)}$$

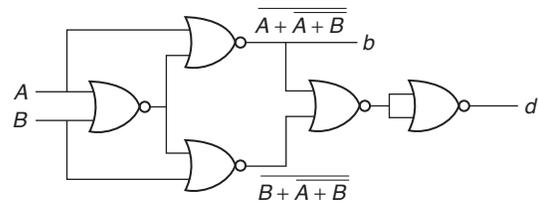


Figure 6 Logic diagram of half subtractor using NOR gate

Full Subtractor

Full subtractor is an arithmetic circuit similar to half subtractor but it performs subtraction with borrow, it involves subtraction of three bits—minuend, subtrahend and borrow-in, and two outputs—difference and borrow. The subtraction of 1 from 0 results in borrow to become logic 1. The presence of odd number of 1's at input lines make difference as logic 1.

Table 4 Truth Table

A	B	b_i	d	b
0	0	0	0	0
0	0	1	1	1
0	1	0	1	1
0	1	1	0	1
1	0	0	1	0
1	0	1	0	0
1	1	0	0	0
1	1	1	1	1

$$d = \bar{A}\bar{B}b_i + \bar{A}B\bar{b}_i + A\bar{B}\bar{b}_i + ABb_i$$

$$= b_i(AB + \bar{A}\bar{B}) + \bar{b}_i(\bar{A}B + A\bar{B})$$

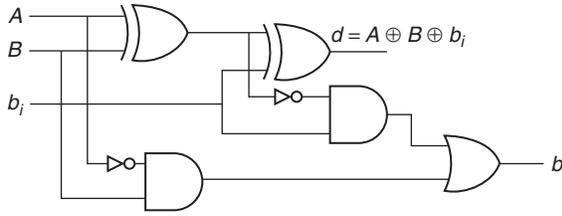
$$= b_i(A \oplus B) + \bar{b}_i(A \oplus B)$$

$$= A \oplus B \oplus b_i$$

and

$$b = \bar{A}\bar{B}b_i + \bar{A}B\bar{b}_i + \bar{A}Bb_i$$

$$= \bar{A}B + (\bar{A} \oplus B)b_i$$



NAND logic

$$d = A \oplus B \oplus b_i$$

$$= \overline{(A \oplus B)(A \oplus B)b_i b_i(A \oplus B)b_i}$$

$$b = \bar{A}B + b_i(\bar{A} \oplus B)$$

$$= \overline{\overline{\bar{A}B} + b_i(\bar{A} \oplus B)}$$

$$= \overline{\bar{A}B + b_i(\bar{A} \oplus B)}$$

$$= \overline{\bar{A}B b_i(A \oplus B)}$$

$$= \overline{B(\bar{A} + \bar{B}) b_i [b_i + (A \oplus B)]}$$

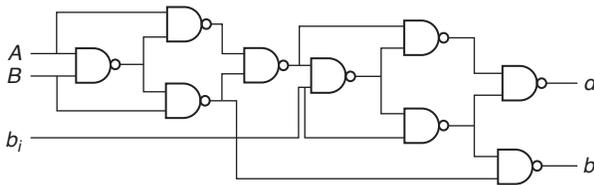


Figure 7 Logic diagram of a full subtractor using NAND logic

NOR logic

Output of full subtractor is also self dual in nature. So, same circuit, with all NAND gates, replaced by NOR gates gives the NOR gate full subtractor. 9 NOR gates required.

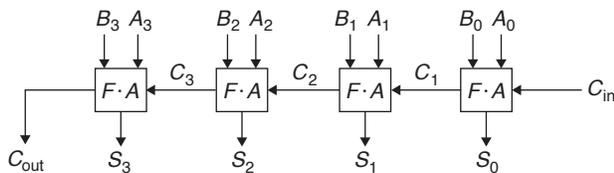
Example 1: How many NAND gates are required for implementation of full adder and full subtractor respectively?
 (A) 11, 10 (B) 11, 11 (C) 9, 9 (D) 9, 10

Solution: (C)

From the circuit diagrams in the previous discussion, full adder requires 9 NAND gates and full subtractor requires 9 NAND gates.

Binary Adder

A binary adder is a digital circuit that produces the arithmetic sum of two binary numbers.



Four bit parallel adder the output carry from each full adder is connected to the input carry of next full adder.

The bits are added with full adders, starting from the LSB position to form the sum bit and carry bit.

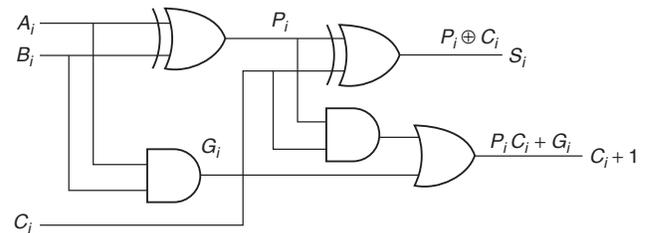
The longest propagation delay time in parallel adder is the time it takes the carry to propagate through the full adders.

For n -bit parallel adders consider t_{pds} is the propagation delay for sum of each full adder and t_{pdc} is the propagation delay of carry.

The total time required to add all n -bits at the n th full adder is

$$T_S = t_{pds} + (n - 1)t_{pdc}$$

So propagation delay increases with number of bits. To overcome this difficulty we use look ahead carry adder, which is the fastest carry adder.



Consider the full adder circuit for i th stage, in parallel adder, with two binary variables A_i, B_i , input carry C_i are:

Carry propagate (P_i) and carry generate (G_i)

$$P_i = A_i \oplus B_i$$

$$G_i = A_i \cdot B_i$$

The output sum and carry can be expressed as

$$S_i = P_i \oplus C_i$$

$$C_{i+1} = P_i C_i + G_i$$

Now, the Boolean functions for each stage can be calculated as substitute $i = 0$

C_0 is input carry

$$C_1 = G_0 + P_0 C_0$$

Substitute $i = 1, 2 \dots$

$$C_2 = G_1 + P_1 C_1 = G_1 + P_1 (G_0 + P_0 C_0)$$

$$= G_1 + P_1 G_0 + P_1 P_0 C_0$$

$$C_3 = G_2 + P_2 C_2 = G_2 + P_2 (G_1 + P_1 G_0 + P_1 P_0 C_0)$$

$$= G_2 + P_2 G_1 + P_2 P_1 G_0 + P_2 P_1 P_0 C_0$$

Since the Boolean function for each output carry is expressed in SOP form, each function can be implemented with AND-OR form or two level NAND gates.

From the above equations we can conclude that this circuit can perform addition in less time as C_3 does not have to wait for C_2 and C_1 to propagate. C_3, C_2, C_1 can have equal time delays.

The gain in speed of operation is achieved at the expense of additional complexity (hardware).

n-bit Comparator

The comparison of two numbers is an operation that determines whether one number is greater than, less than, or equal to the other number.

A magnitude comparator is a combinational circuit that compares two input numbers A and B , and specifies the output with three variables, $A > B$, $A = B$, $A < B$:

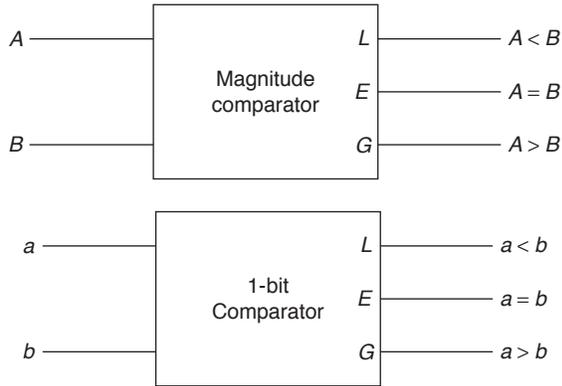


Figure 8 1-bit comparator will have only 1 bit input a, b .

a	b	$a < b$	$a = b$	$a > b$
0	0	0	1	0
0	1	1	0	0
1	0	0	0	1
1	1	0	1	0

By considering minterms for each output.

$$\begin{aligned} (a < b) &= a'b \\ (a = b) &= a'b' + ab = a \odot b \\ (a > b) &= ab' \end{aligned}$$

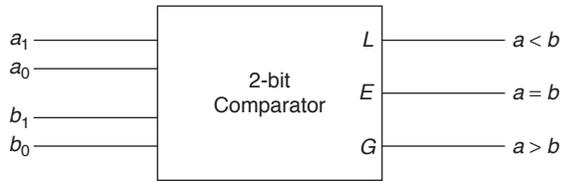
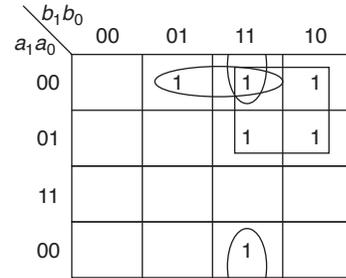


Figure 9 2-bit comparator will have 2-bit inputs a_1, a_0 and b_1, b_0 .

a_1	a_0	b_1	b_0	L $a < b$	E $a = b$	G $a > b$
0	0	0	0	0	1	0
0	0	0	1	1	0	0
0	0	1	0	1	0	0
0	0	1	1	1	0	0
0	1	0	0	0	0	1
0	1	0	1	0	1	0
0	1	1	0	1	0	0
0	1	1	1	1	0	0
1	0	0	0	0	0	1
1	0	0	1	0	0	1
1	0	1	0	0	1	0

1	0	1	1	1	0	0
1	1	0	0	0	0	1
1	1	0	1	0	0	1
1	1	1	0	0	0	1
1	1	1	1	0	1	0

$$\begin{aligned} (a < b) &= \Sigma(1, 2, 3, 6, 7, 11) \\ (a > b) &= \Sigma(4, 8, 9, 12, 13, 14) \\ (a = b) &= \Sigma(0, 5, 10, 15) \end{aligned}$$



$$\begin{aligned} a < b &= \overline{a_1}a_0'b_0 + a_0'b_1b_0 + a_1'b_1 \\ L &= \overline{a_1}b_1 + (a_1 \odot b_1)\overline{a_0}b_0 \end{aligned}$$

Similarly, $a > b = a_0b_1'b_0' + a_1a_0b_0' + a_1b_1'$

$$G = a_1\overline{b_1} + (a_1 \odot b_1)a_0\overline{b_0}$$

$a = b$ is possible when $a_1 = b_1, a_0 = b_0$

$$\text{So } (a = b) = (a_1 \odot b_1)(a_0 \odot b_0)$$

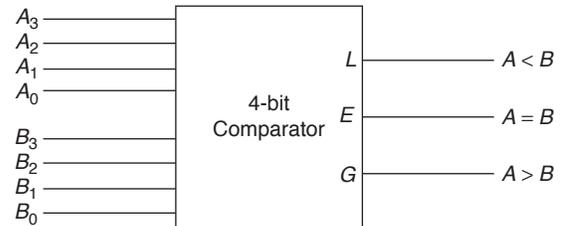


Figure 10 4-bit comparator will compare 2 input numbers each of 4-bits A_3, A_2, A_1, A_0 and B_3, B_2, B_1, B_0 ($A = B$) output will be 1 when each bit of input A is equal to corresponding bit in input B .

So we can write $(A = B) = (A_3 \odot B_3)(A_2 \odot B_2)(A_1 \odot B_1)(A_0 \odot B_0)$.

To determine whether A is greater or less than B , we inspect the relative magnitudes of pairs of significant bits, starting from MSB. If the two bits of a pair are equal, we compare the next lower significant pair of bits. The comparison continues until a pair of unequal bits is reached.

for $A < B, A = 0, B = 1$

for $A > B, A = 1, B = 0$

$$\begin{aligned} A < B &= A_3'B_3 + (A_3 \odot B_3)A_2'B_2 + (A_3 \odot B_3)(A_2 \odot B_2) \\ &\quad \times A_1'B_1 + (A_3 \odot B_3)(A_2 \odot B_2)(A_1 \odot B_1)A_0'B_0 \end{aligned}$$

$$\begin{aligned} A > B &= A_3B_3' + (A_3 \odot B_3)A_2B_2' + (A_3 \odot B_3)(A_2 \odot B_2) \\ &\quad \times A_1B_1' + (A_3 \odot B_3)(A_2 \odot B_2)(A_1 \odot B_1)A_0B_0' \end{aligned}$$

4-bit comparator will have total 8 inputs and $2^8 = 256$ input combinations in truth table.

For 16 combinations ($A = B$) = 1, and for 120 combinations $A < B = 1$.

For remaining 120 combination $A > B = 1$

Parity Bit Generator and Parity Bit Checker

When digital information is transmitted, it may not be received correctly by the receiver. To detect one bit error at receiver we can use parity checker.

For detection of error an extra bit, known as parity bit, is attached to each code word to make the number of 1's in the code even (in case of even parity) or odd (in case of odd parity).

For n -bit data, we use n -bit parity generator at the transmitter end. With 1 parity bit and n -bit data, total $n + 1$ bit will be transmitted. At the receiving end $n + 1$ parity checker circuit will be used to check correctness of the data.

For even parity transmission, parity bit will be made 1 or 0 based on the data, so that total $n + 1$ bits will have even number of 1's. For example, if we want to transmit data 1011 by even parity transmission, then we will use parity bit as 1, so data will have even number of 1's, i.e., data transmitted will be 11011. At the receiving end this data will be received and checked for even number of ones.

To transmit data $B_3B_2B_1B_0$ using even parity, we will transmit sequence $P B_3B_2B_1B_0$, where $P = B_3 \oplus B_2 \oplus B_1 \oplus B_0$. (Equation for parity generator)

At the receiving end we will check data received $P B_3B_2B_1B_0$ for error, $E = P \oplus B_3 \oplus B_2 \oplus B_1 \oplus B_0$ (equation for parity checker). If $E = 0$ (no error), or if $E = 1$ (1 bit error).

We use EX-OR gates for even parity generator/checker as EX-OR of bits gives output 1 if there are odd number of 1's else EX-OR output is 0.

Odd parity generator/checker is complement of even parity generator/checker. Odd parity circuits check for presence of odd number of 1's in data.

CODE CONVERTERS

There are many situations where it is desired to convert from one code to another within a system. For example, the information from output of an analog to digital converter is often in gray code, before it can be processed in arithmetic unit, conversion to binary is required.

Let us consider simple example of 3-bit binary to gray code converter. This will have input lines supplied by binary codes and output lines must generate corresponding bit combination in gray code. The combination circuit code converter performs this transformation by means of logic gates.

The output logic expression derived for code converter can be simplified by using the usual techniques including 'don't-care' if any present. For example, BCD code uses only codes from 0000 to 1001 and remaining combinations are treated as don't-care combinations. Similarly, EXS-3 uses only combinations from 0011 to 1100 and remaining combinations are treated as don't-care.

The relationship between the two codes is shown in the following truth table:

Decimal	B_2	B_1	B_0	G_2	G_1	G_0
0	0	0	0	0	0	0
1	0	0	1	0	0	1
2	0	1	0	0	1	1
3	0	1	1	0	1	0
4	1	0	0	1	1	0
5	1	0	1	1	1	1
6	1	1	0	1	0	1
7	1	1	1	1	0	0

For conversion we have to find out minimized functions of

$$G_2(B_2, B_1, B_0) = \sum m(4, 5, 6, 7)$$

$$G_1(B_2, B_1, B_0) = \sum m(2, 3, 4, 5)$$

$$G_0(B_2, B_1, B_0) = \sum m(1, 2, 5, 6)$$

	B_2			
B_0B_1	00	01	11	10
0		1		1
1		1		1

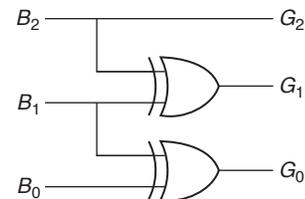
$$G_0(B_2, B_1, B_0) = B'_1 B_0 + B_1 B'_0 = B_1 \oplus B_0$$

	B_2			
B_0B_1	00	01	11	10
0			1	1
1	1	1		

$$G_1(B_2, B_1, B_0) = B'_1 B_2 + B_1 B'_2 = B_2 \oplus B_1$$

	B_2			
B_0B_1	00	01	11	10
0				
1	1	1	1	1

$$G_2(B_2, B_1, B_0) = B_2$$



In similar fashion we can derive n -bit binary to gray code conversion as

$$G_n = B_n$$

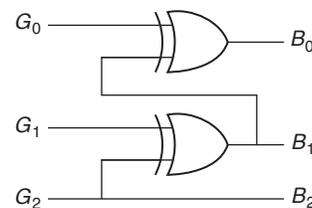
$$G_{n-1} = B_{n-1} \oplus B_n$$

$$G_{i-1} = B_{i-1} \oplus B_i$$

Thus conversion can be implemented by $n - 1$ X-OR gates for n -bits.

For reverse conversion of gray to binary, by following similar standard principle of conversion, we will get

$$B_0 = G_0 \oplus G_1 \oplus G_2, B_1 = G_1 \oplus G_2, B_2 = G_2$$



In general for n -bit gray to binary code conversion

$$B_i = G_n \oplus G_{n-1} \oplus G_{n-2} \dots \oplus G_{i+1} \oplus G_i$$

$B_n = G_n$ (MSB is same in gray and binary). It also requires $n-1$ X-OR gates for n -bits.

Example 2: Design 84-2-1 to XS-3 code converter.

Solution: Both 84-2-1 and XS-3 are BCD codes, each needs 4-bits to represent. The following table gives the relation between these codes. 84-2-1 is a weighted code, i.e., each position will have weight as specified. XS-3 is non-weighted code; the binary code is 3 more than the digit in decimal.

Decimal	84-2-1 $B_3 B_2 B_1 B_0$	XS-3 $X_3 X_2 X_1 X_0$
0	0000	0011
1	0111	0100
2	0110	0101
3	0101	0110
4	0100	0111
5	1011	1000
6	1010	1001
7	1001	1010
8	1000	1011
9	1111	1100

We will consider minterm don't-care combinations as 1, 2, 3, 12, 13, 14. For these combinations 84-2-1 code will not exist and the remaining minterms can be found from truth table.

$$X_0(B_3, B_2, B_1, B_0) = \sum m(0, 4, 6, 8, 10) + \sum \Phi(1, 2, 3, 12, 13, 14) = \overline{B_0}$$

$$X_1(B_3, B_2, B_1, B_0) = \sum m(0, 4, 5, 8, 9, 15) + \sum \Phi(1, 2, 3, 12, 13, 14) = \overline{B_1}$$

$$X_2(B_3, B_2, B_1, B_0) = \sum m(4, 5, 6, 7, 15) + \sum \Phi(1, 2, 3, 12, 13, 14) = B_2$$

$$X_3(B_3, B_2, B_1, B_0) = \sum m(8, 9, 10, 11, 15) + \sum \Phi(1, 2, 3, 12, 13, 14) = B_3$$

DECODER

A binary code of n -bits is capable of representing up to 2^n elements of distinct elements of coded information.

The three inputs are decoded into eight outputs, each representing one of the minterms of the three input variables.

A decoder is a combinational circuit that converts binary information from n input lines to a maximum 2^n unique output lines.

A binary decoder will have n inputs and 2^n outputs.

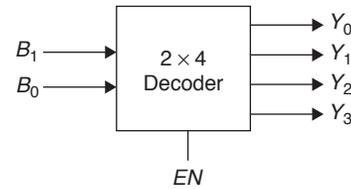
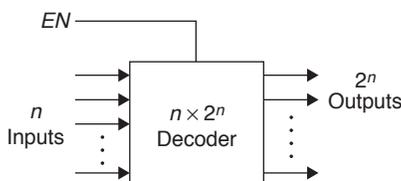


Figure 11 2 x 4 decoder

Table 5 Truth Table

EN	B ₁	B ₀	Y ₃	Y ₂	Y ₁	Y ₀
0	X	X	0	0	0	0
1	0	0	0	0	0	1
1	0	1	0	0	1	0
1	1	0	0	1	0	0
1	1	1	1	0	0	0

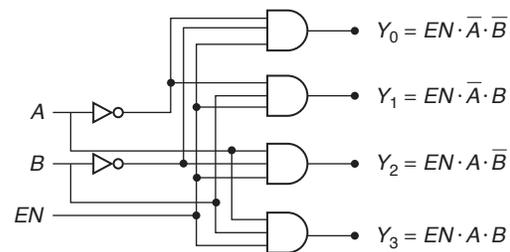


Figure 12 2 x 4 decoder

Decoder outputs are implemented by AND gates, but realization of AND gates at circuit level is done by the NAND gates (universal gates). So, the decoders available in IC form are implemented with NAND gates, i.e., the outputs are in complemented form and outputs are maxterms of the inputs rather than minterms of inputs as in AND gate decoders.

Furthermore, decoders include one or more enable inputs to control the circuit operation. Enable can be either active low/high input.

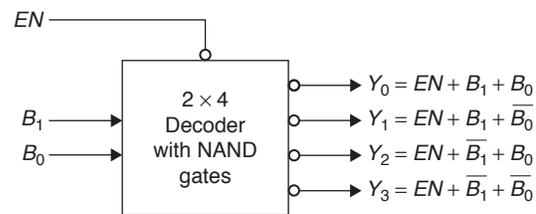


Figure 13 Active low 2 x 4 decoder

Table 6 Truth Table

EN	B ₁	B ₀	Y ₃	Y ₂	Y ₁	Y ₀
1	X	X	1	1	1	1
0	0	0	1	1	1	0
0	0	1	1	1	0	1
0	1	0	1	0	1	1
0	1	1	0	1	1	1

The block diagram shown here is 2×4 decoder with active low output and active low enable input.

The logic diagram is similar to the previous 2×4 decoder, except, all AND gates are replaced by NAND gates and \overline{EN} will have inverter, \overline{EN} is connected to all NAND inputs, as \overline{EN} is active low input for this circuit.

The decoder is enabled when EN is equal to 0.

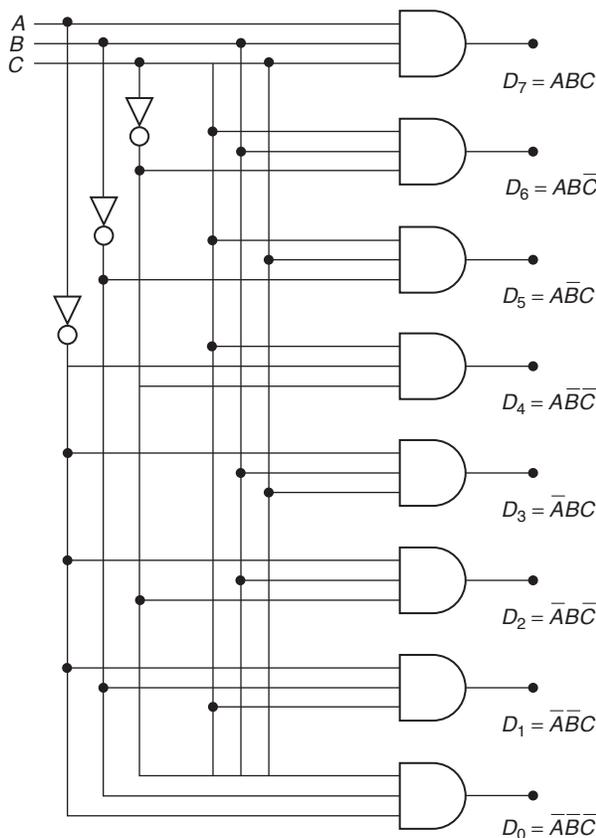
As shown in the truth table, only one output can be equal to 0 at any given time, all other outputs are equal to 1. The output whose value is equal to 0 represents the minterm selected by inputs, enable.

Consider a 3–8 line decoder

Table 7 Truth Table

Inputs			Outputs							
A	B	C	D_0	D_1	D_2	D_3	D_4	D_5	D_6	D_7
0	0	0	1	0	0	0	0	0	0	0
0	0	1	0	1	0	0	0	0	0	0
0	1	0	0	0	1	0	0	0	0	0
0	1	1	0	0	0	1	0	0	0	0
1	0	0	0	0	0	0	1	0	0	0
1	0	1	0	0	0	0	0	1	0	0
1	1	0	0	0	0	0	0	0	1	0
1	1	1	0	0	0	0	0	0	0	1

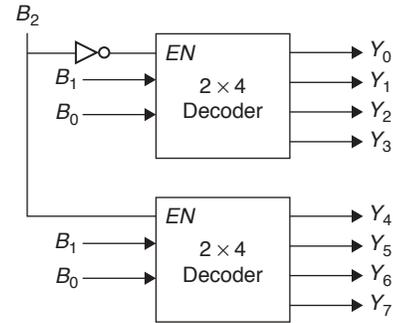
A 3–8 decoder has 3 input lines and 8 output lines, based on the combination of inputs applied for the 3 inputs, one of the 8 output lines will be made logic 1 as shown in the truth table. So, each output will have only one minterm.



Designing High Order Decoders from Lower Order Decoders

Decoder with enable input can be connected together to form larger decoder circuit.

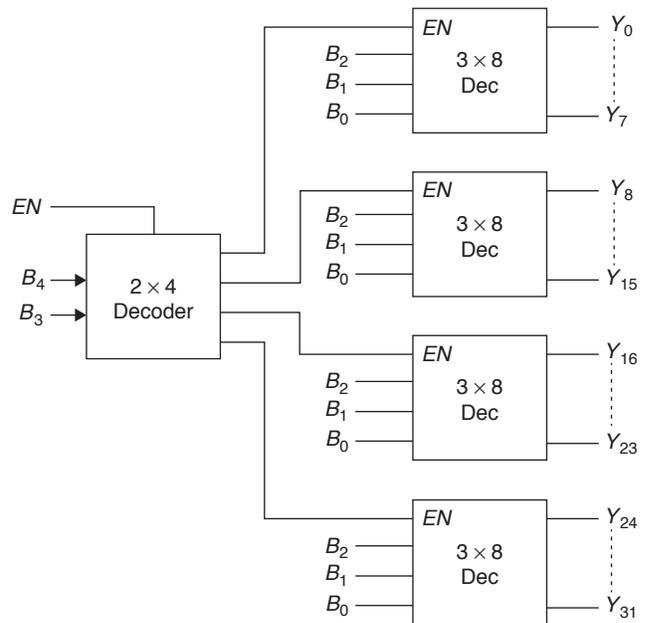
The following configuration shows 3×8 decoder with 2×4 decoders.



When $B_2 = 0$, top decoder is enabled and other is disabled, for 000–011 inputs, outputs are Y_0 – Y_3 , respectively, and other outputs are 0.

For $B_2 = 1$, the enable conditions are reversed.

The bottom decoder outputs generates minterms 100–111, while the outputs of top decoder are all 0's. 5×32 decoder with 3×8 decoders, 2×4 decoders



5×32 decoder will have 5 inputs $B_4 B_3 B_2 B_1 B_0$, 3×8 decoder will have 8 outputs, so 5×32 requires four 3×8 decoders, and we need one of the 2×4 decoders to select one 3×8 decoders and the connections are as shown in the circuit above.

Combinational Logic Implementation

An $n \times 2^n$ decoder provides 2^n minterms of n input variables. Since any Boolean function can be expressed in sum-of-minterms form, a decoder that generates the minterms of

the function, together with an external OR gate that forms their logical sum, provides a hardware implementation of the function.

Similarly, any function with n inputs and m outputs can be implemented with $n \times 2^n$ decoders and m OR gates.

Example 3: Implement full adder circuit by using 2×4 decoder.

$$\text{Sum} = \Sigma (1, 2, 4, 7), \text{Carry} = \Sigma (3, 5, 6, 7)$$

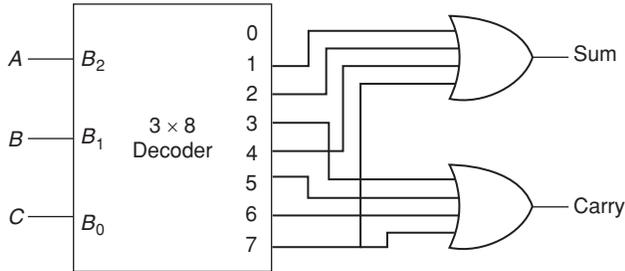
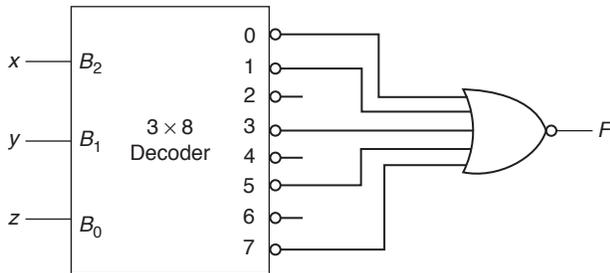


Figure 14 Implementation of full adder circuit with decoder

The 3×8 decoder generates the 8 minterms for A , B , and C . The OR gate for output sum forms the logical sum of minterms 1, 2, 4 and 7. The OR gate for output carry forms the logical sum of minterms 3, 5, 6 and 7.

Example 4: The minimized SOP form of output $F(x, y, z)$ is

- (A) $x'y + z'$ (B) $x'y' + z'$
 (C) $x'y' + z'$ (D) $x' + y'z$



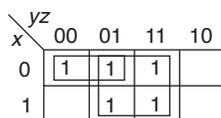
Solution: (C)

The outputs of decoder are in active low state. So, we can express outputs as $\overline{Y_7}, \overline{Y_6} \dots \overline{Y_0}$

Outputs 0, 1, 3, 5, 7 are connected to NAND gate to form function $F(x, y, z)$

$$\begin{aligned} \text{So } F &= \overline{Y_0} \cdot \overline{Y_1} \cdot \overline{Y_3} \cdot \overline{Y_5} \cdot \overline{Y_7} \\ &= Y_0 + Y_1 + Y_3 + Y_5 + Y_7 \\ &= \Sigma(0, 1, 3, 5, 7) \end{aligned}$$

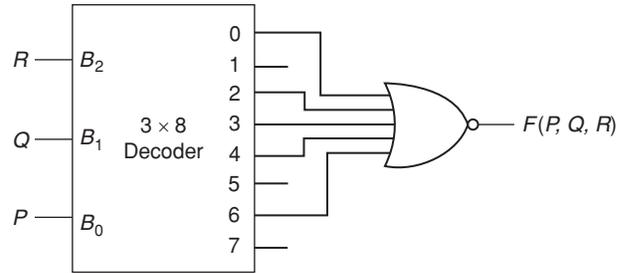
By using K-maps



$$F = z + x'y'$$

Example 5. The minimal POS form of output function $f(P, Q, R)$ is

- (A) $P\overline{Q} + PR$ (B) $P + \overline{Q}R$
 (C) $P(\overline{Q} + R)$ (D) $Q(\overline{P} + R)$



Solution: (C)

The outputs of decoder are in normal form. 0, 2, 3, 4, 6 outputs are connected to NOR gate to form $F(P, Q, R)$

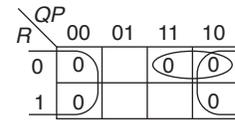
$$\begin{aligned} \text{So } F &= \overline{Y_0 + Y_2 + Y_3 + Y_4 + Y_6} \\ &= \overline{Y_0} \cdot \overline{Y_2} \cdot \overline{Y_3} \cdot \overline{Y_4} \cdot \overline{Y_6} \end{aligned}$$

Y_0, Y_1, \dots, Y_7 indicate minterms, whereas $\overline{Y_0}, \overline{Y_1}, \dots, \overline{Y_7}$ are maxterms.

So $F = \pi(0, 2, 3, 4, 6)$

Here, from the decoder circuit MSB is R , LSB is P .

By using K-map



$$F(P, Q, R) = P(R + \overline{Q})$$

ENCODERS

It is a digital circuit that performs the inverse operation of a decoder.

An encoder has 2^n (or fewer) input lines and n output lines.

It is also known as an octal to binary converter.

Consider an 8–3 line encoder:

Table 8 Truth Table

Inputs								Outputs		
D_0	D_1	D_2	D_3	D_4	D_5	D_6	D_7	A	B	C
1	0	0	0	0	0	0	0	0	0	0
0	1	0	0	0	0	0	0	0	0	1
0	0	1	0	0	0	0	0	0	1	0
0	0	0	1	0	0	1	0	0	1	1
0	0	0	0	1	0	0	1	1	0	0
0	0	0	0	0	1	0	0	1	0	1
0	0	0	0	0	0	1	0	1	1	0
0	0	0	0	0	0	0	1	1	1	1

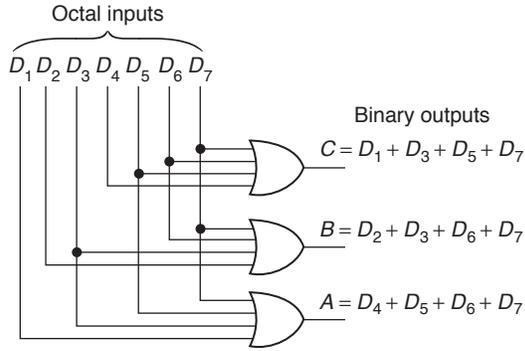


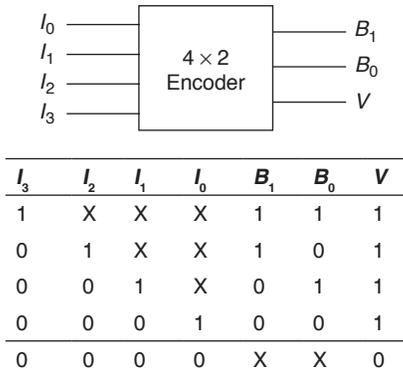
Figure 15 Logic diagram

Priority Encoder

A priority encoder is an encoder circuit that includes the priority function.

When two or more inputs are present, the input with higher priority will be considered.

Consider the 4×2 priority encoder.



I_3-I_0 are inputs and $B_1 B_0$ are binary output bits, valid (V) output is set to 1, when at least one input is present at input (I_3-I_0).

When there is no input present, ($I_3-I_0 = 0000$) then $V = 0$, for this combination the output $B_1 B_0$ will not be considered.

The higher the subscript number, the higher the priority of the input. Input I_3 has the highest priority, I_2 has the next priority level. Input I_0 has lowest priority level. The Boolean expressions for output $B_1 B_0$ are

$$\begin{aligned}
 B_1 &= I_3 + \overline{I_3} I_2 \\
 &= I_3 + I_2 \\
 B_0 &= I_3 + \overline{I_3} \overline{I_2} I_1 \\
 &= I_3 + \overline{I_2} I_1 \\
 V &= I_3 + I_2 + I_1 + I_0
 \end{aligned}$$

MULTIPLEXER

A multiplexer (MUX) is a device that allows digital information from several sources to be converted on to a single line for transmission over that line to a common destination.

The MUX has several data input lines and a single output line. It also has data select inputs that permits digital data on any one of the inputs to be switched to the output line.

Depending upon the binary code applied at the selection inputs, one (out of 2^n) input will be gated to single output. It is one of the most widely used standard logic circuits in digital design. The applications of multiplexer include data selection, data routing, operation sequencing, parallel to serial conversion, and logic function generation.

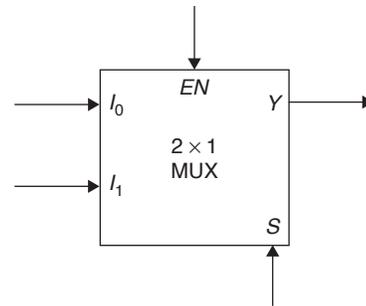
2^n inputs will be controlled by n selection lines and multiplexer will have 1 output, we denote it as $2^n \times 1$ multiplexer (data selector).

In other words, a multiplexer selects 1 out of n input data sources and transmits the selected data to a single output channel, this is called as multiplexing.

Basic 2×1 Multiplexer

The figure shows 2×1 multiplexer block diagram; it will have 2 inputs— I_0 and I_1 , one selection line S , and one output Y . The function table is as shown here.

EN	S	Y
0	x	0
1	0	I_0
1	1	I_1

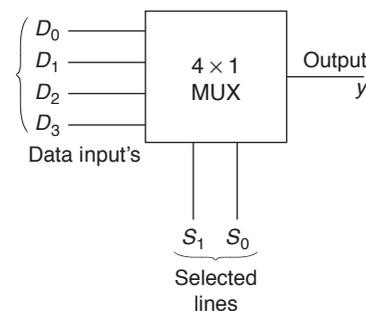


The output equation of 2×1 multiplexer is $Y = EN(I_0 \overline{S} + I_1 S)$.

When enable is 1, the multiplexer will work in normal mode, else the multiplexer will be disabled.

Sometimes enable input will be active low enable \overline{EN} , then $Y = \overline{EN}(I_0 \overline{S} + I_1 S)$.

The 4×1 Multiplexer



If a binary zero $S_1 = 0$ and $S_0 = 0$ as applied to the data select line the data input D_0 appear on the data output line and so on.

S_1	S_0	y
0	0	D_0
0	1	D_1
1	0	D_2
1	1	D_3

$$y = \bar{S}_1 \bar{S}_0 D_0 + \bar{S}_1 S_0 D_1 + S_1 \bar{S}_0 D_2 + S_1 S_0 D_3$$

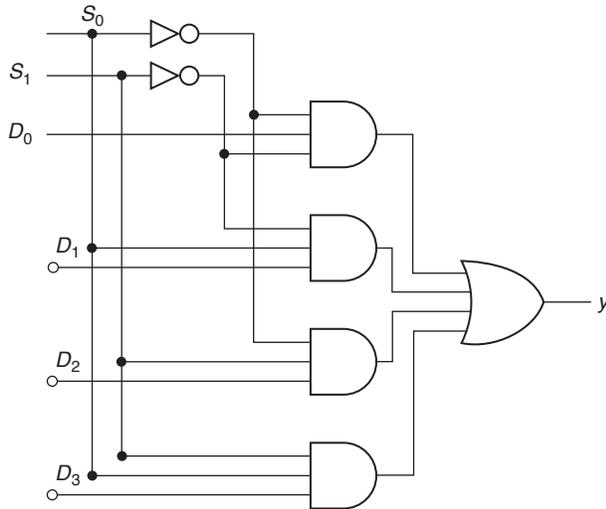


Figure 16 Logic diagram

For 8×1 multiplexer with 8 inputs from I_0-I_7 based on selection inputs S_2, S_1, S_0 , the equation for output

$$Y = I_0 \bar{S}_2 \bar{S}_1 \bar{S}_0 + I_1 \bar{S}_2 \bar{S}_1 S_0 + I_2 \bar{S}_2 S_1 \bar{S}_0 + I_3 \bar{S}_2 S_1 S_0 + I_4 S_2 \bar{S}_1 \bar{S}_0 + I_5 S_2 \bar{S}_1 S_0 + I_6 S_2 S_1 \bar{S}_0 + I_7 S_2 S_1 S_0$$

From multiplexer equation, we can observe, each input is associated with its minterm (in terms of selection inputs).

Basic Gates by Using MUX

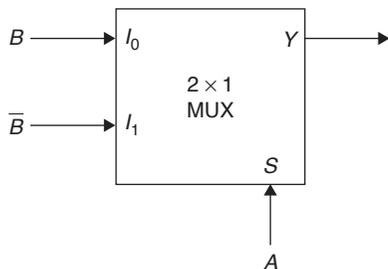


Figure 17 X-OR gate by using 2×1 MUX

$Y = \bar{A}B + A\bar{B} = X\text{-OR gate}$, we can interchange inputs A and B also,

By interchanging inputs I_0 and I_1 , $Y = \bar{A}\bar{B} + AB$, X-NOR gate.

Similarly, we can build all basic gates by using 2×1 multiplexer.

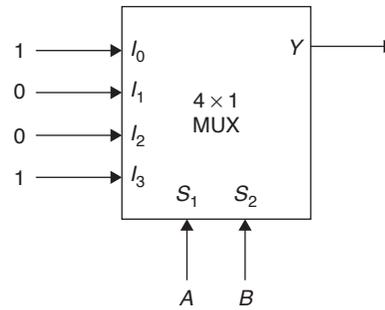
Example 6: If $I_0 = 1, I_1 = 0, S = A$, then Y is

Solution: $Y = (I_0 \bar{S} + I_1 S) = \bar{A}$. It Implements NOT gate.

Example 7: What should be the connections to implement NAND gate by using 2×1 MUX?

Solution: $Y = \bar{A}B = \bar{A} + \bar{B} = \bar{A} + \bar{A}B = 1 \cdot \bar{A} + \bar{B} \cdot A$

By considering $I_0 = 1, I_1 = \bar{B}, S = A$, we can implement NAND gate, or by interchanging A and B also we can get the same answer.



For the above 4×1 multiplexer $Y = \bar{A}B + AB = X\text{-NOR gate}$, similarly to implement 2 input gates by using 4×1 multiplexer, the inputs I_0, I_1, I_2, I_3 should be same as the terms in the truth table of that gate.

Logic Function Implementation by Using Multiplexer

Let us consider a full subtractor circuit (borrow) to be implemented by using multiplexer.

Full subtractor borrow (B) is a function of 3 inputs X, Y, Z . The truth table is

X	Y	Z	B	4×1 MUX	2×1 MUX
0	0	0	0	$B = Z$	$B = Y + Z$
0	0	1	1		
0	1	0	1	$B = 1$	
0	1	1	1		
1	0	0	0	$B = 0$	$B = YZ$
1	0	1	0		
1	1	0	0	$B = Z$	
1	1	1	1		

To implement borrow by using 8×1 multiplexer, connect the three variables X, Y, Z directly to selection lines of the multiplexer, and connect the corresponding values of B to inputs, i.e., for $I_0 = 0, I_1 = 1, I_2 = 1$, etc. as per above truth table.

To implement borrow by using 4×1 multiplexer, connect any two variables to selection lines (in this case X, Y) and write output (B) in terms of other variable, for $XY = 00$, output B is same as Z , so connect $I_0 = Z$, similarly $1, 0, Z$ for remaining inputs.

To implement the function by using 2×1 multiplexer, connect 1 variable as selection line (in this case consider X) and write output (B) in terms of other variables, for $X = 0$,

output B varies as $B = Y + Z$, so connect $I_0 = Y + Z$. For $X = 1$, output B varies as $B = YZ$, connect $I_1 = YZ$.

N -variable function can be implemented by using $2^{N-1} \times 1$ multiplexer without any extra hardware.

Implementation of Higher Order Multiplexer by Using Lower Order Multiplexers

By using lower order multiplexers, we can implement higher order multiplexers, for example by using 4×1 multiplexer, we can implement 8×1 MUX or 16×1 MUX or other higher order multiplexers.

Let us consider implementation of 16×1 MUX by using 4×1 MUX. 16×1 MUX will have inputs I_0-I_{15} and selection lines S_0-S_3 , whereas 4×1 MUX will have only 4 input lines, and 2 selection lines, so we require four 4×1 MUX to consider all inputs I_0-I_{15} , and again to select one of the four outputs of these four multiplexers one more 4×1 multiplexer is needed (for which we will connect higher order selection lines S_2 and S_3). So, total of 5, 4×1 multiplexers are required to implement 16×1 MUX.

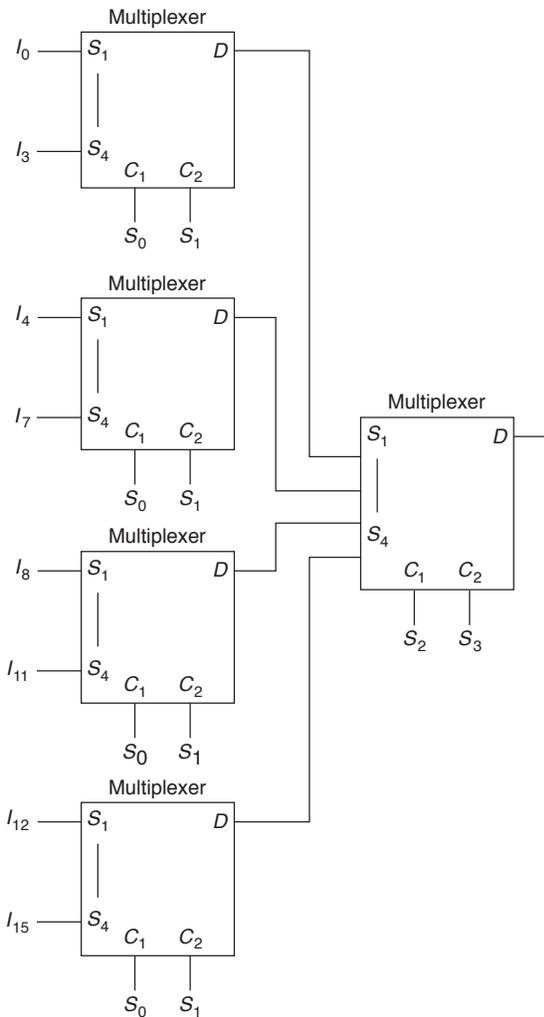


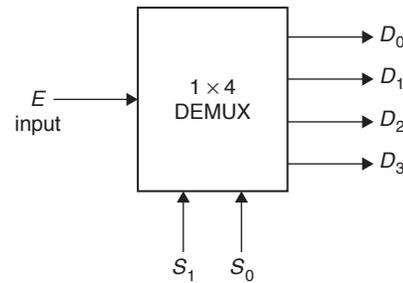
Figure 18 Realization of 16×1 multiplexer by using 4×1 multiplexers

In a similar fashion, to design 4×1 MUX, we require $3, 2 \times 1$ multiplexers, and to design 8×1 multiplexer, we require $7, 2 \times 1$ multiplexers.

DEMULTIPLEXER

The demultiplexer [DeMUX] basically serves opposite of the multiplexing function. It takes data from one line and distributes them to a given number of output lines.

The other name for demultiplexer is data distributor, as it receives information on a single line and distributes it to a possible 2^n output lines, where n is the number of selection lines, and value of n selects the line.



S_1	S_0	D_3	D_2	D_1	D_0
0	0	0	0	0	E
0	1	0	0	E	0
1	0	0	E	0	0
1	1	E	0	0	0

When $S_1 S_0 = 10$; D_2 will be same as input E , and other outputs will be maintained at zero (0).

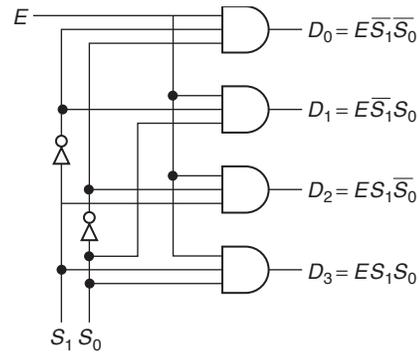
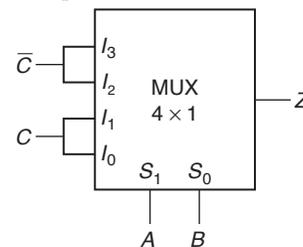


Figure 17 Logic diagram

Solved Examples

Example 1: The multiplexer shown in the figure is a $4 : 1$ multiplexer. The output z is

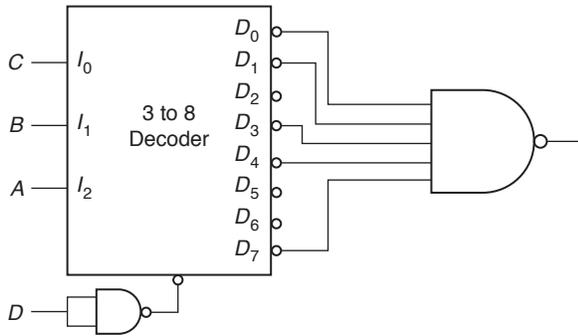


Solution:

A_1	B_0	Z
0	0	C
0	1	C
1	0	\bar{C}
1	1	\bar{C}

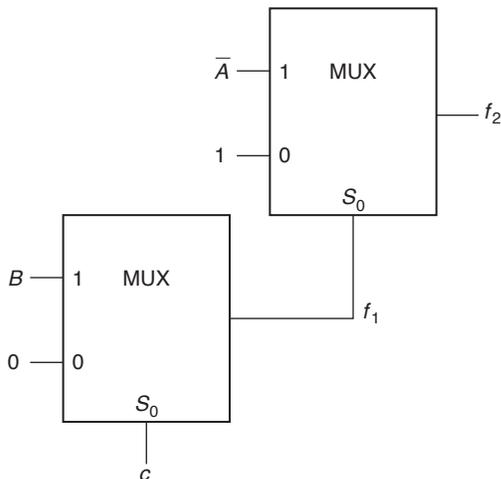
$$\begin{aligned} \therefore Z &= \bar{A}\bar{B}C + \bar{A}BC + A\bar{B}\bar{C} + AB\bar{C} \\ &= \bar{B}(\bar{A}C + A\bar{C}) + B(\bar{A}\bar{C} + A\bar{C}) \\ &= (\bar{A}C + A\bar{C})(\bar{B} + B)(x + \bar{x} = 1) \\ \therefore &= \bar{A}C + A\bar{C} = A \oplus C \end{aligned}$$

Example 2: The logic circuit shown in figure implements



Solution: $z = D(\bar{A}\bar{B}\bar{C} + \bar{A}\bar{B}C + \bar{A}BC + A\bar{B}\bar{C} + ABC)$
 $= D(\bar{A}\bar{B}(C + \bar{C}) + BC(\bar{A} + A) + A\bar{B}\bar{C})$
 $\times D(\bar{B}\bar{A} + \bar{B}\bar{C} + BC)$
 $= D(B \oplus C + \bar{A}\bar{B})$

Example 3. The network shown in figure implements

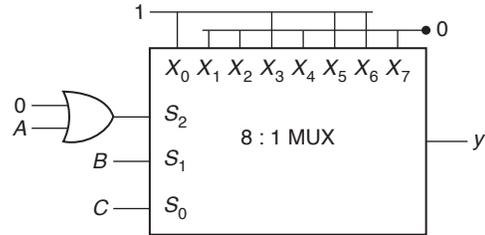


Solution: $f_1 = \bar{C}0 + CB = CB, f_1 = CB$
 $F_2 = \bar{f}_1 + f_1\bar{A} = \bar{A} \cdot CB + \bar{C}B$

$$\begin{aligned} &= \bar{A} + \bar{C}\bar{B} \\ &= \bar{A} + \bar{C} + \bar{B} = \overline{ABC} \end{aligned}$$

\therefore NAND Gate

Example 4: In the TTL circuit in figure, S_2-S_0 are select lines and x_7-x_0 are input lines. S_0 and X_0 are LSBs. The output Y is

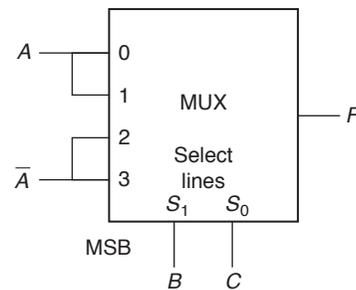


Solution: $S_2 = A, S_1 = B, S_0 = C$

$S_2(A)$	$S_1(B)$	$S_0(C)$	Y
0	0	0	1
0	0	1	0
0	1	0	0
0	1	1	1
1	0	0	0
1	0	1	1
1	1	0	1
1	1	1	0

$$\begin{aligned} Y &= \bar{A}\bar{B}\bar{C} + \bar{A}B\bar{C} + A\bar{B}\bar{C} + \bar{A}BC \\ &= \bar{C}(\bar{A}\bar{B} + AB) + C(\bar{A}\bar{B} + \bar{A}B) \\ Y &= \bar{C}(\bar{A} \oplus B) + C(A \oplus B) = A \oplus B \oplus C \end{aligned}$$

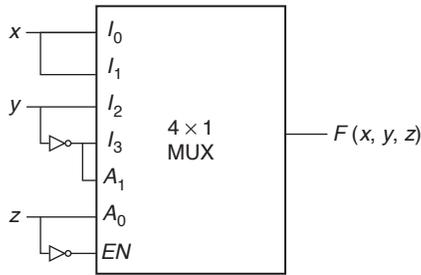
Example 5: The logic realized by the adjoining circuit is



Solution: $F = \bar{B}\bar{C}A + \bar{B}CA + \bar{B}\bar{C}\bar{A} + \bar{B}C\bar{A}$
 $\times \bar{C}(\bar{B}A + B\bar{A}) + C(\bar{B}A + B\bar{A})$
 $\times \bar{A}\bar{B} + \bar{A}B(C + \bar{C}) = A \oplus B$

Example 6: Consider the following multiplexer, where I_0, I_1, I_2, I_3 are four data input lines selected by two address line combinations $A_1A_0 = 00, 01, 10, 11$, respectively and f

is the output of the multiplexer. EN is the enable input, the function $f(x, y, z)$ implemented by the below circuit is



Solution: $A_1 = \bar{y} \cdot A_0 = z, EN = \bar{z}$

A_1	A_0	S	I
0	0	$(y\bar{z})$	x
0	1	(yz)	x
1	0	$(\bar{y}\bar{z})$	y
1	1	$(\bar{y}z)$	\bar{y}

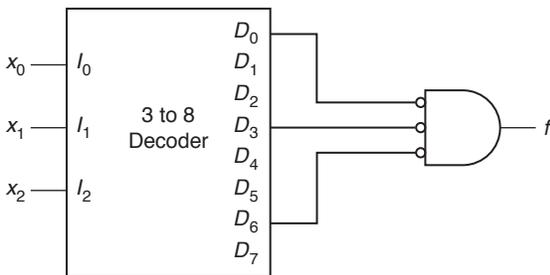
$$f(x, y, z) = S.I = (xy + 0 + \bar{y}z) \cdot EN = xy \cdot \bar{z}$$

EXERCISES

Practice Problems I

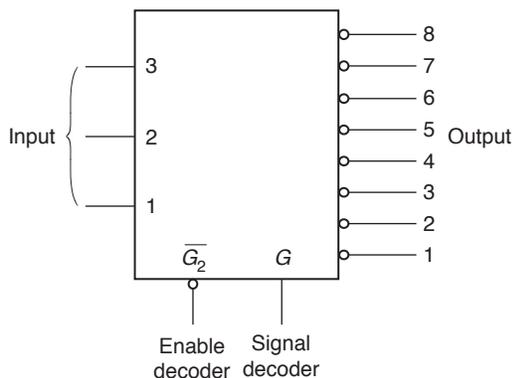
Directions for questions 1 to 21: Select the correct alternative from the given choices.

- The binary number 110011 is to be converted to gray code. The number of gates and type required are
 (A) 6, AND (B) 6, X-NOR
 (C) 6, X-OR (D) 5, X-OR
- The number of 4-to 16-line decoder required to make an 8- to 256-line decoder is
 (A) 16 (B) 17
 (C) 32 (D) 64
- $f(x_2, x_1, x_0) = ?$



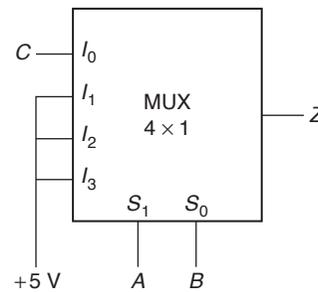
- (A) $\pi(1, 2, 4, 5, 7)$ (B) $\Sigma(1, 2, 4, 5, 7)$
 (C) $\Sigma(0, 3, 6)$ (D) $\pi(0, 2, 3, 6)$

4. A 3-to-8 decoder is shown below



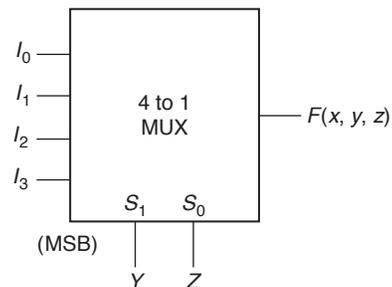
All the output lines of the chip will be high except pin 8, when all the inputs 1, 2, and 3

- (A) are high; and G, G_2 are low
 (B) are high; and G is low G_2 is high
 (C) are high; and G, G_2 are high
 (D) are high; and G is high G_2 is low
5. The MUX shown in figure is 4×1 multiplexer the output z is



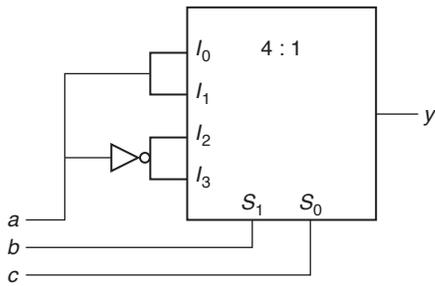
- (A) $A B C$
 (B) $A \oplus B \oplus C$
 (C) $A \ominus B \ominus C$
 (D) $A + B + C$

6. If a 4 to 1 MUX (shown below) realizes a three variable function $f(x, y, z) = xy + x\bar{z}$ then which of the following is correct?



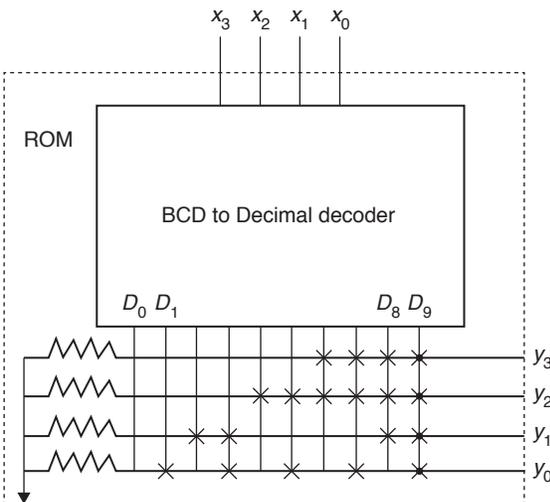
- (A) $I_0 = X, I_1 = 0, I_2 = X, I_3 = X$
 (B) $I_0 = 0, I_1 = 1, I_2 = Y, I_3 = X$
 (C) $I_0 = X, I_1 = 1, I_2 = 0, I_3 = X$
 (D) $I_0 = X, I_1 = 0, I_2 = X, I_3 = Z$

7. The circuit shown in the figure is same as



- (A) two input NAND gate with a and c inputs
- (B) two input NOR gate with a and c inputs
- (C) two input X-OR gates with a and b inputs
- (D) two input X-NOR gate with b and c inputs

8. If the input x_3, x_2, x_1, x_0 to the ROM in the figure are 8421 BCD numbers, then the outputs y_3, y_2, y_1, y_0 are

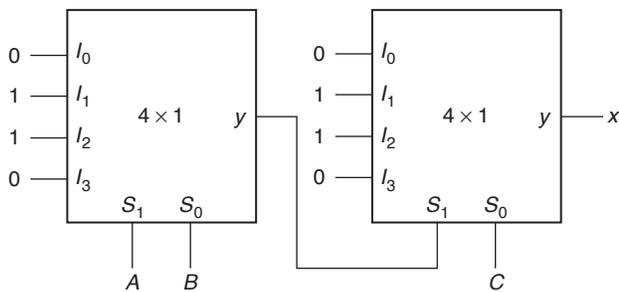


- (A) gray code numbers
- (B) 2421 BCD
- (C) Excess - 3 code numbers
- (D) 84-2-1

9. A 4-bit parallel full adder without input carry requires

- (A) 8 HA, 4 OR gates
- (B) 8 HA, 3 OR gates
- (C) 7 HA, 4 OR gates
- (D) 7 HA, 3 OR gates

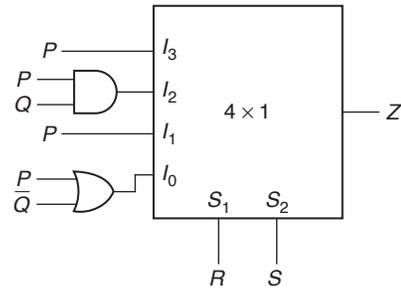
10. In the circuit find X .



- (A) $\bar{A}\bar{B}\bar{C} + \bar{A}B\bar{C} + \bar{A}\bar{B}C + ABC$
- (B) $\bar{A}BC + \bar{A}\bar{B}C + A\bar{B}\bar{C} + \bar{A}\bar{B}\bar{C}$

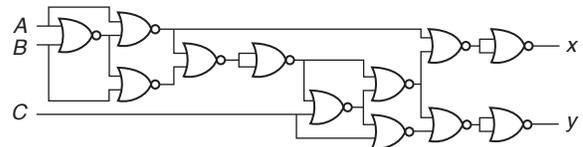
- (C) $AB + BC + AC$
- (D) $\bar{A}\bar{B} + \bar{B}\bar{C} + \bar{A}\bar{C}$

11. Find the function implemented.



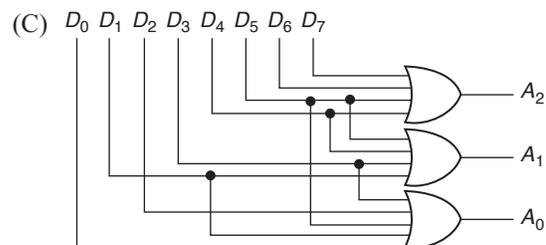
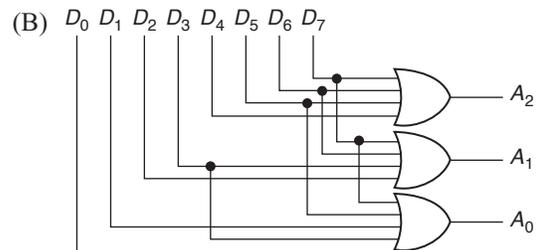
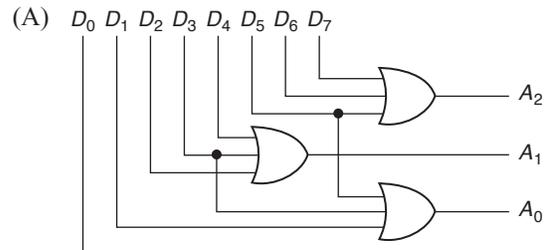
- (A) $PQ + PS + \bar{Q}\bar{R}\bar{S}$
- (B) $P\bar{Q} + PQ\bar{R} + \bar{P}\bar{Q}\bar{S}$
- (C) $P\bar{Q}\bar{R} + \bar{P}QR + PQRS + \bar{Q}\bar{R}\bar{S}$
- (D) $PQ\bar{R} + PQR\bar{S} + P\bar{Q}\bar{R}S + \bar{Q}\bar{R}\bar{S}$

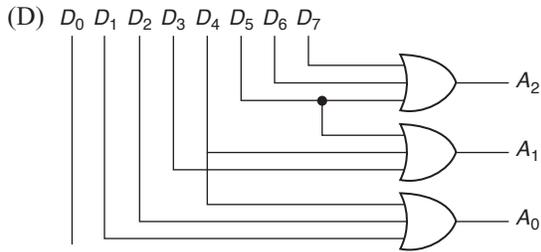
12. Which function is represented by the given circuit?



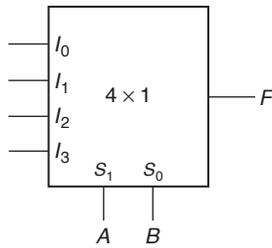
- (A) Full adder
- (B) Full subtractor
- (C) Comparator
- (D) Parity generator

13. Which of the following represents octal to binary encoder?



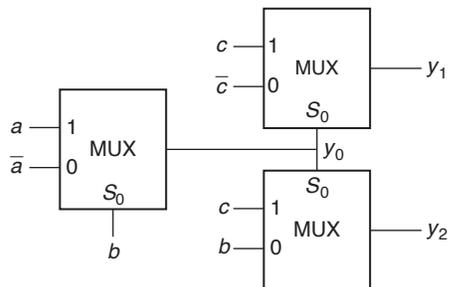


14. For a MUX to function as a full adder what should be the input provided to the I_0, I_1, I_2, I_3 if the A and B are the select lines?



- (A) $I_0 = I_1 = C_{in}; I_2 = I_3 = \overline{C_{in}}$
- (B) $I_0 = I_1 = \overline{C_{in}}; I_2 = I_3 = C_{in}$
- (C) $I_0 = I_3 = C_{in}; I_1 = I_2 = \overline{C_{in}}$
- (D) $I_0 = I_3 = \overline{C_{in}}; I_1 = I_2 = C_{in}$

15. The given circuit act as

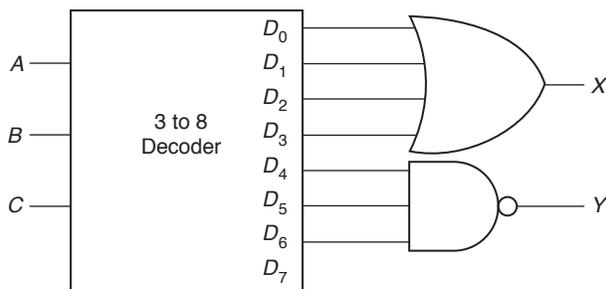


- (A) Full adder
- (B) Half adder
- (C) Full subtractor
- (D) Half subtractor

16. For a 4×16 decoder circuit, the outputs of decoder ($y_0, y_1, y_4, y_5, y_{10}, y_{11}, y_{14}, y_{15}$) are connected to 8 input NOR gate, the expression of NOR gate output is

- (A) $A \oplus D$
- (B) $A \odot D$
- (C) $A \odot C$
- (D) $A \oplus C$

17. The function implemented by decoder is

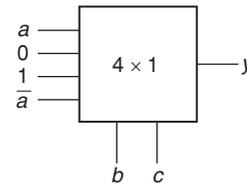


- (A) $X = A'BC' + B'C', y = A + B$
- (B) $X = A'C' + B'C', y = 1$
- (C) $X = \overline{A}, y = 0$
- (D) $X = \overline{A}, y = 1$

18. A relay is to operate with conditions that it should be on when the input combinations are 0000, 0010, 0101, and 0111. The states 1000, 1001, 1010 don't occur. For rest of the status, relay should be off. The minimized Boolean expression notifying the relationship is

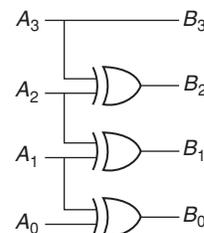
- (A) $BC + ACD$
- (B) $\overline{B}\overline{D} + \overline{A}BD$
- (C) $BD + AC$
- (D) $AB + CD$

19. If a function has been implemented using MUX as shown, implement the same function with a and c as the select lines



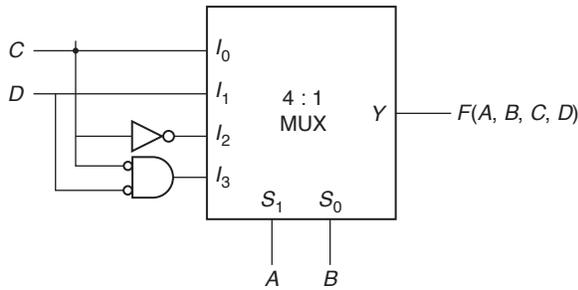
- (A)
- (B)
- (C)
- (D)

20. The circuit is used to convert one code to another. Identify it.



- (A) Binary to gray
- (B) Gray to binary
- (C) Gray to XS-3
- (D) Gray to 8421

21. The Boolean function realised by logic circuit is

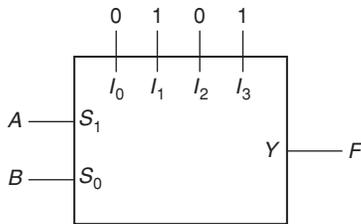


- (A) $F = \Sigma m(0, 1, 3, 5, 9, 10, 14)$
- (B) $F = \Sigma m(2, 3, 5, 7, 8, 12, 13)$
- (C) $F = \Sigma m(1, 2, 4, 5, 11, 14, 15)$
- (D) $F = \Sigma m(2, 3, 5, 7, 8, 9, 12)$

Practice Problems 2

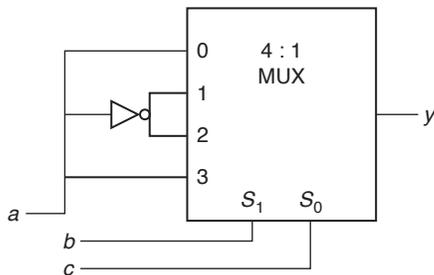
Directions for questions 1 to 21: Select the correct alternative from the given choices.

1. For a binary half subtractor having two input A and B , the correct set of logical expression for the outputs D (A minus B) and X (borrow) are
 - (A) $D = AB + \overline{AB}, X = \overline{AB}$
 - (B) $D = \overline{AB} + \overline{AB}, X = \overline{AB}$
 - (C) $D = \overline{AB} + \overline{AB}, X = \overline{AB}$
 - (D) $D = AB + \overline{AB}, X = \overline{AB}$
2. The function 'F' implemented by the multiplexer chip shown in the figure is



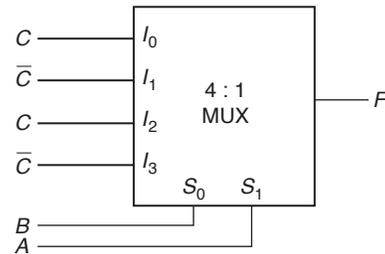
- (A) A
- (B) B
- (C) \overline{AB}
- (D) $\overline{AB} + \overline{AB}$

3 The following multiplexer circuit is equal to



- (A) implementation of sum equation of full adder
- (B) implementation of carry equation of full adder
- (C) implementation of borrow equation of full subtractor
- (D) all of the above

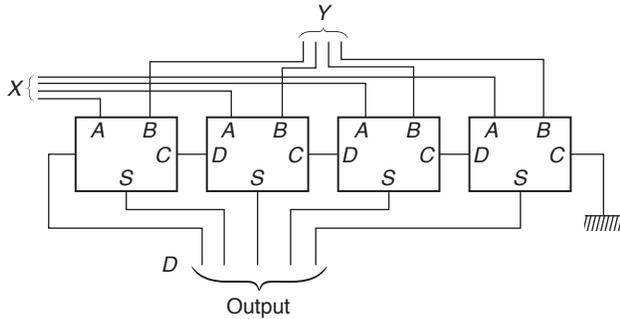
4 The output 'F' of the multiplexer circuit shown in the figure will be



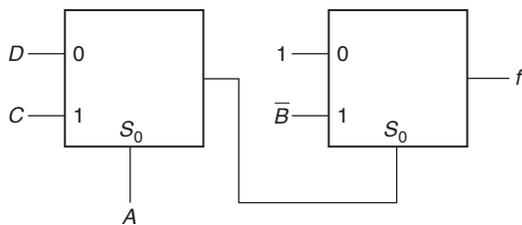
- (A) $AB + B\overline{C} + \overline{C}A + \overline{B}\overline{C}$
- (B) $A \oplus B \oplus C$
- (C) $A \oplus B$
- (D) $B \oplus C$

5. Full subtractor can be implemented by using
 - (A) 3-to-8 line decoder only
 - (B) 3-to-8 line decoder and one OR gate
 - (C) 3-to-8 line decoder and two OR gates
 - (D) None
6. What are the difference and borrow equations for the above circuit?
 - (A) $D = x \ominus y \ominus z, B = x'y + yz + zx'$
 - (B) $D = X \oplus y \oplus z, B = xy + yz + zx$
 - (C) $D = x \oplus y \oplus z, B = x'y + yz + zx'$
 - (D) A and C both
7. Combinational circuits are one in which output depends _____, whereas sequential circuit's output depends _____.
 - (A) only on present input, only on past input
 - (B) only on present input, only on past and future input
 - (C) only on present input, only on present input and past output
 - (D) on present input, on past and present output
8. The sum output of the half adder is given by (assume A and B as inputs)
 - (A) $S = AB(\overline{A+B})$
 - (B) $S = (A+B)\overline{AB}$
 - (C) $S = (A+B)(\overline{AB})$
 - (D) $S = (\overline{A+B})(\overline{AB})$
9. MUX implements which of the following logic?
 - (A) NAND-XOR
 - (B) AND-OR
 - (C) OR-AND
 - (D) XOR-NOT
10. A DeMUX can be used as a
 - (A) Comparator
 - (B) Encoder
 - (C) Decoder
 - (D) Adder

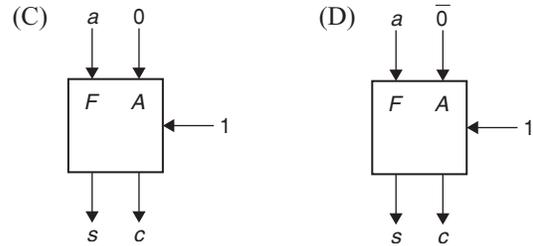
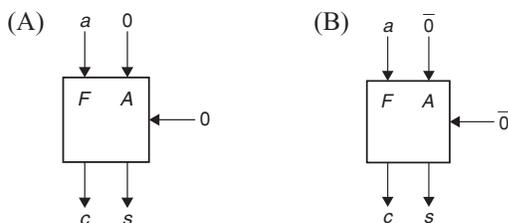
11. If we have inputs as A, B and C and output as S and D . We are given that $S = A \oplus B \oplus C$. $D = BC + \bar{A}B + \bar{A}C$. Which of the circuit is represented by it?



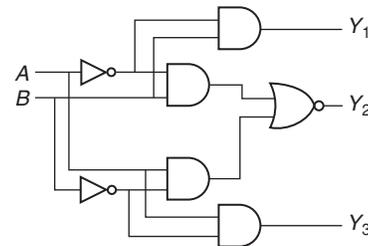
- (A) 4-bit adder giving $X + Y$
 (B) 4-bit subtractor giving $X - Y$
 (C) 4-bit subtractor giving $Y - X$
 (D) 4-bit adder giving $X + Y + S$
12. The Boolean function f implemented in the figure using two input multiplexers is



- (A) $A\bar{C} + \bar{A}\bar{D} + \bar{D}C + \bar{A}\bar{B}D + \bar{A}\bar{B}C$
 (B) $\bar{A} + A\bar{C} + \bar{A}\bar{D} + \bar{D}\bar{C}$
 (C) $\bar{B} + A\bar{C} + \bar{A}\bar{D} + \bar{D}\bar{C}$
 (D) $AC + AD + A + B$
13. The carry generate and carry propagate function of the look ahead carry adder is
- (A) $CG = A + B, CP = A \oplus B$
 (B) $CG = A \oplus B, CP = A + B$
 (C) $CG = AB, CP = A \oplus B$
 (D) $CG = AB, CP = A + B$
14. If we have a comparator and if E represents the condition for equality i.e., $(A_n \oplus B_n)$, if A_n and B_n are to be compared then the expression $A_3\bar{B}_3 + E_3A_2\bar{B}_2 + E_3E_2A_1\bar{B}_1 + E_3E_2E_1A\bar{B}$. represents which of the condition for a 4-bit number?
- (A) $A > B$ (B) $B > A$
 (C) $A = B$ (D) None of these
15. When full adder is used to function as a 1-bit incrementer, which of the circuit configurations must be used?



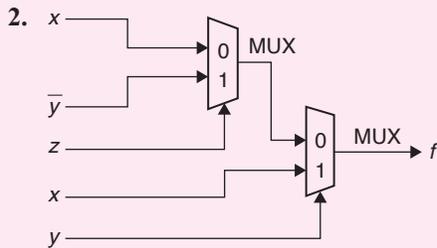
16. Identify the circuit.



- (A) Half adder
 (B) Full adder
 (C) 1-bit magnitude comparator
 (D) Parity generator
17. In order to implement n variable function (without any extra hardware) the minimum order of MUX is
- (A) $2^n \times 1$ (B) $2^n \times 1$
 (C) $(2^n - 1) \times 1$ (D) $(2^{n-1}) \times 1$
18. A full adder circuit can be changed to full subtractor by adding a
- (A) NOR gate (B) NAND gate
 (C) Inverter (D) AND gate
19. The half adder when implemented in terms of NAND logic is expressed as
- (A) $A \oplus B$ (B) $\overline{\overline{A \cdot AB \cdot B \cdot AB}}$
 (C) $\overline{\overline{\overline{A \cdot AB \cdot B \cdot AB}}}$ (D) $\overline{\overline{\overline{A \cdot AB \cdot B \cdot AB}}}$
20. For a DeMUX to act as a decoder, what is the required condition?
- (A) Input should be left unconnected and select lines behave as a input to decoder
 (B) Input should be always 0 and select line behave as inputs to decoder
 (C) Both are same
 (D) Input should become enable and select lines behave as inputs to decoder
21. For a full subtractor, which of the combination will give the difference?
- (A) $\overline{\overline{\overline{(A \oplus B)(A \oplus B)b_i \cdot b_i(A \oplus B)b_i}}}$
 (B) $\overline{\overline{\overline{B \cdot AB \cdot b_i(A \oplus B)}}}$
 (C) $\overline{\overline{\overline{A + B + b_i + A \oplus B}}}$
 (D) None of these

PREVIOUS YEARS' QUESTIONS

1. A 4-bit carry look ahead adder, which adds two 4-bit numbers, is designed using AND, OR, NOT, NAND, NOR gates only. Assuming that all the inputs are available in both complemented and uncomplemented forms and the delay of each gate is one time unit, what is the overall propagation delay of the adder? Assume that the carry network has been implemented using two-level AND-OR logic. **[2004]**
- (A) 4 time units (B) 6 time units
(C) 10 time units (D) 12 time units



Consider the circuit above. Which one of the following options correctly represents $f(x, y, z)$? **[2006]**

- (A) $x\bar{z} + xy + \bar{y}z$ (B) $x\bar{z} + xy + \bar{y}\bar{z}$
(C) $xz + xy + \bar{y}\bar{z}$ (D) $xz + x\bar{y} + \bar{y}\bar{z}$
3. Given two 3-bit numbers $a_2a_1a_0$ and $b_2b_1b_0$ and c , the carry in, the function that represents the carry generate function when these two numbers are added is **[2006]**
- (A) $a_2b_2 + a_2a_1b_1 + a_2a_1a_0b_0 + a_2a_0b_1b_0 + a_1b_2b_1 + a_1a_0b_2b_0 + a_0b_2b_1b_0$
(B) $a_2b_2 + a_2b_1b_0 + a_2a_1b_1b_0 + a_1a_0b_2b_1 + a_1a_0b_2 + a_1a_0b_2b_0 + a_2a_0b_1b_0$
(C) $a_2 + b_2 + (a_2 \oplus b_2)(a_1 + b_1 + (a_1 \oplus b_1)(a_0 + b_0))$
(D) $a_2b_2 + \bar{a}_2a_1b_1 + \bar{a}_2a_1a_0b_0 + \bar{a}_2a_0\bar{b}_1b_0 + a_1\bar{b}_2b_1 + \bar{a}_1a_0\bar{b}_2b_0 + a_0\bar{b}_2\bar{b}_1b_0$
4. We consider the addition of two 2's complement numbers $b_{n-1}b_{n-2} \dots b_0$ and $a_{n-1}a_{n-2} \dots a_0$. A binary adder for adding unsigned binary numbers is used to add the two numbers. The sum is denoted by $c_{n-1}c_{n-2} \dots c_0$ and the carry-out by c_{out} . Which one of the following options correctly identifies the overflow condition? **[2006]**
- (A) $c_{out}(\overline{a_{n-1} \oplus b_{n-1}})$
(B) $a_{n-1}b_{n-1}c_{n-1} + \overline{a_{n-1}b_{n-1}c_{n-1}}$
(C) $c_{out} \oplus c_{n-1}$
(D) $a_{n-1} \oplus b_{n-1} \oplus c_{n-1}$

5. Consider numbers represented in 4-bit gray code. Let $h_3h_2h_1h_0$ be the gray code representation of a number n and let $g_3g_2g_1g_0$ be the gray code of $(n + 1)$ modulo

16 value of the number. Which one of the following functions is correct? **[2006]**

- (A) $g_0(h_3h_2h_1h_0) = \Sigma(1, 2, 3, 6, 10, 13, 14, 15)$
(B) $g_1(h_3h_2h_1h_0) = \Sigma(4, 9, 10, 11, 12, 13, 14, 15)$
(C) $g_2(h_3h_2h_1h_0) = \Sigma(2, 4, 5, 6, 7, 12, 13, 15)$
(D) $g_3(h_3h_2h_1h_0) = \Sigma(0, 1, 6, 7, 10, 11, 12, 13)$

6. How many 3-to-8 line decoders with an enable input are needed to construct a 6-to-64 line decoder without using any other logic gates? **[2007]**

- (A) 7 (B) 8
(C) 9 (D) 10

7. Suppose only one multiplexer and one inverter are allowed to be used to implement any Boolean function of n variables. What is the minimum size of the multiplexer needed? **[2007]**

- (A) 2^n line to 1 line (B) 2^{n+1} line to 1 line
(C) 2^{n-1} line to 1 line (D) 2^{n-2} line to 1 line

8. In a look-ahead carry generator, the carry generate function G_i and the carry propagate function P_i for inputs A_i and B_i are given by:

$$P_i = A_i \oplus B_i \text{ and } G_i = A_i B_i$$

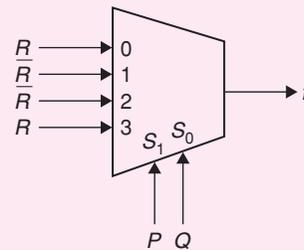
The expressions for the sum bit S_i and the carry bit C_{i+1} of the look-ahead carry adder are given by:

$$S_i = P_i \oplus C_i \text{ and } C_{i+1} = G_i + P_i C_i, \text{ where } C_0 \text{ is the input carry.}$$

Consider a two-level logic implementation of the look-ahead carry generator. Assume that all P_i and G_i are available for the carry generator circuit and that the AND and OR gates can have any number of inputs. The number of AND gates and OR gates needed to implement the look-ahead carry generator for a 4-bit adder with S_3, S_2, S_1, S_0 and C_4 as its outputs are respectively: **[2007]**

- (A) 6, 3 (B) 10, 4
(C) 6, 4 (D) 10, 5

9. The Boolean expression for the output f of the multiplexer shown below is



- (A) $\overline{P \oplus Q \oplus R}$
(B) $P \oplus Q \oplus R$
(C) $P + Q + R$
(D) $\overline{P + Q + R}$

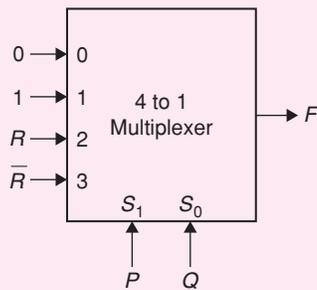
10. The amount of ROM needed to implement a 4-bit multiplier is [2012]
 (A) 64 bits
 (B) 128 bits
 (C) 1K bits
 (D) 2K bits

11. In the following truth table, $V = 1$ if and only if the input is valid.

Inputs				Outputs		
D_0	D_1	D_2	D_3	X_0	X_1	V
0	0	0	0	x	x	0
1	0	0	0	0	0	1
x	1	0	0	0	1	1
x	x	1	0	1	0	1
x	x	x	1	1	1	1

What function does the truth table represent? [2013]

- (A) Priority encoder
 (B) Decoder
 (C) Multiplexer
 (D) Demultiplexer
12. Consider the 4-to-1 multiplexer with two select lines S_1 and S_0 given below.



The minimal sum-of-products form of the Boolean expression for the output F of the multiplexer is [2014]

- (A) $\bar{P}Q + Q\bar{R} + P\bar{Q}R$
 (B) $\bar{P}Q + \bar{P}Q\bar{R} + P\bar{Q}R + P\bar{Q}R$
 (C) $\bar{P}QR + \bar{P}Q\bar{R} + Q\bar{R} + P\bar{Q}R$
 (D) $PQ\bar{R} + \bar{P}QR + \bar{P}Q\bar{R} + Q\bar{R} + P\bar{Q}R$
13. Consider the following combinational function block involving four Boolean variables x, y, a, b , where x, a, b are inputs and y is the output. [2014]

$f(x, y, a, b)$

if $(x \text{ is } 1) y = a;$
 else $y = b;$
 }

Which one of the following digital logic blocks is the most suitable for implementing this function?

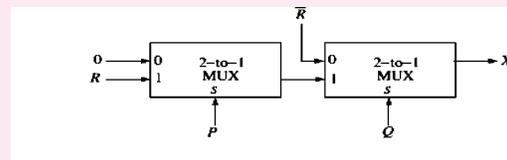
- (A) Full adder
 (B) Priority encoder
 (C) Multiplexer
 (D) Flip-flop
14. Let \oplus denote the Exclusive OR (X-OR) operation. Let '1' and '0' denote the binary constants. Consider the following Boolean expression for F over two variables P and Q :

$F(P, Q) = ((1 \oplus P) \oplus (P \oplus Q)) \oplus ((P \oplus Q) \oplus (Q \oplus 0))$ [2014]

The equivalent expression for F is

- (A) $P + Q$
 (B) $\overline{P + Q}$
 (C) $P \oplus Q$
 (D) $\overline{P \oplus Q}$
15. A half adder is implemented with XOR and AND gates. A full adder is implemented with two half adders and one OR gate. The propagation delay of an XOR gate is twice that of an AND/OR gate. The propagation delay of an AND/OR gate is 1.2 microseconds. A 4-bit ripple-carry binary adder is implemented by using four full adders. The total propagation time of this 4-bit binary adder in microseconds is _____ [2015]

16. Consider the two cascaded 2-to-1 multiplexers as shown in the figure.



minimal sum of products form of the output x is

The minimal sum of products form of the output X is [2016]

- (A) $\bar{P} \bar{Q} + PQR$
 (B) $\bar{P} Q + QR$
 (C) $PQ + \bar{P} \bar{Q} R$
 (D) $\bar{P} \bar{Q} + PQR$
17. When two 8-bit numbers $A_7 \dots A_0$ and $B_7 \dots B_0$ in 2's complement representation (with A_0 and B_0 as the least significant bits) are added using a ripple-carry adder, the sum bits obtained are $S_7 \dots S_0$ and the carry bits are $C_7 \dots C_0$. An overflow is said to have occurred if [2017]
- (A) the carry bit C_7 is 1
 (B) all the carry bits (C_7, \dots, C_0) are 1
 (C) $(A_7 \cdot B_7 \cdot \bar{S}_7 + \bar{A}_7 \cdot \bar{B}_7 \cdot S_7)$ is 1
 (D) $(A_0 \cdot B_0 \cdot \bar{S}_0 + \bar{A}_0 \cdot \bar{B}_0 \cdot S_0)$ is 1

ANSWER KEYS**EXERCISES****Practice Problems 1**

1. D 2. B 3. B 4. D 5. D 6. A 7. C 8. B 9. D 10. A
11. A 12. B 13. B 14. C 15. C 16. D 17. D 18. B 19. C 20. A
21. D

Practice Problems 2

1. C 2. B 3. A 4. D 5. C 6. C 7. C 8. B 9. B 10. C
11. B 12. C 13. C 14. A 15. C 16. C 17. B 18. C 19. C 20. D
21. A

Previous Years' Questions

1. B 2. A 3. A 4. B 5. C 6. C 7. C 8. B 9. B 10. D
11. A 12. A 13. C 14. D 15. 19.2 16. D 17. C

Chapter 4

Sequential Circuits

LEARNING OBJECTIVES

- Sequential circuit
- Basic storage elements
- Latches (SR Latch, D Latch, JK Latch)
- Flip-flops (JK flip-flop, T flip-flop, D flip-flop)
- Counters
- Asynchronous counter design
- Synchronous counter design
- Registers
- Various types of registers
- Application of shift register

SEQUENTIAL CIRCUITS

In sequential circuits the output depends on the input as well as on the previous history of output, i.e., they contain memory elements.

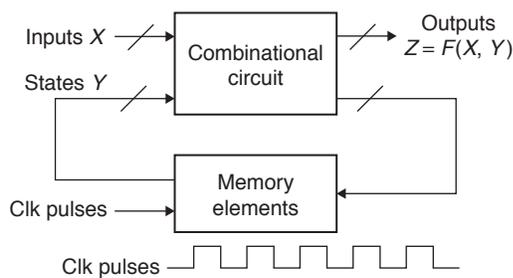


Figure 1 Block diagram of sequential circuit

Table 1 Comparison between combinational and sequential circuits

Combinational Circuits	Sequential Circuits
1. Output at any time depends on the combine set of input applied to it simultaneously at that instant of time	Output depends on the present input as well as on the previous history of output
2. Contains no memory element	Contains at least one memory element
3. Easy to design due to absence of memory	Difficult to design
4. Totally described by the set of output values	Its performance is totally described by the set of subsequent values as well as set of output values
5. Faster in speed because all inputs are primary inputs and applied simultaneously	Slower in speed because secondary inputs are also needed which are applied after delay
6. It need more hardware for realization	Less hardware required
7. Expensive in cost	Cheap in cost

Sequential circuits are of two types:

1. Clocked or synchronous
2. Unclocked or asynchronous

In synchronous sequential circuits the logic circuits action is allowed to occur in synchronization with the input clock pulse from a system clock.

In asynchronous sequential circuits the logic sequential action is allowed to occur at any time.

Basic Storage Elements

Latches and flip-flops

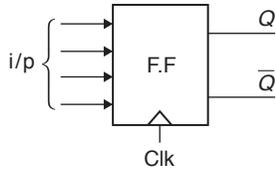
A storage element in digital circuit can maintain a binary state indefinitely until directed by an input signal to switch states. Storage elements that operate with signal levels (i.e., level triggering of signal inputs) are referred to as latches. Those controlled by a clock transition (i.e., edge triggering) are flip-flops.

Latches and flip-flops are related because latches are basic circuits from which all flip-flops are constructed. Latches are useful for storing binary information and for the design of asynchronous sequential circuits. But latches are not practical for use in synchronous sequential circuits, so we use flip-flops.

Flip-flops

They are also known as bistable multivibrators. This is a basic memory element to store 1-bit of information 0 or 1 and is used in storage circuits, counters, shift register, and many other computer applications. It has two stable states: 1 and 0. The high state is called set state and zero as reset.

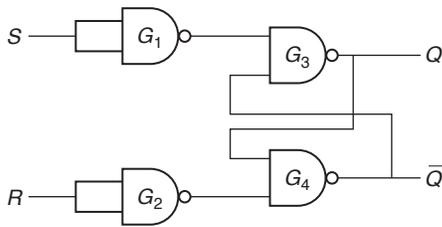
It has two outputs one being the complement of the other usually designated by Q and \bar{Q} .



There are different types of flip-flops S-R flip-flop, D-flip-flop, T-flip-flop, J-K flip-flops, etc.

LATCHES

(i) **S-R Latch:** The simplest latch is called S-R latch. S-R means Set-Reset. It has two outputs Q and \bar{Q} and two inputs S and R , which represent set or reset signal.



Above figure shows two cross coupled gates G_3 and G_4 and inverters G_1 and G_2 . Here output of G_3 is connected to the input of G_4 and output of G_4 is applied to the input of G_3 . $S = 1, R = 0$ output of $G_1 = 0$ and $G_2 = 1$. Since one of the input of G_3 is 0, so its output will be certainly 1 and consequently both input of G_4 will be 1 and the output $\bar{Q} = 0$.

For $S = 1, R = 0, Q = 1, \bar{Q} = 0$. $S = 0, R = 1$ the output will be $Q = 0$ and $\bar{Q} = 1$. The first of the input condition $S = 1$ and $R = 0$ makes $Q = 1$ which referred as the set state and the second condition $S = 0$ and $R = 1$ makes $Q = 0$ which is referred as reset state.

For $S = 0$ and $R = 0$ output of both G_1 and G_2 will be one and hence there will be no change in Q and \bar{Q} .

For $S = R = 1$, both the outputs Q and \bar{Q} will try to become one, which produces invalid results and should not be used for the above latch.

Input		Output		
S	R	Q	\bar{Q}	State
1	0	1	0	Set
0	1	0	1	Reset
0	0	0	0	No change
1	1	?	?	Invalid

(ii) **SR latch by using NAND/NOR gates:** The SR latch is a circuit with two cross-coupled NOR gates or two cross-coupled NAND gates. Two inputs labelled S for set and R for reset. Latch will have two outputs:

- Q : output state in normal form and
- Q' : output state in complemented form.

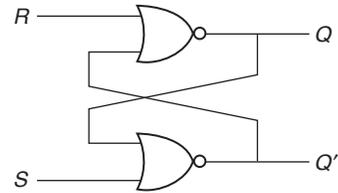


Figure 2 Logic diagram for SR latch

S	R	Q_{n+1}	Q'_{n+1}	
0	0	Q_n	Q'_n	(no change)
0	1	0	1	(Reset)
1	0	1	0	(set)
1	1	0	0	(invalid)

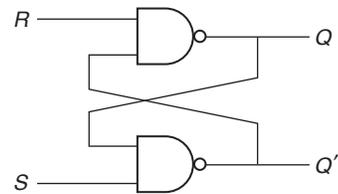
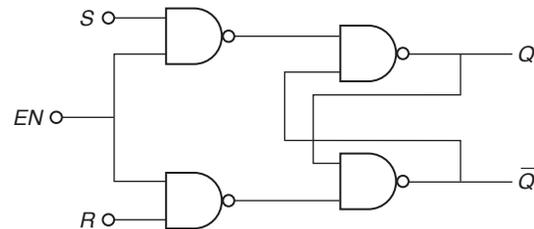


Figure 3 $\bar{S} \bar{R}$ Latch

S	R	Q_{n+1}	\bar{Q}_{n+1}	
0	0	1	1	(Invalid)
0	1	1	0	(Set)
1	0	0	1	(Reset)
1	1	Q_n	Q'_n	(No change)

$\bar{S} \bar{R}$ latch is active low SR latch

(iii) **SR latch with control input:** The working of gated SR latch is exactly the same as SR latch when the EN pulse is present. When the EN pulse is not present (EN pulse = 0) the gates G_1 and G_2 are inhibited and will not respond to the input.

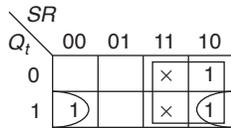


Characteristic table of SR latch shows the operation of latch in tabular form. Q_i stands as the binary state of the latch before the application of latch pulse and referred to as the present state. The S and R columns give the possible values of the inputs and Q_{i+1} is the state of the latch after the application of a single pulse, referred to as next stage. EN input is not included in the characteristic table.

Characteristic table for SR latch is given below:

Q_t	S	R	Q_{t+1}
0	0	0	0
0	0	1	0
0	1	0	1
0	1	1	X
1	0	0	1
1	0	1	0
1	1	0	1
1	1	1	X

Characteristic equation of the latch is derived from the K-map.



$$\therefore Q_{t+1} = S + \bar{R}Q_t$$

This equation specifies the value of the state as a function of the present state and the inputs.

(iv) **Preset and clear inputs:** For the latch/flip-flop, when the power is switched ON, the state of the circuit is uncertain. It can be either $Q = 0$ (reset) or $Q = 1$ (set) state.

In many applications it is desired to set or reset the circuit, so that initial state of the circuit will be known. This is accomplished by using the asynchronous, inputs referred to as preset (Pr) and clear (Clr), inputs.

These inputs can be applied any time, and are not synchronized with EN input/Clr input.

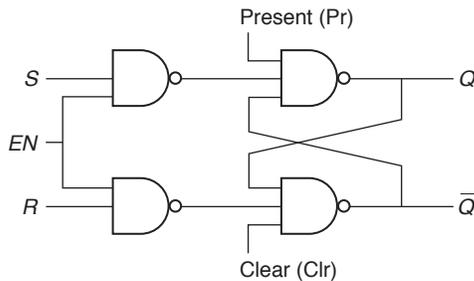


Figure 4 SR latch with Pr and Clr inputs

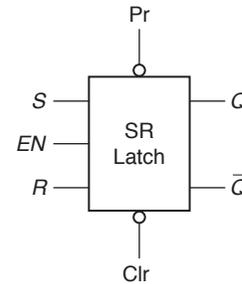
If $Pr = Clr = 1$, the circuit operates as of S-R latch explained previously.

If $Pr = 0, Clr = 1$, the output Q will become 1, which in turn changes $\bar{Q} = 0$.

If $Pr = 1, Clr = 0$ the output \bar{Q} will become 1, which in turn changes $Q = 0$.

If $Pr = Clr = 0$, both Q and \bar{Q} will become 1, which is invalid case, so $Pr = Clr = 0$ condition must not be used.

Pr	Clr	Q_{n+1}
1	1	Q – No change
0	1	1 – Set
1	0	0 – Reset
0	0	X – Invalid



(v) **D latch (Transparent latch):** One way to eliminate the invalid condition of SR latch (when $S = R = 1$) is to ensure that inputs S and R are never equal to 1 at the same time.

By connecting a NOT gate between S and R inputs, i.e., complement of S will be given to R , we can form D latch as shown in block diagram.

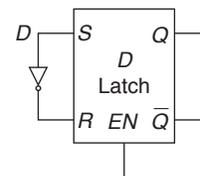


Figure 5 Block diagram for D latch

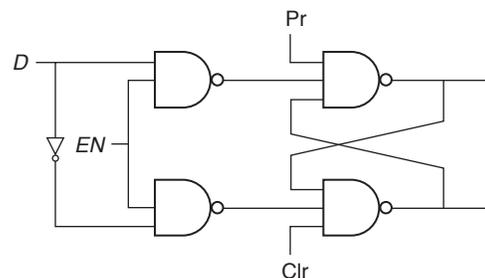


Figure 6 Logic diagram for D latch

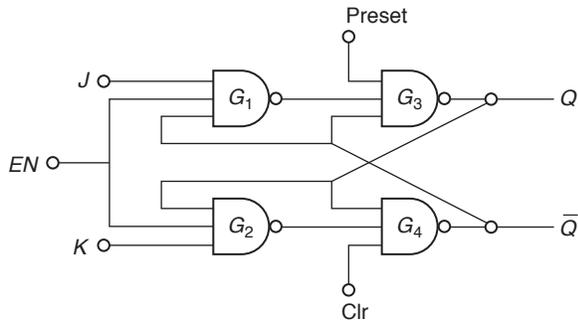
EN	D	Q_{n+1}
0	X	Q_n – No change (Disabled)
1	0	0 – Reset state
1	1	1 – Set state

When $EN = 0$, the circuit will be disabled and input D will not have any effect on output, and output will be same as previous state.

When $EN = 1, D = 0$, i.e., $S = 0, R = 1$ which makes output $Q = 0$ and $\bar{Q} = 1$ (Reset state).

When $EN = 1, D = 1$, i.e., $S = 1, R = 0$ which makes output $Q = 1$, and $\bar{Q} = 0$ (Set state).

(vi) **JK latch:** The function of JK latch is identical to that of SR latch except that it has no invalid state as that of SR latch where $S = R = 1$. In this case the state of the output is changed as complement of previous state.



EN	J	K	Q_{t+}	
1	0	0	Q_t	No change
1	0	1	0	Reset
1	1	0	1	Set
1	1	1	\bar{Q}_t	Toggle
0	x	x	Q_t	No change

JK latch by using SR latch: The uncertainty of SR flip-flop (when $S = 1, R = 1$) can be eliminated by converting it into JK latch.

The data inputs J and K , which are ANDed with \bar{Q} and Q respectively, to obtain S and R inputs.

$$S = J \cdot \bar{Q}, \quad R = K \cdot Q$$

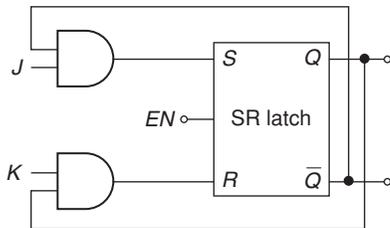
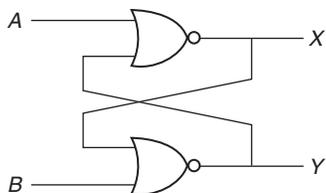


Figure 7 JK latch by using SR latch

J	K	S	R	Q_{n+1}	\bar{Q}_{n+1}
0	0	0	0	Q_n	\bar{Q}_n -No change
0	1	0	Q_n	0	1-Reset
1	0	\bar{Q}_n	0	1	0-Set
1	1	\bar{Q}_n	Q_n	\bar{Q}_n	Q_n -Toggle

Example 1: The following binary values were applied to A and B inputs of NOR gate latch shown in the figure, in the sequence indicated below. $A = 1, B = 0; A = 1, B = 1; A = 0, B = 0$. The corresponding stable X, Y outputs will be



- (A) 10, 01, 10 or 01
- (B) 11, 00, 10
- (C) 01, 00, 10 or 01
- (D) 10, 11, 10 or 01

Solution: (C)

Given circuit is RS latch with NOR gates.

By comparing with RS latch $A = R, B = S$, and $X = Q, Y = \bar{Q}$, so from truth table of RS latch

S/B	R/A	Q/X	\bar{Q}/Y	
0	1	0	1	Reset (Invalid)
1	1	0	0	
0	0	1	0	(Same as previous state)
		0	1	

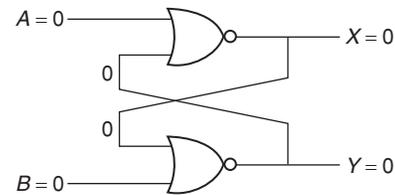
After invalid case $S = 1, R = 1$, i.e., $A = B = 1$,

The output $Q = 0, \bar{Q} = 0$, i.e., $X = Y = 0$

By applying $A = 0, B = 0$

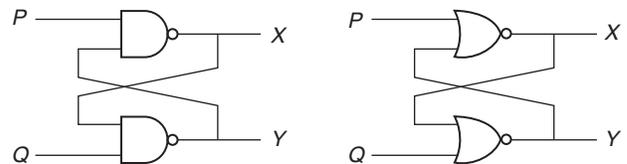
The output X becomes $\overline{(0+0)} = 1$ and which in turn changes

$$Y = \overline{(B + X)} = \overline{(0 + 1)} = 0$$



(or) the output Y becomes $\overline{(0+0)} = 1$ and which in turn changes $X = \overline{(Y + A)} = \overline{(0 + 1)} = 0$. So, output (X, Y) cannot be predicted after the invalid condition. So, $X = 0, Y = 1$ or $X = 1, Y = 0$

Example 2: Refer to the NAND and NOR latches shown in the figure the inputs (P, Q) for both the latches are first made $(1, 0)$ and then after a few seconds, made $(0, 0)$. The corresponding stable outputs (X, Y) are



- (A) NAND: first $(0, 1)$ then $(0, 1)$; NOR: first $(1, 0)$ then $(1, 0)$
- (B) NAND: first $(0, 1)$ then $(1, 1)$; NOR: First $(0, 1)$ then $(0, 1)$
- (C) NAND: first $(1, 0)$ then $(0, 0)$; NOR: first $(1, 0)$ then $(1, 0)$
- (D) NAND: first $(1, 0)$, then $(1, 0)$; NOR: first $(1, 0)$ then $(1, 1)$

Solution: (B)

From the truth table of SR latch and $\bar{S} \bar{R}$ latch SR latch with NOR gates:

For $(P, Q) = (1, 0) = (R, S)$ output $(X, Y) = (Q, \bar{Q}) = (0, 1)$

Then (P, Q) are made $(0, 0)$, i.e., $(R, S) = (0, 0)$, which results in no change at output. So, $(X, Y) = (Q, \bar{Q}) = (0, 1)$ SR latch with NAND gates:

For $(P, Q) = (1, 0) = (S, R)$ output $(X, Y) = (Q, \bar{Q}) = (0, 1)$. Then (P, Q) are made $(0, 0)$, i.e., $(S, R) = (0, 0)$ which is invalid conditions for $\bar{S} \bar{R}$ latch. So, $(X, Y) = (Q, \bar{Q}) = (1, 1)$

(vii) Race around condition: The difficulties of both the inputs $(S = R = 1)$ being not allowed in an SR latch is eliminated in JK latch by using the feedback connection from the output to the input of the gate G_1 and G_2 . In a normal JK latch if $J = K = 1$ and $Q = 0$ and enable signal is applied without RC differentiator, after a time interval Δt (the propagation delay through two NAND gate in series) the output will change to $Q = 1$. Now we have $J = K = 1$ and $Q = 1$ and after another time interval of Δt the output will change back to $Q = 0$. Hence for the duration of (t_p) of the enable signal the output will oscillates back and forth between 0 and 1. At the end of the enable signal the values of Q is uncertain. This situation is referred to as race around condition.

The race around condition can be avoided if enable time period $t_p < \Delta t$ but it may be difficult to satisfy this condition, because of very small propagation delays in ICs. To solve this problem the enable signals are converted to narrow spike using RC differentiator circuit having a short time constant. Its output will be high during the high transmission time of the enable. Another method to avoid this problem is master-slave JK flip-flop.

FLIP-FLOPS

(i) Master-slave JK flip-flop: This is a cascade of 2 SR latches with feedback from the output of the second SR latch to the inputs of the first as shown in the figure below.

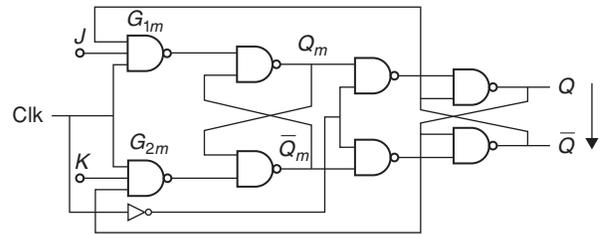
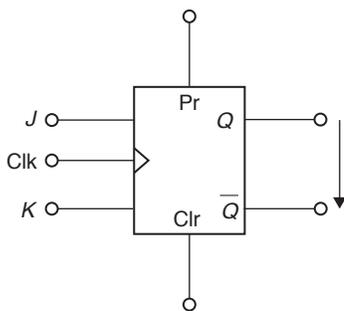
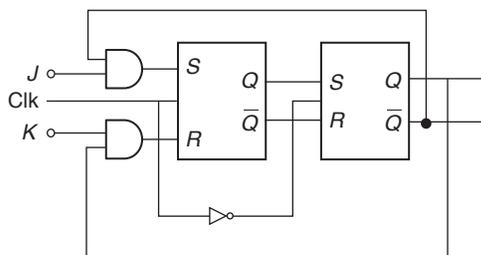


Figure 8 Logic diagram of JK flip-flop

Positive clock pulse is applied to the first latch and the clock pulse will be inverted before its arrival at the second latch. When $Clk = 1$, the first latch is enabled and the outputs Q_m and \bar{Q}_m responds to the inputs J and K , according to the truth table of JK latch. At this time the 2nd latch is inhibited because its clock is low ($Clk = 0$). When the clock goes low ($Clk = 0$), the first latch is inhibited and the second is enabled. Therefore, the outputs Q and \bar{Q} follow the outputs Q_m and \bar{Q}_m , respectively. Since the second latch simply follows the first one, it is referred to as slave and the first one as the master. Hence this configuration is known as master-slave JK flip-flop. In this circuit, the input to the gate G_{1m} and G_{2m} do not change, during the clock pulse levels.

The race around condition does not exist.

Table 2 State/characteristic Table

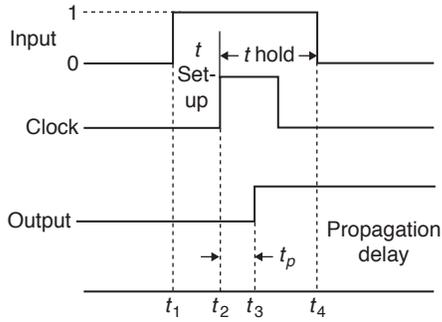
Clk	J	K	Q_t	Q_{t+1}
↓	0	0	0	0
↓	0	0	1	1
↓	1	0	0	1
↓	1	0	1	1
↓	0	1	0	0
↓	0	1	1	0
↓	1	1	0	1
↓	1	1	1	0

	JK			
Q	00	01	11	10
0			1	1
1	1			1

$$Q_{t+1} = J\bar{Q}_t + Q_t\bar{K}$$

(ii) Flip-flop switching time: In designing circuits with flip-flop the following parameters are important:

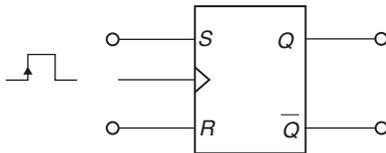
- 1. Set-up time:** The minimum amount of time required for the data input to be present before the clock arrived.
- 2. Hold time:** The minimum amount of time that the data input to be present after the clock trigger arrived.
- 3. Propagation delay:** The amount of time it takes for the output to change states after an input trigger.



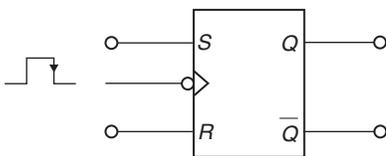
For example, $t_{\text{setup}} = 50 \text{ m sec}$ and $t_{\text{hold}} = 5 \text{ m sec}$, the data bit has to be the input at least 50 m sec before the clock bit arrives and hold at least 5 m sec after the clock edge.

(iii) **Triggering of flip-flop:** The flip-flop can be triggered to set or reset either at one of the edges of the clock pulse. There are three types of triggering as described below:

1. **Positive edge triggering flip-flop:** These set or reset at the positive (rising or leading) edge of the clock pulse depending upon the state of i/p signal and o/p remain steady for 1 clock period. Positive edge triggering is indicated by an arrow head at the clock terminal of the flip-flop.



2. **Negative edge triggered flip-flop:** There are flip-flops those in which state transmissions take place only at the negative edge (falling or trailing) of the clock signal. Negative edge triggering is indicated by arrow head with bubble at the clock terminal.



3. **Level triggering:** Level triggering means the specified action occurs based on the steady state value of the input. That is, when a certain level is reached (0 or 1) the output will change states level triggering will be used in latches.

(iv) **D flip-flop:** It receives the designation from its ability to hold data into its internal storage. An SR/JK flip-flop has two inputs. It requires two inputs S/J and R/K to store 1 bit. This is a serious disadvantage in many application to overcome the difficulty D flip-flop has been developed which has only one input line. A D

flip-flop can be realized using a SR/JK as show in the figure below.

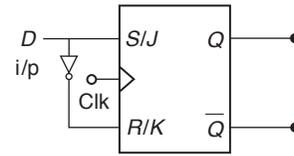


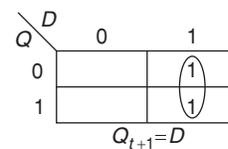
Table 3 Truth Table

clk	D	Q_{t+1}
\nearrow	X	Q_t
\uparrow	0	0
\uparrow	1	1

There is no raising problem with D flip-flop. High or 1 state will set the flip-flop and a low or 0 state will reset the flip-flop. The presence of inverter at the input ensure that S/J and R/K inputs will always be in the opposite state.

Table 4 Characteristic Table of D Flip-flop

Q_t	D	Q_{t+1}
0	0	0
0	1	1
1	0	0
1	1	1



From the characteristic table of D flip-flop, the next state of the flip-flop is independent of the present state since $Q_{t+1} = D$, whether $Q_t = 0$ or 1.

(v) **T flip-flop:** In a JK flip-flop $J = K = 1$ and the resulting flip-flop is referred to as a T flip-flop.

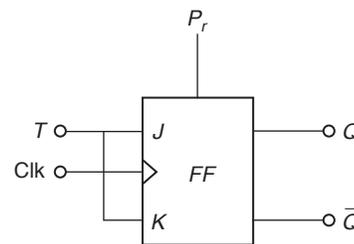


Table 5 Truth Table

Clk	T	Q_{n+1}
\uparrow	0	Q_n
\uparrow	1	\bar{Q}_n
\nearrow	x	Q_n

If $T = 1$, it acts as a toggle switch for every Clk pulse with high input, the Q changes to its opposite state.

Table 6 Characteristic Table

Q_t	T	Q_{t+1}
0	0	0
0	1	1
1	0	1
1	1	0
1	0	1

$T \backslash Q$	0	1
0		1
1	1	

$$Q_{t+1} = T\bar{Q}_t + Q_t\bar{T}$$

(vi) **Excitation table of flip-flops:** The truth table of flip-flop is also referred to as the characteristic table, which specifies the operational characteristic of flip-flop.

Sometimes we come across situations in which present state and the next state of the circuit are known and we have to find the input conditions that must prevail to cause the desired transition of the state.

Consider initially JK flip-flop output $Q_n = 1, \bar{Q}_n = 0$, after clock pulse it changed to $Q_{n+1} = 0, \bar{Q}_{n+1} = 1$,

The input conditions, which made this transition, can be

Toggle – for $J = 1, K = 1, Q_{n+1} = \bar{Q}_n$
or

Reset – for $J = 0, K = 1, Q_{n+1} = 0, \bar{Q}_{n+1} = 1$

From the above conditions we can conclude that for transition $Q_n = 1$ to $Q_{n+1} = 0$ occurs when $J = 0$ (or) 1 (don't care) and $K = 1$.

Similarly, input conditions can be found out for all possible situations.

Table 7 Excitation table of flip-flop.

Present State	Next State	SR Flip-flop	JK Flip-flop	T Flip-flop	D Flip-flop		
Q_n	Q_{n+1}	S	R	J	K	T	D
0	0	0	×	0	×	0	0
0	1	1	0	1	×	1	1
1	0	0	1	×	1	1	0
1	1	×	0	×	0	0	1

These excitation tables are useful in the design of synchronous circuits.

(vii) **State diagrams of flip-flops:** State diagram is a directed graph with nodes connected with directed arcs. State of the circuit is represented by the node, the directed arcs represent the state transitions, from present state (node) to next state (node) at the occurrence of clock pulse.

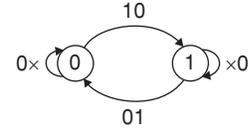


Figure 9 State diagram of SR flip-flop

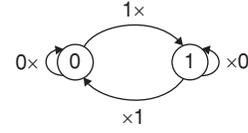


Figure 10 State diagram of JK flip-flop

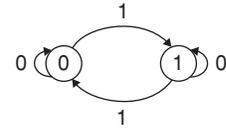


Figure 11 State diagram of T flip-flop

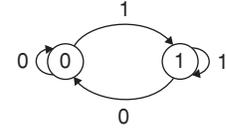


Figure 12 State diagram of D flip-flop

(viii) **Conversion of one flip-flop to other flip-flop**

Conversion of T flip-flop to JK flip-flop

1. Write the characteristic table of required flip-flop (here JK).
2. Write the excitation table of available or given Flip-flop (here T).
3. Solve for inputs of given flip-flop in terms of required flip-flop inputs and output.

Table 8 JK flip-flop characteristic and T flip-flop excitation table

JK Flip-flop Characteristic Table		T Flip-flop Excitation Table		
J	K	Q_n	Q_{n+1}	T
0	0	0	0	0
0	0	1	1	0
0	1	0	0	0
0	1	1	0	1
1	0	0	1	1
1	0	1	1	0
1	1	0	1	1
1	1	1	0	1

$J \backslash KQ_n$	00	01	11	10
0			1	
1	1		1	1

$$T = J\bar{Q}_n + KQ_n$$

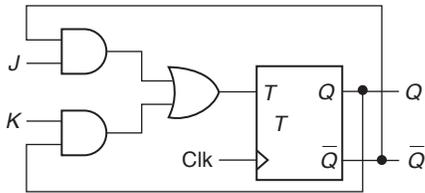


Figure 13 D Flip-flop by using other flip-flops

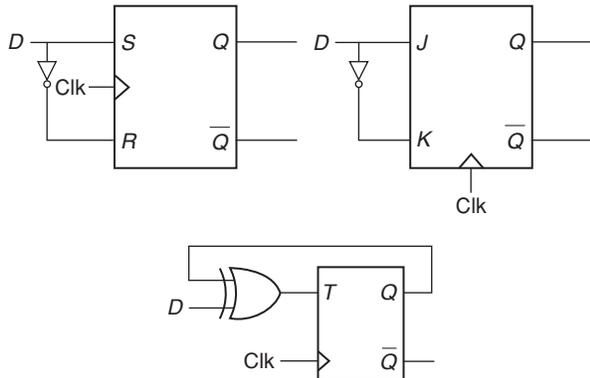
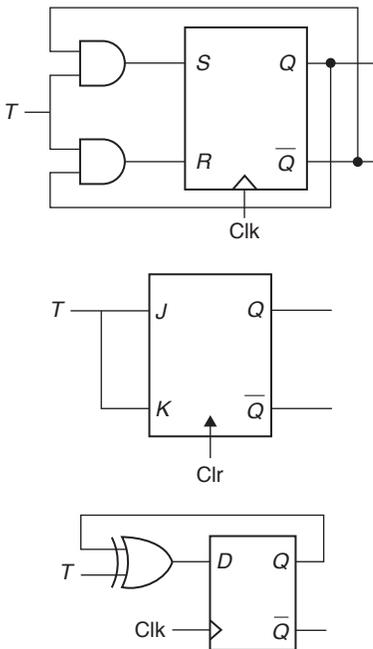
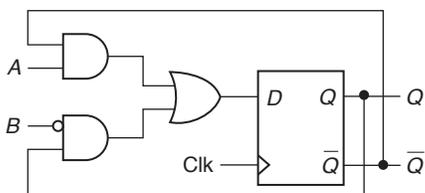


Figure 14 T flip-flop by using other flip-flops



Example 3: A sequential circuit using D flip-flop and logic gates is shown in the figure, where A and B are inputs and Q is output.



The circuit is

- (A) SR flip-flop with inputs $A = S, B = R$
- (B) SR flip-flop with inputs $A = \bar{R}, B = S$
- (C) JK flip-flop with inputs $A = J, B = K$
- (D) JK flip-flop with inputs $A = K, B = \bar{J}$

Solution: (C)

The characteristic equation of D flip-flop is

$$Q_{n+1} = D$$

Here input $D = A\bar{Q}_n + \bar{B}Q_n$

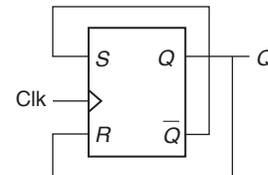
So, output $Q_{n+1} = A\bar{Q}_n + \bar{B}Q_n$

By comparing this equation with characteristic equation of JK

$$Q_{n+1} = J\bar{Q}_n + \bar{K}Q_n$$

If $A = J, B = K$, then this circuit works like JK flip-flop.

Example 4: The input Clk frequency for the flip-flop given is 10 kHz, then the frequency of Q will be



- (A) 10 kHz
- (B) 5 kHz
- (C) 20 kHz
- (D) 2.5 kHz

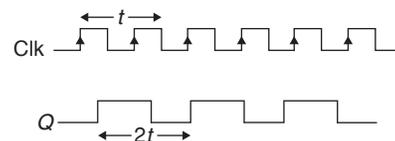
Solution: (B)

Form circuit we can say $S = \bar{Q}_n, R = Q_n$.

If initially $(Q_n, \bar{Q}_n) = (0, 1)$, then inputs $(S, R) = (1, 0)$, by applying clk pulse (Q_{n+1}, \bar{Q}_{n+1}) becomes $(1, 0) \dots$

Clk	Q_n	\bar{Q}_n	S	R	Q_{n+1}	\bar{Q}_{n+1}
1	0	1	1	0	1	0
2	1	0	0	1	0	1
3	0	1	1	0	1	0
4	1	0	0	1	0	1

The output Q_{n+1} toggles for every clock pulse.



So frequency of $Q = \frac{1}{2t} = \frac{f}{2} = \frac{10}{2} = 5$ kHz

Table 9 Identification of up/down Counters

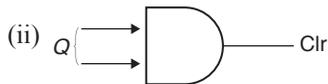
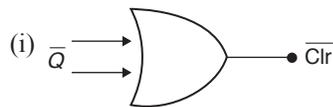
Clock Triggering	Q	Type
+ve edge	\bar{Q}	Up
+ve edge	Q	Down
-ve edge	\bar{Q}	Down
-ve edge	Q	Up

Clock is negative triggering pulse and Q is connected to next level clock, it is acting like a up counter.

Table 10 Identification of GATE to Clear the Flip-flops

Input to the Gate	Output of the Gate	Type of Gate
\bar{Q}	$\bar{\text{Clr}}$	OR
\bar{Q}	Clr	NOR
Q	$\bar{\text{Clr}}$	NAND
Q	Clr	AND

Example:



Example 6: Design and Implement a MOD-6 asynchronous counter using T flip-flops.

Solution: Counting sequence is 00, 001, 010, 011, 100, 101

After Pulses	States Q_3, Q_2, Q_1	Reset R
0	000	0
1	001	0
2	010	0
3	011	0
4	100	0
5	101	0
6	110	1
7	↓↓↓	
	000	0
	111	X

From the Truth table $R = Q_3 Q_2$

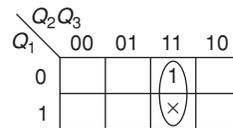
For active Low \bar{R} is used.

$\therefore R = 0$ for 000 to 101

$R = 1$ for 110

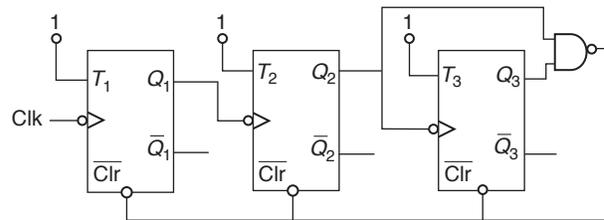
$R = X$ for 111

\therefore K-map is



$\therefore R = Q_2 Q_3$

Logic diagram is



Asynchronous Decode Counter

A ripple counter is an asynchronous sequential circuit, clock is applying only for LSB side. Decade ripple counter it counts from 0 to 9 for up counter.

MOD-10 counter it counts starting from 0000 to 1001. If the NAND gate output is logic '0' at that instant the counter reset to initial state.

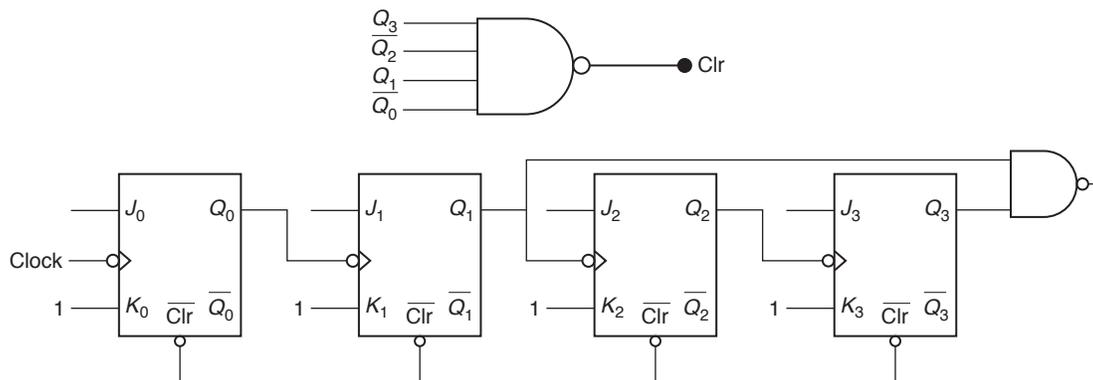


Figure 16 MOD-10 or decade counter

To design a MOD- N counter minimum number of flip-flops required is

$$N \leq 2^n$$

where $N \rightarrow$ MOD

$n \rightarrow$ No. of flip-flops

Example:

MOD-5 counter

$$5 \leq 2^n$$

$\therefore n = 3$

Operating Clock Frequency

(i) **Synchronous counter:**

$$f_{\text{clk}} \leq \frac{1}{t_{pd}}$$

(ii) **Asynchronous counter:**

$$f_{\text{clk}} \leq \frac{1}{nt_{pd}}$$

Output frequency of the MOD- N counter is

$$\Rightarrow f_o = \frac{f_{\text{clk}}}{N}$$

(iii) **Synchronous counter:** When counter is clocked such that each flip-flop in the counter is triggered at the same time, the counter is called as synchronous counter.

- Synchronous counters have the advantage of high speed and less severe decoding problems.
- Disadvantage is having more circuiting than that of asynchronous counter.

Synchronous Series Carry Counters

For normal ring counters to count N sequence total N flip-flops are required.

Unused states in ring counter = $2^N - N$.

Unused states in Johnson ring counter = $2^N - 2N$.

Asynchronous counters are slower than the synchronous counters. By using synchronous series carry adders we can design MOD- N counter with n Flip-flops-only.

For non-binary counters $N \leq 2^n$

3-bit series carry up counter

It counts from initial state 000 to 111.

\therefore MOD = $2^3 = 8$ states

\therefore MOD-8

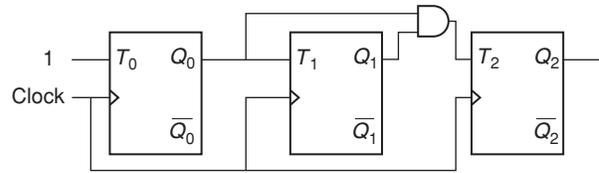


Figure 17 3-bit series carry counter

$$f_{\text{clk}} \leq \frac{1}{t_{pd} + (n-2)t_{pd \text{ AND}}}$$

where

$t_{pd} \rightarrow$ Propagation delay of each flip-flop.

$t_{pd \text{ AND}} \rightarrow$ Propagation delay of AND gate.

$n \rightarrow$ Number of flip-flops.

In this, Q_0 toggles for every clock pulse.

Q_1 toggles when Q_0 is 1.

Q_2 toggles when o/p of AND gate is logic 1.

Note: To design a synchronous series carry down counter. Connect \overline{Q}_0 to the next flip-flop input.

Design of Synchronous Counter

Step I: Determine the required number of flip-flop.

Step II: Draw the state diagram showing all possible states.

Step III: Select the type of flip-flop to be used and write the excitation table.

Step IV: Obtain the minimal expressions for the excitations of the FFs using the K-maps.

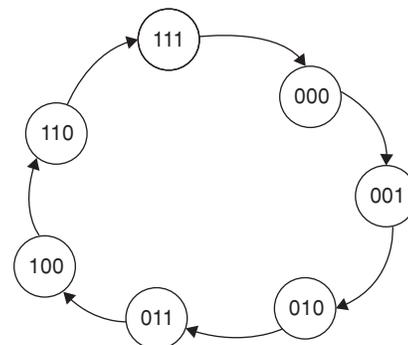
Step V: Draw a logic diagram based on the minimal expression. Let us employ these techniques to design a MOD-8 counter to count in the following.

Example 7: Sequence: 0, 1, 2, 3, 4, 5, 6, and 7. Design a synchronous counter by using JK flip-flops.

Solution:

Step I: Determine the required number of flip-flops. The sequence shows a 3-bit up counter that requires 3 flip-flops.

Step II: Draw the state diagram.

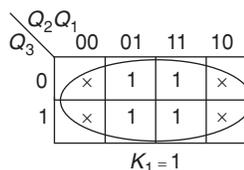
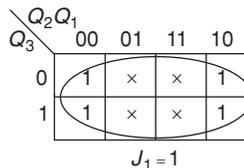
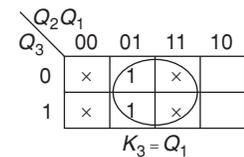
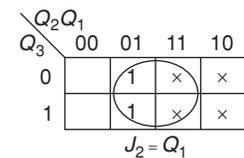
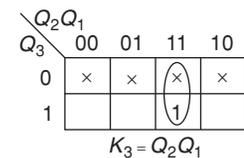
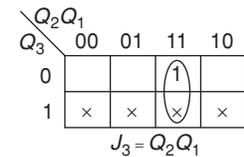


Step III: Select the type of flip-flop to be used and write the excitation table.

JK flip-flop is selected and excitation table of a 3-bit up counter is

PS			NS			Required Excitation					
Q_3	Q_2	Q_1	Q_3	Q_2	Q_1	J_3	K_3	J_2	K_2	J_1	K_1
0	0	0	0	0	1	0	x	0	x	1	x
0	0	1	0	1	0	0	x	1	x	x	1
0	1	0	0	1	1	0	x	x	0	1	x
0	1	1	1	0	0	1	x	x	1	x	1
1	0	0	1	0	1	x	0	0	x	1	x
1	0	1	1	1	0	x	0	1	x	x	1
1	1	0	1	1	1	x	0	x	0	1	x
1	1	1	0	0	0	x	1	x	1	x	1

Step IV: Obtain the minimal expression using K-map.



Step V: Draw the logic diagram based on the minimal expression.

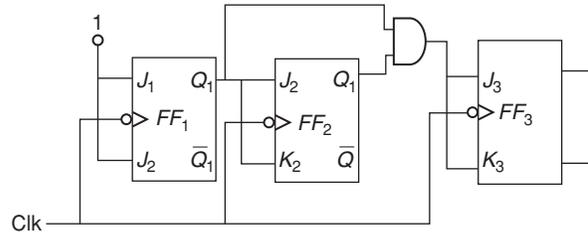


Table 11 Comparison between asynchronous counter and synchronous counter

Asynchronous Counter	Synchronous Counter
1. In this type of counter, flip-flops are connected in such a way that output of first flip-flop drives the clock for the next flip-flop	In this type there is no connection between output of first flip-flop and clock input of the next flip-flop
2. All the flip-flops are not clocked simultaneously	All the flip-flops are clocked simultaneously
3. Logic circuit is very simple even for more number of states	Design involves complex logic circuits as number of state increases
4. Main drawback of these counters is their low speed as the clock is propagated through number of flip-flops before it reaches last flip-flop	As clock is simultaneously given to all flip-flops, there is no problem of propagation delay. Hence they are preferred when number of flip-flops increases in the given design.

The main drawback of ripple counters is their high delays, if propagation delay of each flip-flop is assumed as x , then to get output of the first flip-flop it takes x , i.e., after x seconds the second flip-flop will get its clock pulse from previous stage, and output of second flip-flop will be out after another x seconds, similarly the final output of last flip-flop will be after nx seconds, where n is the number of flip-flops. So the propagation delay of ripple counter is nx , which is directly proportionate to the number of flip-flops.

The maximum frequency of operation of ripple counter is inverse of delay, $f_{max} = \frac{1}{nx}$

Maximum operating frequency is the highest frequency at which a sequential circuit can be reliably triggered. If the clock frequency is above this maximum frequency the flip-flops in the circuit cannot respond quickly and the operation will be unreliable.

In case of synchronous counters (synchronous circuits) as clock is applied simultaneously to all the flip-flops, the output of all the flip-flops change by x seconds (delay of one flip-flop) and this delay is independent of number of flip-flops used in circuit.

The maximum frequency of operation of synchronous counter is inverse of delay $f_{max} = \frac{1}{x}$

Example 8: The maximum operation frequency of a MOD-64 ripple counter is 33.33 kHz, the same flip-flops are used to design a MOD-32 synchronous counter, and then the maximum operating frequency of the new counter is

- (A) 400 kHz
- (B) 200 kHz
- (C) 40 kHz
- (D) 500 kHz

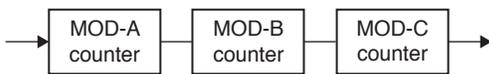
Solution: For ripple counter $f_{max} = \frac{1}{nx}$, given is a MOD-64 ripple counter, i.e., 2^6 states, so $n = 6$ flip-flops are required.

$$x = \frac{1}{33.33K \times 6} = 5\mu S$$

For synchronous counter

$$f_{max} = \frac{1}{x} = \frac{1}{5\mu S} = 0.2\text{ MHz} = 200\text{ kHz}$$

When multiple counters are connected in cascade, then the total number of states of the new counter is $A \times B \times C$, i.e., it will work as MOD- $A \times B \times C$ counter.



For example, decade counter counts from 0 to 9, 10 states – If two such decade counters are connected in cascade, then the total counting states will be $10 \times 10 = 100$, it will work as MOD-100 counter, which counts from 00 to 99.

REGISTERS

A number of flip-flops connected together such that data may be shifted into and shifted out of them is called a shift register. There are four basic types of shift register:

1. Serial-in–serial-out
2. Serial-in–parallel-out
3. Parallel-in–serial-out
4. Parallel-in–parallel-out

(i) Serial-in–serial-out:

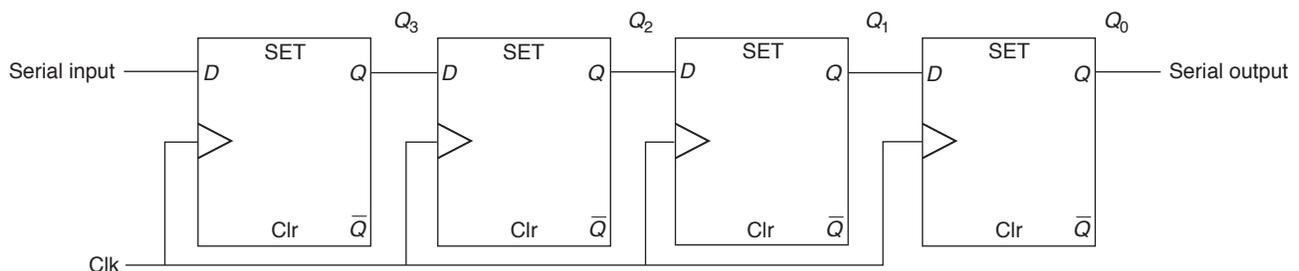
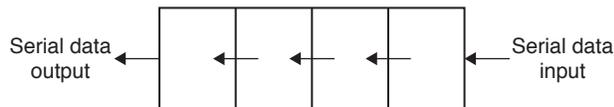
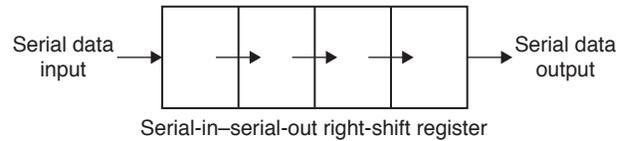
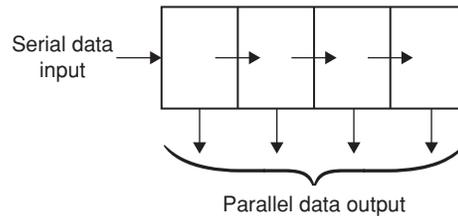


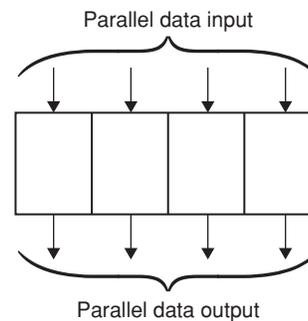
Figure 18 Serial input and serial output register



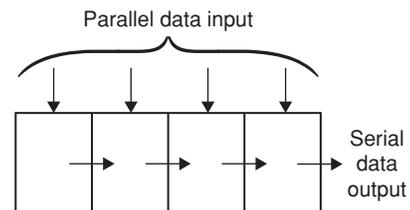
(ii) Serial-in–parallel-out:



(iii) Parallel-in–parallel-out:



(iv) Parallel-in–serial-out:



Serial input and serial output register: This type of shift register accepts data serially, i.e., one bit at a time and also outputs data serially. The logic diagram of 4-bit serial input, serial output, shift-right, shift register is shown in the following figure. With four D flip-flops the register can store up to four bits of data.

If initially, all flip-flops are reset, then by applying serial input 1101, the flip-flop states will change as shown in below table.

Clk	S.I	Q ₃	Q ₂	Q ₁	Q ₀
0	1	0	0	0	0
1	0	1	0	0	0
2	1	0	1	0	0
3	1	1	0	1	0
4		1	1	0	1

The first data bit 1 will appear at serial output after 4 clock pulses.

Application of Shift Registers

1. Delay line: Serial input and serial output shift register can be used to introduce delay in digital signals.

$$\text{Delay} = \text{no. of flip-flops} \times \frac{1}{\text{Clk frequency}} = \text{No. of flip-flops} \times \text{time period of clock pulse}$$
2. Serial to parallel, parallel to serial converter: SIPO, PISO registers used for data conversion.
3. Sequence generator: A circuit, which generates a prescribed sequence of bits, with clock pulses is called as sequence generator
 The minimum number of flip-flops 'n' required to generate a sequence of length 'S' bits is given by $S \leq 2^n - 1$

Shift register counters

One of the applications of the shift register is that they can be arranged to work as ring counters. Ring counters are constructed by modifying the serial-in, serial-out, shift registers. There are two types of ring counters—basic ring counter and twisted ring counter (Johnson counter). The basic ring counter is obtained from SISO shift register by connecting serial output to serial input.

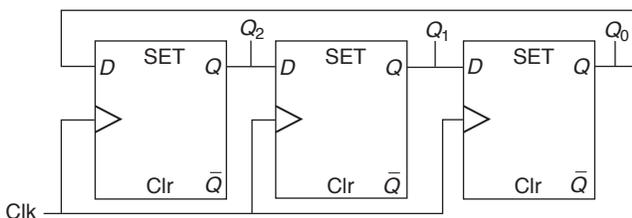
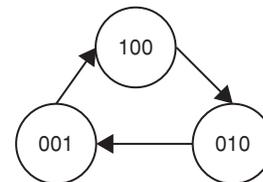


Figure 19 Ring counter

In most instances, only a single 1 or single 0 is in the register and is made to circulate around the register as long as the clock pulses are applied. Consider initially first flip-flop is set, and others are reset. After 3 clock pulses, again we will get initial state of 100. So this is a MOD-3 counter.

Clk	Q ₂	Q ₁	Q ₀
0	1	0	0
1	0	1	0
2	0	0	1
3	1	0	0
4	0	1	0



A ring counter with N flip-flops can count up to N states, i.e., MOD-N counter, whereas, N-bit asynchronous counter can count up to 2^N states. So, ring counter is uneconomical compared to a ripple counter, but has the advantage of requiring no decoder. Since it is entirely synchronous operation and requires no gates for flip-flop inputs, it has further advantage of being very fast.

Twisted ring counter (Johnson counter): This counter is obtained from a SISO shift register by connecting the complement of serial output to serial input as shown in below figure.

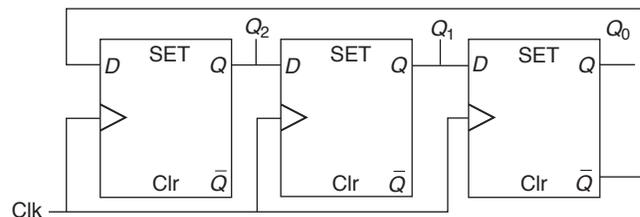
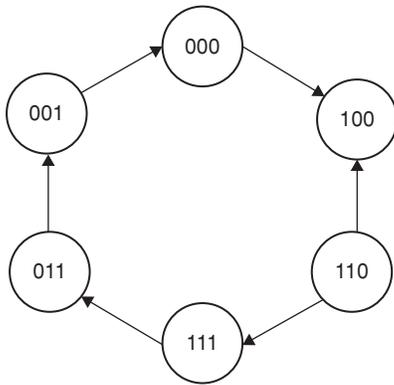


Figure 20 Twisted Ring Counter

Let initially all the FFs be reset, after each clock pulse the complement of last bit will appear as at MSB, and other bits shift right side by 1-bit. After 6 clock pulses the register will come to initial state 000. Similarly, the 3-bit Johnson counter will oscillate between the states 101, 010.

Clk	Q ₂	Q ₁	Q ₀
0	0	0	0
1	1	0	0
2	1	1	0
3	1	1	1
4	0	1	1
5	0	0	1
6	0	0	0



An n -bit Johnson counter can have $2n$ unique states and can count up to $2n$ pulses, so it is a MOD- $2n$ counter. It is more economical than basic ring counter but less economical than ripple counter.

Solved Examples

Example 1: Assume that 4-bit counter is holding the count 0101. What will be the count after 27 clock pulses?

Solution: Total clock pulses: $27 = 16 + 11$
 $0101 + 1011 = 0000$

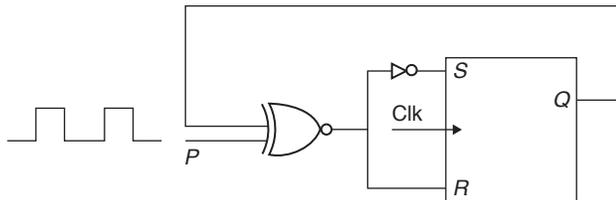
Example 2: A MOD-2 counter followed by MOD-5 counter is

Solution: A decade counter, counts 10 states (5×2).

Example 3: A 4-bit binary ripple counter uses flip-flops with propagation delay time of 25 msec each. The maximum possible time required for change of state will be

Solution: The maximum time = $4 \times 25 \text{ ms} = 100 \text{ ms}$

Example 4: Consider the circuit, the next state Q^+ is



Solution

P	Q	S	R	Q ⁺
0	0	0	1	0
0	1	1	0	1
1	0	1	0	1
1	1	0	1	0

So, $Q^+ = P \oplus Q$

Example 5: A certain JK FF has $t_{pd} = 12 \text{ n}$ sec what is the largest MOD counter, that can be constructed from these FF and still operate up to 10 MHz?

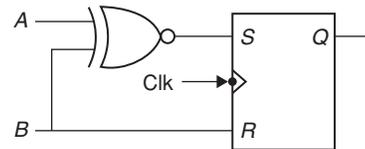
Solution: $N \leq \frac{1}{f_{\max} \cdot t_{pd}}$

$f_{\max} = 10 \text{ MHz}$ $N \leq 8$
 $t_{pd} = 12 \text{ ns}$

$$N \leq \frac{1}{10 \times 10^6 \times 12 \times 10^{-9}}$$

MOD counter is = $2^N = 2^8 = 256$

Example 6: An AB flip-flop is constructed from an SR flip-flop as shown below. The expression for next state Q^+ is

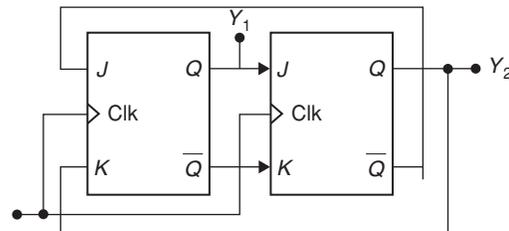


Solution:

A	B	Q	S	R	Q ⁺
0	0	0	1	0	1
0	0	1	1	0	1
0	1	0	0	1	0
0	1	1	0	1	0
1	0	0	0	0	0
1	0	1	0	0	1
1	1	0	1	1	×
1	1	1	1	1	×

$\therefore Q^+ = \overline{A}B + AQ = \overline{A}B + \overline{B}Q$

Example 7: In the circuit shown below, the output y_1 and y_2 for the given initial condition $y_1 = y_2 = 1$ and after four input pulses will be



Solution:

- After 1st pulse $y_1 = 0, y_2 = 1$
- After 2nd pulse $y_1 = 0, y_2 = 0$
- After 3rd pulse $y_1 = 1, y_2 = 0$
- After 4th pulse $y_1 = 1, y_2 = 1$

Example 8: A ripple counter is to operate at a frequency of 10 MHz. If the propagation delay time of each flip-flop in the counter is 10 ns and the storing time is 50 ns, how many maximum stages can the counter have?

Solution: $nt_{pd} + t_s \leq \frac{1}{f}$

where, n = number of stages

t_{pd} = propagation delay time

t_s = strobing time

f = frequency of operation = $10 \times 10^{-9}n + 50 \times 10^{-9}$

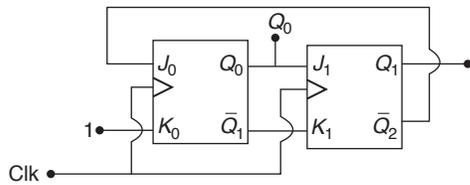
$$\leq \frac{1}{10 \times 10^6}$$

(or) $10n + 50 \leq 100$

(or) $10n \leq 50$

For max stages $n = \frac{50}{10} = 5$

Example 9: In the circuit assuming initially $Q_0 = Q_1 = 0$. Then the states of Q_0 and Q_1 immediately after the 33rd pulse are

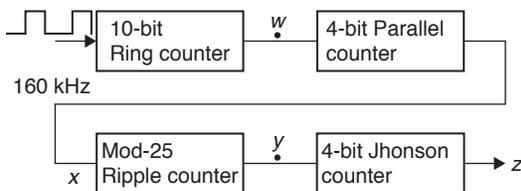


Solution:

J_0	K_0	J_1	K_1	Q_0	Q_1	Count
1	1	0	1	0	0	Initial
1	1	1	0	1	0	1st pulse
0	1	0	1	0	1	2nd
1	1	0	1	0	0	3rd
1	1	1	0	1	0	4th
0	1	0	1	0	1	5th pulse

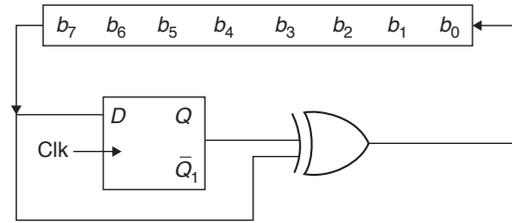
After 4th pulse, output is same as after 1st one, so, sequence gets repeated. So output after 33rd pulse would be same as after 3rd pulse. i.e., (00).

Example 10: The frequency of the pulse at z in the network shown in figure is



Solution: 10-bit ring counter is a MOD-10. So, it divides the 160 kHz input by 10. Therefore, $w = 16$ kHz. The 4-bit parallel counter is a MOD-16. Thus, the frequency at $x = 1$ kHz. The MOD-25 ripple counter produces a frequency at $y = 40$ Hz ($1 \text{ kHz}/25 = 40 \text{ Hz}$). The 4-bit Johnson counter is a MOD-8. The frequency at $Z = 5$ Hz.

Example 11: The 8-bit shift left shift register, and D flip-flop shown in the figure is synchronized with the same clock. The D flip-flop is initially cleared. The circuit acts as



Solution: The output of XOR gate is $Z = b_{i+1} \oplus b_i$ and this output shift the register to left. Initially, $Z = 0$

After 2nd clock $Z = b_7 \oplus 0 = b_7$

After 2nd clock $Z = b_7 \oplus b_6$

3rd clock $Z = b_6 \oplus b_5$

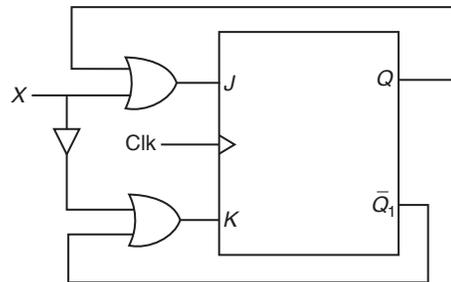
4th clock $Z = b_5 \oplus b_4$

It is a binary to gray code converter.

Example 12: A 4-bit MOD-16 ripple counter uses JK flip-flops. If the propagation delay of each flip-flop is 50 ns sec, the maximum clock frequency that can be used is equal to

Solution: Max = clock frequency = $\frac{1}{4 \times 50 \times 10^{-9}} = 5 \text{ MHz}$

Example 13: What is the state diagram for the sequential circuit shown?

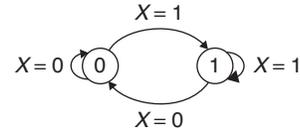


- (A)
- (B)
- (C)
- (D)

Solution: (D)

State diagram of a sequential circuit will have states (output) of all the flip-flops.

Present state	Next state		Q_{n+1}
Q_n	For $x = 0$	For $x = 1$	
0	0	1	
1	0	1	

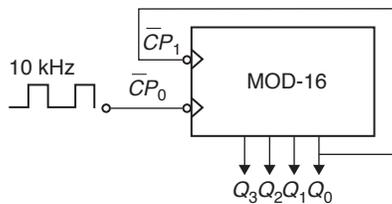


EXERCISES

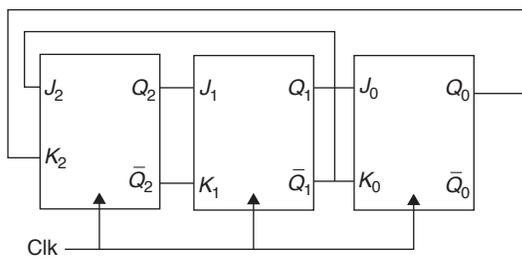
Practice Problems I

Directions for questions 1 to 22: Select the correct alternative from the given choices.

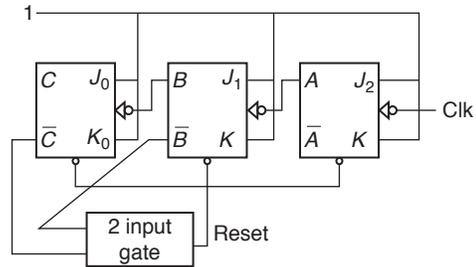
- How many flip-flops are needed for MOD-16 ring counter and MOD-16 Johnson counter?
 - (A) 16, 16
 - (B) 16, 8
 - (C) 4, 3
 - (D) 4, 4
- A 2-bit synchronous counter uses flip-flops with propagation delay time of 25 n sec, each. The maximum possible time required for change of state will be
 - (A) 25 n sec
 - (B) 50 n sec
 - (C) 75 n sec
 - (D) 100 n sec
- For given MOD-16 counter with a 10 kHz clock input determine the frequency at Q_3



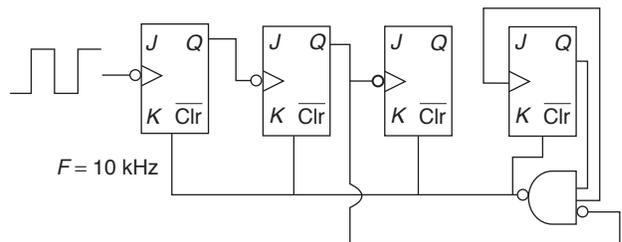
- (A) 625 Hz
 - (B) 10 kHz
 - (C) 2.5 kHz
 - (D) 0 Hz
- A 4-bit ripple counter and a 4-bit synchronous counter are made using flip-flops having a propagation delay of 10 n sec each. If the worst case delay in the ripple counter and the synchronous counter be R and S , respectively, then
 - (A) $R = 10$ ns, $S = 40$ ns
 - (B) $R = 40$ ns, $S = 10$ ns
 - (C) $R = 10$ ns, $S = 30$ ns
 - (D) $R = 30$ ns, $S = 10$ ns
 - The counter shown in the figure has initially $Q_2Q_1Q_0 = 000$. The status of $Q_2Q_1Q_0$ after the first pulse is



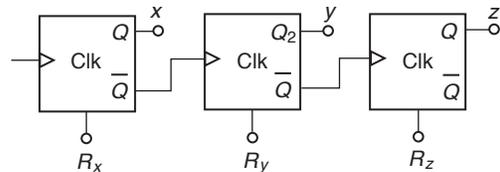
- (A) 001
 - (B) 010
 - (C) 100
 - (D) 101
- 12 MHz clock frequency is applied to a cascaded counter of MOD-3 counter, MOD-4 counter and MOD-5 counter. The lowest output frequency is
 - (A) 200 kHz
 - (B) 1 MHz
 - (C) 3 MHz
 - (D) 4 MHz
 - In the modulo-6 ripple counter shown in the figure below, the output of the 2-input gate is used to clear the JK flip-flops. The 2-input gate is

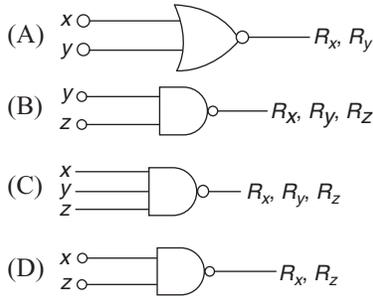


- (A) a NAND gate
 - (B) a NOR gate
 - (C) an OR gate
 - (D) an AND gate
- In figure, J and K inputs of all the 4 flip-flops are made high, the frequency of the signal at output y is

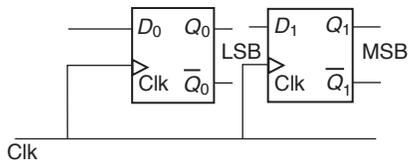


- (A) 0.833 kHz
 - (B) 1.0 kHz
 - (C) 0.91 kHz
 - (D) 0.77 kHz
- In a number system a counter has to recycle to 0 at the sixth count. Which of the connections indicated below will realize this resetting? (a logic '0' at the R inputs resets the counters)

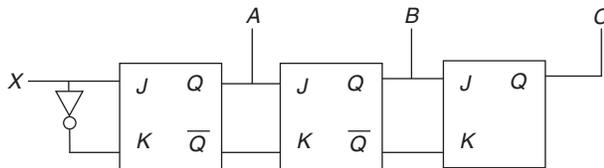




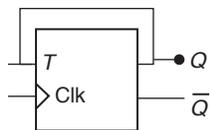
10. Two *D* flip-flops, as shown below, are to be connected as a synchronous counter that goes through the following Q_1Q_0 sequence $00 \rightarrow 01 \rightarrow 11 \rightarrow 10 \rightarrow 00$. The inputs D_0 and D_1 respectively should be connected as



- (A) \bar{Q}_0 and Q_1 (B) \bar{Q}_0 and Q_1
 (C) $\bar{Q}_1\bar{Q}_0$ and \bar{Q}_1Q_0 (D) $\bar{Q}_1\bar{Q}_0$ and \bar{Q}_1Q_0
11. N flip-flops can be used to divide the input clock frequency by
 (A) N (B) $2N$
 (C) 2^N (D) 2^{N-1}
12. For a shift register as shown, $x = 1011$, with initially FF cleared, ABC will have value of after 3 clock pulses



- (A) 101 (B) 011
 (C) 001 (D) 111
13. If a FF is connected as shown what will be the output? (initially $Q = 0$)



- (A) 11111 (B) 0000
 (C) 1010 (D) 0101
14. The excitation table for a FF whose output conditions are if $AB = 00$, no change of state occurs
 $AB = 01$, FF becomes 1 with next clock pulse
 $AB = 10$, FF becomes 0 with next clock pulse
 $AB = 11$, FF changes its state

(A)

Q_n	Q_{n+1}	A	B
0	0	0	x
0	1	1	x
1	0	x	1
1	1	x	0

(B)

Q_n	Q_{n+1}	A	B
0	0	1	x
0	1	0	x
1	0	x	0
1	1	x	1

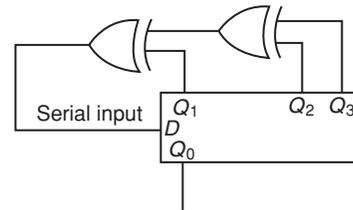
(C)

Q_n	Q_{n+1}	A	B
0	0	x	0
0	1	x	1
1	0	x	x
1	1	0	x

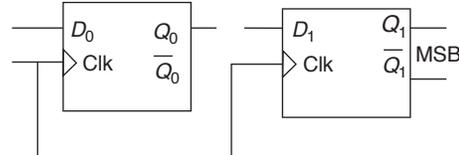
(D)

Q_n	Q_{n+1}	A	B
0	0	x	0
0	1	1	x
1	0	x	1
1	1	0	x

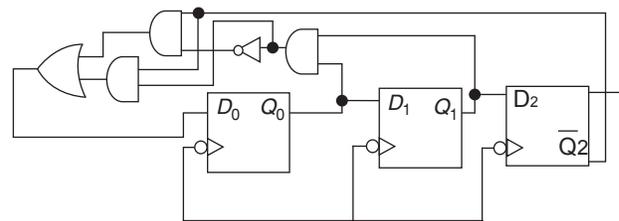
15. A shift register that shift the bits 1 position to the right at each clock pulse is initialized to 1100 (Q_0, Q_1, Q_2, Q_3). The outputs are combined using an XOR gate circuit and fed to the *D* input. After which clock pulse, will the initial pattern reappear at the output?



- (A) 6th (B) 5th
 (C) 4th (D) 7th
16. If we need to design a synchronous counter that goes through the states $00 \rightarrow 01 \rightarrow 11 \rightarrow 10 \rightarrow 00$ using *D* FF, what should be the input to the FF?



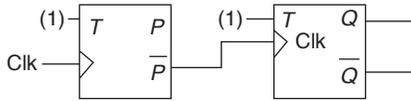
- (A) $D_0 = Q_0, D_1 = \bar{Q}_1$
 (B) $D_0 = \bar{Q}_1, D_1 = Q_0$
 (C) $D_0 = \bar{Q}_1 \cdot Q_0, D_1 = \bar{Q}_1 \bar{Q}_0$
 (D) $D_0 = \bar{Q}_0, D_1 = Q_1$
17. Find the counter state sequence (Assume Q_0 as MSB).



- (A) 4, 6, 7, 3, 1, 0, 4 (B) 4, 6, 5, 3, 1, 0, 4
 (C) 4, 5, 6, 7, 0, 4, 5 (D) 4, 6, 7, 1, 0, 4

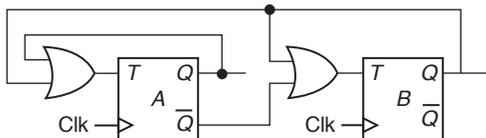
18. If the propagation delay of each FF is 50 ns, and for the AND gate to be 20 ns. What will be the f_{max} for MOD-32 ripple and synchronous counters?
 (A) 14.3 MHz, 4 MHz (B) 14.3 MHz, 5 MHz
 (C) 5 MHz, 14.3 MHz (D) 3.7 MHz, 14.3 MHz

19. For a given counter identify its behavior



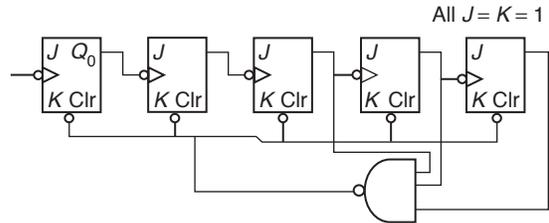
The output is taken from PQ.

- (A) MOD-4 up counter
 (B) MOD-2 down counter
 (C) MOD-4 down counter
 (D) MOD-2 up counter
20. A circuit using T FF is given. Identify the circuit.

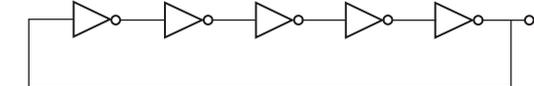


- (A) MOD-2 counter
 (B) MOD-4 counter
 (C) MOD-3 counter
 (D) MOD-2 generate 00, 10, 00

21. The MOD number of asynchronous counter shown



- (A) 24 (B) 48
 (C) 29 (D) 28
22. For the oscillator, find the fundamental frequency if propagation delay of each inverter is 1000 psec.



- (A) 100 MHz (B) 10 MHz
 (C) 1 GHz (D) 10 GHz

Practice Problems 2

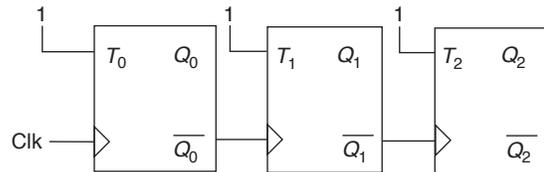
Directions for questions 1 to 30: Select the correct alternative from the given choices.

1. Match List 1 (operation) with List 2 (associated device) and select the correct answer using the codes given below:

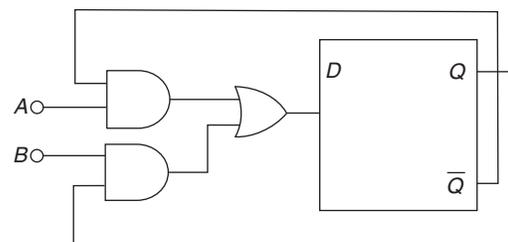
List 1	List 2
(a) Frequency Ddivision	(1) ROM
(b) Decoding	(2) Multiplexer
(c) Data selection	(3) Demultiplexer
(d) Code conversion	(4) Counter

- (A) a-3, b-4, c-2, d-1
 (B) a-3, b-4, c-1, d-2
 (C) a-4, b-3, c-1, d-2
 (D) a-4, b-3, c-2, d-1
2. A MOD-5 synchronous counter is designed by using JK flip-flop, the number of counts skipped by it will be
 (A) 2 (B) 3
 (C) 5 (D) 0
3. A counter starts off in the 0000 state, then clock pulses are applied. Some time later the clock pulses are removed and the counter flip-flops read 0011. How many clock pulses have occurred?
 (A) 3 (B) 35
 (C) 51 (D) Any of these

4. Figure below shown as ripple counter using positive edge triggered flip-flops. If the present state of the counters is $Q_2 Q_1 Q_0 = 011$, then its next state ($Q_2 Q_1 Q_0$) will be



- (A) 010 (B) 100
 (C) 111 (D) 101
5. A synchronous sequential circuit is designed to detect a bit sequence 0101 (overlapping sequence include). Every time, this sequence is detected, the circuit produces output of 1. What is the minimum number of states the circuit must have?
 (A) 4 (B) 5
 (C) 6 (D) 7
6. What is represented by digital circuit given below?



- (A) An SR flip-flop with $A = S$ and $B = R$
- (B) A JK flip-flop with $A = K$ and $B = J$
- (C) A JK flip-flop with $A = J$ and $B = \bar{K}$
- (D) An SR flip-flop with $A = R$ and $B = S$

7. In a ripple counter, the state whose output has a frequency equal to $\frac{1}{8}$ th that of clock signal applied to the first stage, also has an output periodically equal to $\frac{1}{8}$ th that of the output signal obtained from the last stage. The counter is
- (A) MOD-8
 - (B) MOD-6
 - (C) MOD-64
 - (D) MOD-16

8. A flip-flop is popularly known as
- (A) Astable multivibrator
 - (B) Bistable multivibrator
 - (C) Monostable multivibrator
 - (D) None of these

9. Which of the following represents the truth table for JK flip-flop?

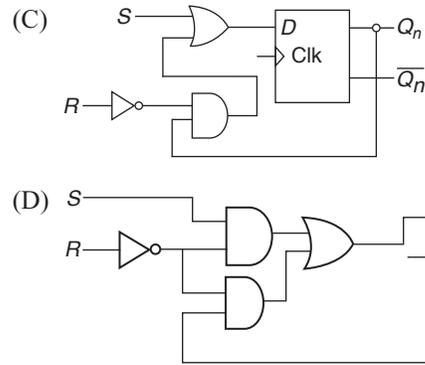
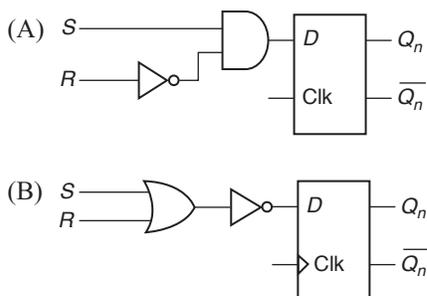
(A)	J	K	Output
	0	0	Q_0
	0	1	0
	1	0	1
	1	1	\bar{Q}_0

(B)	J	K	Output
	0	0	\bar{Q}_0
	0	1	0
	1	0	1
	1	1	Q_0

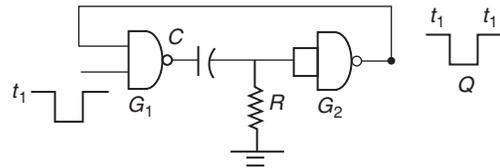
(C)	J	K	Output
	0	0	Q_0
	0	1	0
	1	0	1
	1	1	Invalid

(D)	J	K	Output
	0	0	1
	0	1	0
	1	0	1
	1	1	0

10. One disadvantage of master-slave FF is
- (A) setup time becomes longer.
 - (B) it requires input to be held constant before clock transition.
 - (C) unpredictable output even if input held constant.
 - (D) hold time becomes longer.
11. Which of the following converts D FF to SR FF?



12. Which of the circuit is being represented by the figure?



- (A) NAND gate
- (B) Monostable multivibrator
- (C) Astable multivibrator
- (D) Schmitt trigger

13. Hold time is
- (A) time for which output is held constant.
 - (B) time for which clock is to be held constant on applying input.
 - (C) time for which input should be maintained constant after the triggering edge of clock pulse.
 - (D) time for which input should be maintained constant prior to the arrival of triggering edge of clock pulse.

14. Shift registers are made up of
- (A) MOS inverters
 - (B) FF
 - (C) Latches
 - (D) None of these

15. Data from a satellite is received in serial form. If the data is coming at 8 MHz rate, how long will it take to serially load a word in 40-bit shift register?
- (A) $1.6 \mu\text{s}$
 - (B) $5 \mu\text{s}$
 - (C) $6.4 \mu\text{s}$
 - (D) $12.8 \mu\text{s}$

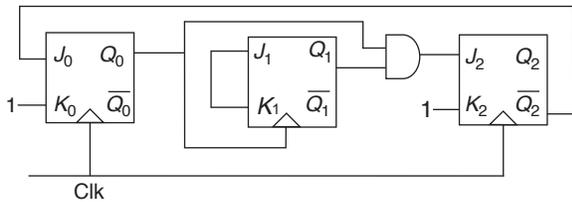
16. A JK FF can be converted into T FF by connecting
- (A) \bar{Q} to 0
 - (B) 0 to \bar{Q}
 - (C) 0 to Q
 - (D) by connecting both J and K inputs to T

17. The flip-flop that is not affected by race around condition
- (A) T FF
 - (B) JK FF
 - (C) SR FF
 - (D) None of these

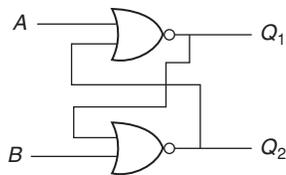
18. The characteristic equation of JK FF is
- (A) $JQ(t) + KQ'(t)$
 - (B) $J'Q(t) + KQ(t)$
 - (C) $JQ'(t) + K'Q(t)$
 - (D) None of these

19. For a D.FF input, the \bar{Q} is connected. What would be the output sequence?
 (A) 0000 (B) 1111
 (C) 010101 (D) 101010
20. In order to implement a MOD-6 synchronous counter we have 3 FF and a combination of 2 input gate(s). Identify the combination circuit.
 (A) One AND gate
 (B) One OR gate
 (C) One AND and one OR gate
 (D) Two AND gates

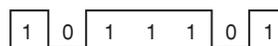
21. Given a MOD-5 counter. The valid states for the counter are (0, 1, 2, 3, 4). The propagation delay of each FF is T_F and that of AND gate is t_A . The maximum rate at which counter will operate satisfactorily



- (A) $\frac{1}{t_F + t_A}$ (B) $\frac{1}{3t_F}$
 (C) $\frac{1}{2t_F + t_A}$ (D) $\frac{1}{3t_F - t_A}$
22. For a NOR latch as shown up A and B are made first (0, 1) and after a few seconds it is made (1, 1). The corresponding output (Q_1, Q_2) are

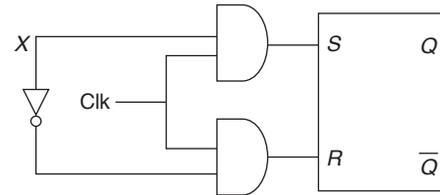


- (A) first (1, 0) then (0, 0)
 (B) first (1, 0) then (1, 0)
 (C) first (1, 0) then (1, 1)
 (D) first (1, 0) then (0, 1)
23. In order to design a pulse generator to generate the wave form using a shift register, what is the number of FF required?

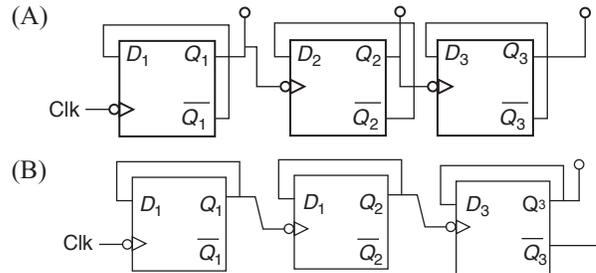


- (A) 3 (B) 4
 (C) 5 (D) 6
24. For what minimum value of propagation delay in each FF will a 10-bit ripple counter skip a count when it is clocked at 5 MHz?
 (A) 10 ns (B) 20 ns
 (C) 25 ns (D) 15 ns

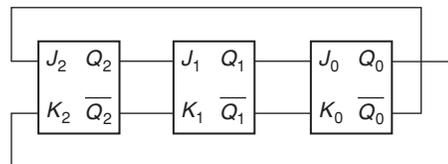
25. A divide by 50 counter can be realized by using
 (A) 5 no. of MOD-10 counter
 (B) 10 no. of MOD-5 counter
 (C) One MOD-5 counter followed by one MOD-10 counter
 (D) 10 no. of MOD-10 counter
26. The following latch is



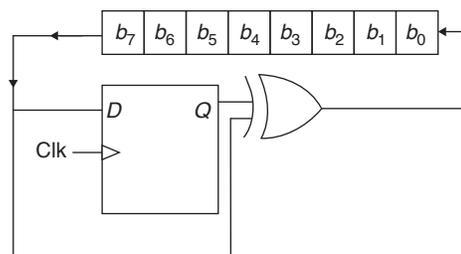
- (A) D latch (B) T latch
 (C) JK latch (D) RS latch
27. Which of the following represent a 3-bit ripple counter using D FF?



- (A) Both (A) and (B)
 (D) None of these
28. For the Johnson counter with initial Q_2, Q_1, Q_0 as 101, the frequency of the output is (Q_2, Q_1, Q_0)



- (A) $\frac{f_c}{8}$ (B) $\frac{f_c}{6}$
 (C) $\frac{f_c}{2}$ (D) $\frac{f_c}{4}$
29. For the given circuit the contents of register ($b_7 - b_0$) are 10010101, what will be the register contents after 8 clock pulses?



- (A) 10010101 (B) 01101010
 (C) 11011111 (D) 01101011

30. A latch is to be build with A and B as input. From the table find the expression for the next state Q^+

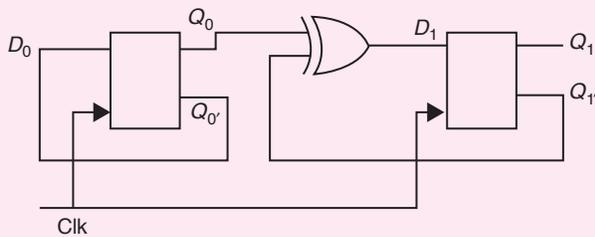
A	B	Q	Q^+
0	0	0	0
0	0	1	0
0	1	0	0
0	1	1	0

1	0	0	1
1	0	1	1
1	1	0	1
1	1	1	1

- (A) A
 (B) B
 (C) $A + \bar{B}$
 (D) $A\bar{B} + AB$

PREVIOUS YEARS' QUESTIONS

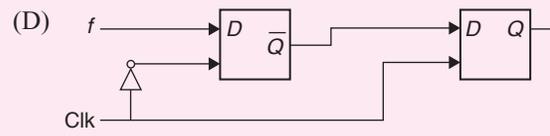
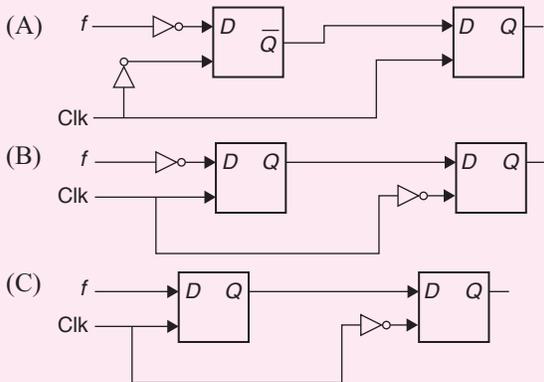
1. Consider the following circuit.



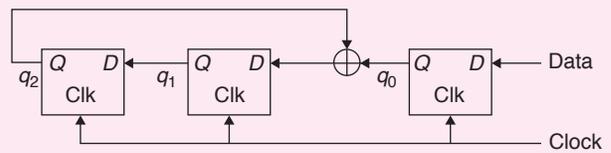
The flip-flops are positive edge triggered D FFs. Each state is designated as a 2-bit string Q_0Q_1 . Let the initial state be 00. The state transition sequence is [2005]

- (A) $00 \rightarrow 11 \rightarrow 01$
 (B) $00 \rightarrow 11$
 (C) $00 \rightarrow 11 \rightarrow 01 \rightarrow 11$
 (D) $00 \rightarrow 11 \rightarrow 01 \rightarrow 10$

2. You are given a free running clock with a duty cycle of 50% and a digital waveform f which changes only at the negative edge of the clock. Which one of the following circuits (using clocked D flip-flops) will delay the phase of f by 180° ? [2006]



3. Consider the circuit in the diagram. The \oplus operator represents Ex-OR. The D flip-flops are initialized to zeros (cleared). [2006]



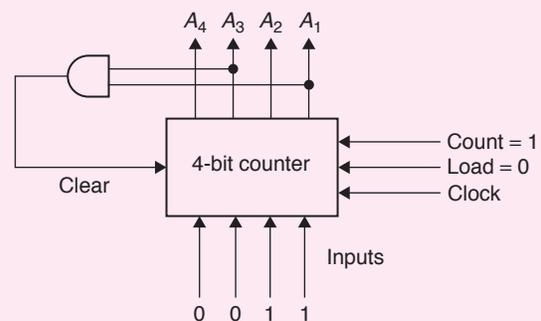
The following data: 100110000 is supplied to the data terminal in 9 clock cycles. After that the values of $q_2q_1q_0$ are

- (A) 000 (B) 001
 (C) 010 (D) 101

4. The control signal functions of a 4-bit binary counter are given below (where X is 'don't care'):

Clear	Clock	Load	Count	Function
1	X	X	X	Clear to 0
0	X	0	0	No change
0	\uparrow	1	X	Load input
0	\uparrow	0	1	Count next

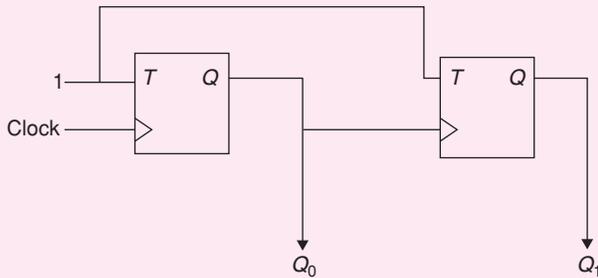
The counter is connected as follows:



Assume that the counter and gate delays are negligible. If the counter starts at 0, then it cycles through the following sequence: [2007]

- (A) 0, 3, 4 (B) 0, 3, 4, 5
 (C) 0, 1, 2, 3, 4 (D) 0, 1, 2, 3, 4, 5

5. In the sequential circuit shown below, if the initial value of the output Q_1Q_0 is 00, what are the next four values of Q_1Q_0 ? [2010]

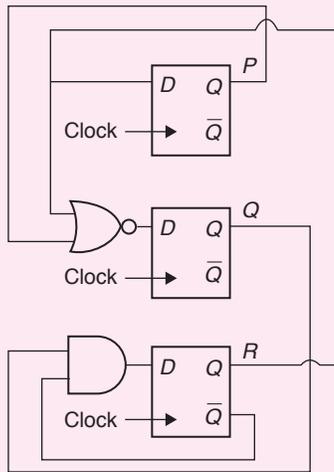


- (A) 11, 10, 01, 00 (B) 10, 11, 01, 00
 (C) 10, 00, 01, 11 (D) 11, 10, 00, 01

6. The minimum number of *D* flip-flops needed to design a MOD-258 counter is [2011]

- (A) 9 (B) 8
 (C) 512 (D) 258

Common Data for Questions 7 and 8: Consider the following circuit involving three D-type flip-flops used in a certain type of counter configuration.



7. If all the flip-flops were reset to 0 at power on, what is the total number of distinct outputs (states) represented by PQR generated by the counter? [2011]

- (A) 3 (B) 4
 (C) 5 (D) 6

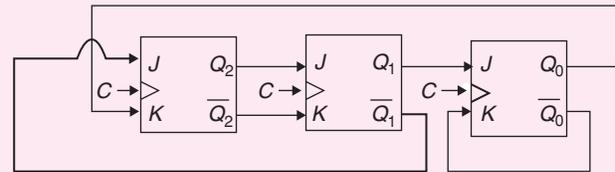
8. If at some instance prior to the occurrence of the clock edge, *P*, *Q*, and *R* have a value 0, 1, and 0, respectively, what shall be the value of PQR after the clock edge? [2011]

- (A) 000 (B) 001
 (C) 010 (D) 011

9. Let $K = 2^n$. A circuit is built by giving the output of an *n*-bit binary counter as input to an *n*-to- 2^n -bit decoder. This circuit is equivalent to a [2014]

- (A) *K*-bit binary up counter
 (B) *K*-bit binary down counter
 (C) *K*-bit ring counter
 (D) *K*-bit Johnson counter

10.



The above synchronous sequential circuit built using *JK* flip-flops is initialized with $Q_2Q_1Q_0 = 000$. The state sequence for this circuit for the next 3 clock cycles is [2014]

- (A) 001, 010, 011 (B) 111, 110, 101
 (C) 100, 110, 111 (D) 100, 011, 001

11. Consider a 4-bit Johnson counter with an initial value of 0000. The counting sequence of this counter is [2015]

- (A) 0, 1, 3, 7, 15, 14, 12, 8, 0
 (B) 0, 1, 3, 5, 7, 9, 11, 13, 15, 0
 (C) 0, 2, 4, 6, 8, 10, 12, 14, 0
 (D) 0, 8, 12, 14, 15, 7, 3, 1, 0

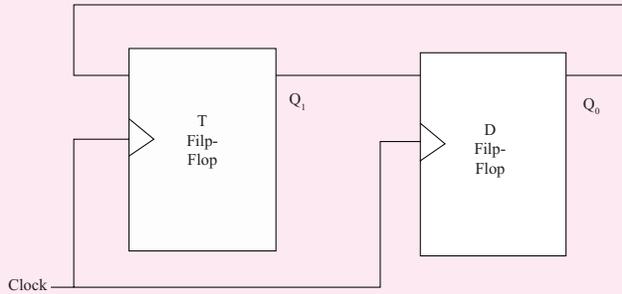
12. A positive edge-triggered D flip-flop is connected to a positive edge-triggered JK flip-flop as follows. The *Q* output of the D flip-flop is connected to both the *J* and *K* inputs of the JK flip-flop, while the *Q* output of the JK flip-flop is connected to the input of the D flip-flop. Initially, the output of the D flip-flop is set to logic one and the output of the JK flip-flop is cleared. Which one of the following is the bit sequence (including the initial state) generated at the *Q* output of the JK flip-flop when the flip-flops are connected to a free-running common clock? Assume that $J = K = 1$ is the toggle mode and $J = K = 0$ is the state-holding mode of the JK flip-flop. Both the flip-flops have non-zero propagation delays. [2015]

- (A) 0110110... (B) 0100100...
 (C) 011101110... (D) 011001100...

13. The minimum number of JK flip-flops required to construct a synchronous counter with the count sequence (0, 0, 1, 1, 2, 2, 3, 3, 0, 0, ...) is [2015]

14. We want to design a synchronous counter that counts the sequence 0-1-0-2-0-3 and then repeats. The minimum number of *J-K* flip-flops required to implement this counter is _____. [2016]

15. Consider a combination of T and D flip-flops connected as shown below. The output of the D flip-flop is connected to the input of the T flip-flop and the output of the T flip-flop is connected to the input of the D flip-flop.



Initially, both Q_0 and Q_1 are set to 1 (before the 1st clock cycle). The outputs [2017]

- (A) $Q_1 Q_0$ after the 3rd cycle are 11 and after the 4th cycle are 00 respectively
- (B) $Q_1 Q_0$ after the 3rd cycle are 11 and after the 4th cycle are 01 respectively
- (C) $Q_1 Q_0$ after the 3rd cycle are 00 and after the 4th cycle are 11 respectively
- (D) $Q_1 Q_0$ after the 3rd cycle are 01 and after the 4th cycle are 01 respectively

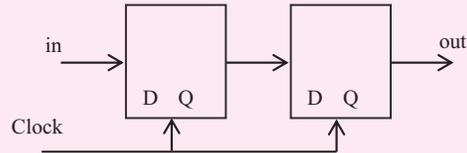
16. The next state table of a 2-bit saturating up-counter is given below.

Q_1	Q_0	Q_1^+	Q_0^+
0	0	0	1
0	1	1	0
1	0	1	1
1	1	1	1

The counter is built as a synchronous sequential circuit using T flip-flops. The expressions for T_1 and T_0 are [2017]

- (A) $T_1 = Q_1 Q_0$, $T_0 = \bar{Q}_1 \bar{Q}_0$
- (B) $T_1 = \bar{Q}_1 Q_0$, $T_0 = \bar{Q}_1 + \bar{Q}_0$
- (C) $T_1 = Q_1 + Q_0$, $T_0 = \bar{Q}_1 + \bar{Q}_0$
- (D) $T_1 = \bar{Q}_1 Q_0$, $T_0 = Q_1 + Q_0$

17. Consider the sequential circuit shown in the figure, where both flip-flops used are positive edge-triggered D flip-flops.



The number of states in the state transition diagram of this circuit that have a transition back to the same state on some value of “in” is _____. [2018]

ANSWER KEYS

EXERCISES

Practice Problems 1

- | | | | | | | | | | |
|-------|-------|-------|-------|-------|-------|-------|-------|-------|-------|
| 1. B | 2. A | 3. A | 4. B | 5. C | 6. A | 7. C | 8. A | 9. B | 10. A |
| 11. C | 12. B | 13. B | 14. C | 15. D | 16. B | 17. A | 18. D | 19. A | 20. C |
| 21. D | 22. A | | | | | | | | |

Practice Problems 2

- | | | | | | | | | | |
|-------|-------|-------|-------|-------|-------|-------|-------|-------|-------|
| 1. D | 2. B | 3. D | 4. B | 5. A | 6. C | 7. C | 8. B | 9. A | 10. B |
| 11. C | 12. B | 13. C | 14. B | 15. B | 16. D | 17. C | 18. C | 19. C | 20. D |
| 21. C | 22. A | 23. D | 24. B | 25. C | 26. A | 27. A | 28. C | 29. C | 30. A |

Previous Years' Questions

- | | | | | | | | | | |
|-------|-------|------------------|-------|-------|-------|-------|------|------|-------|
| 1. D | 2. C | 3. C | 4. C | 5. A | 6. A | 7. B | 8. D | 9. C | 10. C |
| 11. D | 12. A | 13. $3(8 = 2^3)$ | 14. 3 | 15. B | 16. B | 17. 2 | | | |

TEST

DIGITAL LOGIC

Time: 60 min.

Directions for questions 1 to 30: Select the correct alternative from the given choices.

- What is the range of signed decimal numbers that can be represented by 4-bit 1's complement notation?
 (A) -7 to +7 (B) -16 to +16
 (C) -7 to +8 (D) -15 to +16
- Which of the following signed representation have a unique representation of 0?
 (A) Sign-magnitude (B) 1's complement
 (C) 0's complement (D) 2's complement
- Find the odd one out among the following
 (A) EBCDIC (B) GRAY
 (C) Hamming (D) ASCII
- Gray code for number 8 is
 (A) 1100 (B) 1111
 (C) 1000 (D) 1101
- Find the equivalent logical expression for $z = x + \bar{x}y$
 (A) $z = x\bar{y}$ (B) $Z = \bar{x}y$
 (C) $Z = \bar{x} + y$ (D) $Z = x + y$
- The number of distinct Boolean expression of 3 variables is
 (A) 256 (B) 16
 (C) 1024 (D) 65536
- The Boolean expression for the truth table shown is

X	Y	Z	F
0	0	0	0
0	0	1	0
0	1	0	0
0	1	1	1
1	0	0	0
1	0	1	0
1	1	0	1
1	1	1	0

- (A) $Y(X+Z)(\bar{X}+\bar{Z})$ (B) $Y(X+\bar{Z})(\bar{x}+Z)$
 (C) $\bar{Y}(X+\bar{Z})(\bar{x}+Z)$ (D) $\bar{Y}(X+Z)(\bar{X}+Z)$
- The number of essential prime implicants for the Boolean functions shown in the given K-map.

WZ	XY			
	00	01	11	10
00	1	1	0	1
01	1	0	0	1
11	1	0	0	0
10	1	0	0	1

- (A) 4 (B) 5
 (C) 6 (D) 8

- The number of product terms in the minimized SOP from is

1	0	0	1
0	D	0	0
0	0	D	1
1	0	0	1

- (A) 2 (B) 4
 (C) 5 (D) 3
- The minimum number of 2 input NAND gates needed to implement $Z = XY + VW$ is
 (A) 2 (B) 3
 (C) 4 (D) 5
 - The operation $\bar{a} \oplus \bar{b}$ represents
 (A) $ab + \bar{a}\bar{b}$ (B) $\bar{a}b + a\bar{b}$
 (C) $a\bar{b} + \bar{a}b$ (D) $a - \bar{b}$
 - Find the dual of $X + [Y + XZ] + U$
 (A) $X + [Y(X+Z)] + U$ (B) $X(Y+XZ)U$
 (C) $X + [Y(X+Z)]U$ (D) $X[Y(X+Z)]U$
 - The simplified form of given function $AB + BC + \bar{A}\bar{C}$ is equal to
 (A) $AB + \bar{A}\bar{C}$ (B) $\bar{A}\bar{C} + BC$
 (C) $\bar{A}\bar{C} + BC$ (D) $\bar{A}\bar{B} + \bar{A}\bar{C}$
 - Simplify the following

WX	YZ			
	1	1	0	1
1	1	1	0	1
0	0	0	1	1
0	0	0	0	0

- (A) $\bar{W}\bar{Y} + \bar{W}\bar{Z} + WXY$
 (B) $\bar{W}\bar{X} + \bar{W}\bar{Z} + WXY$
 (C) $WY + WYZ + WXY + XY\bar{Z}$
 (D) $\bar{W}\bar{X} + \bar{Y}\bar{Z} + \bar{W}\bar{Z}$
- Simplify the following
 $F = ABCD + \bar{A}\bar{B}CD + \bar{A}C\bar{B}D + \bar{A}BCD$
 (A) CD (B) BC
 (C) AB (D) $\bar{C} + \bar{D}$
 - Find the equivalent Boolean expression for $AC + B\bar{C}$
 (A) $\bar{A}\bar{C} + B\bar{C} + AC$
 (B) $ABC + \bar{A}\bar{B}\bar{C} + ABC\bar{C} + \bar{A}\bar{B}\bar{C}$
 (C) $ABC + \bar{A}\bar{B}\bar{C} + ABC\bar{C} + \bar{A}\bar{B}\bar{C}$
 (D) $\bar{A}\bar{C} + B\bar{C} + \bar{A}\bar{C}$

17. Simplify the following expression

$$\bar{A}\bar{B}\bar{C} + \bar{A}B\bar{C} + \bar{A}BC + ABC$$

- (A) $\bar{A}\bar{C} + B\bar{C} + \bar{A}B$ (B) $A\bar{C} + B\bar{C} + \bar{A}B$
 (C) $\bar{A}\bar{C} + \bar{B}C + \bar{A}B$ (D) $\bar{A}\bar{C} + \bar{B}\bar{C} + \bar{A}B$

18. If $A = 1$ in the logic equation $[A + C\{\bar{B} + (\bar{C} + A\bar{B})\}]$
 $[\bar{A} + \bar{C}(A + B)] = 1$, then

- (A) $B = C$ (B) $B = \bar{C}$
 (C) $C = 1$ (D) $C = 0$

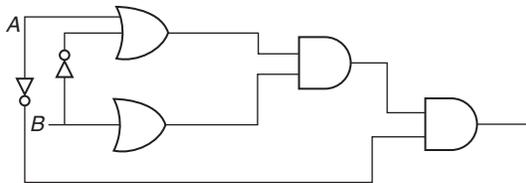
19. Which is the odd function with 3 Boolean variables in it

- (A) $\sum(0, 3, 5, 6)$ (B) $\sum(0, 2, 4, 6)$
 (C) $\sum(1, 2, 4, 7)$ (D) $\sum(1, 3, 5, 7)$

20. Which of the following expressions is/are incorrect?

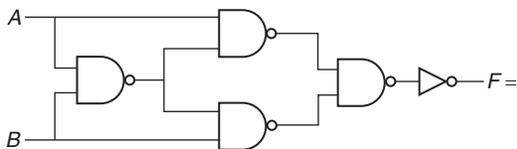
- (A) $\overline{a+b} = \bar{a}\bar{b}$ (B) $\overline{\bar{a}+b} = \bar{a}\bar{b}$
 (C) $\overline{\bar{a}\bar{b}} = \bar{a} + \bar{b}$ (D) $\overline{a+b} = \bar{a}\bar{b}$

21. The simplified form of logic circuit is



- (A) $A + B$ (B) $\bar{A}\bar{B}$
 (C) $\bar{A} + \bar{B}$ (D) $\bar{A}\bar{B}$

22. The circuit shown in figure is equivalent to — gate.



- (A) X-OR gate (B) EX-NOR gate
 (C) Half adder (D) Half subtractor

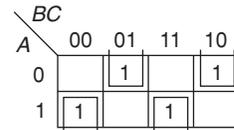
23. The truth table of the circuit shown in figure

A	B	C	Z
0	0	0	1
0	0	1	0
0	1	0	1
0	1	1	1
1	0	0	1
1	0	1	1
1	1	0	1
1	1	1	1

The Boolean expression for Z

- (A) $\overline{\overline{(A+B)}(B+C)}$ (B) $\overline{\overline{(A+B)}(B+C)}$
 (C) $\overline{\overline{(A+B)}(B+C)}$ (D) All of the above

24. A combinational circuit has input A, B and C and its K-map is as shown in figure. The output of the circuit is given by



- (A) $(\bar{A}B + A\bar{B})\bar{C}$ (B) $(AB + \bar{A}\bar{B})\bar{C}$
 (C) $\bar{A}\bar{B}\bar{C}$ (D) $A \oplus B \oplus C$

25. Which of the following two 2-input gates will realize the Boolean expression $X(P, Q, R) = \pi(0, 5)$

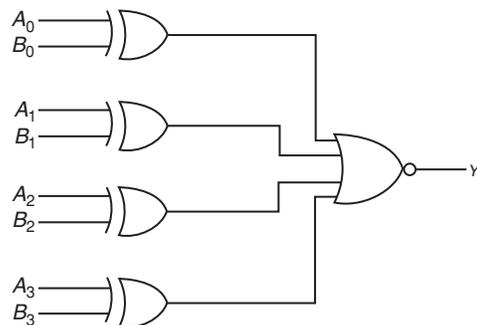
- (A) AND and OR (B) NAND and OR
 (C) AND and X-OR (D) OR and X-OR

26. Simplify the given function

$$f(x, y, z) = \sum m(0, 2, 3, 4, 5, 7)$$

- (A) $\bar{x}y + \bar{y}\bar{z} + xz$ (B) $\bar{x}\bar{z} + x\bar{y} + yz$
 (C) Both (A) and (B) (D) $\bar{x}\bar{z} + \bar{x}y + x\bar{y} + xz$

27. Figure below shows a digital circuit, which compares two numbers $A_0 A_1 A_2 A_3, B_0 B_1 B_2 B_3$. Choose the pair of correct input number to get output $Y = 0$.



- (A) 1100, 1100 (B) 0110, 0110
 (C) 1011, 0010 (D) 1011, 1011

28. How many 3 to 8 line decoders with an enable input are required to build 6 of 34 decoder?

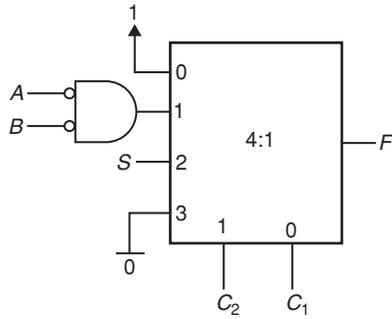
- (A) 6 (B) 2
 (C) 9 (D) 4

29. It is required to construct a 2^n to 1 multiplexer by using 2-to-1 multiplexer only. How many of 2-to-1 multiplexer are needed?

- (A) n (B) 2^{2n}
 (C) 2^{n-1} (D) $2^n - 1$

1.82 | Unit 1 • Digital Logic

30. Consider the following circuit



Which one of the following give the function implemented by the MUX based digital circuit?

- (A) $F = C_2 \cdot \overline{C_1}S + \overline{C_2}C_1(\overline{A} + \overline{B})$
- (B) $F = \overline{C_2} \cdot \overline{C_1} + C_2C_1 + C_2\overline{C_1}S + \overline{C_2}C_1\overline{AB}$
- (C) $F = \overline{AB} + S$
- (D) $F = \overline{C_2} \cdot \overline{C_1} + C_2 \cdot \overline{C_1}S + \overline{C_2}C_1\overline{A} \cdot \overline{B}$

ANSWER KEYS

- | | | | | | | | | | |
|-------|-------|-------|-------|-------|-------|-------|-------|-------|-------|
| 1. A | 2. D | 3. C | 4. A | 5. D | 6. A | 7. A | 8. A | 9. A | 10. B |
| 11. C | 12. D | 13. B | 14. A | 15. A | 16. B | 17. A | 18. D | 19. C | 20. D |
| 21. D | 22. B | 23. B | 24. D | 25. D | 26. C | 27. C | 28. C | 29. D | 30. D |

Computer Organization and Architecture

Chapter 1: Machine Instructions, Addressing Modes	2.3
Chapter 2: ALU and Data Path, CPU Control Design	2.16
Chapter 3: Memory Interface, I/O Interface	2.33
Chapter 4: Instruction Pipelining	2.47
Chapter 5: Cache and Main Memory, Secondary Storage	2.60

U

N

I

T

2

This page is intentionally left blank

Chapter 1

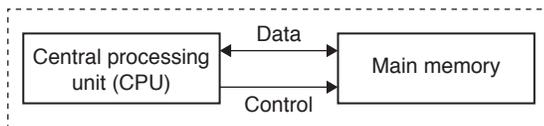
Machine Instructions, Addressing Modes

LEARNING OBJECTIVES

- Computer
- Computer system
- Computer component
- Machine instruction
- Instruction types
- Types of operands
- Types of operations
- Procedure call instruction
- Addressing modes
- Computer performance

COMPUTER

A computer is a data-processing machine which is operated automatically under the control of a list of instructions (called a program) stored in its main memory.



Computer System

- A computer system consists usually of a computer and its peripherals.
- Computer peripherals include input devices, output devices and secondary memories.

Computer Architecture: Computer Architecture refers to those attributes of a system visible to a programmer, i.e., the attributes that have direct impact on the logical execution of a program.

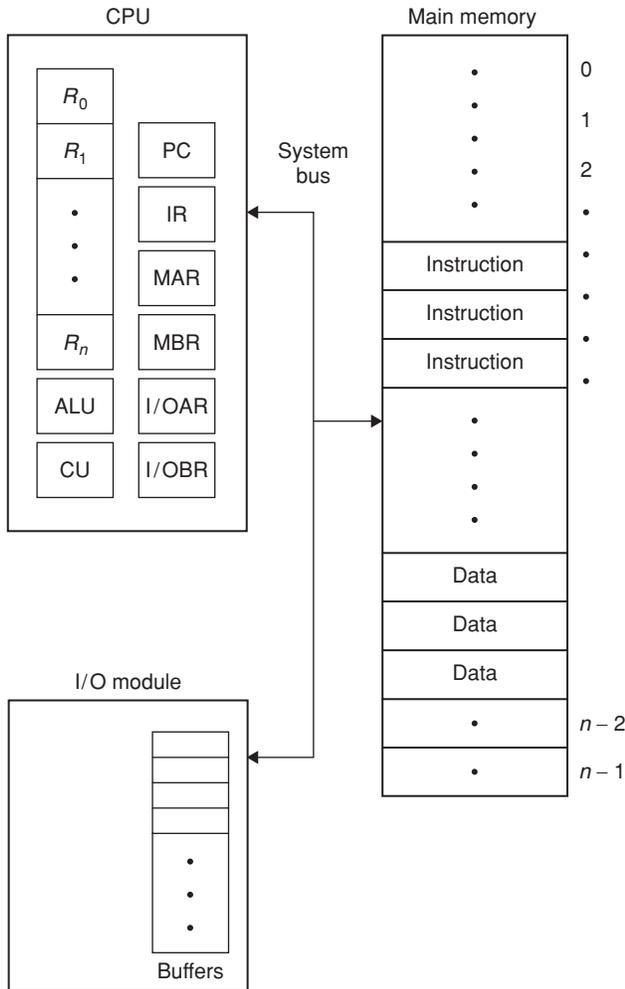
Example: Whether a computer will have a multiply instruction or not.

Computer Organization: Computer organization refers to operational units and their interconnections that realize the architectural specifications.

Example: Whether the multiply instruction will be implemented by a special multiply unit or by a mechanism that makes repeated use of the add unit of the system.

Computer Components

- $R_1, R_2 \dots R_n$: General Purpose Registers.
- PC:** Program counter. Holds address of next instruction to be executed. $PC = PC + I$ (I = instruction length)
- IR:** Instruction Register. It holds the instruction which is fetched from memory.
- MAR:** Memory Address Register: MAR specifies the address in memory for the next read or write.
- MBR:** Memory Buffer Register. It contains the data to be written into memory or receives the data read from memory.
- Input-outputAR:** Input-output Address Register. It specifies a particular input-output device.
- Input-outputBR:** Input-output Buffer Register. Used for the exchange of data between an input-output module and the CPU.
- ALU:** Arithmetic and Logic Unit. Used to perform arithmetic and logical operations.



- **CU:** Control Unit. It causes operations to happen within the processor. Also generates timing signals.
- **Memory:** It consists of set of locations, defined by sequential numbered address.
- **Input-output module:** Transfer data from external device to CPU and memory and vice versa.
- **System bus:** A bus that connects major computer components is called a system bus.

MACHINE INSTRUCTIONS

- The operation of the CPU is determined by the instructions it executes. These instructions are called machine instructions or computer instructions.
- The collection of different instructions that the CPU can execute is referred to as the CPU's instruction set.

Elements of Machine Instructions

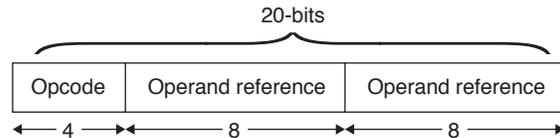
Each instruction must contain the information required by the processor for execution. The elements of a machine instruction are

1. Operation code: Specifies the operation
2. Source operand reference: Inputs for operation.

3. Result operand reference
4. Next instruction reference

Instruction representation

- Each instruction is represented by a sequence of bits.
- Example: 20-bit instruction format:



Instruction Types

Number of addresses

Most of the instructions have one, two or three operand addresses, with the address of next instruction being implicit.

- (i) **3-address instructions:** Computers with 3-address instruction formats can use each address field to specify either a processor register or a memory operand.

Example: 3-address instruction format for the evaluation of $X = (P + Q) \times (R + S)$ is

```
ADD R1, P, Q
ADD R2, R, S
MUL X, R1, R2.
```

Here R_1, R_2 are processor registers.

Advantage: Shorter programs when evaluating arithmetic expressions.

Disadvantage: The binary coded instructions required too many bits to specify three addresses.

- (ii) **2-address instructions:** These are most common in commercial computers. Each address field can specify either a processor register or a memory word.

Example: For evaluating $X = (P + Q) \times (R + S)$,

The 2-address instructions are

```
MOV R1, P
ADD R1, Q
MOV R2, R
ADD R2, S
MUL R1, R2
MOV X, R1.
```

(The first symbol of instruction is both source and destination)

- (iii) **One-address instructions:** Use an implied accumulator (AC) register for all data manipulations.

Example: 1-address instructions to evaluate $R = (P + Q) \times (R + S)$.

```
LOAD P
ADD Q
```

STORE T
 LOAD R
 ADD S
 MUL T
 STORE X .

Here ' T ' is a temporary memory location required to store the intermediate result.

- (iv) **Zero-address instructions:** A stack organized computer does not use an address field for the instructions ADD and MUL. The push and pop instructions require an address field to specify the operand that communicates with the stack.

Example: Zero-address instructions for the evaluation of $X = (P + Q) \times (R + S)$

PUSH P
 PUSH Q
 ADD
 PUSH R
 PUSH S
 ADD
 MUL
 POP X .

- (v) **RISC instructions:** The instruction set of a reduced instruction set computer (RISC) processor is restricted to the use of load and store instruction when communicating between memory and CPU. All other instructions are executed with in the register of the CPU without referring to memory.

Example: RISC instruction to evaluate,
 $X = (P + Q) \times (R + S)$

LOAD R_1, P
 LOAD R_2, Q
 LOAD R_3, R
 LOAD R_4, S
 ADD R_1, R_1, R_2
 ADD R_3, R_3, R_4
 MUL R_1, R_1, R_3
 STORE X, R_1

Types of Operands

Machine instructions operate on data. The most important general categories of data are

- Addresses
- Numbers
- Characters
- Logical data.

Types of Operations

The number of different opcodes varies widely from machine to machine. A useful and typical categorization is the following

1. Data transfer
2. Arithmetic

3. Logic
4. Conversion
5. Input-output
6. System control
7. Transfer of control

- (i) **Data transfer operations:** This type of instructions transfers data from one location to another.

Example: move, store, load, exchange, clear, set, push, pop.

- (ii) **Arithmetic operations:** Perform some function in ALU.

Example: add, subtract, multiply, divide, absolute, negate, increment, decrement

- (iii) **Logical operations:** Perform some logical operation in ALU and set condition codes and flags.

Example: AND, OR, NOT, EX-OR, Test, Compare, set control variables, shift, Rotate.

Let $R_1 = 10100101$, $R_2 = 00001111$ then

$(R_1) \text{ AND } (R_2) = 00000101$.

AND is also called mask operation.

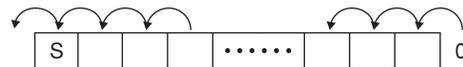
$(R_1) \text{ OR } (R_2) = 10101111$.

$\text{NOT } (R_1) = 01011010$.

$(R_1) \text{ EX-OR } (R_2) = 10101010$.

- (iv) **Shifting and rotating operations:** The operations are:

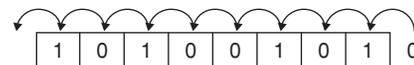
- (a) **Logical left shift:**



Here the bits of a word are shifted left. The left most bits is lost and 0 is shifted in right most bit position (i.e., bit empty).

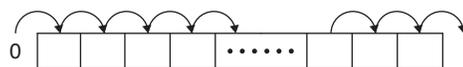
Example: $R_1 = 1010\ 0101$

Logical left shift R_1 :



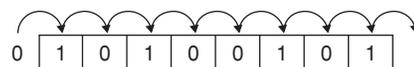
After left shift $R_1 = 0100\ 1010$.

- (b) **Logical right shift:** Here the bits of a word are shifted right. The right most bit lost and '0' is shifted in left most bit position.



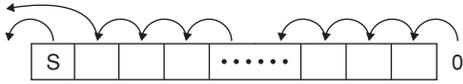
Example: $R_1 = 1010\ 0101$

Logical right shift $R_1 = 0101\ 0010$

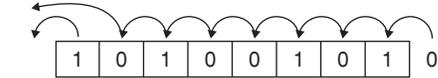


Logical shift operations are useful primarily for isolating fields within a word and also used to displace unwanted information.

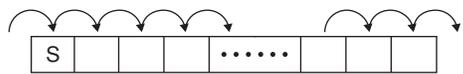
(c) **Arithmetic left shift:** Arithmetic shift operation treats the data as a signed integer and does not shift the sign bit. In Arithmetic left shift, a logical left shift is performed on all bits but the sign bit, which is retained.



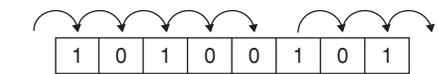
Example: $R_1 = 1010\ 0101$
Arithmetic Left shift $R_1 = 1100\ 1010$.



(d) **Arithmetic right shift:** Here, the sign bit is replicated into the bit position to its right.



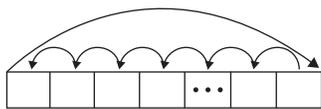
Example: $R_1 = 1010\ 0101$
Arithmetic Right shift $R_1 = 1101\ 0010$



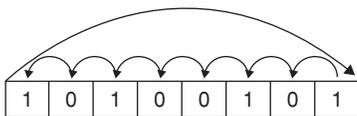
Notes:

1. With numbers in 2's complement form, a right arithmetic shift corresponds to a division by 2, with truncation for odd numbers.
2. Both arithmetic left shift and logical left shift correspond to a multiplication by 2 when there is no overflow.

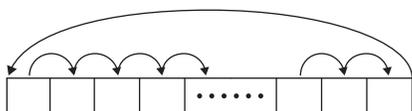
(e) **Left rotate (Cyclic left shift):** Rotate operations preserve all the bits being operated on. Here the bits from LSB will move one bit position to the left and MSB will be placed in LSB position.



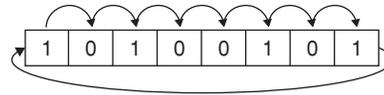
Example: $R_1 = 1010\ 0101$
Left Rotate $R_1 = 0100\ 1011$



(f) **Right rotate (Cyclic right shift):** Here the bits from MSB will be shifted to one bit position right and LSB is placed in MSB.



Example: $R_1: 1010\ 0101$
Right Rotate $R_1 = 1101\ 0010$.



(v) **Transfer of control:** This type of operations updates the program counter. Used for subroutine call/return, manage parameter passing and linkage.

Example: Jump, jump unconditional, return, execute, skip, skip conditional, Halt, Wait, NOP, etc.

(vi) **Input–output operations:** These are used to issue a command to input–output module.

Example: input, output, start input–output, test input–output etc.

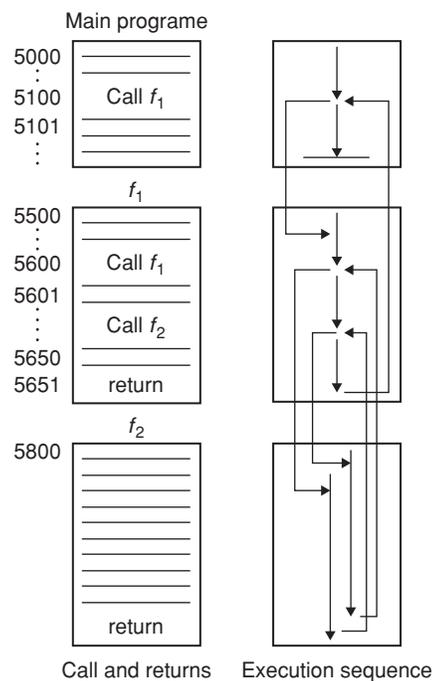
(vii) **Conversion operations:** These are similar to arithmetic and logical operations. May also involve special logic to perform conversion.

Procedure Call Instruction

A procedure is a self-contained computer program that is incorporated into a large program. It allows us to use the same piece of code many times. The procedure mechanism involves two basic instructions:

1. A call instruction that branches from the present location to the procedure.
2. A return instruction that returns from the procedure to the place from where it was called.

Example:

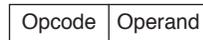


We can call a procedure from a variety of points, so the processor must somehow save the return address to return

Here operand in accumulator is implied in the definition of instruction.

- All register-reference instructions that use an accumulator are implied mode instructions.
- Zero-address instructions in a stack-oriented computer are implied-mode instructions.

(ii) **Immediate mode:** The operand is specified in the instruction itself. Instruction format in immediate mode is

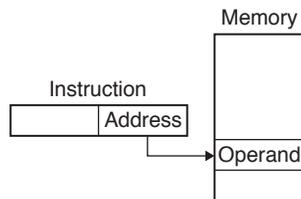


Example: Move A, 50.

- These are useful for initializing registers to constant value or to set initial values of variables.
- No memory reference is required other than the instruction fetch.
- The size of number is restricted to the size of the address (operand) field.

(iii) **Direct mode:**

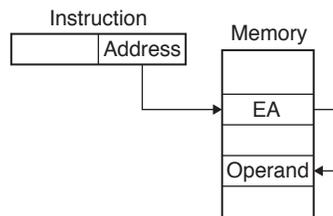
- Here the address of the operand is equal to the address part of the instruction.



- Required only one memory reference and no special calculation required.
- Limitation is limited address space.

(iv) **Indirect mode:**

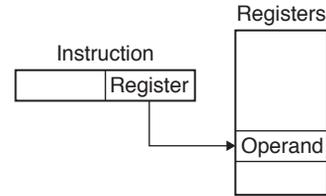
- The address field of the instruction gives the address of the operand which is stored in memory.
- The advantage of this approach is that for a word length of N , an address space of 2^N is available.
- The disadvantage is that the instruction execution requires two memory references to fetch the operand.



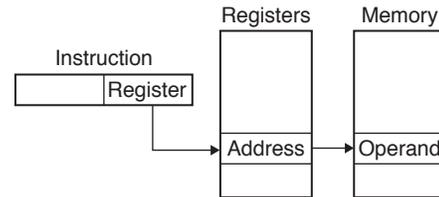
(Here EA is effective address of operand)

(v) **Register mode:**

- Here the operands are in registers that reside within the CPU.
- Only small address field required in instructions.
- No time consuming memory references are required.
- Address space is very limited.

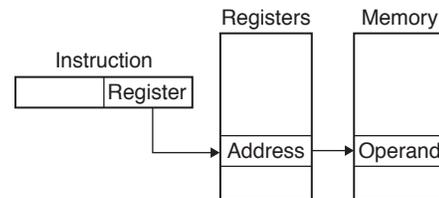


(vi) **Register indirect mode:** In this mode the instruction specifies a register in the CPU whose contents give the address of the operand in memory. Address field of the instruction uses fewer bits to select a register than would have required to specify a memory address directly.



Effective address: The effective address is defined to be the memory address obtained from the computation based on the addressing mode, consists the actual address of the operand.

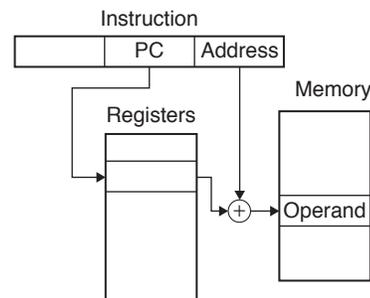
(vii) **Auto increment and auto decrement mode:** This is similar to register indirect mode except that the register is incremented or decremented after (or before) its value is used to access memory.



- After fetch operand, increment or decrement address.
- Used to access table of data.

(viii) **PC-relative mode:** Here the content of the program counter is added to the address part of the instruction to get the effective address.

The address part of the instruction is usually a signed number which can be either positive or negative. The effective address will be a displacement relative to address of the instruction.



Effective address = PC + address part

- (C) Relative (3) $Z = Y + X$
 (D) Indexed (4) $Z = Y$
 (A) $A - 4, B - 1, C - 2, D - 3$
 (B) $A - 3, B - 1, C - 2, D - 4$
 (C) $A - 4, B - 2, C - 1, D - 3$
 (D) $A - 3, B - 2, C - 1, D - 4$
3. A computer has 32-bit instruction and 12-bit addresses. If there are 250 two-address instructions, how many one-address instructions can be formulated?
 (A) 6 (B) 256
 (C) 12,288 (D) 24,576
4. The memory unit of a computer has 256k words of 32-bits each. The computer has an instruction format with four fields:
 1. An operation field
 2. A mode field to specify one of 8 addressing modes
 3. A Register address field to specify one of 120 processor registers.
 4. A memory address.
 Then what is the number of bits in each field respectively if the instruction is in one memory word?
 (A) 4, 3, 7, 18 (B) 3, 4, 7, 18
 (C) 2, 1, 6, 23 (D) 3, 7, 4, 18

Common data for questions 5 to 7: A relative mode branch type of instruction is stored in memory at an address equivalent to decimal 750. The branch is made to an address equivalent to decimal 400.

5. What should be the value of the relative address field of the instruction in decimal?
 (A) 351 (B) -351
 (C) 350 (D) -350
6. The relative address value in binary using 12-bits, will be
 (A) 000101011111 (B) 100101011111
 (C) 111010100000 (D) 111010100001
7. What will be the binary value in PC after the fetch phase (in binary)?
 (A) 001011101111 (B) 001011101110
 (C) 111011101111 (D) 110100010001

Common data for questions 8 to 10: Consider a 16-bit processor in which the following appears in main memory, starting at location 200:

200	Load to AC	Mode
201	500	
202	Next to instruction	

The first part of the first word indicates that this instruction loads a value into an accumulator. The mode field specifies an addressing mode or a source register, R_1 , which has a value 400. There is a base register that

contain the value 100. The value 500 in location 201, may be the part of address calculation. Assume that location 399 contains the value 999, location 400 contains the value 1000 and so on.

8. What will be the effective address and operand to be loaded by using Register indirect mode?
 (A) 200, 400 (B) 400, 1000
 (C) 400, 500 (D) 200, 1000
9. What will be the effective address using indirect addressing mode?
 (A) 200 (B) 201
 (C) 500 (D) Present in 500 location
10. What will be the effective address using immediate addressing mode?
 (A) 202 (B) 201
 (C) 500 (D) 400
11. A CPU of a computer has 48-bit instructions. A program starts at address $(600)_{10}$. Which one of the following is a legal program counter value in decimal?
 (A) 610 (B) 650
 (C) 672 (D) 693
12. Consider a new instruction named branch- on-bit-reset (bbr). The instruction ' $BBR R_1, I, label$ '
 Jumps to label, and if bit in position I of register operand, R_1 is zero. The registers of the computer are 16-bits wide and are numbered 0 to 15, position 0 being LSB. Consider the following implementation of this instruction on a processor that does not have *BBR* implemented.
 $Temp \leftarrow R_1$ and mask
 Branch to label if temp is zero.
 The variable 'temp' is a temporary register. For correct implementation, the variable 'mask' must be generated by
 (A) $mask \leftarrow 0 \times 1 \ll I$
 (B) $mask \leftarrow 0 \times FFFFFFFF \gg I$
 (C) $mask \leftarrow I$
 (D) $mask \leftarrow 0 \times F$.

13. Consider a hypothetical processor with an instruction of type

$LW R_1, 40(R_2)(R_3)$.

Which during execution reads a 16-bit word from memory and stores it in a 16-bit register R_1 . The effective address of the memory location is obtained by the addition of constant 40, contents of R_2 and R_3 registers. Which of the following best reflects the addressing mode implemented by this instruction for the operand in memory?

- (A) Immediate addressing
 (B) Register addressing
 (C) Register indirect scaled addressing
 (D) Base with index and displacement addressing

14. Which of the following is true of base-register addressing mode?
- It is useful in creating self-relocating code.
 - If it is included in an instruction set architecture, then an additional ALU is required for effective address calculation.
 - The amount of displacement depends on the content of base register.
- (A) (i) only (B) (ii) only
(C) (i) and (ii) only (D) (ii) and (iii) only
15. Which of the following addressing modes are suitable for program relocation at run time?
- Direct addressing
 - Based register addressing
 - PC-relative addressing
 - Index register addressing
- (A) (i) and (ii) (B) (ii) and (iii)
(C) (iii) and (iv) (D) (ii), (iii) and (iv)
16. In which of the following addressing mode, the address of the operand is inside the instruction?
- (A) Implied mode
(B) Absolute addressing mode
(C) Immediate addressing mode
(D) Register addressing mode
17. A certain processor supports only the immediate and the direct addressing modes. Which of the following programming language features can be on this processor?
- Pointers
 - Arrays
 - Initialization
- (A) (i) and (ii) (B) (i) and (iii)
(C) (ii) and (iii) (D) (iii) only
18. In which of the following situation, relative addressing mode is useful?
- (A) Coroutine writing
(B) Position-independent code writing
(C) Sharable code writing
(D) Interrupt handlers
19. In indexed addressing mode with scaling, the effective address is calculated as
- (A) Index + scaling + signed displacement
(B) (Index * scaling) + signed displacement
(C) Index + (scaling * displacement)
(D) (Index + scaling) * displacement
20. Which of the following addressing modes require more number of memory accesses?
- (A) DIRECT
(B) IMMEDIATE
(C) INDIRECT
(D) IMPLIED

Practice Problems 2

Directions for questions 1 to 20: Select the correct alternative from the given choices.

- The addressing mode that facilitates access to an operand whose location is defined relative to the beginning of the data structure in which it appears is

(A) Direct (B) Indirect
(C) Immediate (D) Index
- Stack addressing is same as

(A) Direct addressing
(B) Indirect addressing
(C) Zero addressing
(D) Relative addressing
- The Register which contain the Instruction to be executed is called

(A) Instruction register
(B) Memory address register
(C) Index register
(D) Memory data register
- The Register which keeps track of the execution of a program and which contains the memory address of the next instruction to be executed is called

(A) Instruction register
(B) Program counter
(C) Index register
(D) Memory address register
- A stack pointer is

(A) A 16-bit register in the microprocessor that indicate the beginning of the Stack Memory
(B) A register that decodes and execute 16-bit arithmetic operation.
(C) The first memory location where a subroutine address is stored
(D) A register in which flag bits are stored.
- Function of Control Unit in the CPU is

(A) To transfer data to primary storage
(B) To store program instruction
(C) To perform logic operations
(D) To generate timing signals
- When a subroutine is called the address of the instruction following the CALL instruction stored in the

(A) Stack (B) Accumulator
(C) Program counter (D) Stack pointer
- In Immediate addressing mode the operand is placed

(A) In the CPU Register
(B) After the OP Code in the instruction
(C) In the memory
(D) In the stack memory

9. When the RET instruction at the end of subroutine is executed
- The information where the stack is initialized is transferred to the stack pointer.
 - The memory address of the RET instruction is transferred to the program counter.
 - Two data bytes stored in the top two locations of the stack are transferred to the program counter.
 - Two data bytes stored in the top two location of the stack are transferred to the stack pointer.
10. Match the following:
- | List I | List II |
|------------------------------|--------------|
| P. Indirect Addressing | 1. Loops |
| Q. Auto decrement Addressing | 2. Constants |
| R. Immediate Addressing | 3. Pointers |
- P – 1, Q – 3, R – 2
 - P – 3, Q – 1, R – 2
 - P – 2, Q – 1, R – 3
 - P – 3, Q – 2, R – 1
11. An instruction used to set the carry flag in a computer can be
- Data control
 - Process control
 - Logical
 - Data transfer
12. The addressing mode in which the address of the location of the operand is given explicitly as part of the instruction is
- Direct addressing mode
 - Indirect addressing mode
 - Immediate addressing mode
 - Register addressing mode
13. The unit that is used to supervise each instructions in the CPU is
- Control register
 - Control logic unit
 - ALU
 - Address register
14. The address of the location to or from which data are to be transferred is called
- Memory data register
 - Memory address register
 - Program counter
 - Index register
15. Which register is used as a working area in CPU?
- Program counter
 - Accumulator
 - Stack pointer
 - Instruction register
16. Which of the following statement is false about the PC relative addressing mode?
- It allows indexing of array element with same instruction.
 - It enables reduced instruction size.
 - It enables faster address calculations than indirect addressing.
 - It enables easy relocation of data.
17. Which of the following is not an application of logic operations?
- Insert new bit values into a register
 - Change bit value
 - Delete a group of bits
 - Shift bit values in a register
18. In which of the following addressing mode, less number of memory references are required?
- Immediate
 - Register
 - Implied
 - All of the above
19. Which of the following is not involved in a memory write operation?
- MDR
 - MAR
 - PC
 - Data bus
20. In ____ addressing mode the instruction contains 8-bit signed offset, address register A_n and index register R_K .
- Basic index
 - Full index
 - Basic relative
 - Full relative

PREVIOUS YEARS' QUESTIONS

Common Data for Questions 1 to 3: Consider the following program segment. Here R_1 , R_2 and R_3 are the general purpose registers.

Instruction	Operation	Instruction Size (No. of Words)
MOV $R_1, (3000)$	$R_1 \leftarrow M[3000]$	2
LOOP: MOV $R_2, (R_3)$	$R_2 \leftarrow M[R_3]$	1
ADD R_2, R_1	$R_2 \leftarrow R_1 + R_2$	1
MOV $(R_3), R_2$	$M[R_3] \leftarrow R_2$	1
INC R_3	$R_3 \leftarrow R_3 + 1$	1

DEC R_1	$R_1 \leftarrow R_1 - 1$	1
BNZ LOOP	Branch on not zero	2
HALT	Stop	1

Assume that the content of memory location 3000 is 10 and the content of the register R_3 is 2000. The content of each of the memory locations from 2000 to 2010 is 100. The program is loaded from the memory location 1000. All the numbers are in decimal.

1. Assume that the memory is word addressable. The number of memory references for accessing the data

- in executing the program completely is: **[2007]**
 (A) 10 (B) 11
 (C) 20 (D) 21
2. Assume that the memory is word addressable. After the execution of this program, the content of memory location 2010 is: **[2007]**
 (A) 100 (B) 101
 (C) 102 (D) 110
3. Assume that the memory is byte addressable and the word size is 32 bits. If an interrupt occurs during the execution of the instruction 'INC R_3 ', what return address will be pushed on to the stack? **[2007]**
 (A) 1005 (B) 1020
 (C) 1024 (D) 1040
4. Which of the following is/are true of the auto-increment addressing mode?
 (i) It is useful in creating self-relocating code
 (ii) If it is included in an Instruction Set Architecture, then an additional ALU is required for effective address calculation
 (iii) The amount of increment depends on the size of the data item accessed **[2008]**
 (A) (i) only (B) (ii) only
 (C) (iii) only (D) (ii) and (iii) only
5. Consider a hypothetical processor with an instruction of type LW $R_1, 20(R_2)$, which during execution reads a 32-bit word from memory and stores it in a 32-bit register R_1 . The effective address of the memory location is obtained by the addition of a constant 20 and the contents of register R_2 . Which of the following best reflects the addressing mode implemented by this instruction for the operand in memory? **[2011]**
 (A) Immediate addressing
 (B) Register addressing
 (C) Register indirect scaled addressing
 (D) Base indexed addressing
6. Consider two processors P_1 and P_2 executing the same instruction set. Assume that under identical conditions, for the same input, a program running on P_2 takes 25% less time but incurs 20% more CPI (Clock cycles per instructions) as compared to the program running on P_1 . If the clock frequency of P_1 is 1GHz, then the clock frequency of P_2 (in GHz) is _____. **[2014]**
7. A machine has a 32-bit architecture with 1-word long instructions. It has 64 registers, each of which is 32 bits long. It needs to support 45 instructions, which have an immediate operand in addition to two register operands. Assuming that the immediate operand is an unsigned integer, the maximum value of the immediate operand is _____. **[2014]**
8. Consider a new instruction named branch-on-bit-set (mnemonic bbs). The instruction 'bbs reg, pos,

label' jumps to label if bit in position pos of register operand reg is one. A register is 32 bits wide and the bits are numbered 0 to 31, bit in position 0 being the least significant. Consider the following emulation of this instruction on a processor that does not have bbs implemented.

temp \leftarrow reg & mask

Branch to label if temp is non-zero.

The variable temp is a temporary register. For correct emulation, the variable mask must be generated by **[2006]**

- (A) mask $\leftarrow 0 \times 1 \ll \text{pos}$
 (B) mask $\leftarrow 0 \times \text{ffffff} \gg \text{pos}$
 (C) mask $\leftarrow \text{pos}$
 (D) mask $\leftarrow 0 \times \text{f}$

9. Consider a processor with byte-addressable memory. Assume that all registers, including Program Counter (PC) and Program Status Word (PSW), are of size 2 bytes. A stack in the main memory is implemented from memory location $(0100)_{16}$ and it grows upward. The stack pointer (SP) points to the top element of the stack. The current value of SP is $(016E)_{16}$. The CALL instruction is of two words, the first word is the op-code and the second word is the starting address of the subroutine (one word = 2 bytes). The CALL instruction is implemented as follows:
- Store the current value of PC in the stack
 - Store the value of PSW register in the stack
 - Load the starting address of the subroutine in PC

The content of PC just before the fetch of a CALL instruction is $(5FA0)_{16}$. After execution of the CALL instruction, the value of the stack pointer is **[2015]**

- (A) $(016A)_{16}$ (B) $(016C)_{16}$
 (C) $(0170)_{16}$ (D) $(0172)_{16}$

10. A processor has 40 distinct instructions and 24 general purpose registers. A 32-bit instruction word has an opcode, two register operands and an immediate operand. The number of bits available for the immediate operand field is _____. **[2016]**
11. Suppose the functions F and G can be computed in 5 and 3 nanoseconds by functional units U_F and U_G , respectively. Given two instances of U_F and two instances of U_G , it is required to implement the computation $F(G(X_i))$ for $1 \leq i \leq 10$. Ignoring all other delays, the minimum time required to complete this computation is _____ nanoseconds. **[2016]**
12. Consider a processor with 64 registers and an instruction set of size twelve. Each instruction has five distinct fields, namely, opcode, two source register identifiers, one destination register identifier, and a twelve-bit immediate value. Each instruction must be stored in

memory in a byte - aligned fashion. If a program has 100 instructions, the amount of memory (in bytes) consumed by the program text is ____.

[2016]

13. Consider the C struct defined below:

```
struct data {
    int marks [100];
    char grade;
    int cnumber;
};
struct data student;
```

The base address of student is available in register R1. The field student.grade can be accessed efficiently using

[2017]

- (A) Post-increment addressing mode, $(R1)+$
 (B) Pre-decrement addressing mode, $-(R1)$
 (C) Register direct addressing mode, R1
 (D) Index addressing mode, $X(R1)$, where X is an offset represented in 2's complement 16-bit representation.
14. Consider a RISC machine where each instruction is exactly 4 bytes long. Conditional and unconditional branch instructions use PC-relative addressing mode. Offset specified in bytes to the target location of the branch instruction. Further the Offset is always with respect to the address of the next instruction in the

program sequence. Consider the following instruction sequence.

Instr. No.	Instruction
i :	add R2, R3, R4
i+1 :	sub R5, R6, R7
i + 2 :	cmp R1, R9, R10
i + 3 :	beq R1, Offset

If the target of the branch instruction is i, then the decimal value of the Offset is ____.

[2017]

15. A processor has 16 integer registers ($R0, R1, \dots, R15$) and 64 floating point registers ($F0, F1, \dots, F63$). It uses a 2-byte instruction format. There are four categories of instructions: Type-1, Type-2, Type-3, and Type-4. Type-1 category consists of four instructions, each with 3 integer register operands ($3R_s$). Type-2 category consists of eight instructions, each with 2 floating point register operands ($2F_s$). Type-3 category consists of fourteen instructions, each with one integer register operand and one floating point register operand ($1R + 1F$). Type-4 category consists of N instructions, each with a floating point register operand ($1F$).

The maximum value of N is ____.

[2018]

ANSWER KEYS

EXERCISES

Practice Problems 1

1. A 2. A 3. D 4. A 5. B 6. D 7. A 8. B 9. D 10. B
 11. C 12. A 13. D 14. A 15. B 16. B 17. B 18. B 19. B 20. C

Practice Problems 2

1. D 2. C 3. A 4. B 5. A 6. D 7. A 8. B 9. C 10. B
 11. B 12. A 13. B 14. B 15. B 16. A 17. D 18. B 19. D 20. B

Previous Years' Questions

1. D 2. A 3. C 4. C 5. D 6. 1.6 7. 16383 8. A 9. D 10. 16
 11. 28 12. 500 13. D 14. -16 15. 32

Chapter 2

ALU and Data Path, CPU Control Design

LEARNING OBJECTIVES

- Arithmetic and logic unit
- Fixed-point arithmetic operation
- Floating point arithmetic operation
- BCD
- Data path
- CPU control design
- Instruction cycle
- Control unit
- Control of processor
- Function of control unit
- Design of control unit
- Types of micro-instructions
- Micro-instruction sequencing
- RISC and CISC
- RISC characteristic
- CISC characteristic

ALU (ARITHMETIC AND LOGIC UNIT)

ALU performs arithmetic and logical operations on data (see Figure 1).

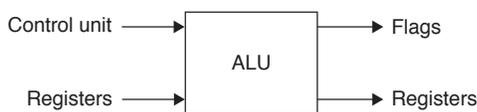


Figure 1 ALU inputs and outputs

- Data are presented to ALU in registers and the results of an operation are stored in registers.
- Registers are temporary storage locations within the processor that are connected by signal paths to ALU.
- The control unit provides signals that control the operation of ALU and the movement of data into and out of the ALU.
- Here we will discuss
 - Fixed-point arithmetic operations
 - Floating-point arithmetic operations
 - BCD data arithmetic operations

Fixed-point Arithmetic Operations

Fixed-point representation

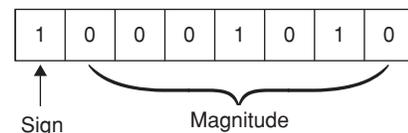
The numbers may be positive, zero or negative. So we have two types of numbers:

Unsigned numbers Only zero and positive integers can be represented. All bits represent magnitude and no need of sign.

Signed numbers In signed representation, the most significant bit represents the sign. If the number is positive, the MSB is 0 and remaining bits represent magnitude. If the number is negative, we have three techniques to represent that number:

- Signed magnitude representation:** In signed magnitude representation, the MSB represents sign and remaining bits represent magnitude. If the number is negative then the MSB is 1.

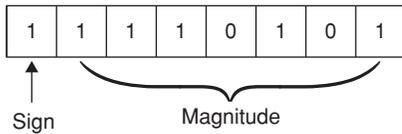
Example: Signed magnitude representation of $-10 =$



- Signed 1's complement representation:** In signed 1's complement representation, the MSB bit is 1. The remaining bits of its signed magnitude bits are inverted i.e., convert 0's to 1's and 1's to 0's to obtain 1's complement.

Example:

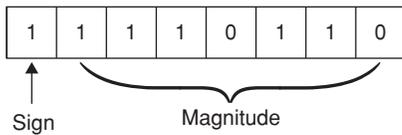
Signed 1's complement (-10) =



- 3. Signed 2's complement representation:** To get signed 2's complement representation, add 1 to the signed 1's complement of that number.

Example:

Signed 2's complement (-10) =



Fixed-point arithmetic operations

We will discuss the following operations using signed magnitude data and signed 2's complement data.

1. Addition
2. Subtraction
3. Multiplication
4. Division

Addition and subtraction using signed magnitude data Consider two numbers whose magnitude is represented as A and B . When the signed numbers are added or subtracted, there are eight different conditions to consider, depending on the sign of the numbers and operation performed.

Operation	Add Magnitudes	Subtract Magnitudes ($A > B$)
$(+A) + (+B)$	$+(A + B)$	
$(+A) + (-B)$		$+(A - B)$
$(-A) + (+B)$		$-(A - B)$
$(-A) + (-B)$	$-(A + B)$	
$(+A) - (+B)$		$+(A - B)$
$(+A) - (-B)$	$+(A + B)$	
$(-A) - (+B)$	$-(A + B)$	
$(-A) - (-B)$		$-(A - B)$

Algorithm for addition (subtraction): When the signs of A and B identical (different), add the two magnitudes and attach the sign of A to the result. When the signs of A and B are different (identical), compare the magnitudes and subtract the smaller number from the larger. Choose the sign of result based on magnitudes of A and B .

Example: All eight cases for the numbers $A = 5, B = 2$.

$$\begin{aligned}
 (+A) + (+B) &= (+5) + (+2) \\
 &= 0101 + 0010 = 0111 = +7 \\
 (+A) + (-B) &= (+5) + (-2) \\
 &= 0101 + 1010
 \end{aligned}$$

Take 2's complement of -2 and add it to 5

$$101$$

$$110$$

$$-----$$

$$1] 011$$

$$\uparrow$$

Discard

$$\therefore \text{result} = +3 \ (A > B)$$

$$(-A) + (+B) = (-5) + (+2)$$

$$= 1101 + 0010$$

add 2's complement of -5 to 2

$$011$$

$$010$$

$$101$$

As MSB is 1 take 2's complement to get original number i.e., 011.

$$\text{Result} = 011 = -3 \ (\because A > B)$$

$$(-A) + (-B) = (-5) + (-2)$$

$$= 1101 + 1010$$

$$101$$

$$010$$

$$\text{Result} = 1111 = -7$$

Similarly we can perform the subtractions using signed magnitude data.

Hardware implementation:

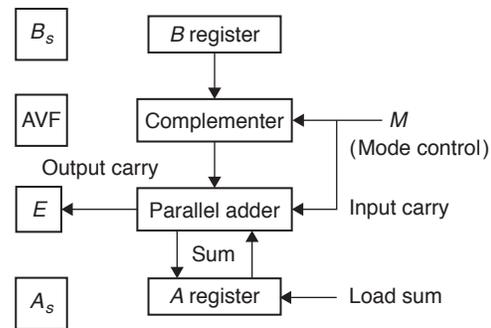


Figure 2 Hardware implementation for addition and subtraction.

Figure 2 shows the hardware implementation for addition and subtraction operations. It consists of registers A and B and sign flip-flops A_s and B_s . Subtraction is done by adding A to the 2's complement of B . The output carry is transferred to flip-flop E and add overflow flip-flop AVF holds the overflow bit when A and B are added. The addition is done through the parallel adder. The output of adder is sent to 'A' register. The complements provides an output of B or complement of B depending on the state of mode control M . When $M = 0$, the output equal to $A + B$, when $M = 1$, the output equal to $A + \bar{B} + 1$, i.e., $A - B$.

Addition and Subtraction with signed 2's complement data

Addition: In 2's complement representation, addition proceeds as if the two numbers were unsigned integers. If the result of the operation is positive, we get a positive number

in 2's complement form, which is same as in unsigned integer form. If the result of the operation is negative, we get a negative number in 2's complement form.

Example:

+5 = 0101
 +2 = 0010
0111 = +7
 +5: 0101
 -2: 1110
10011 = +3
 -5: 101
 +2: 0010
1101

As the result is negative, take 2's complement of result to get original number, i.e., 0110 + 1 = 0111 and the answer is -3.

-5: 1011
 -2: 1110
11001

As the result is negative take 2's complement to get original number, i.e., 0110 + 1 = 0111.

∴ Answer is -7.

Note: If two numbers are added and they are both positive or both negative, then overflow occurs if the result has the opposite sign.

Subtraction: To subtract subtrahend from minuend, take the 2's complement of subtrahend and add it to the minuend.

Example:

+5: 0101
 +2: 0010
 To subtract these two numbers add 2's complement of 2 to 5.
 +5: 0101
 -2: 1110
10011 = +3
 +5: 0101
 -2: 1110
 2's complement of -2 = 0010.
 +5: 0101
 +2: 0010
0111 = +7

Similarly for the other cases we can perform the subtraction.

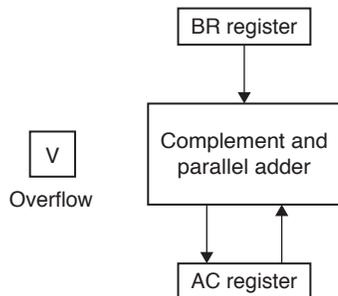


Figure 3 Hardware implementation for signed 2's complement addition and subtraction:

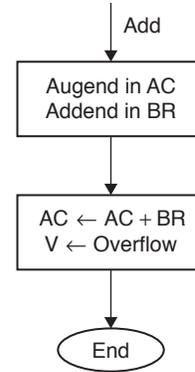


Figure 4 Flowchart for addition in 2's complement form

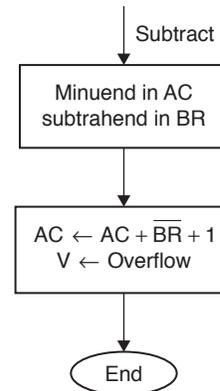


Figure 5 Flowchart for subtraction of 2's complement data

Multiplication of signed magnitude data

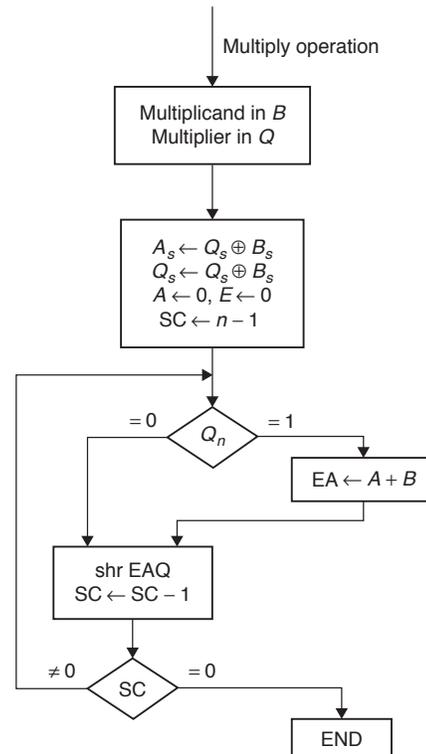


Figure 6 Flowchart for multiplication of signed magnitude data

Multiplication of two fixed point binary numbers in signed magnitude representation is a process of successive shift and add operations (see Figure 6).

Example 1: Multiply the two numbers -7 and $+8$, using 5-bit registers.

$-7 = 10111$

$+8 = 01000$

By excluding sign-bits, the multiplicand, $B = 0111$ and multiplier $Q = 1000$. Initially $A = 0000$, SC is sequence counter contains number of bits in multiplier magnitude.

Here $SC = 4$

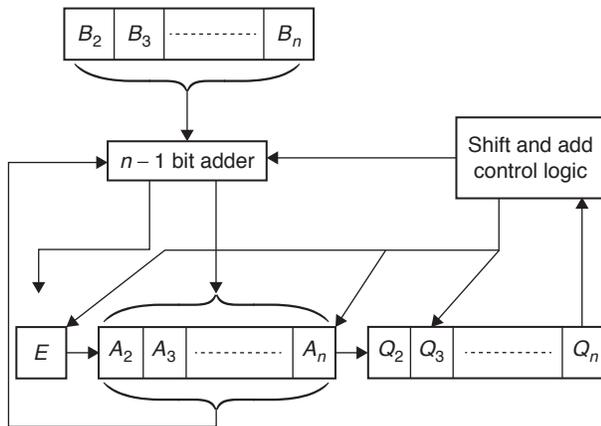
Multiplicand B = 0111	E	A	Q	SC
Multiplier in Q	0	0000	1000	4
Last bit of Q, $Q_n = 0 \Rightarrow \text{Shr } EAQ$	0	0000	0100	3
$Q_n = 0 \Rightarrow \text{Shr } EAQ$	0	0000	0010	2
$Q_n = 0 \Rightarrow \text{Shr } EAQ$	0	0000	0001	1
$Q_n = 1 \Rightarrow \text{Add } B \text{ to } A$	0	0000 0111 0111	0001	
Shr EAQ	0	0011	1000	0

$B \times Q = 00111000 = 56$

Sign = $Q_s \oplus B_s = 1 \oplus 0 = 1$

\therefore Result = -56

Hardware for signed magnitude data multiplication: B_1, A_1, Q_1 represent the respective signs of the registers B, A, Q . Final result will be in AQ , which consist of $2n$ -bits. (Here each register has n -bits).



Multiplication of Signed 2's complement data The straight forward multiplication will not work if either the multiplicand or the multiplier is negative. There are number of ways to perform multiplication of signed 2's complement data. One such a technique is *Booth's multiplication algorithm*. The following flowchart depicts about Booth's algorithm (see figure 7).

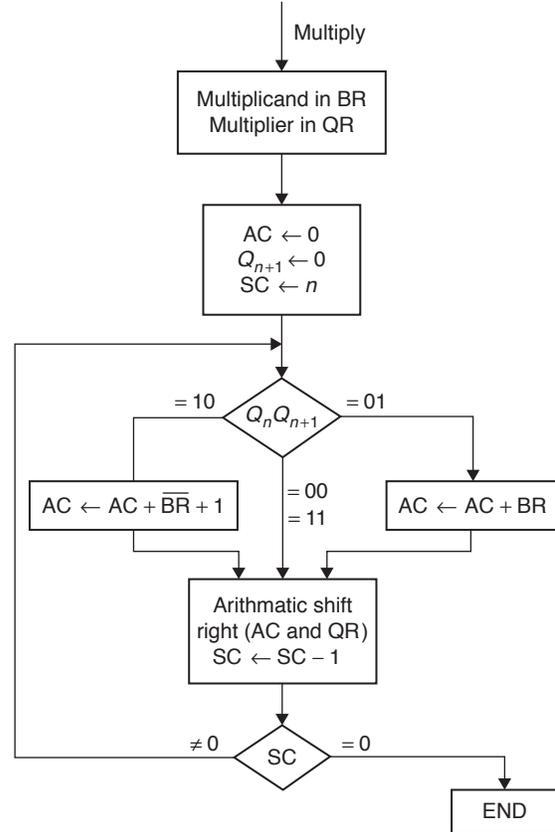


Figure 7 Booth's multiplication algorithm

An additional 1-bit register placed logically to the right of the LSB (Q_n) of Q register designated Q_{n+1} .

Example 2: Multiply the two numbers -7 and $+8$ using booth's algorithm, using 5 bits.

$BR = -7 = 11001$

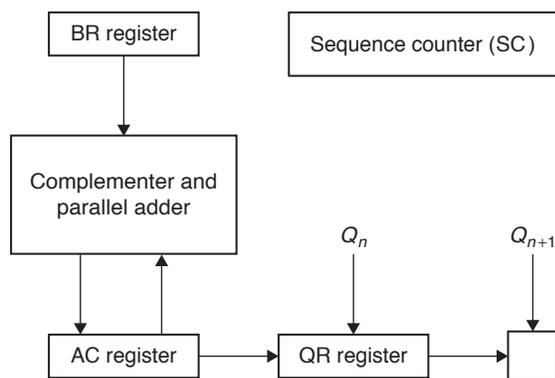
$QR = +8 = 01000$

Initially $AC = 00000, Q_{n+1} = 0, SC = 5$

$Q_n Q_{n+1}$	BR = 11001 BR + 1 = 00111	AC	QR	Q_{n+1}	SC
Initial		00000	01000	0	5
00	ashr(AC and QR)	00000	00100	0	4
00	ashr(AC and QR)	00000	00010	0	3
00	ashr(AC and QR)	00000	00001	0	2
10	subtract BR ashr(AC and QR)	00000 00111 00111 00011	10000	1	1
01	add BR ashr(AC and QR)	00011 11001 11100 11110	01000	0	0

\therefore Result = $1111001000 = -56$

Hardware implementation for Booth's algorithm:



Note: These two multiplication algorithms are sequential but we can also do the operation by means of a combinational circuit that forms the product bits all at once. The circuit consist of AND gates and adders.

Division algorithms

Division of signed magnitude data: Division of signed magnitude data is a process of successive compare, shift and subtract operations.

Example:

Dividend = 0111000000
 Divisor = 10001
 10001) 0111000000 (11010
 -10001
 010110
 -10001
 0010100
 -10001
 000110

Hardware implementation:

- The hardware implementation of division is same as multiplication, instead of shifting the divisor to the right, the dividend or partial remainder is shifted to the left, thus leaving the two numbers in the required relative position.
- The divisor is stored in the *B* register and the double-length dividend is stored in registers *A* and *Q*. The dividend is shifted to the left and the divisor is subtracted by adding its 2's complement value. The information about the relative magnitude is available in *E*.
- If $E = 1$, it signifies that $A \geq B$. *A* quotient bit 1 is inserted into Q_n and the partial remainder is shifted to the left to repeat the process.
- If $E = 0$, it signifies that $A < B$ so the quotient is Q_n remains a 0. The value of *B* is added to restore the partial remainder in *A* to its previous value. The partial remainder is shifted to the left and the process is repeated again until all quotient bits are formed.
- Finally, the quotient is in *Q* and remainder is in *A*. This method is called *restoring method*.

Divide overflow:

- A divide overflow condition occurs if the high-order half bits of the dividend constitute a number greater than or equal to the divisor.
- A division by zero must be avoided.

Other algorithms for division: Two other methods are available for dividing numbers:

Comparison method: To divide the two numbers *A* and *B* in comparison method, they are compared prior to the subtraction operation. If $A \geq B$, *B* is subtracted from *A*. If $A < B$ nothing is done. The partial remainder is shifted left and the numbers are compared again.

Non-restoring method: To divide two numbers *A* and *B* is non-restoring method, *B* is not added if the difference is negative but instead, the negative difference is shifted left and then *B* is added.

Floating-point Arithmetic Operations

Floating-point representation

Fixed-point representation allows representation of numbers with fractional component as well. But this approach has limitations. It is not possible to represent very large numbers and very small numbers in fixed point representation.

In floating-point representation, the numbers can be represented in the form,

$$\pm S \times B^{\pm E}$$

The three fields are

Sign: plus or minus

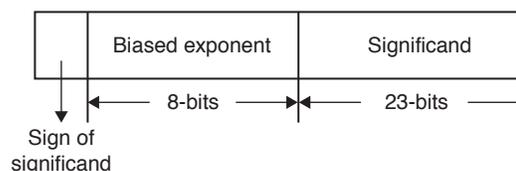
Significand: *S*

Exponent: *E*

are stored in a binary word.

The base *B* is implicit and need not be stored because it is same for all the numbers. It is assumed that the radix point is to the right of the left most or most significant bit of the significand, i.e., there is one bit to the left of the radix point.

32-bit floating-point format: The left most bit stores the sign of the number. The exponent value is stored in next 8-bits. This is represented in biased representation.



Biased representation: In biased representation, a fixed value, called the bias, is subtracted from the exponent field to get the true exponent value.

Bias = $(2^{k-1} - 1)$, where k = number of bits in binary exponent. In IEEE 32-bit floating point representation, bias = $2^7 - 1 = 127$.

And the range of true exponents is -127 to $+128$.

The advantage of biased representation is that non-negative floating point numbers can be treated as integers for comparison purposes.

The last portion of word is the significand.

Normalized numbers:

- To simplify operation on floating point numbers, it is typically required that they must be normalized.
- A normalized number is one in which the most significant digit of significand is non-zero.
- For base-2 representation, a normalized number is therefore one in which the MSB of the significand is one.
- Normalized non-zero number is one in the form $\pm 1.bbb \dots b \times 2^{\pm E}$, where b is either binary digit 0 or 1.

- As the MSB is always one, it is unnecessary to store this bit. Thus the 23-bit field is used to store a 24-bit significand with a value in the half open interval $(1, 2)$.
- A number may be normalized by shifting the radix point to the right of the leftmost 1 bit and adjusting the exponent accordingly.

Example: In 32-bit floating representation of $-1.6328125 \times 2^{-20}$,

Sign = 1 (as the number is negative)
 $(.6328125)_{10} = (.1010001\dots)_2$
 Exponent = -20
 Biased exponent = $127 - 20 = 107$
 = 1101011
 $\therefore -1.6328125 \times 2^{-20}$
 = 1 01101011 1010001000000000000000

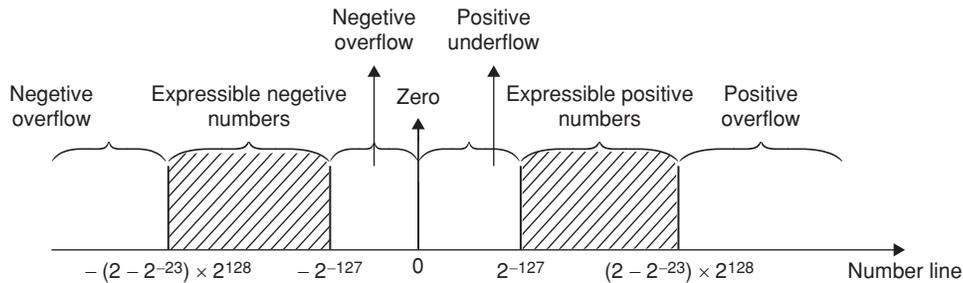
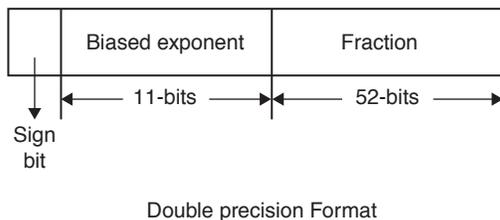
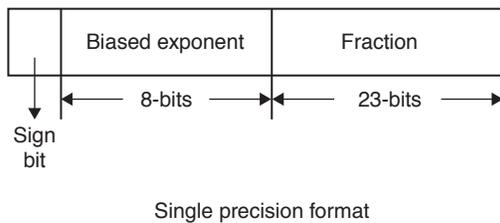


Figure 8 Range of expressible numbers in a 32-bit floating point format

IEEE standard for binary floating-point representation



Example:

$$(123 \times 10^0) + (234 \times 10^{-2}) = 123 \times 10^0 + 2.34 \times 10^0 = 125.34 \times 10^0$$

(ii) **Multiplication:** The steps to multiply two floating point numbers are

1. Check for zeros
2. Add the exponents
3. Multiply the significands
4. Normalize the result

(iii) **Division:** The steps to divide two floating point numbers are

1. Check for zeros
2. Initialize registers and evaluate the sign
3. Align the dividend
4. Subtract the exponents
5. Divide the significands

Floating-point arithmetic

(i) **Addition and subtraction:** The algorithm consists the following phases:

1. Check for zeros
2. Align the significands/mantissas
3. Add or subtract the significands
4. Normalize the result

Binary-Coded Decimal (BCD) Arithmetic Operations

Computers capable of performing decimal arithmetic must store the data in binary-coded form.

Example: BCD of 239 = 0010 0011 1001

BCD addition

In BCD each digit do not exceed 9, so the sum of two BCD digits cannot be greater than $9 + 9 + 1 = 19$, the 1 in the sum being an input carry. When the binary sum of two BCD digits is greater than 1001, we obtain a non-valid BCD representation. The addition of binary 6 (0110) to the binary sum converts it to the correct BCD representation and also produces an output carry as required.

Example:

$$\begin{array}{r}
 239 = 0010\ 0011\ 1001 \\
 426 = 0100\ 0010\ 0110 \\
 665 = 0110\ 0101\ 1111 > 1001 \\
 \quad \quad \quad \underline{0110} \\
 \quad \quad \quad 0110\ 0110\ 0101 \\
 = 665
 \end{array}$$

BCD subtraction

- Perform the subtraction by taking the 9’s or 10’s complement of the subtrahend and adding it to the minuend.
- The 9’s complement of a decimal digit represented in BCD can be obtained by complementing the bits in the coded representation of the digit, provided a correction is included.

There are two possible correction methods:

1. Binary 1010 is added to each complemented digit and the carry discarded after each addition.

Example:

$$\begin{array}{r}
 9\text{'s complement of } 7 = 2 \\
 7 \text{ in BCD} = 0111. \\
 \text{Complement of } 7 = 1000 \\
 \text{Add } 1010 \quad \quad = \underline{1010} \\
 \quad \quad \quad \quad \quad 1] 0010 = 2 \\
 \quad \quad \quad \quad \quad \uparrow \\
 \quad \quad \quad \quad \quad \text{Discard}
 \end{array}$$

2. Binary 0110 is added before the digit is complemented.

Example:

$$\begin{array}{r}
 \text{BCD of } 7 = 0111 \\
 \text{Add } 0110 = \underline{0110} \\
 \quad \quad \quad \underline{1101} \\
 \text{Complement} = 0010 = 2
 \end{array}$$

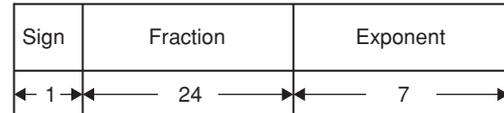
Example 4: Which of the following multiplier bit pattern of Booth’s multiplication algorithm gives worst case performance?

- (A) 01010101....0101
- (B) 000000....0000
- (C) 11111111....1111
- (D) 011110111110....01110

Solution: (A)

Booth’s multiplication algorithm works well with consecutive 0’s or 1’s. But it gives worst case performance when the multiplier consists of alternative 0’s and 1’s. (As 01, 10, pattern leads to addition and subtractions).

Example 5: Consider the following 32-bit floating point representation scheme as shown in the format below:



A value is specified by three fields:

Sign field: 1 bit (0 for positive and 1 for negative values)
 Fraction: 24-bits (with binary point being at the left end of fraction bits)

Exponent: 7-bits (in excess-64 signed integer representation)

The base of exponentiation is 16. The sign bit is in MSB. Then the normalized floating point representation of -6.5 is

- (A) E8000042
- (B) E1000012
- (C) D8000841
- (D) D0000042

Solution: (D)

Here sign = 1 as the number is negative.

$$\begin{aligned}
 (-6.5)_{10} &= (-0110.1)_2 \\
 &= (-1.101 \times 2^2)
 \end{aligned}$$

$$\text{Fraction} = 101000000000000000000000$$

$$\text{Exponent} = \text{Excess} - 64 \text{ exponent}$$

$$= 64 + 2 = 66 = 1000010$$

$$\therefore (-6.5)_{10}$$

$$= 1\ 1010000000000000000000001000010$$

$$= D0000042.$$

DATA PATH

Data path consists of the components of the processor that performs arithmetic operations.

Components of data path: ALU is just one data path building block. Other components are

1. Computational Components, which consist of combinational circuits (output follow inputs)

Example: ALU.

2. State components, which consists of sequential circuits (output changes on clock edge)

Example: Registers.

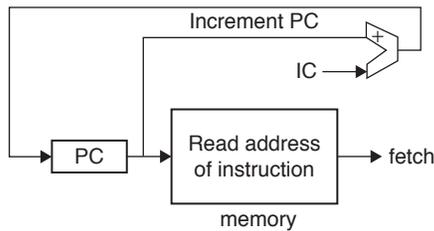
Example: The sequence of steps for the addition of two registers content are

1. $R1_{out}, X_{in}$
2. $R2_{out}$, Choose X, ADDITION, Y_{in}
3. $Y_{out}, R3_{in}$

(Each step executed in a single clock cycle).

- Data path and control unit forms the processing unit of a computer. The Data path includes ALU, multiplexers, all registers (like PC, IR) etc.

Example Data Path Design:

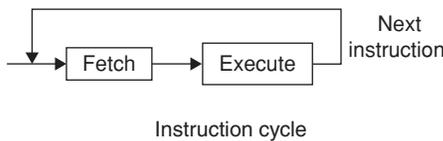


CPU CONTROL DESIGN

Instruction Cycle

A program residing in the memory unit of the computer consists of a sequence of instructions. The program is executed in the computer by going through a cycle for each instruction. Each instruction cycle in turn is subdivided into a sequence of sub cycles. For example, the phases of instruction cycle may be

1. Fetch
2. Decode
3. Read effective address
4. Execute, etc.



The cycle will be repeated, till all the instructions are executed.

Each phase is made up of more fundamental operations, called micro-operations.

Example micro-operations: Transfer between registers, simple ALU operation, etc.

Control Unit

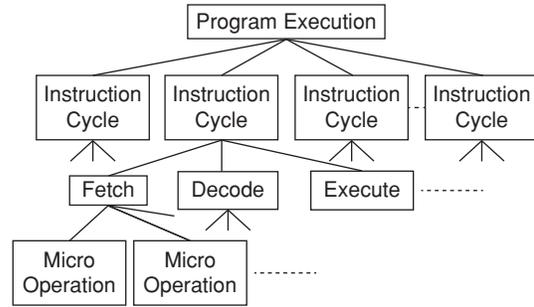
The control unit of a processor performs two tasks:

1. It causes the processor to execute micro-operations in the proper sequence, determined by the program being executed.
2. It generates the control signals that cause each micro-operation to be executed.

We now discuss the micro-operations of various phases of instruction cycle.

Micro-operation

- These are the functional or atomic operations of a processor.



Fetch Cycle: ‘Fetch’ stage of an instruction occurs at the beginning of each instruction, which causes an instruction to be fetched from memory. The micro-operations involved in fetch phase are

- t_1 : $MAR \leftarrow PC$ (move contents of PC to MAR)
- t_2 : $MBR \leftarrow \text{memory}; PC \leftarrow (PC) + I$ (move contents of MAR location to MBR and increment PC by I)
- t_3 : $IR \leftarrow (MBR)$ (move contents of MBR to IR)

Here I is instruction length.

Each micro-operation can be performed within the time of a single time unit.

Execute Cycle: For a machine with N different opcodes, there will be N different sequence of micro-operations. For the execution of following instruction.

Add R_1, X , the micro-operations will be

- t_1 : $MAR \leftarrow (IR(\text{Address}))$
- t_2 : $MBR \leftarrow \text{memory}$
- t_3 : $R_1 \leftarrow (R_1) + (MBR)$

Control of the Processor

(i) Functional requirements of control unit: Let us consider the following concepts to the characterization of a CU.

1. Define the basic elements of the processor.
2. Describe the micro-operations that the processor performs.
3. Determine the functions that the control unit must perform to cause the micro-operations to be performed.

(ii) Basic elements of processor:

- ALU
- Registers
- Internal data path: Used to move data between registers and between register and ALU.
- External data path: Used to link registers to memory and *input-output* modules, often by means of a system bus.
- Control unit: Causes operations to happen within the processor.

(iii) Micro-operations of processor:

- Transfer data from one register to another.
- Transfer data from a register to an external interface.

- Transfer data from an external interface to a register.
- Perform an arithmetic or logic operation, using registers for input and output.

(iv) Control unit tasks:

- **Sequencing:** The control unit causes the processor to step through a series of micro-operations in the proper sequence, based on the program being executed.
- **Execution:** The control unit causes each micro-operation to be executed.

(v) Control signals: For the control unit to perform its function, it must have inputs that allow it to determine the state of the system and outputs that allows it to control the behaviour of the system. These are external specifications of the control unit.

Internally, the control unit must have the logic required to perform its sequencing and execution functions.

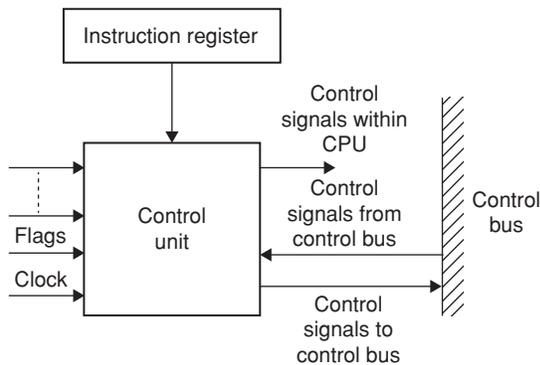


Figure 9 Block diagram of control unit

- (a) Clock:** This is how the control unit ‘keeps time.’ The control unit causes one micro-operation to be performed for each clock pulse. This is referred as processor cycle time or clock cycle time.
- (b) Instruction registers:** The opcode of current instruction is used to determine which micro-operations to perform during the execute cycle.
- (c) Flags:** Used to determine the status of the processor and outcome of previous ALU operations.
- (d) Control signals from control bus:** The control bus portion of system bus provides signals to the control unit.
- (e) Control signals within the processor:**
 1. Those that cause data to be moved from register to another.
 2. Those that activate specific ALU functions.
- (f) Control signals to control bus:**
 1. Control signals to memory.
 2. Control signals to input–output modules.

Totally, there are three types of control signals:

1. Those that activate ALU function.
2. Those that activate a data path.
3. Those that are signals on the external system bus.

Functions of Control Unit

- The control unit directs the entire computer system to carry out stored program instructions.
- The control unit must communicate with both the Arithmetic Logic Unit and Main memory.
- The control unit instructs the arithmetic logic unit by which, logical or arithmetic operation is to be performed.
- The control unit coordinate the activities of the other two units as well as all peripheral and auxiliary storage devices linked to the computer.

Design of Control Unit

Control unit generates control signals using one of the two organizations

- (1) Hardwired control unit
- (2) Micro-programmed control unit.

Hardwired control unit

- It is implemented as logic circuits (gates, flip-flops, decoders, etc.) in the hardware.
- It is very complicated if we have a large control unit.
- In this organization, if the design has to be modified or changed. It requires changes in wiring among the various components. Thus the modification of all the combinational circuits may be very difficult.

Architecture of hardwired control unit An example hardwired control unit is shown in Figure 9.

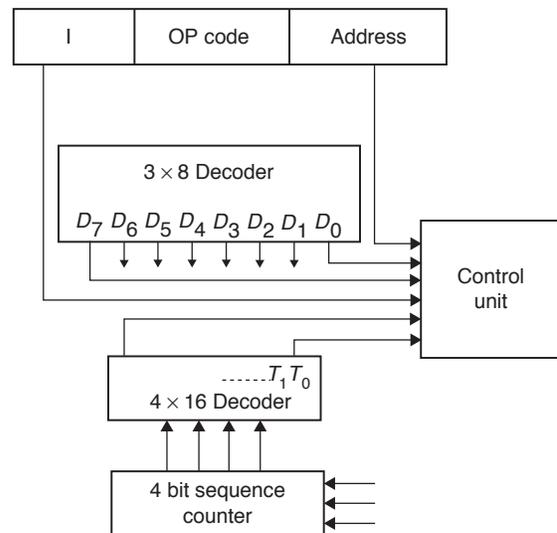


Figure 10 Hardwired control unit

The above control unit consists of:

- Instruction Register
- Number of control logic gates

- Two decoders
- 4-bit sequence counter.
- An instruction read from memory is placed in the instruction register (IR)
- The instruction register is divided into three parts: the I bit, operation code and Address part.
- First 12-bits (0-11) to specify an address, next 3-bits specify the operation code (op code) field of the instruction and last left most bit specify the addressing mode I .
 $I = 0$ for direct address
 $I = 1$ for indirect address
- First 12-bits are applied to the control logic gates.
- The Opcode bits (12-14) are decoded with 3×8 decoder.
- The eight outputs (D_0 through D_7) from a decoder go to the control logic gates to perform specific operation.
- Last bit 15 is transferred to a I flip flop designated by symbol I .
- The 4-bit sequence counter SC can count in binary from 0 through 15.
- The counter output is decoded into 16 timing pulses T_0 through T_{15} .
- The sequence counter can be incremented by INR input or clear by CLR input synchronically.

Advantages:

- Hardwired control unit is fast because control signals are generated by combinational circuits.
- The delay in generation of control signals depends upon the number of gates.

Disadvantages:

- More is the control signal required by CPU, more complex will be the design of control unit.

- Modifications in control signal are very difficult. That means it requires rearranging of wires in the hardware circuit.
- It is difficult to correct mistake in original design or adding new features.

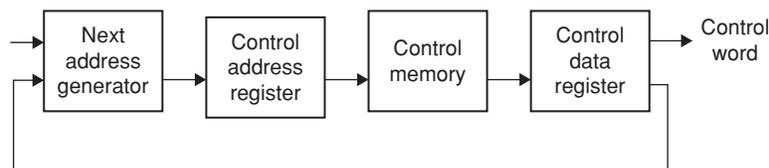
Micro-programming control unit

- A micro-programmed Control unit is implemented using programming approach. A sequence of micro-operations are carried out by executing a program consisting of microinstructions.
- Micro-program, consisting of micro instructions is stored in the control memory of the control unit.
- Execution of micro-instruction is responsible for generation of a set of control signals.

A micro-instruction consists of:

- One or more micro-instructions to be executed.
 - Address of next micro-instruction to be executed.
- (a) Micro-operations:** The operations performed on the Data stored inside the registers are called Micro-operations.
 - (b) Micro-programs:** Micro-programming is the concept for generating control signals using programs. These programs are called Micro-programs.
 - (c) Micro-instructions:** The instructions that make Micro-programs are called micro-instructions.
 - (d) Micro-code:** Micro-program is a group of micro-instructions. Micro-program can also be termed as micro-code.
 - (e) Control memory:** Micro-programs are stored in the read-only memory (ROM). That memory is called control memory.

(f) Architecture of Micro-Programmed Control Unit:



- The address of micro-instruction that is to be executed is stored in the control address register (CAR).
- Micro-instruction corresponding to the address stored in CAR is fetched from control memory and is stored in the control data register (CDR).
- This micro-instruction contains control word to execute one or more micro-operations.
- After the execution of all micro-operations of micro-instructions, the address of next micro-instructions is located.

Advantages:

- The design of micro-program control unit is less complex because micro-programs are implemented using software routines.

- The micro-programmed control unit is more flexible because design modifications, correction and enhancement is easily possible.
- The new or modified instruction set of CPU can be easily implemented by simply rewriting or modifying the contents of control memory.
- The fault can be easily diagnosed in the micro-program control unit using diagnostic tools by maintaining the contents of flags, registers and counters.

Disadvantages:

- The micro-program control unit is slower than hardwired control unit. That means to execute an instruction in micro-program control unit requires more time.

- The micro-program control unit is expensive than hardwired control unit in case of limited hardware resources.
- The design duration of micro-program control unit is more than hardwired control unit for smaller CPU.

Types of Micro-instructions

Micro-instructions can be classified as

Horizontal micro-instruction

- Individual bits in horizontal micro-instructions correspond to individual control lines.
- These are long and allow maximum parallelism since each bit controls a single control line.
- No decoding needed.

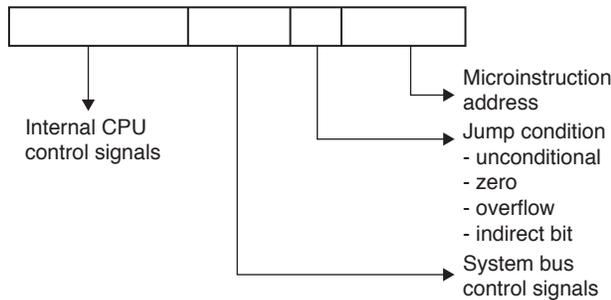


Figure 11 Horizontal micro-instruction format

Vertical micro-instruction

- Here, control lines are coded into specific fields within a micro-instruction.
- Decoders are needed to map a field of k -bits to 2^k possible combinations of control lines.

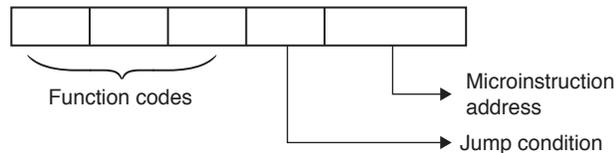


Figure 12 Vertical micro-instruction format

Example: A 3-bit field in a micro-instruction could be used to specify any one of eight possible lines.

- Hence these instructions are much shorter than horizontal ones.
- Control fields encoded in the same field cannot be activated simultaneously. Therefore vertical micro-instructions allow only limited parallelism.
- Decoding is necessary.

Micro-instruction Sequencing

Two concerns are involved in the design of a micro-instruction sequencing technique:

1. The size of micro-instruction: Minimizing size of control memory reduces the cost of that component.
2. The address-generation time:

A desire to execute micro-instructions as fast as possible. In executing a micro program, the address of next micro-instruction to be executed is in one of these categories.

1. Determined by IR
2. Next sequential address
3. Branch

Micro-instructions Execution

The Micro-instruction cycle has two parts:

1. Fetch
2. Execution

The effect of execution of a micro-instruction is to generate control signals. Some of the signals control points internal to the processor. The remaining signals go to the external control bus or other external interface.

Micro-instructions can be classified in a variety of ways.

1. Vertical/horizontal
2. Packed/unpacked
3. Hard/soft micro-programming
4. Direct/indirect encoding.

RISC AND CISC

One of the important aspects of computer architecture is the design of the instruction set for the processor. The instruction set chosen for a particular computer determines the way that machine language programs are constructed. There are two categories of computers based on instructions:

1. Complex instruction set computer (CISC)
2. Reduced instruction set computer (RISC)

CISC: A computer with a large number of instructions is classified as a complex instruction set computer.

RISC: A computer which has fewer instructions with simple constructs, so they can be executed much faster with in the CPU without having to use memory as often. This type of computer is classified as RISC.

CISC characteristics

- CISC provides a single machine instruction for each statement, That is written in a high level language so that compilation process is simplified and the over all computer performance improved.
- It has variable length instruction formats.
- It provides direct manipulation of operands residing in memory.
- Some instructions that perform specialized tasks and are used infrequently.
- A large variety of addressing modes

Drawback of CISC architecture As more instructions and addressing modes are incorporated into a computer, the more hardware logic is needed to implement and support them and hence this causes the computations to slow down.

RISC characteristics

- Reduce execution time by simplifying the instruction set of the computer.
- Fewer numbers of instructions
- Relatively fewer addressing modes
- Memory access is limited to load and store instructions.
- All operations are done with in the register of the CPU.
- Fixed - length, easily decoded instruction format.
- Single-cycle instruction execution.
- Hardwired rather than micro-programmed control.
- Relatively large number of registers.
- Uses overlapped register windows to speed - up procedure call and return.
- Efficient instruction pipeline.
- Efficient translation of high - level language programs into machine language programs by the compiler.

Example 6: An instruction set of a processor has 200 signals which can be divided into 5 groups of mutually exclusive signals as follows.

Group 1: 30 Signals
Group 2: 90 Signals

Group 3: 20 Signals

Group 4: 10 Signals

Group 5: 50 Signals

How many bits of the control words can be saved by using vertical micro-programming over horizontal microprogramming?

- (A) 27 (B) 173
(C) 200 (D) 227

Solution: Horizontal micro-programming requires 200 signals. But vertical micro-programming uses encoding. So

Group 1 requires 5-bits ($\therefore 2^5 = 32$)

Group 2 requires 7-bits ($\therefore 2^7 = 128$)

Group 3 requires 5-bits ($\therefore 2^5 = 32$)

Group 4 requires 4-bits ($\therefore 2^4 = 16$)

Group 5 requires 6-bits ($\therefore 2^6 = 64$)

\therefore Total bits required using vertical micro programming = 27

\therefore Number of bits saved = $200 - 27 = 173$

EXERCISE**Practice Problems I**

Directions for questions 1 to 20: Select the correct alternative from the given choices.

- Using two's complement arithmetic the resultant of $111100001111 - 110011110011$ is
(A) 0010 0001 1111
(B) 0011 0000 1100
(C) 0010 0001 1101
(D) 0010 0001 1100
- IEEE 32-bit floating point format of 384 is
(A) 0 10001111 000000000000000000000000
(B) 0 10001111 100000000000000000000000
(C) 0 00001000 000000000000000000000000
(D) 0 00001000 100000000000000000000000
- Consider the following IEEE 32-bit floating point number:
 $0\ 01111110\ 101000000000000000000000$.
What is the decimal value equivalent to given number?
(A) 0.25 (B) 3.25
(C) 0.8125 (D) 0.9375
- What would be the bias value for a base-8 exponent in a 7-bit field?
(A) 8 (B) 16
(C) 63 (D) 64
- The normalized value of the resultant of $8.844 \times 10^{-3} - 2.233 \times 10^{-1}$ is
(A) -2.144×10^{-1} (B) -0.2144
(C) -2×10^{-1} (D) -0.2

- Which of the following is the correct sequence of micro-operations to add a number to the AC when the operand is a direct address operand and store the final result to AC?

- (A) $MAR \leftarrow (IR(\text{address}))$
 $MBR \leftarrow \text{memory}$
 $R_1 \leftarrow (AC) + (MBR)$
- (B) $MAR \leftarrow IR(\text{address})$
 $MBR \leftarrow MAR$
 $R_1 \leftarrow (MBR)$
 $R_2 \leftarrow (AC) + (R_1)$
 $AC \leftarrow R_2$
- (C) $MAR \leftarrow (IR(\text{address}))$
 $MBR \leftarrow \text{Memory}(MAR)$
 $R_1 \leftarrow (MBR)$
 $R_2 \leftarrow (AC) + (R_1)$
 $AC \leftarrow (R_2)$
- (D) $MAR \leftarrow (IR(\text{address}))$
 $MBR \leftarrow \text{Memory}(MAR)$
 $AC \leftarrow (AC) + (MAR)$

Statement for linked answer questions 7 to 9: Assume that the control memory is 24 bits wide. The control portion of the micro-instruction format is divided into two fields. A micro-operation field of 13-bits specifies the micro-operation to be performed. An address selection field specifies a condition, based on the flags, that will cause a micro-instruction branch. There are eight flags.

7. How many bits are there in address selection field?
 (A) 1 (B) 2
 (C) 3 (D) 4
8. How many bits are there in address field?
 (A) 8 (B) 9
 (C) 13 (D) 24
9. What is the size of control memory in bits?
 (A) 256 (B) 768
 (C) 3328 (D) 6144
10. A simple processor has 3 major phases to its instructions cycle:
 1. Fetch
 2. Decode
 3. Execute
 Two 1-bit flags are used to specify the current phase in hardwired implementation. Will these flags required in micro-programming also?
 (A) Yes
 (B) No
 (C) Cannot predict
 (D) Depends on clock cycle time
11. In a 3-bus data path, the micro instructions format will be Opcode src1, src2, desti; The number of operations supported are 8 and the src1, src2 and desti require 20, 16 and 20 bits respectively.
 The total number of horizontal microinstructions specified will be
 (A) 2^{64} (B) 2^8
 (C) 2^{56} (D) 2^{61}
12. What is the smallest positive normalized number represented using IEEE single precision floating point representation?
 (A) 2^{-128} (B) $1 - 2^{-127}$
 (C) 2^{-127} (D) 2^{-126}
13. A micro program control unit is required to generate a total of 30 control signals. Assume that during any micro instruction, almost two control signals are active. Minimum number of bits required in the control word to generate the required control signals will be
 (A) 2 (B) 2.5
 (C) 10 (D) 12
14. What is the fraction field of the single-precision floating point representation of 6.25?
 (A) 1110 1000 0000 0000 0000 000
 (B) 1001 0000 0000 0000 0000 000
 (C) 1100 0000 0000 0000 0000 000
 (D) 0110 0100 0000 0000 0000 000
15. Let the total number of control signals generated are n , then what is the number of bits allocated in control field of vertical micro programming?
 (A) $n/2$ (B) n
 (C) 2^n (D) \log_2^n
16. In a micro programmed control unit, a control field of one address control instruction has to support two groups of control signals. In group1 it is required to generate either one or none of the 32 control signals. In group 2 at most 5 from the remaining, what will be the number of bits needed for the control field?
 (A) 8 (B) 10
 (C) 35 (D) 37
17. Assume that the exponent e is constrained to lie in the range $0 \leq e \leq x$, with a bias of q , that the base is b and that the significant is P -digits in length.
 What is the largest positive value that can be written is normalized floating point?
 (A) $b^{x-q}(1 - b^{-p})$ (B) b^{-q-1}
 (C) b^{-q-p} (D) $b^{x-q}(b^{-p} - 1)$
18. By using Booth's Multiplication algorithm. Below two numbers are multiplied:
 Multiplicand: 0111 0111 1011 1101
 Multiplier: 0101 1010 1110 1110
 How many additions/subtractions are required for the multiplication of the above two numbers?
 (A) 8 (B) 10
 (C) 13 (D) 7
19. Let us assume, we are multiplying two positive integers 1101 and 1011. The multiplicand M is 1101 and Multiplier Q is 1011. What is partial product after second cycle?
 (A) 0110 1101 (B) 1001 1110
 (C) 0100 1111 (D) 1000 1111
20. The decimal representation of the 2's complement number 1101011 is
 (A) 21 (B) -21
 (C) 219 (D) 91

Practice Problems 2

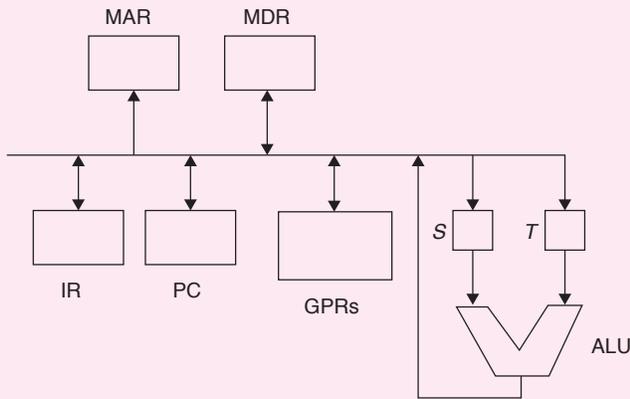
Directions for questions 1 to 20: Select the correct alternative from the given choices.

1. A microprogrammed control unit
 (A) is faster than a hard-wired control unit
 (B) facilitates easy implementation of new instructions.
 (C) is useful when very small programs are to be run.
 (D) usually refers to the control unit of a micro-processor.
2. Micro-program is
 (A) the name of a source program in micro computers.
 (B) a primitive form of macros used in assembly language programming.
 (C) a program of a very small size.
 (D) the set of instruction indicating the basic elemental commands which directly control the operation of a system.

3. Programming that actually controls the path of signal or data within the computer is called
 - (A) System programming
 - (B) Micro-programming
 - (C) High-level language programming
 - (D) Assembly language programming
4. The instruction cycle time in a generic microprocessor is
 - (A) Longer than the machine cycle time
 - (B) Shorter than the machine cycle time
 - (C) Same as the machine cycle time
 - (D) Double the machine cycle time
5. Microprocessor unit or central processor unit consist of
 - (A) Control circuitry
 - (B) ALU
 - (C) Memory
 - (D) All of these
6. The exponent of a floating point number is represented in excess-N code so that
 - (A) the dynamic range is large
 - (B) overflow is avoided
 - (C) the precision is high
 - (D) the smallest number is represented efficiently
7. Using Booth's algorithm for Multiplication, the Multiplier -14 is coded as
 - (A) 11110
 - (B) 01110
 - (C) 10010
 - (D) 00010
8. Data Path consists of
 - (A) Registers
 - (B) ALU
 - (C) Bus
 - (D) All of these
9. A floating point number that has a '0' in MSB of mantissa is said to have ____
 - (A) Overflow
 - (B) Underflow
 - (C) Normalization
 - (D) Positive exponent
10. Let the Binary sum after BCD addition is stored in K , Z_8 , Z_4 , Z_2 , and Z_1 Then the condition for a correction and output carry can be expressed as $C =$
 - (A) $K + Z_8 Z_4 + Z_8 Z_2$
 - (B) $K + Z_8 Z_4 + Z_4 Z_2$
 - (C) $K + Z_8 Z_2 + Z_8 Z_1$
 - (D) $K + Z_4 Z_2 + Z_2 Z_1$
11. Which of the following is an advantage of biased exponents?
 - (A) Convenient way to represent exponents
 - (B) Useful for conversion
 - (C) Convenient for comparison purposes
 - (D) All of these
12. Booth multiplication skips over runs of zeros and ones which reduces the number of add and subtract steps needed to multiply two n -bit numbers to n to a variable number whose average value n_{avg} is less than n what will be n_{avg} ?
 - (A) $n/3$
 - (B) $n/4$
 - (C) $n/2$
 - (D) n
13. The sequence of events that happen during a fetch operation is:
 - (A) $PC \rightarrow \text{memory} \rightarrow IR$
 - (B) $PC \rightarrow MAR \rightarrow \text{memory} \rightarrow IR$
 - (C) $PC \rightarrow MAR \rightarrow \text{memory} \rightarrow MDR \rightarrow IR$
 - (D) $PC \rightarrow \text{memory} \rightarrow MDR \rightarrow IR$
14. Micro-programming is a technique for
 - (A) Programming input or output routines
 - (B) Programming the microprocessors
 - (C) Programming the control steps of a computer
 - (D) Writing small programs
15. In a micro program ____ specifies the address of Micro-instructions to be executed.
 - (A) AR
 - (B) PC
 - (C) SP
 - (D) CAR
16. Which one of the following statements is correct?
 - (A) Micro-programmed control unit is costlier and slow.
 - (B) Micro-programmed control unit are cheap and slow.
 - (C) Micro-programmed control unit is costlier and fast.
 - (D) Micro-programmed control unit are fast and cheaper.
17. Horizontal micro-instructions have
 - (A) High degree parallelism, more encoding of control information.
 - (B) High degree parallelism, little encoding of control information.
 - (C) Low degree parallelism, more encoding of control information.
 - (D) Low degree parallelism, little encoding of control information.
18. A vertical micro-instruction have _____.
 - (A) Short formats and considerable encoding of control information
 - (B) Long formats and considerable encoding of control information
 - (C) Short formats and little encoding of control information
 - (D) Long formats and little encoding of control information
19. Guard bits are used to
 - (A) avoid unnecessary loss of MSB
 - (B) avoid unnecessary loss of LSB
 - (C) the loss of MSB
 - (D) the loss of LSB
20. Which of the following is not the essential element of a number represented in floating-point notation?
 - (A) Exponent
 - (B) Significand
 - (C) Sign
 - (D) Normalization

PREVIOUS YEARS' QUESTIONS

Common data for questions 1 and 2: Consider the following data path of a CPU.

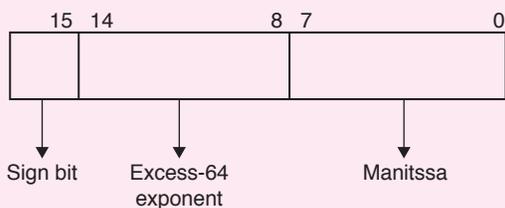


The ALU, the bus and all the registers in the data path are of identical size. All operations including incrementation of the PC and the GPRs are to be carried out in the ALU. Two clock cycles are needed for memory read operation—the first one for loading address in the MAR and the next one for loading data from the memory bus into the MDR.

- The instruction 'add R_0, R_1 ' has the register transfer interpretation $R_0 \leftarrow R_0 + R_1$. The minimum number of clock cycles needed for execution cycle of this instruction is
 (A) 2 (B) 3
 (C) 4 (D) 5
- The instruction 'call Rn, sub ' is a two word instruction. Assuming that PC is incremented during the fetch cycle of the first word of the instruction, its register transfer interpretation is
 $Rn \leftarrow PC + 1;$
 $PC \leftarrow M[PC]$

The minimum number of CPU clock cycles, needed during the execution cycle of this instruction is
 (A) 2 (B) 3
 (C) 4 (D) 5

Data for question 3: Consider the following floating-point format.



Mantissa is a pure fraction in sign-magnitude form.

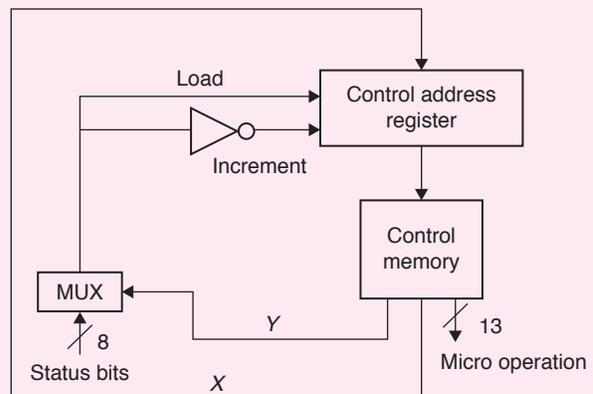
- The normalized representation for the above format is specified as follows. The mantissa has an implicit 1

preceding the binary (radix) point. Assume that only 0's are padded in while shifting a field.

The normalized representation of the above number (0.239×2^{13}) is: [2005]

- (A) 0A 20 (B) 11 34
 (C) 49 D0 (D) 4A E8

- In the IEEE floating point representation the hexadecimal value $0x00000000$ corresponds to [2008]
 (A) The normalized value 2^{-127}
 (B) The normalized value 2^{-126}
 (C) The normalized value $+0$
 (D) The special value $+0$
- P is a 16-bit signed integer. The 2's complement representation of P is $(F87B)_{16}$. The 2's complement representation of $8 * P$ is [2010]
 (A) $(C3D8)_{16}$ (B) $(187B)_{16}$
 (C) $(F878)_{16}$ (D) $(987B)_{16}$
- The decimal value 0.5 in IEEE single precision floating point representation has [2012]
 (A) fraction bits of 000...000 and exponent value of 0
 (B) fraction bits of 000...000 and exponent value of -1
 (C) fraction bits of 100...000 and exponent value of 0
 (D) no exact representation
- The smallest integer that can be represented by an 8-bit number in 2's complement form is [2013]
 (A) -256 (B) -128
 (C) -127 (D) 0
- Let $A = 1111\ 1010$ and $B = 0000\ 1010$ be two 8-bit 2's complement numbers. Their product in 2's complement is [2004]
 (A) 1100 0100 (B) 1001 1100
 (C) 1010 0101 (D) 1101 0101
- The microinstructions stored in the control memory of a processor have a width of 26 bits. Each microinstruction is divided into three fields, a micro-operation field of 13 bits, a next address field (X), and a MUX select field (Y), there are 8 status bits in the inputs of the MUX [2004]



How many bits are there in the X and Y fields, and what is the size of the control memory in number of words?

- (A) 10, 3, 1024 (B) 8, 5, 256
(C) 5, 8, 2048 (D) 10, 3, 512

10. Consider the following sequence of micro-operations.

$MBR \leftarrow PC$

$MAR \leftarrow X$

$PC \leftarrow Y$

$Memory \leftarrow MBR$

Which one of the following is a possible operation performed by this sequence? [2013]

- (A) Instruction fetch
(B) Operand fetch
(C) Conditional branch
(D) Initiation of interrupt service

11. For computers based on three-address instruction formats, each address field can be used to specify which of the following: [2015]

- (S_1) A memory operand
(S_2) A processor register
(S_3) An implied accumulator register

- (A) Either S_1 or S_2
(B) Either S_2 or S_3
(C) Only S_2 and S_3
(D) All of S_1 , S_2 and S_3

12. Let X be the number of distinct 16-bit integers in 2's complement representation. Let Y be the number of distinct 16-bit integers in sign magnitude representation. They $x - y$ is _____. [2016]

13. The n -bit fixed-point representation of an unsigned real number X uses f bits for the fraction part. Let $i = n - f$. The range of decimal values for X in this representation is [2017]

- (A) 2^{-f} to 2^i (B) 2^{-f} to $(2^i - 2^{-f})$
(C) 0 to 2^i (D) 0 to $(2^i - 2^{-f})$

- 14 Consider the C code fragment given below.

```
typedef struct node {
    int data;
    node* next;
} node;
void join (node* m, node* n) {
```

```
node* p = n;
while (p ->next != NULL) {
    p = p ->next;
}
p ->next = m;
}
```

Assuming that m and n point to valid NULL-terminated linked lists, invocation of join will [2017]

- (A) append list m to the end of list n for all inputs.
(B) either cause a null pointer dereference or append list m to the end of list n .
(C) cause a null pointer dereference for all inputs.
(D) append list n to the end of list m for all inputs.

15. The representation of the value of a 16-bit unsigned integer X in hexadecimal number system is BCA9. The representation of the value of X in octal number system is [2017]

- (A) 571244 (B) 736251
(C) 571247 (D) 136251

16. Consider the following processor design characteristics.

- I. Register-to-register arithmetic operations only
II. Fixed-length instruction format
III. Hardwired control unit

Which of the characteristics above are used in the design of a RISC processor? [2018]

- (A) I and II only (B) II and III only
(C) I and III only (D) I, II and III

17. Consider the unsigned 8-bit fixed point binary number representation below:

$$b_7 b_6 b_5 b_4 b_3 \cdot b_2 b_1 b_0$$

where the position of the binary point is between b_3 and b_2 . Assume b_7 is the most significant bit. Some of the decimal numbers listed below cannot be represented exactly in the above representation:

- (i) 31.500 (ii) 0.875
(iii) 12.100 (iv) 3.001

Which one of the following statements is true?

- [2018]

- (A) None of (i), (ii), (iii), (iv) can be exactly represented
(B) Only (ii) cannot be exactly represented
(C) Only (iii) and (iv) cannot be exactly represented
(D) Only (i) and (ii) cannot be exactly represented

ANSWER KEYS

EXERCISES

Practice Problems I

1. D 2. B 3. C 4. C 5. A 6. C 7. C 8. A 9. D 10. B
11. A 12. D 13. C 14. B 15. D 16. B 17. A 18. B 19. B 20. B

Practice Problems I

1. B 2. D 3. B 4. C 5. D 6. D 7. C 8. D 9. B 10. A
11. C 12. C 13. C 14. C 15. D 16. A 17. B 18. A 19. B 20. D

Previous Years' Questions

1. B 2. B 3. D 4. D 5. A 6. B 7. B 8. A 9. A 10. D
11. A 12. 1 13. D 14. B 15. D 16. D 17. C

Chapter 3

Memory Interface, I/O Interface

LEARNING OBJECTIVES

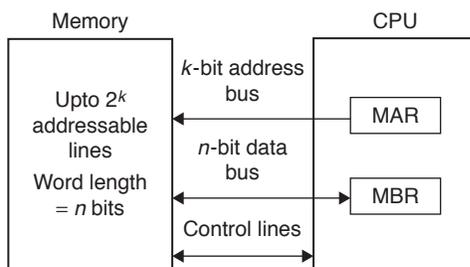
- Memory interface
- RAM
- ROM
- Memory interfacing
- Input-output interfacing
- Handshaking
- Data transfer mode
- Design techniques for interrupts
- Direct memory access
- Input-output processor

MEMORY INTERFACE

Basic Concepts

Computer memory is used to store programs and data. The maximum size of a memory that can be used in any computer is determined by the addressing scheme.

Example: If the memory address has 16-bits, then the size of memory will be 2^{16} Bytes.



If MAR is k -bits long and MDR is n -bits long, then the memory may contain up to 2^k addressable locations and the n -bits of data are transferred between the processor and memory. This transfer takes place over processor bus. The processor bus has

1. Address line
2. Data line
3. Control line

Control line is used for coordinating data transfer.

Processor reads the data from the memory by loading the address of the required memory location into MAR and setting the R/\bar{W} line to 1.

The memory responds by placing the data from the addressed location onto the data lines and confirms the actions. Upon confirmation, the processor loads the data onto the data lines, into MDR register. The processor writes the data into the memory location by loading the address of this location into MAR and loading the data into MDR sets the R/\bar{W} line to 0.

- **Memory Access Time:** It is the time that elapses between the initiation of an operation and the completion of that operation.
- **Memory Cycle Time:** It is the minimum time delay that required between the initiations of two successive memory operations.

RAM (Random Access Memory)

In RAM, if any location that can be accessed for a read/write operation in fixed amount of time, it is independent of the location's address:

- Memory cells are usually organized in the form of array, in which each cell is capable of storing one bit of information.
- Each row of cells constitutes a memory word and all cells of a row are connected to a common line called as word line.
- The cells in each column are connected to sense/write circuit by two bit lines.

The data input and data output of each sense/write circuit are connected to a single bidirectional data line that can be connected to a data bus.

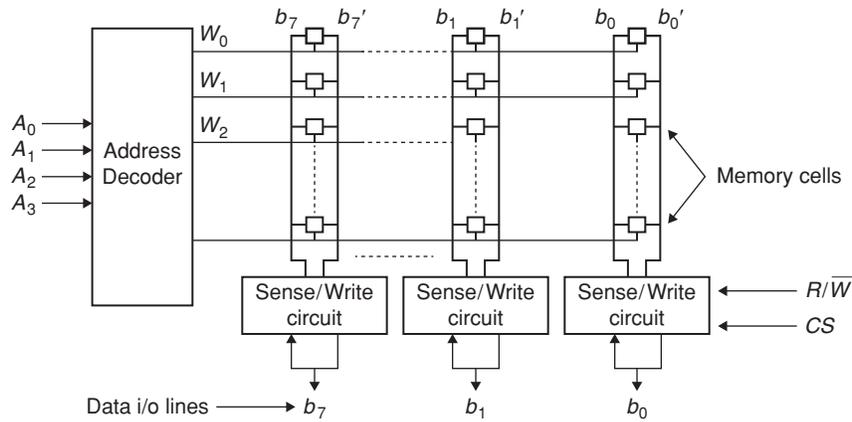


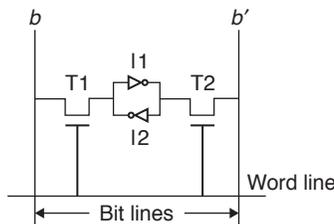
Figure 1 Organization of bit cells in a memory chip

- R/\overline{W} : Specifies the required operation.
- CS: Chip select input selects a given chip in the multi-chip memory system.

Static memories

Memories that consist of circuits capable of retaining their state as long as power is applied are known as static memories.

SRAM (static RAM) SRAM consists of two inverters, two transistors. In order to read the state of the SRAM cell, the word line is activated to close switches T1 and T2.



Advantages of SRAM:

1. It has low power consumption, because the current flows in the cell only when the cell is being activated or accessed.
2. SRAM can be accessed quickly.

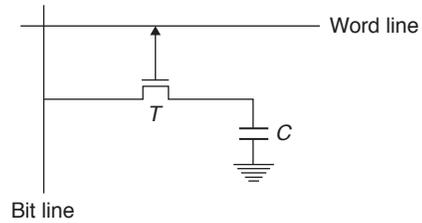
Disadvantages of SRAM: SRAMs are said to be volatile memories, because their contents are lost when the power is interrupted.

DRAM (Dynamic RAM) Less expensive RAMs can be implemented if simplex cells are used, such cells cannot retain their state indefinitely. Hence they are called dynamic RAMs.

The information stored in a dynamic memory cell in the form of a charge on a capacitor and this charge can be maintained only for tens of milliseconds.

The contents must be periodically refreshed by restoring the capacitor charge to its full value.

Example: Single-transistor dynamic memory cell:



If charge on capacitor > threshold value, then bit line will have '1'. If charge on capacitor < threshold value, then bit line will have '0'.

DRAM	SRAM
1. Volatile	1. Volatile
2. Simple to build and slower than SRAM	2. Faster than DRAM
3. Need refresh circuitry	3. More expensive to build
4. Favoured for large memory units	4. Favoured for cache memory units

Latency It is the amount of time it takes to transfer a word of data to or from the memory.

- For the transfer of a single word, the latency provides the complete indication of memory performance.
- For a block transfer, the latency denotes the time it takes to transfer the first word of data.

Bandwidth It is defined as the number of bits or bytes that can be transferred in one second.

Note: All dynamic memories have to be refreshed.

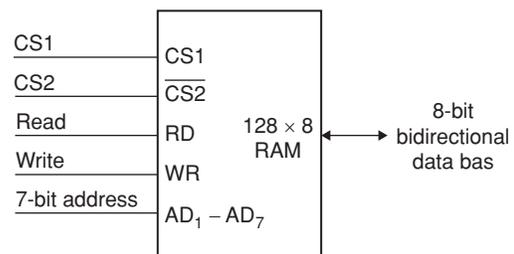


Figure 2 RAM chip block diagram

Read-only Memory (ROM)

Both SRAM and DRAM chips are volatile, which means that they lose the stored information if power is turned off. If the normal operation involves only reading of stored data, use ROM memory.

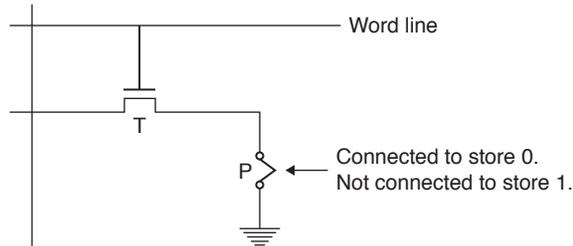


Figure 3 ROM cell

Types of ROM

Different types of non-volatile ROM are:

- PROM (Programmable ROM):**
 - Allows the data to be loaded by the user.
 - Less expensive, faster, flexible.
- EPROM (Erasable PROM):**
 - Allows the stored data to be erased and new data to be loaded.
 - Flexible, retain information for a long time.
 - Contents erased by UV light.
- EEPROM (Electrically Erasable PROM):**
 - Programmed and erased electrically.
 - Allows the erasing of all cell contents selectively.
 - Requires different voltage for erasing, writing and reading of stored data.
- Flash memory:** Allows to read the contents of a single cell but it is only possible to write the entire contents of a block.

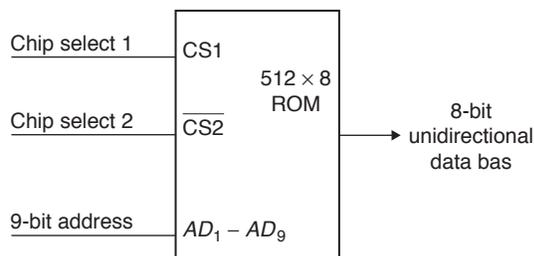


Figure 4 Block diagram of ROM chip

Memory Interfacing

The interfacing circuit enables the access of processor to memory. The function of memory interfacing is that the processor should be able to read from and write into a given

register of a memory chip. To perform this, the microprocessor should be

- able to select the chip.
- identify the register.
- enable the appropriate buffer.

INPUT-OUTPUT INTERFACING

Basic Concepts of I/O Module

I/O module contains logic for performing a communication function between the peripherals and the bus. The peripherals are not connected to the system bus directly. The reasons for this are

- Peripherals are electromechanical and electromagnetic devices and their manner of operation is different from the operation of the CPU and memory, which are electronic devices. So a conversion of signal values may be required.
- The data transfer rate of peripherals is usually slower than the transfer rate of the CPU and hence a synchronization mechanism may be needed.
- Data codes and formats in peripherals differ from the word format in the CPU and memory.
- The operating modes of peripherals are different from each other and each must be controlled so as not to disturb the operation of other peripherals connected to the CPU.

To resolve these differences, computer systems include special hardware components between the CPU and peripherals to supervise and synchronize all input and output transfers. These components are called 'interface' units.

By using this interfacing,

- interface to the processor and memory via the system bus or central switch.
- interface to one or more peripheral devices by tailored data links.

Input-output devices

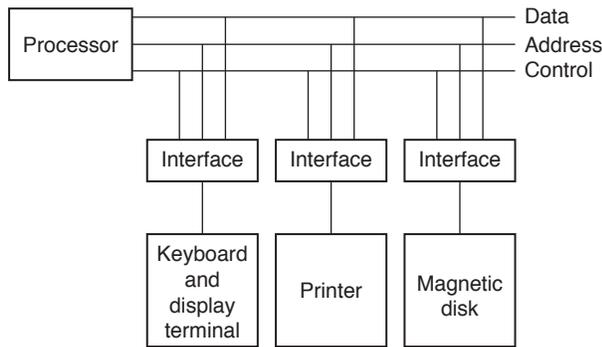
- Input and Output devices provide a means for people to make use of a computer.
- Some I/O devices function as an interface between a computer system and other physical system.

Input-output interface

Input/output Interface provides a method for transferring information between internal storage (such as memory and CPU Register) and external I/O devices. It resolves the difference between the computer and peripheral devices.

Input–output bus and interface modules

Each peripheral has an interface module associated with it. The interface module decodes the device address (device code), decodes signals for the peripheral controller, synchronizes the data flow and supervises the transfer rate between peripheral and CPU or memory.



Function of buses

- 1. Memory bus:** It is used for information transfer between CPU and main memory.
- 2. I/O bus:** It is used for information transfers between CPU and I/O devices through their I/O interface.

Isolated versus memory mapped I/O

- 1. Isolated I/O:**
 - Separate I/O read/write control lines in addition to memory read/write control lines.
 - Separate (isolated) memory and I/O address space
 - Distinct input and output instructions.
- 2. Memory-mapped I/O:**
 - A single set of Read/write control lines (i.e., no distinction between memory and I/O transfer).
 - Memory and I/O address share the common address space (reduces memory address range available).
 - No specific input or output instruction.
 - The same memory reference instructions can be used for I/O transfer.
 - Considerable flexibility in handling I/O operations.

Asynchronous serial transfer

In serial data transmission, each bit in the message is sent in sequence one at a time. Serial transmission can be synchronous or asynchronous.

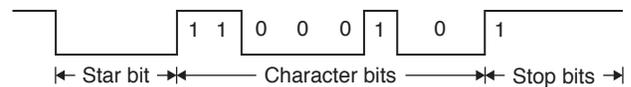
In synchronous transmission, the two units share a common clock frequency and bits are transmitted continuously at the rate dictated by the clock pulses.

In asynchronous transmission, binary information is sent only when it is available and the line remains idle when there is no information to be transmitted.

In serial asynchronous transmission technique, each character consists of three parts:

1. start bits
2. character bits
3. stop bits

Example:



A transmitted character can be detected by the receiver from the knowledge of the transmission rules:

1. When a character is not being sent, the line is kept in the 1-state.
2. The initiation of a character transmission is detected from the start bit, which is always 0.
3. The character bits always follow the start bit.
4. After the last bit of the character is transmitted, a stop bit is detected when the line returns to the 1-state for at least one bit time.
 - The baud rate is defined as the rate at which serial information is transmitted and is equivalent to the data transfer in bits per second.

Strobe control Employs a single control line to time each transfer. Strobe may be activated by either the source or the destination units.

(a) Source initiated transfer:

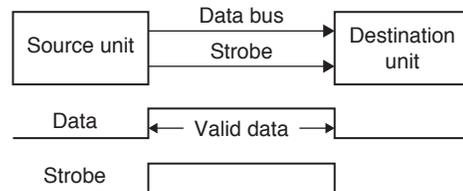
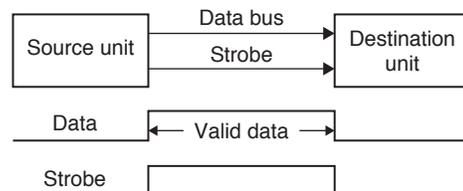


Figure 5 Source initiated strobe for data transfer

- The data bus carries the binary information from source unit to the destination unit.
- The strobe is a single line that informs the destination unit when a valid data word is available in the bus.

(b) Destination initiated strobe for data transfer:



- The destination unit activates the strobe pulse, informing the source to provide the data. The source unit responds by placing the requested binary information on the data bus.
- The data must be valid and remain in the bus long enough for the destination unit to accept it.
- The falling edge of the strobe pulse can be used again to trigger a destination register. The destination unit then disables the strobe.

Handshaking

Disadvantage of strobe method: Source unit which initiated the transfer has no way of knowing whether the destination unit has actually received the data item that was placed in the bus.

The handshake method solves this problem by introducing a second control signal that provides a reply to the unit that initiates the transfer.

Principle of two-wire handshaking: One control line is in the same direction as the data flow in the bus from the source to the destination. It is used by the source unit to inform the destination unit whether there are valid data in the bus.

The other control line is in the other direction from the destination to the source. It is used by the destination unit to inform the source whether it can accept data.

The sequence of control during the transfer depends on the unit that initiates the transfer.

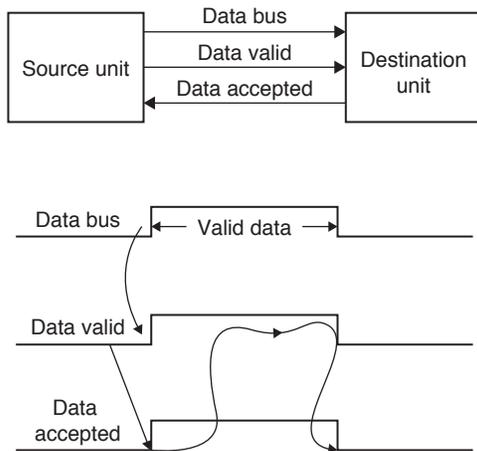


Figure 6 Source initiated transfer using hand shaking

Similarly a destination unit may also initiate the transfer.

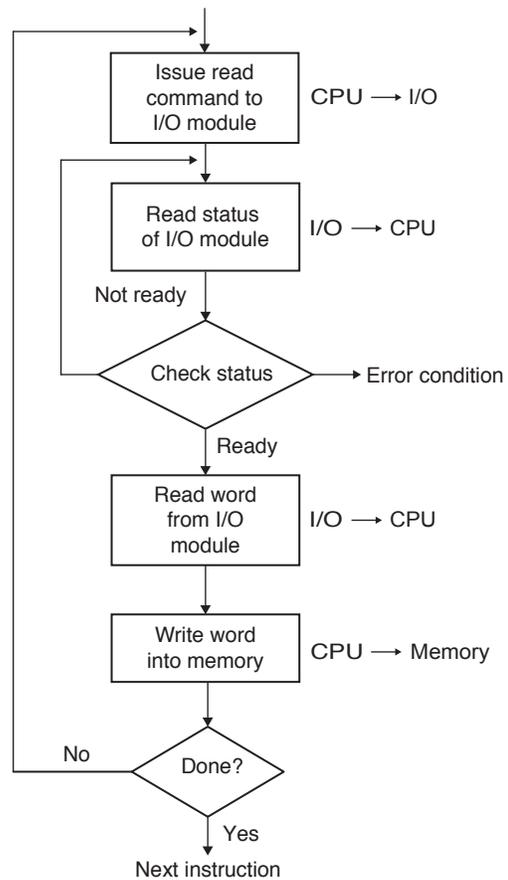
Advantage: Handshaking scheme provides a high degree of flexibility and reliability because the successful completion of a data transfer relies on active participation by both units.

Modes of transfer

There are three different data transfer modes between the central computer (CPU or Memory) and peripherals:

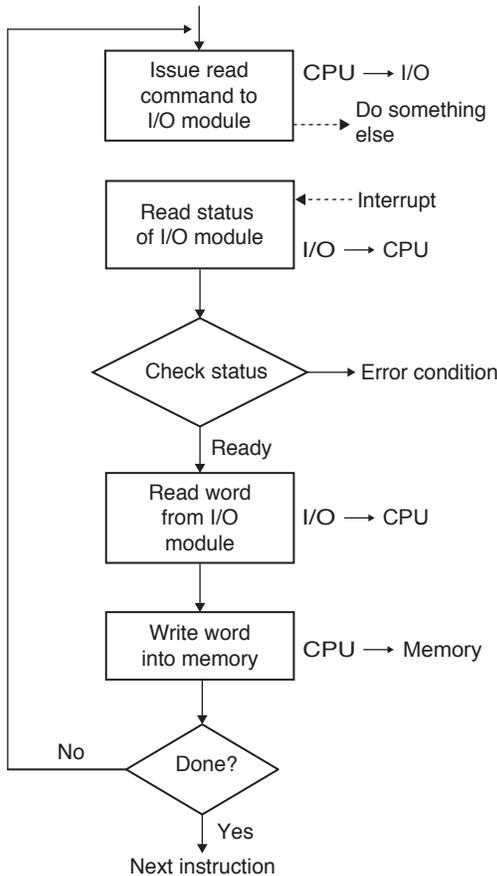
1. Program-controlled I/O
2. Interrupt-initiated I/O
3. Direct memory access

Program-controlled input-output With programmed I/O, the I/O module will perform the requested action and then set the appropriate bits in the I/O status register. The I/O module takes no further action to alert the CPU. In particular it does not interrupt the CPU. Thus, it is the responsibility of the CPU to periodically check the status of the I/O module until it finds that the operation is complete.



Interrupt initiated input-output The problem with programmed I/O is that the CPU has to wait a long time for the I/O module of concern to be ready for either reception or transmission of data. The CPU, while waiting must repeatedly interrogate the status of the I/O module. As a result, the level of the performance of the entire system is severely degraded.

An alternation is for the CPU to issue an I/O command to a module and then go on to do some other useful work. The I/O module will then interrupt the CPU to request service when it is ready to exchange data with the CPU. The CPU then executes the data transfer, as before and then resumes its former processing.



Interrupt Processing: In all computers, there is a mechanism by which the normal processing of the processor is interrupted by other modules like I/O, memory. The interrupts may be of the following class:

1. Program: Generated by some condition that occurs as a result of an instruction execution.
Examples: Arithmetic overflow, division by zero, etc.
2. Timer: Generated by timer within the processor.
3. I/O: Generated by an I/O controller, to signal normal completion of an operation or to signal a variety of error conditions.
4. Hardware failure: Generated by a failure such as power failure or memory parity error.

Interrupts are provided primarily as a way to improve processing efficiency.

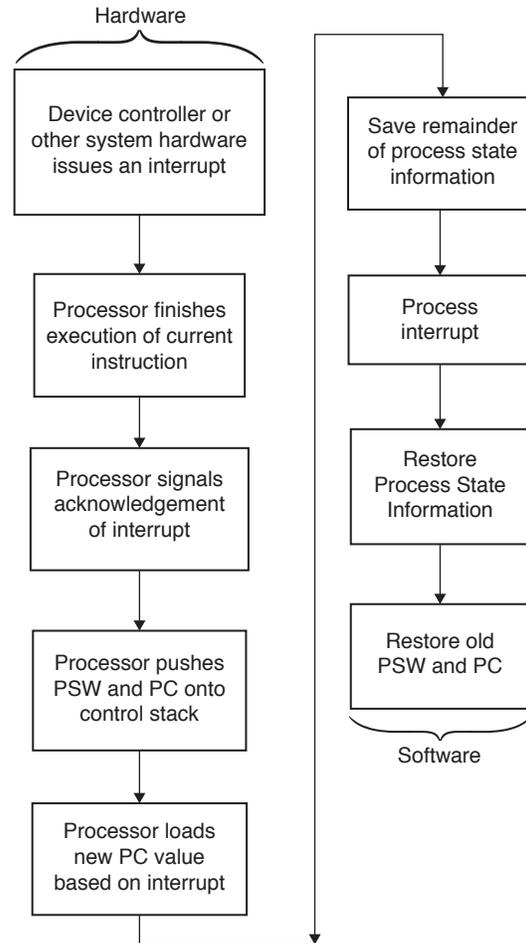


Figure 7 Interrupt Processing Flowchart

Consider the following figures, which show the contents of memory and registers before and after interrupt instruction processing.

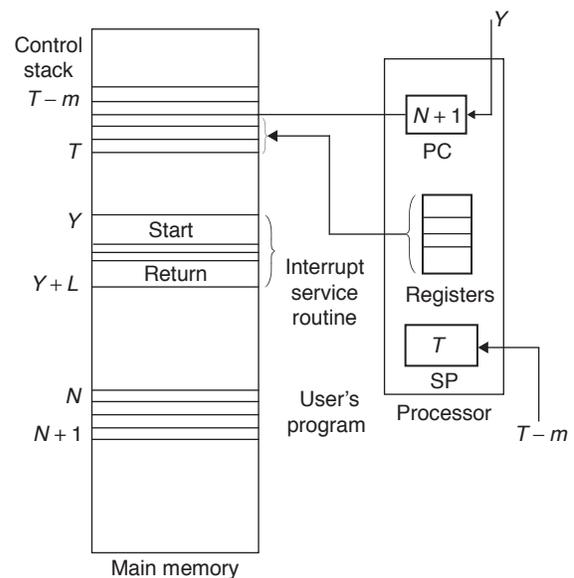


Figure 8 Interrupt occurs after instruction at location N

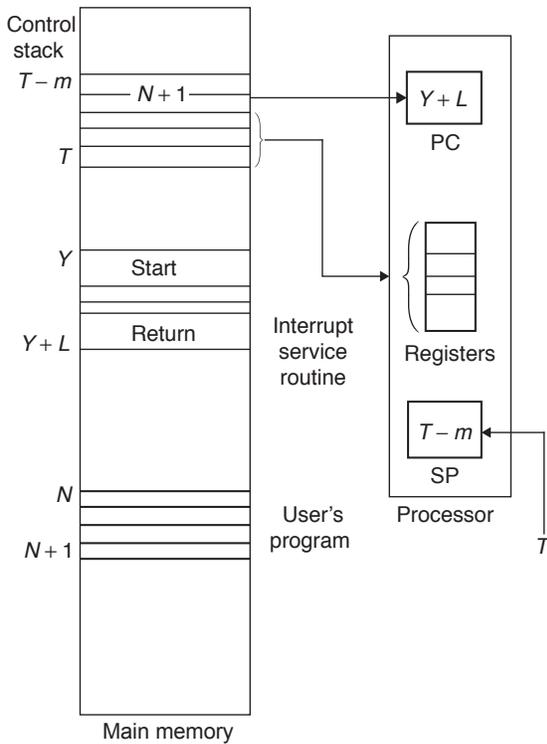


Figure 9 Return from interrupt

Interrupt priority: Priority determines which interrupt is to be served first when two or more requests are made simultaneously. Priority also determines which interrupts are permitted to interrupt the computer while another is being serviced. Higher priority interrupts can make requests while servicing a lower priority interrupt.

Design techniques for interrupts: Two design issues arise in implementing interrupt I/O:

1. Since there will almost invariably be multiple I/O module, how does the CPU determine which device issued the interrupt.
2. If multiple interrupts have occurred, how does the CPU decide which one to process.

Four general categories of techniques are there which are common in use:

- (a) Multiple interrupt lines:** In this technique, multiple interrupt lines are provided between the CPU and the I/O modules. However, it is impractical to dedicate more than a few bus lines or CPU pins to interrupt lines. Consequently, even if multiple lines are used, it is likely that each line will have multiple I/O modules attached to it. Thus, one of the other three techniques must be used on each line.
- (b) Software poll:** When the CPU detects an interrupt, it branches to an interrupt-service routine whose job is

to poll each I/O module to determine which module generated the interrupt. The poll could be in the form of a separate command line. The CPU receives the command and places the address of a particular I/O module on the address lines. The I/O module responds positively if it set the interrupt. Alternatively, each I/O module could contain an addressable status register. The CPU then reads the status register of each I/O module to identify the interrupting module. Once the correct module is identified, the CPU branches to a device service routine specified to that device. It is time consuming.

- (c) Daisy chain:** Daisy chain in effect provides a hardware poll. For interrupts all I/O modules share a common interrupt request line. The interrupt acknowledge line is daisy chained through the modules. When the CPU is interrupted, it sends out an interrupt acknowledgement. This signal propagates through a series of I/O modules until it gets to a requesting module. The requesting module typically responds by placing a word on the data lines. This word is referred to as a vector and is either the address of the I/O module or some other unique identifier. In either case, the CPU uses the vector as a pointer to the appropriate device-service routine. This avoids the need to execute a general interrupt-service routine first. This technique is referred to as a vectored Interrupt.
- (d) Bus arbitration:** Bus arbitration is also another technique which makes use of vectored Interrupts. With bus arbitration, an I/O module must first gain control of the bus before it can raise the interrupt request line. Thus only one module can raise the line at a time. When the CPU detects the interrupt, it responds on the interrupt acknowledge line. The requesting module then places its vector on the data lines.

Direct Memory Access (DMA)

Drawbacks of programmed and interrupt-driven I/O:

1. The I/O transfer rate is limited by the speed with which the processor can test and service a device.
2. The processor is tied up in managing I/O transfer; a number of instructions must be executed for each I/O transfer.

Both methods have an adverse impact on both processor activity and I/O transfer rate.

When large volume of data is to be moved, a more efficient technique is required: Direct Memory Access (DMA).

DMA function: DMA involves an additional module on the system bus.

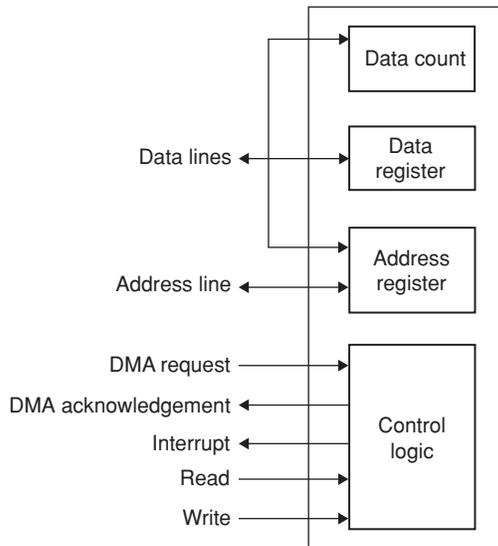


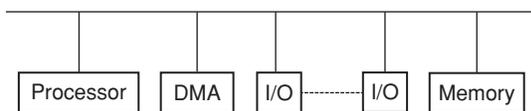
Figure 10 DMA Block Diagram

The DMA module is capable of mimicking the processor and indeed, of taking over control of the system from the processor. It needs to do this to transfer data to and from memory over the system bus. For this purpose, the DMA module must use the bus only when the processor does not need it or it must force the processor to suspend operation temporarily. The latter technique is more common and is referred to as cycle stealing.

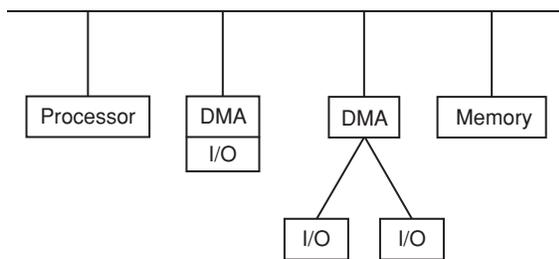
DMA configurations

1. Single bus, detached DMA

- Inexpensive, inefficient
- Each transfer of a word consumes two bus cycles.

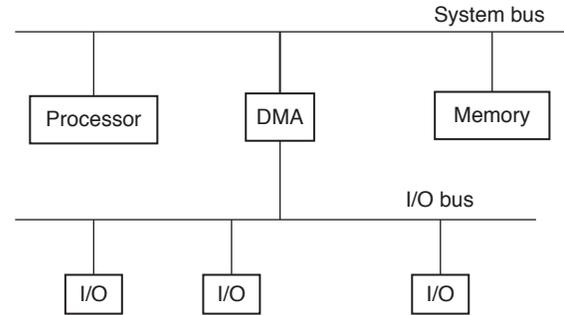


2. Single bus, integrated DMA I/O There is a path between the DMA module and one or more I/O modules that does not include system bus.



3. I/O bus

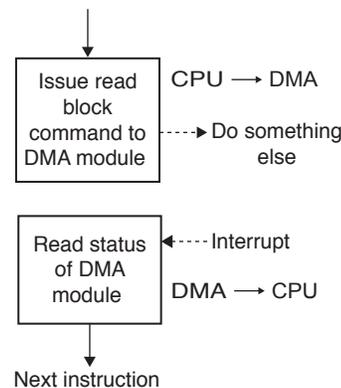
- Reduces the number of I/O interfaces in the DMA module to one.
- Easily expandable configuration.



With DMA, when the CPU wishes to read or write a block of data, it issues a command to the DMA module, by sending the following information to the DMA module.

1. Whether a read or write is requested.
2. The address of the I/O device involved.
3. The starting location in memory to read from or write to.
4. The number of words to be read or written.

The CPU then continues with other work. It has delegates this I/O operation to the DMA module, and that module will take care of it. The DMA module transfers the entire block of data, one word at a time, directly to or from memory, without going through the CPU. When the transfer is complete, the DMA module sends an interrupt signal to the CPU. Thus, the CPU is involved only at the beginning and end of the transfer.



DMA transfer can either happen as:

1. Burst transfer

- A block sequence consisting of a number of memory words is transferred in continuous burst.
- DMA controller is master of memory Buses.
- This mode of transfer is needed for fast devices such as Magnetic Disks, where transmission cannot be stopped or slowed down.

2. Cycle stealing

- CPU is usually much faster than I/O (DMA), thus CPU uses the most of the memory cycles.
- DMA controller steals the memory cycles from CPU.
- For those stolen cycles, CPU remains idle.

- For those slow CPU, DMA Controller may steal most of the memory.
- Cycle stealing, which may cause CPU remain idle long time.

Input–Output Processor (IOP)

An IOP is a processor, having a direct memory access capability, used to communicate with I/O devices.

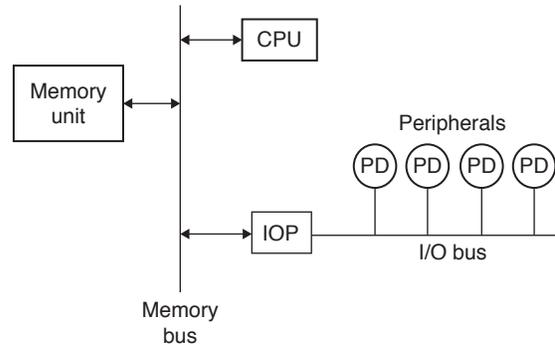
In this configuration, the computer system can be divided into a memory unit and a number of processors comprised of the CPU and one or more IOPs.

Each IOP takes care of input and output tasks, relieving the CPU from the house keeping chores involved in I/O transfers.

- IOP is similar to a CPU except that it is designed to handle the details of I/O processing.

- Unlike DMA, the IOP can fetch and execute its own instructions.

The following figure shows a computer with two processors:



EXERCISES

Practice Problems I

Directions for questions 1 to 20: Select the correct alternative from the given choices.

- Consider a DRAM that must be given a refresh cycle 64 times per ms. Each refresh operation requires 150 ns, a memory cycle requires 250 ns. What is the approximate percentage of the memory's total operating time must be given to refreshes?

(A) 1%	(B) 2%
(C) 9%	(D) 60%
- A DMA controller transfers 16-bit words to memory using cycle stealing. The words are assembled from a device that transmits characters at a rate of 2400 characters per second. The CPU is fetching and executing instructions at an average rate of 1 million instructions per second. By how much time will the CPU be slowed down because of the DMA transfer?

(A) 0.6%	(B) 0.1%
(C) 0.12%	(D) 0.24%
- A system is based on a 16-bit microprocessor and has two I/O devices. The I/O controllers for this system use separate control and status registers. Both devices handle data on a one-byte-at-a time basis. The first device has two status lines and three control lines. The second device has three status lines and four control lines. How many 16-bit I/O control module registers do we need for status reading and control of each device?

(A) 1, 2	(B) 2, 1
(C) 2, 2	(D) 1, 1
- In a programmed I/O technique, the processor is stuck in a wait loop doing status checking of an I/O device. To increase efficiency, the I/O software could be written so that the processor periodically checks the status of the device. If the device is not ready, the processor can jump to other tasks. After some timed interval, the processor comes back to check status again. Let us assume that above scheme is used for outputting data one character at a time to a printer that operates at 10 characters per second (CPS). Which of the following statement is true if its status is scanned every 200 ms?

(A) The printing speed is increased by 5 CPS
(B) The printing rate is slowed to 5 CPS
(C) The printing rate is at 10 CPS only
(D) The printing rate is at 20 CPS
- Consider a system employing interrupt-driven I/O for a particular device that transfers data at an average of 8 KB/s on a continuous basis. The interrupt processing takes about 100 μ s and the I/O device interrupts processor for every byte. Let assume that the device has two 16-byte buffers and interrupts the processor when one of the buffer is full. While executing the ISR, the processor takes about 8 μ s for the transfer of each byte. Then what is the fraction of processor time is consumed by this I/O device?

(A) 8%	(B) 11%
(C) 50%	(D) 65%
- A 32-bit computer has two selector channels one multiplexor channel. Each selector channel supports two magnetic disks and three magnetic tape units. The multiplexor channel has two line printers, two card readers and 10 VDT terminals connected to it. Assume the following transfer rates:

Disk drive: 1000 KB/sec
Magnetic tape drive: 300 KB/sec
Line printer: 6.2 KB/sec
Card reader: 2.4 KB/sec
VDT: 1 KB/sec

- What is the maximum aggregate I/O transfer rate of this system?
- (A) 1625.6 KB/s (B) 1327.2 KB/s
(C) 2027.2 KN/s (D) 2327.2 KB/s
7. Consider a disk drive with 16 surfaces, 512 tracks per surface and 512 sectors per track, 1 kilo bytes per sector and a Rotation speed of 3000 RPM. The disk is operated in cycle stealing mode where by whenever one 4 byte word is ready it is sent to memory; similarly, for writing, the disk interface read a 4 byte word from the memory in each DMA cycle. The memory cycle time is 40 nsec. Find the maximum percentage of time that the CPU gets blocked during DMA operation?
- (A) 2.62% (B) 26.21%
(C) 0.26% (D) 0.52%
8. How many RAM chips of size (256 K × 1-bit) are needed to build a 1 M Byte memory?
- (A) 16 (B) 8
(C) 32 (D) 24
9. Four memory chips of 16 × 4 size have their address bases connected together. The whole system will have a size of
- (A) 16 × 8 (B) 64 × 64
(C) 16 × 16 (D) 256 × 1
10. In which of following I/O techniques, there will be no interrupt?
- (A) Programmed I/O (B) Interrupt-driven I/O
(C) DMA (D) Both (B) and (C)
11. The capacity of a memory unit is defined by the number of words multiplied by the number of bits/word. How many separate address and data lines are needed for a memory 16K × 16?
- (A) 10 address, 4 data lines
(B) 14, 4
(C) 14, 16
(D) 14, 14
12. The main problem of strobe asynchronous data transfer is
- (A) it employs a single control line
(B) it is controlled by clock pulses in the CPU.
(C) the falling edge again used to trigger
(D) no way of knowing whether the destination has received the data item.
13. Which of the following DMA transfer modes and interrupt handling mechanisms will enable the highest I/O bandwidth?
- (A) Block transfer and polling interrupt
(B) Cycle stealing and polling interrupt
(C) Block transfer and vectored interrupt
(D) Transparent DMA and vectored interrupt
14. Which of the following enables peripherals to pass a signal down the bus to the next device on the bus during polling of the device?
- (A) Interrupt vectoring (B) Cycle stealing
(C) DMA (D) Daisy chain
15. What will be the response of the CPU, on receiving an interrupt from an input/output device?
- (A) It hands over the control of address bus and data bus to the interrupting device.
(B) It branches off to the interrupt service routine after completion of the current instruction.
(C) It halts for a predetermined time.
(D) It branches off to the interrupt service routine immediately.
16. What is the bandwidth of memory system that has a latency of 50ns, a pre charge time of 10ns and transfers 2 bytes of data per access?
- (A) 60 B/sec (B) 1.67 B/sec
(C) 1.67×10^7 B/sec (D) 3.33×10^7 B/sec
17. A hard disk is connected to a 50MHz processor through a DMA controller. Assume that the initial set-up of a DMA transfer takes 2000 clock cycles for the processor and also assume that the handling of the interrupt at DMA completion requires 1000 clock cycles for the processor. The hard disk has a transfer rate of 4000 K bytes/sec and average block size transferred is 8 K bytes. What fraction of the processor time is consumed by the disk, if the disk is actively transferring 100% of the time?
- (A) 1% (B) 1.5%
(C) 2% (D) 3%
18. A device with transfer rate of 20KB/sec is connected to a CPU. Data is transferred byte wise. Let the interrupt overhead is 6 micro seconds. The byte transfer time between the device interface register and CPU or memory is negligible. What is minimum performance gain of operating the device under interrupt mode over operating it under program-controlled mode?
- (A) 6 (B) 8
(C) 10 (D) 12
19. A DMA module is transferring characters to main memory from an external device at 76800 bits per second. The processor can fetch instructions at a rate of 2 million instructions per second. How much will the processor be slowed down due to DMA activity? (Express this as a percent of the time from when there is a conflict between DMA and the CPU)
- (A) 0.24% (B) 0.48%
(C) 0.96% (D) 0.50%
20. Let us suppose that we want to read 2048 bytes in programmed I/O mode of CPU. The bus width is 32-bits. Each time an interrupt occurs from Hard disk drive and it taken 4 μsec to service it. How much CPU time is required to read 2048 bytes?
- (A) 512 msec (B) 768 msec
(C) 1024 msec (D) 2048 msec

Practice Problems 2

Directions for questions 1 to 21: Select the correct alternative from the given choices.

1. Memory which is ultraviolet erasable is
(A) RAM (B) EPROM
(C) PROM (D) EEPROM
2. Memory which is electrically erasable is
(A) EPROM (B) EEPROM
(C) ROM (D) PROM
3. The minimum time delay that is required between the initiation of two successive memory operations is called
(A) Memory access time (B) Transmission time
(C) Seek Time (D) Memory cycle
4. The memory that is programmed at the time of manufacture is
(A) RAM (B) PROM
(C) ROM (D) EEPROM
5. The disadvantage of dynamic RAM over static RAM is
(A) High power consumption
(B) Higher bit density
(C) Need to refresh the capacitor charge every once in two milliseconds.
(D) Variable speed
6. If an error is detected, a part of the memory can be erased in
(A) PROM (B) EPROM
(C) EAROM (D) EROM
7. What are sequences of events in source initiated hand shaking transfer?
(A) Source enable data valid, destination enable data accepted, source disable data valid, destination disable data accepted
(B) Source disable data valid, destination enable data accepted, source disable data valid, destination enable data accepted
(C) Source disable data valid, destination Disable data valid, source enable data valid, destination enable data accepted.
(D) Source disable data valid, destination enable data valid, source enable data valid, destination disable data accepted.
8. Processor needs software interrupt to
(A) return from subroutine
(B) implement co-routines
(C) test the interrupt system of the processor
(D) obtain system services which need execution of privileged instructions
9. A microcomputer has primary memory of 512 KB. what is the exact number of bytes contained in this memory?
(A) 512×1000 (B) 512×100
(C) 512×1024 (D) 512×1028
10. The number of address lines required in a microprocessor which has to access 1 K bytes of memory is
(A) 6 (B) 4
(C) 10 (D) 8
11. Software interrupt is
(A) used to stimulate an external device
(B) generated by an external device
(C) Both (A) and (B)
(D) None of these
12. The bus that is used to transfer data from main memory to peripheral devices and vice-versa is
(A) Control bus (B) input bus
(C) output bus (D) DMA bus
13. The bus which is connected between the CPU and the main memory that permits transfer of information between the CPU and main memory is called
(A) memory bus (B) address bus
(C) control bus (D) DMA bus
14. An interrupt in which the external device supplies the interrupt requests as well as its address is called
(A) maskable interrupt
(B) vectored interrupt
(C) designated interrupt
(D) non-maskable interrupt
15. A temporarily ignored interrupt is called
(A) designated interrupt
(B) maskable interrupt
(C) non-maskable interrupt
(D) low priority interrupt
16. Which of the following device is used to connect a peripheral to a bus?
(A) control register
(B) interface
(C) communication protocol
(D) None of these
17. Which of the following is true for the daisy scheme of connecting input/output devices?
(A) It gives non-uniform priority to various devices.
(B) It gives uniform priority to all devices.
(C) It is only useful for connecting slow devices to a processor device.
(D) It requires a separate interrupt pin on the processor for each device.
18. In direct memory access data are directly transferred
(A) from CPU to input/output device and memory
(B) from an input/output device to memory only.
(C) from memory to an input/output device only.
(D) from an input/output device to the memory or vice versa

19. Which one of the following is true for a CPU having a single interrupt request line and a single interrupt grant line?
- (A) Vectored interrupt multiple interrupting devices are always possible.
 - (B) Vectored interrupts are not possible but multiple interrupting devices are possible
 - (C) Vectored interrupts and multiple interrupting devices are sometimes possible
 - (D) Vectored interrupt is possible but multiple interrupting devices are not possible
20. In which of the following I/O, there is a single address space for memory locations and I/O devices
- (A) Isolated I/O
 - (B) Memory mapped I/O
 - (C) DMA
 - (D) Both (A) and (B)
21. ____ signal used to interrupt processor and to execute service routine that takes an error recovery action.
- (A) Strobe
 - (B) Handshaking
 - (C) Polling
 - (D) Time out

PREVIOUS YEARS' QUESTIONS

1. A device with data transfer rate 10 KB/sec is connected to a CPU. Data is transferred byte-wise. Let the interrupt overhead be 4 μ sec. The byte transfer time between the device interface register and CPU or memory is negligible. What is the minimum performance gain of operating the device under interrupt mode over operating it under program-controlled mode? **[2005]**
- (A) 15
 - (B) 25
 - (C) 35
 - (D) 45
2. A computer handles several interrupt sources of which the following are relevant for this question.
- Interrupt from CPU temperature sensor (raises interrupt if CPU temperature is too high)
 - Interrupt from Mouse (raises interrupt if the mouse is moved or a button is pressed)
 - Interrupt from Keyboard (raises interrupt when a key is pressed or released)
 - Interrupt from Hard Disk (raises interrupt when a disk read is completed)

Which one of these will be handled at the HIGHEST priority? **[2011]**

- (A) Interrupt from Hard Disk
- (B) Interrupt from Mouse
- (C) Interrupt from Keyboard
- (D) Interrupt from CPU temperature sensor

The following information pertains to 3 and 4: Consider the following program segment for a hypothetical CPU having three user registers R_1 , R_2 and R_3 .

Instruction	Operation	Instruction Size (in words)
MOV R_1 , 5000;	$R_1 \leftarrow$ Memory [5000]	2
MOV R_2 (R_1);	$R_2 \leftarrow$ Memory [(R_1)]	1
ADD R_2 , R_3 ;	$R_2 \leftarrow R_2 + R_3$	1
MOV 6000, R_2 ;	Memory [6000] $\leftarrow R_2$	2
HALT;	Machine halts	1

3. Consider that the memory is byte addressable with size 32 bits, and the program has been loaded starting from memory location 1000 (decimal). If an interrupt occurs while the CPU has been halted after executing the HALT instruction, the return address (in decimal) saved in the stack will be **[2004]**
- (A) 1007
 - (B) 1020
 - (C) 1024
 - (D) 1028
4. Let the clock cycles required for various operations be as follows:
- Register to/from memory transfer: 3 clock cycles
 ADD with both operands in register: 1 clock cycle
 Instruction fetch and decode: 2 clock cycles per word
 The total number of clock cycles required to execute the program is **[2004]**
- (A) 29
 - (B) 24
 - (C) 23
 - (D) 20
5. A CPU generally handles an interrupt by executing an interrupt service routine **[2009]**
- (A) As soon as an interrupt is raised.
 - (B) By checking the interrupt register at the end of fetch cycle.
 - (C) By checking the interrupt register after finishing the execution of the current instruction.
 - (D) By checking the interrupt register at fixed time intervals.
6. A hard disk with a transfer rate of 10 Mbytes/second is constantly transferring data to memory using DMA. The processor runs at 6000 MHz, and takes 300 and 900 clock cycles to initiate and complete DMA transfer respectively. If the size of the transfer is 20 Kbytes, what is the percentage of processor time consumed for the transfer operation?
- (A) 5.0%
 - (B) 1.0%
 - (C) 0.5%
 - (D) 0.1%
7. On a non-pipelined sequential processor, a program segment, which is a part of the interrupt service routine, is given to transfer 500 bytes from an I/O device to memory.

ANSWER KEYS

EXERCISES

Practice Problems 1

- | | | | | | | | | | |
|-------|-------|-------|-------|-------|-------|-------|-------|-------|-------|
| 1. A | 2. C | 3. D | 4. B | 5. B | 6. C | 7. B | 8. C | 9. C | 10. A |
| 11. C | 12. D | 13. A | 14. D | 15. B | 16. D | 17. D | 18. B | 19. B | 20. D |

Practice Problems 2

- | | | | | | | | | | |
|-------|-------|-------|-------|-------|-------|-------|-------|-------|-------|
| 1. B | 2. B | 3. A | 4. C | 5. A | 6. C | 7. B | 8. D | 9. C | 10. C |
| 11. A | 12. D | 13. A | 14. B | 15. B | 16. B | 17. A | 18. D | 19. B | 20. B |
| 21. D | | | | | | | | | |

Previous Years' Questions

- | | | | | | | | | | |
|---------|-------|--------------|------|------|------|------|------|------|--------|
| 1. B | 2. D | 3. D | 4. B | 5. C | 6. D | 7. A | 8. A | 9. D | 10. 31 |
| 11. 456 | 12. A | 13. 59 to 60 | | | | | | | |

Chapter 4

Instruction Pipelining

LEARNING OBJECTIVES

- ☞ Flynn's classification
- ☞ Pipelining
- ☞ Six-stages of pipelining
- ☞ Pipeline performance
- ☞ Pipeline hazards
- ☞ Structural hazards
- ☞ Data hazards
- ☞ Control hazards
- ☞ Conditional branch
- ☞ Dealing with branches

FLYNN'S CLASSIFICATION

In parallel processing, the system is able to perform concurrent data processing to achieve faster execution time. A classification introduced by M.J. Flynn considers the organization of a computer system by the number of instructions and data items that are manipulated simultaneously. According to this classification, there are four major groups of computers:

- 1. Single instruction stream, single data stream (SISD):**
 - Single computer containing a control unit, a processor unit and a memory unit
 - Instructions executed sequentially; parallel processing may be achieved by multiple functional units or by pipeline processing.
- 2. Single instruction stream, Multiple data stream (SIMD):**
 - Include many processing units under the supervision of a common control unit.
 - All processors receive the same instruction from the control unit but operate on different items of data.
- 3. Multiple instruction stream, Single data stream (MISD):**
 - No practical system has been constructed.
- 4. Multiple instruction stream, Multiple data stream (MIMD):**
 - Capable of processing several programs at the same time.

One type of parallel processing that does not fit Flynn's classification is pipelining.

PIPELINING

Pipelining is a technique of decomposing a sequential process into sub operations, with each subprocess being executed in special dedicated segment that operates with all other segments.

In pipelining, new inputs are accepted at one end before previously accepted inputs appear as outputs at the other end.

Two-stage Pipeline

As a simple approach, consider subdividing instruction processing into two stages:

1. Fetch instruction
2. Execute instruction
 - There are times during the execution of an instruction when main memory is not being accessed. This time can be used to fetch the next instruction in parallel with the execution of the current one. This is called instruction prefetch or fetch overlap.
 - This process will speed up instruction execution. If the fetch and execute stages were of equal duration, the instruction cycle time would be halved.
 - But there are some problems in this technique:
 - (i) The execution time is generally longer than the fetch time.
 - (ii) A conditional branch instruction makes the address of the next instruction to be fetched unknown.
 - These two factors reduce the potential effectiveness of the two-stage pipeline, but some speed up occurs. To gain further speed up, the pipeline must have more stage(s).

Six-stage Pipeline

Let the six stages be

- F: Fetch Instruction
- D: Decode Instruction
- C: Calculate Operand Address
- O: Operand Fetch

E: Execute instruction
W: Write operand

Then the time line diagram for seven instructions is shown below:

Clock cycle	1	2	3	4	5	6	7	8	9	10	11	12
Instruction <i>i</i>	F	D	C	O	E	W						
<i>i</i> + 1		F	D	C	O	E	W					
<i>i</i> + 2			F	D	C	O	E	W				
<i>i</i> + 3				F	D	C	O	E	W			
<i>i</i> + 4					F	D	C	O	E	W		
<i>i</i> + 5						F	D	C	O	E	W	
<i>i</i> + 6							F	D	C	O	E	W

Execution Time for the seven instructions with pipelining = $\left(\frac{t_{ex}}{6}\right) \times 12 = 2 * t_{ex}$. Where t_{ex} in the execution time required for each instruction.

- A deeper pipeline means that there are more stages in the pipeline. This generally means that the processor’s frequency can be increased as the cycle time is lowered. This happens because there are fewer components in each stage of the pipeline, so the propagation delay is decreased for the overall stage.
- An instruction pipeline is said to be fully pipelined if it can accept a new instruction in every clock cycle.
- A pipeline that is not fully pipelined has wait cycle that delays the progress of the pipeline.

Advantages of pipelining:

- The cycle time of the processor is reduced, thus increasing instruction issue rate in most cases.
- Some combinational circuits such as adders or multipliers can be made faster by adding more circuitry. If pipelining is used instead it can save circuitry versus a more complex combinational circuit.

Limitations of pipelining:

1. If the stages are not of equal duration, there will be some waiting involved at various stages.
2. Conditional branch instruction may invalidate several instruction fetches.
3. The contents of one stage may depend on the contents of other stages of previous instructions, which is still in pipeline.

PIPELINE PERFORMANCE

The cycle time τ of an instruction pipeline is the time needed to advance a set of instructions one stage through the pipeline.

$$\text{Cycle time} = \max[\tau_i] + d = \tau_m + d, 1 \leq i \leq K$$

where τ_i = Time delay of the circuitry in the i^{th} stage of the pipeline.

τ_m = maximum stage delay.

K = number of stages in instruction pipeline

d = time delay of a latch, needed to advance signals and data from one stage to the next.

Suppose that n instructions are processed without any branches. Let $T_{k,n}$ be the total time required for a pipeline with K stages to execute n instructions. Then

$$T_{k,n} = [K + (n - 1)]\tau$$

Example 1: Let $n = 7, K = 6, \tau = 1$. Then $T_{k,n} = [6 + (7 - 1)] \times 1 = 12$ cycles.

Now consider a processor with equivalent functions but no pipeline and assume that the instruction cycle time is $k\tau$. The speed up factor for the instruction pipeline compared to execution without the pipeline is defined as

$$S_k = \frac{T_{1,n}}{T_{k,n}} = \frac{nk\tau}{[(k + (n - 1))\tau]} = \frac{nk}{k + (n - 1)}$$

Note: Larger the number of stages, greater the potential for speed up. But practically, the potential gains of additional pipeline stages are countered by increase in cost, delays between stages, and the fact that branches will be encountered requiring the flushing of the pipeline.

Arithmetic pipeline:

- Pipeline arithmetic units are usually found in very high-speed computers.
- These are used to implement floating point operations, multiplication of fixed point numbers and similar computations encountered in scientific problems.

PIPELINE HAZARDS

- Pipeline hazards are situations that prevent the next instruction in the instruction stream from executing during its designated clock cycle. The instruction is said to be stalled. When an instruction is stalled, all instructions later in the pipeline than the stalled instructions are also stalled. Instructions earlier than the stalled one can continue. No new instructions are fetched during the stall.

Note: Keeping a pipeline at its maximal rate is prevented by pipeline hazard.

Different types of Hazards:

1. Structural Hazards
2. Data Hazards
3. Control Hazards

Structural Hazards

Structural hazards occur when a certain resource is requested by more than one instruction at the same time.

Example 2: Instruction MVI *B, X* fetches in the *O* stage operand *X* from memory. The memory does not accept another access during that cycle.

Clock cycle	1	2	3	4	5	6	7	8	9	10	11
MVI B, X	F	D	C	O	E	W					
Instruction $i+1$		F	D	C	O	E	W				
$i+2$			F	D	C	O	E	W			
$i+3$				stall	F	D	C	O	E	W	
$i+4$						F	D	C	O	E	W

Penalty: 1 cycle

Certain resources are duplicated in order to avoid structural hazards (ALU, floating-point unit) can be pipelined themselves in order to support several instructions at a time. A classical way to avoid hazards at memory access is by providing separate data and instruction caches.

Note: Structural hazards are due to resource conflict.

Data Hazards

In a pipeline execution of two instructions I_1 and I_2 a certain stage of the pipeline I_2 needs the result produced by I_1 , but this result has not yet been generated, then we have a data hazard.

Example 3: $I_1: \text{ADD } R_3, R_2 \quad R_3 \leftarrow R_3 + R_2$
 $I_2: \text{MUL } R_1, R_3 \quad R_1 \leftarrow R_1 * R_3$

Clock cycle	1	2	3	4	5	6	7	8	9	10	11	12
ADD R_3, R_2	F	D	C	O	E	W						
MUL R_1, R_3		F	D	C	stall	stall	O	E	W			
Instruction $i+2$			F	D			C	O	E	W		

Penalty: 2 cycles Before executing the O stage (operand fetch stage), the MUL instruction is stalled until the ADD instruction has written the result into R_3 .

Data dependencies

Data dependency exists between two instructions if the data used by an instruction depends on the data created by other instructions.

Two type of dependencies exist between instructions:

1. True data dependency
2. Name dependencies
 - (i) Anti-dependency
 - (ii) Output dependency

True data dependency

- This is also called as Read-After-Write Hazard (RAW).
- This type of dependency occurs when the value produced by an instruction is required by a subsequent instruction.

- This is also known as a flow dependency because dependency is due to flow of data in a program.

Example: $\text{ADD } R_3, R_2, R_1; R_3 \leftarrow R_2 + R_1$
 $\text{SUB } R_4, R_3, 1; R_4 \leftarrow R_3 - 1$

- Here R_3 is read before it is written by 'ADD' instruction.
- In RAW hazard $(i+1)^{\text{st}}$ instruction tries to read a source before it is written by ' i^{th} ' instruction. So $(i+1)^{\text{st}}$ instruction incorrectly gets the old value.
- This kind of hazard can be reduced by using forwarding (or Bypassing).

Name dependencies

1. Anti-dependency:

- This is also called as Write-After-Read hazard
- This kind of dependency occurs when an instruction writes to a location which has been read by a previous instruction.
- Here $(i+1)^{\text{st}}$ instruction tries to write an operand before it is read by i^{th} instruction. So i^{th} instruction incorrectly gets the new value.

Example: $I_1: \text{ADD } R_3, R_2, R_1; R_3 \leftarrow R_2 + R_1$
 $I_2: \text{SUB } R_2, R_5, 1; R_2 \leftarrow R_5 - 1$

I_2 must not produce its result in R_2 before I_1 read R_2 , otherwise I_1 would use the value produced by I_2 rather than the previous value of R_2 .

2. Output dependency:

- This is also called as Write - After - Write (WAW) hazard.
- This dependency occurs when a location is written by two instructions.
- i.e., $(i+1)^{\text{th}}$ instruction tries to write an operand before it is written by i^{th} instruction.

The writes end up being performed in the wrong order.

Example: $I_1: \text{ADD } R_3, R_2, R_1; R_3 \leftarrow R_2 + R_1$
 $I_2: \text{SUB } R_2, R_3, 1; R_2 \leftarrow R_3 - 1$
 $I_3: R_3, R_2, R_5; R_3 \leftarrow R_2 + R_5$

There is a possibility of WAW hazard between I_1 and I_3 .

Handling data dependency There are ways to handle data dependency.

1. Hardware interlocks
2. Operand forwarding
3. Delayed load

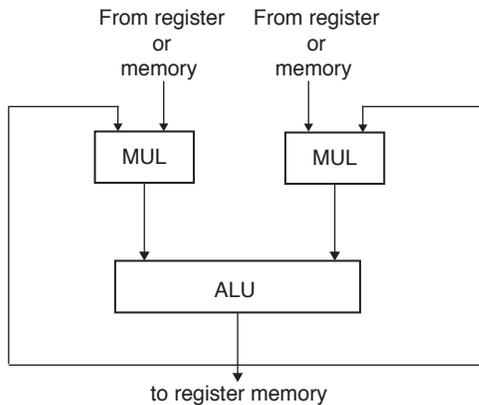
1. Hardware interlocks:

- To avoid data dependency, insert hardware interlock.
- An interlock is a circuit that detects instructions whose source operands are destinations of instructions farther up in the pipeline.

2. Operand forwarding:

- Uses special hardware to detect a conflict and then avoid it by routing the data through special paths between pipeline segments.

- Some of the penalty produced by data hazards can be avoided using a technique called forwarding (Bypassing).
- The ALU result is always fed back to the ALU input. If the hardware detects that the value needed for the current operation is the one produced by the previous operation (but which has not yet been written back). It selects the forwarded result as the ALU input, instead of the value read from register or memory.



Clock cycle	1	2	3	4	5	6	7	8
ADD R ₃ , R ₂	F	D	C	O	E	W		
MUL R ₁ , R ₃		F	D	C	stall	O	E	W

Penalty: 1 cycle After the E stage of the MUL instruction the result is available by forwarding. Therefore the penalty is reduced to one cycle.

Delayed Load: Here the compiler of a computer will detect the data conflicts and reorder the instructions as necessary to delay the loading of the conflicting data by inserting no-operation instruction.

Control Hazards

Control hazards are produced by Branch Instructions.

Unconditional branch

- Jump loop
- Loop

Clock cycle	1	2	3	4	5	6	7	8	9	10	11
Loop	F	D	C	O	E	W					
Loop		stall	stall	stall	F	D	C	O	E	W	
Loop + 1						F	D	C	O	E	W

Penalty: 3 cycles

- The instruction following the branch is fetched before the D stage is finished in 2nd clock. It is not known that a branch is executed. Later the fetched instruction is discarded.
- After the O stage of the branch instruction the address of the target is known and it can be fetched.

Conditional branch

Example: ADD B; A ← A + B
JZ Loop

Loop: If condition satisfies and branch is taken:

Clock cycle	1	2	3	4	5	6	7	8	9	10	11	12
ADD B	F	D	C	O	E	W						
JZ Loop		F	D	C	O	E	W					
Loop			stall	stall	stall	F	D	C	O	E	W	

Penalty: 3 cycles

At this moment both the condition (set by ADD) and the target address are known.

If condition not satisfied and branch not taken:

Clock cycle	1	2	3	4	5	6	7	8	9	10	11	12
ADD B	F	D	C	O	E	W						
JZ Loop		F	D	C	O	E	W					
Instruction i + 1			F	stall	stall	D	C	O	E	W		

Penalty: 2 cycles

At this moment the condition is known and instruction i + 1 can go on.

- With conditional branch, we have a penalty even if the branch has not been taken. This is because we have to wait until the branch condition is available.

Dealing with branches

One of the major problems in designing an instruction pipeline is the occurrence of branch instructions. A variety of approaches have been taken for dealing with branches

1. Multiple streams
2. Prefetch branch target
3. Branch target buffer
4. Loop buffer
5. Branch prediction
6. Delayed branch

Multiple streams A branch instruction may cause to choose one of two instructions to fetch next, then allow the pipeline to fetch both instructions, making use of streams.

There are two problems with this approach:

1. With multiple pipelines there are contention delays for access to the registers and to memory.
2. Additional branch instructions may enter the pipeline before the original branch decision is resolved.

Prefetch branch target When a conditional branch is recognized, the target of the branches is prefetched, in addition to the instruction following the branch. This target is then saved until the branch instruction is executed.

Branch targets buffer (BTB) BTB is an associative memory included in the fetch segment of the pipeline. Each entry in the BTB consists of the address of a previously executed branch instruction and target instructions for that branch.

It also stores the next few instructions after the branch target instruction.

When the pipeline decodes a branch instruction, it searches the associative memory BTB for the address of the instruction. If it is in BTB, the instruction is available directly and prefetch continues from the new path. If the instruction is not in BTB, the pipeline shifts to a new instruction stream and stores the target instruction in the BTB.

Advantage: Branch instructions that occurred previously are readily available in the pipeline without interruption.

Loop Buffer A loop buffer is a small, very high speed memory maintained by the instruction fetch stage of the pipeline and containing the n most recently fetched instructions in sequence. If a branch is to be taken, the hardware first checks whether the branch target is within the buffer. If so, the next instruction is fetched from the buffer. The Advantages of loop buffer are

1. Loop buffer will contain some instructions sequentially ahead of the current instruction fetch address. Thus instructions fetched in sequence will be available without the usual memory access time.
2. If the branch occurs to a target just a few locations ahead of the address of the branch instruction, the target will already be in the buffer.
3. This strategy is well suited in dealing with loops.

Branch prediction Various techniques can be used to predict whether a branch will be taken or not. The common techniques are

Static { 1. Predict never taken Static
2. Predict always taken
3. Predict by opcode

Dynamic { 4. Taken/not taken switch
5. Branch history table

- The first two approaches are static, i.e., no dependency on execution history. Here always assume that the branch will not be taken and continue to fetch instructions in sequence, or always assume that the branch will be taken and always fetch from the branch target.
- The third approach is also static. Takes the decision based on the opcode of the branch instruction in a program.
- Dynamic branch strategies attempt to improve the accuracy of prediction by recording the history of conditional branch instructions in a program.

(a) Taken/not taken switch:

- Use two bits to record the result of the last two instances of the execution of the associated instruction or record a state in some other fashion.

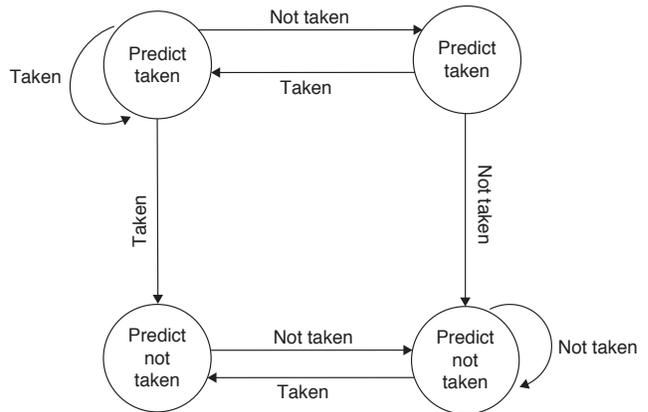


Figure 1 Branch prediction state diagram

- As long as each succeeding conditional branch instruction that is encountered is taken, the decision process predicts that the next branch will be taken.
- If a single prediction is wrong, the algorithm continues to predict that the next branch is taken.
- Only if two successive branches are not taken does the algorithm shift to not taken branch.

Drawback: If the decision is made to take the branch, the target instruction cannot be fetched until the target address, which is an operand in the conditional branch instruction is decoded.

(b) Branch history table: It is a small cache memory associated with the instruction fetch stage of the pipeline.

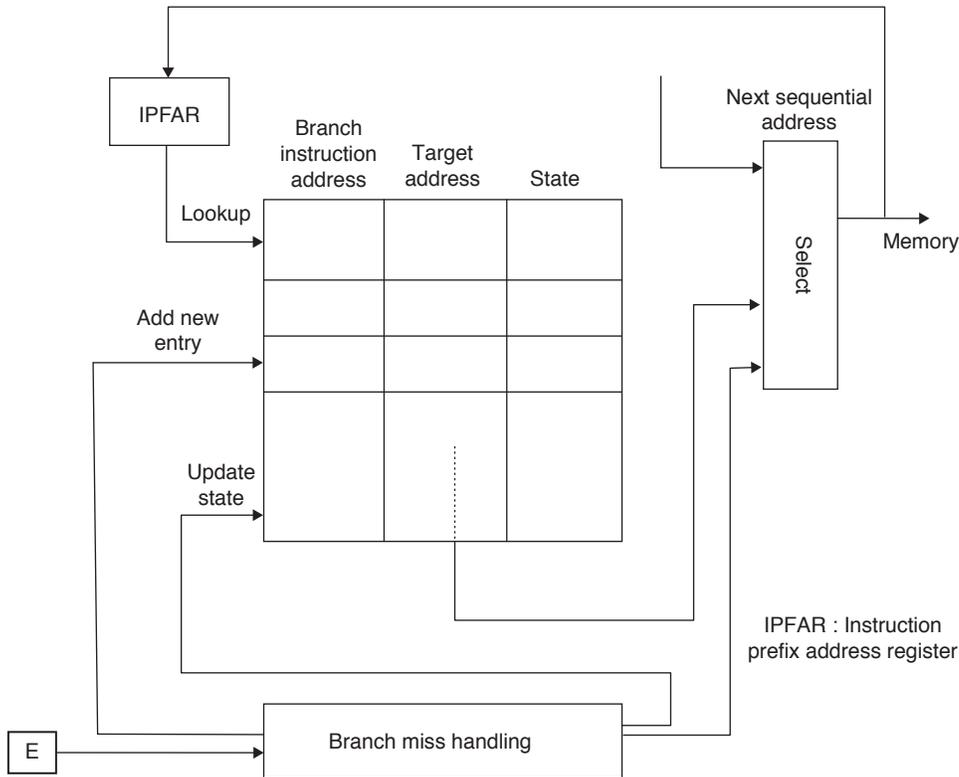


Figure 2 Branch history table

- Each entry in the table consist of three elements:
 1. The address of branch instruction.
 2. Some number of history bits that record the state of use of that instruction.
 3. Information about target instruction.

Delayed branch A compiler detects the branch instructions and rearranges the machine language code sequence by inserting useful instructions that keep a pipeline operating without interruptions.

Exercises

Practice Problems I

Directions for questions 1 to 21: Select the correct alternative from the given choices.

Common data for questions 1 and 2: An unpipelined processor with eight number cycle time and pipeline batches with 1 ns latency is given.

1. Find the cycle times of pipelined versions of the processor with 2, 4, 8 and 16 stages if the Data path logic is evenly divided among the pipeline stages.

(A) 5, 3, 2, 1.5	(B) 4, 2, 1, 0.5
(C) 8, 4, 2, 1	(D) 10, 6, 4, 3
2. What is the latency of each of the pipelined versions of the processor?

(A) 4, 2, 1, 0.5	(B) 10, 6, 4, 3
(C) 5, 3, 2, 1.5	(D) 10, 12, 16, 24
3. A 4-stage pipeline has the stage delays as 110, 120, 130, and 140 nanoseconds respectively. Registers that are used between the stages have a delay of 2 nanoseconds each. Assuming constant clocking rate. Find

the total time taken to process 1000 instructions on this pipeline.

- | | |
|------------|--------------|
| (A) 7.1 ms | (B) 14.24 ms |
| (C) 28 ms | (D) 2000 ms |

4. Consider a pipelined processor with the following four stages:

IF: instruction fetch
 ID: Instruction decode
 EX: Execute
 WB: Write back

The IF, ID and WB stages takes one clock cycle each to complete the operation. The number of clock cycles for EX stage depends on the instruction; for I_1 and I_3 one clock cycle is needed and for I_2 three clock cycles are needed. Find the number of clock cycles taken to complete the following sequence of instructions?

I_1 :	ADD R_0, R_1, R_2 $R_0 \leftarrow R_1 + R_2$
I_2 :	MUL R_2, R_3, R_4 $R_2 \leftarrow R_3 \times R_4$
I_3 :	SUB R_4, R_5, R_6 $R_4 \leftarrow R_5 - R_6$

- (A) 7 (B) 8
(C) 6 (D) 9
5. A CPU has five stage pipelines and runs at 1 GHz frequency. Instruction fetch happens in the first stage of the pipeline. A conditional branch instruction computes the target address and evaluates the condition in the third stage of the pipeline. The processor stops fetching new instructions following a conditional branch until the branch outcome is known. A program executes 10^9 instructions. Out of which 10% are conditional branches. If each instruction takes one cycle to complete on average then find the total execution time of the program?
- (A) 1 sec (B) 1.2 sec
(C) 1.4 sec (D) 1.8 sec
6. Consider a four stage pipeline processor, number of cycles needed by the four instructions I_1, I_2, I_3 and I_4 in stages S_1, S_2, S_3 and S_4 are shown below:

	S_1	S_2	S_3	S_4
I_1	2	1	1	1
I_2	1	3	2	2
I_3	2	1	1	3
I_4	1	2	2	2

What is the number of cycles needed to execute the instructions in the order:

$I_1 : I_2 : I_3 : I_4$

- (A) 8 (B) 12
(C) 14 (D) 15
7. A non-pipelined system takes 50 ns to process a task; the same task can be processed in a six-segment pipeline with a clock cycle of 10 ns. Speedup ratio of 100 tasks for pipeline is
- (A) 1.62 (B) 3.21
(C) 4.76 (D) 8.21
8. Consider a pipelined processor with the following four stages:
IF: Instruction fetch
ID: Instruction decode
EX: Execute
WB: Write back
- The IF, ID and WB stages takes one clock cycle each to complete the operation. The number of clock cycles for EX stage depends on the instruction; for I_1 and I_3 one clock cycle is needed and for I_2 three clock cycles are needed. The number of clock cycles taken to complete the following sequence of instructions is

I_1 :	ADD $R_0, R_1, R_2,$	$R_0 \leftarrow R_1 + R_2$
I_2 :	MUL $R_2, R_0, R_4,$	$R_2 \leftarrow R_0 \times R_4$
I_3 :	SUB $R_4, R_5, R_2,$	$R_4 \leftarrow R_5 - R_2$

- (A) 7 (B) 8
(C) 9 (D) 10
9. Following are the sequence of stages in a pipeline CPU:
(1) IF: Instruction fetch from instruction memory
(2) RD: Instruction decode and register read
(3) EX: Execute ALU operation for data and address computation
(4) MA: Data memory access, for write access, the register read at RD stage is used.
(5) WB: Register write back
- Consider the following sequence of instructions:
LOAD $R_1, M[\text{loc}]$
ADD R_1, R_1, R_1
ADD R_2, R_1, R_2
- Let each stage take one clock cycle.
- What is the number of clock cycles taken to complete the above sequence of instructions starting from the fetch of first instruction?
- (A) 18 (B) 15
(C) 13 (D) 10
10. Which of the following can cause a hazard for a pipelined CPU with a single ALU?
- (i) The $(j + 1)^{\text{st}}$ instruction uses the result of the j^{th} instruction as an operand.
(ii) The j^{th} and $(j + 1)^{\text{st}}$ instructions require the ALU at the same time.
(iii) The execution of a conditional jump instruction.
(iv) The execution of non-conditional jump instruction.
- (A) (i) and (ii) (B) (ii) and (iii)
(C) (i), (ii) and (iii) (D) (i), (ii), (iii) and (iv)
11. Given an unpipelined processor with a 10 ns cycle time and pipeline latches with 0.5 ns latency, how many stages of pipelining are required to achieve cycle time of 2 ns?
- (A) 5.5 (B) 6.67
(C) 7 (D) 6
12. In a 4-stage pipeline,
IF – instruction fetches
ID – instruction decode and fetch operands
EX – Execute
WB – write back
ADD, SUB take one clock cycle, MUL take three clock cycles. Then for
ADD $R_2, R_1, R_0, R_2 \leftarrow R_1 + R_0$
MUL $R_4, R_3, R_2, R_4 \leftarrow R_3 * R_2$
SUB $R_6, R_5, R_4, R_6 \leftarrow R_5 - R_4$
Number of clock cycles required using operand forwarding technique are
- (A) 8 (B) 12
(C) 10 (D) 14
13. Consider an instruction sequence of length 'n' that is streaming through a K-stage instructions pipeline. Let P be the probability of encountering a conditional or

unconditional branch instruction and let q be the probability that execution of a branch instruction I_B causes a jump to a non-consecutive address. Assume that each such jump requires the pipeline to be cleared, destroying all ongoing instruction processing, when I_B emerges from the last stage. Also assume that T is the cycle time. Then which of the following expression correctly specifies the time required for this pipeline?

- (A) $pqnk\tau + (1 - pq)[K + (n - 1)]\tau$
 (B) $(1 - pq)[k + (n - 1)]\tau + pqn\tau$
 (C) $pqnk\tau + (1 - pq)n[k + (n - 1)]\tau$
 (D) $pqn + (1 - pq)n[k + (n - 1)]\tau$
14. If T_m is maximum stage delay of an m -stage pipeline with time delay of the latch is d then cycle time is
 (A) T_m/d (B) $T_m + d$
 (C) $2T_m + d$ (D) $T_m \times d$
15. Pipelining is a general technique for increasing processor _____ without requiring large amounts of extra hardware.
 (A) turnaround time (B) waiting time
 (C) latency (D) throughput
16. A 4-stage instruction pipeline executes a 100 instruction program. The probability of occurrence of a conditional or unconditional branch is 0.4 and the probability of execution of a branch instruction I_B causing a jump to a non-consecutive address is 0.1. Then the speed up factor for the instruction pipeline compared to execution without pipeline is
 (A) 2.14 (B) 6.23
 (C) 3.21 (D) 3.48
17. A non-pipelined processor has a clock rate of 2.5GHz and an average cycles per instruction of 4. An upgrade to the processor introduces a five stage pipeline. However, due to internal pipeline delays, such as latch delay, the clock rate of the new processor has to be reduced to 2GHz. What is the MIPS rate for each of these processors respectively.
 (A) 625, 400 MIPS (B) 625, 2000 MIPS
 (C) 3125, 2000 MIPS (D) 3125, 400 MIPS
18. Consider the following sequence of instructions:
 I_1 : MUL R_1, R_2 $R_1 \leftarrow R_1 * R_2$
 I_2 : SUB $R_3, 1$ $R_3 \leftarrow R_3 - 1$

I_3 : ADD R_3, R_4 $R_3 \leftarrow R_3 + R_4$
 I_4 : BEZ Target Branch if zero
 I_5 : MOVE $R_3, 10$ $R_3 \leftarrow 10$
 :

Target:

Which of the following instruction will be placed in delayed slot to reduce penalty in a 6-stage pipeline? (Assume that the branch outcome will be known during 5th stage)

- (A) I_1 (B) I_2
 (C) I_3 (D) I_5

19. Consider the following sequence of instructions:

ADD R_1, R_2 $R_1 \leftarrow R_1 + R_2$
 BEZ Target Branch if Zero
 MUL R_3, R_4 $R_3 \leftarrow R_3 * R_4$
 MOVE $R_1, 10$ $R_1 \leftarrow 10$
 :

Target:

Assume that this program executed on a 6-stage pipelined processor and each stage required 1 clock cycle.

Let us suppose that “branch not taken” Prediction is used but the prediction is not fulfilled, then the penalty will be (branch outcome is known at 5th stage)

- (A) 1 clock cycle (B) 2 clock cycles
 (C) 3 clock cycles (D) 4 clock cycles
20. Suppose 40% of the instructions are loads and half the time they are followed by instruction that depends on value loaded. If this hazard causes single cycle delay, how must faster is ideal pipelined machine (CPI = 1) than real one? (Ignore other stalls)
 (A) 1 time (B) 1.2 times
 (C) 1.5 times (D) 1.15 times
21. Assume that a pipelined processor has three categories of instructions: Branch, load/store, other. If it is a branch instruction it will take 3 clock cycles, if it is a load/store instruction it will take 4 clock cycles and all other instructions require 6 clock cycles. A program consisting of 10% Branch instructions, 10% of load/store instructions is executed on this processor. Then the number of clock cycles required for the execution of the program is
 (A) 2.45 (B) 3.61
 (C) 4.66 (D) 5.5

Practice Problems 2

Directions for questions 1 to 20: Select the correct alternative from the given choices.

1. The time required for the five functional units, which operated in each of the five cycles are 10 ns, 7 ns, 10 ns, 10 ns and 8 ns. Assume that pipelining add 1 ns of overhead. The speed up of pipeline compared to unpipeline is
 (A) 4.5 times (B) 1.1 times
 (C) 4.1 times (D) 2.4 times

2. Which of the following is a technique of decomposing a sequential process into sub operations with each sub-process being executed in a special dedicated segment that operates concurrently with each other?

- (A) Straight line sequencing
 (B) Random sequencing
 (C) Pipelining
 (D) Serial execution

3. Which of the following statements is incorrect?
- Latency is the number of time units between two initiations in a pipelined architecture.
 - If initiations are of different but fixed reservation tables, the architecture is known as static pipelined configuration.
 - A collision in a pipelined architecture is an attempt by two different initiations to use the same stage at the same time.
 - None of the above
4. Which of the following technique is used in a pipelined processor, when there is a conditional branch?
- Loop butter
 - Branch prediction
 - Delayed Branch
 - All of the above
5. Which of the following cases, leads to a pipelined computer architecture?
- The evaluation of each basic function is relatively independent of the previous one.
 - The sub-functions are closely related to each other.
 - The evaluation of each sub function requires approximately the same sequence.
 - All of the above
6. The performance of a pipelined processors is degraded if
- the pipeline stages have different delays.
 - consecutive instructions are dependent on each other.
 - the pipeline stages share hardware resources.
 - All of the above
7. The following is a limit on how much the performance of a processor can be improved using pipelining:
- the number of pipeline stages
 - data dependencies
 - branch delays
 - All of the above
8. A pipeline processor consists of a sequence of 'm' data processing circuits called ____, which collectively perform a single operation on a stream of data operands passing through them.
- stages
 - pipelines
 - latches
 - None of the above
9. A five-stage pipelined CPU has the following sequence of stages:
- IF – Instruction fetch from memory
 RD – Decode instruction
 EX – Execute
 MA – Data memory access
 WB – Register write back
- Consider the following instruction sequence:
- I_1 : Load R_0 $R_0 \leftarrow M$
 I_2 : ADD R_1 , $R_1 R_1 \leftarrow R_1 + R_1$
 I_3 : SUB R_2 , $R_3 R_2 \leftarrow R_2 - R_3$

Each stage takes one clock cycle.

Number of clock cycles to execute the program is

- 8
- 10
- 7
- 15

Common data for questions 10 and 11: Given an unpipelined processor with a 10 number cycle time and pipeline latches with 0.5 ns latency.

10. Which are the cycle times of pipelined versions of the processors with 2, 4, 8 and 16 stages if the data path logic is evenly divided among the pipeline stages?
- 5.0, 3.0, 1.5, 1.0
 - 5.5, 3.0, 1.75, 1.125
 - 4.0, 5.0, 6.0, 7.0
 - None of the above
11. What is the latency of each of the pipelined versions of the processor with 2, 4, 8 and 16 stages?
- 10, 11, 12, 14 ns
 - 10, 10, 11, 11 ns
 - 11, 12, 14, 18 ns
 - None of the above
12. Assume an unpipelined processor has a 1 ns clock cycle and it uses 5 cycles for ALU operations and branches. And 6 clock cycles for memory operations. A program has 40%, 30%, and 20% of ALU operations, branch instructions and memory operations respectively. If we are using pipelining it adds 0.2 ns overhead. Then what is the speedup of pipelining compared to unpipelined processor?
- 1.2
 - 3.91
 - 4.7
 - 2.5
13. Consider a five-stage pipeline processor in which each instructions on an average has 2 clock cycle stalls. Then the speed up of this pipelined processor compared to an unpipelined processor is
- 2.5
 - 1.67
 - 0.4
 - 5
14. Pipelining strategy is called to implement
- instruction execution
 - instruction prefetch
 - instruction decoding
 - instruction manipulation
15. If an Instruction 'j' tries to read a source operand before instruction 'i' writes it. Then it is a ____ type of hazard.
- WAR
 - RAW
 - WAW
 - None of these
16. What is the average instruction processing time of a five-stage instruction pipeline for 32 instructions if conditional branch instructions occur as follows: $I_2, I_5, I_7, I_{25}, I_{27}$.
- 1.97
 - 1.67
 - 1.75
 - 1.25

17. Consider the execution of 1000 instructions on a five-stage pipeline machine. Then the speed-up due to the use of pipelining given that the probability of an instruction being a branch is 0.2.
- (A) 1.77 (B) 2.6
(C) 2.77 (D) 3.2
18. If an instructions following a branch (taken or not taken) have a dependency on the branch and cannot be executed until the branch is executed, then the dependency is
- (A) True data dependency
(B) Procedural dependency
(C) Resource conflict
(D) Output dependency
19. 'A two-stage instruction pipeline unlikely to cut the instruction cycle time in half, compared with the use of no pipeline.' The statement is
- (A) Always true (B) Always False
(C) Can't predict (D) Some times true
20. Write after read dependency is also known as
- (A) True dependency (B) Anti-dependency
(C) Output dependency (D) Inverse dependency

PREVIOUS YEARS' QUESTIONS

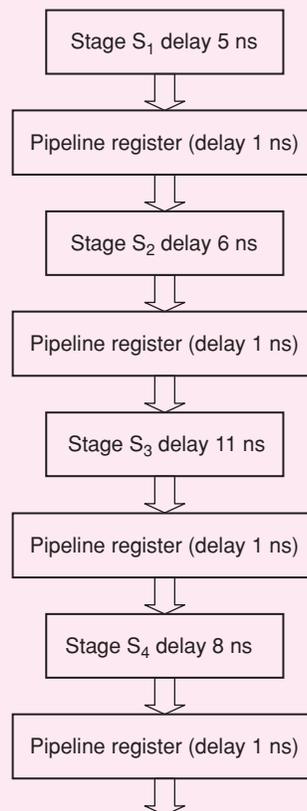
1. Consider a 6-stage instruction pipeline, where all stages are perfectly balanced. Assume that there is no cycle time overhead of pipelining. When an application is executing on this 6-stage pipeline, the speedup achieved with respect to non-pipelined execution if 25% of the instructions incur 2 pipeline stall cycles is _____. [2014]
2. Consider the following processors (ns stands for nano-seconds). Assume that the pipeline registers have zero latency.
- P1: Four-stage pipeline with stage latencies 1 ns, 2 ns, 2 ns, 1 ns.
 P2: Four-stage pipeline with stage latencies 1 ns, 1.5 ns, 1.5 ns, 1.5 ns.
 P3: Five-stage pipeline with stage latencies 0.5 ns, 1 ns, 1 ns, 0.6 ns, 1 ns.
 P4: Five-stage pipeline with stage latencies 0.5 ns, 0.5 ns, 1 ns, 1 ns, 1.1 ns.

Which processor has the highest peak clock frequency? [2014]

- (A) P1 (B) P2
(C) P3 (D) P4
3. An instruction pipeline has five stages, namely, instruction fetch (IF), instruction decode and register fetch (ID/RF), instruction execution (EX), memory access (MEM), and register write back (WB) with stage latencies 1 ns, 2.2 ns, 2 ns, 1 ns, and 0.75 ns, respectively (ns stands for nano seconds). To gain in terms of frequency, the designers have decided to split the ID/RF stage into three stages (ID, RF1, RF2) each of latency 2.2/3 ns. Also, the EX stage is split into two stages (EX1, EX2) each of latency 1 ns. The new design has a total of eight pipeline stages. A program has 20% branch instructions which execute in the EX stage and produce the next instruction pointer at the end of the EX stage in the old design and at the end of the EX2 stage in the new design. The IF stage stalls after fetching a branch

instruction until the next instruction pointer is computed . All instructions other than the branch instruction have an average CPI of one in both the designs. The execution times of this program on the old and the new design are P and Q nanoseconds, respectively. The value of P/Q is _____. [2014]

4. Consider an instruction pipeline with four stages (S_1, S_2, S_3 and S_4) each with combinational circuit only. The pipeline registers are required between each stage and at the end of the last stage. Delays for the stages and for the pipeline registers are as given in the figure.



What is the approximate speed up of the pipeline in steady state under ideal conditions when compared to the corresponding non-pipeline implementation?

[2011]

- (A) 4.0 (B) 2.5
(C) 1.1 (D) 3.0

5. Register renaming is done in pipelined processors [2012]

- (A) as an alternative to register allocation at compile time
(B) for efficient access to function parameters and local variables
(C) to handle certain kinds of hazards
(D) as part of address translation

6. Consider an instruction pipeline with five stages without any branch prediction: Fetch Instruction (FI), Decode Instruction (DI), Fetch Operand (FO), Execute Instruction (EI) and Write Operand (WO). The stage delays for FI, DI, FO, EI and WO are 5 ns, 7 ns, 10 ns, 8 ns and 6 ns, respectively. There are intermediate storage buffers after each stage and the delay of each buffer is 1 ns. A program consisting of 12 instructions $I_1, I_2, I_3, \dots, I_{12}$ is executed in this pipelined processor. Instruction I_4 is the only branch instruction and its branch target is I_9 . If the branch is taken during the execution of this program, the time (in ns) needed to complete the program is [2013]

- (A) 132 (B) 165
(C) 176 (D) 328

Common data for Questions 7 and 8: Delayed branching can help in the handling of control hazards

7. For all delayed conditional branch instructions, irrespective of whether the condition evaluates to true or false [2008]
- (A) The instruction following the conditional branch instruction in memory is executed
(B) The first instruction in the fall through path is executed
(C) The first instruction in the taken path is executed
(D) The branch takes longer to execute than any other instruction

8. The following code is to run on a pipelined processor with one branch delay slot:

I_1 : ADD $R_2 \leftarrow R_7 + R_8$

I_2 : SUB $R_4 \leftarrow R_5 - R_6$

I_3 : ADD $R_1 \leftarrow R_2 + R_3$

I_4 : STORE Memory [R_4] $\leftarrow R_1$

BRANCH to Label if $R_1 = 0$

Which of the instructions I_1, I_2, I_3 or I_4 can legitimately occupy the delay slot without any other program modification? [2008]

- (A) I_1 (B) I_2
(C) I_3 (D) I_4

9. A 5 stage pipelined CPU has the following sequence of stages:

IF - Instruction fetch from instruction memory

RD - Instruction decode and register read

EX - Execute: ALU operation for data and address computation

MA - Data memory access - for write access, the register read at RD state is used

WB - Register write back

Consider the following sequence of instructions.

I_1 : L $R_0, \text{loc}_1; R_0 \leq M[\text{loc}_1]$

I_2 : A $R_0, R_0; R_0 \leq R_0 + R_0$

I_3 : S $R_2, R_0; R_2 \leq R_2 - R_0$

Let each stage take one clock cycle.

What is the number of clock cycles taken to complete the above sequence of instructions from the fetch of I_1 ?

- (A) 8 (B) 10
(C) 12 (D) 15

10. Consider a non-pipelined processor with a clock rate of 2.5 gigahertz and average cycles per instruction of four. The same processor is upgraded to a pipelined processor with five stages; but due to the internal pipeline delay, the clock speed is reduced to 2 gigahertz. Assume that there are no stalls in the pipeline. The speed up achieved in this pipelined processor is _____. [2015]

11. Consider the sequence of machine instructions given below:

MUL R_5, R_0, R_1

DIV R_6, R_2, R_3

ADD R_7, R_5, R_6

SUB R_8, R_7, R_4

In the above sequence, R_0 to R_8 are general purpose registers. In the instructions shown, the first register stores the result of the operation performed on the second and the third registers. This sequence of instructions is to be executed in a pipelined instruction processor with the following 4 stages: (1) Instruction Fetch and Decode (IF), (2) Operand Fetch (OF), (3) Perform Operation (PO) and (4) Write back the result (WB). The IF, OF and WB stages take 1 clock cycle each for any instruction. The PO stage takes 1 clock cycle for ADD or SUB instruction, 3 clock cycles for MUL instruction and 5 clock cycles for DIV instruction. The pipelined processor uses operand forwarding from the PO stage to the OF stage. The number of clock cycles taken for the execution of the above sequence of instructions is _____. [2015]

12. Consider the following reservation table for a pipeline having the stages S_1 , S_2 and S_3 .

	Time →				
	1	2	3	4	5
S_1	X				X
S_2		X		X	
S_3			X		

The minimum average latency (MAL) is _____ [2015]

13. Consider the following code sequence having five instructions I_1 to I_5 . Each of these instructions has the following format. [2015]

$OP\ Ri,\ Rj,\ Rk$

Where operation OP is performed on contents of registers Rj and Rk and the result is stored in register Ri .

I_1 : ADD R_1, R_2, R_3

I_2 : MUL R_7, R_1, R_3

I_3 : SUB R_4, R_1, R_5

I_4 : ADD R_3, R_2, R_4

I_5 : MUL R_7, R_8, R_9

Consider the following three statements.

S_1 : There is an anti-dependence instructions between instructions I_2 and I_5

S_2 : There is an anti-dependence between Instructions I_2 and I_4

S_3 : |Within an instruction pipeline an anti-dependence always creates one or more stalls

Which one of the above statements is/are correct?

- (A) Only S_1 is true
- (B) Only S_2 is true
- (C) Only S_1 and S_3 are true
- (D) Only S_2 and S_3 are true

14. The stage delays in a 4 - stage pipeline are 800, 500, 400 and 300 picoseconds. The first stage (with delay 800 picoseconds) is replaced with a functionally equivalent design involving two stages with respective

delays 600 and 350 picoseconds. The throughput increase of the pipeline is _____ percent. [2016]

15. Consider a 3 GHz (gigahertz) processor with a three - stage pipeline and stage latencies τ_1, τ_2 and τ_3 such that $\tau_1 = 3\tau_2 / 4 = 2\tau_3$. If the longest pipelines stage is split into two pipeline stages of equal latency, the new frequency is _____ GHz, ignoring delays in the pipeline registers. [2016]

16. Instruction execution in a processor is divided into 5 stages, *Instruction Fetch* (IF), *Instruction Decode* (ID), *Operand Fetch* (OF), *Execute* (EX), and *Write Back* (WB). These stages take **5, 4, 20, 10, and 3 nanoseconds (ns)** respective. A pipelined implement action of the processor requires buffering between each pair of consecutive stages with a delay of **2 ns**. Two pipelined implementations of the processor are contemplated:

- (i) A navie pipeline implementation (NP) with 5 stages and
- (ii) An efficient pipeline (EP) where the OF stage is divided into stages OF1 and OF2 with execution times of **12 ns** and **8 ns** respectively.

The speedup (correct to two decimal places) achieved by EP over NP in executing 20 independent instructions with no hazards is _____. [2017]

17. The instruction pipeline of a RISC processor has the following stages: *Instruction Fetch* (IF), *Instruction Decode* (ID), *Operand Fetch* (OF), *Perform Operation* (PO) and *Write back* (WB). The IF, ID, OF and WB stages take 1 clock cycle each for every instruction. Consider a sequence of 100 instructions. In the PO stage, 40 instructions take 3 clock cycles each, 35 instructions take 2 clock cycles each, and the remaining 25 instructions take 1 clock cycle each. Assume that there are no data hazards and no control hazards.

The number of clock cycles required for completion of execution of the sequence of instructions is _____. [2018]

ANSWER KEYS**EXERCISES****Practice Problems 1**

- | | | | | | | | | | |
|-------|-------|-------|-------|-------|-------|-------|-------|-------|-------|
| 1. A | 2. D | 3. B | 4. B | 5. B | 6. D | 7. C | 8. D | 9. C | 10. D |
| 11. A | 12. A | 13. A | 14. B | 15. D | 16. D | 17. B | 18. A | 19. B | 20. B |
| 21. D | | | | | | | | | |

Practice Problems 2

- | | | | | | | | | | |
|-------|-------|-------|-------|-------|-------|-------|-------|-------|-------|
| 1. C | 2. C | 3. B | 4. D | 5. D | 6. D | 7. D | 8. A | 9. C | 10. B |
| 11. C | 12. B | 13. B | 14. B | 15. B | 16. C | 17. C | 18. B | 19. A | 20. B |

Previous Years' Questions

- | | | | | | | | | | |
|--------|-------|---------|-----------------|-------|----------|---------|------|------|---------|
| 1. 4 | 2. C | 3. 1.54 | 4. B | 5. C | 6. B | 7. A | 8. D | 9. A | 10. 3.2 |
| 11. 13 | 12. 3 | 13. B | 14. 33.0 : 34.0 | 15. 4 | 16. 1.51 | 17. 219 | | | |

Chapter 5

Cache and Main Memory, Secondary Storage

LEARNING OBJECTIVES

- ☞ Characteristics of memory system
- ☞ Memory hierarchy
- ☞ Locality of reference
- ☞ Cache memory
- ☞ Basic operation of cache
- ☞ Elements of cache design
- ☞ Replacement algorithm
- ☞ Secondary storage
- ☞ Disk
- ☞ Diskette
- ☞ Magnetic tape
- ☞ Optimal memory

CHARACTERISTICS OF MEMORY SYSTEM

1. **Location:** The term refers to whether memory is internal or external to the computer. The location of memory may be
 - Processor
 - Internal (main)
 - External (secondary)
2. **Capacity:** The capacity of internal memory is expressed in terms of bytes. The capacity specified using
 - Word size
 - Number of words
3. **Unit of transfer**
 - For internal memory, the unit of transfer is equal to the number of data lines into and out of the memory module. The unit of transfer need not equal a word or an addressable unit.
 - For external memory, data are often transferred in much larger units than a word, and these are referred to as blocks.
4. **Access method:** The various methods of accessing units of data are
 - (i) **Sequential access:** Memory is organized into units of data, called records.
Example: Magnetic tapes

- (ii) **Direct access:** Individual blocks or records have a unique address based on physical location.

Example: Magnetic disks

- (iii) **Random access:** Each addressable location in memory has a unique, physically wired-in addressing mechanism. The time to access a given location is independent of the sequence of prior accesses and is constant.

Example: Main memory

- (iv) **Associative:** This is a random access type of memory that enables one to make a comparison of desired bit locations within a word for a specified match.

5. **Performance:** Three performance parameters are:

- (i) **Access time (latency):**

- For random access memory, this is the time it takes to perform a read or write operation.
- For non-random-access memory, access time is the time it takes to position the read-write mechanism at the desired location.

- (ii) **Memory cycle time:** For a random access memory it consists of the access time plus any additional time required before a second access can commence.

- (iii) **Transfer rate:** This is rate at which data can be transferred into or out of memory unit.

For Random access memory,

$$\text{Transfer rate} = \frac{1}{\text{Cycle Time}}$$

For non-random access memory, $T_N = T_A + \frac{N}{R}$

Where, T_N = Average time to read or write N -bits.

T_A = Average access time

N = Number of bits

R = Transfer rate in bits per second

6. **Physical type:** The physical type of a memory will be
 - i. Semiconductor
 - ii. Magnetic
 - iii. Optical
 - iv. Magneto-optimal
7. **Physical characteristics:** The memory may be
 - Volatile/non-volatile
 - Erasable/non-erasable
8. **Organization:** There is a trade-off among the three key characteristic of memory.
 - i. Cost
 - ii. Capacity
 - iii. Access time

(i) **Registers:** The fastest, smallest and most expensive type of memory consists of the registers internal to the processor.

(ii) **Main memory:** The principal internal memory system of the computer is main memory. Each location in main memory has a unique address.

(iii) **Cache:** Main memory is usually extended with a higher speed, smaller cache. The cache is not visible to the programmer or, indeed, to the processor. It is a device for staging the movement of data between main memory and processor registers to improve performance.

These three forms of memory are volatile and employ semi conductor technology.

(iv) **Magnetic tapes and disks:** Data are stored more permanently on external mass storage devices, of which the most common are hard disk and removable media.

- External, non-volatile memory is also referred to as secondary or auxiliary memory.

- Used to store program and data files, which are visible to the programmer in the form of files and records.

MEMORY HIERARCHY

Consider the following memory hierarchy, which shows the various memory components.

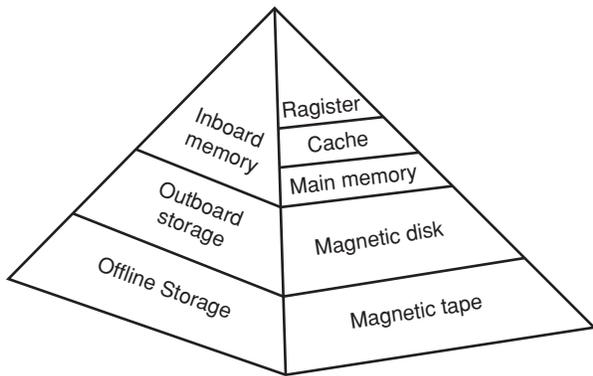


Figure 1 Memory Hierarchy

As one goes down the hierarchy, the following occur:

1. Decreasing cost per bit
2. Increasing capacity
3. Increasing access time
4. Decreasing frequency of access of the memory by the processor.

Locality of Reference

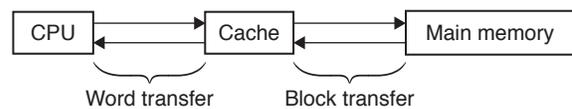
During the course of execution of a program, memory references by the processor, for both instructions and data, tend to cluster. This is referred as principal of locality.

CACHE MEMORY

The locality of reference property states that over a short interval of time, the address generated by a typical program refer to a few localized areas of memory repeatedly, while the remainder of memory is accessed relatively infrequently (Because of frequent loops and subroutine calls).

If the active portions of the program and data are placed in a fast small memory, the average memory access time can be reduced, thus reducing the total execution time of the program. Such a fast small memory is referred to as a cache memory.

Cache memory is intended to give memory speed approaching that of the fastest memories available and at the same time provide a large memory size at the price of less expensive types of semiconductor memories. The following figure shows the structure of cache/main memory system.



The fundamental idea of cache organization is that by keeping the most frequently accessed instructions and data in the fast memory, the average memory access time will approach the access time of cache.

Basic Operation of Cache

- When the CPU need to access memory, the cache is examined. If the word is found in cache, it is read otherwise main memory is accessed to read the word.

- The performance of cache memory is measured in terms of hit ratio.
- When the CPU refers to memory and find the word in cache, it is called hit.
- If the word is not found in cache and is in Main Memory, it is called miss.

$$\text{Hit ratio} = \frac{\text{hits}}{\text{hits} + \text{misses}}$$

$$\text{Average access time} = hc + (1 - h)(c + m)m$$

Where, $c \rightarrow$ Cache access time

$m \rightarrow$ Main memory access time

$h \rightarrow$ hit ratio

- Let main memory consists of up to 2^n addressable words, with each word having a unique n -bit address.
- For mapping purposes, this memory is considered to consist of a number of fixed length blocks of K words each.

$$\therefore \text{Number of blocks } (M) = \frac{2^n}{K}$$

- The cache consists of C lines.
- Each line contains K words, plus a tag of a few bits.
- The number of words in a line is referred to as the line size.
- The number of lines is considerably less than the number of main memory blocks i.e., $C \ll M$.
- Each line includes a tag that identifies which particular block is currently being stored.

The tag is usually a portion of the main memory address.

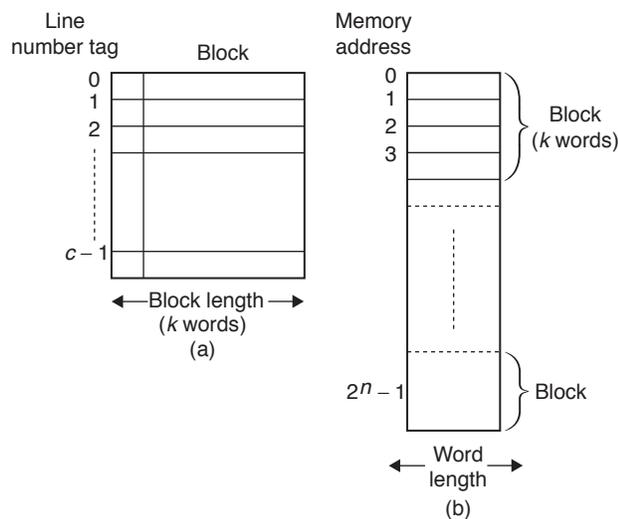


Figure 2 (a) Cache, (b) Main memory

Elements of Cache Design

1. Cache size
2. Mapping function
 - Direct

- Associative
 - Set-associative
3. Replacement algorithm
 4. Write policy
 - Write through
 - Write back
 - Write once
 5. Line size
 6. Number of caches
 - Single or two level
 - Unified or split

Cache size

The size of the cache to be small enough so that the overall average cost per bit is close to that of main memory alone and large enough so that the overall average access time is close to that of the cache alone.

Mapping function

Because there are fewer cache lines than main memory blocks, an algorithm is needed for mapping main memory blocks into cache lines. Three techniques can be used for mapping.

- (i) Direct
- (ii) Associative
- (iii) Set-associative

Direct mapping Maps each block of main memory into only one possible cache line. Figure 2 illustrates the general mechanism. The mapping is expressed as

$$i = j \text{ modulo } m, \text{ where}$$

- $i =$ cache line number
- $j =$ main memory block number
- $m =$ number of lines in the cache

For purpose of cache access, each main memory address can be viewed as consisting of three fields.

- The least significant w bits identify a unique word or byte within a block of main memory.
- The remaining s -bits specify one of the 2^s blocks of main memory. The cache logic interpret these s -bits as a tag of $s-r$ bits. (most significant portion)
- A line field of r -bits, to identify one of 2^r lines.

To summarize,

$$\text{Address length} = (s + w) \text{ bits}$$

$$\text{Number of Addressable units} = 2^{s+w} \text{ words or bytes}$$

$$\text{Block size} = \text{line size} = 2^w \text{ words or bytes}$$

$$\text{Number of blocks in main memory} = \frac{2^{s+w}}{2^w} = 2^s \text{ Number}$$

$$\text{of lines in cache} = M = 2^r.$$

$$\text{Size of Tag} = (s - r) \text{ bits}$$

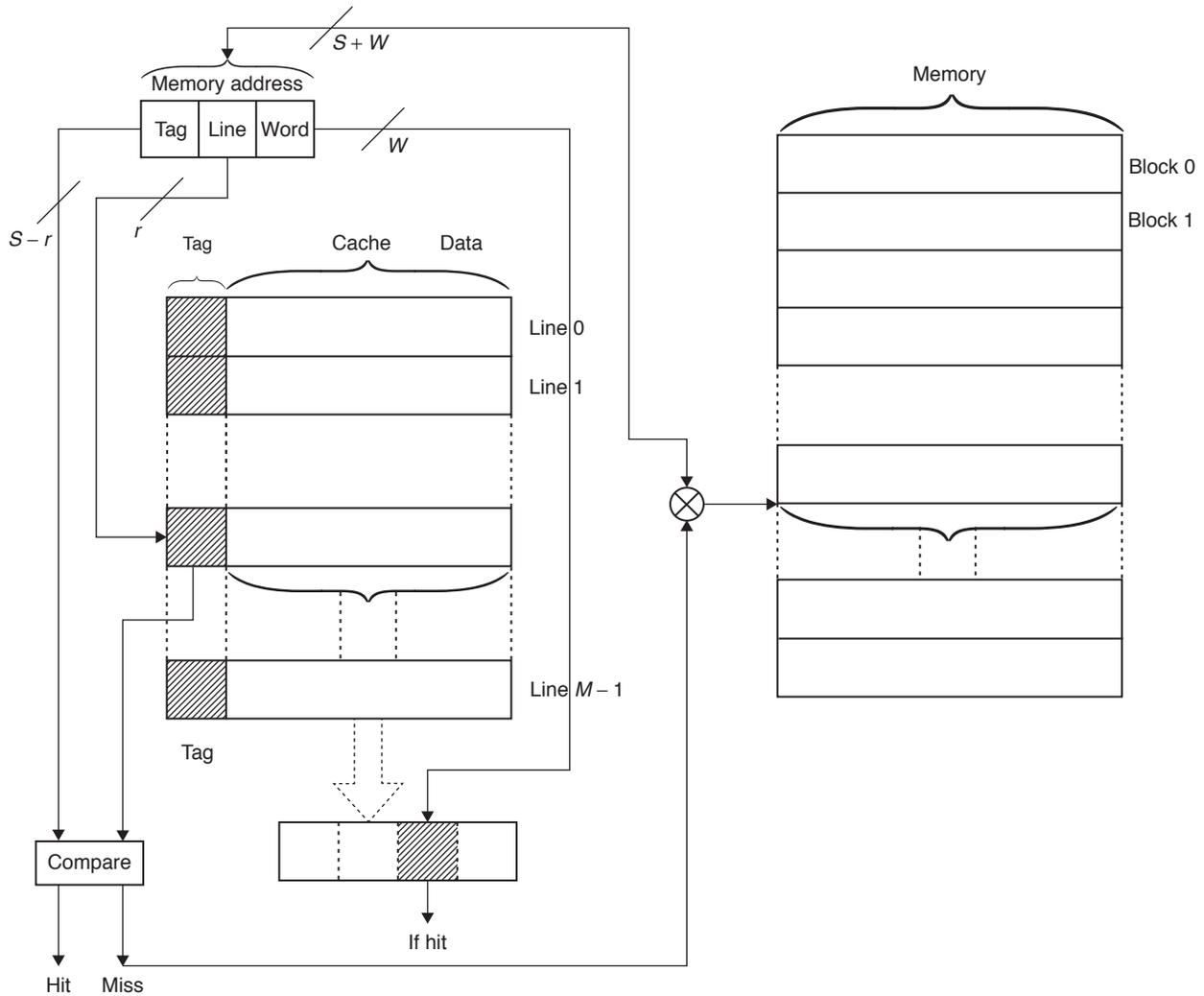


Figure 3 Direct Mapping

The effect of this mapping is that blocks of main memory are assigned to lines of the cache as follows:

Cache Line	Main Memory Blocks Assigned
0	0, m, 2m, ... 2 ^s - m
1	1, m+1, 2m+1, ... 2 ^s - m + 1
⋮	⋮
⋮	⋮
⋮	⋮
⋮	⋮
m - 1	m - 1, 2m - 1, 3m - 1, ... 2 ^s - 1

Example 1: Let cache capacity = 64 KB
 Line size = 4 B
 Main memory capacity = 16 MB = 2²⁴ B
 Using direct mapping, Address length = s + w = 24-bits
 Block size = 2² B

$$\text{Number of blocks in main memory} = \frac{2^{24}}{2^2} = 2^{22}$$

$$\text{Number of lines in cache} = m = 2^r = \frac{2^{16}}{2^2} = 2^{14}$$

$$\therefore \text{Size of tag} = s - r = 22 - 14 = 8$$

$$\therefore \text{Main memory address} =$$

Tag	Line	Word
8	14	2

The mapping becomes

Cache Line	Starting Memory Address of Blocks (Hexa)
0	00000, 010000, ... FF0000
1	000004, 010004, ... FF0004
⋮	⋮
⋮	⋮
⋮	⋮
⋮	⋮
2 ¹⁴ - 1	00FFFC, 01FFFC, ... FFFFFC

Note: No two blocks that map into the same line number have the same tag number.

Advantages:

- Simple and cheap
- The tag field is short; only those bits have to be stored which are not used to address the cache.
- Access is very fast.

Disadvantages: A given block fits into a fixed cache location, i.e., a given cache line will be replaced whenever there is a reference to another memory block which fits to the

same line, regardless what the status of the other cache line is.

This can produce a low hit ratio, even if only a very small part of the cache is effectively used.

Associative mapping This technique overcomes the disadvantage of direct mapping by permitting each main memory block to be loaded into any line of the cache. Here the cache control logic interprets a memory address as two fields.

1. Tag
2. Word

Figure shows associative mapping technique:

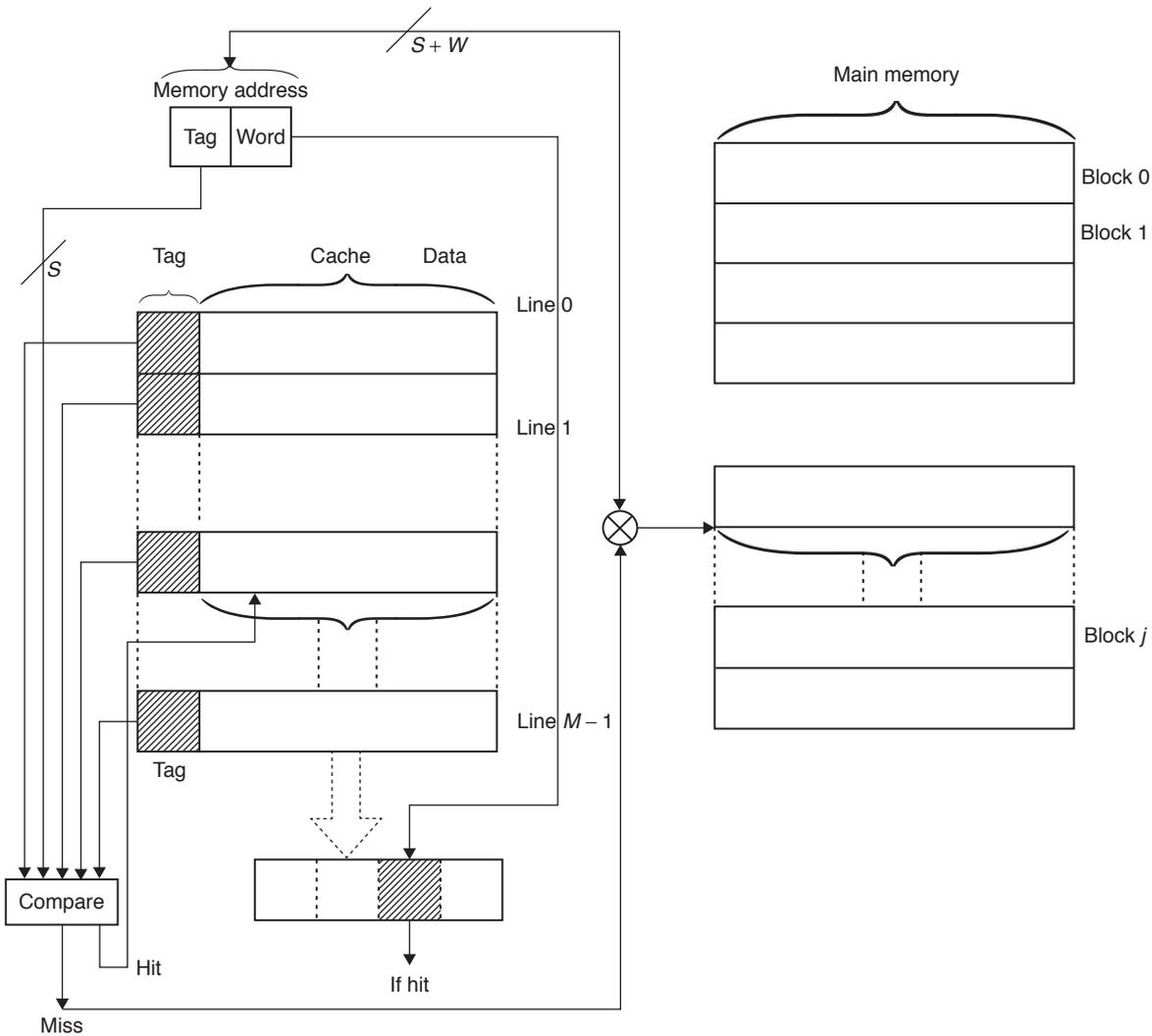


Figure 4 Associative mapping

To determine whether a block is in the cache, the cache control logic must simultaneously examine every line tag for a match. No field in the address corresponds to line number, so that the number of lines in the cache is not determined by the address format.

To summarize,

Address length = $(s + w)$ bits

Number of addressable units = 2^{s+w} words or bytes

Block size = line size = 2^w words or bytes

Number of blocks in main memory = $\frac{2^{s+w}}{2^w} = 2^s$

Number of lines in cache = undetermined

Size of tag = s -bits

Example 2: Cache size = 64 KB

Line size = 4 B

Main memory capacity = 16 MB

$$\begin{aligned} \text{Number of blocks in main memory} &= \frac{2^{24}}{2^2} \\ &= 2^{22}. \end{aligned}$$

∴ Size of tag = 24 – 2 = 22-bits

For example, the tag of the hexadecimal main memory address 16339C is 058CE7

Main memory address =

Tag	Word
22	2

Advantages: Associative mapping provides the highest flexibility concerning the line to be replaced when a new block is read into a cache.

Disadvantages:

- Complex
- The tag field is long
- Fast access can be achieved only using high performance associative memories for the cache, which is difficult and expensive.

Set-associative mapping: It exhibits strengths of both the direct and associative approaches and reduces their disadvantages.

Here the cache is divided into V sets, each of which consists of K lines

$$\text{i.e., } m = V \times K$$

$$i = j \text{ modulo } V$$

Where i = cache set number

j = main memory block number

m = number of lines in cache

As there are K lines in each set, this is referred as K -way set associative mapping. The cache control logic interprets a memory address simply as three fields.

1. Tag
2. Set
3. Word

The d set bits specify one of $V = 2^d$ sets. The S -bits of the tag and the set fields specify one of the 2^S blocks of main memory.

Figure 3 shows Set-associative mapping.

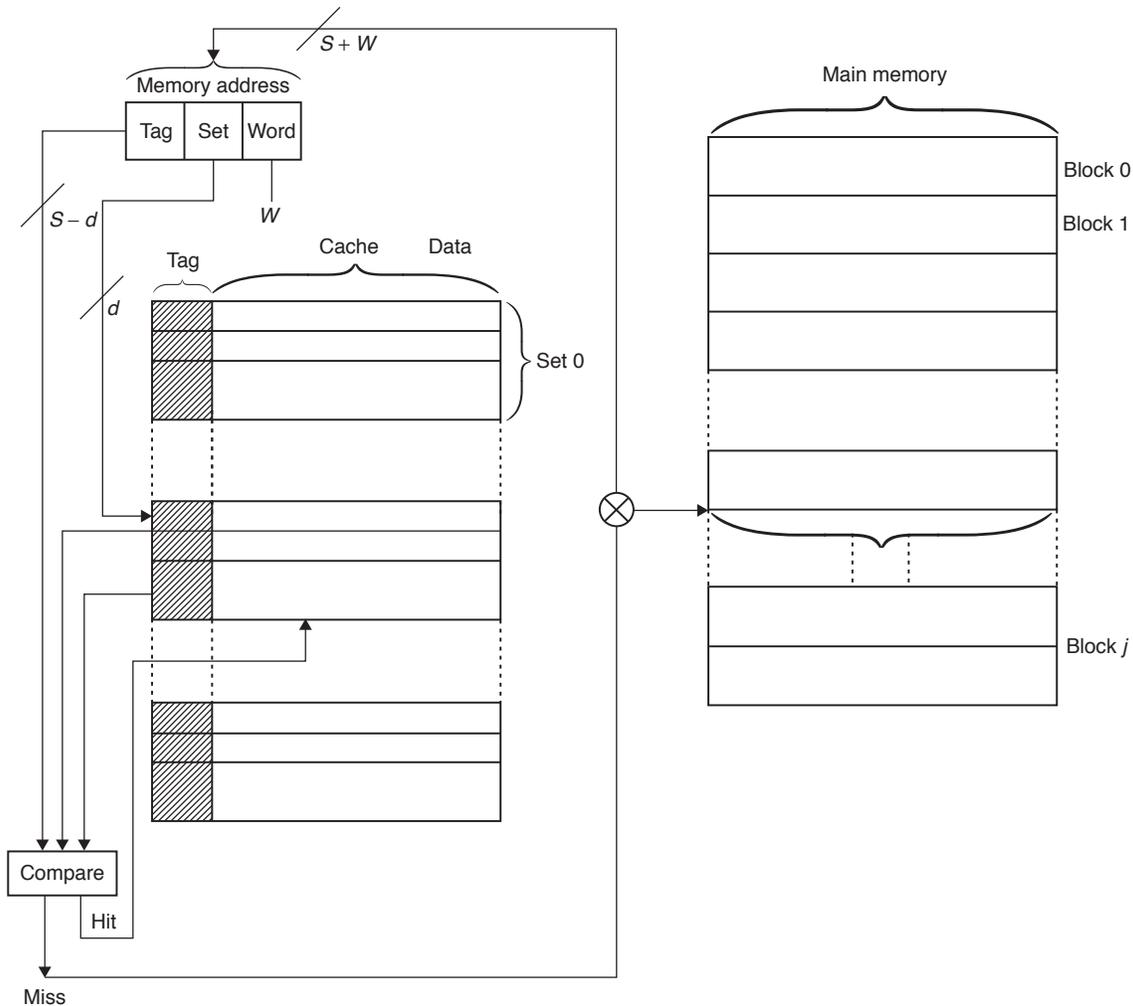


Figure 5 K-way set associative cache

Here the tag in a memory address is much smaller and is only compared to the K tags within a single set. To summarize,

Address length = $(s + w)$ bits

Number of Addressable units = 2^{s+w} words or bytes.

Block size = line size = 2^w words or bytes.

Number of blocks in main memory = $\frac{2^{s+w}}{2^w} = 2^s$

Number of lines in set = K

Number of sets $V = 2^d$

Number of lines in cache = $KV = K \times 2^d$

\therefore Size of tag = $(s - d)$ bits

Example 3: Cache capacity = 64 KB

Block size = 4 B

Main memory capacity = 16 MB

Number of blocks in main memory = $\frac{2^{24}}{2^2} = 2^{22}$

For 2-way set associative mapping,

Number of lines in a set $K = 2$

Number of sets = $V = 2^d$

Number of lines in cache = $K \times 2^d = \frac{2^{16}}{2^2}$

= 2^{14}

$\Rightarrow 2 \times 2^d = 2^{14}$

$\Rightarrow 2^d = 2^{13}$

$\Rightarrow d = 13$

\therefore Size of Tag = $22 - 13 = 9$

Hence main memory address =

Tag	Set	Word
9	13	2

In practice, 2 and 4-way set associative mapping are used with very good results. Larger sets do not produce further significant performance improvement.

If a set consist of a single line, i.e., it is direct mapping; If there is one single set consisting of all lines i.e., it is associative mapping.

Replacement algorithms

Once the cache has been filled, when a new block is brought into the cache, one of the existing blocks must be replaced. For direct mapping, there is only possible line for any particular block, and no choice is possible.

For associative and set associative techniques, a replacement algorithm is needed. Four of the most common replacement algorithms are

- (i) **LRU** (Least recently used): Replaces the block in the set that has been in the cache longest with no reference to it.

- (ii) **FIFO** (First-in-first-out): Replace the block in the set that has been in the cache longest.

- (iii) **LFU** (Least frequently used): Replace the block in the set that has experienced the fewest references.

- (iv) **Random**

Write policy

When a block that is resident in the cache is to be replaced, there are two cases to consider.

- (i) If the old block in the cache has not been altered, then it may be over-written with a new block without first writing out the old block.
- (ii) If at least one write operation has been performed on a word in that line of the cache, then main memory must be updated by writing the line of cache out of the block of memory before bringing in the new block.

The write policies are

- (a) **Write through:** All write operations are made to main memory as well as to the cache, ensuring that main memory is always valid.

Drawback: Creates substantial memory traffic

- (b) **Write back:** This technique minimizes memory writes. It updates are made only in the cache. When a block is replaced it is written back to main memory if and only if it is updated.

Drawback: There some portions of main memory are invalid and hence accesses by I/O modules can be allowed only through the cache.

Line size

Larger blocks reduce the number of blocks that fit into a cache. Because each block fetch overwrites older cache contents, a small number of blocks results in data being over written shortly after they are fetched.

As a block becomes larger, each additional word is farther from the requested word, therefore less likely to be needed in the near future.

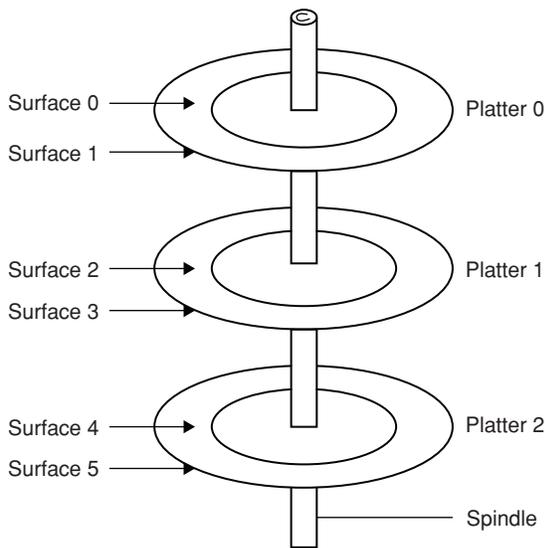
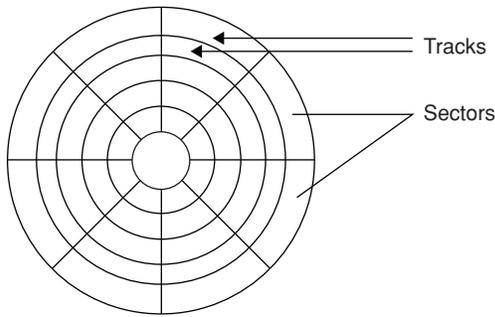
Number of caches

Multilevel caches We may have on-chip cache as well as external cache. This is a two level cache organization, with the internal cache designated as level 1, and external cache designated as level 2.

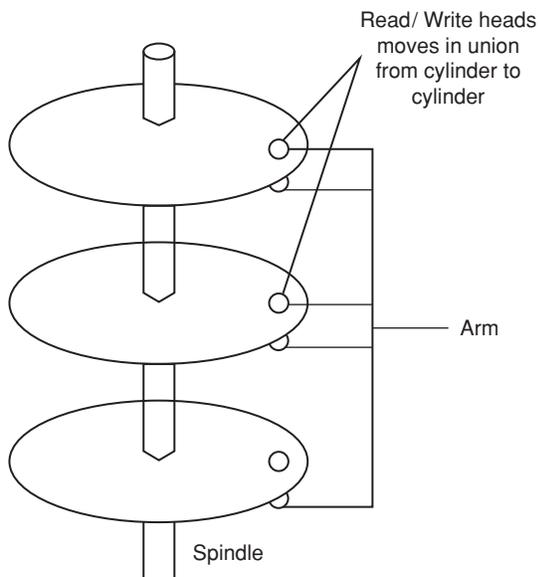
SECONDARY STORAGE

Disk

Disk consists of platters, each with two surfaces. Each surface consists of concentric rings called tracks. Each track consists of sectors separated by gaps.



Disk operation: The disk surface spins at a fixed rotational rate. There is a read/write head attached to the end of the arm and flies over the disk surface on a thin cushion of air. By moving radially the arm can position the read/write head over any track.



Disk access time: Average time to access some target sector:

$$T_{ae} = T_{avg\ seek} + T_{avg\ rotation} + T_{avg\ transfer}$$

Where $T_{avg\ seek}$ is typical 9 ms.

$$T_{avg\ rotation} = \frac{1}{2} \times \frac{1}{RPM} \times 60 \text{ Sec/1min}$$

$$T_{avg\ rotation} = \frac{1}{RPM} \times \frac{1}{(avg\ sector/track)} \times 60 \text{ Sec/1min}$$

Notes:

1. Seek time is the Time to position heads over cylinder containing target sectors ($T_{avg\ seek}$).
 2. Rotational Latency is the time waiting for first bit of target sector to pass under read/write head. ($T_{avg\ rotation}$).
 3. Transfer Time is the time to read the bits in the target sector ($T_{avg\ transfer}$).
- Data are recorded on the surface of a hard disk made of metal coated with magnetic material.
 - The disks and the drive are usually built together and encased in an air tight container to protect the disk from pollutants such as smoke particle and dust. Several disks are usually started on a common drive shaft with each disk having its own read/write head.

Diskette

Data are recorded on the surface of a floppy disk made of polyester coated with magnetic material.

A special diskette drive must be used to access data stored in the floppy disk. It works much like a record turntable of Gramophone.

Main features

- Direct access
- Cheap
- Portable, convenient to use

Main Standards

- 5¹/₄ inch capacity \cong 360 KB/ disk
- 3¹/₂ inch capacity \cong 1.44 MB/disk (about 700 pages of A₄ text)

Magnetic Tape

Magnetic tape is made up from a layer of plastic which is coated with iron oxide. The oxide can be magnetized in different directions to represent data. The operation uses a similar principle as in the case of a tape recorder.

Main features

- Sequential access (access time about 1.55)
- High value of storage (50 MB/tape)
- Inexpensive

It is often used for Batch up or archive purpose.

Optimal Memory

CD-ROM (Compact disk ROM): The disk surface is imprinted with microscopic holes which record digital information. When a low-powered power beam shines on the surface, the intensity of the reflected light changes as it encounters a hole. The change is detected by a photo sensor and converted into a digital signal.

- Huge capacity: 775 MB/disk (≈ 550 diskette)
- Inexpensive replication, cheap production.
- Removable, read only.
- Long access time (could be half a second)

WORM (Write Once Read Memory) CD: A lower beam of modest intensity equipped in the disk drive is used to imprint the hole pattern.

- Good for archival storage by providing a permanent record of large volumes of data.

Erasable Optical Disk: Combination of Laser technology and magnetic surface technique.

- Can be repeatedly written and overwritten.
- High reliability and longer life than magnetic disks.

EXERCISE

Practice Problems I

Directions for questions 1 to 20: Select the correct alternative from the given choices.

- Find the number of bits in the cache index and tag for a direct mapped cache of size 32 KB with block size of 32 bytes. The CPU generates 48-bit addresses.
(A) 33,15 (B) 15,10
(C) 10,33 (D) 15,33
- Given the cache access time is 200 ns and the memory access time is 400 ns. If the effective access time is 20% greater than the cache access time, what is the hit ratio?
(A) 80% (B) 20%
(C) 40% (D) 100%
- A computer system has an L_1 cache, an L_2 cache and a main memory unit connected as shown below. The block size in L_1 cache is 2 words. The block size in L_2 cache is 8 words. The memory access times are 2 nanoseconds, 20 nanoseconds and 200 nanoseconds for L_1 cache, L_2 cache and main memory unit, respectively.



When there is a miss in L_1 cache and a hit in L_2 cache, a block is transferred from L_2 cache to L_1 cache. What is the time taken for this transfer?

- (A) 22 ns (B) 44 ns
(C) 66 ns (D) 88 ns
- In direct memory management, CPU references address of 15-bits. Main memory size is $512 * 8$ and cache memory size is $128 * 8$. Tag and line are respectively
(A) 9, 7 (B) 7, 9
(C) 15, 7 (D) 7, 15
 - Consider a cache with 64 blocks and a block size of 16 bytes. The byte address of 1200 maps to ____ block number.
(A) 10 (B) 11
(C) 64 (D) 16

- In a cache memory, cache line is 64 bytes. The main memory has latency of 32 ns and bandwidth of 1 GB/sec. Then the time required to fetch the entire cache line from main memory is

(A) 32 ns (B) 64 ns
(C) 96 ns (D) 128 ns

- A set associative cache consists of 64 lines or slots divided into four-line sets. Main memory contains 4K blocks of 128 word each. Then the number of bits present in tag, set and word fields are respectively.

(A) 7, 6, 7 (B) 6, 7, 7
(C) 4, 8, 7 (D) 8, 4, 7

- A 2-way set-associative cache has lines of 32 bytes and a total size of 16 KB. The 32 MB main memory is byte addressable. Then which of the following two memory addresses mapped to same set?

(A) 10D6A32, 035C3A2
(B) 2A36D01, 2A3C530
(C) 10D63A2, 035C3A0
(D) 2A36D08, 0A3C538

- Let the cache memory capacity is 64 KB and main memory capacity is 16 MB. Let block size is 4 bytes. Then the tag, line, word fields in hexadecimal notation for the main memory address ccccc using direct mapped cache will be

(A) cc, ccc, c (B) cc, 3333, 0
(C) cc, ccc, 0 (D) cc, 333, 30

- Consider a 32-bit microprocessor that has an on-chip 16 KB four-way set associative cache. Assume that the cache has a line size of four 32-bit words. Then the word in the memory location ABCDE8F8 will be mapped to

(A) 143rd set (B) 815th set
(C) 255th set (D) 0th set

- Given the following specifications for an external cache memory:

Four-way set associative, Line size of two 16-bit words; Able to accommodate a total of 4K 32-bit words from

main memory. Used with a 16-bit processor that issues 24-bit address. Then the number of bits used to represent set field is

- (A) 2-bits (B) 10-bits
(C) 12-bits (D) 14-bits

12. Consider a machine with a byte addressable main memory of 2^{16} bytes and block size of 8 bytes. Assume that a direct mapped cache consisting of 32 lines is used with this machine. Then in what line would bytes with the address 1100 0011 0011 0100 is stored?
(A) slot 3 (B) slot 4
(C) slot 6 (D) slot 12
13. A computer system contains a main memory of 32 K 16-bit words. It also has a 4 K-word cache divided into four line sets with 64 words per line. The processor fetches words from locations 0, 1, 2, ..., 4351 in that order. It then repeats this fetch sequence 10 more times. The cache is 10 times faster than main memory. Then the improvement resulting from the use of the cache is (assume an LRU policy is used for block replacement)
(A) 0.63 (B) 0.45
(C) 1.21 (D) 2.18
14. Consider an L_1 cache with an access time of 1ns and a hit ratio of $H = 0.95$. Suppose that we can change the cache design such that we increase H to 0.98, but increase access time to 1.5ns. Which of the following condition is met for this change to result in improved performance?
(A) Next level memory access time must be less than 16.67
(B) Next level memory access time must be greater than 16.67
(C) Next level memory access time must be less than 50
(D) Next level memory access time must be greater than 50
15. Consider a single-level cache with an access time of 2.5 ns and a line size of 64 bytes and a hit ratio of $H = 0.95$. Main memory uses a block transfer capability that

has a first-word (4 bytes) access time of 50 ns and an access time of 5ns for each word thereafter. What is the access time when there is a cache miss?

- (A) 130 ns (B) 149.4 ns
(C) 2.375 ns (D) 8.875 ns

16. The tag, block and word fields of main memory address using direct mapping technique for 2048 main memory blocks, 128 blocks of cache memory and block size of 16:
(A) 4, 7, 4 (B) 7, 4, 4
(C) 11, 7, 4 (D) Data insufficient
17. Let H_1 is level 1 cache hit ratio, H_2 is level 2 cache hit ratio, C_1 is the time required to access Level 1 cache, C_2 is the time required to access Level 2 cache and M is the time required to access Main memory. Then the average access time required by the processor is
(A) $H_1 C_1 + (1 - H_1) H_2 (C_2) + (1 - H_1) (1 - H_2) (M)$
(B) $H_1 C_1 + (1 - H_1) H_2 (C_1 + C_2) + (1 - H_1) (1 - H_2) (C_1 + C_2 + M)$
(C) $H_1 C_1 + H_1 H_2 (C_1 + C_2) + H_1 H_2 (C_1 + C_2 + M)$
(D) $H_1 C_1 + (1 - H_1) H_2 (C_1 \cdot C_2) + (1 - H_1) (1 - H_2) (C_1 \cdot C_2 \cdot M)$
18. If $p = 2^m$ be the number of lines in cache and $b = 2^n$ be the size of each block, then total words that can be stored in cache memory is given by
(A) 2^{m+n} (B) 2^{m-n}
(C) $m + n$ (D) $p + b$
19. Cache memory enhances
(A) memory capacity
(B) memory access time
(C) secondary storage capacity
(D) secondary storage access time
20. Which of the following property allows the processor to execute a number of clustered locations?
(A) Spatial (B) Temporal
(C) Inclusion (D) Coherence

Practice Problems 2

Directions for questions 1 to 20: Select the correct alternative from the given choices.

1. If average access time of CPU is 20 ns, access time of main memory is 110 ns and the cache access time is 10 ns. What is the hit ratio?
(A) 100% (B) 90%
(C) 80% (D) 70%
2. A hard disk spins at 180 revolutions per minute. What is the average rotational latency?
(A) 0.16 sec (B) 0.32 sec
(C) 0.2 sec (D) 0.4 sec
3. A disk pack have 16 surfaces, with 128 tracks per surface and 256 sectors per track. 512 bytes of data are stored in a bit serial manner in a sector. The number of bits required to specify a particular sector in the disk is
(A) 4 (B) 7
(C) 11 (D) 19
4. A disk has 19456 cylinders, 16 heads and 63 sectors per track. The disk spins at 5400 rpm. Seek time between adjacent tracks is 2 ms. Assuming the read/write head is already positioned at track 0, how long does it take to read the entire disk?
(A) 48 min (B) 58 min
(C) 64 min (D) 72 min

5. A certain moving arm disk storage with one head has following specifications:
 Number of tracks/recording surface = 200
 Disk Rotation Speed = 2400 rpm
 Track storage capacity = 62500-bits
 The average latency time (assuming that head can move from one track to another only by traversing the entire track) is
 (A) 0.125 sec (B) 1.25 sec
 (C) 0.0125 sec (D) 12.5 sec
6. In Memory management system, cache memory access time is 100 ns and main memory access time is 200 ns. Number of CPU references is 100 and number of hits is 10. Average access time is
 (A) 150 ns (B) 100 ns
 (C) 190 ns (D) 280 ns
7. The seek time of disk is 40 m sec. It rotates at the rate of 40 rps. The capacity of each track is 400 words. The access time is
 (A) 50 m sec (B) 53 m sec
 (C) 60 m sec (D) 63 m sec
8. An Associated cache and one million word main memory are divided into 256 word blocks. How many blocks are there?
 (A) 2^8 (B) 2^{12}
 (C) 2^{20} (D) 2^{28}
9. The average access time of a disk is
 (A) Seek time + Rotational latency time
 (B) Seek time
 (C) Rotational latency + transfer time + seek time
 (D) Rotation latency + transfer time.
10. What will be the size of the memory whose last memory location is FFFF?
 (A) 64 k (B) 32 k
 (C) 10 k (D) 24 k
11. Data from a cassette tape is obtained by ____ accessing method.
 (A) Parallel (B) Serial
 (C) Sequential (D) Random
12. For a memory system, the desirable characteristics is/are
 (A) Speed and reliability
 (B) Durability and compactness
 (C) Low power consumption
 (D) All of these
13. The memory that has the shortest access time is
 (A) Magnetic bubble (B) Magnetic core memory
 (C) Cache memory (D) RAM
14. Cache memory
 (A) has greater capacity than RAM.
 (B) enhances secondary storage access time.
 (C) is faster to access than registers.
 (D) is faster to access than main memory
15. Consider a disk pack with 16 surfaces 128 tracks per surface and 256 sectors per track. 512 bytes of data are stored in bit and serial manner, Then the capacity of the disk is
 (A) 256 MB (B) 256 KB
 (C) 512 MB (D) 64 MB
16. Principle of locality justifies the use of
 (A) Cache (B) DMA
 (C) Disk (D) RAM
17. The main memory of a computer has $2ab$ blocks while cache has $2a$ blocks. If the cache uses the set associative mapping scheme with two blocks per set, then block k of main memory maps to the set:
 (A) $(k \bmod b)$ of the cache (B) $(k \bmod a)$ of cache
 (C) $(k \bmod 2a)$ of cache (D) $(k \bmod 2ab)$ of cache
18. Which of the following factors do not affect the hit ration of cache?
 (A) Block replacement algorithms.
 (B) Block frame size
 (C) Cycle counts
 (D) Main memory size
19. In which of the following mapping function, there is no need of replacement algorithm?
 (A) Direct Mapping
 (B) Set-associative mapping
 (C) Full associative mapping
 (D) Both (A) and (B)
20. In a direct mapping, the index field equals to
 (A) Sum of tag and word fields
 (B) Sum of block and word fields
 (C) Sum of tag an block fields
 (D) Same as block field

PREVIOUS YEARS' QUESTIONS

1. Consider a small two-way set-associative cache memory, consisting of four blocks. For choosing the block to be replaced, use the least recently used (LRU) scheme. The number of cache misses for the following sequence of block addresses is 8, 12, 0, 12, 8 [2004]
 (A) 2 (B) 3
 (C) 4 (D) 5
2. Consider a direct mapped cache of size 32 KB with block size 32 bytes. The CPU generates 32 bit addresses. The number of bits needed for cache indexing and the number of tag bits are respectively. [2005]
 (A) 10, 17 (B) 10, 22
 (C) 15, 17 (D) 5, 17

Common data for questions 3 and 4: Consider two cache organizations: The first one is 32 KB 2-way set associative with 32-byte block size. The second one is of the same size but direct mapped. The size of an address is 32 bits in both cases. A 2-to-1 multiplexer has a latency of 0.6 ns while a k-bit comparator has a latency of $k/10$ ns. The hit latency of the set associative organization is h_1 while that of the direct mapped one is h_2 .

3. The value of h_1 is: [2006]
 (A) 2.4 ns (B) 2.3 ns
 (C) 1.8 ns (D) 1.7 ns
4. The value of h_2 is: [2006]

Data for question 5: Consider a machine with a byte addressable main memory of 2^{16} bytes. Assume that a direct mapped data cache consisting of 32 lines of 64 bytes each is used in the system. A 50×50 two-dimensional array of bytes is stored in the main memory starting from memory location 1100 H. Assume that the data cache is initially empty. The complete array is accessed twice. Assume that the contents of the data cache do not change in between the two accesses.

5. Which of the following lines of the data cache will be replaced by new blocks in accessing the array for the second time? [2007]
 (A) line 4 to line 11 (B) line 4 to line 12
 (C) line 0 to line 7 (D) line 0 to line 8
6. For inclusion to hold between two cache levels L_1 and L_2 in a multi-level cache hierarchy, which of the following are necessary? [2008]
 (i) L_1 must be a write-through cache
 (ii) L_2 must be a write-through cache
 (iii) The associativity of L_2 must be greater than that of L_1
 (iv) The L_2 cache must be atleast as large as the L_1 cache
 (A) (iv) only (B) (i) and (iv) only
 (C) (i), (ii) and (iv) only (D) (i), (ii), (iii) and (iv)

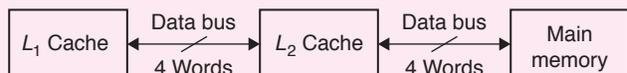
Common data for questions 7, 8 and 9: Consider a machine with a 2-way set associative data cache of size 64K-bytes and block size 16-bytes. The cache is managed using 32-bit virtual addresses and the page size is 4Kbytes. A program to be run on this machine begins as follows:

```
double ARR [1024] [1024];
int i, j;
/* Initialize array ARR to 0.0 */
for (i = 0; i < 1024; i++)
  for (j = 0; j < 1024; j++)
    ARR [i] [j] = 0.0;
```

The size of double is 8 Bytes. Array ARR is located in memory starting at the beginning of virtual page 0XFF000 and stored in row major order. The cache is initially empty and no pre-fetching is done. The only data memory references made by the program are those to array ARR.

7. The total size of the tags in the cache directory is [2008]
 (A) 32K-bits (B) 34K-bits
 (C) 64K-bits (D) 68K-bits
8. Which of the following array elements has the same cache index as ARR [0] [0]? [2008]
 (A) ARR [0] [4] (B) ARR [4] [0]
 (C) ARR [0] [5] (D) ARR [5] [0]
9. The cache hit ratio for this initialization loop is [2008]
 (A) 0% (B) 25%
 (C) 50% (D) 75%
10. Consider a 4-way set associative cache (initially empty) with total 16 cache blocks. The main memory consists of 256 blocks and the request for memory blocks is in the following order: [2009]
 0, 255, 1, 4, 3, 8, 133, 159, 216, 129, 63, 8, 48, 32, 73, 92, 155.
 Which one of the following memory block will NOT be in cache if LRU replacement policy is used?
 (A) 3 (B) 8
 (C) 129 (D) 216

Common data questions 11 and 12: A computer system has an L_1 cache, an L_2 cache, and a main memory unit connected as shown below. The block size in L_1 cache is 4 words. The block size in L_2 cache is 16 words. The memory access times are 2 nanoseconds, 20 nanoseconds and 200 nanoseconds, for L_1 cache, L_2 cache and main memory unit respectively.

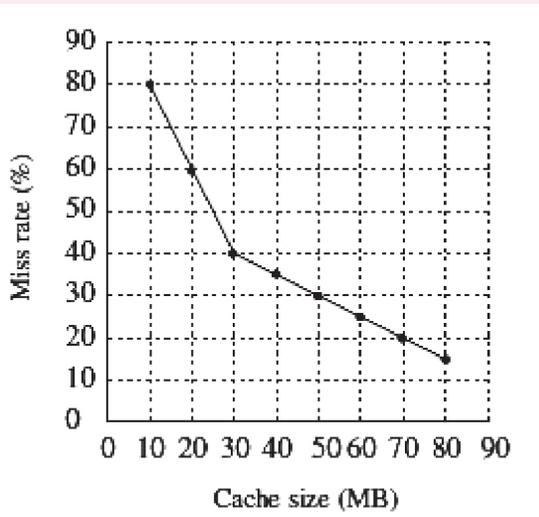


11. When there is a miss in L_1 cache and a hit in L_2 cache, a block is transferred from L_2 cache to L_1 cache. What is the time taken for this transfer? [2010]
 (A) 2 ns (B) 20 ns
 (C) 22 ns (D) 88 ns
12. When there is a miss in both L_1 cache and L_2 cache, first a block is transferred from main memory to L_2 cache, and then a block is transferred from L_2 cache to L_1 cache. What is the total time taken for these transfers? [2010]
 (A) 222 ns (B) 888 ns
 (C) 902 ns (D) 968 ns
13. An 8 KB direct-mapped write-back cache is organized as multiple blocks, each of size 32 bytes. The processor generates 32-bit addresses. The cache controller maintains the tag information for each cache block comprising of the following.
 1 Valid bit
 1 Modified bit
 As many bits as the minimum needed to identify the memory block mapped in the cache.
 What is the total size of memory needed at the cache controller to store meta-data (tags) for the cache? [2011]
 (A) 4864 bits (B) 6144 bits
 (C) 6656 bits (D) 5376 bits
- Common data for questions 14 and 15:** A computer has a 256 KB, 4-way set associative, write back data cache with block size of 32 bytes. The processor sends 32 bit addresses to the cache controller. Each cache tag directory entry contains, in addition to address tag, 2 valid bits, 1 modified bit and 1 replacement bit.
14. The number of bits in the tag field of an address is [2012]
 (A) 11 (B) 14
 (C) 16 (D) 27
15. The size of the cache tag directory is [2012]
 (A) 160 K-bits (B) 136 K-bits
 (C) 40 K-bits (D) 32 K-bits
 (A) 2.4 ns (B) 2.3 ns
 (C) 1.8 ns (D) 1.7 ns
16. In a k -way set associative cache, the cache is divided into v sets, each of which consists of k lines. The lines of a set are placed in sequence one after another. The lines in set s are sequenced before the lines in set $(s + 1)$. The main memory blocks are numbered 0 onwards. The main memory block numbered j must be mapped to any one of the cache lines from [2013]
 (A) $(j \bmod v) * k$ to $(j \bmod v) * k + (k - 1)$
 (B) $(j \bmod v)$ to $(j \bmod v) + (k - 1)$
 (C) $(j \bmod k)$ to $(j \bmod k) + (v - 1)$
 (D) $(j \bmod k) * v$ to $(j \bmod k) * v + (v - 1)$
17. An access sequence of cache block addresses is of length N and contains n unique block addresses. The number of unique block addresses between two consecutive accesses to the same block address is bounded above by K . What is the miss ratio if the access sequence is passed through a cache of associativity $A \geq K$ exercising least-recently used replacement policy? [2014]
 (A) n/N (B) $1/N$
 (C) $1/A$ (D) K/n
18. A 4-way set -associative cache memory unit with a capacity of 16 KB is built using a block size of 8 words. The word length is 32-bits. The size of the physical address space is 4 GB. The number of bits for the TAG field is _____. [2014]
19. In designing a computer's cache system, the cache block (or cache line) size is an important parameter. Which one of the following statements is correct in this context? [2014]
 (A) A smaller block size implies better spatial locality.
 (B) A smaller block size implies a smaller cache tag and hence lower cache tag overhead.
 (C) A smaller block size implies a larger cache tag and hence lower cache hit time.
 (D) A smaller block size incurs a lower cache miss penalty.
20. Consider a main memory system that consists of 8 memory modules attached to the system bus, which is one word wide. When a write request is made, the bus is occupied for 100 nanoseconds (ns) by the data, address, and control signals. During the same 100 ns, and for 500 ns thereafter, the addressed memory module executes one cycle accepting and storing the data. The (internal) operation of different memory modules may overlap in time, but only one request can be on the bus at any time. The maximum number of stores (of one word each) that can be initiated in 1 millisecond is _____. [2014]
21. If the associativity of processor cache is doubled while keeping the capacity and block size unchanged, which one of the following is guaranteed to be NOT affected? [2014]
 (A) Width of tag comparator
 (B) Width of set index decoder
 (C) Width of way selection multiplexer
 (D) Width of processor to main memory data bus
22. The memory access time is 1 nanosecond for a read operation with a hit in cache, 5 nanoseconds for a read operation with a miss in cache, 2 nanoseconds for a write operation with a hit in cache and 10 nanoseconds for a write operation with a miss in cache. Execution of a sequence of instructions involves 100 instruction fetch operations. 60 memory operand read operations

and 40 memory operand write operations. The cache - hit ratio is 0.9. The average memory access time (in nanoseconds) in executing the sequence of instructions is _____. [2014]

23. Assume that for a certain processor, a read request takes 50 nanoseconds on a cache miss and 5 nanoseconds on a cache hit. Suppose while running a program, it was observed that 80% of the processor's read requests result in a cache hit. The average read access time in nanoseconds is _____. [2015]
24. A computer system implements a 40-bit virtual address, page size of 8 kilobytes, and a 128-entry translation look-aside buffer (TLB) organized into 32 sets each having four ways. Assume that the TLB tag does not store any process id. The minimum length of the TLB tag in bits is _____. [2015]
25. Consider a machine with a byte addressable main memory of 2^{20} bytes, block size of 16 bytes and a direct mapped cache having 2^{12} cache lines. Let the addresses of two consecutive bytes in main memory be $(E201F)_{16}$ and $(E2020)_{16}$. What are the tag and cache line address (in hex) for main memory address $(E201F)_{16}$? [2015]
 (A) E, 201 (B) F, 201
 (C) E, E20 (D) 2, 01F
26. The width of the physical address on a machine is 40 bits. The width of the tag field in a 512KB 8-way set associative cache is ____ bits. [2016]
27. A file system uses an in - memory cache to cache disk blocks. The miss rate of the cache is shown in the figure. The latency to read a block from the cache is 1 ms and to read a block from the disk is 10ms. Assume that the cost of checking whether a block exists in the cache is negligible. Available cache sizes are in multiples of 10MB.

The smallest cache size required to ensure an average read latency of less than 6 ms is ____ MB. [2016]



28. Consider a two-level cache hierarchy with L1 and L2 caches. An application incurs 1.4 memory accesses per instruction on average. For this application, the miss rate of L1 cache is 0.1; the L2 cache experiences on average 7 misses per 1000 instructions. The miss rate of L2 expressed correct to two decimal places is _____. [2017]
29. Consider a 2-way set associative cache with 256 blocks and uses LRU replacement. Initially the cache is empty. Conflict misses are those misses which occur due to contention of multiple blocks for the same cache set. Compulsory misses occur due to first time access to the block. The following sequence of accesses to memory blocks (0, 128, 256, 128, 0, 128, 256, 128, 1, 129, 257, 129, 1, 129, 257, 129) is repeated 10 times. The number of *conflict misses* experienced by the cache is _____. [2017]
30. A cache memory unit with capacity of N words and block size of B words is to be designed. If it is designed as a direct mapped cache, the length of the TAG field is 10 bits. If the cache unit is now designed as a 16-way set-associative cache, the length of the TAG field is ____ bits. [2017]
31. In a two-level cache system, the access times of L_1 and L_2 caches are 1 and 8 clock cycles, respectively. The miss penalty from the L_2 cache to main memory is 18 clock cycles. The miss rate of L_1 cache is twice that of L_2 . The average memory access time (AMAT) of this cache system is 2 cycles. The miss rates of L_1 and L_2 respectively are: [2017]
 (A) 0.111 and 0.056 (B) 0.056 and 0.111
 (C) 0.0892 and 0.1784 (D) 0.1784 and 0.0892

32. The read access times and the hit ratios for different caches in a memory hierarchy are as given below.

Cache	Read access time (in nanoseconds)	Hit ratio
I-cache	2	0.8
D-cache	2	0.9
L2-cache	8	0.9

The read access time of main memory is 90 nanoseconds. Assume that the caches use the referred- word-first read policy and the write back policy. Assume that all the caches are direct mapped caches. Assume that the dirty bit is always 0 for all the blocks in the caches. In execution of a program, 60% of memory reads are for instruction fetch and 40% are for memory operand fetch. The average read access time in nanoseconds (up to 2 decimal places) is _____. [2017]

33. Consider a machine with a byte addressable main memory of 2^{32} bytes divided into blocks of size 32

bytes. Assume that a direct mapped cache having 512 cache lines is used with this machine. The size of the tag field in bits is _____.

[2017]

34. The size of the physical address space of a processor is 2^P bytes. The word length is 2^W bytes. The capacity of cache memory is 2^N bytes. The size of each cache

block is 2^M words. For a K -way set-associative cache memory, the length (in number of bits) of the tag field is:

[2018]

- (A) $P - N - \log_2 K$
 (B) $P - N + \log_2 K$
 (C) $P - N - M - W - \log_2 K$
 (D) $P - N - M - W + \log_2 K$

ANSWER KEYS

EXERCISES

Practice Problems 1

1. D 2. A 3. D 4. A 5. B 6. C 7. D 8. C 9. B 10. A
 11. B 12. C 13. D 14. B 15. A 16. A 17. B 18. A 19. B 20. A

Practice Problems 2

1. B 2. A 3. D 4. B 5. C 6. D 7. B 8. B 9. C 10. A
 11. C 12. D 13. C 14. D 15. A 16. A 17. B 18. D 19. A 20. B

Previous Years' Questions

1. C 2. A 3. A 4. D 5. C 6. B 7. D 8. B 9. C 10. D
 11. D 12. D 13. D 14. C 15. A 16. A 17. A 18. 20 19. D
 20. 10000 21. D 22. 1.68 23. 14 24. 22 25. A 26. 24 27. 30 28. 0.05
 29. 76 30. 14 31. A 32. 4.72 33. 18 34. B

TEST

COMPUTER ORGANIZATION AND ARCHITECTURE

Time: 60 min.

Directions for questions 1 to 30: Select the correct alternative from the given choices.

- Which of the following register keeps track of instruction execution sequence?
(A) Accumulator (B) Program counter
(C) Stack pointer (D) Instruction register

- Consider the following Register Transfer Language:

$$R_1 \leftarrow R_1 + M[R_2 + R_3]$$

Where R_1 , R_2 and R_3 are the CPU registers and 'M' is a memory location in primary memory, which addressing mode is suitable for above register transfer language?

- Indirect (B) Direct
(C) Indexed (D) Displacement
- Which of the following is/are advantage(s) of using a multiple-bus architecture over a single-bus architecture?
(i) Multiple-bus architecture reduces propagation delay.
(ii) Multiple-bus architecture reduces bottleneck effects.
(A) (i) only (B) (ii) only
(C) Both (i) and (ii) (D) Neither (i) nor (ii)
 - Which of the following statement is false with respect to Booth's Multiplication Algorithm?
(i) Corrections required for the final result.
(ii) Sign bit is protected due to internal arithmetic shift.
(iii) More space required to maintain the sum.
(A) (i), (ii) only (B) (ii), (iii) only
(C) (i), (iii) only (D) (i), (ii), (iii)
 - After selective complement of $A = 1100$ with $B = 0101$, the resultant A will be
(A) 0000 (B) 1100
(C) 0101 (D) 1001
 - Which type of shift operation always keeps the sign bit unchanged?
(A) Logical shift (B) Arithmetic shift
(C) Circular shift (D) Any right shift
 - Consider the register transfer language instructions:
 $AC \leftarrow M[R_1];$
 $R_1 \leftarrow R_1 + 1;$
Which addressing mode is specified by the instructions?
(A) Register addressing mode
(B) Register indirect mode
(C) Auto-increment mode
(D) Relative mode

- Which of the following statement is true?
(A) Floating point representation is better than fixed point representation.
(B) Fixed point representation is better than floating point representation.
(C) Datapath is same as ALU.
(D) Both (A) and (C)
- Which of the following statements correctly specifies about overflow?
(i) When adding two unsigned numbers the carry-out, from the MSB position serves as the overflow indicator.
(ii) Overflow can occur only by adding two signed numbers that have the same sign.
(A) (i) only (B) (ii) only
(C) Both (i) and (ii) (D) Neither (i) nor (ii)
- A certain processor supports only the immediate and direct addressing modes. Which of the following programming language features cannot be implemented on this processor?
(A) Pointers (B) Arrays
(C) Records (D) All of these
- The special purpose storage location(s) used by both ALU and CU are
(A) Decoders (B) Demultiplexers
(C) Registers (D) Buffers
- Which of the following is a component of the datapath of Von Neumann machine?
(i) Registers
(ii) ALU input bus
(iii) ALU I/O registers
(A) (i), (ii) only (B) (ii), (iii) only
(C) (i), (ii), (iii) (D) None of these
- Which of the following is/are false with respect to single-bus datapath?
(i) It is simplest and least expensive.
(ii) No limit on the amount of data transfer in a single clock cycle.
(A) (i) only (B) (ii) only
(C) Both (i) and (ii) (D) Neither (i) nor (ii)
- If we have shifted the significant to the right by a single position, then
(A) Add one to the exponent
(B) Subtract one from the exponent
(C) Don't change the exponent
(D) Data insufficient

15. In which addressing mode, the effective address of the operand is generated by adding a constant value to the content of a register?
 (A) Absolute mode (B) Indirect mode
 (C) Immediate mode (D) Index mode
16. What is the number of instructions required to add ' n ' numbers and store the result in memory using only one-address instructions?
 (A) n (B) $n - 1$
 (C) $n + 1$ (D) independent of n
17. Which unit of a computer system executes program, communicates with and often controls the operation of other subsystems?
 (A) CPU (B) ALU
 (C) I/O module (D) DMA
18. The multiplicand register and multiplier register of a hardware circuit implementing booth's algorithm have 1001 and 1100 respectively. The resultant will be
 (A) 10011100 (B) 00011100
 (C) 01101100 (D) 00010010
19. A floating point number has sign bit 0, Excess-64 exponent is 1010100 and fractional part is 0000000000011011. After converting this number to normalized form, the exponent (in decimal) will be
 (A) 20 (B) 9
 (C) 31 (D) 0
20. In IEEE floating point single precision representation, the number of bits in the fractional part is
 (A) 24
 (B) 23
 (C) 32
 (D) Depends on the architecture
21. After multiplying the binary numbers 010111 and 110110 using booth's multiplication algorithm, the resultant will be
 (A) -1242 (B) 1242
 (C) 230 (D) -230
22. The IEEE standard 754 single precision floating point representation of $(0.000000110110100101)_2$ is.
 (A) 0 10001111 11011010010100000000000
 (B) 0 01111001 11011010010100000000000
 (C) 0 10001110 1011010011010000000000000
 (D) 0 01111000 10110100101000000000000
23. How many clock cycles are required to perform two-operand operations using one bus datapath?
 (A) 1 (B) 2
 (C) 3 (D) Can't be determined
24. Which of the following is a rounding mode in IEEE754 standard?
 (i) round to 0
 (ii) round towards $+\infty$
 (iii) round towards $-\infty$
 (iv) round to nearest representable number
 (A) (i), (iv) only (B) (ii), (iii) only
 (C) (i), (ii) only (D) (i), (ii), (iii), (iv)
25. What is the normalized form of $0.00000110 \times 16^{101}$?
 (A) 1.10×16^{107} (B) 1.10×16^{95}
 (C) 0.110×16^{94} (D) 0.110×16^{106}
26. What is the biased representation of -7 , using 4-bits for the bias?
 (A) 0111 (B) 1111
 (C) 0000 (D) 1001
27. What is the total resultant after adding $A = -7$ and $B = -6$ using signed two's complement representation?
 (A) 0100 (B) 11101
 (C) 1101 (D) Overflow occurs
28. What is the total number of additions and subtractions required using Booths multiplication algorithm for the multiplier 00011110?
 (A) 1 (B) 2
 (C) 30 (D) Can't be determined
- Common data questions 29 and 30:** Consider a 12-bit floating point format in which base $b = 2$, a 5-bit exponent e with a bias = 16 and 6-bit normalized mantissa m . Given two floating point numbers:
 $A = 0\ 10001\ 011011$
 $B = 1\ 01111\ 101010$
29. After adding A and B , the resultant will be
 (A) 1 10001 000000
 (B) 1 10001 000001
 (C) 0 10001 000000
 (D) 0 10001 000001
30. After subtracting B from A , the resultant will be
 (A) 1 10001 110101 (B) 0 10001 110101
 (C) 1 10001 110110 (D) 0 10001 110110

ANSWERS KEYS

- | | | | | | | | | | |
|-------|-------|-------|-------|-------|-------|-------|-------|-------|-------|
| 1. B | 2. C | 3. C | 4. C | 5. D | 6. B | 7. C | 8. A | 9. C | 10. D |
| 11. C | 12. C | 13. B | 14. A | 15. D | 16. C | 17. A | 18. B | 19. B | 20. B |
| 21. D | 22. D | 23. B | 24. D | 25. B | 26. C | 27. D | 28. B | 29. C | 30. D |

Programming and Data Structures

PART A

Chapter 1: Programming in C	3.3
Chapter 2: Functions	3.14
Chapter 3: Arrays, Pointers, and Structures	3.30
Chapter 4: Linked Lists, Stacks, and Queues	3.47
Chapter 5: Trees	3.60

U

N

I

T

3

This page is intentionally left blank

Chapter 1

Programming in C

LEARNING OBJECTIVES

- Basic concepts
- Character set
- Identifier
- Declaring a variable
- Visualization of declaration
- Constants
- Single character constants
- String constants
- Using const keyword
- Precedence decreases as we move from top to bottom
- Type conversion
- Documentation section
- Preprocessing
- Global declaration
- Control statements
- Selection/Decision making statement
- Looping statements
- Unconditional jump statements

BASIC CONCEPTS

Character Set

A character refers to an alphabet, digit or a special symbol.
Alphabets: $A - Z, a - z$

Digits: 0 -9

Special symbols:

~ ! # % ^ and * () - + { } [] - < > , . | ? \ | ; ; " ' White space

Identifier

Identifier is a user-defined name used for naming a variable or a function.

Rules for naming an identifier

- Consists only letters, digits and underscore
- Starts only with an alphabet or underscore
- Keywords cannot be used.
- Can be as long as you like, first 31 characters are significant.

Example: Valid identifiers: RollNo, Roll_No, _Roll_No
rollno, Name2;
Invalid: 2name, Roll No.

Variable

The name itself represents value, is not constant. Variable is a data name whose value varies/changes during program execution. Variable name is a name given to memory cell (may be one or multiple bytes).

DATA TYPES

Represents type of data and set of operations to perform on data .

Data Type			
Primitive/Basic	Derived	User defined	Valueless
- Char	- Array	- Structure	
- float	- pointer	- union	- void
- double		Enumeration	
- integer			

Type	Keyword	Number of Bytes
Integer	int	2
Floating	float	4
Double	double	8
Character	char	1

Declaring a Variable

- Before using a variable, you must give some information to compiler about the variable. i.e., you must declare it.
- Declaration statement includes the type and variable name.

Syntax:

Datatype Var_name;

Example:

```
int roll_no;
char ch;
float age;
```

- When we declare a variable
 - memory space is allocated to hold a value of specified type.
 - space is associated with variable name
 - space is associated with a unique Address.

Table 1 Visualization of declaration

	roll no
int roll no;	garbage
	2002
	marks
int marks = 10;	10
	3008
	diameter
float diameter = 5.9	5.9
	4252
	ch → variable name
char ch : 'A'	A → value
	2820 → address

Note: The default value is garbage, i.e., an unknown value is assigned randomly.

Renaming data types with typedef Typedef is a keyword, which can form complex types from the basic type, and will assign some simpler names for such combinations. This is more helpful when some declaration is very tough, confusing or varies from one implementation to another.

For example, the data type unsigned long int is redefined as LONG as follows:

```
typedef unsigned long int LONG;
```

Uses of enumerated data types Enumerated data types are most useful when one is working over small, discrete set of values, in which each is having a meaning and it is not a number.

A best example can be given on months jan, feb, mar, ..., dec, which are 12 in number, with assigning consecutive numbers for it.

The main advantages are storage efficiency, the *c*-code can become readable

Constants

A constant value is one which does not change during the execution of a program.

C supports several types of constants:

1. Integer constants
2. Real constants
3. Single character constants
4. Strings constants

Integer constants

An integer constant is a sequence of digits. It consists of a set of digits 0 to 9 preceded by an optional + or – sign spaces, commas, and non-digit characters are not permitted between digits.

Examples for valid decimal integer constants are

```
123
-31
0
562321
+78
```

Examples for invalid integer constants are

```
20,000
₹1000
```

Real constants

Real constants consist of a fractional part in their representation. Integer constants are inadequate to represent quantities that vary continuously.

Examples of real constants are

```
0.0026
-0.97
435.29
+487.0
```

Single character constants

A single character constant represents a single character which is enclosed in a pair of quotation symbols.

Examples for character constants are

```
'5'
'x'
','
```

String constants

A string constant is a set of characters enclosed in double quotation marks. The characters in a string constant sequence may be alphabet, number, special character and blank space.

Examples of string constants are

```
“VISHAL”
“1234”
“C language”
“!....?”
```

Naming constants

A name given to a constant value. Value of name does not change during program execution.

Using const keyword

When we use 'const' with data type, memory will be allocated to variable and the initialized value does not change.

```
const int x = 10;
const float pi = 3.141;
```

Using # define

```
# define x 10
# define pi 3.141
```

Where '# define' is instruction to preprocessor so memory is allocated. The preprocessor replace each occurrence of name with value in program before execution.

OPERATOR

An operator is a symbol which performs operations on given data elements.

Table 2 Precedence and Associativity

() Parenthesis	L - R
[] Index	
→ Member of	
• Member of	
Pre ++, -- (unary) -, &(address of) * (Indirection)	R - L
Arithmetic *, /, %	L - R
Arithmetic: +, -	L - R
Bitwise shift: <<, >>	L - R
Relational: <, >, =, >, >=, !=	L - R
Bitwise ex -OR : ^	L - R
Logical AND : &&	L - R
Logical OR :	L - R
Conditional: ? :	R - L
Assignment & compound Assignment =, +=, -=, *=, /= ; %=	R - L
Separation operator: , (comma)	L - R

Note: For Assignment operator, only a variable is allowed on its left.

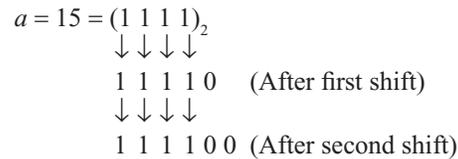
Precedence Decreases as We Move from Top to Bottom

Examples:

```
1. int a,b,c;
   a = b = c = 0;
   Assigns '0' to a,b,c;
```

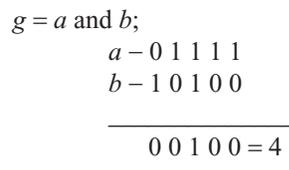
```
2. int a, b = 55, c = 10;
   initializes 'b' with 55 and 'c' with '10'.
   b + c = a; // Invalid
   Only variable is allowed on left side of assignment.
```

```
3. int a = 15, b = 20, c = 2, d = 5, e = 10, f, g, h, i;
   f = a << c;
   'a' is left shifted for 'c' times and result stored in 'f'
   i.e.,
```



One left shift multiplies 15 by 2 = 30
Again the 2nd left shift multiplies 30 by 2 = 60
Thus $15 \times 2^2 = 60$, where the power of 2 is the number of times shift is made. Value of 'f' becomes 60.

Note: Left shift multiplies the value by 2. Right shift divides the value by 2.

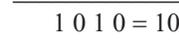


'&' performs bitwise AND. So 'g' value is '4'.
 $h = a \& b$; $a - 0\ 1\ 1\ 1\ 1$
 $b - 1\ 0\ 0\ 0\ 0$



'|' performs bitwise 'OR'. R value is '31'.

```
i = a ^ d: a - 1 1 1 1
          b - 0 1 0 1
```



'^' performs bit-wise ex - OR. i value is '10'.

```
4. int a = 100, b = 200, c = 300, x;
   x = (a > b) ? ((a > c) ? a : c) : ((b > c) ? b : c);
   false c
```

```
x = c
so, x = 300
```

```
5. int i = 10, j = 10, x, y;
   x = i++++i+i++++i++++i
   executes as
   ++ i }
   ++ i } pre-increments
   ++ i }
   X = i + i + i + i + i
   i ++ ; } post-increments
   i ++ ; }
   so x = 65, i = 15.
   y = j - - - + - - j + j - - + - - j +
   - - j
```

```

executes as
- - j; }
- - j; } pre decrements
- - j; }
y = j + j + j + j + j ;
j - -; }
j - -; } post decrements.
y = 35; j = 5
    
```

```

6. int i = 10;
printf("%d%d%d%d%d", i++, ++i, ++i,
i++, ++i);
evaluates the values in printf from
right to left.
So
i++, ++i, ++i, i++, ++i
←
Prints 14 14 13 11 11
Printf ("%d", i)
Prints 15:
    
```

TYPE CONVERSION

‘C’ allows mixed mode operations, i.e., variables of different type may appear in same expression. To perform the operation, the data need to convert into compatible type.

The conversion takes place in two ways:

Implicit

C automatically converts any intermediate values to proper type so that the expression can be evaluated without losing any significance.

For mixed mode operations, generally the ‘lower’ type is automatically converted to ‘higher’ type before the operation proceeds.

Explicit

‘C’ allows programmer to use type conversion operator to convert a data value to the required type.

Syntax:

```
V1 = (type) V2;
```

Type in parenthesis represents the destination type.

```
Example: int a = 3, b = 2, float x, y;
```

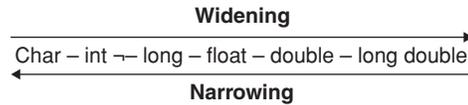
```
Case I: x = a/b;
results x = 1.000000
```

```
Case II: y = (float) a/b;
results y = 1.500000.
```

Because, in case 1, the integer division is performed and so returns an integer by division operator. While assigning the integer value implicitly converted to 1.000000, then assigns to float variable *x* where as in case 2, (float)a converts value of ‘a’ to float, the second variable ‘b’ is integer. The compiler implicitly converts integer to float. Then it performs float division. So 1.500000 is stored into floating variables.

Notes: ‘C’ allows both implicit and explicit type conversion. Type conversion is of two types:

1. **Narrowing:** Conversion of ‘higher’ type to ‘lower’ type.
2. **Widening:** Conversion of ‘lower’ type to ‘higher’ type.



Note: Narrowing causes loss of data.

Input/output Functions	
Function	Purpose
printf	prints formatted string
scanf	reads formatted string
getchar	reads character
putchar	displays a character
gets	reads a string
puts	displays a string

Format Specifier	Purpose
%c	single character
%d	decimal integer
%e	floating point
%f	floating point
%h	short int
%o	octal integer
%x	hexa decimal
%s	string
%u	unsigned decimal integer

Note: scanf(“%s”, string_var); does not read string which contains white space. Hence to read multi word string use gets(string_var);

Example 1: Which of following comment regarding the reading of a string using scanf() and gets () is true?

- (A) Both can be used interchangeably
- (B) scanf is delimited by end of line, gets is delimited by blank space
- (C) scanf is delimited by blank, gets is delimited by end of line
- (D) None of these

Ans: (C)

PROGRAM STRUCTURE

```

/* Documentation section */
Preprocessor commands;
Global declaration;
main ()
    
```

```
{
Body of main;
}
User defined function area;
```

Documentation section/comments Ignored by compiler, provides additional information to user to improve readability.

Preprocessing Tells the compiler to do pre-processing before doing compilation. For example

#include <stdio.h> tells to include stdio header file.

Global declaration It contains variable declarations, these are accessible in more than one function.

Function Functions are main building blocks of 'C' program. Every 'C' program contains one or more functions. A mandatory function called 'main()' instructs the compiler to start execution from here.

User defined area Here user can define his own functions.

CONTROL STATEMENTS

The statement that controls the execution sequence of a program is called "control statement".

The control statements are classified as:

1. Selection statement: if, switch
2. Iterative/looping statement: While, do-while, for
3. Unconditional jump statements: break, continue, return, goto

Selection/Decision-making Statement

Makes a decision to select and execute statement(s) based on the condition. 'C' supports if and switch selection statements.

The if statement "if" is called two-way selection statement.

Syntax:

```
if (expression) // simple-if
    statement(s);
if (expression) // if-else
{
    statement1(s);
}
else
{
    statements(s);
}
if (expression) // ladder else-if.
{
    Statement1(s);
}
else if (expression2)
```

```
{
Statement2(s);
}
else
{
Statement3(s);
}
Nested if:
    if (expression1)
{
    Statement(s)1;
    if (expression(s)2)
else
    Statement(s)3;
}
else
Statement(s)4;
```

Note: If the expression evaluates to true then the statements of if block gets executed otherwise else block statements will execute.

Example 2: Consider the following program segment:

```
if (a > b) printf ("a > b");
else
printf ("else part");
printf ("a < = b");
a < = b will be printed if
```

- (A) $a > b$ (B) $a < b$
(C) $a = b$ (D) all of these

Ans: (D)

Because the statement, printf("a < = b"); is not the part of either if block or else block.

The switch statement Switch is a multi-way (n-way) selection statement.

Syntax:

```
switch (var_name/exp)
{
case const1: stmts1;
            break;
case const2: stmts2;
            break;
            .
            .
            .
case constn: stmts n;
            break;
default: statements;
}
```

Notes:

- For switch only the integral (integer/char) type variables or expression evaluates to integral allowed.
- Absence of break after case statements leads continuation execution of all case statements followed by matching case block.

Example 3:

```
main ()
{
int i = 10;
switch(i)
{
case 10 : printf ("case 10");
case 15 : printf ("case 15");
case 20 : printf ("case 20");
default : printf ("default case");
}
}
```

Output: Case 10 case 15 case 20 default case

Reason: Missing break after each case, leads to execution of all the cases from matching case.

Example 4:

```
main ()
{
int i = 10;
switch (i)
{
case 10 : printf("case 10");
break ;
case 8 + 2 : printf("case 8+2");
break;
default : printf(" No matching case");
}
}
```

Program raises an error called ‘Duplicate case’ while compiling because the expression ‘8 + 2’ evaluates to ‘10’.

Looping Statements

Sometimes, there is a situation to execute statement(s) repeatedly for a number of times or until the condition satisfies. C supports following looping statements: while, do-while, for.

While Statement

Syntax: while (condition)

```
{
Statement(s);
}
```

If the condition is true the block of statements will execute and control returns to condition, i.e., the statement(s) executes till the condition becomes false.

Notes:

- ‘While’ executes the block either ‘0’ or more times.
- ‘While’ is called entry control loop.

Do-while Statement.

Syntax:

```
do
{
Statement(s);
} while (condition);
```

do-while is same as ‘while; except that the statement(s) will execute for at least once.

Notes:

- The condition will not be evaluate to execute the block for first time.
- ‘do-while’ is called exit-control loop.

Example 5:

```
main ( )
{
int i = 0;
while (i! = 0)
{
printf("%d", i);
i++;
}
}
```

No output, because the condition is false for the first time.

```
main ( )
{
int i = 0;
do
{
printf("%d", i);
i++;
} while (i! = 0);
}
```

Output: Displays 0 to 32767 and–32768 to–1

The for loop ‘for’ provides more concise loop control structure.

Syntax:

```
for(exp1; exp2; exp3)
{
Statement(s);
}
```

Expression 1: Initialization expression may contain multiple initializations. It executes only once before executing the loop for first time.

Expression 2: Condition expression. Only one condition expression is allowed. That may be single or compound condition, evaluates before every execution.

Expression 3: Modification statement may contain multiple statements. It executes on completion of loop body for every iteration.

Note: All the expressions in parenthesis are optional. Two semi-colons (;) are compulsory even though there are no expressions.

Odd loops In the for loop, while loop, the condition specifies the number of times a loop can be executed. Sometimes a user may not know, about the number of times a loop is to be executed. If we want to execute a loop for unknown number of times, then the concept of odd loops should be implemented, these can be done using the for, while (or) do-while loops. Let us illustrate odd-loop with a program

```
# include <stdio.h>
main()
{
int num, x;
num = 1;
while (num == 1)
{
printf ("enter a number");
scanf ("%d", & x);
if((x % 2) == 0)
printf("number is even");
else
printf("number is odd");
printf("do u want to test any num.");
printf("for yes-enter '1', No-enter '0'");
scanf("%d",& num);
}
}
```

Unconditional Jump Statements

- “C” language permits to jump from one statement to another.
- ‘C’ supports break, continue, return and goto jump statements.

Break statement Breaks the execution sequence. That is when the break statement executes in a block (loop) it’ll come out from block (loop).

Syntax:

```
break;
```

Continue statement Used to skip a part of the loop under certain conditions.

Syntax:

```
continue;
```

Return statement Terminates the execution of a function and returns the control to the calling function.

Syntax:

```
return [exp/value];
```

Goto statement Jumps from one point to another with in a function.

Syntax:

```
label1:                goto label2:
Statement(s);         Statement(s);
goto label1;          label2;
reverse jump          forward jump
```

Reverse jump, executes the statements repeatedly where as in forward jump, the statements are skipped from execution.

Example 6:

```
main( )
{
int i ;
for (i=1; i<=10; i++)
{
        if (i == 5)
            break;
        printf("%d" , i);
}
}
```

output: 1 2 3 4
if $i = 5$, then the loop will break.

Example 7:

```
main( )
{
        int i ;
        for (i = 1; i<=10; i++)
        {
                if (i == 5)
                    continue;
                printf("%d" , i);
        }
}
```

o/p: 1 2 3 4 6 7 8 9 10

if $i = 5$, the loop statements skipped for that iteration. So it does not print ‘5’.

Example 8:

Output for the following program segment

```
for (i = 1, j = 10 ; i < 6; ++i, --j)
printf("\n %d %d", i, j);
```

Output:

```
1 10
2 9
3 8
4 7
5 6
```

Note: Since for statement allows multiple initialization and multiple update statements, expression 1 and expression 3, does not raise any error.

EXERCISES

Practice Problems I

Directions for questions 1 to 15: Select the correct alternative from the given choices.

1. What will be the output of the following program?

```
void main()
{
    int i;
    char a[ ] = " \0 ";
    if (printf("%s\n", a))
        printf ("ok \n");
    else
        printf("program error \n");
}
```

- (A) ok (B) program error
(C) no output (D) compilation error

2. Output of the following will be

```
# define FALSE-1
# define TRUE 1
# define NULL 0
main( )
{
    if(NULL)
        puts("NULL");
    else if(FALSE)
        puts("TRUE");
    else
        puts("FALSE");
}
```

- (A) NULL (B) TRUE
(C) FALSE (D) 1

3. main()

```
{
    printf("%x",-1 << 4) ;
}
```

For the above program output will be

- (A) FFF0 (B) FF00
(C) 00FF (D) 0FFF

4. For the following program

```
# define sqr (a) a*a
main()
{
    int i;
    i = 64 / sqr(4);
    printf( "%d", i);
}
```

output will be

- (A) 4 (B) 16
(C) 64 (D) compilation error

5. #define clrscr () 1000

```
main ( )
{
    clrscr();
```

```
printf ( "%d \n", clrscr());
}
```

Output of the above program will be?

- (A) error (B) No output
(C) 1000 (D) 1

6. Output of the following program is

```
main( )
{
    int i = -2;
    +i;
    printf("i = %d, +i = %d\n", i, +i);
```

- (A) error (B) -2,+2
(C) -2,-2 (D) -2,2

7. main()

```
{
    int n;
    printf("%d", scanf ("%d", &n));
}
```

For the above program if input is given as 20. What will be the output?

- (A) 20 (B) 1
(C) 2 (D) 0

8. How many times will the following code be executed?

```
{
    x = 10;
    while (x = 1)
        x ++;
}
```

- (A) Never
(B) Once
(C) 15 times
(D) Infinite number of times

9. The following statement

```
printf("%d", 9%5); prints
```

- (A) 1.8 (B) 1.0
(C) 4 (D) 2

10. int a;

```
printf("%d", a);
```

What is the output of the above code fragment?

- (A) 0 (B) 2
(C) Garbage value (D) 3

11. printf("%d", printf("time"));

- (A) syntax error
(B) outputs time 4
(C) outputs garbage
(D) prints time and terminates abruptly

12. The following program

```
main( )
{
    int i = 2;
    {
        int i = 4, j = 5;
```

```
printf ("%d%d", i, j);
}
printf ("%d%d", i, j);
}
```

- (A) Compiler error: unrecognised symbol *j*;
 (B) Prints 2545
 (C) Print 4525
 (D) None of the above

13. What is the output of the following program fragment?

```
for (i = 3; i < 15; i += 3);
printf ("%d", i);
```

- (A) a syntax error (B) an execution error
 (C) prints 12 (D) prints 15

14. What is the output of the following program segment?

```
int a = 4, b = 6;
printf ("%d", a = b);
```

- (A) Outputs an error message
 (B) Prints 0
 (C) Prints 1
 (D) None of these

15. The statements:

```
a = 7;
printf ("%d", (a++));
prints
```

- (A) Value of 8 (B) Value of 7
 (C) Value of 0 (D) None of the above

Practice Problems 2

Directions for questions 1 to 12: Select the correct alternative from the given choices.

1. If the condition is missing in a FOR loop of a C program then

- (A) It is assumed to be present and taken to be false
 (B) It is assumed to be present and taken to be true
 (C) It results in syntax error
 (D) Execution will be terminated abruptly

2. Which of the following operators in 'C' does not associate from the right?

- (A) = (B) +=
 (C) postfix++ (D) >

3. In a C programming language $x -= y + 1$ means

- (A) $x = -x - y - 1$ (B) $x = x - y + 1$
 (C) $x = x - y - 1$ (D) $x = -x + y + 1$

4. Minimum number of temporary variables needed to swap two variables is

- (A) 1 (B) 2
 (C) 3 (D) 0

5. A preprocessor command

- (A) need not start on a new line
 (B) need not start on the first column
 (C) has # as the first character
 (D) comes after the first executable statement

6. `printf ("0d", printf ("0d", printf ("time4kids")));`

- (A) Outputs time (B) Syntax error
 (C) Outputs 9 (D) None of the above

7. `for (i = 1; i < 5; i++)`

```
if (i!=3)
printf ("%d", i);
```

Outputs:

- (A) 12345 (B) Error
 (C) 1245 (D) 0000

8. Which operand in 'C' takes only integer operands?

- (A) * (B) /
 (C) % (D) +

9. An unrestricted use of 'goto' statement is harmful because

- (A) it results in increasing the executing time of the program
 (B) it increases the memory of the program
 (C) it decreases the readability and testing of program
 (D) None of the above

10. What will be the output?

```
main ()
{
int i = 0, j = 0;
if (i && j++)
printf ("%d..%d", i++, j);
printf ("%d..%d", i, j);
}
```

- (A) 1..1 (B) 2..2
 (C) 0..0 (D) 1..1, 1..1

11. What is the output?

```
main ()
{
int a = 0;
int b = 20;
char x = 1;
char y = 10;
if (a, b, x, y);
printf ("hello");
}
```

- (A) logical error (B) Garbage value
 (C) hello (D) 20

12. What will be the value of count after executing the below program:

```
main ( ) {
int count = 10, digit = 0;
while (digit <= 9) {
printf ("%d\n", ++count);
++digit;
}
```

- (A) 10 (B) 11
 (C) 20 (D) 21

PREVIOUS YEARS' QUESTIONS

1. Which one of the following are essential features of an object-oriented programming language?
- Abstraction and encapsulation
 - Strictly-typedness
 - Type-safe property coupled with sub-type rule
 - Polymorphism in the presence of inheritance

[2005]

- (A) (i) and (ii) only
 (B) (i) and (iv) only
 (C) (i), (ii) and (iv) only
 (D) (i), (iii) and (iv) only

2. Which of the following are true?

- A programming language which does not permit global variables of any kind and has no nesting of procedures/functions, but permits recursion can be implemented with static storage allocation
- Multi-level access link (or display) arrangement is needed to arrange activation records only if the programming language being implemented has nesting of procedures/functions
- Recursion in programming languages cannot be implemented with dynamic storage allocation
- Nesting procedures/functions and recursion require a dynamic heap allocation scheme and cannot be implemented with a stack-based allocation scheme for activation records
- Programming languages which permit a function to return a function as its result cannot be implemented with a stack-based storage allocation scheme for activation records

[2008]

- (A) (ii) and (v) only (B) (i), (iii) and (iv) only
 (C) (i), (ii) and (v) only (D) (ii), (iii) and (v) only

3. What will be the output of the following C program segment?

```
char inChar = 'A';
switch(inChar) {
case 'A': printf("choice A\n");
case 'B':
case 'C': printf("choice B");
case 'D':
case 'E':
default: printf("No Choice");}
```

[2012]

- (A) No choice
 (B) Choice A
 (C) Choice A
 Choice B No choice
 (D) Program gives no output as it is erroneous

4. Suppose n and p are unsigned int variables in a C program. We wish to set p to nC_3 . If n is large, which one of the following statements is most likely to set p correctly?

[2014]

- (A) $p = n * (n - 1) * (n - 2) / 6;$
 (B) $p = n * (n - 1) / 2 * (n - 2) / 3;$
 (C) $p = n * (n - 1) / 3 * (n - 2) / 2;$
 (D) $p = n * (n - 1) * (n - 2) / 6.0;$

5. The secant method is used to find the root of an equation $f(x) = 0$. It is started from two distinct estimates x_a and x_b for the root. It is an iterative procedure involving linear interpolation to a root. The iteration stops if $f(x_b)$ is very small and then x_b is the solution. The procedure is given below. Observe that there is an expression which is missing and is marked by ?. Which is the suitable expression that is to put in place of ? so that it follows all steps of the secant method?

[2015]

Secant

Initialize: x_a, x_b, ϵ, N // ϵ = convergence indicator
 // N = maximum no. of

iterations

 $f_b = f(x_b)$ $i = 0$ while ($i < N$ and $|f_b| > \epsilon$) do $i = i + 1$ // update counter $x_t = ?$ // missing expression for
// intermediate value $x_a = x_b$ // reset x_a $x_b = x_t$ // reset x_b $f_b = f(x_b)$ // function value at new x_b

end while

if $|f_b| > \epsilon$ then // loop is terminated with $i = N$
write "Non-convergence"

else

write "return x_b ."

end if

(A) $x_b - (f_b - f(x_a)) f_b / (x_b - x_a)$ (B) $x_a - (f_a - f(x_a)) f_a / (x_b - x_a)$ (C) $x_b - (x_b - x_a) f_b / (f_b - f(x_a))$ (D) $x_a - (x_b - x_a) f_a / (f_b - f(x_a))$

6. Consider the following C program:

```
#include<stdio.h>
int main( )
{
int i, j, k = 0;
j = 2 * 3 / 4 + 2.0 / 5 + 8 / 5;
k -= --j;
for (i = 0; i < 5; i ++ )
{
switch(i + k)
```

```

{
case 1:
case 2: printf("\n%d", i + k);
case 3: printf("\n%d", i + k);
default: printf("\n%d", i + k);
}
}
return 0;
}

```

The number of times printf statement is executed is _____.

[2015]

7. Consider the C program fragment below which is meant to divide x by y using repeated subtractions. The variables x, y, q and r are all unsigned int.

```

while (r >= y) {
    r = r - y;
    q = q + 1;
}

```

Which of the following conditions on the variables x, y, q and r before the execution of the fragment will

ensure that the loop terminates in a state satisfying the condition $x == (y * q + r)$? [2017]

- (A) $(q == r) \ \&\& \ (r == 0)$
 (B) $(x > 0) \ \&\& \ (r == x) \ \&\& \ (y > 0)$
 (C) $(q == 0) \ \&\& \ (r == x) \ \&\& \ (y > 0)$
 (D) $(q == 0) \ \&\& \ (y > 0)$

8. Consider the following C Program.

```

#include<stdio.h>
int main () {
    int m = 10;
    int n, n1 ;
    n = ++m;
    n1 = m++;
    n--;
    --n1;
    n -= n1;
    printf ("%d", n) ;
    return 0;
}

```

The output of the program is _____.

[2017]

ANSWER KEYS

EXERCISES

Practice Problems 1

1. A 2. B 3. A 4. C 5. C 6. C 7. B 8. D 9. C 10. C
 11. B 12. A 13. D 14. D 15. B

Practice Problems 2

1. B 2. D 3. C 4. D 5. C 6. D 7. C 8. C 9. C 10. C
 11. C 12. C

Previous Years' Questions

1. B 2. D 3. C 4. B 5. C 6. 10 7. C 8. 0

Functions

LEARNING OBJECTIVES

- Functions
- Library functions
- User defined functions
- Defining user defined functions
- Recursion
- Parameter passing
- Pass by value
- Pass by address
- Scope
- Life time
- Binding

FUNCTIONS

A function is a block of code that performs a specific task. It has a name and is reusable, i.e., it can be executed from as many different parts in a program as required.

Functions make possible top down modular programming. In this style of programming, the high-level logic of overall problem is solved first, whereas the detail of each lower-level function is addressed later. This approach reduces the complexity in writing program.

1. Every C program can be thought of collection of functions.
2. `main()` is also a function.

Types of Functions

Library functions

These are the in-built functions of 'C' library. These are already defined in header files.

Example 1: `printf()`; is a function which is used to print at output. It is defined in 'stdio.h' file.

User-defined functions

Programmers can create their own function in 'C' to perform specific tasks.

Example 2:

```
# include <stdio.h>
main( )
{
message( );
}
message( )
{
printf("Hello");
}
```

- A function receives zero (or) more parameters, performs a specific task, and returns zero or one value.
- A function is invoked by its name and parameters.
- No two functions have the same name in a single C program.
- The communication between the function and invoker is through the parameter and the return value.
- A function is independent.
- It is "completely" self-contained.
- It can be called at any place of your code and can be ported to another program.
- Functions make programs reusable and readable.

Example 3: Return the largest of two integers.

```
int maximum (int a, int b)
{
if (a > b)
return a;
else
return b;
}
```

Note: Function calls execute with the help of execution stack. Execution of 'C program' starts with `main()` function. `Main()` is a user-defined function.

Defining User-defined Functions

In order to work with user-defined functions, it requires the following concepts about functions:

- Declaration of a function
- Definition of a function
- Function call

Declaration specifies what

- is the name of the function
- are the parameters to pass (type, order and number of parameters).
- it returns on completion of execution

Example 4: `int maximum (int, int); int maximum (int a, int b);`

Syntax:

```
Return_type Function_Name(Parameter_list);
```

- Names of parameters are optional in declaration.
- Default return type for 'C' functions is 'int'.
- A function, whose return type is `void` returns nothing.
- Empty parenthesis after a function name in declaration says, function does not accept any parameters.

Definition specifies *how*

- to perform the specified task
- to accept passed parameters
- to process parameters or execute instruction to produce required results (return value).

Function definition is a self-contained block of instructions, will be executed on call:

Syntax:

```
Return_type Function -Name(paralist)
{
  Local declaration(s);
  Executable statements(s);
}
int maximum (int a, int b)
{
  if (a > b)
    return a;
  else
    return b;
}
```

Function call specifies

1. where to execute the function
2. when to execute the function

Note: If the function definition provided before use (call), the declaration is optional.

The following example describes control flow during function call:

```
void hello(); // Declaration
void main()
{
  printf("\n function");
  hello();
  printf("\n Main after call to hello")
void hello()//Definition
{
```

```
printf ("\n function Hello");
return;
}
```

A *return* statement has two important uses:

1. first, it causes an immediate exit from the function.
2. second, it may be used to return a value.

If a function does not return any value, then the return statement is optional.

RECURSION

In general, programmers use two approaches to write repetitive algorithms. One approach using *loops*, the other is *recursion*.

Recursive functions typically implement recurrence relations, which are mathematical formula in which the desired expression (function) involving a positive integer, *n*, is described in terms of the function applied to corresponding values for integers less than 'n'.

1. The function written in terms of itself is a recursive case.
2. The recursive case must call the function with a decreasing 'n'.
3. Recursion in computer programming is exemplified when a function defined in terms of itself.
4. Recursion is a repetitive process in which a function calls itself.

Note: Recursive function must have an *if* condition to force the function to return without recursive call being executed. If there is no such condition, then the function execution falls into infinite loop.

Rules for designing recursive function

1. Determine base case
2. Determine general case
3. Combine, base and general case into a function

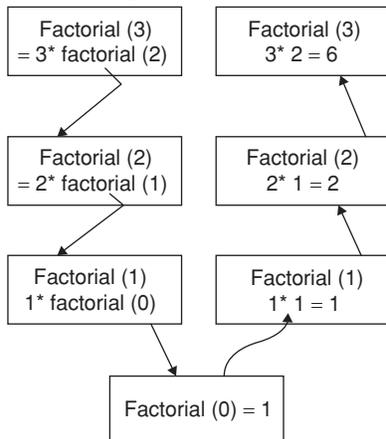
Example 5: Recursive factorial function

```
1. int factorial (int n)
2. {
3. if (n == 0)
4. return 1;
5. else
6. return (n * factorial (n - 1));
7. }
```

The statement 3 is a base condition, which stops the recursive call of function.

The statement 6 reduces the size of problem by recursively calling the factorial with $(n - 1)$.

Execution sequences for factorial (3):



Disadvantages:

1. Recursive programs increase the execution time of program.
2. Recursive programs typically use a large amount of computer memory and the greater the recursion, the more memory is used.
3. Recursive programs can be confusing to develop and extremely complicated to debug.

PARAMETER PASSING

There are two ways of passing parameters to functions in ‘C’ language.

1. Pass-by-value: When parameters are passed by value, create copies in called function. This mechanism is used when we do not want to change the value of actual parameters.
2. Pass-by-address: In this case, only the addresses of parameters are passed to the called function. Therefore, manipulation of formal parameters affects actual parameters.

Examples 6:

```
void swap1(int, int); /* function - to swap two numbers by passing values */
```

```
void swap2 (int *, int *); /* function to swap two numbers by passing Address * /
void main ()
{
int a = 10, b = 15, c = 5, d = 25;
printf("value of a and b before swapping :%d, %d" a , b );
swap1(a, b);
printf("values of a and b after swapping : %d, %d", a, b);
printf ("values of c and d before swapping :%d%d", c,d );
Swap2(&c, &d);
printf("values of c and d after swapping %d, %d", c, d);
}
void swap1(int x, int y )
{
int temp;
temp = x;
x = y;
y = temp;
}
void swap2 (int *x, int *y)
{
int temp;
temp = *x;
*x = *y;
*y = temp;
}
```

Output:

Value of *a* and *b* before swapping: 10, 15
 Value of *a* and *b* after swapping: 10, 15
 Value of *c* and *d* before swapping: 5, 25
 Value of *c* and *d* after swapping: 25, 5

Solved Examples

Example 1: Consider the program below:

```
#include<stdio.h>
int fun (int n, int *fp)
```

Table 1 Comparison of pass-by-value and pass-by-address

Pass-by-value	Pass-by-address
1. Also known as call-by-value	1. Also known as call-by-address or call by-reference
2. Pass the values of actual parameters	2. Pass the address of actual parameters
3. Formal parameters act as duplicates or as a copy to actual parameters	3. Formal parameters acts as references to the actual parameters
4. Operations on formal parameter does not affect actual parameters	4. Operations on formal parameters affect actual parameters
5. Passing of parameters is time consuming as the data size increases	5. The size of parameters does not affect the time for transferring references.
6. Actual parameters are secured	6. Helps to return multiple parameters

```

{
int t, f;
if (n <= 1)
{
*fp=1;
return 1;
}
t = fun(n-1, fp);
f = t+ *fp;
*fp = t;
return f;
}
int main ()
{
int x = 15;
printf ("%d\n", fun(5, &x));
return 0;
}
    
```

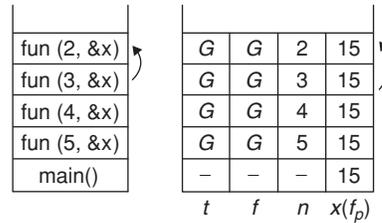
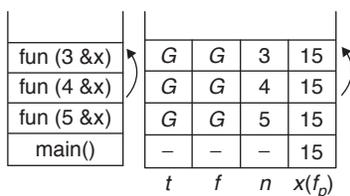
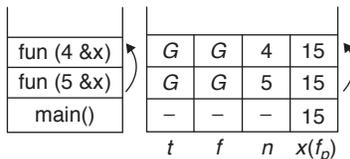
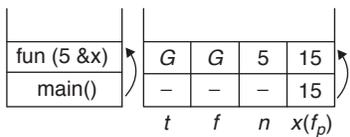
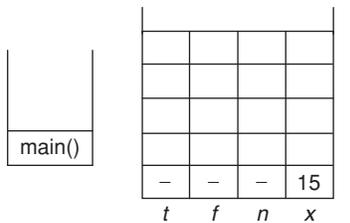
What is the output?

- (A) 2
- (B) 4
- (C) 8
- (D) 16

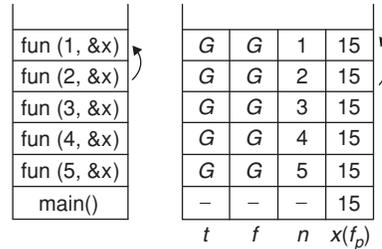
Solution: (C)

Execution stack

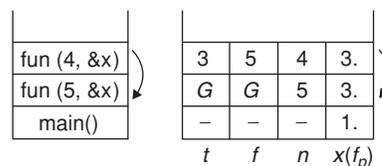
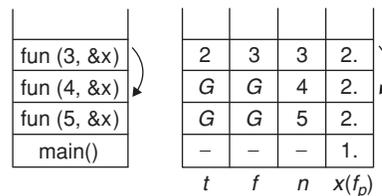
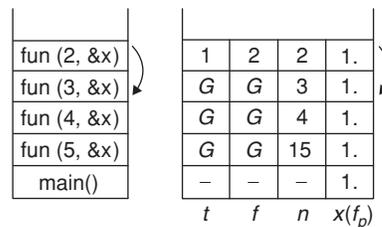
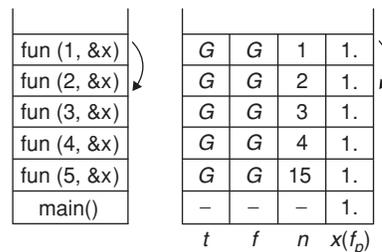
Function call sequence Corresponding values of *t*, *f*, *n*, *x*

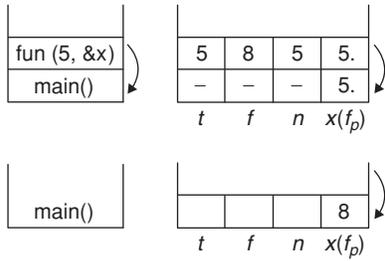


Note: ‘-’ indicates no memory allocated to variable. ‘G’ indicates garbage value.



For the function call fun(1, &x) condition ($n \leq 1$) is true. So Assigns ‘1’ to fp and returns ‘1’.





Finally, x contains '8', so printf prints '8'.

Example 2: What does the following program prints?

```
#include <stdio.h>
void f (int *p, int *q)
{
    p=q;
    *p=12;
}
int i = 0, j=1;
int main()
{
    f(&i, &j);
    printf(" %d%d ", i, j);
    return 0 ;
}
```

- (A) 2 12
- (B) 12 1
- (C) 0 1
- (D) 0 12

Solution: (D)

```
main( )
f (&i, &j)
address of 'i' is stored in to p.
and address of 'j' is stored into 'q'.
i.e., *p and*q refers i and j.
The statement:
p = q; updates pointer 'p', so that both
pointers refer to parameter 'j'.
*p = 12
Changes value of 'j' to '12' But 'i' does
not effected. So, prints 0 12.
```

Example 3: What is the value printed by the following program?

```
# include <stdio.h>
int f(int *a, int n)
{
    if (n<=0) return 0 ;
    else if(*a%2 == 0)
    return *a + f(a+1, n-1);
    else
    return *a - f(a+1, n-1);
}
int main ( )
{
    int G [ ] = { 12, 7, 13, 4, 11, 6};
    printf("%d", f (a,b));
    return 0;
}
```

- (a) -9
- (b) 12
- (c) 15
- (d) 20

Solution: (C)

	0	1	2	3	4	5
a	12	7	13	4	11	6

$f(a, 6)$ is the first call to function $f()$.

The array `_name` refers to base address of array, i.e., address of first element.

Thus,

$F(a, 6)$

$12 \% 2 = 0$. So,

$$12 + f(a+1, \frac{n-1}{5}) [*a \text{ is even}]$$

↓ 7

$$12 + \left(7 - f(a+1, \frac{n-1}{4}) [*a \text{ is odd}] \right)$$

↓ 13

$$12 + \left(7 - \left(13 - f(a+1, \frac{n-1}{3}) [*a \text{ is odd}] \right) \right)$$

↓ 4

$$12 + \left(7 - \left(13 - \left(4 + f(a+1, \frac{n-1}{2}) [*a \text{ is even}] \right) \right) \right)$$

↓ 11

$$12 + \left(7 - \left(13 - \left(4 + \left(11 - f(a+1, \frac{n-1}{1}) [*a \text{ is odd}] \right) \right) \right) \right)$$

↓ 6

$$12 + \left(7 - \left(13 - \left(4 + \left(11 - \left(6 + \frac{f(a+1, \frac{n-1}{0}}{0}) [*a \text{ is even}] \right) \right) \right) \right) \right)$$

↓ 0

$$12 + (7 - (13 - (4 + (11 - (6 + 0)))))) = 15$$

SCOPE, LIFETIME AND BINDING

Storage classes specify the scope, lifetime and binding of variables. To fully define a variable, one needs to mention not only its 'type' but also its 'storage class'.

A variable name identifies some physical location within computer memory where a collection of bits are allocated for storing value of variable.

Storage class tells us:

1. Where the variable would be stored (either in memory or CPU registers)?
2. What will be the initial value of a variable, if no value is specifically initialized?
3. What is the scope of a variable (where it can be accessed)?
4. What is the life of a variable?

Scope

The scope defines the visibility of an object. It defines where an object can be referenced/accessed; generally, the scope of variable is local or global.

1. The variables defined within a block have *local scope*. They are *visible only to the block* in which they are defined.
2. The variables defined in global area are visible from their definition until the end of program. It is *visible everywhere* in program.

Lifetime

The lifetime of a variable defines the duration for which the computer allocates memory for it (the duration between allocation and deallocation of memory).

In *C*, variable can have automatic, static or dynamic lifetime.

1. Automatic: Variables with automatic lifetime are created each time their declaration are encountered and are destroyed each time their blocks are exited.
2. Static: A variable is created when the declaration is executed for the first time and destroyed when the execution stops/terminates.
3. Dynamic: The variable's memory is allocated and deallocated through memory management functions.

Binding

Binding finds the corresponding binding occurrence (declaration/definition) for an applied occurrence (usage) of an identifier. For Binding.

1. Scope of variables should be known. What is the block structure? In which block the identifier is variable?
2. What will happen if we use same identifier name again? '*C* forbids use of same identifier name in the same scope'. Same name can be used in different scopes.

Examples:

1.

```
double f,y;
int f( ) // error
{
.
.
.
}
double y; // error
```
2.

```
double y;
int f( )
{
double f;// legal
int y; //legal
}
```

There are four storage classes in *C*.

Storage class	Storage Area	Default Initial Value	Lifetime	Scope	Keyword
Automatic	Memory	Till the control remains in block	Till the control remains in block	Local	auto
Register	CPU register	An unpredictable value (or) garbage value	Till the control remains in block	Local	register
Static	Memory	Zero	Value of variable persist between function calls	Local	static
External	Memory	Unpredictable or garbage value	Throughout program execution	Global	extern

Note: Default storage class is auto.

Example 4: What will be the output for the program?

```
int i = 33;
main( )
{
    extern int i;
    {
        int i = 22;
        {
            const volatile unsigned i
            = 11;
            printf (" %d ", i);
        }
        printf (" %d ", i);
    }
    printf ("%d ", i) ;
}
```

- (A) error
 (B) 11 22 33
 (C) 11 22 garbage
 (D) 11 11 11

Solution: (B)

'{' introduces new block and thus new scope. In the innermost block, *i* is declared as const volatile unsigned which is a valid declaration. *i* is assumed of type int. So printf prints 11. In the next block, *i* has value 22 and so printf prints 22. In the outermost block, *i* is declared as extern, so no storage space is allocated for it. After compilation is over, the linker resolves it to global variable, *i* since it is the only variable visible there. So it prints its value as 33.

3.20 | Unit 3 • Programming and Data Structures

Example 5: Consider the following C program:

```
int f(int n)
{
    static int r;
    if (n<=0) return 1;
    if (n> 3)
    {
        r=n;
        return (f(n-2)+2));
    }
    return f(n-1) + r;
}
```

What is the value of $f(5)$?

- (a) 15
- (b) 17
- (c) 18
- (d) 19

Solution: (C)

Call Sequence	r	Return Sequence
$f(5)$	5	18
$f(3)+2$	5	16+2
$f(2)+r$	5	11+5
$f(1)+r$	5	6+5
$f(0)+r$	5	1+5

Common data for questions 6 and 7: Consider the following recursive 'C' function that takes two arguments.

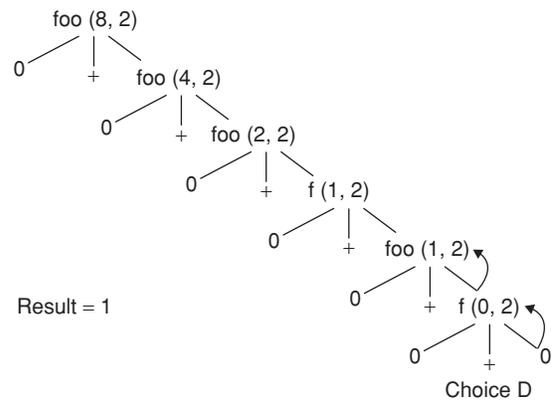
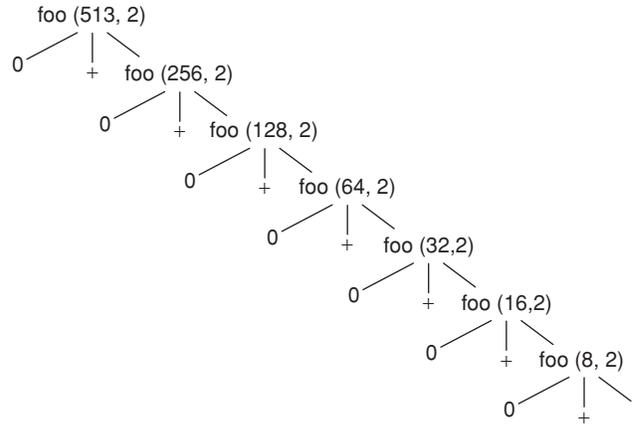
unsigned int foo (unsigned int n , unsigned int r)

```
{
if (n>0)
return ((n%r)+ foo(n/r,r));
else
return 0;
}
```

Example 6: What is the return value of the function foo when it is called as $foo(512,2)$?

- (A) 9
- (B) 8
- (C) 2
- (D) 1

Solution: (D)



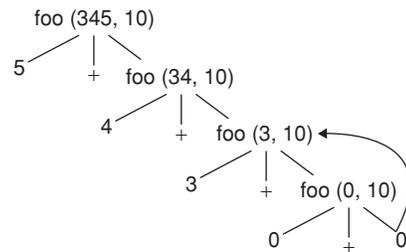
Result = 1

Result = 1

Example 7: What is return value for the function call $foo(345, 10)$?

- (A) 345
- (B) 12
- (C) 5
- (D) 3

Solution: (B)



result $5 + 4 + 3 = 12$

EXERCISES

Practice Problems I

Directions for questions 1 to 15: Select the correct alternative from the given choices.

1. What will be the output of the following program?

```
main( )
{
    main( );
}
```

- (A) overflow error (B) syntax error
(C) returns 0 (D) returns 1

2. Output of the following program is

```
main( )
{
    static int var = 6;
    printf("%d\t", var--);
    if(var)
        main( );
}
```

- (A) 5 4 3 2 1 0 (B) 6 6 6 6 6 6
(C) 6 5 4 3 2 1 (D) Error

3. Which of the following will be the output of the program?

```
main( )
{
    char str[ ] = "Hello";
    display( str );
}
void display (char *str)
{
    printf ( "%s", str );
}
```

- (A) compilation error (B) hello
(C) print null string (D) no output

4. Consider the following C function

```
int fun (int n)
{
    static int x = 0;
    if (n<=0) return 1;
    if (n>3)
    {
        x = n;
        return fun(n-2)+3;
    }
    return fun(n-1)+ x;
}
```

What is the value of fun(5)?

- (A) 4 (B) 15
(C) 18 (D) 19

5. For the following C function

```
void swap (int a, int b)
{
    int t;
```

```
    t = a;
    a = b;
    b = t;
```

```
}
```

In order to exchange the values of two variables *w* and *z*,

- (A) call swap (*w*, *z*)
(B) call swap (and *w*, and *z*)
(C) swap (*w*, *z*) cannot be used as it does not return any value
(D) swap (*w*, *z*) cannot be used as the parameters are passed by value

6. Choose the correct option to fill? *x* and? *y* so that the program below prints an input string in reverse order. Assume that the input string is terminated by a new line character:

```
void Rev(void) {
    int a;
    if (?x) Rev( );
    ?y
}
```

```
main( ) {
    printf("Enter the text");
    printf(" \n");
    Rev( );
    printf("\n");
}
```

- (A) ? *x* is (getchar() != '\n')
 ? *y* is getchar (A);
(B) ? *x* is ((A = getchar()) != '\n')
 ? *y* is getchar(A) ;
(C) ? *x* is (A! = '\n')
 ? *y* is putchar (A);
(D) ? *x* is (A = getchar ()) != '\n')
 ? *y* is putchar(A) ;

7. main ()

```
{
    extern int a;
    a = 30;
    printf ("%d", a);
}
```

What will be the output of the above program?

- (A) 30 (B) Compiler error
(C) Runtime error (D) Linker error

8. Which of the following will be the output of the program?

```
void main ( )
{
    int n = ret(sizeof(float));
    printf("\n value is %d ", ++n);
}
int ret(int ret)
{
```

3.22 | Unit 3 • Programming and Data Structures

```
ret += 2.5;
return (ret);
}
```

- (A) Value is 6 (B) Value is 6.5
(C) Value is 7 (D) Value is 7.5

9. The following program

```
main( )
{
    pt( ); pt( );pt( );
}
pt( )
{
    static int a;
    printf("%d", ++a) ;
}
prints
```

- (A) 0 1 2
(B) 1 2 3
(C) 3 consecutive, but unpredictable numbers
(D) 1 1 1

10. What is the output of the following program?

```
main( ) {
    int i = 0;
    while (i < 4) {
        sum(i);
        i++;
    }
}
void sum(int i) {
    static int k;
    printf ("%d", k + i);
    k++;
}
```

- (A) 0 2 4 6 (B) 0 1 2 3
(C) 0 2 0 0 (D) 1 3 5 7

11. What will be the output of following code?

```
# include <stdio.h>
aaa() {
    printf("hi");
}
bbb() {
    printf("hello");
}
ccc()
{
    printf("bye");
}
main ( )
{
    int *ptr[3]( );
    ptr[0] = aaa;
    ptr[1] = bbb;
    ptr[2] = ccc;
```

```
ptr[2]();
}
```

- (A) hi (B) hello
(C) bye (D) Garbage value

12. What is the output?

```
void main()
{
    static int i = 5;
    if(--i)
    {
        main();
        printf("%d", i);
    }
}
```

- (A) 5 (B) 5 5 5 5
(C) 0 0 0 0 (D) 1 1 1 1

13. If the following function gets compiled, what error would be raised?

```
double fun(int x, double y)
{
    int x;
    x = 100;
    return y;
}
```

- (A) Function should be defined as int fun(int x, double y)
(B) Missing parenthesis in return
(C) Redclaration of x
(D) All of these

14. Consider the following function:

```
fun(int x)
{
    if ((x/2) != 0)
    return (fun (x/2) 10 + x%2);
    else return 1;
}
```

What will happen if the function 'fun' called with value 16 i.e., as fun(16).

- (A) Infinite loop
(B) Random value will be returned
(C) 11111
(D) 10000

15. What is the output of the following program?

```
void main( )
{
    static int x = 5;
    printf("%d", x - - );
    if (x != 0)
    main( );
}
```

- (A) error:main() cannot be called from main()
(B) Infinite loop
(C) 5 4 3 2 1
(D) 0

Practice Problems 2

Directions for questions 1 to 15: Select the correct alternative from the given choices.

- An external variable
 - is globally accessible by all functions
 - has a declaration “extern” associated with it when declared within a function
 - will be initialized to 0, if not initialized
 - all of the above
- The order in which actual arguments are evaluated in a function call
 - is from the left
 - is from the right
 - is unpredictable
 - none of the above
- In C language, it is necessary to declare the type of a function in the calling program if the function
 - returns an integer
 - Returns a float
 - both (A) and (B)
 - none of the above

- What is the output?

```
void main()
{
    int k = ret(sizeof(int));
    printf("%d", ++k);
}
int ret (int ret)
{
    ret + = 2.5;
    return (ret);
}
```

- 3.5
 - 5
 - 4
 - logical error
- When a recursive function is called, all its automatic variables are
 - maintained in stack
 - retained from last execution
 - initialized during each call of function
 - none of these

- Consider the following program segment:

```
int fun(int x, int y)
{
    if(x > 0)
        return ((x % y) + fun(x/y, y));
    else
        return 0;
}
```

What will be the output of the program segment if the function is called as fun(525, 25)?

- 25
 - 12
 - 21
 - 42
- Consider the following C program segment:

```
int fun (int x)
{
    static int i = 0;
    if (x < = 0)
```

```
return 1;
else if (x > 5)
{
    i = x;
    return fun (x - 3) +2;
}
return fun (x - 2) + i;
}
```

What is the value of fun(7)?

- 17
- 10
- 11
- 9

- Consider the following C program:

```
void rearrange( )
{
    char ch;
    if (X)
        rearrange( );
    Y;
}
void main ( )
{
    printf("\n enter text to print reverse
order :");
    rearrange( ) ;
}
```

Choose the correct option to fill X and Y, so that the program prints the entered text in reverse order. Assume that input string terminates with new line.

- X: (getchar(ch) == '\n')
- Y: putchar(ch);
- X: (getchar(ch) != '\n')
- Y: ch = putchar();
- X: ((ch = getchar()) != '\n')
- Y: putchar(ch);
- X: ((ch = getchar()) == '\n')
- Y: putchar (ch);

- Consider the following C function:

```
int f(int n)
{
    static int i = 1;
    if (n > = 5) return n;
    n = n + i;
    i++ ;
    return f(n);
}
```

The value returned by f(1) is

- 5
- 6
- 7
- 8

- Consider the following C function:

```
int incr (int i)
{
    static int count = 0;
    count = count + i;
    return (count);
}
```

```

}
main ( )
{
int i, j;
for (i = 0; i < =4; i++)
j = incr (i);
}

```

The *j* value will be

- (A) 10
- (B) 4
- (C) 6
- (D) 7

11. The following function

```

int Trial (int a, int b, int c)
{
if ((a > = b) && (c < b))
return b;
else if(a > = b)
return Trail(a, c, b);
else return Trail (b, a, c);
}

```

- (A) finds the maximum of *a, b, c*
- (B) finds the middle value of *a, b, c* after sorting
- (C) finds the minimum of *a, b, c*
- (D) none of the above

12. Consider the following pseudo code

```

f(a, b)
{
while(b! = 0)
{
t = b;
b = a % b;
a = t;
}
return a;
}

```

- (A) The above code computes HCF of two numbers *a* and *b*
- (B) The above code computes LCM of *a* and *b*
- (C) The above code computes GCD of *a* and *b*
- (D) None of the above

13. 1. main ()
2. {int a = 10, *j;
3. void *k;
4. j = k = &a;
5. j++;
6. k++;
7. printf("\n %u, %u", j, k);
8. }

Which of the following is true in reference to the above code?

- (A) The above code will compile successfully
- (B) Error on line number 6
- (C) Error on line number 3
- (D) Error on line number 4

14. Aliasing in the context of programming language refers to

- (A) multiple variables having the same memory location
- (B) multiple variables having the same value
- (C) multiple variables having the same identifier
- (D) multiple uses of the same variable

15. Match the following:

X: <i>m</i> = malloc (5); <i>m</i> = NULL;	1: Using dangling pointers
Y: free (<i>n</i>); <i>n</i> value = 5;	2: Using un initialized pointers
Z: char * <i>p</i> ; * <i>p</i> = 'a';	3: Lost memory

- (A) X – 1 Y – 3 Z – 2
- (B) X – 3 Y – 1 Z – 2
- (C) X – 3 Y – 2 Z – 1
- (D) X – 2 Y – 1 Z – 3

PREVIOUS YEARS' QUESTIONS

1. In the following C function, let $n \geq m$.

```
int gcd(n,m)
{
if (n%m ==0) return m;
n = n%m;
return gcd(m,n);
}
```

How many recursive calls are made by this function?

[2007]

- (A) $\Theta(\log_2 n)$ (B) $\Omega(n)$
 (C) $\Theta(\log_2 \log_2 n)$ (D) $\Theta(\sqrt{n})$
2. What is the time *complexity* of the following recursive function?

```
int DoSomething (int n) {
if (n <= 2)
return 1;

else
return(DoSomething(floor(sqrt(n)))+ n);}

```

[2007]

- (A) $\Theta(n^2)$ (B) $\Theta(n \log_2 n)$
 (C) $\Theta(\log_2 n)$ (D) $\Theta(\log_2 \log_2 n)$
3. Choose the correct option to fill ? 1 and ? 2 so that the program below prints an input string in reverse order. Assume that the input string is terminated by a newline character.

```
void reverse (void) {
int c;
if (?1) reverse( );
?2
}

main ( ) {
printf ("Enter Text") ; printf ("\ n");
reverse ( ); printf ("\ n") ;
}

```

[2008]

- (A) ?1 is (getchar() != '\n')
 ?2 is getchar(c);
 (B) ?1 is (c = getchar()) != '\n'
 ?2 is getchar(c);
 (C) ?1 is (c != '\n')
 ?2 is putchar(c);
 (D) ?1 is ((c = getchar()) != '\n')
 ?2 is putchar(c);
4. Consider the program below:

```
# include < stdio.h >
int fun(int n, int * f_p) {
int t, f;
if (n <=1) {
```

```
*f_p =1;
return 1;
}
t = fun (n-1, f_p);
f = t+*f_p;
*f_p = t;
return f;
}

int main( ) {
int x = 15;
printf ("%d\n", fun(5,&x));
return 0;
}
```

The value printed is

[2009]

- (A) 6 (B) 8
 (C) 14 (D) 15

5. What is the value printed by the following C program?

```
#include <stdio.h>
int f(int *a, int n)
{
if (n <= 0) return 0;
else if(*a % 2 == 0) return * a + f(a+1, n-1);
else return *a-f(a+1, n-1);
}

int main( )
{
int a[ ] = {12, 7, 13, 4, 11, 6};
printf("%d", f(a,6));
return 0;
}

```

[2010]

- (A) -9 (B) 5
 (C) 15 (D) 19

Common data for questions 6 and 7: Consider the following recursive C function that takes two arguments.

```
unsigned int foo (unsigned int n, unsigned int r)
{
if ( n > 0 ) return ((n % r) + foo (n /r, r));
else return 0;
}
```

6. What is the return value of the function foo when it is called as foo (513, 2)?

[2011]

- (A) 9 (B) 8
 (C) 5 (D) 2

7. What is the return value of the function foo when it is called as foo (345, 10)? [2011]
 (A) 345 (B) 12
 (C) 5 (D) 3

Common data for questions 8 and 9: Consider the following C code segment

```
int a, b, c = 0;
void prtFun (void);
main ( )
{ static int a = 1;
  prtFun ( );
  a+ = 1;
  prtFun ( );
  printf("\n %d %d", a, b);
}
void prtFun(void)
{static int a = 2;
  int b = 1;
  a+ = ++b;
  printf("\n %d %d", a, b);
}
```

8. What output will be generated by the given code segment if: [2012]

Line 1 is replaced by **auto int a = 1;**

Line 2 is replaced by **register int a = 2;**

(A)	(B)	(C)	(D)
3 1	4 2	4 2	4 2
4 1	6 1	6 2	4 2
4 2	6 1	2 0	2 0

9. What output will be generated by the given code segment? [2012]

(A)	(B)	(C)	(D)
3 1	4 2	4 2	3 1
4 1	6 1	6 2	5 2
4 2	6 1	2 0	5 2

10. What is the return value of $f(p, p)$, if the value of p is initialized to 5 before the call? Note that the first parameter is passed by reference, whereas the second parameter is passed by value.

```
int f(int &x, int c) {
  c = c - 1;
  if (c == 0) return 1;
  x = x + 1;
  return f(x, c) * x;
}
```

- [2013]
 (A) 3024 (B) 6561
 (C) 55440 (D) 161051

11. Consider the following pseudo code. What is the total number of multiplications to be performed? [2014]

```
D = 2
for i = 1 to n do
  for j = i to n do
    for k = j + 1 to n do
      D = D * 3
```

- (A) Half of the product of the three consecutive integers.
 (B) One-third of the product of the three consecutive integers.
 (C) One-sixth of the product of the three consecutive integers.
 (D) None of the above.
12. Consider the function func shown below:

```
int func (int num) {
  int count = 0;
  while (num) {
    count ++;
    num>>=1;
  }
  return (count);
}
```

The value returned by func(435) is _____ [2014]

13. Consider the following function

```
double f (double X)
if (abs(X*X - 3) < 0.01) return X;
else return f(X/2 + 1.5/X);
}
```

Give a value q (to two decimals) such that $f(q)$ will return q : _____ [2014]

14. Consider the following pseudo code, where x and y are positive integers [2015]

```
begin
  q := 0
  r := x
  while r ≥ y do
    begin
      r := r - y
      q := q + 1
    end
  end
```

The post condition that needs to be satisfied after the program terminates is

- (A) $\{r = qx + y \wedge r < y\}$
 (B) $\{x = qy + r \wedge r < y\}$
 (C) $\{y = qx + r \wedge 0 < r < y\}$
 (D) $\{q + 1 < r - y \wedge y > 0\}$

15. Consider the following C function [2015]

```
int fun(int n) {
    int x = 1, k;
    if (n == 1) return x;
    for (k = 1; k < n; ++k)
        x = x + fun(k) * fun(n
- k);
    return x;
}
```

The return value of fun(5) is _____

16. Consider the following recursive C function

```
void get (int n)
{
    if (n < 1) return;
    get (n - 1);
    get (n - 3);
    printf("%d", n);
}
```

If get (6) function is being called in main () then how many times will the get () function be invoked before returning to the main ()?

- (A) 15 (B) 25
(C) 35 (D) 45

17. Consider the following C program [2015]

```
#include<stdio.h>
int f1(void);
int f2(void);
int f3(void);
int x = 10;
int main ( )
{
    int x = 1;
    x += f1 ( ) + f2 ( ) + f3 ( ) + f2 (
);
    printf("%d", x);
    return 0;
}
int f1 ( ) { int x = 25;
x++; return x;}
int f2 ( ) { static int x
= 50; x++; return x;}
int f3 ( ) { x *= 10; return
x};
```

The output of the program is _____

18. Suppose $c = \langle c[0], \dots, c[k-1] \rangle$ is an array of length k , where all the entries are from the set $\{0, 1\}$. For any positive integers a and n , consider the following pseudo code. [2015]

DOSOMETHING (c, a, n)

```
z ← 1
for  $i \leftarrow 0$  to  $k - 1$ 
do  $z \leftarrow z^2 \bmod n$ 
if  $c[i] = 1$ 
then  $z \leftarrow (z \times a) \bmod n$ 
return z
```

If $k = 4$, $c = \langle 1, 0, 1, 1 \rangle$, $a = 2$ and $n = 8$, then the output of DOSOMETHING(c, a, n) is _____

19. What will be the output of the following C program? [2016]

```
void count (int n) {
    static int d = 1;
    printf("%d", n);
    printf("%d", d);
    d++;
    if (n > 1) count (n - 1);
    printf("%d", d);
}
void main ( ) {
    count (3);
}
```

- (A) 3 1 2 2 1 3 4 4 4
(B) 3 1 2 1 1 1 2 2 2
(C) 3 1 2 2 1 3 4
(D) 3 1 2 1 1 1 2

20. The following function computes X^Y for positive integers X and Y . [2016]

```
int exp (int X, int Y)
{
    int res = 1, a = X, b = Y;
    while (b != 0)
    {
        if (b%2 == 0) {a = a*a; b = b/2;}
        else {res = res *a; b = b -1;}
    }
    return res;
}
```

Which one of the following conditions is TRUE before every iteration of the loop?

- (A) $X^Y = a^b$
(B) $(res * a)^Y = (res * X)^b$
(C) $X^Y = res * a^b$
(D) $X^Y = (res * a)^b$

21. Consider the following two functions.

```
void fun1 (int n) {
    if (n == 0) return;
    printf ("%d", n);
    fun2 (n - 2);
    printf ("%d", n);
}

void fun2 (int n) {
    if (n == 0) return;
    printf ("%d", n);
    fun1(++n);
    printf ("%d", n);
}
```

The output printed when fun1 (5) is called is [2017]

- (A) 53423122233445 (B) 53423120112233
(C) 53423122132435 (D) 53423120213243

22. Consider the C functions foo and bar given below:

```
int foo (int val) {
    int x = 0;
    while (val > 0) {
        x = x + foo (val--);
    }
    return val;
}

int bar (int val) {
    int x = 0;
    while (val > 0) {
        x = x + bar (val - 1);
    }
    return val;
}
```

Invocations of foo (3) and bar (3) will result in: [2017]

- (A) Return of 6 and 6 respectively.
(B) Infinite loop and abnormal termination respectively.
(C) Abnormal termination and infinite loop respectively.
(D) Both terminating abnormally.

23. The output of executing the following C program is_____.

```
# include <stdio.h>
int total (int v) {
    static int count = 0;
    while (v) {
        count + = v&1;
        v >> = 1;
    }
    return count;
}
void main ( ) {
```

```
static int x = 0;
int i = 5;
for (; i > 0,i--) {
    x = x + total (i);
}
printf ("%d\n", x);
}
```

[2017]

24. Consider the following C program:

```
#include <stdio.h>
int counter = 0;
int calc (int a, int b) {
    int c;
    counter++;
    if (b==3) return (a*a*a);
    else {
        c = calc (a, b/3);
        return (c*c*c);
    }
}
int main () {
    calc (4, 81);
    printf ("%d", counter);
}
```

The output of this program is _____. [2018]

25. Consider the following program written in pseudo-code. Assume that x and y are integers.

```
Count (x,y) {
    if (y != 1) {
        if (x != 1) {
            print ("*");
            Count (x/2, y);
        }
        else {
            y = y-1;
            Count (1024, y);
        }
    }
}
```

The number of times that the print statement is executed by the call count (1024, 1024) is _____. [2018]

ANSWER KEYS**EXERCISES****Practice Problems 1**

1. A 2. C 3. A 4. D 5. D 6. D 7. D 8. C 9. B 10. A
11. C 12. C 13. C 14. D 15. C

Practice Problems 2

1. D 2. C 3. B 4. B 5. C 6. C 7. A 8. C 9. C 10. A
11. B 12. C 13. B 14. A 15. B

Previous Years' Questions

1. C 2. - 3. -D 4. B 5. C 6. D 7. B 8. D 9. C 10. B
11. C 12. 9 13. 1.72 to 1.74 14. B 15. 51 16. B 17. 230 18. 0 19. A
20. C 21. A 22. C 23. 23 24. 4 25. 10230

Chapter 3

Arrays, Pointers and Structures

LEARNING OBJECTIVES

- Arrays
- Array initialization
- Passing array elements to function
- Two dimensional arrays
- Syntax for 3D array declaration
- Pointers
- Pointer to pointer
- Pointer to void (generic pointer)
- Array of pointers
- Pointer to function
- Dynamic memory management
- Memory allocation function
- Realloc
- Structures
- Nesting of structures
- Array of structures
- Structures & functions
- Union
- Declaration
- Bit fields

ARRAYS

In C we have the following derived data types:

- Arrays
- Pointers
- Structures
- Unions

Imagine a problem that requires to read, process, and print 10 integers. We can declare 10 variables, each with different name. Having 10 different names creates a problem; we need 10 read and 10 write statements each for different variable.

Definition

An array is a collection of elements of same data type. Array is a sequenced collection. So, we can refer to the elements in array as 0th element, 1st element, and so on, until we get the last element. The array elements are individually addressed with their subscripts/indices along with array name. We place subscript value in square brackets ([]) followed by array name. This notation is called indexing.

There is a powerful programming construct, loop, that makes array processing easy. It does not matter if there are 1, 10, 100 or 1000 elements.

We can also use a variable name in subscript, as the value of variable changes; it refers different elements at different times.

Syntax of Array Declaration

```
Data_type array_name [size];
```

Here, data type says the type of elements in collection, array_name is the name given to collection of elements and size says the number of elements in array.

Example: `int marks[6];`

Here, 'int' specifies the type of variable, marks specifies name of variable. The number 6 tells the dimension/size. The '[']' tells the compiler that we are dealing with array.

Accessing array elements: All the array elements are numbered, starting from 0, thus marks [3] is not the third, but the fourth element.

Example: marks[2] – 3rd element

marks[0] – 1st element

We can use the variable as index.

Thus marks[i] – ith element. As the value of i changes, refers different elements in array.

Summary about Arrays

- An array is a collection of similar elements.
- The first element in array is numbered 0, and the last element is one less than the total size of the array.
- An array is also known as subscripted variable.
- Before using an array, its type and dimension must be declared.
- How big an array is, its elements are always stored in contiguous memory locations.
- Individual elements accessed by index indicating relative position in collection.
- Index of an array must be an integer.

Array Initialization

Syntax

```
Data_type array_name[size] = {values};
```

Example:

```
int n[6]= {2,4,8,12,20,25}; // Array initialized with list of values
int num[10] = {2,4,12,20,35};
// remaining 5 elements are initialized with 0
// values
int b[10] = {0}; // Entire array elements initialized with 0.
```

Note:

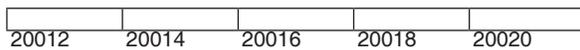
- Till the array elements are not given any specific value, they are supposed to contain garbage values.
- If the number of elements used for initialization is lesser than the size of array, then the remaining elements are initialized with zero.
- Where the array is initialized with all the elements, mentioning the dimension is optional.

Array Elements in Memory

Consider the following declaration – `int num[5]`.

What happens in memory when we make this declaration?

- 10 bytes get reserved in memory, 2 bytes each for 5 integers.
- Since array is not initialized, all five values in it would be garbage. This happens because the default storage class is `auto`. If it is declared as `static`, all the array elements would be initialized with 0.



Note: In *C*, the compiler does not check whether the subscript used for array exceeds the size of array.

Data entered with a subscript exceeding the array size will simply be placed in memory out size the array, and there will be no error/warning message to warn the programmer.

Passing array elements to function

Array elements can be passed to a function by value or by reference.

Example: A program to pass an array by value:

```
void main( )
{
void display(int[ ]);// Declaration
int marks[ ] = {10,15,20,25,30};
display (marks);// function call
}
void display(int n[ ] )// function definition
{
int i;
for(i = 0 ; i < 5 ; i++)
printf("%d ", n[ i ] );
}
```

Output:

```
10 15 20 25 30
```

Here, we are passing the entire array by name. The formal parameter to receive is declared as an array, so it receives entire array elements.

To pass the individual elements of an array, we have to use index of element with array name.

Example: `display (marks[i])`; sends only the *i*th element as parameter.

Example: A program to demonstrate call by reference:

```
void main( )
{
void display (int *);
int marks[ ] = {5, 10 15, 20, 25};
display(&marks[0]);
}
void display(int *p)
{
int i;
for(i = 0; i < 5; i++)
printf("%d ",*(p+i));
}
```

Output:

```
5 10 15 20 25
```

Here, we pass the address of very first element. Hence, the variable in which this address is collected (*p*) is declared as a pointer variable.

Note: Array elements are stored in contiguous memory location, by passing the address of the first element; entire array elements can be accessed.

TWO-DIMENSIONAL ARRAYS

In *C* a two-dimensional array looks like an array of arrays, i.e., a two-dimensional array is the collection of one-dimensional arrays.

Example: `int x[4][2]`;

	0	1
0		
1		
2		
3		

By convention, first dimension says the number of rows in array and second dimension says the number of columns in each row.

In memory, whether it is one-dimensional or a two-dimensional array, the array elements are stored in one continuous chain.

The arrangement of array elements of a two-dimensional array in memory is shown below:

X[0][0]	X[0][1]	X[1][0]	X[1][1]	X[2][0]	X[2][1]	X[3][0]	X[3][1]
6000	6002	6004	6006	6008	6010	6012	6014

Initialization

We can initialize two-dimensional array as one-dimensional array:

```
int a[4][2] = {0,1,2,3,4,5,6,7}
```

The nested braces can be used to show the exact nature of array, i.e.,

```
int a[4][2] = {{0,1},{2,3},{4,5},{6,7}}
```

Here, we define each row as a one-dimensional array of two elements enclosed in braces.

Note: If the array is completely initialized with supplied values, then we can omit the size of first dimension of an array (the left most dimension).

- For accessing elements of multi-dimensional arrays, we must use multiple subscripts with array name.
- Generally, we use nested loops to work with multi-dimensional array.

MULTIDIMENSIONAL ARRAYS

C allows array of two or more dimensions and maximum numbers of dimensions a C program can have depends on the compiler, we are using. Generally, an array having one dimension is called 1D array; array having two dimensions is called 2D array and so on.

Syntax:

type array-name[d1][d2][d3][d4]...[dn];
where dn is the size of last dimension.

Example:

```
int table[5][5][20];
float arr[5][6][5][6][5];
```

In our example array “table” is a 3D. (A 3D array is an array of array of array)

Declaration and Initialization of 3D array

A 3D array can be assumed as an array of arrays; it is an array of 2D arrays and as we know 2D array itself is an array of 1D arrays. A diagram can help you to understand this.

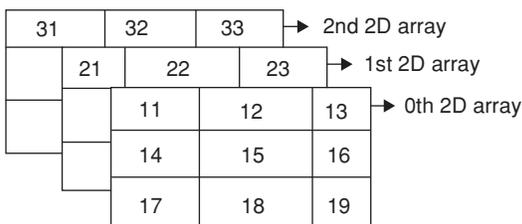


Figure 1 3D array conceptual view

Example:

```
void main( )
{
int i, j, k;
int arr [3] [3] [3] =
{
{11, 12, 13},
{14, 15, 16},
{17, 18, 19}
},
{21, 22, 23},
{24, 25, 26},
{27, 28, 29}
},
{31, 32, 33},
{34, 35, 36},
{37, 38, 39}
}
};
printf("3D Array Elements \n");
for (i = 0; i<3; i++)
{
for(j =0; j <3; j++)
{
for (k= 0; k<3; k++)
{
printf ("% d\t", arr[i][j][k]);
}
printf ("\n");
}
printf ("\n");
}
}
```

Output: 3D Array Elements

11	12	13
14	15	16
17	18	19
21	22	23
24	25	26
27	28	29
31	32	33
34	35	36
37	38	39

Syntax for 3D Array Declaration

data-type array-name [table][row][column];

To store values in any 3D array, first point to table number, row number and lastly to column number.

POINTERS

Pointer is a variable which contains address of another variable. C's clever use of pointers makes it the excellent language.

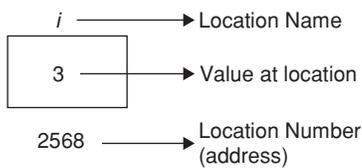
Consider the declaration:

```
int i = 3;
```

The declaration tells the C compiler to:

- Reserve space in memory to hold in integer value.
- Associate the name *i* with this memory location.
- Store the value 3 at this location.

Memory map is:



Computer may choose different location at different times for same variable. The important point is the address is a number.

The expression '&i' gives the address of variable 'i'.

```
p = &i;
```

Assigns the address of 'i' to variable 'p'.

The variable 'p' is declared as:

```
int *p;
```

* tells the compiler that variable 'p' is an address variable.

Memory map of i, *p is –

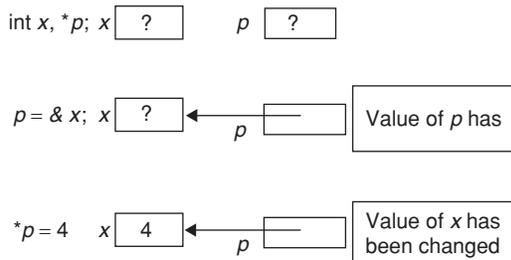
<i>*p</i>	<i>i</i>
2568	3
2720	2568

Now, pointer 'p' is referring to the variable 'i'.

The variable 'i' can be accessed in two ways:

- By using the name of variable.
- By using the pointer variable referring to location 'i'.

The operator '*' can also be used along with pointer variable in expressions. The operator '*' acts as indirection operator.



Usage of 'p' refers to value of 'p', where as '*p' refers to value at the address stored in 'p', i.e., value of 'i'.

```
Example: int *p;
float *x;
char *ch ;
```

Here, *p*, *x* and *ch* are pointer variables, i.e., variables capable of holding address. Since addresses are always whole numbers, pointers would always contain whole numbers.

The declaration `float *x` does not mean that *x* contains floating value, *x* will contain address of floating point variable. Similarly, 'ch' contains address of char value.

Pointer to Pointer

We know, pointer is a variable that contains address of another variable. Now this variable address might be stored in another pointer. Thus, we now have a pointer that contains address of another pointer, known as pointer to pointer.

Example:

```
void main()
{
int i = 3, *p, **q;
p = &i;
q = &p;
printf("\n Address of i = %u", &i);
printf("\n Address of i = %u", p);
printf("\n Address of i = %u", *q);
printf("\n Address of p= %u", &p);
printf("\n Address of p= %u, &q);
printf("\n Address of q = %u", &q)
printf("\n value of i= %d", i);
printf("\n value of i= %d", *(&i));
printf("\n value of i= %d", *p);
printf("\n value of i= %d", **q);
}
```

If the memory map is

<i>**q</i>	<i>*p</i>	<i>i</i>
2010	2000	3
2050	2010	2000

Then the output is:

- Address of *i* = 2000
- Address of *i* = 2000
- Address of *i* = 2000
- Address of *p* = 2010
- Address of *p* = 2010
- Address of *q* = 2050
- Value of *i* = 3

Note: We can extend pointer to a pointer to pointer. In principal, there is no limit on how far we can go on extending this definition.

Pointers for Inter-function Communication

We know that functions can be called by value and called by reference.

- If the actual parameter should not change in called function, pass the parameter-by value.

- If the value of actual parameter should get changed in called function, then use pass-by reference.
- If the function has to return more than one value, return these values indirectly by using call-by-reference.

Example: The following program demonstrates how to return multiple values.

```
void main( )
{
void areaperi(int, int *, int *);
int r;
float a,p;
printf("\n Enter radius of a circle");
scanf("%d", &r);
areaperi(r, &a, &p);
printf("Area = %f", a);
printf("\n Perimeter = %f", p);
}
void areaperi(int x, int *p, int *q)
{
*p = 3.14*x*x;
*q = 2 * 3.14*x;
}
```

Output:

```
Enter radius of circle 5
Area = 78:500000
Perimeter = 31.400000
```

Compatibility: Pointers have a type associated with them. They are not just pointer types, but rather are pointers to a specific type. The size of all pointers is same, which is equal to size of int. Every pointer holds the address of one memory location in computer, but size of variable that the pointer references can be different.

Pointer to Void (Generic Pointer)

A pointer to void is a generic type; this can point to any type. Its limitation is that the pointed data cannot be referenced directly. Since void pointer has no object type, so its length is undetermined; it cannot be dereference unless it is cast.

Example: The following example demonstrates generic pointer.

```
void main ( )
{
int a = 10;
float x = 5.7;
void *p;
p = &a;
printf("\n value of a = %d", *((int*)p));
p = &x;
printf ("\n value of x = % f", *((float *)p));
}
```

Output:

```
value of a = 10
value of x = 5.700000
```

Operations can be Performed on Pointers

1. Addition of a number to a pointer.

Example: `int i = 4, *j, *k;`
`j = &i;`
`j = j + 1;`
`k = j + 5;`

2. Subtraction of a number from a pointer.

Example: `int i = 4, *j, *k;`
`j = &i; j = j - 1;`
`k = j - 3;`

3. Subtraction of one pointer from another. One pointer variable can be subtracted from another (provided both variables point to same array elements). The resulting value indicates the number of bytes (elements) separating (the corresponding array elements).

Example:

```
void main ( )
{
int a[ ] = {5,10,15,20,25} , *i, *j;
i = &a[0];
j = &a[4];
printf("%d, %d", j-i, *j-*i);
}
```

Output: 4, 20

The expression `j-i` prints 4 but not 8. because `j` and `i` pointing to integers that are 4 integers apart.

4. Comparison of two pointer variables. Pointer variables can be compared provided both pointing to the same data type.

Notes: Do not attempt the following operations on pointers:

1. Addition of two pointers.
2. Multiplication of a pointer with a number or another pointer.
3. Division of a pointer with a number or another pointer.

Important points about pointer arithmetic

- A pointer when incremented always points to an immediately next location.
- A pointer when decremented always points to an element precedes the current element.

Notice the difference with:

`(*p)++`

Here, the expression would have been evaluated as the value pointed by `p` increased by one. The value of `p` would not be modified if we write

`*p++ = *q++;`

Because `++` has a higher precedence than `*`, both `p` and `q` are increased, but because both increase operators (`++`) are used as postfix and not prefix, the value assigned to `*p` is

*q before both p and q are increased. And then both are increased, it would be equivalent to

```
*p = *q;
++p;
++q;
```

Implementation of arrays in C

Array name is the pointer to the first element in array. The following discussion explains how pointers are used for implementing arrays in C.

```
int n[ ] = {10, 20, 30, 40, 50};
```

n	10	20	30	40	50
	5512	5514	5516	5518	5520

- We know that mentioning the array name gets the base address.

```
int *p = n;
```

Now 'p' points to 0th element of array 'n'.

- 0th element can be accessed as *array_name.

```
int x = *n;
stores n[0] into 'x'.
```

- we can say that *array_name and *(array_name+0) are same. This indicates the following are same.

```
num[i]
*(num + i)
*(i+num)
```

(num is an array; i is an index)

ARRAY OF POINTERS

The way there can be an array of ints or array of floats, similarly there can be an array of pointers. An array of pointers is the collection of addresses.

These arrays of pointers may point to isolated elements or an array element.

Example 1: Array of pointers pointing to isolated elements:

```
int i = 5, j=10, k =15;
int *ap[3];
ap[0] = &i; ap[1] = &j; ap[2] = &k;
```

Example 2: Array of pointers pointing to elements of an array:

```
int a[ ] = {0, 20, 45, 50, 70};
int *p[5], i ;
for(i = 0; i < 5 ; i++ )
p[i] = &a[i] ;
```

Example 3: Array of pointers pointing to elements of different arrays;

```
int a[ ] = {5, 10, 20, 25};
int b[ ] = {0, 100, 200, 300, 400};
int c[ ] = {50, 150, 250, 350, 450};
int *p[3];
p[0] = a; p[1] = b; p[2]=c;
```

Example 4: Array of pointers pointing to 0th element of each row of a two-dimensional array

```
int a[3][2] = {{1,2} {3,4}, {5,6}};
int *p[3];
p[0] = a[0]; p[1] =a[1];p[2] = a[2];
```

POINTER TO FUNCTION

Function is a set of instructions stored in memory, so the function also contains the base address. This address can hold by using a pointer called pointer to function.

Syntax:

```
return_type (*function_pointer)(parameter - list );
```

Example: int (*fp)(float, char, char);

Example:

```
// pointer to functions
# include <iostream>
Using name space std;
int addition(int a, int b)
{
return (a + b);
}
int subtraction(int a, int b)
{
return (a - b) ;
}
int operation (int x, int y, int (*funtocall)
(int, int))
{
int g;
g = (*funtocall)(x, y);
return (g);
}
int main( )
{
int m, n;
int (*minus)(int, int) = subtraction;
m = operation(7, 5, addition);
n = operation(20, m, minus);
cout << n;
return 0;
}
```

In the example, minus is a pointer to a function that has two parameters of type int. It is immediately assigned to point to the function subtraction, all in a single line.

Example: Program to demonstrate function pointer

```
int add(int, int);
int sub(int, int);
void main( )
{
Int (*fp) (int, int);
fp = add;
```

```
printf("\n 4+5=%d", fp(4,5));
fp = sub;
printf ("\n 4 - 5 = %d", fp(4,5));
}
int add(int x, int y)
{
return x + y;
}
int sub(int x,int y)
{
return x - y;
}
```

Output: 4 + 5 = 9
4 - 5 = -1

Pointer to structure The main usage of pointer to structure is we can pass structure as parameter to function as call by reference.

The other usage is to create linked lists and other dynamic data structures which depend on dynamic allocation.

Consider the declaration

struct employee

```
{
char name[20];
Int age;
float salary;
};
```

struct employee *p;

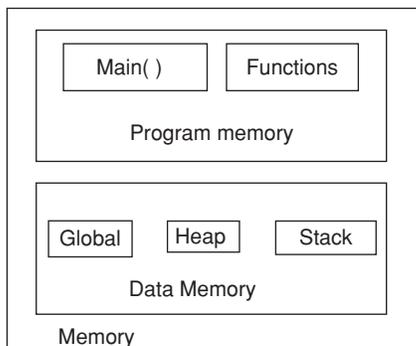
Variable of structures can be accessed using ‘.’ Operator (or) → operator that is

```
(*p).age = 20 ; (or) p → age = 20;
(*p).salary = 40, 231.0; (or) p → salary = 40,231.0;
```

DYNAMIC MEMORY MANAGEMENT

We can allocate the memory to objects in two ways—static and dynamic allocation. Static memory allocation requires declaration and definition of memory fully specified in the source program. The number of bytes required cannot be changed during run time. Dynamic memory allocation uses predefined functions to allocate and de-allocate memory for data dynamically during the execution of program.

We can refer to dynamically allocated memory only through pointers. Conceptual view of memory:



Memory Allocation Function

- Static memory allocation uses stack memory for variables.
- Dynamic memory management allocates memory from heap.

The following are the four memory management functions available in alloc.h and stdlib.h.

- 1. Malloc (Block memory allocation):** Malloc function allocates block of memory that contained the number of bytes specified in parenthesis. It returns ‘void’ pointer to the first byte of allocated memory. The allocated memory is not initialized. If the memory allocation is not successful then it return NULL pointer.

Declaration

```
void *malloc (size_t size);
```

The type size_t is defined as unsigned int in several header files including `stdio.h`.

Syntax: pointer = (type*) malloc(size);

- 2. Calloc (contiguous memory allocation):** Calloc is primarily used to allocate memory for arrays. It initializes the allocated memory with null characters.

Declaration: void *calloc (size_t ele_count, size_t ele_size);

Syntax: ptr = (type*)calloc(ele-count,ele-size);

- 3. Realloc (reallocation of memory):** The realloc function is highly inefficient. When given a pointer to a previously allocated block of memory, realloc changes the size of block by deleting or extending the memory at the end of block. If the memory cannot be extended, then realloc allocates completely new block, copies the contents from existing memory location to new location, and deletes the old location.

Declaration: void *realloc (void *ptr, size_t new_size);

Syntax: ptr = (type*)realloc(ptr, new_size);

- 4. Free (Releasing memory):** When the memory allocated by malloc, calloc or realloc is no longer needed, they can be freed using the function free().

Declaration: void free(void *ptr);

Syntax: free(ptr);

Free function de-allocates complete memory referenced by the pointer. Part of the memory block cannot be de-allocated.

STRUCTURES

Arrays are used to store large set of data and manipulate them but the disadvantage is that all the elements stored in an array are to be of the same data type. When we require using a collection of different data items of different data types, we can use a structure.

- Structure is a method of packing data of different types.
- A structure is a convenient method of handling a group of related data items of different data types.

Syntax for declaration

```
struct sturct_name
{
Data_type_1 var1;
Data_type_2 var2;
:
Data_type_n varn;
};
```

Example:

```
struct lib - books
{
char title [20];
char author[15];
int pages;
float price;
};
```

The keyword struct declares a structure to hold the details of four fields namely title, author, pages and price, these are members of the structures.

We can declare structure variables using the tag name anywhere in the program.

Example: struct lib – books book1, book2, book3;

- Declares book1, book2, book3 as variables of type struct lib _ books, each declaration has four elements of the structure lib _ books.

Memory map of book1:

Book1	Title	20 bytes
	Author	15 bytes
	Pages	2 bytes
	Price	4 bytes

- Memory will not be allocated to the structure until it is instantiated. i.e., till the declaration of a variable to structure.
- To access the members of a structure variable, C provides the member of (.) operator.

Example: To access author of book 1 – book1. author

Syntax: structure_var.member_name;

- The structures can also be initialized as any other variable of C.

Example: struct lib-books book4={"Let us C", "yashwanth", 450, 200.95};

Note: The values must provide in the same order as they appear in structure declaration.

- One structure variable can be assigned to another structure variable.
- Structure variables cannot be compared.

Example:

```
# include <stdio.h>
void main( )
{
Struct s1{
int id_no;
char name[20];
```

```
char address[20];
char combination[3];
    int age;
    } newstudent;
printf (" Enter student Information");
printf ("Enter student id - no");
scanf ("%d", &newstudent.id_no);
printf (" Enter the name of the student");
scanf ("%s", &newstudent.name);
printf (" Enter the address of the student");
scanf ("%s", &newstudent.address);
printf("Enter the combination of the student");
scanf("%s", &newstudent.combination);
printf (" Enter the age of student);
scanf ("%d ", &newstudent.age");
printf (" student information");
printf (" student id-no = %d", newstudent.
id - no);
printf("student name = %s", newstudent.
name);
printf("student address = %s", newstudent.
address);
printf ("students combination = %s", newstu-
dent. combination);
printf("Age of student = %d", newstudent.
age);
}
```

Nesting of Structures

The structures can be nested in two ways:

- Placing the structure variable as a member in another structure declaration.
- Declaration of the entire structure in another structure.

Example:

```
struct date
{
int day;
int month;
int year;
};
struct student
{
int id_no;
char name[20];
char address [20];
int age;
structure date doa;
} oldstudent, newstudent;
```

The structure 'student' contains another structure date as one of its members.

To access the day of date of admission (doa) of old student – oldstudent.doa.day.

Example:

```
struct outer
{
```

```
int o1;
float o2;
struct inner
{
int i1;
float i2;
};
} out1, out2;
```

The innermost members in a nested structure can be accessed by chaining all the concerned structure variables, from outermost to innermost; accessing `i1` for `out1-out1.inner.i1`;

Array of Structures

It is possible to define an array of structures. For example, if we are maintaining information of all the students in the college and if 100 students are studying in the college, we need to use an array than single variables.

Example:

```
structure information
{
int id_no;
char name[20];
char address[20];
char combination[3];
int age;
}
student[100];
```

Example:

```
# include <stdio.h>
{
struct info
{
int id_no;
char name[20];
char address[20];
char combination[3];
int age;
}
struct info std[100];
int, i ,n;
printf (" Enter the number of students");
scanf ("%d", &n);
scanf("Enter id_no, name, address, combination and age");
for (i = 0; i<n; i++)
scanf(" %d %s %s %s %d", &std[i].id_no,
std[i].name, std[i].address,
std[i]. combination,&std [i].age);
printf("student information");
for (i = 0 ; i < n; i++)
printf("%d %s %s % s % d", std[i].id_no,
std[i].name, std[i].address, std[i]. combi-
nation, std[i]. age);
```

Structures and Functions

- An entire structure can be passed as a parameter like any other variable.
- A function can also return a structure variable.

Example:

```
# include <stdio.h>
struct employee
{
int emp_id;
char name[25];
char department[10];
float salary;
};
void main( )
{
static struct employee emp1 = {
12, "shyam", "computer", 7500.00};
/* sending entire employee structure */
display(emp1);
}
/* function to pass entire structure vari-
able */
display(empf)
struct employee empf
{
printf (" %d %s % s %f", empf.empid, empf.
name, empf.department, empf.salary);
}
```

UNION

Union, like structure contains members whose individual data types may differ from one another. The members that compose union all share the same storage area within the computer's memory whereas each member within a structure is assigned its own unique storage area. Thus, unions are used to conserve memory.

Declaration

```
union item
{
int m;
float p;
char c;
}Code;
```

This declares a variable code of type union item.

The union contains three members each with a different data type. However, we can use only one of them at a time. The compiler allocates a piece of storage that is large enough to access a union member; we can use the same syntax that we use to access structure members, i.e.,

```
Code.m
Code.p
Code.c
```

are all valid member variables. During accessing, we should make sure that we are accessing the member whose value is currently stored.

Example:

```
union marks
{
float perc;
char grade;
}
main( )
{
union marks student1;
student1.perc = 98.5;
printf("marks are %f address is %16lu", stu-
dent1.perc, &student1.perc);
student1.grade = 'c';
printf("grade is %c address is %16lu", stu-
dent1.grade, &student1.grade);
}
```

Example:

```
# include <stdio.h>
void main ( )
{
Union u_example
{
float decval;
int p_num;
double my_value;
}U1;
U1.my_value = 125.5;
U1.pnum = 10;
```

```
U1.decval = 1000.5f;
printf("decval = %f pnum = %d my_value = % lf
", U1.decval, U1.pnum, U1.my_value);
printf(" U1 size = %d decval size =%d,
pnum size = %d my-value size = % d",
sizeof (U1), sizeof (U1.decval), sizeof
(U1.pnum), sizeof (U1.my_value));
}
```

Bit Fields

When a program variable 'x' is declared as int, then 'x' takes the values from (-2^{15}) to $(2^{15} - 1)$, if x in the program takes only two values, 1 and 0, which requires only one bit, then the remaining 15 bits are waste.

In order to not to have this wastage, we can use bit fields with the several variables with the small enough maximal values, which can pack into a single memory location

Example:

```
struct student
{
Int gender : 1 ; // gender takes only 0,1
values
Int marriage : 2 ; // marriage takes 4(0, 1,
2, 3) values
Int marks : 7 ; // marks takes values from
0 - 127
}
```

EXERCISES**Practice Problems I**

Directions for questions 1 to 15: Select the correct alternative from the given choices.

1. Output of the following C program is

```
intF(int x, int *py, int **pz)
{
int y, z;
** pz+= 1;
z = *pz;
*py+= 2;
y = *py;
x+ = 3;
return x+y+z;
}
void main( )
{
int c, *b, **a ;
c = 4;
b = &c;
a = &b;
printf( "%d", F(c, b, a));
}
```

- (A) 30 (B) 22
(C) 20 (D) Error

2. main()

```
{
char *ptr;
ptr = "Hello World";
printf("%c\n", *&*ptr);
}
```

Output of the above program is

- (A) Garbage value
(B) Error
(C) H
(D) Hello world

3. #include <stdio.h>

```
main( )
{
register a =10;
char b[ ] = "Hi";
printf("%s %d ", b, a);
}
```

Output is

- (A) Hi 10 (B) Error
(C) Hi (D) Hi garbage value

4. main ()

```

{
    int fun( ) ;
    (*fun)( ) ;
}
int fun( )
{ printf("Hello") ;
}

```

- (A) Hello (B) Error
(C) No output (D) H

5. Let B be a two-dimensional array declared as $B : \text{array}[1\dots 10][1\dots 15]$ of integer;

Assuming that each integer takes one memory location the array is stored in row major order and the first element of the array is stored at location 100, what is the address of the element $B[i][j]$?

- (A) $15i + 10j + 84$ (B) $15i + j - 16$
(C) $15i + j$ (D) $15i + j + 84$

6. Consider the following C program which is supposed to compute the transpose of a given 4×4 matrix M . Note that, there is a Y in the program which indicates some missing statements. Choose the correct option to replace Y in the program.

```

#include <stdio.h>
int M[4][4] = { 8, 10, 9, 16, 12, 13, 11,
15, 14, 7, 6, 3, 4, 2, 1, 5 };
main( )
{
    int i, j, temp;
    for (i = 0; i<4; ++i)
    {
        Y
    }
    for (i=0; i<4; ++i)
    for (j=0; j<4; ++j)
        printf("%d", M[i][j]);
}

```

- (A) for (j=0; j<4; ++j)
{
M[j][i] = temp;
temp = M[j][i];
M[j][i] = M[i][j];
}
- (B) for (j=0; j<4; ++j)
{
temp = M[j][i];
M[i][j] = M[j][i];
M[j][i] = temp;
}
- (C) for (j=i; j<4; ++j)
{
temp = M[i][j];
M[i][j] = M[j][i];
M[j][i] = temp;
}

```

(D) for (j=i; j<4; ++j)
{
    M[i][j] = temp;
    temp = M[j][i];
    M[j][i] = M[i][j] ;
}

```

7. Consider the C program shown below:

```

#include <stdio.h>
#define print(a) printf("%d", a)
int a;
void z(int n)
{
    n += a;
    print (n);
}
void x(int *p)
{
    int a = *p+2;
    z(a) ;
    *p = a;
    print(a);
}
main(void)
{
    a = 6;
    x(&a);
    print(a);
}

```

The output of this program is

- (A) 14 8 6 (B) 16 6 6
(C) 8 6 6 (D) 22 11 12

8. Consider the program below:

```

#include <stdio.h>
int fun(int n, int *p)
{
    int x,y;
    if (n<=1)
    {
        *p = 1;
        return 1;
    }
    x = fun(n-1, p);
    y = x +p;
    *p = x;
    return y;
}
int main( )
{
    int a =15;
    printf( "%d\n", fun(5, &a));
    return 0;
}

```

The output value is

- (A) 14 (B) 15
(C) 8 (D) 95

9. Consider the following C program segment

```
char p[20] ;
int i;
char *s = "string" ;
int l = strlen(s);
for (i=0; i<l; i++)
p[i] = s[l - i] ;
printf("%s", p) ;
```

The output of the program is

- (A) string
 (B) gnirt
 (C) gnirts
 (D) No output is printed
10. # include <stdio.h>
 main()
 {
 struct AA
 {
 int A = 5;
 char name[] = "ANU";
 };
 struct AA *p = malloc(sizeof(struct
 AA));
 printf("%d",p->A);
 printf("%s",p->name);
 }
 Output of the program is
 (A) 5 ANU
 (B) Runtime error
 (C) Compiler error
 (D) Linker error
11. The declaration
 union u_tag {
 int ival;
 float fval;
 char sval;
 } u;
 denotes *u* is a variable of type *u_tag* and
 (A) *u* can have a value of int, float and char
 (B) *u* can represent either integer value, float value or
 character value at a time
 (C) *u* can have a value of float but not integer
 (D) None of the above
12. If the following program is run from command line as
 myprog 1 2 3, what would be the output?

```
main (int argc, char *argv[ ])
{
int i;
i = argv [1] + argv [2] - argv [3];
printf ("%d", i);
}
```

- (A) 123
 (B) 6
 (C) 0
 (D) Error

13. The following C program is run from the command line
 as

myprog one two;
 what will be the output?

```
main (int argc, char *argv [ ])
{
printf ("%c",**++argv);
}
```

- (A) m
 (B) o
 (C) myprog
 (D) one

14. The following program

```
change(int *);
main( ) {
int a = 4;
change(a);
printf ("%d", a);
}
change(a)
int a;
{
printf("%d", a);
}
```

Outputs

- (A) 44
 (B) 55
 (C) 34
 (D) 22

15. What is the output of the following program:

```
main( )
{
const int x = 10;
int *ptrx;
ptrx = &x;
*ptrx = 20;
printf ("%d", x);
}
```

- (A) 5
 (B) 10
 (C) Error
 (D) 20

Practice Problems 2

Directions for questions 1 to 11: Select the correct alternative from the given choices.

- The following program segment


```
int *i;
*i = 10;
```

 (A) Results in run time error
 (B) Is a dangling reference
 (C) Results in compilation error
 (D) Assigns 10 to *i*
- A $m \times n$ matrix is stored in column major form. The expression which accesses the (*ij*)th entry of the same matrix is
 (A) $n \times (j - 1) + i$
 (B) $m \times (j - 1) + i$
 (C) $n \times (m - 1) + ij$
 (D) $m \times (n - 1) + j$
- `int * S[a]` is 1D array of integers, which of the following refers to the third element in the array?
 (A) $*(S+2)$ (B) $*(S+3)$
 (C) $S+2$ (D) $S+3$
- If an array is declared as `char a[10][12]`; what is referred to by `a[5]`?
 (A) Pointer to 3rd Row
 (B) Pointer to 4th Row
 (C) Pointer to 5th Row
 (D) Pointer to 6th Row
- The following code is run from the command line as `myprog 1 2 3`. What would be the output?


```
main(int argc, char *argv[ ])
{
    int i, j = 0;
    for (i = 1; i < argc; i++)
        j = j + atoi (argv [i]);
    printf ("%d", j);
}
```

 (A) 123 (B) 6
 (C) Error (D) "123"
- What will be the following C program output?


```
main (int argc, char *argv[ ], char *env
[ ]) {
    int i;
    for(i = 1; i < argc; i++)
        printf ("%s", env[i]);
}
```

 (A) List of all arguments
 (B) List of all path parameters
 (C) Error
 (D) List of environment variables
- The declaration


```
enum colors {
    red,
```

```
    blue,
    yellow = 1,
    green
};
```

assigns the value 1 to

- (A) Red and Yellow
 (B) Blue
 (C) Red and blue
 (D) Blue and yellow
- What would be the output of the following program?


```
sum = 0;
for (i = -10; i < 0; i++)
    sum = sum + abs(i);
printf ("%d", sum);
```

 (A) 100 (B) -505
 (C) 55 (D) -55
 - An integer occupies 2 bytes of memory, float occupies 4 bytes and character occupies 1 byte. A structure is defined as:


```
struct tab {
    char a;
    int b;
    float c;
} table [10];
```

 Then the total memory requirement (in bytes) is
 (A) 14 (B) 70
 (C) 40 (D) 100
 - What are the values of `u1` and `u2`?


```
int u1, u2;
int x = 2;
int *ptr;
u1 = 2*(x + 10);
ptr = &x;
u2 = 2>(*ptr + 10);
```

 (A) $u1 = 8, u2 = 16$
 (B) $u1 = 23, u2 = 24$
 (C) $u1 = 24, u2 = 24$
 (D) None of the above
 - What is the output?


```
func(a, b)
int a, b;
{
    return (a = (a == b));
}
main ()
{
    int process(), func();
    printf("The value of process is %d", pro-
cess (func,3,6));
}
process (pf, val1, val2)
int (*pf) ();
```

```
int val1, val2;
{
return ((*pf) (val1, val2));
}
```

- (A) The value of process is 0
- (B) The value of process is 3
- (C) The value of process is 6
- (D) Logical error

PREVIOUS YEARS' QUESTIONS

1. Consider the following program in C language:

```
# include < stdio. h>
main ()
{
int i;
int *pi = &i;
scanf ("%d", pi);
printf("%d\n", i + 5);
}
```

Which one of the following statement is TRUE?

[2014]

- (A) Compilation fails
 - (B) Execution results in a run-time error
 - (C) On execution, the value printed is 5 more than the address of variable *i*.
 - (D) On execution, the value printed is 5 more than the integer value entered.
2. Consider the following C function in which size is the number of elements in the array *E*:

```
int MyX (int *E, unsigned int size)
{
int Y = 0;
int Z;
int i, j, k;
for (i = 0; i < size; i++)
Y = Y + E[i];
for (i = 0; i < size; i++)
for (j = 1; j < size; j++)
{
Z = 0;
for (k = i; k <= j; k++)
Z = Z + E[k];
if (Z > Y)
Y = Z;
}
return Y;
}
```

The value returned by the function My X is the **[2014]**

- (A) maximum possible sum of elements in any sub-array of array *E*.
 - (B) maximum element in any sub-array of array *E*.
 - (C) sum of the maximum elements in all possible sub-arrays of array *E*.
 - (D) the sum of all the elements in the array *E*.
3. The output of the following C program is _____

[2015]

```
void f1 (int a, int b) {
int c;
c=a; a=b; b=c;
}
void f2(int *a, int *b) {
int c;
c=*a; *a=*b; *b=c;
}
int main ( ) {
int a=4, b=5, c=6;
f1 (a, b);
f2 (&b, &c);
printf("%d", c-a-b);
}
```

4. What is the output of the following C code? Assume that the address of *x* is 2000 (in decimal) and an integer requires four bytes of memory. **[2015]**

```
int main ( ) {
unsigned int x[4] [3] =
{ {1, 2, 3}, {4, 5, 6}, {7, 8, 9},
{10, 11, 12}};
printf ("%u, %u, %u", x + 3, *(x +
3), *(x + 2) + 3);
}
```

- (A) 2036, 2036, 2036
- (B) 2012, 4, 2204
- (C) 2036, 10, 10
- (D) 2012, 4, 6

5. Consider the following function written in the C programming language. **[2015]**

```
void foo(char *a {
if ( *a && *a != ' '){
foo(a + 1);
putchar(*a);
}
}
```

The output of the above function on input "ABCD EFGH" is

- (A) ABCD EFGH
- (B) ABCD
- (C) HGFE DCBA
- (D) DCBA

6. Consider the following C program segment. **[2015]**

```
#include <stdio.h>
int main()
{
char s1[7] = "1234", *p;
```

```

p = s1 + 2;
*p = '0';
printf("%s", s1);
}

```

What will be printed by the program?

- (A) 12 (B) 120400
(C) 1204 (D) 1034

7. Consider the following C program [2015]

```

#include<stdio.h>
int main ( )
{
    static int a[ ] = {10, 20, 30, 40,
50};
    static int *p[ ] = {a, a+3, a+4,
a+1, a+2};
    int **ptr = p;
    ptr++;
    printf("%d%d", ptr-p, **ptr);
}

```

8. Consider the following C program. [2016]

```

void f(int, short);
void main()
{
    int i = 100;
    short s = 12;
    short *p = &s;
    ____; // call to f()
}

```

Which one of the following expressions, when placed in the blank above, will **NOT** result in a type checking error?

- (A) $f(s,*s)$ (B) $i=f(i,s)$
(C) $f(i,*s)$ (D) $f(i,*p)$

9. Consider the following C program. [2016]

```

# include<stdio.h>
void mystery (int *ptra, int *ptrb) {
    int *temp;
    temp = ptrb;
    ptrb = ptra;
    ptra = temp;

}
int main ( ) {
    int a = 2016, b = 0, c = 4, d = 42;
    mystery (&a, &b);
    if (a < c)

```

```

    mystery(&c, &a);
    mystery (&a, &d);
    printf(“%d\n”, a)
}

```

The output of the program is _____.

10. The following function computes the maximum value contained in an integer array $p []$ of size n ($n \geq 1$). [2016]

```

int max (int *p, int n) {
    int a = 0, b = n - 1;
    while (____) {
        if (p [a] <= p [b]) {a = a+1;}
        else { b = b - 1;}
    }
    return p[a];
}

```

The missing loop condition is

- (A) $a != n$
(B) $b != 0$
(C) $b > (a + 1)$
(D) $b != a$

11. The value printed by the following program is _____. [2016]

```

void f(int* p, int m) {
    m = m + 5;
    *p = *p + m;
    return;
}
void main () {
    int i = 5, j = 10;
    f(&i, j);
    printf(“%d”, i + j);
}

```

12. Consider the following program: [2016]

```

int f(int *p, int n)
{   if (n <= 1) return 0;
    else return max (f(p + 1, n - 1), p [0] - p [1] );
}
int main ()
{
    int a[ ] = {3,5,2,6,4};
    printf(“%d”, f(a,5));
}

```

Note: $\max(x,y)$ returns the maximum of x and y .

The value printed by this program is _____

13. Consider the following C code:

```
# include <stdio.h>
int *assignval (int *x, int val) {
    *x = val;
    return x;
}
void main ( ) {
    int *x = malloc (sizeof (int));
    if (NULL == x) return;
    x = assignval (x, 0);
    if (x) {
        x = (int *) malloc
            (sizeof (int));
        if (NULL == x) return;
        x = assignval (x, 10);
    }
    printf ("%d\n", *x);
    free (x);
}
```

The code suffers from which one of the following problems: [2017]

- (A) compiler error as the return of malloc is not type-cast appropriately
- (B) compiler error because the comparison should be made as $x == \text{NULL}$ and not as shown
- (C) compiles successfully but execution may result in dangling pointer
- (D) compiles successfully but execution may result in memory leak

14. Consider the following C program.

```
# include <<stdio.h>
# include <<string.h>
void printlength (char *s, char *t)
{
    unsigned int c = 0;
    int len = ((strlen(s) - strlen
        (t)) > c) ? strlen (s) : strlen
        (t);
    printf ("%d\n", len);
}
void main ( ) {
    char *x = "abc";
    char *y = "defgh";
    printlength (x, y);
}
```

Recall that strlen is defined in string.h as returning a value of type size_t, which is an unsigned int. the output of the program is _____. [2017]

15. Given the following binary number in 32-bit (single precision) IEEE-754 format:

00111110011011010000000000000000

The decimal value closest to this floating-point number is [2017]

- (A) 1.45×10^1
- (B) 1.45×10^{-1}
- (C) 2.27×10^{-1}
- (D) 2.27×10^1

16. Match the following:

(P) static char var;	(i) Sequence of memory locations to store addresses
(Q) $m = \text{malloc}(10);$ $m = \text{NULL};$	(ii) A variable located in data section of memory
(R) char *ptr[10];	(iii) Request to allocate a CPU register to store data
(S) register int var1;	(iv) A lost memory which cannot be freed

[2017]

- (A) P → (ii), Q → (iv), R → (i), S → (iii)
- (B) P → (ii), Q → (i), R → (iv), S → (iii)
- (C) P → (ii), Q → (iv), R → (iii), S → (i)
- (D) P → (iii), Q → (iv), R → (i), S → (ii)

17. Consider the following function implemented in C:

```
void printxy (int x, int y) {
    int *ptr;
    x = 0;
    ptr = &x;
    y = *ptr;
    *ptr = 1;
    printf ("%d, %d" x, y);
}
```

The output of invoking printxy (1, 1) is [2017]

- (A) 0, 0
- (B) 0, 1
- (C) 1, 0
- (D) 1, 1

18. Consider the following snippet of a C program. Assume that swap (&x, &y) exchanges the contents of x and y.

```
int main () {
    int array[] = {3, 5, 1, 4, 6, 2};
    int done = 0;
    int i;
    while (done == 0) {
        done = 1;
        for (i=0; i <=4; i++) {
            if (array[i] < array[i+1]) {
                swap(&array[i], &array[i + 1]);
                done = 0;
            }
        }
        for (i=5; i >=1; i--) {
            if (array[i] > array[i-1]) {
                swap(&array[i], &array[i-1]);
                done = 0;
            }
        }
    }
}
```

```
printf("%d", array[3]);
}
```

The output of the program is _____. [2017]

19. Consider the following C Program.

```
#include<stdio.h>
#include<string,h>
int main () {
    char* c = "GATECSIT2017";
    char* p = c;
    printf("%d",
        (int) strlen(c+2[p]-6[p]-1)) ;
    return 0;
}
```

The output of the program is _____. [2017]

20. Consider the following C program.

```
#include<stdio.h>
struct Ournode {
    char x, y, z;
};
int main () {
    struct Ournode p = {'1', '0', 'a'+2};
    struct Ournode *q = &p;
    printf ("%c, %c", *( (char*)q+1),
        * ( (char*)q+2) );
    return 0;
}
```

The output of this program is: [2018]

- (A) 0, c
- (B) 0, a+2
- (C) '0', 'a+2'
- (D) '0', 'c'

21. Consider the following C program:

```
#include<stdio.h>
void fun1 (char *s1, char * s2) {
    char *tmp;
    tmp = s1;
    s1 = s2
    s2 = tmp;
}
void fun2 (char **s1, char **s2) {
    char *tmp;
    tmp = *s1;
    *s1 = *s2;
    *s2 = tmp;
}
int main () {
    char *str1 = "Hi", *str2 = "Bye";
    fun1 (str1, str2);
    printf ("%s %s ", str1, str2);
    fun2 (&str1, &str2);
    printf ("%s %s", str1, str2);
    return 0;
}
```

The output of the program above is: [2018]

- (A) Hi Bye Bye Hi
- (B) Hi Bye Hi Bye
- (C) Bye Hi Hi Bye
- (D) Bye Hi Bye Hi

ANSWER KEYS

EXERCISES

Practice Problems 1

1. B 2. C 3. A 4. A 5. D 6. C 7. A 8. C 9. D 10. C
11. B 12. D 13. B 14. A 15. D

Practice Problems 2

1. B 2. B 3. A 4. D 5. B 6. D 7. D 8. C 9. B 10. C
11. A

Previous Years' Questions

1. D 2. A 3. -5 4. A 5. D 6. C 7. 140 8. D 9. 2016 10. D
11. 30 12. 3 13. D 14. 3 15. C 16. A 17. C 18. 3 19. 2 20. A
21. A

Chapter 4

Linked Lists, Stacks and Queues

LEARNING OBJECTIVES

- ☞ Data structure
- ☞ Linked list
- ☞ Single-linked list
- ☞ Double-linked list
- ☞ Circular linked list
- ☞ Double circular-linked list
- ☞ Stack
- ☞ Queue
- ☞ Double-ended queue
- ☞ Circular queue
- ☞ Priority queue
- ☞ Array implementation
- ☞ Linked list implementation
- ☞ Linked list implementation of priority queue

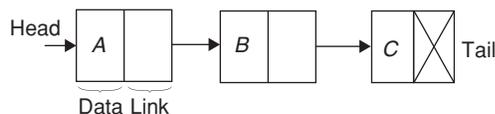
DATA STRUCTURE

Data structure represents the logical arrangement of data in computer memory for easily accessing and maintenance.

LINKED LIST

A linked list is a data structure that consists of a sequence of nodes, each of which contains data field and a reference (i.e., link) to next node in sequence.

- Generally node of linked list is represented as self-referential structure.
- The linked list elements are accessed with special pointer(s) called head and tail.



- The principal benefit of a linked list over a conventional array is that the list elements can easily be added or removed without reallocation or reorganization of the entire structure because the data items need not be stored contiguously in memory or on disk.
- Linked lists allow insertion and removal of nodes at any point in the list.
- Finding a node that contains a given data, or locating the place where a new node should be inserted may require scanning most or all of the list elements.
- The list element does not have to occupy contiguous memory.
- Adding, insertion or deletion of list elements can be accomplished with the minimal disruption of neighbouring nodes.

SINGLE-LINKED LIST

List in which each node contains only one link field.

Node structure

```
struct
{
int ele;
struct node * next;
};
typedef struct node Node;
```

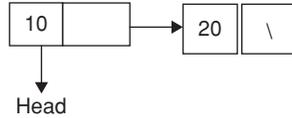
Creating a linked list with two nodes of type list node

Creating a linked list with 2 nodes

```
struct node
{
Int ele;
struct node * next ;
};
typedef struct node Node ;
Node * ptr1, * ptr2;
ptr1 = getnode ();
ptr2 = getnode ();
if((ptr1) && (ptr2))
{
Printf("No memory");
exit(1);
}
Ptr1 → ele = 10;
```

```
Ptr1 → next = ptr2;
Ptr2 → ele = 20;
Ptr2 → next = NULL;
Head = ptr1;
```

the linked list appears as below



Operations on SLL (single-linked list)

- Insert at Head
- Insert at Tail
- Insert in Middle
- Delete Head
- Delete Tail
- Delete Middle
- Search
- Display

Declare two special pointers called head and tail as follows:

```
Node *Head, *Tail;
Head = Tail = NULL;
```

Head or tail is NULL represents list is empty.

Steps for Insertion:

1. Allocate memory
2. Read data
3. Adjust references

Insert head element

```

1. void ins _ Head (int x)
2. {
3. Node *temp;
4. temp = (Node *) malloc(sizeof (Node));
5. temp → ele = x;
6. temp → next = Head;
7. Head = temp;
8. if (Tail == NULL)
9. Tail = Head;
10 }
```

- Step 4 allocates memory
- Step 5 read data
- Steps from 6 to 9 adjust reference
- 'if' condition represents first insertion

Insert tail element

```

1. void ins_tail (int x)
2. {
3. Node *temp;
4. temp = (Node *) malloc (sizeof (Node));
5. temp → ele = x;
6. temp → next = NULL;
7. Tail = temp;
8. if (Head == NULL)
9. Head = Tail;
10. }
```

- Step 4 allocates memory
- Step 5 read data
- Steps from 6 to 9 adjust reference
- 'if' condition represents first insertion

Insert in middle/random position of list

```

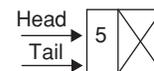
1. void ins _ mid (int n, int pos)
2. {
3. int i = 1;
4. Node * temp, N, P; //N,P represent
   previous //& next nodes
5. if (Head == NULL)
6. {
7. ins _ head(n);
8. }
9. temp = (Node *) malloc(sizeof(Node));
10. temp → ele = n;
11. P = head;
12. while (i < pos -1)
13. {
14. P = P → next;
15. i++;
16. }
17. N = P → next;
18. temp → next = N;
19. P → next = temp;
20. }
```

- step 4 checks, whether the insertion is into an empty list.
- If list is empty, invokes ins-head() function.
- If list is not empty, then step 9 allocates memory.
- Step 10 reads data.
- Steps from 11 to 14 make the reference to the previous and next nodes of new node to be inserted.
- Steps 15 and 16 create the reference to new node from previous node and from new node to next node.

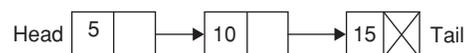
Example 1: Head = Tail = NULL

$n = 5, P = \text{NULL};$

Here the list is empty. So,



Example 2:



Insert element (n) 20 at position(pos) 3.

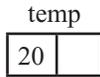
In current list, element 5 is the first element, 10 is the second and 15 is the third element.

To insert an element at $\text{pos} = 3$, the new node has to be placed between elements 10 and 15.

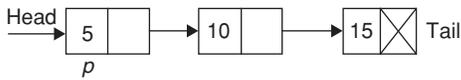
Condition in step 4 is false so step 9 executes and allocates memory.



On completion of step 10 –



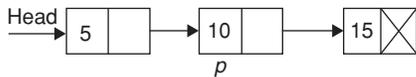
Step 11



Step 12, 13

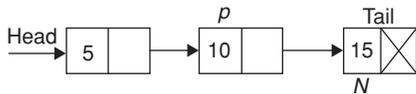
```

While (i < pos - 1)
{
P = P → next;
i++;
}
i < pos
1 < 2
Condition true, so
    
```

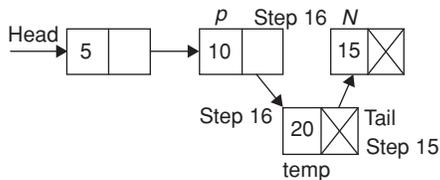


i becomes 2,
2 < 2 // condition false

Step 14 makes a reference to next of previous element.



Steps 15 and 16 execute as follows:



Now the element 20 becomes the 3rd element in the list.

Deletion

- Identify the node
- Adjust the links, such that deallocation of that node does not make the list as unconnected components.
- Return/display element to delete.
- Deallocate memory.

Delete head element

```

1. void del _ head()
2. {
3. int x;
   Node * temp;
4. if (Head == NULL)
    
```

```

5. {
6. printf("List empty");
7. return;
8. }
9. x = Head → ele;
10. temp = Head;
11. if (Head == Tail)
12. Head = Tail = NULL;
13. else
14. Head = Head → next;
15. printf ("Deleted element %d", x);
16. free(temp);
17. }
    
```

- Step 4 – Checks for list empty
- Step 9 – Reads element to delete
- Step 10 – Head referred by temp pointer
- Step 11 – Checks for last deletion
- Step 14 – Moves the head pointer to next element in the list
- Step 15 – Displays element to delete
- Step 16 – Deallocates memory

Delete tail element

```

1. void del _ tail()
2. {
3. int x;
4. Node * temp;
5. if (Head == NULL)
6. {
7. printf("\n list empty")
8. return ;
9. }
10. temp = Head;
11. while(temp → next != Tail)
12. temp = temp → next;
13. x = Tail → ele;
14. Tail = temp;
15. temp = temp → next;
16. Tail → next = NULL;
17. printf("\n Deleted element : %d", x)
18. free (temp);
19. }
    
```

- Step 4 – Checks for list empty
- Step 10, 11, 12 – Move the temp pointer to last but one node of the list
- Step 13 – Reads tail element to delete
- Step 14 – Moves tail pointer to last but one node
- Step 15 – Moves the temp pointer to last node of the list
- Step 16 – Removes the reference from tail node to temp node, i.e., tail node becomes the last element
- Step 17 – Displays elements to delete
- Step 18 – Deallocate memory

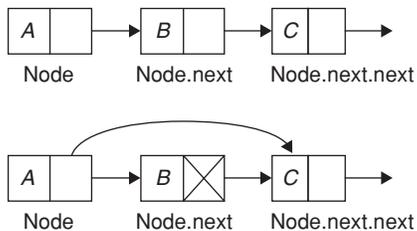
Delete middle element

```

1. void del _ mid (int pos)
2. {
3. int i = 1, x;
4. Node * temp P, N;
5. if(Head == NULL)
6. {
7. printf ("\n list empty")
8. return;
9. }
10. P = head;
11. while (i < pos -1)
12. {
13. P = P → next;
14. i++;
15. }
16. temp = P → next;
17. N = temp → next;
18. P → next = N;
19. x = temp → ele;
20. printf("\n Element to Delete %d", x);
21. free(temp);
22. }

```

- Step 5 – Checks for empty list
 Step 10, 11, 12 – Move previous pointer *P* to previous node of node to delete.
 Step 13 – Temp points to node to delete
 Step 14 – *N* points to temp next
 Step 15 – Creates link from *P* to *N*
 Steps 16, 17, 18 – Read and display elements to delete and deallocate memory.

**Linked list using dynamic variables**

Node in the linked list contains data part that is *ele* and link part which points to the next node, and some other external pointer will be pointing to this as these take some storage, a programmer when creating a list, should check with the available storage. For this we make use of `getnode ()`

Function which is defined as follows:

```

struct node
{
int ele
struct node * next ;
};
typedef struct node Node;
Node getnode ()

```

```

{
Node ptr;
ptr = (Node *) malloc (size of (struct node));
return (ptr);
}

```

If *ptr* returns NULL, then it is underflow (there is no available memory) otherwise, it returns start address of memory location.

Search an element

```

1. void search (int x)
2. {
3. Node * temp = head;
4. int c = 1;
5. while (temp != NULL)
6. {
7. if (temp → ele == x)
8. {
9. printf("\n Element found at % d", c);
10. break;
11. }
12. c++;
13. }
14. if (temp == NULL)
15. printf("\n search unsuccessful");
16. }

```

- Step 7 – Checks temp data with search element. Repeats this step until the element is found or reaches the last node
 Step 9 – Displays the position of search element in the list, if found
 Step 14, 15 – Represents search element not exists in list

Display

```

1. void display ( )
2. {
3. Node *temp = Head;
4. printf("\n list elements: ");
5. while (temp != NULL)
6. {
7. printf("%d", temp → ele);
8. temp = temp → next;
9. }
10. }

```

Step 7 – Displays temp data

Step 8 – Moves temp pointer to next node

Algorithm to reverse direction of all links of singly linked list

Consider a linked list '*L*' with head as pointer pointing to the first node contains data element '*ele*' and a pointer called '*next*' which points to the next node.

Reverse is the routine which will reverse the list, there are three node pointers *P*, *Q*, *R* with *P* pointing to the first node, *Q* pointing to NULL.

```

1. START
2. if (P = NULL)
    1. print ("List is null");
    2. Exit
3. While (P)
4. R = Q;
5. Q = P;
6. P = P → next;
7. Q → next = R
8. End While
9. Head = Q;
10. STOP

```

Double-linked List (DLL)

Double-linked list is a linked list in which, each node contains data part and two link fields.

Node structure:

```

struct Dnode
{
struct Dnode *prev;
int ele;
struct Dnode *next;
};

```

- prev – points to previous node in list
- next – points to next node in list
- The operations which can be performed in SLL can also be performed on DLL.
- The major difference is that we have to adjust double reference as compared to SLL.
- We can traverse or display the list elements in forward as well as in reverse direction.

Example:



Circular-linked List (CLL)

Circular-linked list is completely same as SLL, except, in CLL the last (Tail) node points to first (Head) node of list.

So, the Insertion and Deletion operation at Head and Tail are little different from SLL.

Double Circular-linked List (DCL)

Double circular-linked list can be traversed in both directions again and again. DCL is very similar to DLL, except the last node's next pointer points to first node of list and first node's previous pointer points to last node of list.

So, the insertion and deletion operations at head and tail in DCL are little different in adjusting the reference as compared to DLL.

Storing ordered table as linked list: The table is stored as a linked list, it is retrieved and stored with two pointers, one pointer will point to node holding a record having the smallest key and other pointer performs the search.

Stack

A stack is a last in first out (LIFO) abstract data type and data structure. A stack can have any abstract data type as an element, but is characterized by only two fundamental operations.

√ PUSH

√ POP

- The PUSH operation adds an item to the top of the stack, hiding any items already on the stack or initializing the stack if it is empty.
- The POP operation removes an item from the top of the stack, and returns the popped value to the caller.
- Elements are removed from the stack in the reverse order to the order of their insertion. Therefore, the lower elements are those that have been on the stack for longest period.

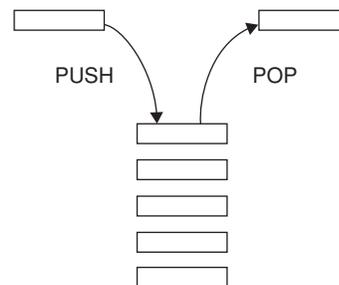


Figure 1 Simple representation of a stack

Implementation

A stack can be easily implemented either through an array or a linked list. The user is only allowed to POP or PUSH items onto the array (or) linked list.

1. **Array Implementation:** Array implementation aims to create an array where the first element inserted is placed `st[0]` which will be deleted last.

The program must keep track of position `top` (last) element of stack.

Operations

Initially `Top = -1`; //represents stack empty

(i) Push (`S, N, TOP, x`)

```

{
if (TOP == N - 1)
printf("overflow");
else
TOP = TOP + 1;
S[TOP] = x;
}

```

(ii) POP (`S, N, TOP, x`)

```

{
if (TOP == -1)
printf("underflow");
else
x = S[TOP]
TOP = TOP - 1
return x;
}

```

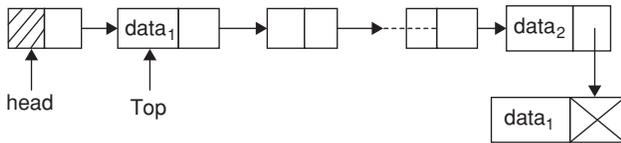
2. Dynamic Implementation: The Array implementation is also called static implementation, because the stack size is fixed.

The stack implementation using linked list is called dynamic implementation, because the stack size can grow and shrink as the elements added or removed from the stack.

- The PUSH operation on stack is same as insert head in SLL.
- The POP operation is same as delete head in SLL.

Algorithm to add and delete to a link stack and link queue

Link stack:



The linked stack with head and top pointers is shown above

The algorithm to push the elements into stack is given below, the method push (item)

Steps:

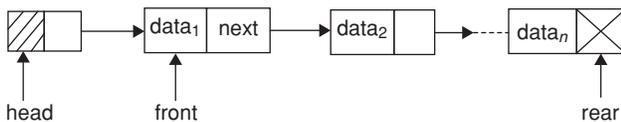
1. ptr = getnode (Node)
2. ptr.data = item
3. ptr.next = Top
4. Top = new
5. Head.next = Top
6. Stop.

for deletion of elements from stack, its algorithm is pop(), it is given below

Steps:

1. if (Top = NULL)
 1. print "stack is empty"
 2. exit
2. Else
 1. ptr = Top.next
 2. item = Top.data
 3. Head.next = ptr
 4. Top = ptr
3. End if
4. Stop.

Linked queue representation



The linked queue with head, front and rear point is shown above.

The algorithm to enqueue the elements into queue is given below, the method enqueue (item)

Steps:

1. ptr = getNode (Node)
2. ptr.data = item
3. ptr.next = NULL
4. if (front = NULL)
 - front = ptr
 - else
 - rear.next = ptr;
5. end if
6. rear = ptr
7. Stop

For deletion of elements from queue that is ptr dequeue () is given below

Steps:

1. if (front = NULL)
 1. print "underflow".
 2. exit
2. ptr = front;
3. front = ptr.next
4. Head.next = front
5. item = ptr.data
6. free(ptr)
7. end.

USES OF STACK

- Function calls: When a function is called all local storage for the function is allocated on system 'stack', and return address also pushed on to system stack.
- Recursion stacks can be used to implement recursion if the programming language does not provide recursion facility.
- Reversing a list
- Parsing: Stacks are used by compilers to check the syntax of program.
- For evaluating expressions.

Expression Notations

Infix expression: Here binary operator comes between the operands.

Postfix expression: Here the binary operator comes after both the operands.

Example: ab+

Prefix expression: Here the binary operator comes before both the operands.

Example: +ab

Infix to postfix conversion

- If operand, output to postfix expression
- If operator, push it onto stack
- In case of parenthesis, when an opening parenthesis is read, it is pushed onto stack and when a closing parenthesis is read, all operators up to the first opening parenthesis must be popped from the stack into the post fix notation.

Example: $(A + (B - C)) * D$

i/p	Postfix notation	Stack
((
A	A	(
+	A	(+)
(A	(+(
B	AB	(+(
-	AB	(+(-
C	ABC	(+(-
)	ABC-	(+
)	ABC-+	-
*	ABC-+	*
D	ABC-+D	*
	ABC-+D*	

Evaluation of postfix expression

We use operand stack for evaluation. Scan the post fix expression,

- When an operand encounters while scanning, push on to stack.
- While scanning post fix expression, if operator found then
 - Pop top two operands from stack
 - Perform the operation on those two operands
 - Push, result on to stack top
- Finally, the stack contains only one value, which represents result of the expression.

Example: $6\ 2\ 3\ +\ -\ 3\ 8\ 2\ /\ +\ 2\ 3\ +$

Symbol	OP1	OP2	Value	Operand stack
6				6
2				6, 2
3				6, 2, 3
+	2	3	5	6, 5
-	6	5		1
3				1, 3
8				1, 3, 8
2				1, 3, 8, 2
/	8	2	4	1, 3, 4
+	3	4	7	1, 7
*	1	7	7	7
2				7, 2
*	7	2	14	14
3				14, 3
+	14	3		17

Result is 17.

Performing add, delete operations on stack (multiple stack)

Let us consider an array whose size is 'max' with multiple stack A, B having top A and top B, push and pop operations on one stack A is given below.

Algorithm for push A(x)

Initially A[Max], top A = -1, top B = MAX;

1. if (top A = top B)
 - a. print "overflow"
 - b. exit
2. top A = top A + 1
3. A[top A] = x
4. stop

Algorithm for pop A(x)

1. if (top A = - 1)
 - a. print "underflow"
 - b. exit
2. y = A[top A]
3. top A = top A - 1
4. return y
5. stop

Algorithm for push B(x)

1. if (top B - 1 = top A)
 - a. print "overflow"
 - b. exit
2. top B = top B - 1
3. A[top B] = x
4. stop

Algorithm for pop B(x)

1. if (top B = max)
 - a. print "underflow"
 - b. exit
2. y = A [top B]
3. top B = top B - 1
4. return y
5. stop

QUEUE

A queue is an ordered collection of items from which items may be deleted at one end (called that front of queue) and into which items may be inserted at the other end (called rear of queue).

Queue is a linear data structure maintains the data in first in-first out (FIFO) order.

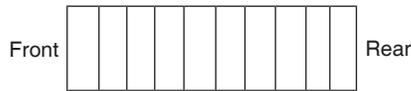
Implementation

Queue can be implemented in the following ways:

1. Array static implementation: queue cannot be extended beyond the array size.
2. Linked list dynamic implementation: Queue size increases as the elements added/inserted to queue. Queue shrinks when an element deleted from queue.

Array Implementation

```
const int SIZE = 10;
int q[SIZE];
int f = -1, r = -1; //f = r = -1 represents queue empty
```



Insertion

```
1. void insert (int x)
2. {
3.   if (r == SIZE - 1)
4.   {
5.     printf("Q FULL")
6.     return;
7.   }
8.   r++;
9.   q[r] = x;
10.  if (f == -1)
11.   f = r;
12. }
```

- Step 3 – Checks for queue full
- Step 8 – Increments rear (*r*)
- Step 9 – Inserts ‘x’ into queue
- Step 10 – Checks whether insertion is first
- Step 11 – If first insertion, updates front (*f*)

Deletion

```
1. void deletion()
2. {
3.   int x;
4.   if (f == -1)
5.   {
6.     printf("\n Q Empty");
7.     return;
8.   }
9.   x = q[f];
10.  if (f == r)
11.   f = r = -1;
12.  else
13.   f++;
14.  printf("\n deleted element %d", x);
15. }
```

- Step 4 – Checks for queue empty
- Step 9 – Deletes ‘q’ front element
- Step 10 – Checks whether queue having only one element
- Step 11 – Rear and front initializes to -1, if queue is having only one element
- Step 13 – Queue front points to next element
- Step 14 – Deleted element is printed

Display

```
1. void display( )
2. {
3.   int i = f;
4.   if (f == -1)
5.   {
```

```
6.   printf("Queue Empty");
7.   return;
8. }
9. printf ("\n Queue Elemetns");
10. for(; i <= r; i++)
11.  printf(" %d", q[i]);
12. }
```

- Step 4 – Checks for ‘q’ empty
- Step 10 and 11 – Display ‘q’ elements

Double-ended Queue

A double-ended queue (deque) is an abstract data structure that implements a queue for which elements can only be added to or removed from the front (head) (or) rear (tail) end.



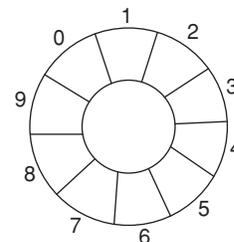
Insertions and deletions are possible at both ends.

Linked List Implementation Double-ended Queue

- Insert – Front is same as insert – Head
- Insert – Rear is same as insert – Tail
- Delete front is same as delete – Head
- Delete – Rear is same as delete – Tail

Circular Queue

As the items from a queue get deleted, the space for that item is reclaimed. Those queue positions continue to be empty. This problem is solved by circular queues. Instead of using a linear approach, a circular queue takes a circular approach; this is why a circular queue does not have a beginning or end.



The advantage of using circular queue over linear queue is efficient usage of memory.

Algorithm to implement addition and deletion from circular queue

Circular Queue Insertion:

To add an element ‘X’ to a Queue ‘Q’ of size ‘N’ with front and rear pointers as ‘F’ and ‘R’ is done with insert (X), Initially $F = R = 0$.

Insert (X)

Steps:

1. if ((R = N) && (F = 1)) or ((R + 1) = F)
 - a. print "overflow"
 - b. exit
2. if (R = N)
 - then R = 0;
 - Else
 - R = R + 1;
3. Q[R] = x;
4. if (F = 0)
 - F = 1
5. Stop.

To delete an element we implement an algorithm delete (). 'y' contains the deleted element.

delete()

Steps:

1. if (F = 0)
 - a. print "underflow"
 - b. exit
2. y = Q[F]
3. if (F = R)
 - F = R = 0
 - else
 - If (F = N)
 - F = 1
 - Else
 - F = F + 1
4. Return y
5. Stop.

Priority Queue

In priority queue, the intrinsic ordering of elements does determine the results of its basic operations.

There are two types of priority queues.

- Ascending priority queue is a collection of items in which items can be inserted arbitrarily and from which only the smallest items can be removed.
- Descending priority queue is similar but allows deletion of the largest item.

Array Implementation

- The insertion operation on priority queue selects the position to the element to insert.
- Makes the position empty/free by moving the existing element (if required).
- Place the element in required position.
- Deletion operation simply deletes front of queue.

Linked-list Implementation

- Insertion operation create a node
- Reads element into node
- Find out the location
- Insert the node into list, by adjusting the reference
- Deletion operation simply deletes head elements, making the head next as head element

Linked-list Implementation of Priority Queue

- Insertion in queue is same as insert-tail of queue
- Deletion from queue is same as delete head

EXERCISES

Practice Problems I

Directions for questions 1 to 16: Select the correct alternative from the given choices.

1. If the array representation of a circular queue contains only one element then

(A) front = rear	(B) front = rear + 1
(C) front = rear - 1	(D) front = rear = NULL
2. The five items *P*, *Q*, *R*, *S* and *T* are pushed in a stack, one after another starting from *P*. The stack is popped four times, and each element is inserted in a queue. The two elements are deleted from the queue and pushed back on the stack. Now one item is popped from the stack. The popped item is _____.

(A) <i>P</i>	(B) <i>Q</i>
(C) <i>R</i>	(D) <i>S</i>
3. What are the contents of the stack (initially the stack is empty) after the following operations?

PUSH (A)
 PUSH (B)
 PUSH (C)
 POP

PUSH(D); POP; POP;
 PUSH(E)
 PUSH(F)
 POP

- | | |
|---------|----------|
| (A) ABE | (B) AE |
| (C) A | (D) ABCE |

4. Consider the below code, which deletes a node from the beginning of a list:

```
void deletefront()
{
  if(head == NULL)
    return;
  else
  {
    .....
    .....
    .....
  }
}
```

Which lines will correctly implement else part of above code?

- (A) if (head → next == NULL)
head = head → next;
- (B) if (head == tail)
head = tail = NULL;
else
head = head → next;
- (C) if (head == tail == NULL)
head = head → next;
- (D) head = head → next;
5. When a new element is inserted in the middle of linked list, then the references of _____ to be adjusted/ updated.
- (A) those nodes that appear after the new node
(B) those nodes that appear before the new node
(C) head and tail nodes
(D) those nodes that appear just before and after the new node

6. The following C function takes double-linked list as an argument. It modifies the list by moving the head (first) element to tail of the list.

```
typedef struct node
{
    struct node *p;
    int data;
    struct node *n;
} Node;
Node * Move - to - last (Node *head)
{
    Node * temp, * prev, * next;
    if (head == NULL) || (head → n == NULL))
        return head;
    temp = head;
    prev = head;
    head = head → n;
    while (prev → n != NULL)
    {
        X;
    }
    Y;
    return head;
}
```

- (A) X: prev = prev → n;
Y: prev → n = temp;
temp → p = prev;
temp → n = NULL;
head → p = NULL;
- (B) X: next = prev → n;
Y: prev → n = temp;
temp → p = prev;
- (C) X: prev = prev → n;
Y: prev → n = temp;
temp → n = NULL;
head → p = NULL;
- (D) X: next = prev → n;
prev = prev → n;
Y: prev → n = Next;

```
next → n = head;
temp → n = NULL;
```

7. Which of the following program segment correctly inserts an element at the front of the linked list. Assume that Node represents linked list node structure, value is the element to be inserted.

- (A) temp = (Node *)malloc (sizeof (Node));
temp → data = value;
temp → next = head;
head = temp;
- (B) temp = (Node *)malloc(sizeof (Node*))
);
temp → data = value;
temp → next = head;
head = temp;
- (C) temp = (Node *)malloc (sizeof (Node));
head = temp;
temp → next = head;
temp → data = value;
- (D) temp = (Node *)malloc (sizeof (Node*))
);
temp → data = value;
head = temp;
temp → next = head;

8. Consider the following program segment:

```
struct element
{
    int x;
    struct element *link;
}
void shuffle(struct element *head)
{
    struct *p, *q;
    int t;
    if (!head || !head → link) return;
    p = head ; q = head → link;
    while(q)
    {
        t = p → x;
        p → x = q → x;
        q → x = t;
        p = q → link;
        q = p? p : 0;
    }
}
```

The function called with list containing 10, 15, 20, 25, 30, 35, 40 in given order. What will the order of elements of the list, after executing the function shuffle?

- (A) 10 15 20 25 30 35 40
(B) 40 35 30 25 20 15 10
(C) 20 15 10 25 40 35 30
(D) 15 10 25 20 35 30 40

9. Primary ADT's are

- (A) Linked list only (B) Stack only
(C) Queue only (D) All of these

10. Linked list uses NULL pointers to signal
 (A) end of list (B) start of list
 (C) Either (A) or (B) (D) Neither (A) nor (B)
11. Which of the following is essential for converting an infix to postfix form efficiently?
 (A) Operator stack (B) Operand stack
 (C) Both (A) and (B) (D) Parse tree
12. Stacks cannot be used to
 (A) Evaluate postfix expression
 (B) Implement recursion
 (C) Convert infix to postfix
 (D) Allocate resource like CPU by the operating system
13. Linked list can be sorted
 (A) By swapping data only
 (B) By swapping address only
 (C) Both (A) and (B)
 (D) None of these
14. Linked list are not suitable for implementing
 (A) Insertion sort
 (B) Binary search
 (C) Radix sort
 (D) Polynomial manipulation
15. Insertion of node in a double-linked list requires how many changes to previous (prev) and next pointers?
 (A) No changes (B) 2 next and 2 prev
 (C) 1 next and 1 prev (D) 3 next and 3 prev
16. Minimum number of stacks required to implement a queue is
 (A) 1 (B) 2
 (C) 3 (D) 4

Practice Problems 2

Directions for questions 1 to 11: Select the correct alternative from the given choices.

1. Stack is useful for implementing _____.
 (A) radix sort
 (B) breadth first search
 (C) quick sort
 (D) recursion
2. Which is true about linked list?
 (A) A linked list is a dynamic data structure.
 (B) A linked list is a static structure.
 (C) A stack cannot be implemented by a linear linked list.
 (D) None of the above
3. The process of accessing the data stored in a tape is similar to manipulating data on a _____.
 (A) stack (B) list
 (C) queue (D) heap
4. Which of the following is used to aid in evaluating a prefix expression?
 (A) Queue (B) Heap
 (C) Stack (D) Hash
5. Select the statement which best completes the sentence — ‘Abstract data type is...’
 (A) a data type which is abstract in nature
 (B) a kind of data type
 (C) data structure
 (D) a mathematical model together with a set of operations defined on it
6. Which of the following data structures may give an overflow error, even through the current number of elements in it is less than its size?
 (A) Simple queue (B) Circular queue
 (C) Stack (D) None of these
7. In a circular linked list, insertion of a record involves the modification of _____.
 (A) no pointer (B) four pointers
 (C) two pointers (D) All of the above
8. Among the following, which one is not the right operation on a stack?
 (A) Remove the item that is inserted latest into the stack.
 (B) Add an item to the stack.
 (C) Remove the first item that is inserted into the stack, without deleting other elements.
 (D) None of the above
9. Among the following which one is not the right operation on dequeue?
 (A) Inserting an element in the middle of a dequeue.
 (B) Inserting an element at the front of a dequeue.
 (C) Inserting an element at the rear of a dequeue.
 (D) None of the above
10. A linear list in which elements can be added or removed at either end but not in the middle is _____.
 (A) queue
 (B) dequeue
 (C) array
 (D) tree
11. The post fix notation of $A/B ** C + D * E - A * C$ is
 (A) $ABC ** /DE * + AC * -$
 (B) $ABC ** D/E * + AC + -$
 (C) $ABC ** /DE * AC + -$
 (D) $ABC ** /DE * + AC + -$

PREVIOUS YEARS' QUESTIONS

1. An abstract data type (ADT) is **[2005]**
 (A) same as an abstract class.
 (B) a data type that cannot be instantiated.
 (C) a data type for which only the operations defined on it can be used, but none else.
 (D) All of the above

2. An implementation of a queue Q , using two stacks S_1 and S_2 , is given below:

```
void insert (Q, x) {
    push (S1, x);
}
void delete (Q) {
    if (stack-empty (S2)) then
        if (stack-empty (S1)) then {
            print ("Q is empty");
            return;
        }
    else while (!(stack-empty (S1)))
    {
        x = pop (S1);
        push (S2, x);
    }
    x = pop (S2);
}
```

Let n insert and m (n) delete operations be performed in an arbitrary order on an empty queue Q . Let x and y be the number of push and pop operations performed respectively in the process. Which one of the following is true for all m and n ? **[2006]**

- (A) $n + m \leq x < 2n$ and $2m \leq y < n + m$
 (B) $n + m \leq x < 2n$ and $2m \leq y < 2n$
 (C) $2m \leq x < 2n$ and $2m \leq y < n + m$
 (D) $2m \leq x < 2n$ and $2m \leq y < 2n$
3. The following postfix expression with single digit operands is evaluated using a stack:
 $8\ 2\ 3\ \wedge\ / \ 2\ 3\ * \ + \ 5\ 1\ * \ -$
 Note that \wedge is the exponentiation operator. The top two elements of the stack after the first $*$ is evaluated are: **[2007]**
 (A) 6 and 1 (B) 5 and 7
 (C) 3 and 2 (D) 1 and 5
4. The following C function takes a single-linked list of integers as a parameter and rearranges the elements of the list. The function is called with the list containing the integers 1, 2, 3, 4, 5, 6, 7 in the given order. What will be the contents of the list after the function completes execution?

```
struct node {
    int value;
    struct node *next;
};
void rearrange (struct node *list) {
    struct node *p, *q;
```

```
int temp;
if (!list || !list -> next) return;
p = list; q = list -> next;
while (q) {
    temp = p -> value; p -> value = q ->
    value;
    q -> value = temp; p = q -> next;
    q = p? p -> next : 0;
}
}
```

[2008]

- (A) 1, 2, 3, 4, 5, 6, 7 (B) 2, 1, 4, 3, 6, 5, 7
 (C) 1, 3, 2, 5, 4, 7, 6 (D) 2, 3, 4, 5, 6, 7, 1

5. Suppose a circular queue of capacity $(n - 1)$ elements is implemented with an array of n elements. Assume that the insertion and deletion operations are carried out using REAR and FRONT as array index variables, respectively. Initially, REAR = FRONT = 0. The conditions to detect queue full and queue empty are **[2012]**

- (A) Full: $(\text{REAR} + 1) \bmod n == \text{FRONT}$
 Empty: REAR == FRONT
 (B) Full: $(\text{REAR} + 1) \bmod n == \text{FRONT}$
 Empty: $(\text{FRONT} + 1) \bmod n == \text{REAR}$
 (C) Full: REAR == FRONT
 Empty: $(\text{REAR} + 1) \bmod n == \text{FRONT}$
 (D) Full: $(\text{FRONT} + 1) \bmod n == \text{REAR}$
 Empty: REAR == FRONT

6. Consider the C program below **[2015]**

```
#include <stdio.h>
int *A, stkTop;
int stkFunc (int opcode, int val)
{
    static int size=0, stkTop=0;
    switch (opcode) {
        case -1: size = val; break;
        case 0: if (stkTop < size)
            A[stkTop++] =
                val; break;
        default: if (stkTop) return A[
            --stkTop];
    }
    return -1;
}
int main ( )
{
    int B[20]; A = B; stkTop = -1;
    stkFunc (-1, 10);
    stkFunc (0, 5);
    stkFunc (0, 10);
    printf ("%d\n", stkFunc(1, 0) +
        stkFunc(1, 0));
}
```

The value printed by the above program is _____

7. The result of evaluating the postfix expression $10\ 5\ +\ 60\ 6\ * \ 8\ -$ is [2015]

- (A) 284 (B) 213
(C) 142 (D) 71

8. Let Q denote a queue containing sixteen numbers and S be an empty stack.

$Head(Q)$ returns the element at the head of the queue Q **without** removing it from Q . Similarly $Top(S)$ returns the element at the top of S **without** removing it from S .

Consider the algorithm given below.

```

while  $Q$  is not Empty do
  if  $S$  is Empty OR  $Top(S) \leq Head(Q)$ 
  then
     $x := Dequeue(Q)$ 
    Push( $S, x$ );
  else
     $x := Pop(S)$ ;
    enqueue( $Q, x$ );
  end
end

```

The maximum possible number of iterations of the while loop in the algorithm is _____. [2016]

9. The attributes of three arithmetic operators in some programming language are given below.

Operator	Precedence	Associativity	Arity
+	High	Left	Binary
-	Medium	Right	Binary
*	Low	Left	Binary

The value of the expression

$2 - 5 + 1 - 7 * 3$ in this language is _____. [2016]

10. A circular queue has been implemented using a singly linked list where each node consists of a value and a single pointer pointing to the next node. We maintain exactly two external pointers **FRONT** and **REAR** pointing to the front node and the rear node of the queue, respectively. Which of the following statements is/are CORRECT for such a circular queue, so that insertion and deletion operations can be performed in $O(1)$ time?

- I. Next pointer of front node points to the rear node.
II. Next pointer of rear node points to the front node.

[2017]

- (A) I only (B) II only
(C) Both I and II (D) Neither I nor II

ANSWER KEYS

EXERCISES

Practice Problems 1

1. A 2. C 3. B 4. B 5. D 6. A 7. A 8. D 9. D 10. A
11. A 12. D 13. C 14. B 15. B 16. B

Practice Problems 2

1. D 2. A 3. C 4. C 5. D 6. A 7. C 8. C 9. A 10. B
11. A

Previous Years' Questions

1. C 2. A 3. A 4. B 5. 6. 15 7. C 8. 256 9. 9 10. B

Trees

LEARNING OBJECTIVES

- Tree
- 2-Tree
- Binary tree
- Properties of binary trees
- Complete binary tree
- Full binary tree
- Binary tree representation
- Linked representation
- Binary search tree
- Binary tree traversing methods
- AVL tree
- Binary heap
- Max-heap
- Min-heap
- Expression tree

TREE

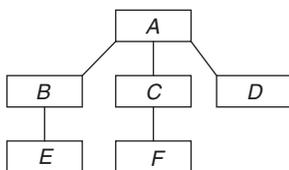
Tree is non-linear data structure designated at a special node called root and elements are arranged in levels without containing cycles.

(or)

The tree is

1. Rooted at one vertex
2. Contains no cycles
3. There is a sequence of edges from any vertex to any other
4. Any number of elements may connect to any node (including root)
5. A unique path traverses from root to any node of tree
6. Tree stores data in hierarchical manner
7. The elements are arranged in layers

Example:



- Root node is *A*.
- *A*'s children are *B*, *C* and *D*.
- *E*, *F* and *D* are leaves.
- Nodes *B*, *C* are called as intermediate nodes.
- *A* is parent of *B*, *C* and *D*.

- *B* is parent of *E* and *C* is parent of *F*.
- Number of children of a node is called degree of node.

2-TREE

A tree in which every node contains either 0 or 2 children.

BINARY TREE

It is a special type of tree where each node of tree contains either 0 or 1 or 2 children.

(or)

Binary Tree is either empty, or it consists of a root with two binary trees called left-sub tree and right sub-tree of root (left or right or both the sub trees may be empty).

Properties of binary tree

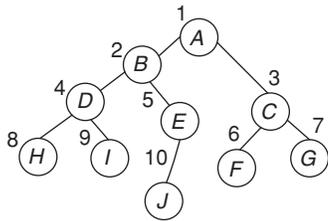
- Binary tree partitioned into three parts.
- First subset contains root of tree.
- Second subset is called left subtree.
- Another subset is called right subtree.
- Each subtree is a binary tree.
- Degree of any node is 0/1/2.
- The maximum number of nodes in a tree with height '*h*' is $2^{h+1} - 1$.
- The maximum number of nodes at level '*i*' is 2^{i-1} .
- For any non-empty binary tree, the number of terminal nodes with n_2 , nodes of degree 2 is $N_0 = n_2 + 1$
- The maximum number of nodes in a tree with depth *d* is $2^d - 1$.

Types of binary tree

Complete binary tree It is a binary tree, in which at every level, except possibly the last, is completely filled and all nodes at the last level are as left as possible.

Example:

Level	Height	Depth
1	3	1
2	2	2
3	1	3
4	0	4



For the given tree:

- Having 4 levels
- Height of the tree is 3
- Depth of the tree is 4
- The numbers at each node represents level order index.
- The level order Index, are assigned to nodes in the following manner
- Root of the tree is '1'
- For a node 'x', the LOI is (2 * LOI (parent)), if 'x' is left child of its parent.
- For a node 'y', the LOI (2 * LOI (Parent) + 1), if 'y' is right child of its parent.

Now complete binary tree can be defined as a binary tree, which contains a sequence of numbers to its nodes as LOI's without any break in sequence.

Full binary tree It is a binary tree, for which all leaf nodes are at same level and all intermediate nodes contains exactly 2 children.

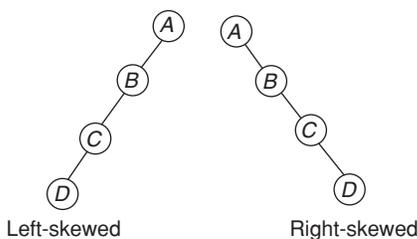
(or)

A tree with depth 'K' contains exactly $2^K - 1$ nodes.

Strictly binary tree A binary tree in which every node contains exactly 0 or 2 children.

Skewed binary tree A binary tree in which elements are added only in one direction.

Example:



Application

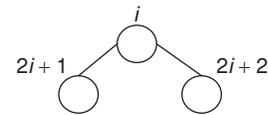
- A binary tree is useful data structure when two way decisions must be made at each point of process.

Binary tree representation

The binary trees can be represented in two ways.

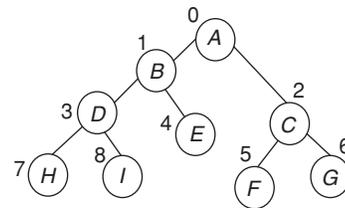
- Array
- Linked list

Array representation The elements of a binary tree are placed in an array using the level order index of each element.



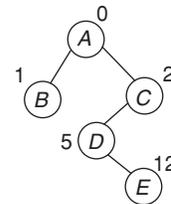
When LOI of Root is 0:

Example 1:



0	1	2	3	4	5	6	7	8
A	B	C	D	E	F	G	H	I

Example 2:

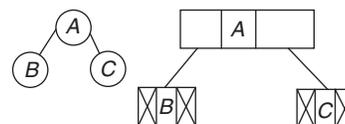


0	1	2	3	4	5	6	7	8	9	10	11	12
A	B	C			D							E

Linked representation Each node contains one data field and two link fields. First link point to the left child and another point to the right child.

In absence of any child, corresponding link field contains NULL.

Example:



Trade-off's between array and linked, representations

- Array representation is somewhat simpler. It must ensure elements are placed in array at proper position.
- Linked representation requires pointer to its left and right child.
- Array representation saves memory for almost complete binary trees.
- Linked representation allocates the number and nodes equal to the number of elements in tree.
- Array representation does not work efficiently for skewed binary trees.
- Array representation limits the size of binary tree to the array size.
- In linked representation, tree can be extended by adding an element dynamically and can be shrunk by deleting an element dynamically.

Binary search tree

It is a special type of binary tree that satisfies the following properties.

- All the elements of left sub tree of root are smaller than root.
- All the elements of right sub tree of root are greater than root.
- The above two properties satisfy for each subtree.

Example:

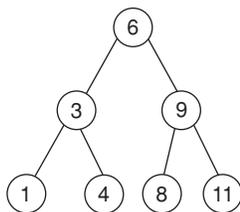


Figure 1 A data structure to encode binary search tree

The binary search tree node contains three fields, data field, left child, right child. Left child is a pointer which points to the predecessor of the node and right child is a pointer which points to the successor of the node.

A data structure to encode binary search tree is

Left child	Data	Right child
------------	------	-------------

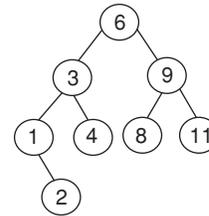
The declaration is

```
Struct node
{
Struct node * left child;
Int data;
Struct node * Right child;
};
```

Insertion If a value to be inserted is smaller than the root, value, it must go in the left subtree, if larger it must go in the right subtree. This reasoning applies recursively until we

reach a node where the required subtree does not exist and that is where we place the new value.

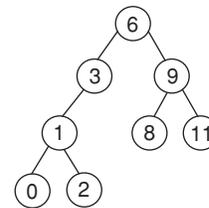
Example: It must go in 6's left subtree, 3's left subtree, 1's right subtree, 1 has no right subtree, so we make a singleton with 2 and it becomes 1's right subtree.



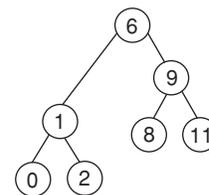
Deletion:

1. If a leaf node has to be deleted, just delete it and the rest of the tree is exactly as it was, so it is still a BST.
2. Suppose the node we are deleting has only one sub tree

Example, In the following tree, '3' has only one sub-tree

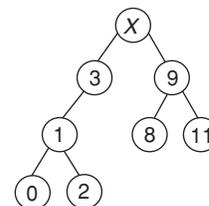


To delete a node with 1 subtree, we just 'link past' the node, i.e., connect the parent of the node directly to the node's only subtree. This always works, whether the one subtree is on the left or on the right. Deleting 3 gives us.



3. Deletion of node which has 2 subtrees

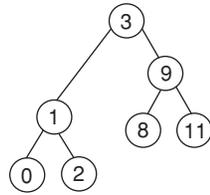
Example: Delete 6.



Choose value 'X'

1. Everything in the left subtree must be smaller than X.
2. Everything in the right subtree must be bigger than X.

We must choose X to be the largest value in the left subtree. In our example, 3 is the largest value in the left subtree. So we replace root node 6 with 3.



Note: We could do the same thing with the right subtree. Just use the smallest value in the right subtree.

Notes:

- The largest element in left subtree is the right most element.
- The smallest element in right subtree is the left most element.

Binary tree traversing methods

The binary tree contains 3 parts:

- V – root
- L – Left subtree
- R – Right subtree

Pre-order: (V, L, R)

- Visit root of the tree first
- Traverse the left - subtree in pre-order
- Traverse the right - subtree in preorder

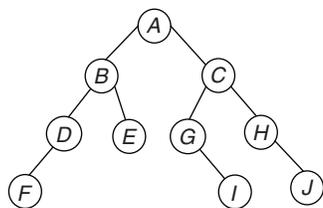
In-order: (L, V, R)

- Traverse the left – subtree in in-order
- Visit Root of the tree
- Traverse right - sub tree in in-order

Post-order: (L, R, V)

- Traverse the left subtree in post-order.
- Traverse the Right - subtree in post-order
- Visit root of the tree

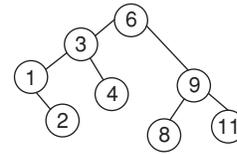
Example 1:



Pre-order: $ABDFECGIHJ$
 In-order: $FDBEAGICHJ$
 Post-order: $FDEBIGJHCA$

Pre-order, In-order and post-order uniquely identify the tree.

Example 2:



Pre-order: 6 3 1 2 4 9 8 11
 In-order: 1 2 3 4 6 8 9 11
 Post-order: 2 1 4 3 8 11 9 6

Points to remember

- Pre-order traversal contains root element as first element in traverse list.
- Post-order traversal contains root element as last in traversal list.
- For BST, in-order traversal is a sorted list.
- A unique binary tree can constructed if either pre-order or post-order traversal list provided with In order traversal list.
- If either pre-order or post-order only given then BST cannot be constructed.

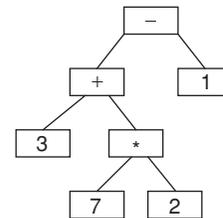
Applications

1. Binary trees can represent arithmetic expressions.
 - An infix expression will have a parent operator and two children operands.

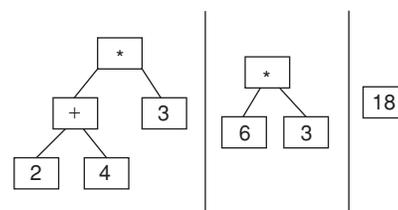
Consider the expression $((3 + (7 * 2)) - 1)$

Each parenthesised expression becomes a tree.

Each operand is a leaf, each operator is an internal node.

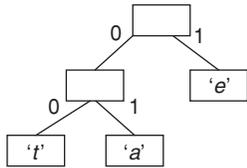


2. To evaluate the expression tree:
 - Take any two leaves
 - Apply the parents operator to them
 - Replace the operator with the value of the sub expression.



3. Binary trees in a famous file compression algorithm Huffman coding tree

- Each character is stored in a leaf
- The code is found by following the path 0 go left, 1 go right.
- *a* is 01
- *e* is 1



AVL Tree

An AVL tree is a self-balancing binary search tree, in which the heights of the two child subtrees of any node differ by at most one.

Insertions and deletions may require the tree to be rebalanced by one or more tree rotations.

- The balance factor of a node is the height of its left subtree minus the height of its right subtree (sometimes opposite) and a node with balance factor -1 , 0 or $+1$ is considered balanced. A node with any other balance factor is considered unbalanced and requires rebalancing the tree.
- The balance factor is either stored directly at each node or computed from the heights of the subtrees.

Insert operations

Step I: Insert a node into the AVL tree as it is inserted in a BST.

Step II: Examine the search path to see if there is a pivot node.

Three cases may arise

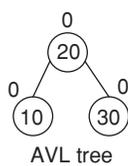
Case I: There is no pivot node. No adjustment required.

Case II: The pivot node exists and the subtree of the pivot node to which the new node is added has smaller height. No adjustment required.

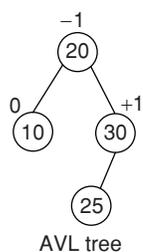
Case III: The pivot node exists and the subtree to which the new node is added has the larger height, Adjustment required.

Example: The numbers at each node represents balance factor.

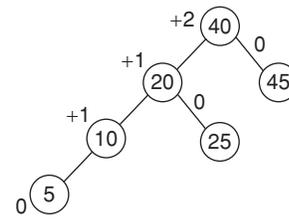
Example 1:



Example 2:



Example 3:



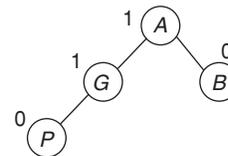
Not an AVL tree

Example 3 is not an AVL tree, because the balance factor of root node is $+2$.

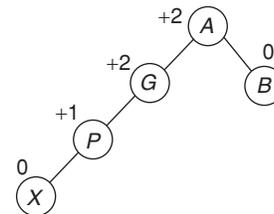
AVL tree becomes height in-balanced tree in following cases:

1. Left-Left case: An insertion in left subtree of left child of pivot node.

Example:

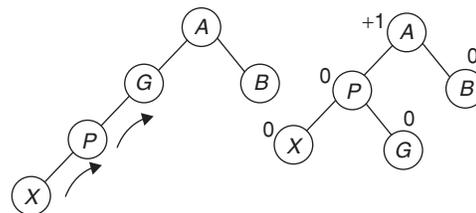


Insert '*X*' as left to node '*P*'. Here '*G*' is pivot node.



Solution:

To make the tree as balanced tree, perform **Left-Left Rotation** as follows:

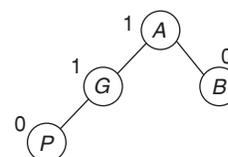


In left-left rotation

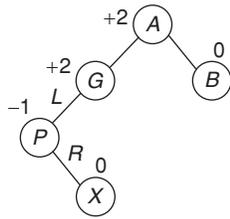
- Intermediate node '*P*' becomes root of subtree.
- Root of subtree '*G*' (pivot) becomes right subtree.
- New node '*X*' remains same as left child of '*P*'.

Left-Right Case

An insertion of left subtree of right child of pivot node.



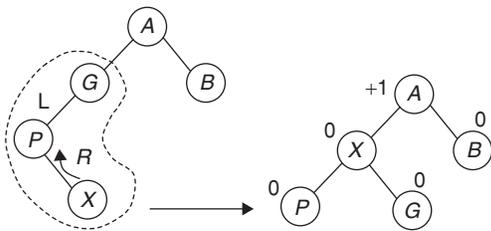
Example 1: Insert 'X' as right child of 'P'.



Is not an AVL tree. Height in-balance at node 'G'.

Solution:

Perform Left-Right Rotation, to balance the height of tree.



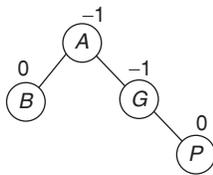
In Left-Right rotation:

- New node 'X' becomes root of subtree.
- Root of subtree 'G' (pivot) becomes right child of 'X'.
- Intermediate node 'P' becomes left child of new node.

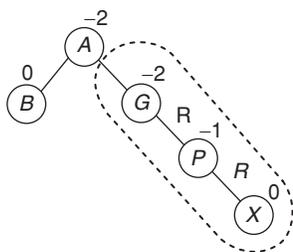
Right-Right case

An insertion of right subtree of right child of pivot node.

Example:



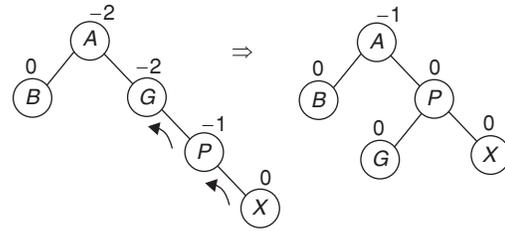
Insert 'X' as right child of 'P'



Is not an AVL tree, because of height in-balance at node 'G'.

Solution:

To make the tree as balanced tree, perform the right-right rotation as follows:

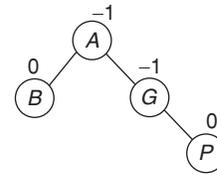


In Right-Right rotation:

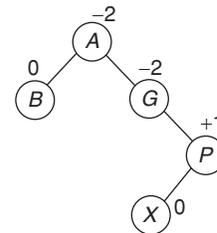
- Intermediate node 'P' becomes root of subtree.
- Root of subtree 'G' (pivot) becomes left child of 'P'.
- New node 'X' remains as right child to 'P'.

Right-Left case

An insertion of right subtree of left child of pivot node.



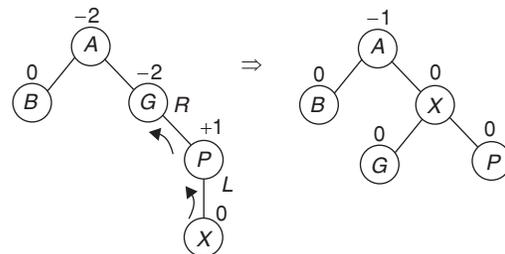
Insert 'X' as left child of 'P'



Is not AVL tree, because height in-balance at node 'G'.

Solution:

To make the above tree as balanced, perform Right-Left rotation as follows:



In Right-Left Rotation:

- New node 'X' becomes root of subtree.
- Root of subtree 'G' (pivot) becomes left child of 'X'.
- Intermediate node 'P' becomes right of 'X'.

Note: Left-Right and Right-Left rotation are also called as double rotations.

BINARY HEAP

A binary heap is a heap data structure created using a binary tree. It can be seen as a binary tree with two additional constraints.

The shape property: The tree is a complete binary tree; that is, all levels of the tree, except possibly the last one (deepest) level of the tree is not complete, the nodes of that level are filled, from left to right.

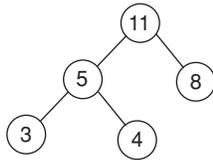
Max-Heap

A heap in which each node is greater than or equal to its children is called max-heap. Max-Heap generally used for heap sort.

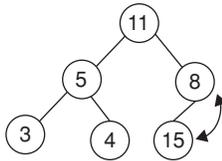
Min-Heap

A heap in which, each node is smaller than or equal to its children is called Min-Heap. Min-heap generally used to implement priority queue.

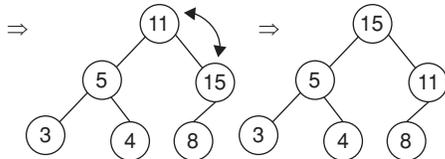
Note: By default heap represent Max-Heap:



Insert 15:

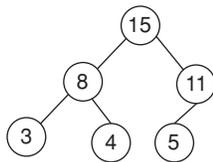


Is not satisfying heap property. So Heapify

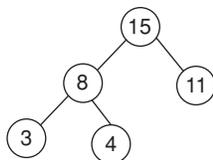


Delete 5: Deletion of a node from heap is always deletes a leaf node.

So interchange the value of last leaf node with node 5.



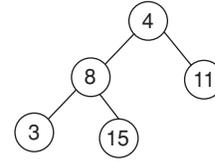
Now delete node '5'



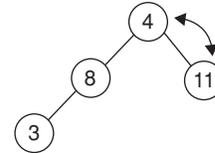
Is satisfying heap property.

Delete 15:

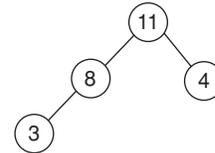
Interchange 4 and 15



Now delete Node '15'



Is not satisfying heap property. So heapify



Note: Insertion or deletion operation on a heap may require heapify process.

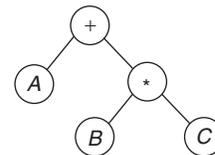
Expression Tree

The expressions can also represented by using a binary tree called expression tree.

Expression tree contains:

- Operators as intermediate nodes.
- Operands as leaf nodes (or) childs to operator nodes.
- The operator at lowest level will be having highest priority.

Example: $A + B * C$



Traversing:

Pre-order: $+ A * B C$

In-order: $A + B * C$

Post-order: $A B C * +$

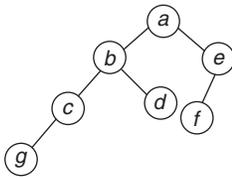
Note: In-order traversal of expression tree generates In-fix expression. Similarly pre-order and post-order generates prefix and postfix, respectively.

EXERCISES

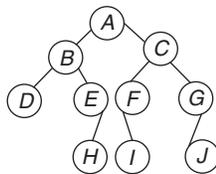
Practice Problems I

Directions for questions 1 to 15: Select the correct alternative from the given choices.

- A binary tree T has n leaf nodes. The number of nodes of degree two in T is _____.
 (A) n (B) $n - 1$
 (C) $\log n$ (D) $n + 1$
- How many numbers of binary tree can be created with 3 nodes which when traversed in post-order gives the sequence C, B, A ?
 (A) 3 (B) 5
 (C) 8 (D) 15
- A binary search tree contains the values 3, 6, 10, 22, 25, 30, 60, 75. The tree is traversed in pre-order and the values are printed out. Which of the following sequence is a valid output?
 (A) 25 6 3 10 22 60 30 75
 (B) 25 6 10 3 22 75 30 60
 (C) 25 6 75 60 30 3 10 22
 (D) 75 30 60 22 10 3 6 25
- Figure shows a balanced tree. How many nodes will become unbalanced when a node is inserted as a child of the node 'g'?



- (A) 7 (B) 2
 (C) 3 (D) 8
- A full binary tree with n non-leaf nodes contains
 (A) $2n$ nodes (B) $\log_2 n$ node
 (C) $n + 1$ nodes (D) $2n + 1$ nodes
 - Which of the following list of nodes corresponds to a post order traversal of the binary tree shown below?



- (A) $ABCDEFGHIJ$ (B) $JIHGFEDCBA$
 (C) $DHEBIFJGCA$ (D) $DEHFIGJBCA$
- Which of the following sequence of array elements forms as heap?

- (A) {23, 17, 14, 6, 13, 10, 1, 12, 7, 5}
 (B) {23, 17, 14, 6, 13, 10, 1, 5, 7, 12}
 (C) {23, 17, 14, 7, 13, 10, 1, 5, 6, 12}
 (D) {23, 17, 14, 7, 13, 10, 1, 12, 5, 6}

- What is the maximum height of any AVL tree with 7 nodes? Assume that the height of a tree with a single node is 0.
 (A) 2 (B) 3
 (C) 4 (D) 5
- A binary search tree is generated by inserting in order the following integers:
 55, 15, 65, 5, 25, 59, 90, 2, 7, 35, 60, 23.
 The number of nodes in the left subtree and right subtree of the root respectively are
 (A) 8, 3 (B) 7, 4
 (C) 3, 8 (D) 4, 7
- In a complete binary tree of n nodes, how far are the most distant two nodes? Assume each in the path counts as 1.
 (A) about $\log_2 n$ (B) about $2\log_2 n$
 (C) about $3\log_2 n$ (D) about $4\log_2 n$
- A complete binary tree of level 5 has how many nodes?
 (A) 20 (B) 63
 (C) 30 (D) 73

Common data for questions 12 and 13: A 3-ary max-heap is like a binary max-heap, but instead of 2 children, nodes have 3 children. A 3-ary heap can be represented by an array as follows:

The root is stored in the first location, $a[0]$, nodes in the next level from left to right is stored from $a[1]$ to $a[3]$ and so on. An item x can be inserted into a 3-ary heap containing n items by placing x in the location $a[n]$ and pushing it up the tree to satisfy the heap property.

- Which one of the following is a valid sequence of elements in an array representing 3-ary max-heap?
 (A) 1, 3, 5, 6, 8, 9 (B) 9, 6, 3, 1, 8, 5
 (C) 9, 3, 6, 8, 5, 1 (D) 9, 5, 6, 8, 3, 1
- Suppose the elements 7, 2, 10 and 4 are inserted, in that order, into the valid 3-ary max-heap found in the above question. Which one of the following is the sequence of items in the array representing the resultant heap?
 (A) 10, 7, 9, 8, 3, 1, 5, 2, 6, 4
 (B) 10, 9, 8, 7, 6, 5, 4, 3, 2, 1
 (C) 10, 9, 4, 5, 7, 6, 8, 2, 1, 3
 (D) 10, 8, 6, 9, 7, 2, 3, 4, 1, 5
- Consider the nested representation of binary trees : $(X Y Z)$ indicated Y and Z are the left and right subtrees respectively, of node X (Y and Z may be null (or) further nested) which of the following represents a valid binary tree?

- (A) (1 2 (4 5 6 7)) (B) (1(2 3 4)5 6)7
 (C) (1(2 3 4) (5 6 7)) (D) (1(2 3 NULL)(4 5))

15. A scheme for storing binary trees in an array X is as follows:

Indexing of X starts at 1 instead of 0. The root is stored at $X[1]$. For a node stored at $X[i]$, the left child, if any,

is stored in $X[2i]$ and the right child, if any, in $X[2i + 1]$. To store any binary tree on ' n ' vertices the minimum size of X should be

- (A) $2n$ (B) n
 (C) $3n$ (D) n^2

Practice Problems 2

Directions for questions 1 to 15: Select the correct alternative from the given choices.

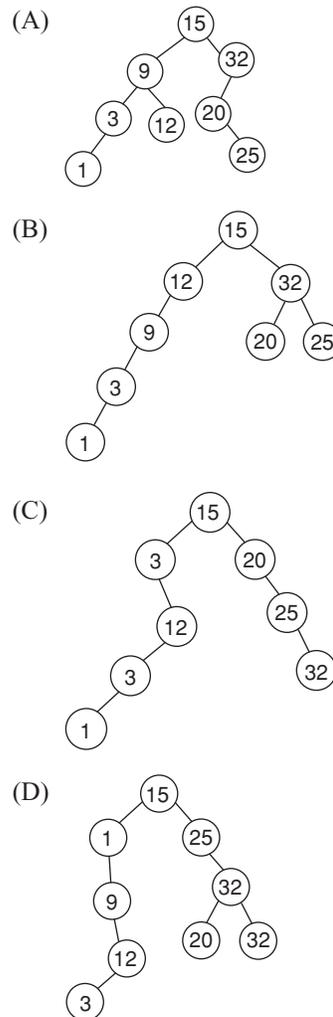
- A binary search tree contains the values 1, 2, 3, 4, 5, 6, 7, 8. The tree is traversed in pre-order and the values are printed out. Which of the following is a valid output?
 (A) 53124786 (B) 53126487
 (C) 53241678 (D) 53124768
- A binary search tree is generated by inserting in order the following integers : 50, 15, 62, 5, 20, 58, 91, 3, 8, 37, 60, 24. The number of nodes in the left subtree and right subtree of the root respectively are:
 (A) (4, 7) (B) (7, 4)
 (C) (8, 3) (D) (3, 8)
- A full binary tree (with root at level 0) of height h has a total number of nodes equal to:
 (A) 2^h (B) $2^{h+1} - 1$
 (C) $2^h - 1$ (D) 2^{h-1}
- The number of null pointers of a binary tree of n nodes is :
 (A) $n + 1$ (B) $n(n + 1)$
 (C) n^2 (D) $2n$
- Which of the following is false?
 (A) A tree with n nodes has $(n - 1)$ edges.
 (B) A labeled rooted binary tree can be uniquely constructed, given its post-order, in-order traversal results.
 (C) The complete binary tree with n internal nodes has $(n + 1)$ leaves.
 (D) The maximum number of nodes in a binary tree of height h is $(2^{h+1} - 1)$.
- The maximum number of nodes in a binary tree at level i is
 (A) 2^i (B) $2^i - 1$
 (C) $2^i + 1$ (D) $\log_2 i + 1$
- The number of leaf nodes in a rooted tree of n nodes, with each node having 0 or 3 children is
 (A) $\frac{n}{3}$ (B) $\frac{(n-1)}{3}$
 (C) $\frac{(n-1)}{2}$ (D) $\frac{(2n+1)}{3}$

8. A complete n -ary tree is one in which every node has 0 or n children. If x is the number of internal nodes of a complete n -ary tree, the number of leaves in it is given by

- (A) $x(n - 1) + 1$
 (B) $xn + 1$
 (C) $xn - 1$
 (D) $x(n + 1) - 1$

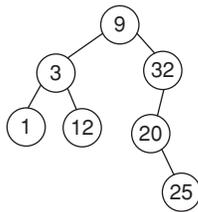
Common data for questions 9 and 10:

9. Insert the keys into a binary search tree in the order specified 15, 32, 20, 9, 3, 25, 12, 1. Which one of the following is the binary search tree after insertion of all elements?

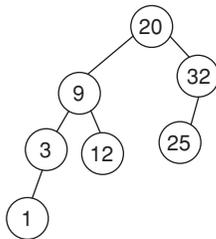


10. Which of the following is the binary tree after deleting 15?

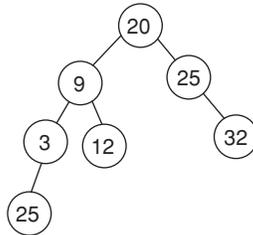
(A)



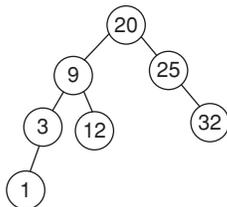
(B)



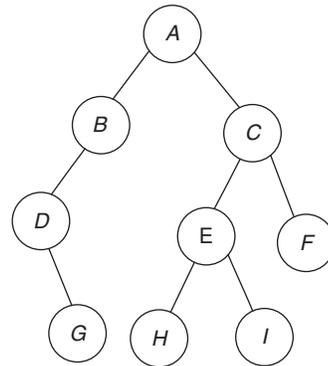
(C)



(D)



For questions 11, 12 and 13 below, use this figure



11. What is the post-order expression?

- (A) ABDGCEJHIF (B) GDBHIEFCA
(C) DGBAHEICF (D) ABHIEFCDG

12. What is the pre-order expression?

- (A) ABDGCEHIF (B) ABHIEFCDG
(C) DGBAHEIFCF (D) GDBHIEFCA

13. What is the in-order expression?

- (A) ABDGCEHIF (B) GDBHIEFCA
(C) DGBAHEICF (D) ABHIEFCDG

14. In a 3-ary tree every internal node has exactly 3 children. The number of leaf nodes in such a tree with 6 internal nodes will be

- (A) 13 (B) 12 (C) 11 (D) 10

15. Minimum number of swaps needed to convert the array 89, 19, 14, 40, 17, 12, 10, 2, 5, 7, 11, 6, 9, 70 into a max heap

- (A) 2 (B) 3 (C) 1 (D) 0

PREVIOUS YEARS' QUESTIONS

1. In a binary tree with n nodes, every node has an odd number of descendants. Every node is considered to be its own descendant. What is the number of nodes in the tree that have exactly one child? [2010]

- (A) 0 (B) 1
(C) $(n-1)/2$ (D) $n-1$

2. The following C function takes a singly-linked list as input argument. It modifies the list by moving the last element to the front of the list and returns the modified list. Some part of the code is left blank.

```

typedef struct node {
    int value;
    struct node *next;
} Node;
Node *move_to_front(Node *head) {
    Node *p, *q;
    if ((head == NULL || (head->next ==
    NULL)) return head;
  
```

```

    q = NULL; p = head;
    while (p->next !=NULL) {
        q = p;
        p = p->next;
    }
  
```

```

    return head;
}
  
```

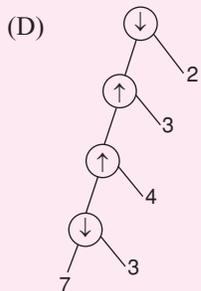
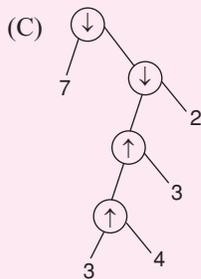
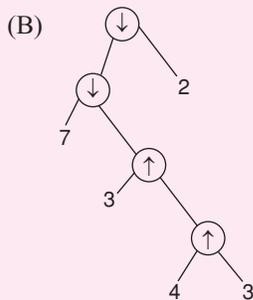
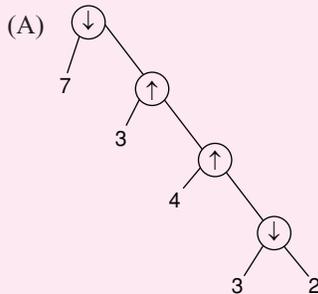
Choose the correct alternative to replace the blank line. [2010]

- (A) $q = \text{NULL}; p \rightarrow \text{next} = \text{head}; \text{head} = p;$
(B) $q \rightarrow \text{next} = \text{NULL}; \text{head} = p; p \rightarrow \text{next} = \text{head};$
(C) $\text{head} = p; p \rightarrow \text{next} = q; q \rightarrow \text{next} = \text{NULL};$
(D) $q \rightarrow \text{next} = \text{NULL}; p \rightarrow \text{next} = \text{head}; \text{head} = p;$

3. Consider two binary operators ' \uparrow ' and ' \downarrow ' with the

precedence of operator \downarrow being lower than that of the operator \uparrow . Operator \uparrow is right associative while operator \downarrow is left associative. Which one of the following represents the parse tree for expression $(7 \downarrow 3 \uparrow 4 \uparrow 3 \downarrow 2)$?

[2011]



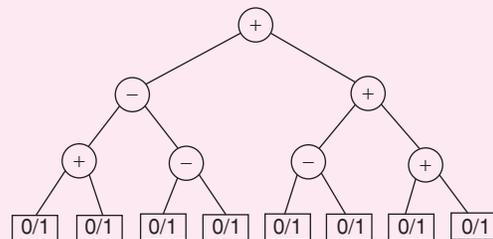
4. The height of a tree is defined as the number of edges on the longest path in the tree. The function shown in the pseudocode below is invoked as “height(root)” to compute the height of a binary tree rooted at the tree pointer “root”.

```
int height(treeptr n)
{if (n == NULL) return -1;
if (n → left ==NULL)
if (n → right == NULL) return 0;
else return B1 ; //Box 1
else {h1 = height (n → left);
if (n → right== NULL) return (1 + h1);
else {h2 = height (n → right);
return B2 ; //Box 2
}
}
}
```

The appropriate expressions for the two boxes B_1 and B_2 are [2012]

- (A) $B_1: (1 + \text{height}(n \rightarrow \text{right}))$
 $B_2: (1 + \max(h_1, h_2))$
- (B) $B_1: (\text{height}(n \rightarrow \text{right}))$
 $B_2: (1 + \max(h_1, h_2))$
- (C) $B_1: \text{height}(n \rightarrow \text{right})$
 $B_2: \max(h_1, h_2)$
- (D) $B_1: (1 + \text{height}(n \rightarrow \text{right}))$
 $B_2: \max(h_1, h_2)$

5. Consider the expression tree shown. Each leaf represents a numerical value, which can either be 0 or 1. Over all possible choices of the values at the leaves, the maximum possible value of expression represented by the tree is _____. [2014]



6. Consider the pseudocode given below. The function **Dosomething()** takes as argument a pointer to the root of an arbitrary tree represented by the *leftMostChild-rightSibling* representation. Each node of the tree is of type **treeNode**. [2014]

```
type def struct treeNode* treeptr;
struct treeNode
{
treeptr leftMostChild, rightSibling;
};
int Dosomething (treeptr tree)
```

```

{
int value = 0;
if (tree != NULL) {
if (tree -> leftMostChild == NULL)
value = 1;
else
value = Dosomething (tree -> leftMostChild);
value = value + Dosomething (tree -> right
Sibling);
}
return (value);
}

```

When the pointer to the root of a tree is passed as the argument to **DoSomething**, the value returned by the function corresponds to the

- (A) Number of internal nodes in the tree
 (B) Height of the tree
 (C) Number of nodes without a right sibling in the tree
 (D) Number of leaf nodes in the tree
7. The height of a tree is the length of the longest root-to-leaf path in it. The maximum and minimum number of nodes in a binary tree of height 5 are [2015]
- (A) 63 and 6, respectively
 (B) 64 and 5, respectively
 (C) 32 and 6, respectively
 (D) 31 and 5, respectively
8. Which of the following is/are correct inorder traversal sequence(s) of binary search tree(s)? [2015]
- I. 3, 5, 7, 8, 15, 19, 25
 II. 5, 8, 9, 12, 10, 15, 25
 III. 2, 7, 10, 8, 14, 16, 20
 IV. 4, 6, 7, 9, 18, 20, 25
- (A) I and IV only (B) II and III only
 (C) II and IV only (D) II only
9. Consider a max heap, represented by the array: 40, 30, 20, 10, 15, 16, 17, 8, 4 [2015]

Array Index	1	2	3	4	5	6	7	8	9
Value	40	30	20	10	15	16	17	8	4

Now consider that a value 35 is inserted into this heap. After insertion, the new heap is

- (A) 40, 30, 20, 10, 15, 16, 17, 8, 4, 35
 (B) 40, 35, 20, 10, 30, 16, 17, 8, 4, 15
 (C) 40, 30, 20, 10, 35, 16, 17, 8, 4, 15
 (D) 40, 35, 20, 10, 15, 16, 17, 8, 4, 30

10. A binary tree T has 20 leaves. The number of nodes in T having two children is _____ [2015]
11. Consider a binary tree T that has 200 leaf nodes. Then, the number of nodes in T that have exactly two children are _____. [2015]
12. While inserting the elements 71, 65, 84, 69, 67, 83 in an empty binary search tree (BST) in the sequence shown, the element in the lowest level is [2015]
 (A) 65 (B) 67
 (C) 69 (D) 83
13. Consider the following New-order strategy for traversing a binary tree: [2016]
- Visit the root;
 - Visit the right subtree using New-order;
 - Visit the left subtree using New-order;
- The New-order traversal of the expression tree corresponding to the reverse polish expression $3\ 4\ * \ 5 - 2 \wedge 6\ 7\ * \ 1 + -$ is given by:
 (A) $+ - 1\ 6\ 7\ * \ 2 \wedge 5 - 3\ 4\ *$
 (B) $- + 1\ * \ 6\ 7 \wedge 2 - 5\ * \ 3\ 4$
 (C) $- + 1\ * \ 7\ 6 \wedge 2 - 5\ * \ 4\ 3$
 (D) $1\ 7\ 6\ * + 2\ 5\ 4\ 3\ * - \wedge -$
14. Let T be a binary search tree with 15 nodes. The minimum and maximum possible heights of T are: [2017]
Note: The height of a tree with a single node is 0.
 (A) 4 and 15 respectively
 (B) 3 and 14 respectively
 (C) 4 and 14 respectively
 (D) 3 and 15 respectively
15. The pre-order traversal of a binary search tree is given by 12,8,6,2,7,9,10,16,15,19,17,20. Then the post-order traversal of this tree is: [2017]
 (A) 2, 6, 7, 8, 9, 10, 12, 15, 16, 17, 19, 20
 (B) 2, 7, 6, 10, 9, 8, 15, 17, 20, 19, 16, 12
 (C) 7, 2, 6, 8, 9, 10, 20, 17, 19, 15, 16, 12
 (D) 7, 6, 2, 10, 9, 8, 15, 16, 17, 20, 19, 12
16. The postorder traversal of a binary tree is 8, 9, 6, 7, 4, 5, 2, 3, 1. The inorder traversal of the same tree is 8, 6, 9, 4, 7, 2, 5, 1, 3. The height of a tree is the length of the longest path from the root to any leaf. The height of the binary tree above is _____. [2018]
17. The number of possible min-heaps containing each value from $\{1, 2, 3, 4, 5, 6, 7\}$ exactly once is _____. [2018]

ANSWER KEYS

EXERCISES

Practice Problems 1

- | | | | | | | | | | |
|-------|-------|-------|-------|-------|------|------|------|------|-------|
| 1. B | 2. B | 3. A | 4. C | 5. D | 6. C | 7. C | 8. B | 9. B | 10. B |
| 11. B | 12. D | 13. A | 14. C | 15. A | | | | | |

Practice Problems 2

- | | | | | | | | | | |
|-------|-------|-------|-------|-------|------|------|------|------|-------|
| 1. D | 2. B | 3. B | 4. A | 5. C | 6. B | 7. D | 8. A | 9. A | 10. B |
| 11. B | 12. A | 13. C | 14. A | 15. B | | | | | |

Previous Years' Questions

- | | | | | | | | | | |
|---------|-------|-------|-------|-------|-------|--------|------|------|--------|
| 1. A | 2. D | 3. B | 4. A | 5. 6 | 6. D | 7. A | 8. A | 9. B | 10. 19 |
| 11. 199 | 12. B | 13. C | 14. B | 15. B | 16. 4 | 17. 80 | | | |

TEST

PROGRAMMING AND DATA STRUCTURE (PART I)

Time: 60 min.

Directions for questions 1 to 30: Select the correct alternative from the given choices.

1. What is the behavior of following code?

```
auto int I;
int main()
{
}
```

- (A) Compiler error (B) No error
(C) Linker error (D) Runtime error

2. What is the output of following code:

```
#define scanf "%S is a string"
int main()
{
printf(scanf, scanf);
}
```

- (A) is a string is a string
(B) %S is a string is a string
(C) %S is a string %S is a string
(D) Syntax error

3. void rec_fun (int n, int sum)

```
{
int k = 0, j = 0;
if(n == 0) return;
k = n%10; j = n/10;
sum+ = k;
rec_fun(j, sum);
printf("%d\t", k);
}
```

main ()

```
{
int a = 2048; sum = 0;
rec_fun(a, sum);
printf("%d", sum);
}
```

What does the above program print?

- (A) 8 4 0 2 14
(B) 8 4 0 2 0
(C) 2 0 4 8 18
(D) 2 0 4 8 0

4. int fun (int * p, int n)

```
{
if (n <= 0) return 0;
```

```
else
if(*P % 2 == 0)
return *p + fun(p + 1, n - 1);
else
return *p - fun(p + 1, n -1);
}
main( )
{
int arr[ ] = {56, 48, 55, 10, 49, 14};
printf("% d", fun(arr, 6));
}
```

Which of the following is the output of above function?

- (A) 110 (B) -122
(C) 114 (D) 108

5. Consider the function:

```
int fun (int n)
{
static int i = 1;
if (n >= 5) return n;
n = n + i;
i ++;
return fun (n);
}
```

The value returned by $f(3)$ is

- (A) 6 (B) 7
(C) 8 (D) 9

6. Which of the following code will change a lower case letter to an upper case?

- (A) $\text{char } C_2 = (C_1 >= 'A' \& C_1 < 'Z') ? 'a' + 'C_1' - 'A' : C_1;$
(B) $\text{char } C_2 = (C_1 >= 'a' \& C_1 <= 'z') ? 'A' - 'a' + 'C_1' : C_1;$
(C) $\text{char } C_2 = (C_1 >= 'a' \& \& C_1 <= 'z') ? 'A' + 'C_1' - 'a' : C_1;$
(D) $\text{char } C_2 = (C_1 >= 'A' \& \& C_1 <= 'Z') ? 'A' - 'C_1' + 'a' : C_1;$

7. The following is a program to find the average length of several lines of text. What should be the lines of code corresponding to 'SecA' and 'SecB'.

```
main( )
{int n, count = 0, sum = 0;
float avg;
secA
```

```

{
count++;
sum + = n;
}
avg = (float) sum/count;
}
int linecount (void)
{
char line [80];
int count = 0;
while (secB)
{
count ++;
}
return (count);
}

```

- (A) Sec A: while ($n > 0$)
Sec B: line [count]! = 0
- (B) Sec A: while (linecount () > 0)
Sec B: (line [count] = getch () != '\n')
- (C) Sec A: while (($n = \text{linecount}()$) > 0)
Sec B: (line [count] = getch () != '\n')
- (D) None of these
8. What is the meaning of following declaration? $\text{int}(*f_1)$ ();
- (A) f_1 is a function which returns a pointer to an integer number.
- (B) f_1 is a pointer to a function which returns an integer number.
- (C) f_1 is a function which takes an integer pointer.
- (D) f_1 is a pointer to an integer number.
9. Which of the following represents the statement: “x is a pointer to a group of one dimensional 20 element arrays”.
- (A) $\text{int } *x[20]$; (B) $\text{int } *x[10][20]$;
- (C) $\text{int } **x[20]$; (D) $\text{int } (*x)[20]$;
10. What will be the output of the following code?
- ```

int number[] = {18, 20, 22, 24};
main()
{
int *q;
q = number;
q + = 4;
printf("%d", *q);
}

```
- (A) 24                                      (B) Compiler error
- (C) syntax error                        (D) -24

```

11. main()
{
char fname[] = "TIME 4 EDUCATION";
time4 (fname);
}
time4(char fname[5])
{
fname + = 7;
printf("%s", fname);
fname = 2;
printf("%s", fname);
}

```

What is the output of above code?

- (A) 4 EDUCATION
- (B) EDUCATION
- (C) TIME 4 EDUCATION
- (D) EDUCATION ME 4 EDUCATION
12. What will be output of following code?

```

int main()
{
extern int x;
x = 13;
printf ("%d", x);
return 0;
}

```

- (A) 13
- (B) Vary from compiler
- (C) Linker error
- (D) Undefined symbol
13. What is the output of the following code?

```

main()
{
int x = y = z = 100;
int i;
i = x > y < z;
printf("% d", i);
}

```

- (A) 0                                      (B) 1
- (C) error                                (D) No output
14. In the following code, how many times ‘while’ loop will be executed?
- ```

int count = 0;
while (count < 32767)
count++;

```
- (A) 32767
- (B) 32766
- (C) infinite times
- (D) varies from compiler to compiler

```
15. #define SUM(x)x + x * x
#define DIF(x)x * x - x
int main( )
{
float y = SUM(5) /DIF(5);
printf("%f", y);
}
(A) 1.5 (B) 0
(C) 2 (D) 1
```

16. Which of the following is a correct description of void (* ptr[10]) ();
- (A) ptr is an array of 10 pointers to functions returning type void.
 - (B) ptr is an array of 10 functions returning pointers of type void.
 - (C) ptr is an array of 10 functions returning void*.
 - (D) ptr is an array of data elements of type void.

```
17. int main( )
{
int p = 2, g = 2;
printf(" %d%d", p << g, p >> g);
}
```

Output of above code is

- (A) 1 16 (B) 4 0
- (C) 16 4 (D) 8 0

18. Which of the following is equivalent expression for $P = a * 16 + b/8$;
- (A) $P = (a << 4) + (b >> 2)$
 - (B) $P = (a >> 4) - (b << 2)$
 - (C) $P = (a << 4) + (b >> 3)$
 - (D) $P = (a << 4) + (b << 3)$

19. What is the output of the below code?

```
main( )
{
int I = 0, *j = &I;
f1(j);
*j = *j + 10;
printf(" %d%d", I, j);
}
f1(int *k)
{
*k+ = 15;
}
```

- (A) 20 55 (B) 25 25
- (C) 45 55 (D) 35 35

```
20. int rec_f2(int r)
{
if (r == 1 || r == 0)
```

```
return 1;
if (r % 2 == 0)
return(rec_f2 (r/2)+ 2);
else return ((rec_f2(r - 1) + 3);
}
main( )
{
printf ("%d", rec_f2(7));
}
```

Which of the following is the output of above program?

- (A) 10 (B) 11
- (C) 13 (D) 0

```
21. int guest(int a);
int host(int b);
main( )
{
int p = 50, q = 100, r ;
for (r = 0; r < 2; r++)
{
q = guest(p) + host(p);
printf(" %d", q);
}
}
int guest (int a)
{
int y;
y = host(a);
return(y);
}
int host (int a)
{
static int y = 0;
y = y + 1;
return (a + y);
}
```

The output of above code will be

- (A) 103 107 (B) 107 103
- (C) 110 107 (D) 107 110

22. Choose the best matching between groups A and B:

Group A	Group B
1 Volatile	P Queue
2 Function pointer	Q Auto
3 Default	R Guest and host
4 FIFO	S Switch-case

- (A) 1 - S, 2 - Q, 3 - R, 4 - P
- (B) 1 - Q, 2 - P, 3 - R, 4 - S

- (C) 1 – Q, 2 – R, 3 – S, 4 – P
 (D) 1 – R, 2 – Q, 3 – P, 4 – S

23. Which of the following expression represents the statement: “*P* is a function that accepts a pointer to a character array”.

- (A) `int p(char *a[]);` (B) `int (*p)(char (*a)[]);`
 (C) `int *p(char *a);` (D) `int p(char (*a) []);`

24. `void arr_fun(int [] [3]);`

```
main( )
{
int x[3][3] = {{10, 20, 30}, {40, 50, 60},
{70, 80, 90}};
arr_fun(x);
printf("%d", x[2][1]);
}
void arr_fun (int y[ ][3])
{
++y;
y[1][1] = 9;
}
```

What is the output of the above code?

- (A) 80 (B) 90
 (C) 70 (D) None of these

25. `#define F - 1`

```
#define T 1
#define N 0
main( )
{
if(N)
printf(" %s", "GOOD");
else
if(F)
printf(" %s", "MORNING");
else
printf(" %s", "GOOD NIGHT");
}
```

Output of above code will be

- (A) GOOD
 (B) MORNING
 (C) GOOD MORNING
 (D) GOOD NIGHT

26. An external variable

- (i) is defined once and declared in other functions.
 (ii) is globally accessible by all functions.
 (iii) cannot be static
 (iv) is defined after `main()`
 (A) (i) and (ii) (B) (i), (ii) and (iii)
 (C) (ii), (iii) and (iv) (D) (i), (ii), (iii), (iv)

Common Data for Questions 27 and 28:

Consider the following code:

```
int f(int P)
{
if(P <= 0) return 1;
if(P% 10 == 0)
return f(P - 2); //X
else
return f(P - 3); //Y
}
main( )
{
printf(" %d", f(30));
}
```

27. What will be the output of above code?

- (A) 10 (B) 50
 (C) 100 (D) 1

28. What will be the output if the lines labeled *X* and *Y* are changed as follows:

X: `return 3 + f(P/2);`

Y: `return 2 + f(P/3);`

- (A) 4 (B) 6
 (C) 8 (D) 10

Common data for Questions 29 and 30:

1. `float fn1(float n, int a)`

2. {

3. `float P, S;`

`int I;`

4. `for(S = x, P = 1, I = 1; P = P * x * x;`

5. `I < a; I ++)`

`S = S + P/(1);`

6. `return (S);`

7. }

8. `int f(int x)`

9. {

10. `int f = 1;`

11. `for (int i = 1; i <= n; i++)`

12. `f = f * i;`

13. `return (f);`

14. }

29. Output of above code for `fn1(2.1, 20)` is

(A) $x + \frac{x^2}{2!} + \frac{x^4}{4!} + \dots + a + 2.1$

(B) $\frac{x-x^3}{3!} + a + 2.1$

(C) $1 + \frac{x}{1!} + \frac{x^2}{2!} +$

(D) $1 + \frac{x}{1!} + \frac{x^3}{3!} + \frac{x^5}{5!} + \dots + a + 2.1$

30. When the statement numbered 4, 5, 6, 7 are replaced by
 {for (S = 1, P = -x, I = 1; I < a; i++)
 P = P * x * x - 1;

```
S = S + P/f(I);  
}
```

What will be the approximation of $f(x)$?

(A) $1 + x - x^2 + x^3 \dots$

(B) $1 - \frac{x}{1!} + \frac{x^3}{3!} -$

(C) $1 - \frac{x}{1!} + \frac{x^2}{2!} + \frac{x^3}{3!} + \dots$

(D) $1 - \frac{x}{1!} + \frac{x^2}{2!} - \frac{x^3}{3!} + \dots$

ANSWER KEYS

- | | | | | | | | | | |
|-------|-------|-------|-------|-------|-------|-------|-------|-------|-------|
| 1. A | 2. B | 3. A | 4. C | 5. A | 6. C | 7. C | 8. B | 9. D | 10. B |
| 11. D | 12. C | 13. B | 14. C | 15. D | 16. A | 17. D | 18. C | 19. B | 20. B |
| 21. A | 22. C | 23. B | 24. A | 25. B | 26. B | 27. D | 28. D | 29. A | 30. D |

This page is intentionally left blank

Algorithms

PART B

Chapter 1: Asymptotic Analysis	3.81
Chapter 2: Sorting Algorithms	3.98
Chapter 3: Divide-and-conquer	3.107
Chapter 4: Greedy Approach	3.116
Chapter 5: Dynamic Programming	3.135

U

N

I

T

3

This page is intentionally left blank

Chapter 1

Asymptotic Analysis

LEARNING OBJECTIVES

- Algorithm
- Recursive algorithms
- Towers of Hanoi
- Time complexity
- Space complexity
- SET representation
- TREE representation
- Preorder traversal
- Post-order traversal
- In order traversal
- Data structure
- Worst-case and average-case analysis
- Asymptotic notations
- Notations and functions
- Floor and ceil
- Recurrence
- Recursion-tree method
- Master method

ALGORITHM

An algorithm is a finite set of instructions that, if followed, accomplishes a particular task.

All algorithms must satisfy the following.

- Input: Zero or more quantities are externally supplied.
- Output: Atleast one quantity is produced.
- Definiteness: Each instruction should be clear and unambiguous.
- Finiteness: The algorithm should terminate after finite number of steps.
- Effectiveness: Every instruction must be very basic.

Once an algorithm is devised, it is necessary to show that it computes the correct answer for all possible inputs. This process is called algorithm validation. Analysis of algorithms refers to the task of determining how much computing time and storage an algorithm requires.

Analyzing Algorithms

The process of comparing 2 algorithms rate of growth with respect to time, space, number of registers, network, bandwidth etc is called analysis of algorithms.

This can be done in two ways

1. **Priori Analysis:** This analysis is done before the execution; the main principle behind this is frequency count of fundamental instruction.

This analysis is independent of CPU, OS and system architecture and it provides uniform estimated values.

2. **Posterior analysis:** This analysis is done after the execution. It is dependent on system architecture, CPU, OS etc. it provides non-uniform exact values.

Recursive Algorithms

A recursive function is a function that is defined in terms of itself. An algorithm is said to be recursive if the same algorithm is invoked in the body.

Towers of Hanoi

There was a diamond tower (labeled *A*) with 64-golden disks. The disks were of decreasing size and were stacked on the tower in decreasing order of size bottom to top. Besides this tower there were 2 other diamond towers (labeled *B* and *C*) we have to move the disks from tower *A* to tower *B* using tower *C* for intermediate storage. As the disks are very heavy, they can be moved only one at a time. No disk can be on top of a smaller disk.

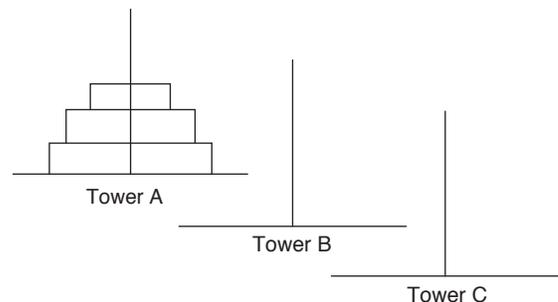
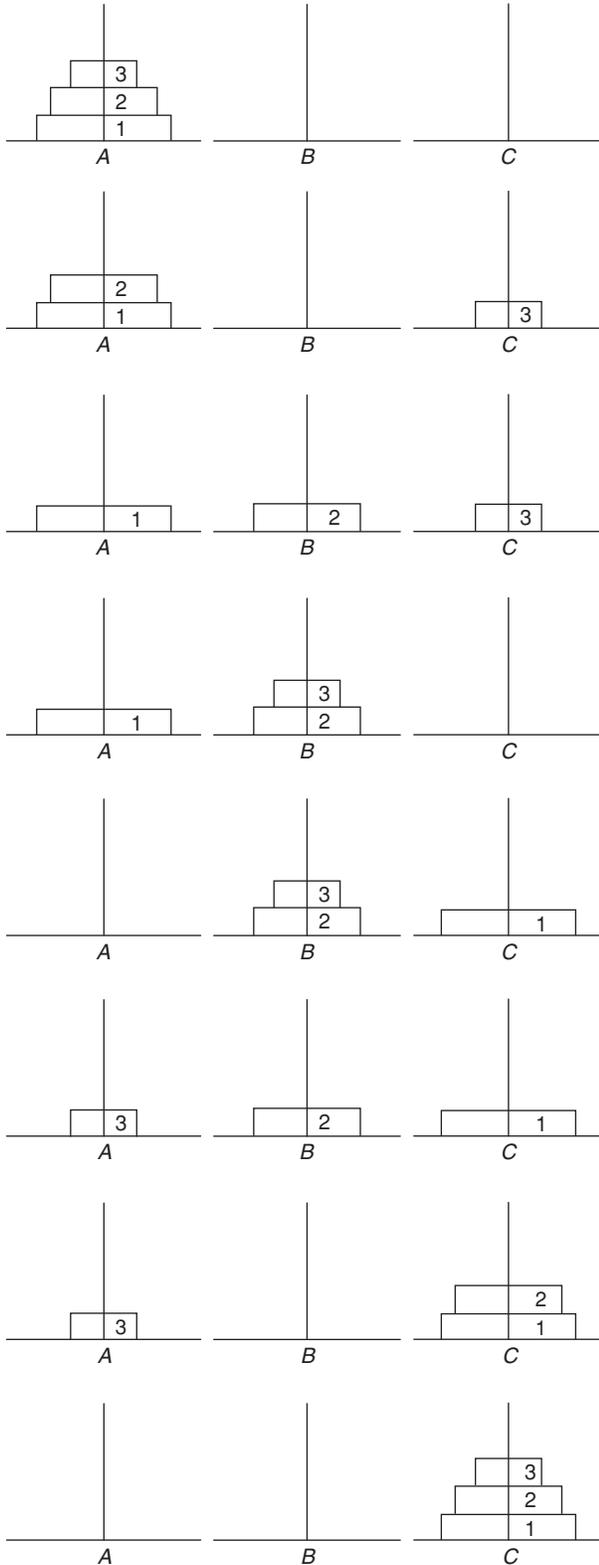


Figure 1 Towers of Hanoi

3.82 | Unit 3 • Algorithms

Assume that the number of disks is 'n'. To get the largest disk to the bottom of tower B, we move the remaining (n - 1) disks to tower C and then move the largest to tower B. Now move the disks from tower C to tower B.

Example:



To move '3' disks from tower A to tower 'C' requires 7 disk movements

∴ For 'n' disks, the number of disk movements required is $2^n - 1 = 2^3 - 1 = 7$

Time complexity

$$T(n) = 1 + 2T(n - 1)$$

$$T(n) = 1 + 2(1 + 2(T(n - 2)))$$

$$T(n) = 1 + 2 + 2^2 T(n - 2)$$

$$T(n) = 1 + 2 + 2^2 (1 + 2T(n - 3))$$

$$T(n) = 1 + 2 + 2^2 + 2^3 + T(n - 3)$$

$$T(n) = 1 + 2 + 2^2 + \dots + 2^{i-1} + 2^i T(n - i)$$

$$T(n) = \sum_{i=0}^{n-1} 2^i$$

The time complexity is exponential, it grows as power of 2.

$$\therefore T(n) \cong O(2^n)$$

Space complexity

The space complexity of an algorithm is the amount of memory it needs to run to completion. The measure of the quantity of input data is called the size of the problem. For example, the size of a matrix multiplication problem might be the largest dimension of the matrices to be multiplied. The size of a graph problem might be the number of edges. The limiting behavior of the complexity as size increases is called the asymptotic time complexity.

- It is the asymptotic complexity of an algorithm which ultimately determines the size of problems that can be solved by the algorithm.
- If an algorithm processes inputs of size 'n' in time cn^2 for some constant c, then we say that the time complexity of that algorithm is $O(n^2)$, more precisely a function $g(n)$ is said to be $O(f(n))$ if there exists a constant c such that $g(n) \leq c(f(n))$ for all but some finite set of non-negative values for n.
- As computers become faster and we can handle larger problems, it is the complexity of an algorithm that determines the increase in problem size that can be achieved with an increase in computer speed.
- Suppose we have 5 algorithms Algorithm 1 – Algorithm 5 with the following time complexities.

Algorithm	Time Complexity
Algorithm – 1	n
Algorithm – 2	n log n
Algorithm – 3	n ²
Algorithm – 4	n ³
Algorithm – 5	2n

The time complexity is, the number of time units required to process an input of size 'n'. Assume that input size 'n' is 1000 and one unit of time equals to 1 millisecond.

The following figure gives the sizes of problems that can be solved in one second, one minute, and one hour by each of these five algorithms.

Algorithm	Time Complexity	Maximum Problem Size		
		1 sec	1 min	1 hour
Algorithm – 1	n	1000	6×10^4	3.6×10^6
Algorithm – 2	$n \log n$	140	4893	2.0×10^5
Algorithm – 3	n^2	31	244	1897
Algorithm – 4	n^3	10	39	153
Algorithm – 5	$2n$	9	15	21

From the above table, we can say that different algorithms will give different results depending on the input size. Algorithm – 5 would be best for problems of size $2 \leq n \leq 9$, Algorithm – 3 would be best for $10 \leq n \leq 58$, Algorithm – 2 would be best for $59 \leq n \leq 1025$, and Algorithm – 1 is best for problems of size greater than 1024.

SET REPRESENTATION

A common use of a list is to represent a set, with this representation the amount of memory required to represent a set is proportional to the number of elements in the set. The amount of time required to perform a set operation depends on the nature of the operation.

- Suppose A and B are 2 sets. An operation such as $A \cap B$ requires time atleast proportional to the sum of the sizes of the 2 sets, since the list representing A and the list representing B must be scanned atleast once.
- The operation $A \cup B$ requires time atleast proportional to the sum of the set sizes, we need to check for the same element appearing in both sets and delete one instance of each such element.
- If A and B are disjoint, we can find $A \cup B$ in time independent of the size of A and B by simply concatenating the two lists representing A and B .

GRAPH REPRESENTATION

A graph $G = (V, E)$ consists of a finite, non-empty set of vertices V and a set of edges E . If the edges are ordered pairs (V, W) of vertices, then the graph is said to be directed; V is called the tail and W the head of the edge (V, W) . There are several common representations for a graph $G = (V, E)$. One such representation is adjacency matrix, a $|V| \times |V|$ matrix M of 0's and 1's, where the ij_{th} element, $m[i, j] = 1$, if and only if there is an edge from vertex i to vertex j .

- The adjacency matrix representation is convenient for graph algorithms which frequently require knowledge of whether certain edges are present.
- The time needed to determine whether an edge is present is fixed and independent of $|V|$ and $|E|$.

- Main drawback of using adjacency matrix is that it requires $|V|^2$ storage even if the graph has only $O(|V|)$ edges.
- Another representation for a graph is by means of lists. The adjacency list for a vertex v is a list of all vertices W adjacent to V . A graph can be represented by $|V|$ adjacency lists, one for each vertex.

Example:

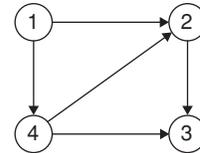


Figure 2 Directed graph

	1	2	3	4
1	0	1	0	1
2	0	0	1	0
3	0	0	0	0
4	0	1	1	0

Figure 3 Adjacency matrix

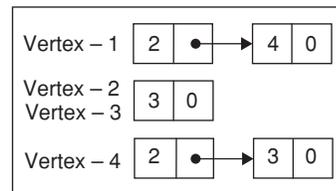


Figure 4 Adjacency lists

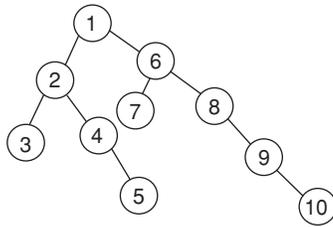
There are edges from vertex – 1 to vertex – 2 and 4, so the adjacency list for 1 has items 2 and 4 linked together in the format given above.

- The adjacency list representation of a graph requires storage proportional to $|V| + |E|$, this representation is used when $|E| < |V|^2$.

TREE REPRESENTATION

A directed graph with no cycles is called a directed acyclic graph. A directed graph consisting of a collection of trees is called a forest. Suppose the vertex ‘ v ’ is root of a sub tree, then the depth of a vertex ‘ v ’ in a tree is the length of the path from the root to ‘ v ’.

- The height of a vertex ‘ v ’ in a tree is the length of a longest path from ‘ v ’ to a leaf.
- The height of a tree is the height of the root
- The level of a vertex ‘ v ’ in a tree is the height of the tree minus the depth of ‘ v ’.



	Left child	Right child
1	2	6
2	3	4
3	0	0
4	0	5
5	0	0
6	7	8
7	0	0
8	0	9
9	0	10
10	0	0

Figure 5 A binary tree and its representation

- Vertex 3 is of depth '2', height '0' and the level is 2 (Height of tree – depth of '3' = 4 – 2 = 2).
- A binary tree is represented by 2 arrays: left child and right child.
- A binary tree is said to be complete if for some integer k , every vertex of depth less than k has both a left child and a right child and every vertex of depth k is a leaf. A complete binary tree of height k has exactly $(2^{k+1} - 1)$ vertices.
- A complete binary tree of height k is often represented by a single array. Position 1 in the array contains the root. The left child of the vertex in position ' i ' is located at position ' $2i$ ' and the right child at position ' $2i + 1$ '.

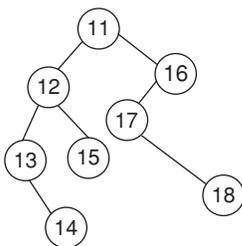
Tree Traversals

Many algorithms which make use of trees often traverse the tree in some order. Three commonly used traversals are pre-order, postorder and inorder.

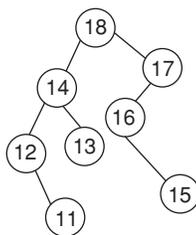
Pre-order Traversal

A pre-order traversal of T is defined recursively as follows:

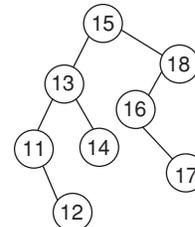
1. Visit the root.
2. Visit in pre-order the sub trees with roots $v_1, v_2 \dots v_k$ in that order.



(a)



(b)



(c)

Figure 6 (a) Pre-order, (b) Post-order (c) In-order

Post-order traversal

A post-order traversal of T is defined recursively as follows:

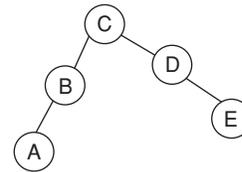
1. Visit in post-order the sub trees with roots $v_1, v_2, v_3, \dots v_k$ in that order.
2. Visit the root r .

In-order Traversal

An in-order traversal is defined recursively as follows:

1. Visit in in-order the left sub tree of the root ' r '.
2. Visit ' r '.
3. Visit in inorder the right sub tree of r .

Example: Consider the given tree



What are the pre-order, post-order and in-order traversals of the above tree?

Solution: Pre-order – CBADE
 Post-order – ABEDC
 In-order – ABCDE

DATA STRUCTURE

A data structure is a way to store and organize data in-order to facilitate access and modifications. No single data structure works well for all purposes, it is important to know the strengths and limitations of several data structures.

Efficiency

Algorithms devised to solve the same problem often differ dramatically in their efficiency. Let us compare efficiencies of Insertion sort and merge sort; insertion sort, takes time equal to $C_1 n^2$ to sort ' n ' elements, where C_1 is a constant that does not depend on ' n '. It takes time proportional to n^2 , merge sort takes time equal to $C_2 n \log n$, C_2 is another constant that also does not depend on ' n '. Insertion sort has a smaller constant factor than merge sort ($C_1 < C_2$) constant factors are far less significant in the running time.

Merge sort has a factor of ‘log n ’ in its running time, insertion sort has a factor of ‘ n ’, which is much larger. Insertion sort is faster than merge sort for small input sizes, once the input size ‘ n ’ becomes large enough, merge sort will perform better. No matter how much smaller C_1 is than C_2 . There will always be a crossover point beyond which merge sort is faster.

Example: Consider 2 computers, computer A (faster computer), B (slower computer). Computer A runs insertion sort and computer B runs merge sort. Each computer is given 2 million numbers to sort. Suppose that computer A executes one billion instruction per second and computer B executes only 10 million instructions per second, computer A is 100 times faster than computer B ($C_1 = 4$, $C_2 = 50$). How much time is taken by both the computers?

Solution: Insertion sort takes $C_1 * n^2$ time
Merge sort takes $C_2 * n * \log n$ time
 $C_1 = 4$, $C_2 = 50$
Computer A takes

$$\frac{4 \times (2 \times 10^6)^2 \text{ instructions}}{10^9 \text{ instructions/second}} \cong 4000 \text{ seconds}$$

Computer B takes

$$\begin{aligned} &= \frac{50 \times 2 \times 10^6 \times \log(2 \times 10^6) \text{ instructions}}{10^7 \text{ instructions/second}} \\ &= 209 \text{ seconds} \end{aligned}$$

By using an algorithm whose running time grows more slowly, even with an average compiler, computer B runs 20 times faster than computer A . The advantage of merge sort is even more pronounced when we sort ten million numbers. As the problem size increases, so does the relative advantage of merge sort.

Worst-case and average-case analysis

In the analysis of insertion sort, the best case occurs when the array is already sorted and the worst case, in which the input array is reversely sorted. We concentrate on finding the worst-case running time, that is the longest running time for any input of size ‘ n ’.

- The worst-case running time of an algorithm is an upper bound on the running time for any input. It gives us a guarantee that the algorithm will never take any longer.
- The ‘average-case’ is as bad as the worst-case. Suppose that we randomly choose ‘ n ’ numbers and apply insertion sort. To insert an element $A[j]$, we need to determine where to insert in sub-array $A[1 \dots J-1]$. On average half the elements in $A[1 \dots J-1]$ are less than $A[j]$ and half the elements are greater. So $t_j = j/2$. The average-case running time turns out to be a quadratic function of the input size.

ASYMPTOTIC NOTATIONS

Asymptotic notations are mostly used in computer science to describe the asymptotic running time of an algorithm. As an example, an algorithm that takes an array of size n as input and runs for time proportional to n^2 is said to take $O(n^2)$ time.

5 Asymptotic Notations:

- O (Big-oh)
- θ (Theta)
- Ω (Omega)
- o (Small-oh)
- ω

How to Use Asymptotic Notation for Algorithm Analysis?

Asymptotic notation is used to determine rough estimates of relative running time of algorithms. A worst-case analysis of any algorithm will always yield such an estimate, because it gives an upper bound on the running time $T(n)$ of the algorithm, that is $T(n) \leq g(n)$.

Example:

$a \leftarrow 0$	1 unit	1 time
for $i \leftarrow 1$ to n do{	1 unit	n times
for $j \leftarrow 1$ to i do{	1 unit	$n(n+1)/2$ times
$a \leftarrow a + 1$	1 unit	$n(n+1)/2$ times

Where the times for the inner loop have been computed as follows: For each i from 1 to n , the loop is executed i times, so the total number of times is $1 + 2 + 3 + \dots + n = \sum_{i=1}^n i = n(n+1)/2$

Hence in this case

$$T(n) = 1 + n + 2n(n+1)/2 = n^2 + 2n + 1$$

If we write $g(n) = n^2 + 2n + 1$, then $T(n) \in \theta(g(n))$,

That is $T(n) \in \theta(n^2 + 2n + 1)$, we actually write $T(n) \in \theta(n^2)$, as recommended by the following rule:

- Although the definitions of asymptotic notation allow one to write, for example, $T(n) \in O(3n^2 + 2)$.

We simplify the function in between the parentheses as much as possible (in terms of rate of growth), and write instead $T(n) \in O(n^2)$

For example: $T(n) \in \theta(4n^3 - n^2 + 3)$

$$T(n) \in \theta(n^3)$$

For instance $O\left(\sum_{i=1}^n i\right)$, write $O(n^2)$ after computing the sum.

- In the spirit of the simplicity rule above, when we are to compare, for instance two candidate algorithms A and B having running times ($T_A(n) = n^2 - 3n + 4$ and $T_B(n) = 5n^3 + 3$, rather than writing $T_A(n) \in O(T_B(n))$, we write $T_A(n) \in \theta(n^2)$, and $T_B(n) \in \theta(n^3)$, and then we conclude that A

is better than B , using the fact that n^2 (quadratic) is better than n^3 (cubic) time, since $n^2 \in O(n^3)$.

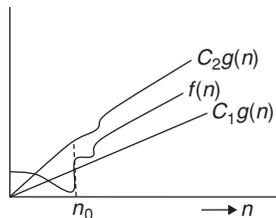
Order of Growth

In the rate of growth or order of growth, we consider only the leading term of a formula. Suppose the worst case running time of an algorithm is $an^2 + bn + c$ for some constants a, b and c . The leading term is an^2 . We ignore the leading term's constant coefficient, since constant factors are less significant than the rate of growth in determining computational efficiency for large inputs. Thus we can write, the worst-case running time is $\theta(n^2)$.

We usually consider one algorithm to be more efficient than another if its worst-case running time has a lower order of growth. Due to constant factors and lower order terms, this evaluation may be in error for small inputs. But for large inputs, $\theta(n^2)$ algorithm will run more quickly in the worst-case than $\theta(n^3)$ algorithm.

θ -Notation

A function $f(n)$ belongs to the set $\theta(g(n))$ if there exists a positive constant C_1 and C_2 such that it can be “sandwiched” between $C_1g(n)$ and $C_2g(n)$ for sufficiently large n . We write $f(n) \in \theta(g(n))$ to indicate that $f(n)$ is a member of $\theta(g(n))$ or we can write $f(n) = \theta(g(n))$ to express the same notation.



The above figure gives an intuitive picture of functions $f(n)$ and $g(n)$, where we have that $f(n) = \theta(g(n))$, for all the values of ‘ n ’ to the right of n_0 , the value of $f(n)$ lies at or above $C_1g(n)$ and at or below $C_2g(n)$. $g(n)$ is asymptotically tight bound for $f(n)$. The definition of $\theta(g(n))$ requires that every member $f(n) \in \theta(g(n))$ be asymptotically non-negative, that is $f(n)$ must be non-negative whenever ‘ n ’ is sufficiently large.

The θ -notation is used for asymptotically bounding a function from both above and below. We would use θ (theta) notation to represent a set of functions that bounds a particular function from above and below.

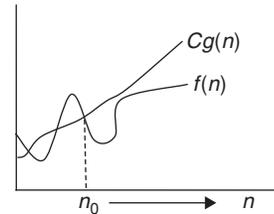
Definition: We say that a function $f(n)$ is theta of $g(n)$ written as $f(n) = \theta(g(n))$ if such exists positive constants C_1, C_2 and n_0 such that $0 \leq C_1g(n) \leq f(n) \leq C_2g(n), \forall n \geq n_0$.

Example: Let $f(n) = 5.5n^2 - 7n$, verify whether $f(n)$ is $\theta(n^2)$. Lets have constants $c_1 = 9$ and $n_0 = 2$, such that $0 \leq f(n) \leq C_1n^2, \forall n \geq n_0$. From example, 4 we have constants $C_2 = 3$, and $n_0 = 2.8$, such that $0 \leq C_2n^2 \leq f(n), \forall n \geq n_0$. To show $f(n)$ is $\theta(n^2)$, we have got hold of two constants C_1 and C_2 . We fix the n_0 for θ as maximum $\{2, 2.8\} = 2.8$.

- The lower order terms of an asymptotically positive function can be ignored in determining asymptotically tight bounds because they are insignificant for large n .
- A small fraction of the highest order term is enough to dominate the lower order term. Thus setting C_1 to a value that is slightly smaller than the coefficient of the highest order term and setting C_2 to a value that is slightly larger permits the inequalities in the definition of θ -notation to be satisfied. If we take a quadratic function $f(n) = an^2 + bn + c$, where a, b and c are constants and $a > 0$. Throwing away the lower order terms and ignoring the constant yields $f(n) = \theta(n^2)$.
- We can express any constant function as $\theta(n^0)$, or $\theta(1)$ we shall often use the notation $\theta(1)$ to mean either a constant or a constant function with respect to some variable.

O-Notation

We use O-notation to give an upper bound on a function, within a constant factor.



The above figure shows the intuition behind O-notation. For all values ‘ n ’ to the right of n_0 , the value of the function $f(n)$ is on or below $g(n)$. We write $f(n) = O(g(n))$ to indicate that a function $f(n)$ is a member of the set $O(g(n))$.

$f(n) = \theta(g(n))$ implies $f(n) = O(g(n))$. Since θ notation is stronger notation than O-notation set theoretically, we have $\theta(g(n)) \subseteq O(g(n))$. Thus any quadratic function $an^2 + bn + c$, where $a > 0$, is in $\theta(n^2)$ also shows that any quadratic function is in $O(n^2)$ when we write $f(n) = O(g(n))$, we are claiming that some constant multiple of $g(n)$ is an asymptotic upper bound on $f(n)$, with no claim about how tight an upper bound it is.

The O-notation is used for asymptotically upper bounding a function. We would use O (big-oh) notation to represent a set of functions that upper bounds a particular function.

Definition We say that a function $f(n)$ is big oh of $g(n)$ written as $f(n) = O(g(n))$ if there exists positive constants C and n_0 such that

$$0 \leq f(n) \leq Cg(n), \forall n \geq n_0$$

Solved Examples

Example 1: let $f(n) = n^2$

Then $f(n) = O(n^2)$

$f(n) = O(n^2 \log n)$

$f(n) = O(n^{2.5})$

$f(n) = O(n^3)$

$f(n) = O(n^4) \dots$ so on.

Example 2: Let $f(n) = 5.5n^2 - 7n$, verify whether $f(n)$ is $O(n^2)$

Solution: Let C be a constant such that

$$5.5n^2 - 7n \leq Cn^2, \text{ or } n \geq \frac{7}{c-5.5}$$

Fix $C = 9$, to get $n \geq 2$

So our $n_0 = 2$ and $C = 9$

This shows that there exists, positive constants $C = 9$ and $n_0 = 2$ such that

$$0 \leq f(n) \leq Cn^2, \forall n \geq n_0$$

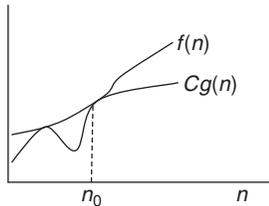
Example 3:

$$h(n) = 3n^3 + 10n + 1000 \log n \in O(n^3)$$

$$h(n) = 3n^3 + 10n + 1000 \log n \in O(n^4)$$

- Using O -notation, we can describe the running time of an algorithm by inspecting the algorithm's overall structure. For example, the doubly nested loop structure of the insertion sort algorithm yields an $O(n^2)$ upper bound on the worst-case running time. The cost of each iteration of the inner loop is bounded from above by $O(1)$ (constant), the inner loop is executed almost once for each of the n^2 pairs.
- $O(n^2)$ bound on worst-case running time of insertion sort also applies to its running time on every input.
- The $\theta(n^2)$ bound on the worst-case running time of insertion sort, however, does not imply a $\theta(n^2)$ bound on the running time of insertion sort on every input, when the input is already sorted, insertion sort runs in $\theta(n)$ time.

Ω (omega)-notation



The Ω -notation is used for asymptotically lower bounding a function. We would use Ω (big-omega) notation to represent a set of functions that lower bounds a particular function.

Definition We say that a function $f(n)$ is big-omega of $g(n)$ written as $f(n) = \Omega(g(n))$ if there exists positive constants C and n_0 such that

$$0 \leq Cg(n) \leq f(n), \forall n \geq n_0$$

The intuition behind Ω -notation is shown in the above figure. For all values 'n' to the right of n_0 , the value of $f(n)$

is on or above $Cg(n)$. For any 2 functions $f(n)$ and $g(n)$ we have $f(n) = \theta(g(n))$ if $f(n) = O(g(n))$ and $f(n) = \Omega(g(n))$. From the above statement we can say that, $an^2 + bn + c = \theta(n^2)$ for any constants a, b and c , where $a > 0$, immediately implies that

$$\therefore an^2 + bn + c = \Omega(n^2)$$

$$\therefore an^2 + bn + c = O(n^2)$$

Example 4: Let $f(n) = 5.5n^2 - 7n$.

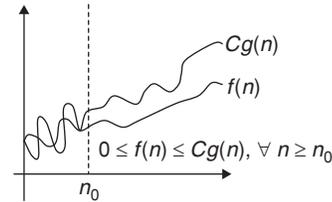
Verify whether $f(n)$ is $\Omega(n^2)$

Solution: Let C be a constant such that $5.5n^2 - 7n \geq Cn^2$ or

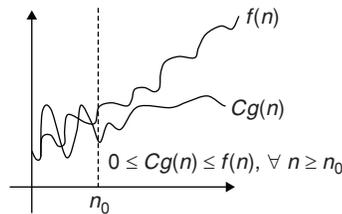
$$n \geq \frac{7}{5.5-C}$$

Fix $C = 3$, to get $n \geq 2.8$. So, our $n_0 = 2.8$ and $C = 3$

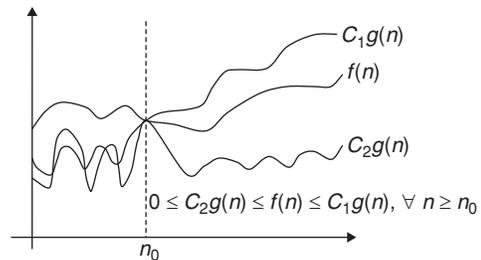
This shows that there exists positive constants $C = 3$ and $n_0 = 2.8$, such that $0 \leq Cn^2 \leq f(n), \forall n \geq n_0$.



(a) $f(n) = O(g(n))$



(b) $f(n) = \Omega(g(n))$



(c) $f(n) = \theta(g(n))$

Figure 7 A diagrammatic representation of the asymptotic notations O, Ω and θ

- Ω -notation describes a lower bound; it is used to bound the best-case running time of an algorithm. The best-case running time of insertion sort is $\Omega(n)$. The running time of insertion sort falls between $\Omega(n)$ and $O(n^2)$, since it falls anywhere between a linear function of 'n' and a quadratic function of 'n'.

- When we say that the running time of an algorithm is $\Omega(g(n))$, we mean that no matter what particular input of size 'n' is chosen for each value of n, the running time on that input is at least a constant times $g(n)$, for sufficiently large 'n'.

O-notation

The asymptotic upper bound provided by O-notation may or may not be asymptotically tight. The bound $2n^3 = O(n^3)$ is asymptotically tight, but the bound $2n = O(n^2)$ is not. We use O-notation to denote an upper bound that is not asymptotically tight.

ω-notation

By analogy, ω-notation is to Ω-notation as o-notation is to O-notation. We use ω-notation to denote a lower bound that is not asymptotically tight.

It is defined as

$$f(n) \in \omega(g(n)) \text{ if and only if } g(n) \in o(f(n))$$

Comparison of functions

Transitivity

- $f(n) = \theta(g(n))$ and $g(n) = \theta(h(n))$
 $\Rightarrow f(n) = \theta(h(n))$
- $f(n) = O(g(n))$ and $g(n) = O(h(n))$
 $\Rightarrow f(n) = O(h(n))$
- $f(n) = \Omega(g(n))$ and $g(n) = \Omega(h(n))$
 $\Rightarrow f(n) = \Omega(h(n))$
- $f(n) = o(g(n))$ and $g(n) = o(h(n))$
 $\Rightarrow f(n) = o(h(n))$
- $f(n) = \omega(g(n))$ and $g(n) = \omega(h(n))$
 $\Rightarrow f(n) = \omega(h(n))$

Reflexivity

- $f(n) = \theta(f(n))$
- $f(n) = O(f(n))$
- $f(n) = \Omega(f(n))$

Symmetry

$$f(n) = \theta(g(n)) \text{ if and only if } g(n) = \theta(f(n))$$

Transpose symmetry

- $f(n) = O(g(n))$ if and only if $g(n) = \Omega(f(n))$
- $f(n) = o(g(n))$ if and only if $g(n) = \omega(f(n))$

NOTATIONS AND FUNCTIONS

Floor and Ceil

For any real number 'x', we denote the greatest integer less than or equal to x by $\lfloor x \rfloor$ called as floor of x and the least integer greater than or equal to x by $\lceil x \rceil$ called as ceiling of x.

$$x - 1 < \lfloor x \rfloor \leq x \leq \lceil x \rceil < x + 1 \text{ for any integer } n,$$

$$\left\lceil \frac{n}{2} \right\rceil + \left\lfloor \frac{n}{2} \right\rfloor = n,$$

For any real number $n \geq 0$ and integer $a, b > 0$

$$\left\lceil \frac{\left\lfloor \frac{n}{a} \right\rfloor}{b} \right\rceil = \left\lfloor \frac{n}{ab} \right\rfloor$$

$$\left\lfloor \frac{\left\lceil \frac{n}{a} \right\rceil}{b} \right\rfloor = \left\lceil \frac{n}{ab} \right\rceil$$

Polynomials

Given a non-negative integer k, a polynomial in n of degree 'k' is a function $p(n)$ of the form $p(n) = \sum_{i=0}^k a_i n^i$

Where the constants a_0, a_1, \dots, a_k are the coefficients of the polynomial and $a_k \neq 0$.

For an asymptotically positive polynomial $p(n)$ of degree k, we have $p(n) = \theta(n^k)$

Exponentials

For all real $a > 0$, m and n, we have the following identities:

$$a^0 = 1$$

$$a^1 = a$$

$$a^{-1} = \frac{1}{a}$$

$$(a^m)^n = a^{mn}$$

$$(a^m)^n = (a^n)^m$$

$$a^m a^n = a^{m+n}$$

$$e^x = 1 + x + \frac{x^2}{2!} + \frac{x^3}{3!} + \dots = \sum_{i=0}^{\infty} \frac{x^i}{i!}$$

- For all real x, we have inequality $e^x \geq 1 + x$
- If $x = 0$, we have $1 + x \leq e^x \leq 1 + x + x^2$

Logarithms

$\lg n = \log_2 n$ (binary logarithm)

$\ln n = \log_e n$ (natural logarithm)

$\lg^k n = (\log n)^k$ (exponentiation)

$\lg \lg n = \lg(\lg n)$ (composition)

For all real $a > 0, b > 0, c > 0$ and n,

$$\log_c(ab) = \log_c a + \log_c b$$

$$\log_b a^n = n \log_b a$$

$$\log_b a = \frac{\log_c a}{\log_c b}$$

$$\log_b (1/a) = -\log_b a$$

$$\log_b a = \frac{1}{\log_a b}$$

$$a^{\log_b c} = c^{\log_b a}$$

Factorials

$n!$ is defined for integers $n \geq 0$ as

$$n! = \begin{cases} 1 & \text{if } n = 0 \\ (n-1)! * n & n > 0 \end{cases}$$

A weak upper bound on the factorial function is $n! \leq n^n$ since each of the n terms in the factorial product is almost n .

$$n! = o(n^n)$$

$$n! = \omega(2^n)$$

$$\lg(n!) = \theta(n \log n)$$

Iterated Logarithm

The notation $\lg^* n$ is used to denote the iterated logarithm. Let ' $\lg^{(i)} n$ ' be as defined above, with $f(n) = \lg n$. The logarithm of a non-positive number is undefined, ' $\lg^{(i)} n$ ' is defined only if $\lg^{(i-1)} n > 0$;

The iterated logarithm function is defined as $\lg^* n = \min \{i \geq 0 : \lg^{(i)} n \leq 1\}$. This function is a very slowly growing function.

$$\lg^* 2 = 1$$

$$\lg^* 4 = 2$$

$$\lg^* 16 = 3$$

$$\lg^* 65536 = 4$$

$$\lg^*(2^{65536}) = 5$$

RECURRENCES

When an algorithm contains a recursive call to itself, its running time can often be described by a recurrence. A recurrence is an equation that describes a function in terms of its value on smaller inputs. For example, the worst-case running time $T(n)$ of the merge-sort can be described as

$$T(n) = \begin{cases} \theta(1) & \text{if } n = 1 \\ 2T(n/2) + \theta(n) & \text{if } n > 1 \end{cases}$$

The time complexity of merge-sort algorithm in the worst-case is $T(n) = \theta(n \log n)$

There are 3 methods to solve recurrence relations:

1. Substitution method
2. Recursion-tree method
3. Master method

Substitution Method

In this method one has to guess the form of the solution. It can be applied only in cases when it is easy to guess the form of the answer. Consider the recurrence relation

$$T(n) = 2T(n/2) + n$$

We guess that the solution is $T(n) = O(n \log n)$ we have to prove that

$$T(n) \leq c n \log n \quad (\because c > 0)$$

Assume that this bound holds for $\lfloor n/2 \rfloor$

$$T(n/2) \leq c(n/2) \log(n/2) + n$$

$$T(n) \leq 2(c(n/2) \log(n/2) + n) + n$$

$$\leq cn \log n - cn \log 2 + n$$

$$\leq cn \log n - cn + n$$

$$\leq cn \log n \quad (\because c \geq 1)$$

Recursion-tree Method

In a recursion-tree, each node represents the cost of single sub problem somewhere in the set of recursive function invocations. We sum the costs within each level of the tree to obtain a set of per-level costs, and then we sum all the per-level costs to determine the total cost of all levels of the recursion. Recursion trees are useful when the recurrence describes the running time of a divide-and-conquer algorithm.

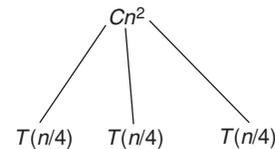
Example:

Consider the given recurrence relation

$$T(n) = 3T(n/4) + \theta(n^2)$$

We create a recursion tree for the recurrence

$$T(n) = 3T(n/4) + Cn^2$$



The Cn^2 term at the root represents the cost at the top level of recursion, and the three sub trees of the root represent the costs incurred by the sub problems of size $n/4$.

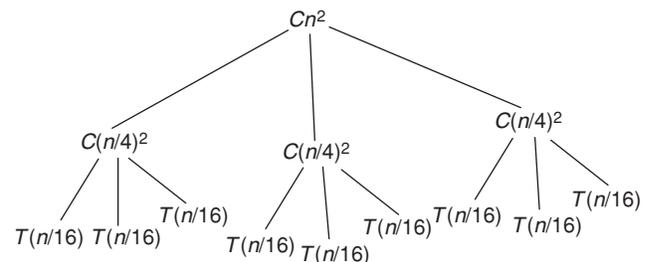


Figure 8 Recursion tree for $T(n) = 3T(n/4) + cn^2$

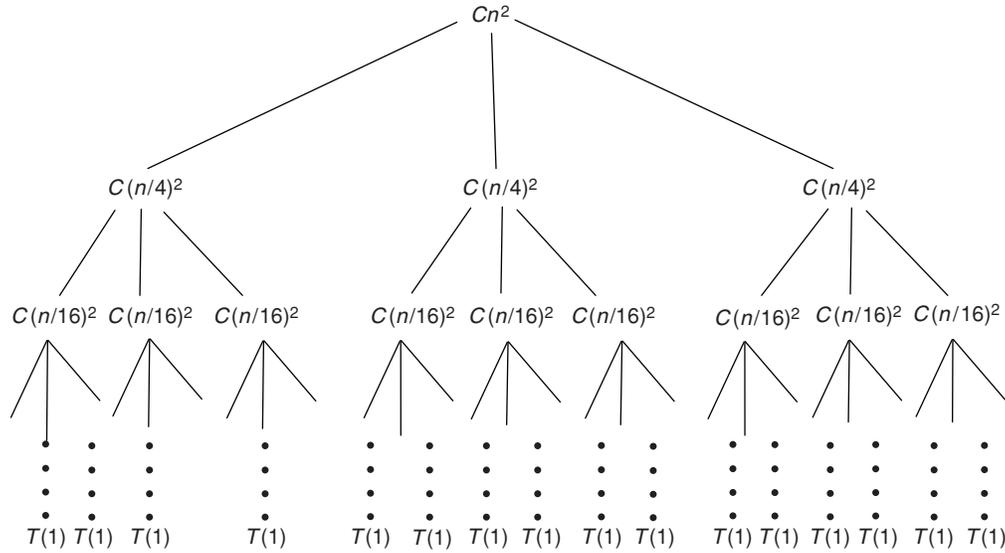


Figure 9 Expanded Recursion tree with height $\log_4 n$ (\therefore levels $\log_4 n + 1$)

The sub-problem size for a node at depth ‘ i ’ is $n/4^i$, at this depth, the size of the sub-problem would be $n = 1$, when $n/4^i = 1$ or $i = \log_4 n$, the tree has $\log_4 n + 1$ levels.

- We have to determine the cost at each level of the tree. Each level has 3 times more nodes than the level above, so the number of nodes at depth ‘ i ’ is 3^i .
- Sub problem sizes reduce by a factor of ‘4’ for each level we go down from the root, each node at depth i , for $i = 0, 1, 2 \dots \log_4 n - 1$, has a cost of $c(n/4^i)^2$.

Total cost over all nodes at depth i , for $i = 0, 1, \dots \log_4 n - 1$

$$= 3^i * c \left(\frac{n}{4^i} \right)^2 = \left(\frac{3}{16} \right)^i cn^2$$

The last level, at depth $\log_4 n$ has 3^i nodes $= 3^{\log_4 n} = n^{\log_4 3}$ each contributing cost $T(1)$, for a total cost of $n^{\log_4 3} T(1)$, which is $\theta(n^{\log_4 3})$ cost of the entire tree is equal to sum of costs over all levels.

$$\begin{aligned} T(n) &= cn^2 + \frac{3}{16}cn^2 + \left(\frac{3}{16}\right)^2 cn^2 + \dots + \\ &\left(\frac{3}{16}\right)^{\log_4 n - 1} cn^2 + \dots + \theta(n^{\log_4 3}) \\ &= \sum_{i=0}^{\log_4 n - 1} \left(\frac{3}{16}\right)^i cn^2 + \theta(n^{\log_4 3}) \\ &< \sum_{i=0}^{\infty} \left(\frac{3}{16}\right)^i cn^2 + \theta(n^{\log_4 3}) \\ &= \frac{1}{1 - \left(\frac{3}{16}\right)} cn^2 + \theta(n^{\log_4 3}) \\ &= \frac{16}{13} cn^2 + \theta(n^{\log_4 3}) = O(n^2) \end{aligned}$$

Master Method

Let $a \geq 1$ and $b > 1$ be constants, let $f(n)$ be a function and let $T(n)$ be defined on the non-negative integers by the recurrence

$$T(n) = aT(n/b) + f(n)$$

$T(n)$ can be bounded asymptotically as follows

1. If $f(n) = O(n^{\log_b a - \epsilon})$ for some constant $\epsilon > 0$, then $T(n) = \theta(n^{\log_b a})$
2. If $f(n) = \theta(n^{\log_b a})$ then $T(n) = \theta(n^{\log_b a} \log n)$
3. If $f(n) = \Omega(n^{\log_b a + \epsilon})$ for some constant $\epsilon > 0$, and if $af(n/b) \leq cf(n)$ for some constant $c < 1$ and all sufficiently large n , then $T(n) = \theta(f(n))$.

Note: In the first case, not only must $f(n)$ be smaller than $n^{\log_b a}$, it must be polynomially smaller. That is, $f(n)$ must be asymptotically smaller than $n^{\log_b a}$ by a factor of n^ϵ , for some constant $\epsilon > 0$.

In the third case, not only must $f(n)$ be larger than $n^{\log_b a}$, it must be polynomially larger and in addition satisfy the regularity condition $af(n/b) \leq Cf(n)$.

Example: Consider the given recurrence relation $T(n) = 9T(n/3) + n$.

To apply master theorem, the recurrence relation must be in the following form:

$$T(n) = aT(n/b) + f(n)$$

$$a = 9, b = 3, f(n) = n$$

$$n^{\log_b a} = n^{\log_3 9} = n^2$$

Since $f(n) = O(n^{\log_3 9 - \epsilon})$, where $\epsilon = 1$

We can apply case 1 of the master theorem and the solution is $T(n) = \theta(n^2)$.

EXERCISES

Practice Problems 1

Directions for questions 1 to 15: Select the correct alternative from the given choices.

1. What is the time complexity of the recurrence relation

$$T(n) = 2T\left(\frac{n}{2}\right) + n^2?$$

- (A) $\theta(n^2)$ (B) $\theta(n)$
(C) $\theta(n^3)$ (D) $\theta(n \log n)$

2. What is the time complexity of the recurrence relation

by using masters theorem $T(n) = 2T\left(\frac{n}{2}\right) + n$?

- (A) $\theta(n^2)$ (B) $\theta(n)$
(C) $\theta(n^3)$ (D) $\theta(n \log n)$

3. What is the time complexity of the recurrence relation

by using master theorem, $T(n) = 2T\left(\frac{n}{4}\right) + n^{0.51}$

- (A) $\theta(n^2)$ (B) $\theta(n)$
(C) $\theta(n^3)$ (D) $(n^{0.51})$

4. What is the time complexity of the recurrence relation

using master theorem, $T(n) = 7T\left(\frac{n}{3}\right) + n^2$?

- (A) $\theta(n^2)$ (B) $\theta(n)$
(C) $\theta(n^3)$ (D) $(\log n)$

5. Time complexity of $f(x) = 4x^2 - 5x + 3$ is

- (A) $O(x)$ (B) $O(x^2)$
(C) $O(x^{3/2})$ (D) $O(x^{0.5})$

6. Time complexity of $f(x) = (x^2 + 5 \log_2 x)/(2x + 1)$ is

- (A) $O(x)$ (B) $O(x^2)$
(C) $O(x^{3/2})$ (D) $O(x^{0.5})$

7. For the recurrence relation, $T(n) = 2T\left(\left\lfloor \sqrt{n} \right\rfloor\right) + \lg n$, which is tightest upper bound?

- (A) $T(n) = O(n^2)$ (B) $T(n) = O(n^3)$
(C) $T(n) = O(\log n)$ (D) $T(n) = O(\lg n \lg \lg n)$

8. Consider $T(n) = 9T(n/3) + n$, which of the following is TRUE?

- (A) $T(n) = \theta(n^2)$ (B) $T(n) = \theta(n^3)$
(C) $T(n) = \Omega(n^3)$ (D) $T(n) = O(n)$

9. If $f(n)$ is $100 * n$ seconds and $g(n)$ is $0.5 * n$ seconds then

- (A) $f(n) = g(n)$ (B) $f(n) = \Omega(g(n))$
(C) $f(n) = w(g(n))$ (D) None of these

10. Solve the recurrence relation using master method:

$$T(n) = 4T(n/2) + n^2$$

- (A) $\theta(n \log n)$ (B) $\theta(n^2 \log n)$
(C) $\theta(n^2)$ (D) $\theta(n^3)$

11. Arrange the following functions according to their order of growth (from low to high):

(A) $\sqrt[3]{n}$, $0.001n^4 + 3n^3 + 1$, 3^n , 2^{2n}

(B) 3^n , 2^{2n} , $\sqrt[3]{n}$, $0.001n^4 + 3n^3 + 1$

(C) 2^{2n} , $\sqrt[3]{n}$, 3^n , $0.001n^4 + 3n^3 + 1$

(D) $\sqrt[3]{n}$, 2^{2n} , 3^n , $0.001n^4 + 3n^3 + 1$

12. The following algorithm checks whether all the elements in a given array are distinct:

Input: array $A[0 \dots n - 1]$

Output: true (or) false

For $i \leftarrow 0$ to $n - 2$ do

For $j \leftarrow i + 1$ to $n - 1$ do

if $A[i] = A[j]$ return false

return true

The time complexity in worst case is

- (A) $\theta(n^2)$ (B) $\theta(n)$
(C) $\theta(\log n)$ (D) $\theta(n \log n)$

13. The order of growth for the following recurrence relation is $T(n) = 4T(n/2) + n^3$, $T(1) = 1$

- (A) $\theta(n)$ (B) $\theta(n^3)$
(C) $\theta(n^2)$ (D) $\theta(\log n)$

14. Time complexity of $T(n) = 2T\left(\frac{n}{4}\right) + \sqrt{3}$ is

(A) $\theta(\sqrt{n} \log n)$ (B) $\theta(\sqrt{n} \log \sqrt{n})$

(C) $\theta(\sqrt{n})$ (D) $\theta(n^2)$

15. Consider the following three claims

(I) $(n + k)^m = \theta(n^m)$, where k and m are constants

(II) $2^{n+1} = O(2^n)$

(III) $2^{2n+1} = O(2^n)$

Which one of the following is correct?

- (A) I and III (B) I and II
(C) II and III (D) I, II and III

Practice Problems 2

Directions for questions 1 to 15: Select the correct alternative from the given choices.

1. Arrange the order of growth in ascending order:

(A) $O(1) > O(\log n) > O(n) > O(n^2)$

(B) $O(n) > O(1) > O(\log n) > O(n^2)$

(C) $O(\log n) > O(n) > O(1) > O(n^2)$

(D) $O(n^2) > O(n) > O(\log n) > O(1)$

2. $\sqrt{n} = \Omega(\log n)$ means

(A) To the least \sqrt{n} is $\log n$

(B) \sqrt{n} is $\log n$ always

(C) \sqrt{n} is at most $\log n$

(D) None of these

3. Which of the following is correct?

(i) $\theta(g(n)) = O(g(n)) \cap \Omega(g(n))$

(ii) $\theta(g(n)) = O(g(n)) \cup \Omega(g(n))$

- (A) (i) is true (ii) is false (B) Both are true
 (C) Both are false (D) (ii) is true (i) is false
4. $2n^2 = x(n^3)$, x is which notation?
 (A) Big-oh (B) Small-oh
 (C) Ω – notation (D) θ – notation
5. Master method applies to recurrence of the form $T(n) = aT(n/b) + f(n)$ where
 (A) $a \geq 1, b > 1$ (B) $a = 1, b > 1$
 (C) $a > 1, b = 1$ (D) $a \geq 1, b \geq 1$
6. What is the time complexity of the recurrence relation using master method?

$$T(n) = 4T\left(\frac{n}{2}\right) + n$$

- (A) $\theta(n^2)$ (B) $\theta(n)$
 (C) $\theta(\log n)$ (D) $\theta(n \log n)$
7. Use the informal definitions of O, θ, Ω to determine these assertions which of the following assertions are true.
 (A) $n(n+1)/2 \in O(n^3)$ (B) $n(n+1)/2 \in O(n^2)$
 (C) $n(n+1)/2 \in \Omega(n)$ (D) All the above
8. Match the following:

(i)	Big-oh	(A)	\geq
(ii)	Small-o	(B)	\leq
(iii)	Ω	(C)	$=$
(iv)	θ	(D)	$<$
(v)	ω	(E)	$>$

- (A) (i) – D, (ii) – A, (iii) – C, (iv) – B, (v) – E
 (B) (i) – B, (ii) – D, (iii) – A, (iv) – C, (v) – E
 (C) (i) – C, (ii) – A, (iii) – B, (iv) – E, (v) – D
 (D) (i) – A, (ii) – B, (iii) – C, (iv) – D, (v) – E
9. Which one of the following statements is true?
 (A) Both time and space efficiencies are measured as functions of the algorithm input size.
 (B) Only time efficiencies are measured as a function of the algorithm input size.
 (C) Only space efficiencies are measured as a function of the algorithm input size.
 (D) Neither space nor time efficiencies are measured as a function of the algorithm input size.
10. Which of the following is true?
 (A) Investigation of the average case efficiency is considerably more difficult than investigation of the worst case and best case efficiencies.
 (B) Investigation of best case is more complex than average case.

- (C) Investigation of worst case is more complex than average case.
 (D) None of these

11. Time complexity of $T(n) = T(n/3) + T(2n/3) + O(n)$ is
 (A) $O(1)$
 (B) $O(n \log n)$
 (C) $O(\log n)$
 (D) $O(n^2)$

12. Solve the recurrence relation to find $T(n)$: $T(n) = 4(n/2) + n$
 (A) $\theta(n^2)$ (B) $\theta(\log_2 n)$
 (C) $\theta(n^2 \log_2 n)$ (D) $\theta(n^3)$

13. What is the worst case analysis for the given code?

```
int search (int a[ ], int x, int n)
{
    int i;
    for (i = 0 ; i < n; i ++ )
        if (a [i] == x)
            return i;
    return -1;
}
```

- (A) $O(n)$ (B) $O(n \log n)$
 (C) $O(\log n)$ (D) $O(n^2)$

14. Find the time complexity of the given code.

```
void f (int n)
{
    if (n > 0)
    {
        f (n/2);
        f (n/2);
    }
}
```

- (A) $\theta(n^2)$
 (B) $\theta(n)$
 (C) $\theta(n \log n)$
 (D) $\theta(2^n)$

15. The running time of the following algorithm procedure

```
A(n)
if n ≤ 2
return (1)
else
return (A(√n))
is described by
```

- (A) $O(\sqrt{n} \log n)$
 (B) $O(\log n)$
 (C) $O(\log \log n)$
 (D) $O(n)$

PREVIOUS YEARS' QUESTIONS

- The median of n elements can be found in $O(n)$ time. Which one of the following is correct about the complexity of quick sort, in which median is selected as pivot? [2006]
 - $\theta(n)$
 - $\theta(n \log n)$
 - $\theta(n^2)$
 - $\theta(n^3)$
- Given two arrays of numbers $a_1 \dots a_n$ and $b_1 \dots b_n$ where each number is 0 or 1, the fastest algorithm to find the largest span (i, j) such that $a_i + a_{i+1} + \dots + a_j = b_i + b_{i+1} + \dots + b_j$, or report that there is not such span, [2006]
 - Takes $O(3^n)$ and $\Omega(2^n)$ time if hashing is permitted
 - Takes $O(n^3)$ and $\Omega(n^{2.5})$ time in the key comparison model
 - Takes $\Theta(n)$ time and space
 - Takes $O(\sqrt{n})$ time only if the sum of the $2n$ elements is an even number
- Consider the following segment of C-code:


```
int j, n;
j = 1;
while (j <= n)
j = j*2;
```

 The number of comparisons made in the execution of the loop for any $n > 0$ is: [2007]
 - $\lceil \log_2 n \rceil + 1$
 - n
 - $\lceil \log_2 + n \rceil$
 - $\lfloor \log_2 n \rfloor + 1$
- In the following C function, let $n \geq m$.


```
int gcd(n, m)
{
if (n%m == 0) return m;
n = n%m;
return gcd(m, n);
}
```

 How many recursive calls are made by this function? [2007]
 - $\Theta(\log_2 n)$
 - $\Omega(n)$
 - $\Theta(\log_2 \log_2 n)$
 - $\Theta(\sqrt{n})$
- What is the time complexity of the following recursive function:


```
int DoSomething (int n) {
if (n <= 2)
return 1;
else
return(DoSomething (floor(sqrt(n))+ n);} [2007]
```

 - $\Theta(n^2)$
 - $\Theta(n \log_2 n)$
 - $\Theta(\log_2 n)$
 - $\Theta(\log_2 \log_2 n)$
- An array of n numbers is given, where n is an even number. The maximum as well as the minimum of

these n numbers needs to be determined. Which of the following is TRUE about the number of comparisons needed? [2007]

- At least $2n - c$ comparisons, for some constant c , are needed.
 - At most $1.5n - 2$ comparisons are needed.
 - At least $n \log_2 n$ comparisons are needed.
 - None of the above.
7. Consider the following C code segment:

```
int IsPrime(n)
{
int i, n;
for(i=2; i<= sqrt(n); i++)
if (n%i == 0)
{printf("Not Prime\n"); return 0;}
return 1;
}
```

Let $T(n)$ denote the number of times the *for* loop is executed by the program on input n . Which of the following is TRUE? [2007]

- $T(n) = O(\sqrt{n})$ and $T(n) = \Omega(\sqrt{n})$
 - $T(n) = O(\sqrt{n})$ and $T(n) = \Omega(1)$
 - $T(n) = O(n)$ and $T(n) = \Omega(\sqrt{n})$
 - None of the above
8. The most efficient algorithm for finding the number of connected components in an undirected graph on n vertices and m edges has time complexity [2008]
- $\Theta(n)$
 - $\Theta(m)$
 - $\Theta(m + n)$
 - $\Theta(mn)$

9. Consider the following functions:

$$f(n) = 2^n$$

$$g(n) = n!$$

$$h(n) = n^{\log n}$$

Which of the following statements about the asymptotic behavior of $f(n)$, $g(n)$, and $h(n)$ is true? [2008]

- $f(n) = O(g(n)); g(n) = O(h(n))$
 - $f(n) = \Omega(g(n)); g(n) = O(h(n))$
 - $g(n) = O(f(n)); h(n) = O(f(n))$
 - $h(n) = O(f(n)); g(n) = \Omega(f(n))$
10. The minimum number of comparisons required to determine if an integer appears more than $n/2$ times in a sorted array of n integers is [2008]
- $\Theta(n)$
 - $\Theta(\log n)$
 - $\Theta(\log * n)$
 - $\Theta(1)$
11. We have a binary heap on n elements and wish to insert n more elements (not necessarily one after another) into this heap. The total time required for this is [2008]

- (A) $\Theta(\log n)$ (B) $\Theta(n)$
 (C) $\Theta(n \log n)$ (D) $\Theta(n^2)$

12. The running time of an algorithm is represented by the following recurrence relation: [2009]

$$T(n) = \begin{cases} n & n \leq 3 \\ T\left(\frac{n}{3}\right) + cn & \text{otherwise} \end{cases}$$

Which one of the following represents the time complexity of the algorithm?

- (A) $\theta(n)$ (B) $\theta(n \log n)$
 (C) $\theta(n^2)$ (D) $\theta(n^2 \log n)$
13. Two alternative packages A and B are available for processing a database having 10^k records. Package A requires $0.0001 n^2$ time units and package B requires $10n \log_{10} n$ time units to process n records. What is the smallest value of k for which package B will be preferred over A ? [2010]
- (A) 12 (B) 10
 (C) 6 (D) 5
14. An algorithm to find the length of the longest monotonically increasing sequence of numbers in an array $A[0 : n - 1]$ is given below.

Let L denote the length of the longest monotonically increasing sequence starting at index in the array.

Initialize $L_{n-1} = 1$,

For all i such that $0 \leq i \leq n - 2$

$L_i = 1 + L_{i+1}$, if $A[i] < A[i + 1]$,

1 otherwise

Finally the length of the longest monotonically increasing sequence is $\text{Max}(L_0, L_1, \dots, L_{n-1})$

Which of the following statements is TRUE? [2011]

- (A) The algorithm uses dynamic programming paradigm.
 (B) The algorithm has a linear complexity and uses branch and bound paradigm.
 (C) The algorithm has a non-linear polynomial complexity and uses branch and bound paradigm.
 (D) The algorithm uses divide and conquer paradigm.
15. Which of the given options provides the increasing order of asymptotic complexity of functions f_1, f_2, f_3 and f_4 ? [2011]

$f_1(n) = 2^n$

$f_2(n) = n^{3/2}$

$f_3(n) = n \log_2 n$

$f_4(n) = n^{\log_2 n}$

- (A) f_3, f_2, f_4, f_1 (B) f_3, f_2, f_1, f_4
 (C) f_2, f_3, f_1, f_4 (D) f_2, f_3, f_4, f_1

16. Let $W(n)$ and $A(n)$ denote respectively, the worst-case and average-case running time of an algorithm

executed on input of size n . Which of the following is ALWAYS TRUE? [2012]

- (A) $A(n) = \Omega(W(n))$ (B) $A(n) = \Theta(W(n))$
 (C) $A(n) = O(W(n))$ (D) $A(n) = o(W(n))$

17. The recurrence relation capturing the optimal execution time of the Towers of Hanoi problem with n discs is [2012]

- (A) $T(n) = 2T(n - 2) + 2$
 (B) $T(n) = 2T(n - 1) + n$
 (C) $T(n) = 2T(n/2) + 1$
 (D) $T(n) = 2T(n - 1) + 1$

18. A list of n strings, each of length n , is sorted into lexicographic order using the merge sort algorithm. The worst-case running time of this computation is [2012]

- (A) $O(n \log n)$ (B) $O(n^2 \log n)$
 (C) $O(n^2 + \log n)$ (D) $O(n^2)$

19. Consider the following function:

```
int unknown (int n) {
    int i, j, k = 0;
    for (i = n/2; i <= n; i++)
        for (j = 2; j <= n; j = j*2)
            k = k + n/2;
    return (k);
}
```

The return value of the function is [2013]

- (A) $\Theta(n^2)$ (B) $\Theta(n^2 \log n)$
 (C) $\Theta(n^3)$ (D) $\Theta(n^3 \log n)$

20. The number of elements that can be sorted in $\Theta(\log n)$ time using heap sort is [2013]

- (A) $\Theta(1)$
 (B) $\Theta(\sqrt{\log n})$
 (C) $\Theta\left(\frac{\log n}{\log \log n}\right)$
 (D) $\Theta(\log n)$

21. Which one of the following correctly determines the solution of the recurrence relation with $T(1) = 1$

$T(n) = 2T\left(\frac{n}{2}\right) + \log n$? [2014]

- (A) $\theta(n)$ (B) $\theta(n \log n)$
 (C) $\theta(n^2)$ (D) $\theta(\log n)$

22. An algorithm performs $(\log N)^{1/2}$ find operations, N insert operations, $(\log N)^{1/2}$ delete operations, and $(\log N)^{1/2}$ decrease-key operations on a set of data items with keys drawn from a linearly ordered set. For a delete operation, a pointer is provided to the record that must be deleted For the decrease – key operation, a pointer is provided to the record that has its key decreased Which one of the following data structures is the most suited for the algorithm to use, if the

goal is to achieve the best total asymptotic complexity considering all the operations? [2015]

- (A) Unsorted array
- (B) Min-heap
- (C) Sorted array
- (D) Sorted doubly linked list

23. Consider the following C function.

```
int fun1(int n) {
    int i, j, k, p, q=0;
    for (i=1; i<n; ++i) {
        p=0;
        for (j=n; j>1; j=j/2)
            ++p;
        for (k=1; k<p; k=k*2)
            ++q;
    }
    return q;
}
```

Which one of the following most closely approximates the return value of the function fun1? [2015]

- (A) n^3
- (B) $n(\log n)^2$
- (C) $n \log n$
- (D) $n \log(\log n)$

24. An unordered list contains n distinct elements. The number of comparisons to find an element in this list that is neither maximum nor minimum is [2015]

- (A) $\theta(n \log n)$
- (B) $\theta(n)$
- (C) $\theta(\log n)$
- (D) $\theta(1)$

25. Consider a complete binary tree where the left and the right subtrees of the root are max-heaps. The lower bound for the number of operations to convert the tree to a heap is [2015]

- (A) $\Omega(\log n)$
- (B) $\Omega(n)$
- (C) $\Omega(n \log n)$
- (D) $\Omega(n^2)$

26. Consider the equality $\sum_{i=0}^n i^3 =$ and the following choices for X

1. $\theta(n^4)$
2. $\theta(n^5)$
3. $O(n^5)$
4. $\Omega(n^3)$

The equality above remains correct if X is replaced by [2015]

- (A) Only 1
- (B) Only 2
- (C) 1 or 3 or 4 but not 2
- (D) 2 or 3 or 4 but not 1

27. Consider the following array of elements

<89, 19, 50, 17, 12, 15, 2, 5, 7, 11, 6, 9, 100>

The minimum number of interchanges needed to convert it into a max-heap is [2015]

- (A) 4
- (B) 5
- (C) 2
- (D) 3

28. Let $f(n) = n$ and $g(n) = n^{(1 + \sin n)}$, where n is a positive integer. Which of the following statements is/are correct? [2015]

I. $f(n) = O(g(n))$

II. $f(n) = \Omega(g(n))$

- (A) Only I
- (B) Only II
- (C) Both I and II
- (D) Neither I nor II

29. A queue is implemented using an array such that ENQUEUE and DEQUEUE operations are performed efficiently. Which one of the following statements is **CORRECT** (n refers to the number of items in the queue)? [2016]

- (A) Both operations can be performed in $O(1)$ time.
- (B) At most one operation can be performed in $O(1)$ time but the worst case time for the other operation will be $\Omega(n)$.
- (C) The worst case time complexity for both operations will be $\Omega(n)$.
- (D) Worst case time complexity for both operations will be $\Omega(\log n)$.

30. Consider a carry look ahead adder for adding two n -bit integers, built using gates of fan-in at most two. The time to perform addition using this adder is [2016]

- (A) $\Theta(1)$
- (B) $\Theta(\log(n))$
- (C) $\Theta(\sqrt{n})$
- (D) $\Theta(n)$

31. N items are stored in a sorted doubly linked list. For a *delete* operation, a pointer is provided to the record to be deleted. For a *decrease-key* operation, a pointer is provided to the record on which the operation is to be performed.

An algorithm performs the following operations on the list in this order: $\Theta(N)$ *delete*, $O(\log N)$ *insert*, $O(\log N)$ *find*, and $\Theta(N)$ *decrease-key*. What is the time complexity of all these operations put together? [2016]

- (A) $O(\log^2 N)$
- (B) $O(N)$
- (C) $O(N^2)$
- (D) $\Theta(N^2 \log N)$

32. In an adjacency list representation of an undirected simple graph $G = (V, E)$, each edge (u, v) has two adjacency list entries: $[v]$ in the adjacency list of u , and $[u]$ in the adjacency list of v . These are called twins of each other. A twin pointer is a pointer from an adjacency list entry to its twin. If $|E| = m$ and $|V| = n$, and the memory size is not a constraint, what is the time complexity of the most efficient algorithm to set the twin pointer in each entry in each adjacency list? [2016]

- (A) $\Theta(n^2)$
- (B) $\Theta(n + m)$
- (C) $\Theta(m^2)$
- (D) $\Theta(n^4)$

ANSWER KEYS**EXERCISES****Practice Problems 1**

1. A 2. D 3. D 4. A 5. B 6. A 7. D 8. A 9. A 10. B
11. A 12. A 13. B 14. A 15. B

Practice Problems 2

1. A 2. A 3. A 4. B 5. A 6. A 7. D 8. B 9. A 10. A
11. B 12. A 13. A 14. B 15. C

Previous Years' Questions

1. B 2. C 3. A 4. C 5. 6. B 7. B 8. C 9. D 10. A
11. B 12. A 13. C 14. A 15. A 16. C 17. D 18. B 19. B 20. C
21. A 22. A 23. D 24. D 25. A 26. C 27. D 28. D 29. A 30. B
31. C 32. B 33. B 34. B 35. C 36. B 37. B

Sorting Algorithms

LEARNING OBJECTIVES

- ☞ *Sorting algorithms*
- ☞ *Merge sort*
- ☞ *Bubble sort*
- ☞ *Insertion sort*
- ☞ *Selection sort*
- ☞ *Selection sort algorithm*
- ☞ *Binary search trees*
- ☞ *Heap sort*
- ☞ *Sorting—performing delete max operations*
- ☞ *Max-heap property*
- ☞ *Min-heap property*
- ☞ *Priority queues*

SORTING ALGORITHMS

Purpose of sorting

Sorting is a technique which reduces problem complexity and search complexity.

- Insertion sort takes $\theta(n^2)$ time in the worst case. It is a fast inplace sorting algorithm for small input sizes.
- Merge sort has a better asymptotic running time $\theta(n \log n)$, but it does not operate in place.
- Heap sort, sorts ' n ' numbers inplace in $\theta(n \log n)$ time, it uses a data structure called heap, with which we can also implement a priority queue.
- Quick sort also sorts ' n ' numbers in place, but its worst – case running time is $\theta(n^2)$. Its average case is $\theta(n \log n)$. The constant factor in quick sort's running time is small, This algorithm performs better for large input arrays.
- Insertion sort, merge sort, heap sort, and quick sort are all comparison based sorts; they determine the sorted order of an input array by comparing elements.
- We can beat the lower bound of $\Omega(n \log n)$ if we can gather information about the sorted order of the input by means other than comparing elements.
- The counting sort algorithm, assumes that the input numbers are in the set $\{1, 2, \dots, k\}$. By using array indexing as a tool for determining relative order, counting sort can sort n numbers in $\theta(k + n)$ time. Thus counting sort runs in time that is linear in size of the input array.
- Radix sort can be used to extend the range of counting sort. If there are ' n ' integers to sort, each integer has ' d ' digits, and each

digit is in the set $\{1, 2, \dots, k\}$, then radix sort can sort the numbers in $\theta(d(n + k))$ time. Where ' d ' is constant. Radix sort runs in linear time.

- Bucket sort, requires knowledge of the probabilistic distribution of numbers in the input array.

MERGE SORT

Suppose that our division of the problem yields ' a ' sub problems,

each of which is $\left(\frac{1}{b}\right)$ th size of the original problem. For merge

sort, both a and b are 2, but sometimes $a \neq b$. If we take $D(n)$ time to divide the problem into sub problems and $C(n)$ time to combine the solutions of the sub problems into the solution to the original problem. The recurrence relation for merge sort is

$$T(n) = \begin{cases} \theta(1) & \text{if } n \leq c, \\ aT(n/b) + D(n) + C(n) & \text{otherwise} \end{cases}$$

Running time is broken down as follows:

Divide: This step computes the middle of the sub array, which takes constant time $\theta(1)$.

Conquer: We solve 2 sub problems of size $(n/2)$ each recursively which takes $2T(n/2)$ time.

Combine: Merge sort procedure on an n -element sub array takes time $\theta(n)$.

- Worst case running time $T(n)$ of merge sort

$$T(n) = \begin{cases} 0(1) & \text{if } n \leq 1 \\ aT(n/2) + \theta(n) & \text{if } n > 1 \end{cases}$$

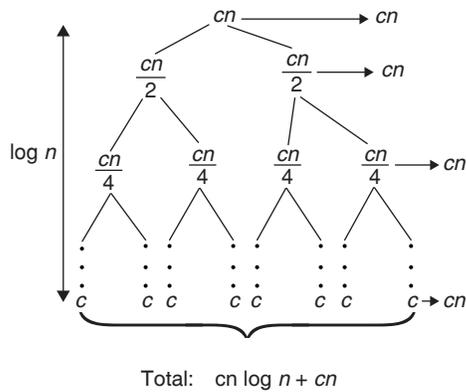


Figure 1 Recurrence tree

The top level has total cost ‘ cn ’, the next level has total cost $c(n/2) + c(n/2) = cn$ and the next level has total cost $c(n/4) + c(n/4) + c(n/4) + c(n/4) = cn$ and so on. The i th level has total cost $2^i c (n/2^i) = cn$. At the bottom level, there are ‘ n ’ nodes, each contributing a cost of ‘ c ’, for a total cost of ‘ cn ’. The total number of levels of the ‘recursion tree’ is $\log n + 1$.

There are $\log n + 1$ levels, each costing cn , for a total cost of $cn (\log n + 1) = cn \log n + cn$ ignoring the low-order term and the constant c , gives the desired result of $\theta(n \log n)$.

BUBBLE SORT

Bubble sort is a simple sorting algorithm that works by repeatedly stepping through the list to be sorted, comparing each pair of adjacent items, and swapping them if they are in the wrong order. The pass through the list is repeated until no swaps are needed, which indicates that the list is sorted. The algorithm gets its name from the way smaller elements ‘bubble’ to the top of the list.

Example: Take the array of numbers ‘5 1 4 2 8’ and sort the array from lowest number to greatest number using bubble sort algorithm. In each step, elements underlined are being compared.

First pass:

(5 1 4 2 8) → (1 5 4 2 8), here algorithm compares the first 2 elements and swaps them
 (1 5 4 2 8) → (1 4 5 2 8), swap (5 > 4)
 (1 4 5 2 8) → (1 4 2 5 8), swap (5 > 2)
 (1 4 2 5 8) → (1 4 2 5 8), since these elements are already in order, algorithm does not swap them.

Second pass:

(1 4 2 5 8) → (1 4 2 5 8)
 (1 4 2 5 8) → (1 2 4 5 8), swap since (4 > 2)
 (1 2 4 5 8) → (1 2 4 5 8)
 (1 2 4 5 8) → (1 2 4 5 8)

The array is already sorted, but our algorithm does not know if it is completed. The algorithm needs one whole pass without any swap to know it is sorted.

Third pass:

(1 2 4 5 8) → (1 2 4 5 8)
 (1 2 4 5 8) → (1 2 4 5 8)
 (1 2 4 5 8) → (1 2 4 5 8)
 (1 2 4 5 8) → (1 2 4 5 8)

Finally the array is sorted, and the algorithm can terminate.

Algorithm

```
void bubblesort (int a [ ], int n)
{
    int i, j, temp;
    for (i=0; i < n-1; i++)
    {
        for (j=0; j < n - 1 - i; j++)
        {
            if (a [j] > a [j + 1])
            {
                temp = a [j + 1];
                a [j + 1] = a [j];
                a [j] = temp;
            }
        }
    }
}
```

INSERTION SORT

Insertion sort is a comparison sort in which the sorted array is built one entry at a time. It is much less efficient on large lists than more advanced algorithms such as quick sort, heap sort, (or) merge sort. Insertion sort provides several advantages.

- Efficient for small data sets.
- Adaptive, i.e., efficient for data set that are already substantially sorted. The complexity is $O(n + d)$, where d is the number of inversions.
- More efficient in practice than most other simple quadratic, i.e., $O(n^2)$ algorithms such as selection sort (or) bubble sort, the best case is $O(n)$.
- Stable, i.e., does not change the relative order of elements with equal keys.
- In-place i.e., only requires a constant amount $O(1)$ of additional memory space.
- Online, i.e., can sort a list as it receives it.

Algorithm

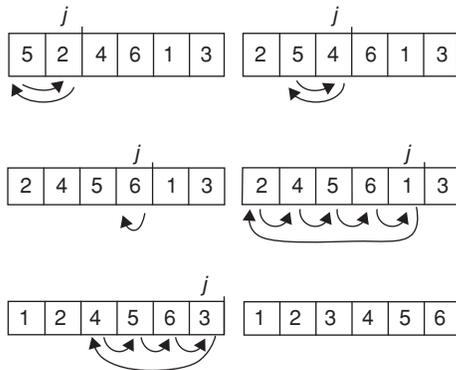
```
Insertion sort (A)
For (j ← 2) to length [A]
Do key ← A [j]
i ← j - 1;
While i > 0 and A [i] > key
{
    Do A [i + 1] ← A [i]
    i ← i - 1
}
A [i + 1] ← key
```

Every repetition of insertion sort removes an element from the input data, inserting it into the correct position in the already sorted list, until no input element remains. Sorting is typically done in-place. The resulting array after K iterations has the property where the first $k + 1$ entries are sorted. In each iteration the first remaining entry of the input is removed, inserted into the result at the correct position, with each element greater than X copied to the right as it is compared against X .

Performance

- The best case input is an array that is already sorted. In this case insertion sort has a linear running time (i.e., $\theta(n)$).
- The worst case input is an array sorted in reverse order. In this case every iteration of the inner loop will scan and shift the entire sorted subsection of the array before inserting the next element. For this case insertion sort has a quadratic running time ($O(n^2)$).
- The average case is also quadratic, which makes insertion sort impractical for sorting large arrays, however, insertion sort is one of the fastest algorithms for sorting very small arrays even faster than quick sort.

Example: Following figure shows the operation of insertion sort on the array $A = (5, 2, 4, 6, 1, 3)$. Each part shows what happens for a particular iteration with the value of j indicated. j indexes the ‘Current card’ being inserted.



Read the figure row by row. Elements to the left of $A[j]$ that are greater than $A[j]$ move one position to the right and $A[j]$ moves into the evacuated position.

SELECTION SORT

Selection sort is a sorting algorithm, specifically an in-place comparison sort. It has $O(n^2)$ complexity, making it inefficient on large lists.

The algorithm works as follows:

1. Find the minimum value in the list.
2. Swap it with the value in the first position.
3. Repeat the steps above for the remainder of the list (starting at the second position and advancing each time).

Analysis

Selection sort is not difficult to analyze compared to other sorting algorithms, since none of the loops depend on the data in the array selecting the lowest element requires scanning all n elements (this takes $n - 1$ comparisons) and then swapping it into the first position. Finding the next lowest element requires scanning the remaining $n - 1$ elements and so on, for $(n - 1) + (n - 2) + \dots + 2 + 1 = n(n - 1)/2 \in \theta(n^2)$ comparisons.

Each of these scans requires one swap for $n - 1$ elements (the final element is already in place).

Selection sort Algorithm

First, the minimum value in the list is found. Then, the first element (with an index of 0) is swapped with this value. Lastly, the steps mentioned are repeated for rest of the array (starting at the 2nd position).

Example 1: Here’s a step by step example to illustrate the selection sort algorithm using numbers.

- Original array:** 6 3 5 4 9 2 7
- 1st pass → 2 3 5 4 9 6 7 (2 and 6 were swapped)
 - 2nd pass → 2 3 5 4 9 6 7 (no swap)
 - 3rd pass → 2 3 4 5 9 6 7 (4 and 5 were swapped)
 - 4th pass → 2 3 4 5 6 9 7 (6 and 9 were swapped)
 - 5th pass → 2 3 4 5 6 7 9 (7 and 9 were swapped)
 - 6th pass → 2 3 4 5 6 7 9 (no swap)

Note: There are 7 keys in the list and thus 6 passes were required. However, only 4 swaps took place.

Example 2: Original array: LU, KU, HU, LO, SU, PU

- 1st pass → HU, KU, LU, LO, SU, PU
- 2nd pass → HU, KU, LU, LO, SU, PU
- 3rd pass → HU, KU, LO, LU, SU, PU
- 4th pass → HU, KU, LO, LU, SU, PU
- 5th pass → HU, KU, LO, LU, PU, SU

Note: There were 6 elements in the list and thus 5 passes were required. However, only 3 swaps took place.

BINARY SEARCH TREES

Search trees are data structures that support many dynamic, set operations, including SEARCH, MINIMUM, MAXIMUM, PREDECESSOR, SUCCESSOR, INSERT and DELETE. A search tree can be used as a dictionary and as a priority Queue. Operations on a binary search tree take time proportional to the height of the tree. For a complete binary tree with ‘ n ’ nodes, basic operations run in $\theta(\log n)$ worst-case time. If the tree is a linear chain of ‘ n ’ nodes, the basic operations take $\theta(n)$ worst-case time.

A binary search tree is organized, in a binary tree such a tree can be represented by a linked data structure in which each node is an object. In addition to key field, each node contains fields left, right and P that point to the nodes corresponding to its left child, its right child, and its parent,

respectively. If the child (or) parent is missing, the appropriate field contains the value NIL. The root node is the only node in the tree whose parent field is NIL.

Binary search tree property

The keys in a binary search tree are always stored in such a way as to satisfy the binary search tree property.

Let 'a' be a node in a binary search tree. If 'b' is a node in the left sub tree of 'a', $\text{key}[b] \leq \text{key}[a]$

If 'b' is a node in the right sub tree of 'a' then $\text{key}[a] \leq \text{key}[b]$.

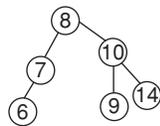


Figure 2 Binary search tree.

The binary search tree property allows us to print out all keys in a binary search tree in sorted order by a simple recursive algorithm called an inorder tree.

Algorithm

INORDER-TREE-WALK (root [T])
 INORDER-TREE-WALK (a)

1. If $a \neq \text{NIL}$
2. Then INORDER-TREE-WALK (left [a])
3. Print key [a]
4. INORDER-TREE-WALK (right [a])

It takes $\theta(n)$ time to walk an n -node binary search tree, since after the initial call, the procedure is called recursively twice for each node in the tree.

Let $T(n)$ denote the time taken by IN-ORDER-TREE-WALK, when it is called on the root of an n -node subtree.

INORDER-TREE-WALK takes a small, constant amount of time on an empty sub-tree (for the test $x \neq \text{NIL}$).

So $T(1) = C$ for some positive constant C .

For $n > 0$, suppose that INORDER-TREE-WALK is called on a node 'a' whose left subtree has k nodes and whose right subtree has $n - k - 1$ nodes.

The time to perform in order traversal is $T(n) = T(k) + T(n - k - 1) + d$.

For some positive constant 'd' that reflects the time to execute in-order (a), exclusive of the time spent in recursive calls $T(n) = (c + d)n + c$.

For $n = 0$, we have $(c + d)0 + c = T(0)$,

For $n > 0$,

$$\begin{aligned} T(n) &= T(k) + T(n - k - 1) + d \\ &= ((c + d)(k + c) + ((c + d)(n - k - 1) + c) + d \\ &= (c + d)n + c - (c + d) + c + d = (c + d)n + c \end{aligned}$$

HEAP SORT

Heap sort begins by building a heap out of the data set, and then removing the largest item and placing it at the end of

the partially sorted array. After removing the largest item, it reconstructs heap, removes the largest remaining item, and places, it in the next open position from the end of the partially sorted array. This is repeated until there are no items left in the heap and the sorted array is full. Elementary implementations require two arrays one to hold the heap and the other to hold the sorted elements.

- Heap sort inserts the input list elements into a binary heap data structure. The largest value (in a max-heap) or the smallest value (in a min-heap) is extracted until none remain, the value having been extracted in sorted order.

Example: Given an array of 6 elements: 15, 19, 10, 7, 17, 16, sort them in ascending order using heap sort.

Steps:

1. Consider the values of the elements as priorities and build the heap tree.
2. Start delete Max operations, storing each deleted element at the end of the heap array.

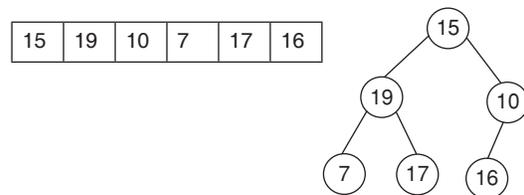
If we want the elements to be sorted in ascending order, we need to build the heap tree in descending order-the greatest element will have the highest priority.

1. Note that we use only array, treating its parts differently,
2. When building the heap-tree, part of the array will be considered as the heap, and the rest part-the original array.
3. When sorting, part of the array will be the heap and the rest part-the sorted array.

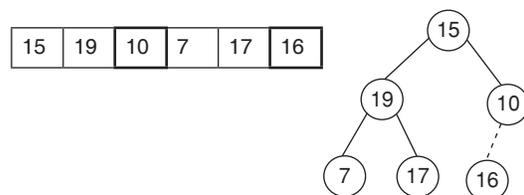
Here is the array: 15, 19, 10, 7, 17, 6.

Building the Heap Tree

The array represented as a tree, which is complete but not ordered.

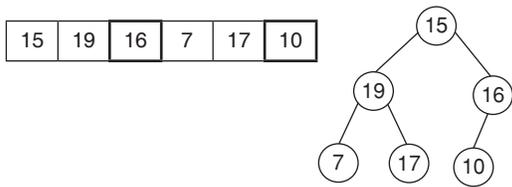


Start with the right most node at height 1 – the node at position 3 = size/2. It has one greater child and has to be percolated down.



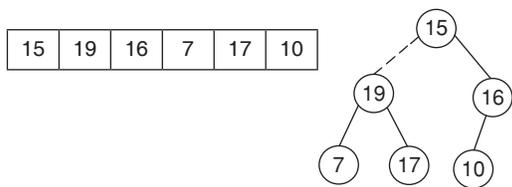
3.102 | Unit 3 • Algorithms

After processing array [3] the situation is:

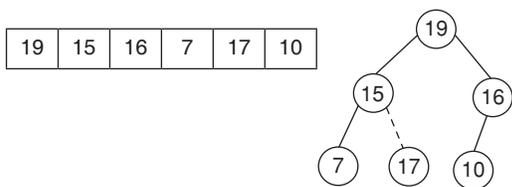


Next comes array [2]. Its children are smaller, so no percolation is needed.

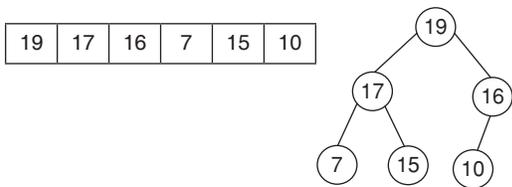
The last node to be processed is array[1]. Its left child is the greater of the children. The item at array [1] has to be percolated down to the left, swapped with array [2].



As a result:



The children of array [2] are greater and item 15 has to be moved down further, swapped with array [5].

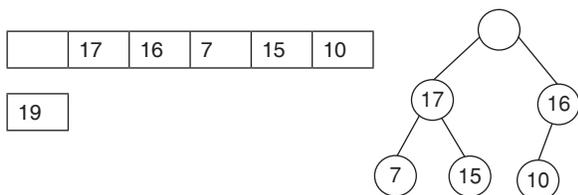


Now the tree is ordered, and the binary heap is built.

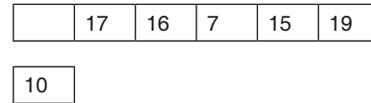
Sorting-performing Delete Max Operations

Delete the top element

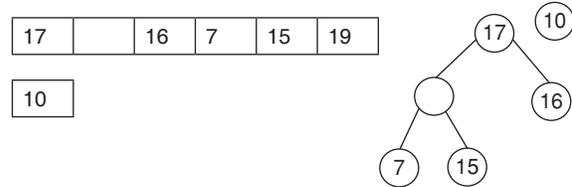
Store 19 in a temporary place, a hole is created at the top.



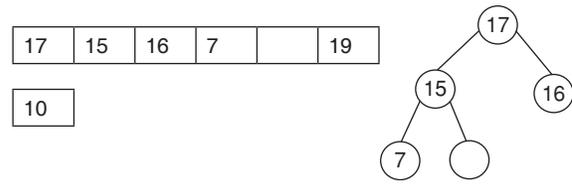
Swap 19 with the last element of the heap. As 10 will be adjusted in the heap, its cell will no longer be a part of the heap. Instead it becomes a cell from the sorted array



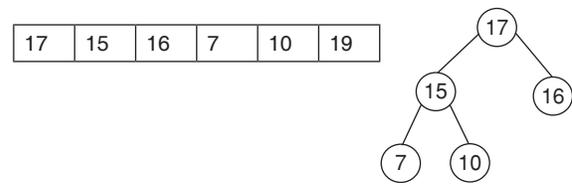
Percolate down the hole



Percolate once more (10 is less than 15, so it cannot be inserted in the previous hole)



Now 10 can be inserted in the hole



Repeat the step B till the array is sorted.

Heap sort analysis

Heap sort uses a data structure called (binary) heap binary, heap is viewed as a complete binary tree. An Array A that represents a heap is an object with 2 attributes: length [A], which is the number of elements in the array and heap size [A], the number of elements in the heap stored within array A.

No element past A [heap size [A]], where heap size [A] ≤ length [A], is an element of the heap.

There are 2 kinds of binary heaps:

1. Max-heaps
2. Min-heaps

In both kinds the values in the nodes satisfy a heap-property.

Max-heap property $A[\text{PARENT}(i)] \geq A[i]$

The value of a node is almost the value of its parent. Thus the largest element in a max-heap is stored at the root, and the sub tree rooted at a node contains values no larger than that contained at the node itself.

Min-heap property For every node 'i' other than the root $[\text{PARENT}(i)] \leq A[i]$. The smallest element in a min-heap is at the root.

Max-heaps are used in heap sort algorithm.

Min-heaps are commonly used in priority queues.

Basic operations on heaps run in time almost proportional to the height of the tree and thus take $O(\log n)$ time

- MAX-HEAPIFY procedure, runs in $O(\log n)$ time.
- BUILD-MAX-HEAP procedure, runs in linear time.
- HEAP SORT procedure, runs in $O(n \log n)$ time, sorts an array in place.
- MAX-HEAP-INSERT
HEAP-EXTRACT-MAX
HEAP-INCREASE-KEY
HEAP-MAXIMUM

All these procedures, run in $O(\log n)$ time, allow the heap data structure to be used as a priority queue.

- Each call to MAX-HEAPIFY costs $O(\log n)$ time, and there are $O(n)$ such calls. Thus, the running time is $O(n \log n)$
- The HEAPSORT procedure takes time $O(n \log n)$, since the call to BUILD-MAX-HEAP takes time $O(n)$ and each of the $(n-1)$ calls to MAX-HEAPIFY takes time $O(\log n)$.

Priority Queues

The most popular application of a heap is its use as an efficient priority queue.

A priority queue is a data structure for maintaining a set S of elements, each with an associated value called a key. A max-priority queue supports the following operations:

INSERT: INSERT(s, x) inserts the element x into the set S . This operation can be written as $S \leftarrow S \cup \{x\}$.

MAXIMUM: MAXIMUM(S) returns the element of S with the largest key

EXTRACT-MAX: EXTRACT-MAX(S) removes and returns the element of S with the largest key.

INCREASE-KEY: INCREASE-KEY(s, x, k) increases the value of element x 's key to the new value k , which is assumed to be at least as large as x 's current key value.

One application of max-priority queue is to schedule jobs on a shared computer.

EXERCISES

Practice Problems I

Directions for questions 1 to 15: Select the correct alternative from the given choices.

1. Solve the recurrence relation $T(n) = 2T(n/2) + k.n$ where k is constant then $T(n)$ is
(A) $O(\log n)$ (B) $O(n \log n)$
(C) $O(n)$ (D) $O(n^2)$
2. What is the time complexity of the given code?

```
Void f(int n)
{
  if (n > 0)
    f (n/2) ;
}
```

(A) $\theta(\log n)$ (B) $\theta(n \log n)$
(C) $\theta(n^2)$ (D) $\theta(n)$
3. The running time of an algorithm is represented by the following recurrence relation;

$$T(n) = \begin{cases} n & n \leq 3 \\ T\left[\frac{n}{3}\right] + cn & \text{otherwise} \end{cases}$$

What is the time complexity of the algorithm?

- (A) $\theta(n)$ (B) $\theta(n \log n)$
(C) $\theta(n^2)$ (D) $\theta(n^2 \log n)$

Common data for questions 4 and 5:

4. The following pseudo code does which sorting?

```
xsort [A, n]
for j ← 2 to n
do key ← A [ i ]
```

$i \leftarrow j - 1$

While $i > 0$ and $A [i] > \text{key}$

do $A [i + i] \leftarrow A [i]$

$i \leftarrow i - 1$

$A [i + 1] \leftarrow \text{key}$

- (A) Selection sort (B) Insertion sort
(C) Quick sort (D) Merge sort

5. What is the order of elements after 2 iterations of the above-mentioned sort on given elements?

8	2	4	9	3	6
---	---	---	---	---	---

(A)

2	4	9	8	3	6
---	---	---	---	---	---

(B)

2	4	8	9	3	6
---	---	---	---	---	---

(C)

2	4	6	3	8	9
---	---	---	---	---	---

(D)

2	4	6	3	8	9
---	---	---	---	---	---

Common data for questions 6 and 7:

6. The following pseudo code does which sort?
 1. If $n = 1$ done
 2. Recursively sort
 $A [1 \dots [n/2]]$ and
 $A [[n/2] + 1 \dots n]$
 3. Combine 2 ordered lists
(A) Insertion sort (B) Selection sort
(C) Merge sort (D) Quick sort

7. What is the complexity of the above pseudo code?
 (A) $\theta(\log n)$ (B) $\theta(n^2)$
 (C) $\theta(n \log n)$ (D) $\theta(2^n)$
8. Apply Quick sort on a given sequence 6 10 13 5 8 3 2
 11. What is the sequence after first phase, pivot is first element?
 (A) 5 3 2 6 10 8 13 11
 (B) 5 2 3 6 8 13 10 11
 (C) 6 5 13 10 8 3 2 11
 (D) 6 5 3 2 8 13 10 11
9. Selection sort is applied on a given sequence:
 89, 45, 68, 90, 29, 34, 17. What is the sequence after 2 iterations?
 (A) 17, 29, 68, 90, 45, 34, 89
 (B) 17, 45, 68, 90, 29, 34, 89
 (C) 17, 68, 45, 90, 34, 29, 89
 (D) 17, 29, 68, 90, 34, 45, 89
10. Suppose there are $\log n$ sorted lists of $\left\lfloor \frac{n}{\log n} \right\rfloor$ elements each. The time complexity of producing sorted lists of all these elements is: (hint: use a heap data structure)
 (A) $\theta(n \log \log n)$ (B) $\theta(n \log n)$
 (C) $\Omega(n \log n)$ (D) $\Omega(n^{3/2})$
11. If Divide and conquer methodology is applied on powering a Number X^n . Which one the following is correct?
 (A) $X^n = X^{n/2} \cdot X^{n/2}$
 (B) $X^n = X^{\frac{n-1}{2}} \cdot X^{\frac{n-1}{2}} \cdot X$
 (C) $X^n = X^{\frac{n+1}{2}} \cdot X^{\frac{n}{2}}$
 (D) Both (A) and (B)
12. The usual $\theta(n^2)$ implementation of insertion sort to sort an array uses linear search to identify the position, where an element is to be inserted into the already

sorted part of the array. If binary search is used instead of linear search to identify the position, the worst case running time would be.

- (A) $\theta(n \log n)$
 (B) $\theta(n^2)$
 (C) $\theta(n(\log n)^2)$
 (D) $\theta(n)$
13. Consider the process of inserting an element into a max heap where the max heap is represented by an array, suppose we perform a binary search on the path from the new leaf to the root to find the position for the newly inserted element, the number of comparisons performed is:
 (A) $\theta(\log n)$ (B) $\theta(\log \log n)$
 (C) $\theta(n)$ (D) $\theta(n \log n)$
14. Consider the following algorithm for searching a given number 'X' in an unsorted array $A[1 \dots n]$ having 'n' distinct values:
 (1) Choose an 'i' uniformly at random from $1 \dots n$
 (2) If $A[i] = x$
 Then stop
 else
 goto(1);
 Assuming that X is present in A, what is the expected number of comparisons made by the algorithm before it terminates.
 (A) n (B) n - 1
 (C) 2n (D) n/2
15. The recurrence equation for the number of additions $A(n)$ made by the divide and conquer algorithm on input size $n = 2^k$ is
 (A) $A(n) = 2A(n/2) + 1$ (B) $A(n) = 2A(n/2) + n^2$
 (C) $A(n) = 2A(n/4) + n^2$ (D) $A(n) = 2A(n/8) + n^2$

Practice Problems 2

Directions for questions 1 to 15: Select the correct alternative from the given choices.

1.

Input Array	Linear Search $W(n)$	Binary search $W(n)$
128 elements	128	8
1024 elements	1024	x

Find x value?

- (A) 10 (B) 11
 (C) 12 (D) 13
2. Choose the correct one
 (i) $\log n$ (ii) n
 (iii) $n \log n$ (iv) n^2

- (a) A result of cutting a problem size by a constant factor on each iteration of the algorithm.
 (b) Algorithm that scans a list of size 'n'.
 (c) Many divide and conquer algorithms fall in this category.
 (d) Typically characterizes efficiency of algorithm with two embedded loops.
 (A) i - b, ii - c, iii - a, iv - d
 (B) i - a, ii - b, iii - c, iv - d
 (C) i - c, ii - d, iii - a, iv - b
 (D) i - d, ii - a, iii - b, iv - c

3. Insertion sort analysis in worst case

- (A) $\theta(n)$
 (B) $\theta(n^2)$
 (C) $\theta(n \log n)$
 (D) $\theta(2^n)$

4. From the recurrence relation. Of merge sort $T(n) = 2T(n/2) + \theta(n)$. Which option is correct?
I. $n/2$ II. $2T$ III. $\theta(n)$
(a) Extra work (divide and conquer)
(b) Sub-problem size
(c) Number of sub-problems
(A) III – b, II – a, I – c (B) I – b, II – c, III – a
(C) I – a, II – c, III – b (D) I – c, II – a, III – b
5. What is the number of swaps required to sort ‘ n ’ elements using selection sort, in the worst case?
(A) $\theta(n)$ (B) $\theta(n^2)$
(C) $\theta(n \log n)$ (D) $\theta(n^2 \log n)$
6. In a binary max heap containing ‘ n ’ numbers, the smallest element can be found in time
(A) $O(n)$ (B) $O(\log n)$
(C) $O(\log \log n)$ (D) $O(1)$
7. What is the worst case complexity of sorting ‘ n ’ numbers using quick sort?
(A) $\theta(n)$ (B) $\theta(n \log n)$
(C) $\theta(n^2)$ (D) $\theta(n!)$
8. The best case analysis of quick sort is, if partition splits the array of size n into
(A) $n/2 : n/m$ (B) $n/2 : n/2$
(C) $n/3 : n/2$ (D) $n/4 : n/2$
9. What is the time complexity of powering a number, by using divide and conquer methodology?
(A) $\theta(n^2)$ (B) $\theta(n)$
(C) $\theta(\log n)$ (D) $\theta(n \log n)$
10. Which one of the following in-place sorting algorithm needs the minimum number of swaps?
(A) Quick sort (B) Insertion sort
(C) Selection sort (D) Heap sort
11. As the size of the array grows what is the time complexity of finding an element using binary search (array of elements are ordered)?
(A) $\theta(n \log n)$ (B) $\theta(\log n)$
(C) $\theta(n^2)$ (D) $\theta(n)$
12. The time complexity of heap sort algorithm is
(A) $n \log n$ (B) $\log n$
(C) n^2 (D) None of these.
13. As part of maintenance work, you are entrusted with the work of rearranging the library books in a shelf in a proper order, at the end of each day. The ideal choices will be _____.
(A) Heap sort (B) Quick sort
(C) Selection sort (D) Insertion sort
14. The value for which you are searching is called
(A) Binary value
(B) Search argument
(C) Key
(D) Serial value
15. To sort many large objects and structures it would be most efficient to _____.
(A) Place them in an array and sort the array
(B) Place the pointers on them in an array and sort the array
(C) Place them in a linked list and sort the linked list
(D) None of the above

PREVIOUS YEARS' QUESTIONS

1. What is the number of swaps required to sort n elements using selection sort, in the worst case? [2009]
(A) $\theta(n)$
(B) $\theta(n \log n)$
(C) $\theta(n^2)$
(D) $\theta(n^2 \log n)$
2. Which one of the following is the tightest upper bound that represents the number of swaps required to sort n numbers using selection sort? [2013]
(A) $O(\log n)$ (B) $O(n)$
(C) $O(n \log n)$ (D) $O(n^2)$
3. Let P be a quick sort program to sort numbers in ascending order using the first element as the pivot. Let t_1 and t_2 be the number of comparisons made by P for the inputs [1 2 3 4 5] and [4 1 5 3 2] respectively. Which one of the following holds? [2014]
(A) $t_1 = 5$ (B) $t_1 < t_2$
(C) $t_1 > t_2$ (D) $t_1 = t_2$
4. The minimum number of comparisons required to find the minimum and the maximum of 100 numbers is _____. [2014]
5. Suppose P, Q, R, S, T are sorted sequences having lengths 20, 24, 30, 35, 50 respectively. They are to be merged into a single sequence by merging together two sequences at a time. The number of comparisons that will be needed in the worst case by the optimal algorithm for doing this is _____. [2014]
6. You have an array of n elements. Suppose you implement quick sort by always choosing the central element of the array as the pivot. Then the tightest upper bound for the worst case performance is [2014]
(A) $O(n^2)$ (B) $O(n \log n)$
(C) $\theta(n \log n)$ (D) $O(n^3)$
7. What are the worst-case complexities of insertion and deletion of a key in a binary search tree? [2015]

- (A) $\theta(\log n)$ for both insertion and deletion
 (B) $\theta(n)$ for both insertion and deletion
 (C) $\theta(n)$ for insertion and $\theta(\log n)$ for deletion
 (D) $\theta(\log n)$ for insertion and $\theta(n)$ for deletion
8. The worst case running times of *Insertion sort*, *Merge sort* and *Quick sort*, respectively, are: [2016]
 (A) $\Theta(n \log n)$, $\Theta(n \log n)$, and $\Theta(n^2)$
 (B) $\Theta(n^2)$, $\Theta(n^2)$, and $\Theta(n \log n)$
 (C) $\Theta(n^2)$, $\Theta(n \log n)$, and $\Theta(n \log n)$
 (D) $\Theta(n^2)$, $\Theta(n \log n)$, and $\Theta(n^2)$
9. An operator delete(i) for a binary heap data structure is to be designed to delete the item in the i-th node. Assume that the heap is implemented in an array and i refers to the i-th index of the array. If the heap tree has depth d (number of edges on the path from the root to the farthest leaf), then what is the time complexity to re-fix the heap efficiently after the removal of the element? [2016]
- (A) $O(1)$ (B) $O(d)$ but not $O(1)$
 (C) $O(2^d)$ but not $O(d)$ (D) $O(d2^d)$ but not $O(2^d)$
10. Assume that the algorithms considered here sort the input sequences in ascending order. If the input is already in ascending order, which of the following are TRUE? [2016]
 I. Quicksort runs in $\Theta(n^2)$ time
 II. Bubblesort runs in $\Theta(n^2)$ time
 III. Mergesort runs in $\Theta(n)$ time
 IV. Insertion sort runs in $\Theta(n)$ time
 (A) I and II only (B) I and III only
 (C) II and IV only (D) I and IV only
11. A complete binary min - heap is made by including each integer in $[1, 1023]$ exactly once. The depth of a node in the heap is the length of the path from the root of the heap to that node. Thus, the root is depth 0. The maximum depth at which integer 9 can appear is _____. [2016]

ANSWER KEYS

EXERCISES

Practice Problems 1

1. B 2. A 3. A 4. B 5. B 6. C 7. C 8. B 9. A 10. B
 11. D 12. A 13. A 14. B 15. A

Practice Problems 2

1. B 2. B 3. B 4. B 5. A 6. A 7. C 8. B 9. C 10. C
 11. B 12. A 13. D 14. C 15. B

Previous Years' Questions

1. A 2. B 3. C 4. 148 5. 358 6. A 7. B 8. D 9. B 10. D
 11. 8

Chapter 3

Divide-and-conquer

LEARNING OBJECTIVES

- ☞ Divide-and-conquer
- ☞ Divide-and-conquer examples
- ☞ Divide-and-conquer technique
- ☞ Merge sort
- ☞ Quick sort
- ☞ Performance of quick sort
- ☞ Recurrence relation
- ☞ Searching
- ☞ Linear search
- ☞ Binary search

DIVIDE-AND-CONQUER

Divide-and-conquer is a top down technique for designing algorithms that consists of dividing the problem into smaller sub problems hoping that the solutions of the sub problems are easier to find and then composing the partial solutions into the solution of the original problem.

Divide-and-conquer paradigm consists of following major phases:

- Breaking the problem into several sub-problems that are similar to the original problem but smaller in size.
- Solve the sub-problem recursively (successively and independently)
- Finally, combine these solutions to sub-problems to create a solution to the original problem.

Divide-and-Conquer Examples

- Sorting: Merge sort and quick sort
- Binary tree traversals
- Binary Search
- Multiplication of large integers
- Matrix multiplication: Strassen's algorithm
- Closest-pair and Convex-hull algorithm

MERGE SORT

Merge sort is a sorting algorithm for rearranging lists (or any other data structure that can only be accessed sequentially, e.g., file streams) into a specified order.

Merge sort works as follows:

1. Divide the unsorted list into two sub lists of about half the size.
2. Sort each of the two sub lists.

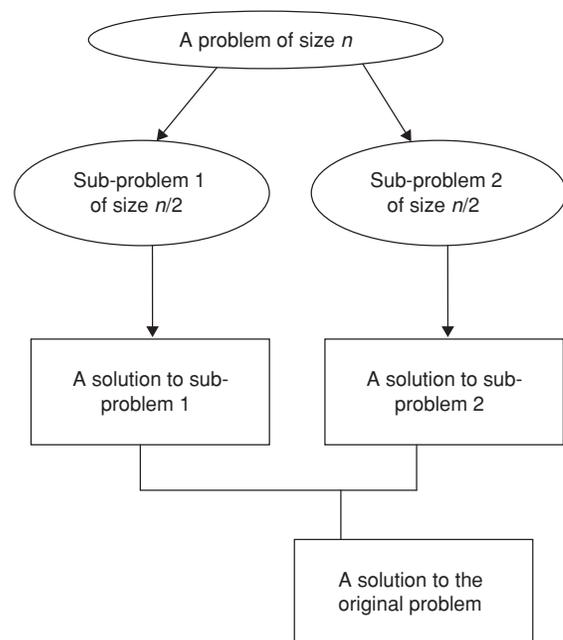


Figure 1 Divide-and-conquer technique.

3. Merge the two sorted sub lists back into one sorted list
4. The key of merge sort is merging two sorted lists into one, such that if we have 2 lists

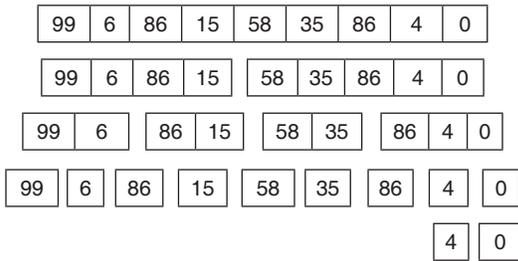
$X(x_1 \leq x_2 \leq x_3 \dots \leq x_m)$ and

$Y(y_1 \leq y_2 \leq y_3 \dots \leq y_n)$ the resulting list is $z(z_1 \leq z_2 \leq \dots \leq z_{m+n})$

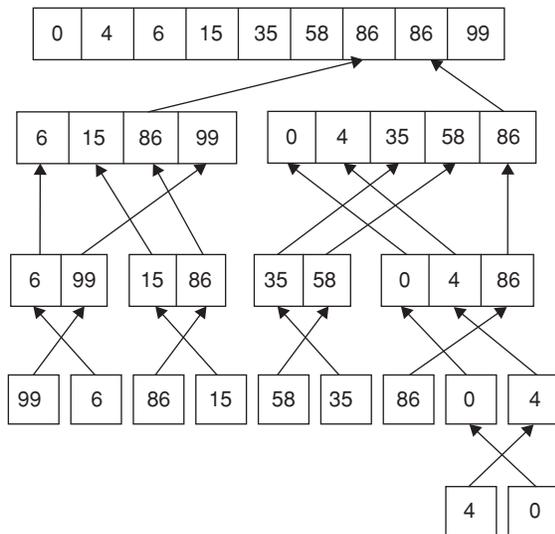
Example 1: $L_1 = \{3, 8, 9\}$, $L_2 = \{1, 5, 7\}$

Merge $(L_1, L_2) = \{1, 3, 5, 7, 8, 9\}$

Example 2:



Merge:



Implementing Merge Sort

Merging is done with a temporary array of the same size as the input array.

Pro: Faster than in-place since the temp array holds the resulting array until both left and right sides are merged into the temp array then the temp array is appended over the input array.

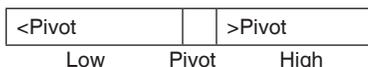
Con: The memory required is doubled. The double memory merge sort runs $O(N \log N)$ for all cases, because of its Divide-and-conquer approach.

$$T(N) = 2T(N/2) + N$$

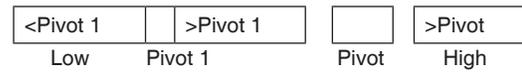
$$= O(N \log N)$$

QUICK SORT

Quick sort is an example of Divide-and-conquer strategy. In Quick sort we divide the array of items to be sorted into two partitions and then call the quick sort procedure recursively to sort the two partitions, i.e., we divide the problem into two smaller ones and conquer by solving the smaller ones. The conquer part of the quick sort routine looks like this



Make bold



Divide: Partition the array A [p - r] into 2 sub arrays A [p - q - 1] and A [q + 1 - r] such that each element of A [p - q - 1] is less than or equal to A[q], which is, in turn, less than or equal to each element of A [q + 1 - r]

Conquer: Sort the 2 sub arrays A [p - q - 1] and A [q + 1 - r] by recursive calls to quick sort.

Combine: Since the sub arrays are sorted inplace, no work is needed to combine them.

Sort left partition in the same way. For this strategy to be effective, the partition phase must ensure that the pivot, is greater than all the items in one part (the lower part) and less than all those in the other (upper) part. To do this, we choose a pivot element and arrange that all the items in the lower part are less than the pivot and all those in the upper part are greater than it. In the general case, the choice of pivot element is first element.

(Here $\lfloor \text{number of elements}/2 \rfloor$ is pivot)

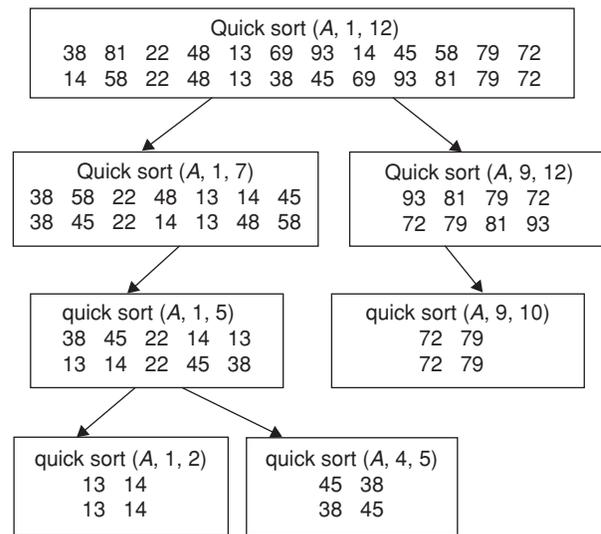


Figure 2 Tree of recursive calls to quick sort.

- Quick sort is a sorting algorithm with worst case running time $O(n^2)$ on an input array of n numbers. Inspite of this slow worst case running time, quick sort is often the best practical choice for sorting because it is efficient on the average: its expected running time is $O(n \log n)$ and the constants hidden in the O -notation are quite small
- Quick sort algorithm is fastest when the median of the array is chosen as the pivot element. This is because the resulting partitions are of very similar size. Each partition splits itself in two and thus the base case is reached very quickly.

Example: Underlined element is pivot.

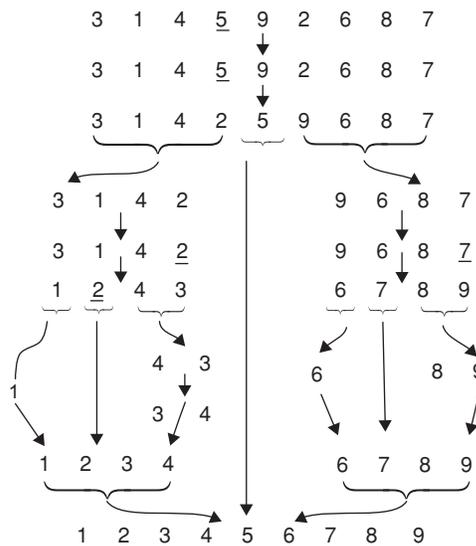


Figure 3 The ideal quick sort on a random array

Performance of Quick Sort

- Running time of quick sort depends on whether the partitioning is balanced or unbalanced, it depends on which elements are used for partitioning. If the partitioning is balanced, the algorithm runs asymptotically as fast as merge sort. If the partitioning is unbalanced, it runs as slowly as insertion sort.
- The worst case of quick sort occurs when the partitioning routine produces one sub-problem with $n - 1$ elements and one with '1' element. If this unbalanced partitioning arises in each recursive call, the partitioning costs $\theta(n)$ time.

Recurrence Relation

$$T(n) = T(n - 1) + T(1) + \theta(n)$$

$$(\because T(0) = \theta(1))$$

$$= T(n - 1) + \theta(n)$$

If we sum the costs incurred at each level of the recursion we get an arithmetic series, which evaluates to $\theta(n^2)$.

- Best case partitioning—PARTITION produces 2 sub-problems, each of size no more than $n/2$, since one is of size $\lfloor n/2 \rfloor$ and one of size $\lceil n/2 \rceil - 1$

The recurrence for the running time is then

$$T(n) \leq 2T(n/2) + \theta(n)$$

The above Recurrence relation has the solution $T(n) = O(n \log n)$ by case 2 of the master theorem.

- The average-case time of quick sort is much closer to the best than to the worst case

For example, that the partitioning algorithm always produces a 8-to-2 proportional split, which at first seems unbalanced. The Recurrence relation will be

$$T(n) \leq T(8n/10) + T(2n/10) + cn$$

The recursion tree for this recurrence has cost 'cn' at every level, until a boundary condition is reached at depth $\log_{10} n = \theta(\log n)$. The recursion terminates at depth $\log_{10/8} n = \theta(\log n)$. The total cost of quick sort is $O(n \log n)$

SEARCHING

Two searching techniques are:

- Linear search
- Binary search

Linear Search

Linear search (or) sequential search is a method for finding a particular value in list that consists of checking every one of its elements, one at a time and in sequence, until the desired one is found. Linear search is a special case of brute force search. Its worst case cost is proportional to the number of elements in the list.

Implementation

```
boolean linear search (int [ ] arr, int target)
{
    int i = 0;
    while (i < arr. length) {
        if (arr [i] == target){
            return true;
        }
        ++ i;
    }
    return false;
}
```

Example:

Consider the array

10	7	1	3	-4	2	20
----	---	---	---	----	---	----

Search for 3

10	7	1	3	-4	2	20
			3?			

Move to next element

10	7	1	3	-4	2	20
				3?		

Move to next element

10	7	1	3	-4	2	20
					3?	

Move to next element

10	7	1	3	-4	2	20
						3?

Element found; stop the search.

Binary Search

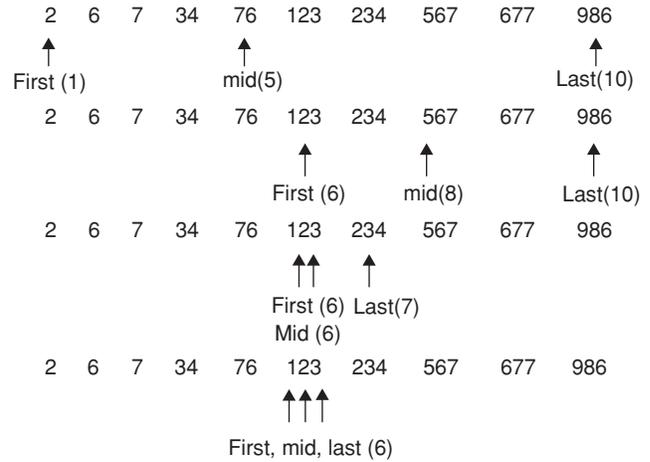
A binary search algorithm is a technique for finding a particular value in a linear array, by ruling out half of the data at each

step; a binary search finds the median, makes comparison, to determine whether the desired value comes before or after it, and then searches the remaining half in the same manner. A binary search is an example of Divide-and-conquer algorithm.

Implementation

```
function binary search (a, value, left, right)
{
  if right < left
    return not found
  mid: = floor ((right -left)/2) + left
  if a [mid] = value
    return mid
  if value < a[mid]
    return binary search (a, value, left, mid -1)
  else return binary search (a, value, mid + 1, right)
}
```

Example: Value being searched 123



EXERCISES

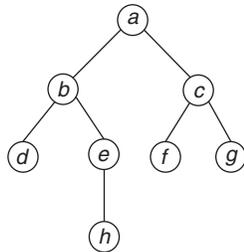
Practice Problems I

Directions for questions 1 to 15: Select the correct alternative from the given choices.

- How many comparisons are required to search an item 89 in a given list, using Binary search?

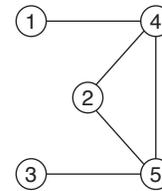
4	8	19	25	34	39	45	48	66	75	89	95
[0]	[1]	[2]	[3]	[4]	[5]	[6]	[7]	[8]	[9]	[10]	[11]

- (A) 3 (B) 4
(C) 5 (D) 6
2. Construct a Binary search tree with the given list of elements:
300, 210, 400, 150, 220, 370, 450, 100, 175, 215, 250
Which of the following is a parent node of element 250?
(A) 220 (B) 150
(C) 370 (D) 215
3. What is the breadth first search order of the given tree?

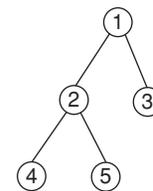


- (A) acbhdefg (B) abcdefgh
(C) adbcefg (D) aebcdfgh

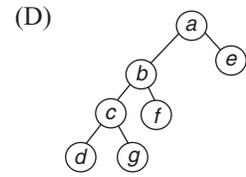
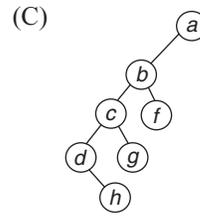
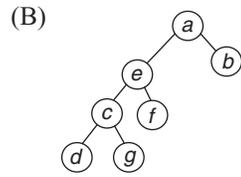
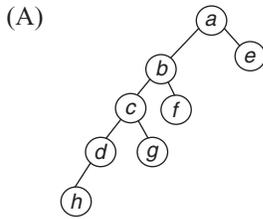
- What is the depth first search order of the given graph?



- (A) 14325 (B) 12435
(C) 14253 (D) 12354
5. When pre-order traversal is applied on a given tree, what is the order of elements?



- (A) 1-2-4-5-3 (B) 1-4-2-5-3
(C) 1-2-4-3-5 (D) 1-2-3-4-5
6. What is the order of post-order traversal and in-order traversals of graph given in the above question?
(A) 4-2-5-1-3 and 4-5-2-3-1
(B) 4-5-2-3-1 and 4-2-5-1-3
(C) 4-5-2-1-3 and 4-2-5-1-3
(D) 4-5-2-3-1 and 4-2-5-3-1



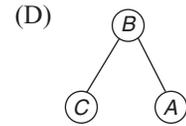
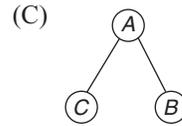
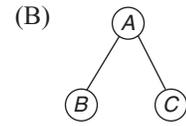
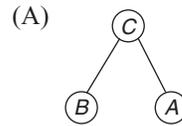
Practice Problems 2

Directions for questions 1 to 15: Select the correct alternative from the given choices.

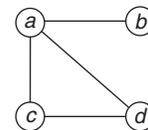
- Which of the following algorithm design technique is used in finding all pairs of shortest distances in a graph?
 - Divide-and-conquer
 - Greedy method
 - Back tracking
 - Dynamic programming
- Let LASTPOST, LASTIN and LASTPRE denote the last vertex visited in a post-order, in-order and pre-order traversals respectively of a complete binary tree. Which of the following is always true?
 - LASTIN = LASTPOST
 - LASTIN = LASTPRE
 - LASTPRE = LASTPOST
 - LASTIN = LASTPOST = LASTPRE
- Match the following:

X : Depth first search
 Y : Breadth first search
 Z : Sorting
 a : Heap
 b : Queue
 c : Stack

 - X – a, Y – b, Z – c
 - X – c, Y – a, Z – b
 - X – c, Y – b, Z – a
 - X – a, Y – c, Z – b
- Let G be an undirected graph, consider a depth first traversal of G , and let T be the resulting DFS Tree. Let ' U ' be a vertex in ' G ' and let ' V ' be the first new (unvisited) vertex visited after visiting ' U ' in the traversal. Which of the following is true?
 - $\{U, V\}$ must be an edge in G and ' U ' is a descendant of V in T .
 - $\{U, V\}$ must be an edge in ' G ' and V is a descendant of ' U ' in T .
 - If $\{U, V\}$ is not an edge in ' G ' then ' U ' is a leaf in T .
 - if $\{U, V\}$ is not an edge in G then U and V must have the same parent in T .
- Identify the binary tree with 3 nodes labeled A, B and C on which preorder traversal gives the sequence C, B, A .



- Consider an undirected unweighted graph G . Let a breadth first traversal of G be done starting from a node r . Let $d(r, u)$ and $d(r, v)$ be the lengths of the shortest paths from r to u and v respectively in ' G '. If u is visited before v during the breadth first travel, which of the following is correct?
 - $d(r, u) < d(r, v)$
 - $d(r, u) > d(r, v)$
 - $d(r, u) \leq d(r, v)$
 - None of these
- In a complete 5-ary tree, every internal node has exactly 5 children. The number of leaves in such a tree with '3' internal nodes are:
 - 15
 - 20
 - 13
 - Can't predicted
- Which of the following algorithm is single pass that is they do not traverse back up the tree for search, create, insert etc.
 - Depth first search
 - Pre-order traversal
 - B-tree traversal
 - Post-order traversal
- Which of the following is the adjacency matrix of the given graph?



(A)
$$\begin{bmatrix} 0 & 1 & 1 & 1 \\ 1 & 0 & 0 & 0 \\ 1 & 0 & 0 & 1 \\ 1 & 0 & 1 & 0 \end{bmatrix}$$

(B)
$$\begin{bmatrix} 1 & 1 & 1 & 1 \\ 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 1 \\ 0 & 0 & 1 & 0 \end{bmatrix}$$

(C)
$$\begin{bmatrix} 1 & 1 & 1 & 1 \\ 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 1 \\ 1 & 0 & 1 & 0 \end{bmatrix}$$

(D)
$$\begin{bmatrix} 1 & 1 & 1 & 1 \\ 0 & 0 & 0 & 0 \\ 1 & 0 & 0 & 1 \\ 0 & 0 & 1 & 0 \end{bmatrix}$$

PREVIOUS YEARS' QUESTIONS

1. Which one of the following is the tightest upper bound that represents the time complexity of inserting an object into a binary search tree of n nodes? [2013]

(A) $O(1)$ (B) $O(\log n)$
(C) $O(n)$ (D) $O(n \log n)$

2. Consider a rooted n node binary tree represented using pointers. The best upper bound on the time required to determine the number of sub trees having exactly 4 nodes is $O(n^a \log^b n)$. The value of $a + 10b$ is _____ [2014]

3. Which one of the following is the recurrence equation for the worst case time complexity of the Quicksort algorithm for sorting $n(\geq 2)$ numbers? In the recurrence equations given in the options below, c is a constant. [2015]

(A) $T(n) = 2T(n/2) + cn$
(B) $T(n) = T(n-1) + T(1) + cn$
(C) $T(n) = 2T(n-1) + cn$
(D) $T(n) = T(n/2) + cn$

4. Suppose you are provided with the following function declaration in the C programming language.

```
int partition (int a [ ], int n);
```

The function treats the first element of $a []$ as a pivot, and rearranges the array so that all elements less than or equal to the pivot is in the left part of the array, and all elements greater than the pivot is in the right part in addition, it moves the pivot so that the pivot is the last element of the left part. The return value is the number of elements in the left part.

The following partially given function in the C programming language is used to find the k th smallest element in an array $a []$ of size n using the partition function. We assume $k \ll n$.

```
int kth_smallest (int a [ ], int n, int k) [2015]
```

```
{
    int left_end = partition(a, n);
    if (left_end+1 == k) {
        return a [left_end];
    }
    if (left_end+1 > k) {
        return kth_smallest (_____);
    } else {
        return kth_smallest (_____);
    }
}
```

The missing argument lists are respectively

- (A) $(a, \text{left_end}, k)$ and $(a+\text{left_end}+1, n-\text{left_end}-1, k-\text{left_end}-1)$
(B) $(a, \text{left_end}, k)$ and $(a, n-\text{left_end}-1, k-\text{left_end}-1)$

(C) $(a+\text{left_end}+1, n-\text{left_end}-1, k-\text{left_end}-1)$ and $(a, \text{left_end}, k)$

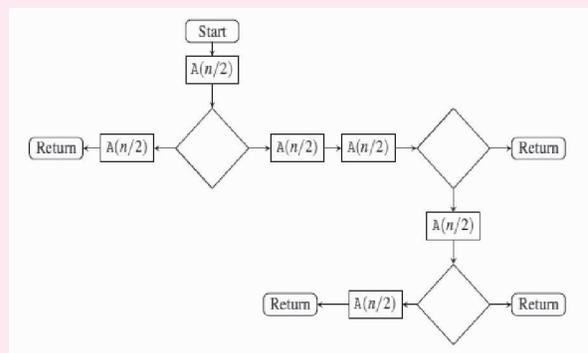
(D) $(a, n-\text{left_end}-1, k-\text{left_end}-1)$ and $(a, \text{left_end}, k)$

5. Assume that a mergesort algorithm in the worst case takes 30 seconds for an input of size 64. Which of the following most closely approximates the maximum input size of a problem that can be solved in 6 minutes? [2015]

(A) 256 (B) 512
(C) 1024 (D) 2048

6. The given diagram shows the flowchart for a recursive function $A(n)$. Assume that all statements, except for the recursive calls, have $O(1)$ time complexity. If the worst case time complexity of this function is $O(n^\alpha)$, then the least possible value (accurate up to two decimal positions) of α is _____. [2016]

Flowchart for Recursive Function $A(n)$



7. Let A be an array of 31 numbers consisting of a sequence of 0's followed by a sequence of 1's. The problem is to find the smallest index i such that $A[i]$ is 1 by probing the minimum number of locations in A . The *worst case* number of probes performed by an *optimal* algorithm is _____. [2017]

8. Match the algorithms with their time complexities:

Algorithm	Time complexity
(P) Towers of Hanoi with n disks	(i) $\Theta(n^2)$
(Q) Binary search given n sorted numbers	(ii) $\Theta(n \log n)$
(R) Heap sort given n numbers at the worst case	(iii) $\Theta(2^n)$
(S) Addition of two $n \times n$ matrices	(iv) $\Theta(\log n)$

- (A) P \rightarrow (iii), Q \rightarrow (iv), R \rightarrow (i), S \rightarrow (ii)
(B) P \rightarrow (iv), Q \rightarrow (iii), R \rightarrow (i), S \rightarrow (ii)
(C) P \rightarrow (iii), Q \rightarrow (iv), R \rightarrow (ii), S \rightarrow (i)
(D) P \rightarrow (iv), Q \rightarrow (iii), R \rightarrow (ii), S \rightarrow (i)

[2017]

ANSWER KEYS**EXERCISES****Practice Problems 1**

1. A 2. A 3. B 4. C 5. A 6. B 7. B 8. A 9. C 10. A
11. C 12. D 13. A 14. C 15. A

Practice Problems 2

1. B 2. B 3. C 4. B 5. A 6. D 7. C 8. C 9. A 10. B
11. B 12. C 13. A 14. B 15. A

Previous Years' Questions

1. C 2. 1 3. B 4. A 5. B 6. 2.2 to 2.4 7. 5 8. C

Chapter 4

Greedy Approach

LEARNING OBJECTIVES

- ☞ Greedy approach
- ☞ Knapsack problem
- ☞ Fractional knapsack problem
- ☞ Spanning trees
- ☞ Prim's algorithm
- ☞ Kruskal's algorithm
- ☞ Tree and graph traversals
- ☞ Back tracking
- ☞ Graph traversal
- ☞ Breadth first traversal
- ☞ Depth first search
- ☞ Huffman codes
- ☞ Task-scheduling problem
- ☞ Sorting and order statistics
- ☞ Simultaneous minimum and maximum
- ☞ Graph algorithms

GREEDY APPROACH

In a greedy method, we attempt to construct an optimal solution in stages.

- At each stage we make a decision that appears to be the best (under some criterion) at the time.
- A decision made at one stage is not changed in a later stage, so each decision should assure feasibility.
- Some problems that use greedy approach are:
 1. Knapsack problem
 2. Minimum spanning tree
 3. Prims algorithm
 4. Kruskals algorithm

KNAPSACK PROBLEM

The knapsack problem is a problem in combinatorial optimization: given a set of items, each with a weight and a value, determine the number of each item to include in a collection so that the total weight is less than or equal to a given limit and the total value is as large as possible. We have n kinds of items, 1 through n . Each kind of item i has a value V_i and a weight W_i , usually assume that all values and weights are non-negative. The maximum weight that we can carry in the bag is W .

Solved Examples

Example 1: (Making change)

Problem:	Accept n dollars, to return a collection of coins with a total value of n dollars.
Configuration:	A collection of coins with a total value of n .
Objective function:	Minimize number of coins returned.
Greedy solution:	Always return the largest coin you can.

- Coins are valued \$.30, \$.020, \$.05, \$.01 use a greedy choice property and make \$.40 by using 3 coins.

Solution: $\$0.30 + \$0.05 + \$0.05 = \0.40

Fractional Knapsack Problem

Given: A set S of n items, with each item i having

- $b_i - a$ positive benefit
- $w_i - a$ positive weight

Goal: Choose items with maximum total benefit but with weight at most W .

If we are allowed to take fractional amounts, then this is the fractional knapsack problem.

- In this case, let x_i denote the amount we take of item i .

• Objective: Maximize $\sum_{i \in S} b_i(x_i/w_i)$

• Constraint: $\sum_{i \in S} x_i \leq W$

Example 2:

Items					
Weight	4 ml	8 ml	2 ml	6 ml	1 ml
Benefit	\$12	\$32	\$40	\$30	\$50
Value (\$ per ml)	3	4	20	5	50



Solution: 1 ml of 5, 2 ml of 3, 6 ml of 4, 1 ml of 2

- Greedy choice: Keep taking item with highest value (benefit to weight ratio).
- Correctness: suppose there is a better solution, there is an item i with higher value than a chosen item j . (i.e., $v_j < v_i$). If we replace some j with i , we get a better solution.

Thus there is no better solution than the greedy one.

$$N = 3, m = 20$$

$$(P_1, P_2, P_3) = (25, 24, 15)$$

$$(W_1, W_2, W_3) = (18, 15, 10)$$

Example 3:

	X_1	X_2	X_3	$\Sigma W_i X_i$	$\Sigma P_i X_i$
1.	1/2	1/3	1/4	$9 + 5 + 2.5 = 16.5$	$12.5 + 8 + 3.75 = 24.25$
2.	1	2/15	0	$18 + 2 + 0 = 20$	$25 + 3.2 + 0 = 28.2$
3.	0	2/3	1	$0 + 10 + 10 = 20$	$0 + 16 + 15 = 31$
4.	0	1	1/2	$0 + 15 + 5 = 20$	$0 + 24 + 7.5 = 31.5$

(1), (2), (3), (4) are feasible ones but (4) is the optimum solution.

SPANNING TREES

A spanning tree of a graph is just a sub graph that contains all the vertices and is a tree. A graph may have many spanning trees.

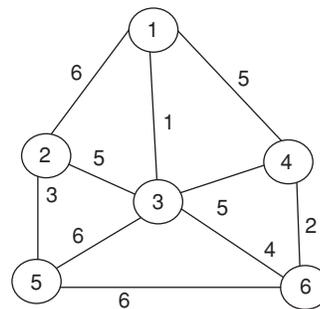
- A sub graph that spans (reaches out to) all vertices of a graph is called a spanning sub graph.
- A sub graph that is a tree and that spans all vertices of the original graph is called a spanning tree.
- Among all the spanning trees of a weighted and connected graph, the one (possibly more) with the least total weight is called a Minimum Spanning Tree (MST).

PRIM'S ALGORITHM

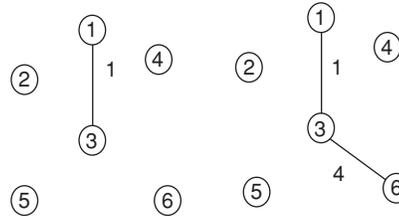
Prim's algorithm is a greedy algorithm that finds a minimum spanning tree for a connected weighted undirected graph. This means it finds a subset of the edges that forms a tree that includes every vertex, where the total weight of all the edges in the tree is minimized. The algorithm continuously increases the size, of a tree, one edge at a time starting with a tree consisting of a single vertex, until it spans all vertices.

- Using a simple binary heap data structure and an adjacency list representation, prim's algorithm can be shown to run in time $O(E \log V)$ where E is the number of edges and V is the number of vertices.

Example:

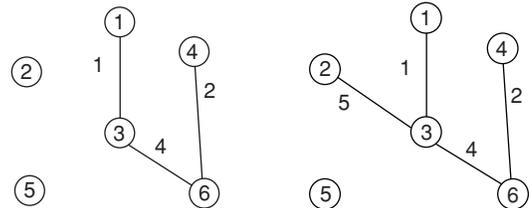


Start:



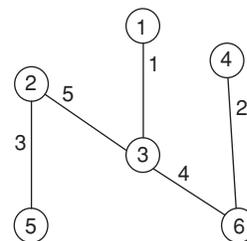
Iteration 1: $U = \{1, 3\}$

Iteration 2: $U = \{1, 3, 6\}$



Iteration 3: $U = \{1, 3, 6, 4\}$

Iteration 4: $U = \{1, 3, 6, 4, 2\}$



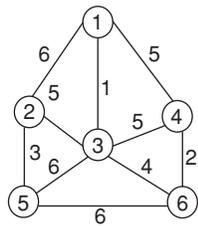
Iteration 5: $U = \{1, 3, 6, 4, 2, 5\}$

Figure 1 An example graph for illustrating prim's algorithm.

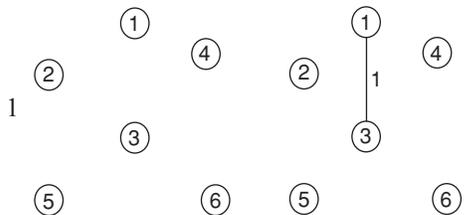
KRUSKAL'S ALGORITHM

Like prim's algorithm, Kruskal's algorithm also constructs the minimum spanning tree of a graph by adding edges to the spanning tree one-by-one. At all points, during its execution the set of edges selected by prim's algorithm forms exactly one tree. On the other hand the set of edges selected by Kruskal's algorithm forms a forest of trees. Kruskal's algorithm is conceptually simple. The edges are selected and added to the spanning tree in increasing order of their weights. An edge is added to the tree only if it does not create a cycle.

Example:

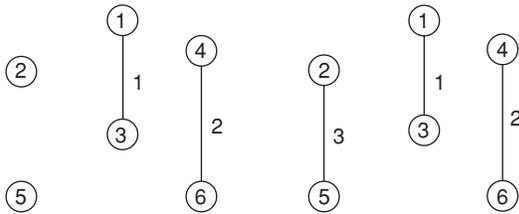


Start:



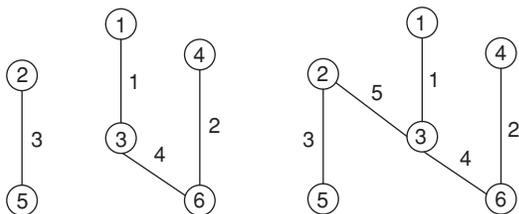
Initial configuration

Setp 1: choose (1, 3)



Setep 2: choose (4, 6)

Setep 3: choose (2, 5)



Setep 4: choose (3, 6)

Setep 6: choose (4, 3)

TREE AND GRAPH TRAVERSALS

Back Tracking

Backtracking is a general algorithm technique that considers searching every possible combination in order to solve an optimization problem.

Backtracking is also known as depth first search (or branch and bound). Backtracking is an important tool for solving constraint satisfaction problems, such as crosswords, verbal arithmetic, sudoku and many other puzzles. It is often the more convenient technique for parsing, for the knapsack problem and other combinational optimization problems.

- The advantage of backtracking algorithm is that they are complete, that is they are guaranteed to find every solution to every possible puzzle.

Graph Traversal

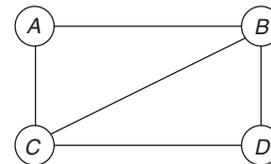
To traverse a graph is to process every node in the graph exactly once, because there are many paths leading from one node to another, the hardest part about traversing a graph is making sure that you do not process some node twice. There are general solutions to this difficulty.

1. When you first encounter a node, mark it as REACHED. When you visit a node, check if it is marked REACHED, if it is, just ignore it. This is the method our algorithms will use.
2. When you process a node, delete it from the graph. Deleting the node causes the deletion of all the arcs that lead to the node, so it will be impossible to reach it more than once.

General traversal strategy

1. Mark all nodes in the graph as NOT REACHED,
2. Pick a starting node, mark it as REACHED, and place it on the READY list.
3. Pick a node on the READY list. Process it remove it from READY. Find all its neighbors, those that are NOT REACHED should be marked as REACHED and added to READY.
4. Repeat 3 until READY is empty.

Example:



Step I: A = B = C = D = NOT REACHED

Step II: READY = {A} . A = REACHED

Step III: Process A. READY = {B, C}.

B = C = REACHED

Step IV: Process C. READY = {B, D}.

D = REACHED

Step V: Process B. READY = {D}

Step VI: Process D. READY = { }

The two most common traversal patterns are

- Breadth first traversal
- Depth first traversal

Breadth First Traversal

In breadth first traversal, READY is a QUEUE, not an arbitrary list. Nodes are processed in the order they are reached (FIFO), this has the effect of processing nodes according to their distance from the initial node. First, the initial node is processed. Then all its neighbors are processed. Then all of the neighbors etc.

- Since a graph has no root, we must specify the vertex at which to start the traversal.
- Breadth first tree traversal first visits all the nodes at depth zero (i.e., the root) then all the nodes at depth 1, and so on.

Procedure

First, the starting vertex is enqueued. Then, the following steps are repeated until the queue is empty.

1. Remove the vertex at the head of the queue and call it vertex.
2. Visit vertex
3. Follow each edge emanating from vertex to find the adjacent vertex and call it 't_o'. If 't_o' has not already been put into the queue, enqueue it.

Notice, that a vertex can be put into the queue at most once. Therefore, the algorithm must some how keep track of the vertices that have been enqueued.

Procedure for BFS for undirected graph G(V, E)

To perform BFS over a graph, the data structures required are queue (Q) and the visited set (Visited), 'V' is the starting vertex.

Procedure for BFS(V)

Steps

1. Visit the vertex 'V'
2. Enqueue the vertex V
3. while (Q is not Empty)
 - (i) V = dequeue ();
 - (ii) for all vertices ϑ adjacent to V
 - (a) if not visited (ϑ)
 - Enqueue (ϑ)
 - Visit the vertex ' ϑ '
 - end if.
 - end for
- end while
- Stop

Example:

- (A) Unexplored vertex
- (A) Visited vertex

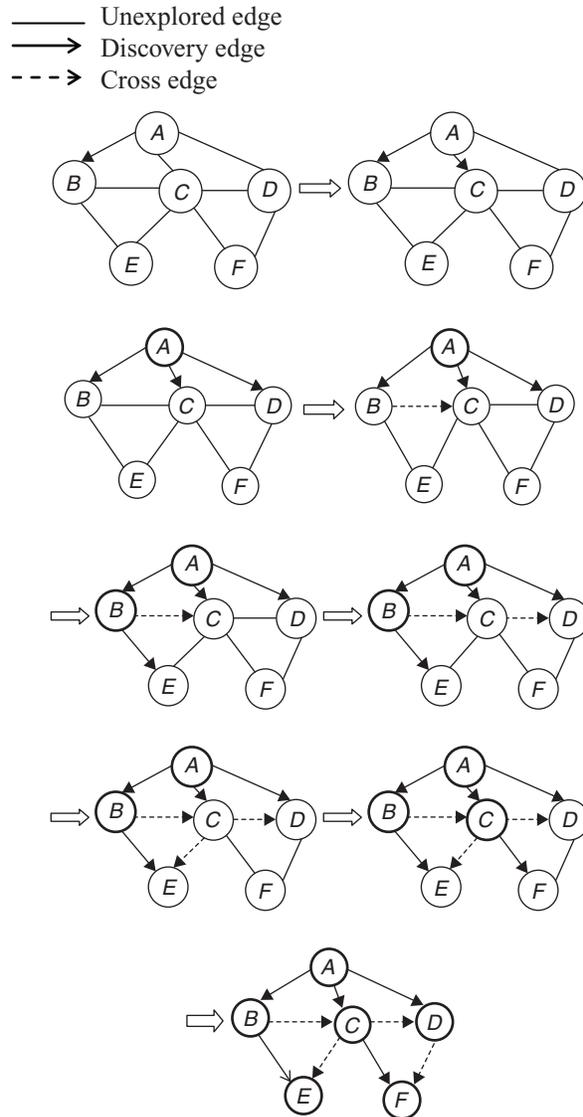


Figure 2 Breadth-first search

Depth First Search

A depth first traversal of a tree always starts at the root of the tree. Since a graph has no root, when we do a depth first traversal, we must specify the vertex at which to begin. A depth first traversal of a tree visits a node and then recursively visits the sub trees of that node similarly, depth first traversal of a graph visits a vertex and then recursively visits all the vertices adjacent to that node. A graph may contain cycles, but the traversal must visit every vertex at most once.

The solution to the problem is to keep track of the nodes that have been visited.

Procedure for DFS for undirected graph G(V, E)

To perform DFS over a graph, the data structures required are stack (S) and the list (visited), 'V' is the start vertex.

Procedure for DFS(*V*)

Steps

1. push the start vertex '*V*' into the stack *S*
2. while (*S* is not empty)
 - (i) pop a vertex *V*
 - (ii) if '*V*' is not visited
 - (a) visit the vertex
 - (b) Store '*V*' in visited
 - (c) push all the adjacent vertices of '*V*' in to visited
 - (iii) End if
3. End while
4. Stop.

Example:

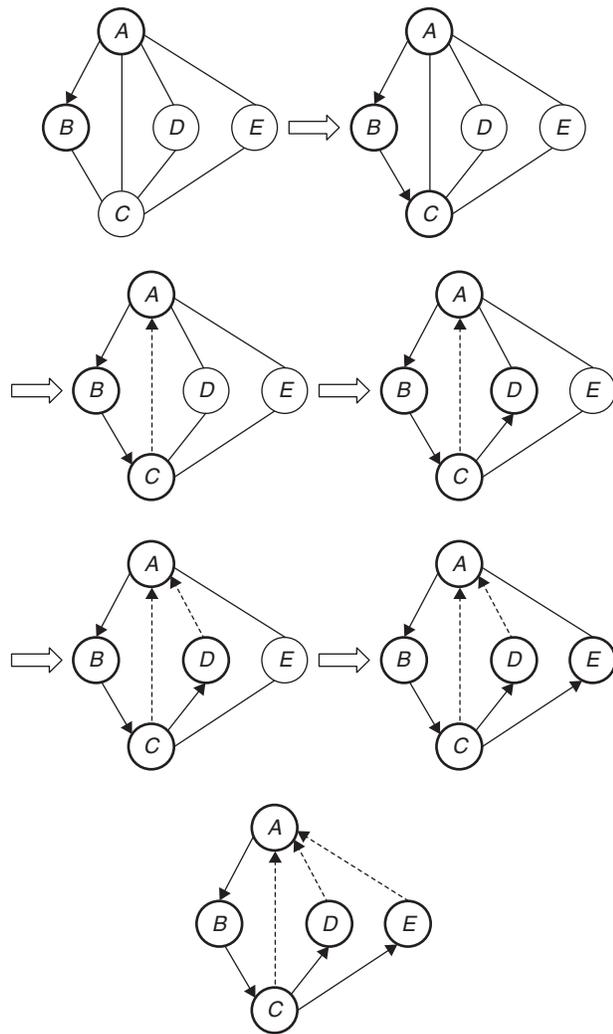
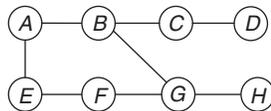


Figure 3 Depth first search

• Let us compare two traversal orders on the following graph:



Initial steps:

READY = [*A*], process A. READY = [*B*, *E*], process *B*.

It is at this point that two traversal strategies differ. Breadth first adds B's neighbors to the back of READY, depth first adds them to the front.

Breadth first

- READY = [*E*, *C*, *G*]
- Process *E*. READY = [*C*, *G*, *F*]
- Process *C*. READY = [*G*, *F*, *D*]
- Process *G*. READY = [*F*, *D*, *H*]
- Process *F*. READY = [*D*, *H*]
- Process *D*. READY = [*H*]
- Process *H*. READY = []

Depth First

- READY = [*C*, *G*, *E*]
- Process *C*. READY = [*D*, *G*, *E*]
- Process *D*. READY = [*G*, *E*]
- Process *G*. READY = [*H*, *F*, *E*]
- Process *H*. READY = [*F*, *E*]
- Process *F*. READY = [*E*]
- Process *E*. READY = []

CONNECTED COMPONENTS

A graph is said to be connected if every pair of vertices in the graph are connected. A connected component is a maximal connected sub graph of '*G*'. Each vertex belongs to exactly one connected component as does each edge.

- A graph that is not connected is naturally and obviously decomposed into several connected components (Figure 4). Depth first search does this handily. Each restart of the algorithm marks a new connected component.
- The directed graph in (Figure 5) is "Connected" Part of it can be "Pulled apart" (so to speak, without "breaking" any edges).
- Meaningful way to define connectivity in directed graph is:

'Two nodes *U* and *V* of a directed graph $G = (V, E)$ connected if there is a path from *U* to *V*', and one from *V* to *U*. This relation between nodes is reflective, symmetric and transitive. As such, it partitions *V* into disjoint sets, called the strongly connected components of the graph. In the directed graph of figure 2 there are four strongly connected components.

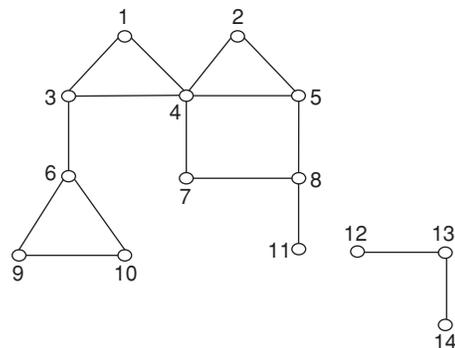


Figure 4 Undirected graph.

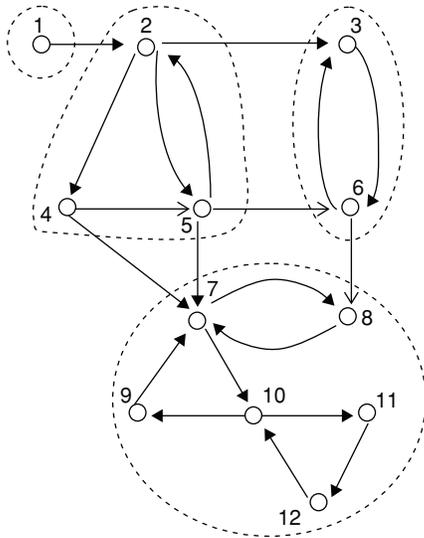
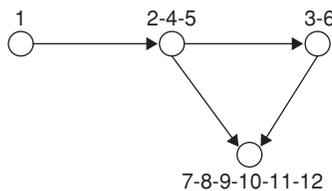


Figure 5 A directed graph and its strongly connected components

If we shrink each of these strongly connected components down to a single node and draw an edge between two of them if there is an edge from some node in the first to some node in the second, the resulting directed graph has to be a directed acyclic graph (DAG) – it has no cycles (figure 6). The reason is simple.

A cycle containing several strongly connected components would merge them all to a single strongly connected component.



Every directed graph is a DAG of its strongly connected components.

HUFFMAN CODES

For compressing data, a very effective and widely used technique is Huffman coding. We consider the data to be a sequence of characters. Huffmans’s greedy algorithm uses a table of the frequencies of occurrence of the characters to build up an optimal way of representing each character as a binary string.

Example: Suppose we have a 1,00,000 – character data file, that we wish to store compactly. The characters in the file occur with the frequencies given below:

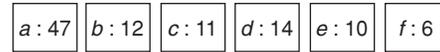
Character	a	b	c	d	e	f
Frequency	47	12	11	14	10	6

Solution: Two methods are used for compression of data are:

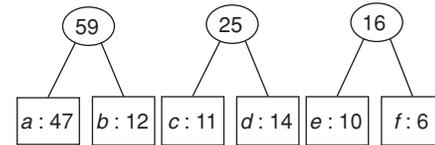
Fixed Length Coding

- Arrange all the characters in sequence (no particular order is followed)
- $a = 47, b = 12, c = 11, d = 14, e = 10, f = 6$

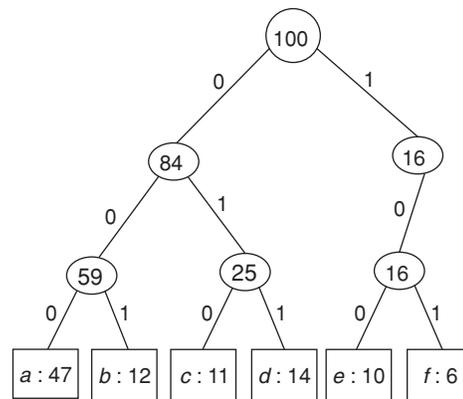
Step I:



Step II:



Step III:



We interpret the binary code word for a character as the path from the root to that character where ‘0’ means ‘go to the left child’, and 1 means ‘go to the right child’.

The above tree is not binary search tree, since the leaves need not appear in sorted order.

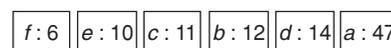
Constructing Huffman Code

This algorithm builds the tree T corresponding to the optimal code in a bottom-up manner. It begins with set of $|C|$ leaves and performs a sequence of $|C|-1$ ‘merging’ operations to create the final tree.

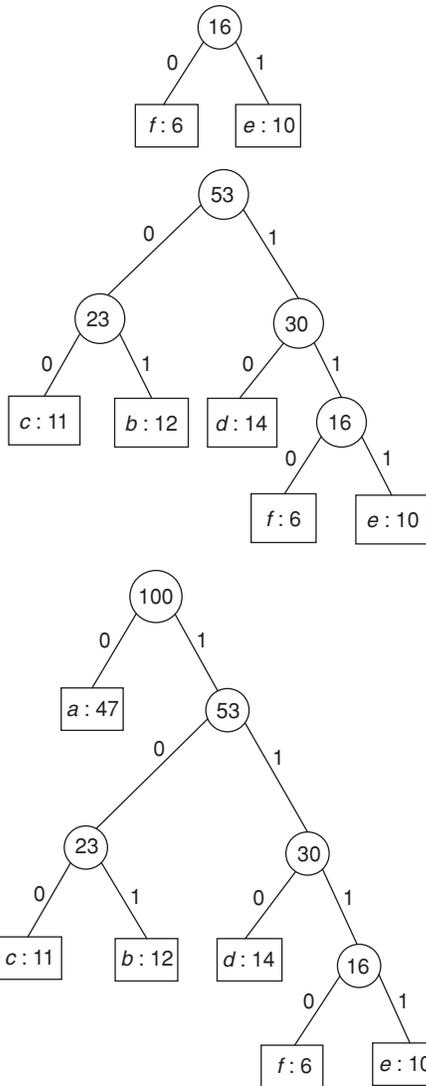
- A min-priority queue Q , keyed on f , is used to identify the 2 least – frequent objects to merge together. The result of the merger of 2 objects is a new object whose frequency is the sum of the frequencies of the 2 objects that were merged.
- In the given example, there are 6 alphabets the initial queue size is $n = 6$, and 5 merge steps are required to build the tree. The final tree represents the optimal prefix code. The code word for a letter is the sequence of edge labels on the path from the root to the letter.

$a = 47, b = 12, c = 11, d = 14, e = 10, f = 6$

Step I: Arrange the characters in non-decreasing order according to their frequencies



Let x and y be 2 characters in C having the lowest frequencies. Then there exists an optimal prefix code for C in which the code words for x and y have the same length and differ only in the last bit



Analysis: The analysis of the running time of Huffman’s algorithm assumes that Q is implemented as a binary min-heap for a set C of ‘ n ’ characters, the initialization of Q can be performed in $O(n)$ time using the BUILD – MIN HEAP procedure.

Each heap operation requires $O(\log n)$ time, and this will be performed exactly $(n - 1)$ times, it contributes to $O(n \log n)$ running time. Thus the total running time of HUFFMAN on a set of ‘ n ’ characters is $O(n \log n)$.

TASK-SCHEDULING PROBLEM

This is the problem of optimally scheduling unit – time tasks on a single processor, where each task has a deadline, along with a penalty that must be paid if the deadline is missed. The problem looks complicated, but it can be solved in simple manner using a greedy algorithm.

- A unit – time task is a job, such as a program to be run on a computer, that requires exactly one unit of time to complete.
- Given a finite set S of unit – time tasks, a schedule for S is a permutation of S specifying the order in which these tasks are to be performed.
- The first task in the schedule begins at time ‘0’ and finishes at time 1, the second task begins at time 1 and finishes at time 2, and so on
- The problem of scheduling unit – time tasks with deadlines and penalties for a single processor has the following inputs:

1. A set $S = \{a_1, a_2, \dots, a_n\}$ of n unit – time tasks:
2. A set of n integer deadlines d_1, d_2, \dots, d_n such that each d_i satisfies $1 \leq d_i \leq n$ and task a_i is supposed to finish by time d_i .
3. A set of n non-negative weights or penalties w_1, w_2, \dots, w_n , such that we incur a penalty of w_i if task a_i is not finished by time d_i and we incur no penalty if a task finishes by its deadline.

Example: Consider the following 7 tasks, $T_1, T_2, T_3, T_4, T_5, T_6, T_7$. Every task is associated with profit and deadline.

Tasks	T_1	T_2	T_3	T_4	T_5	T_6	T_7
Deadline	4	2	4	3	1	4	5
Profit	75	65	55	45	40	35	30

45	65	55	75	30			
T_4	T_2	T_3	T_1	T_7			
0	1	2	3	4	5	6	7

T_1 has highest profit, so it will be executed first and the deadline of T_1 , is ‘4’ so T_1 has to be executed within 4 slots of time, same procedure is applied to other tasks also.

The tasks which are not executed by CPU are T_5 and T_6 .

Profit: sum up the profits made by executing the tasks.
Profit = 45 + 65 + 55 + 75 + 30 = 270

Analysis: We can use a greedy algorithm to find a maximum weight independent set of tasks. We can then create an optimal schedule having the tasks in A as its early tasks.

This method is an efficient algorithm for scheduling unit – time tasks with deadlines and penalties for a single processor. The running time is $O(n^2)$ using GREEDY METHOD, since each of the $O(n)$ independent checks made by that algorithm takes time $O(n)$.

SORTING AND ORDER STATISTICS

Minimum and Maximum

This algorithm determines, how many comparisons are necessary to find minimum or maximum of a set of ‘ n ’ elements. Usually we can obtain maximum or minimum, by performing

$(n - 1)$ comparisons; examine each element of the set in turn and keep track of the smallest element seen so far.

Consider the following procedure.

Assume that the set of elements reside in an array A where length $[A] = n$

```

MINIMUM (A)
Min ← A[1]
For i ← 2 to length [A]
Do if min > A [i]
Then min ← A [i]
Return min.
    
```

Simultaneous minimum and maximum

In some applications, we must find both the minimum and the maximum of a set of ' n ' elements.

We can find the minimum and maximum independently using $(n - 1)$ comparisons for each, for a total of $(2n - 2)$ comparisons.

- In fact, atmost $3\lfloor n/2 \rfloor$ comparisons are sufficient to find both the minimum and the maximum.
- The strategy is to maintain the minimum and maximum elements seen so far.
- Rather than processing each element of the input by comparing it against the current minimum and maximum at a cost of 2 comparisons per element, we process elements in pairs.
- Compare pairs of elements from the input with each other, and then we compare smaller to the current minimum and the larger to the current maximum, at a cost of 3 comparisons for every 2 elements.
- Setting up initial values for the current minimum and maximum depends on whether ' n ' is odd or even. If ' n ' is odd, we set both the minimum and maximum to the value of the first element, and then we process the rest of the elements in pairs.
- If ' n ' is even, we perform 1 comparison on the first 2 elements to determine the initial values of the minimum and maximum and then process the rest of the elements in pairs.

Analysis: If ' n ' is odd the total number of comparisons would be $3\lfloor n/2 \rfloor$.

If ' n ' is even, we need 1 initial comparison followed by $\frac{3(n-2)}{2}$ comparisons, for a total of $\frac{3n}{2} - 2$.

∴ The total number of comparisons is atmost $3\lfloor n/2 \rfloor$

GRAPH ALGORITHMS

Single Source Shortest Path

In a shortest-path problem, we are given a weighted directed graph $G = (V, E)$ with weight function $W : E \rightarrow \mathbb{R}$ mapping edges to real-valued weights. The weight of path $P = \langle V_o, V_1 \dots V_k \rangle$ is the sum of the weights of its constituent edges. Shortest-path weight from U to V is defined by

$$\delta(U, V) = \begin{cases} \min\{W(P) : u \rightarrow v\} & \text{if there is a path from } \\ \infty & \text{'U' to 'V' otherwise} \end{cases}$$

Edge weights can be interpreted as metrics other than distances. They are often used to represent time, cost, penalties, loss, or any other quantity.

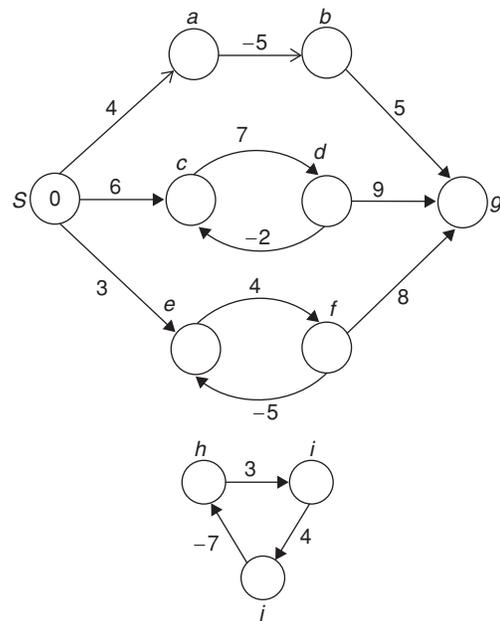
- The breadth first search algorithm is a shortest-path algorithm that works on un weighted graphs, that is, graphs in which each edge can be considered to have unit weight.

Negative-weight edges

Some of the instances of the single-source-shortest-paths problem, there may be edges whose weights are negative.

- If the graph $G = (V, E)$ contains no negative weight cycles reachable from source S , then for all $v \in V$, the shortest – path weight $d(S, V)$ remains well defined, even if it has a negative value.
- If there is a negative-weight cycle reachable from S , shortest-path weights are not well defined.
- No path from ' S ' to a vertex on the cycle can be a shortest path - a lesser weight path can always be found that follows the proposed 'shortest' path and then traverses the negative-weight cycle.
- If there is a negative-weight cycle on some path form ' S ' to ' V ', we define $\delta(S, V) = -\infty$.

Example: Consider the following graph, calculate the shortest distance to all vertices from sources ' S '.



Solution:

- Shortest path from S to a is $\delta(S, a) = W(S, a) = 4$ (because there is only one path from ' S ' to ' a ')
- Similarly, there is only one path from ' s ' to ' b '
 $\delta(S, a) = W(S, a) + W(a, b) = 4 + (-5) = -1$
- Shortest-path from ' s ' to ' c '

There are infinitely many paths from 'S' to 'c'

1. $\langle S, c \rangle$
2. $\langle S, c, d, c \rangle$
3. $\langle S, c, d, c, d, c \rangle$ and so on
 $\delta \langle S, c \rangle = 6$
 $\delta(S, d, d, c) = 6 + 7 - 2 = 11$
 $\delta(S, c, d, c, d, c) = 6 + 7 - 2 + 7 - 2 = 16$
 $\delta(S, c, d, c, d, c, d, c)$
 $= 6 + 7 - 2 + 7 - 2 + 7 - 2 = 21$

The cycle $\langle c, d, c \rangle$ has weight $= 7 + (-2) = 5 > 0$

The shortest path from 'S' to 'c' is $\langle s, c \rangle$ with weight $\delta(S, c) = 6$ similarly, the shortest-path from 'S' to 'd' is $\langle s, c, d \rangle$, with weight $\delta(S, d) = w(S, c) + W(c, d) = 13$

May there are infinitely paths from 'S' to 'e'

1. $\langle s, e \rangle$
2. $\langle s, e, f, e \rangle$
3. $\langle s, e, f, e, f, e \rangle$ and so on

Since the cycle $\langle e, f, e \rangle$ has weight $4 + (-5) = -1 < 0$. However, there is no shortest path from 'S' to 'e' by traversing the negative-weight cycle $\langle e, f, e \rangle$ arbitrarily many times, we can find paths from 's' to 'e' with arbitrarily large negative weights,

So $\delta(S, e) = -\infty$

Similarly, $\delta(S, f) = -\infty$

- The shortest path from 'S' to 'g':
 'g' can be reachable from 'f', we can also find paths with arbitrarily large negative weights from 's' to 'g' and $\delta(s, g) = -\infty$
- Vertices 'h', 'i' and 'j' also form a negative - weight cycle. They are not reachable from 'S' so, $\delta(S, h) = \delta(S, i) = \delta(S, j) = \infty$

Dijkstra's Algorithm

Dijkstra's algorithm solves the single-source shortest-path problem on a weighted, directed graph $G = (V, E)$, for the case in which all edge weights are non-negative.

- The running time of Dijkstra's algorithm is lower than that of the Bellman-Ford algorithm.
- Dijkstra's algorithm maintains a set 's' of vertices whose final shortest-path weights from the source 'S' have already been determined.
- The algorithm repeatedly selects the vertex $u \in (V - S)$ with the minimum shortest-path estimate, adds 'u' to 'S'

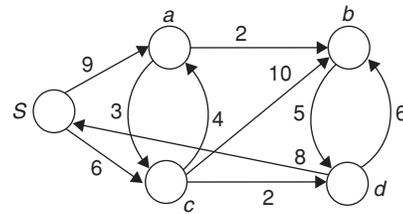
```

DIJKSTRA (G, W, S)
INITIALIZE - SINGLE - SOURCE (G, S)
S ← ∅
S ← V[G]
While Q ≠ ∅
do u ← EXTRACT - MIN(Q)
S ← S ∪ {u}
For each vertex v ∈ Adj[u]
do RELAX (u, v, w)
    
```

The algorithm maintains the invariant that $Q = V - S$ at the start of each iteration of the while loop. Initially the min - priority queue Q contains all the vertices in V. ($\because S = \emptyset$). Each time through the while loop, a vertex 'u' is extracted from $Q = V - S$ and added to set S.

- Each vertex is extracted from Q and added to S exactly once, so the contents of while loop will be executed exactly |V| times.
- Dijkstra's algorithm always chooses the 'closest' or 'lightest' vertex in $(V - S)$ to add to set S, we say that it uses a greedy strategy.

Example: Consider the following graph, what is the shortest path?



Solution:

S	V - S
S	a b c d
S c	a b d
S c d	a b
S c d a	b
S c d a b	∅

Distance from S to all vertices of $(V - S)$

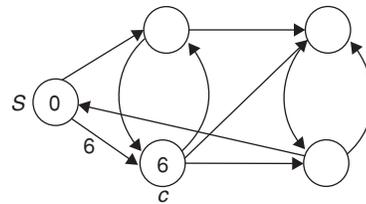
$$d[a] = 9$$

$$d[b] = \infty$$

$$d[c] = 6$$

$$d[d] = \infty$$

9, ∞, 6, ∞ values are given to MIN-PRIORITY Queue 'Q', '6' is returned.

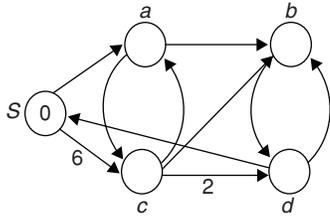


Distance from [Sc] to all vertices of $(V - S)$

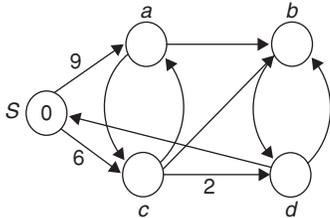
$$d[b] = (S - c - b) = 6 + 10 = 16$$

$$d[a] = \min\{(S - a) = 9, (s - c - a) = 10\} = 9$$

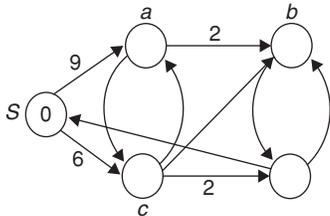
$$d[d] = \min\{\infty, (S - c - d) = 6 + 2 = 8\} = 8$$



Distance from $[s\ c\ d]$ to $[ab]$
 $d[a] = \min\{9, (S - c - d - S - a) = 25\} = 9$
 $d[b] = \min\{16, (S - c - d - b) = 14\} = 14$



$d[a] = \min\{14, (s - a - b) = 9 + 2 = 11\} = 11$



Analysis: It maintains the min-priority queue ‘ Q ’ by calling three priority-queue operations: INSERT, EXTRACT-MIN, and DECREASE-KEY. We maintain the min-priority queue by taking the vertices being numbered 1 to $|v|$. We store $d[v]$ in the v th entry of an array. Each INSERT and DECREASE-KEY operation takes $O(1)$ time, and each EXTRACT-MIN operation takes $O(v)$ time (\therefore we have to search through the entire array) for a total time of $O(v^2 + E) = O(v^2)$.

- If we implement the min - priority queue with a binary min-heap. Each EXTRACT-MIN operation takes time $O(\log V)$, there are $|V|$ such operations.
- The time to build the binary min-heap is $O(v)$. Each DECREASE-KEY operation takes time $O(\log V)$, and there are still atmost $|E|$ such operations. The total running time is $O((V + E) \log V)$, which is $O(E \log V)$ if all vertices are reachable form the source.
- We can achieve a running time of $O(V \log V + E)$ by implementing the min-priority queue with a Fibonacci heap.

Bellman–Ford Algorithm

Bellman–Ford algorithm solves the single-source shortest-path problems in the case in which edge weights may be negative.

- When negative edge lengths are permitted, we require that the graph have no cycles of negative length. This is

necessary to ensure that shortest-paths consist of a finite number of edges.

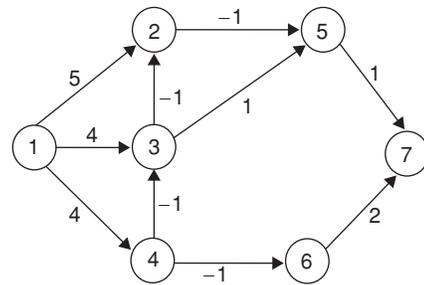
- When there are no cycles of negative length, there is a shortest-path between any two vertices of an n -vertex graph that has atmost $(n - 1)$ edges on it.
- A path that has more than $(n - 1)$ edges must repeat atleast one vertex and hence must contain a cycle
- Let $\text{dist}^x[u]$ be the length of a shortest-path from the source vertex ‘ v ’ to vertex ‘ u ’ under the constraint that the shortest-path contains atmost ‘ x ’ edges. Then $\text{dist}^x[u] = \text{cost}[v, u]$ $1 \leq u \leq n$ when there are no cycles of negative length we can limit our search for shortest-paths to paths with at most $(n - 1)$ edges. Hence, $\text{dist}^{n-1}[u]$ is the length of an unrestricted shortest-path from ‘ v ’ to ‘ u ’.

The Recurrence Relation for dist is:

$$\text{Dist}^k[u] = \min\{\text{dist}^{k-1}[u], \min\{\text{dist}^{k-1}[i] + \text{cost}[i, u]\}$$

This recurrence can be used to compute dist^k from dist^{k-1} , for $k = 2, 3, \dots, n - 1$.

Example: Consider the given directed graph



Find the shortest path from vertex ‘1’ to all other vertices using Bellman–Ford algorithm?

Solution: Source vertex is ‘1’ the distance from ‘1’ to ‘1’ in all ‘6’ iterations will be zero. Since the graph has ‘7’ vertices, the shortest-path can have atmost ‘6’ edges. The following figure illustrates the implementation of Bellman–Ford algorithm:

	1	2	3	4	5	6	7
1	0	5	4	4	∞	∞	∞
2	0	3	3	4	4	3	∞
3	0	2	3	4	2	3	5
4	0	2	3	4	1	3	3
5	0	2	3	4	1	3	2
6	0	2	3	4	1	3	2
7	0	2	3	4	1	3	2

Analysis

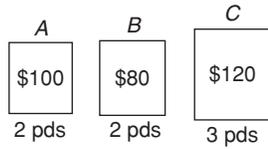
- Each iteration takes $O(n^2)$ time if adjacency matrices are used and $O(e)$ time if adjacency lists are used. Here ‘ e ’ is the number of edges in the graph.
- The time complexity is $O(n^3)$ when adjacency matrices are used and $O(N * E)$ when adjacency lists are used.

EXERCISES

Practice Problems I

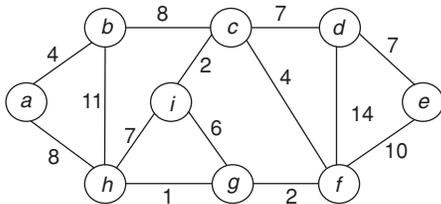
Directions for questions 1 to 14: Select the correct alternative from the given choices.

1. A thief enters a store and sees the following:



His knapsack can hold 4 pounds, what should he steal to maximize profit? (Use 0–1 Knapsack).

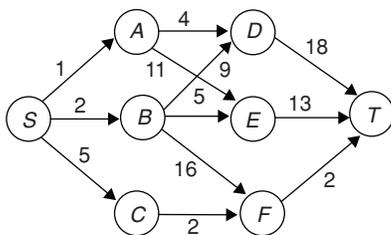
- (A) A and B (B) A and C
 (C) B and C (D) A, B and C
2. By using fractional Knapsack, calculate the maximum profit, for the data given in the above question?
 (A) 180 (B) 170
 (C) 160 (D) 150
3. Consider the below figure:



What is the weight of the minimum spanning tree using Kruskals algorithm?

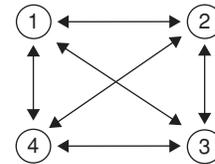
- (A) 34 (B) 35
 (C) 36 (D) 38
4. Construct a minimum spanning tree for the figure given in the above question, using prim's algorithm. What are the first three nodes, added to the solution set respectively (consider 'a' as starting node).
 (A) b, c, i (B) h, b, c
 (C) c, i, b (D) h, c, b

5. Consider the below graph, calculate the shortest distance from 'S' to 'T'?



- (A) 23 (B) 9
 (C) 20 (D) 22
6. Solve the travelling salesman problem, with the given distances in the form of matrix of graph, which of the following gives optimal solution?

C	1	2	3	4
1	0	15	20	25
2	10	0	14	15
3	11	18	0	17
4	13	13	14	0



- (A) 1-2-4-3-1 (B) 2-3-4-1-2
 (C) 1-4-2-3-1 (D) 2-4-3-1-2

7. Calculate the maximum profit using greedy strategy, knapsack capacity is 50. The data is given below:

$n = 3$
 $(w_1, w_2, w_3) = (10, 20, 30)$
 $(p_1, p_2, p_3) = (60, 100, 120)$ (dollars)? (0/1 knapsack)

(A) 180 (B) 220
 (C) 240 (D) 260

Common data for questions 8 and 9: Given that

	a	b	c	d	e	f
Frequency	45	13	12	16	9	5
Fixed length code word	000	001	010	011	100	101

8. Using Huffman code, find the path length of internal nodes.
 (A) 8 (B) 100
 (C) 100×8 (D) $100/8$
9. Using above answer, external path length will be
 (A) 18 (B) 108
 (C) 8 (D) None of these

Common data for questions 10 and 11:



Using 0–1 knapsack select a subset of the three items shown, whose weight must not exceed 50 kg. What is the value?

- (A) 2220 (B) 2100
 (C) 2600 (D) 2180

11. Which of the following gives maximum profit, using fractional knapsack?
 (A) $x_1 = 1, x_2 = 1, x_3 = 0$ (B) $x_1 = 1, x_2 = 1, x_3 = 2/3$
 (C) $x_1 = 1, x_2 = 0, x_3 = 1$ (D) $x_1 = 1, x_2 = 1, x_3 = 1/3$

12. Using dynamic programming find the longest common subsequence (LCS) in the given 2 sub sequences:

$x [1, \dots, m]$
 $y [1, \dots, n]$
 $x : A B C B D A B$
 $y : B D C A B A$

Find longest sequence sets common to both.

- (A) (BDAB, BCAB, BCBA)
 (B) (BADB, BCAB, BCBA)
 (C) (BDAB, BACB, BCBA)
 (D) (BDAB, BCAB, BBCA)
13. Let C_1, C_2, C_3, C_4 represent coins.
 $C_1 = 25$ paisa

$C_2 = 10$ paisa
 $C_3 = 5$ paisa
 $C_4 = 1$ paisa

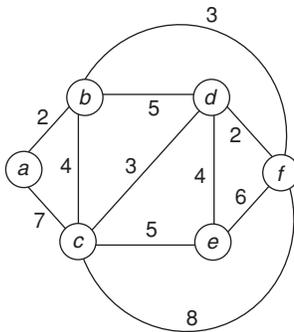
To represents 48 paisa, what is the minimum number of coins used, using greedy approach?

- (A) 6 (B) 7
 (C) 8 (D) 9
14. Worst-case analysis of hashing occurs when
- (A) All the keys are distributed
 (B) Every key hash to the same slot
 (C) Key values with even number, hashes to slots with even number
 (D) Key values with odd number hashes to slots with odd number.

Practice Problems 2

Directions for questions 1 to 15: Select the correct alternative from the given choices.

1. Consider the given graph:

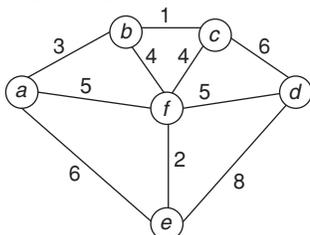


Which one of the following cannot be the sequence of edges added, in that order, to a minimum spanning tree using Kruskal's algorithm?

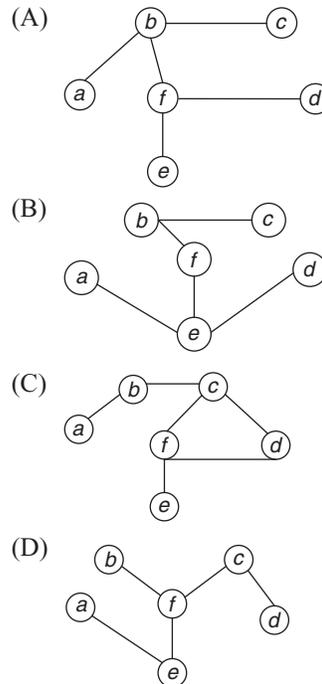
- (A) $(a-b), (d-f), (b-f), (d-c), (d-e)$
 (B) $(a-b), (d-f), (d-c), (b-f), (d-e)$
 (C) $(d-f), (a-b), (d-c), (b-f), (d-e)$
 (D) $(d-f), (a-b), (b-f), (d-e), (d-c)$
2. The worst case height analysis of B-tree is

- (A) $O(n)$
 (B) $O(n^2)$
 (C) $O(\log n)$
 (D) $O(n \log n)$

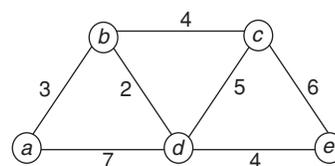
3. Consider the given graph:



Which of the following is the minimum spanning tree. (If we apply Kruskal algorithm).



4. Consider the following graph:



Find the shortest path using Dijkstra's algorithm.

- (A) $a-b-d-e$ (B) $a-b-c-d$
 (C) $a-c-d-e$ (D) $a-b-c-e$
5. Which statement is true about Kruskal's algorithm?
- (A) It is a greedy algorithm for the minimum spanning tree problem.
 (B) It constructs spanning tree by selecting edges in increasing order of their weights.
 (C) It does not accept creation of cycles in spanning tree.
 (D) All the above

6. Dijkstra's algorithm bears similarity to which of the following for computing minimum spanning trees?
 (A) Breadth first search (B) Prim's algorithm
 (C) Both (A) and (B) (D) None of these
7. Which of the following algorithm always yields a correct solution for a graph with non-negative weights to compute shortest paths?
 (A) Prim's algorithm (B) Kruskal's algorithm
 (C) Dijkstra's algorithm (D) Huffman tree
8. Let the load factor of the hash table is number of keys is n , cells of the hash table is m then
 (A) $\infty = n/m$ (B) $\infty = m/n$
 (C) $\infty = \frac{m+1}{n}$ (D) $\infty = \frac{n+1}{m}$
9. To implement Dijkstra's shortest path algorithm on unweighted graphs so that it runs in linear time, the data structure to be used is:
 (A) Queue
 (B) Stack
 (C) Heap
 (D) B-tree
10. The development of a dynamic-programming algorithm can be broken into a sequence of four steps, which are given below randomly.
 I. Construct an optimal solution from computed information.
 II. Compute the value of an optimal solution in a bottom-up fashion.
 III. Characterize the structure of an optimal solution.
 IV. Recursively defines the value of an optimal solution.
 The correct sequence of the above steps is

- (A) I, II, III, IV (B) IV, III, I, II
 (C) IV, II, I, III (D) III, IV, II I

11. Let V stands for vertex, E stands for edges.
 For both directed and undirected graphs, the adjacency list representation has the desirable property that the amount of memory required is
 (A) $\theta(V)$ (B) $\theta(E)$
 (C) $\theta(V + E)$ (D) $\theta(V - E)$
12. Which of the following is false?
 (A) Adjacency-matrix representation of a graph permits faster edge look up.
 (B) The adjacency matrix of a graph requires $\theta(v^2)$ memory, independent of the number of edges in the graph.
 (C) Adjacency-matrix representation can be used for weighted graphs.
 (D) All the above
13. Dynamic programming is a technique for solving problems with
 (A) Overlapped sub problems
 (B) Huge size sub problems
 (C) Small size sub problems
 (D) None of these
14. The way a card game player arranges his cards, as he picks them up one by one is an example of _____.
 (A) Bubble sort (B) Selection sort
 (C) Insertion sort (D) None of the above
15. You want to check whether a given set of items is sorted. Which method will be the most efficient if it is already in sorted order?
 (A) Heap sort (B) Bubble sort
 (C) Merge sort (D) Insertion sort

PREVIOUS YEARS' QUESTIONS

Data for question 1: We are given 9 tasks $T_1, T_2 \dots T_9$. The execution of each task requires one unit of time. We can execute one task at a time. Each task T_i has a profit P_i and a deadline D_i . Profit P_i is earned if the task is completed before the end of the D_i th unit of time.

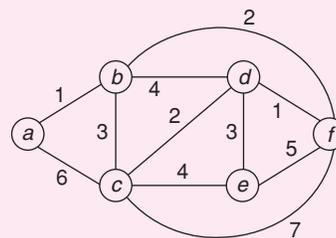
Task	T_1	T_2	T_3	T_4	T_5	T_6	T_7	T_8	T_9
Profit	15	20	30	18	18	10	23	16	25
Deadline	7	2	5	3	4	5	2	7	3

1. What is the maximum profit earned? [2005]
 (A) 147 (B) 165
 (C) 167 (D) 175
2. Consider a weighted complete graph G on the vertex set $\{v_1, v_2, \dots, v_n\}$ such that the weight of the edge (v_i, v_j) is $2|i - j|$. The weight of the minimum spanning tree is: [2006]
 (A) $n - 1$ (B) $2n - 2$
 (C) $\binom{n}{2}$ (D) n^2

3. To implement Dijkstra's shortest path algorithm on unweighted graphs so that it runs in linear time, the data structure to be used is: [2006]

- (A) Queue
 (B) Stack
 (C) Heap
 (D) B-Tree

4. Consider the following graph: [2006]



Which one of the following cannot be the sequence of edges added, in that order, to a minimum spanning tree using Kruskal's algorithm?

- (A) $(a - b), (d - f), (b - f), (d - c), (d - e)$
- (B) $(a - b), (d - f), (d - c), (b - f), (d - e)$
- (C) $(d - f), (a - b), (d - c), (b - f), (d - e)$
- (D) $(d - f), (a - b), (b - f), (d - e), (d - c)$

Common data for questions 5 and 6: A 3-ary max-heap is like a binary max-heap, but instead of 2 children, nodes have 3 children. A 3-ary heap can be represented by an array as follows: The root is stored in the first location, $a[0]$, nodes in the next level, from left to right, is stored from $a[1]$ to $a[3]$. The nodes from the second level of the tree from left to right are stored from $a[4]$ location onward. An item x can be inserted into a 3-ary heap containing n items by placing x in the location $a[n]$ and pushing it up the tree to satisfy the heap property.

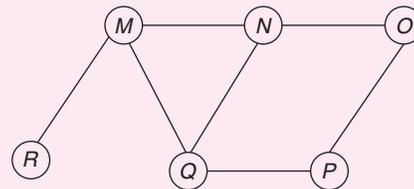
5. Which one of the following is a valid sequence of elements in an array representing 3-ary max-heap? [2006]
 - (A) 1, 3, 5, 6, 8, 9
 - (B) 9, 6, 3, 1, 8, 5
 - (C) 9, 3, 6, 8, 5, 1
 - (D) 9, 5, 6, 8, 3, 1
6. Suppose the elements 7, 2, 10 and 4 are inserted, in that order, into the valid 3-ary max-heap found in the above question, Q-76. Which one of the following is the sequence of items in the array representing the resultant heap? [2006]
 - (A) 10, 7, 9, 8, 3, 1, 5, 2, 6, 4
 - (B) 10, 9, 8, 7, 6, 5, 4, 3, 2, 1
 - (C) 10, 9, 4, 5, 7, 6, 8, 2, 1, 3
 - (D) 10, 8, 6, 9, 7, 2, 3, 4, 1, 5
7. In an unweighted, undirected connected graph, the shortest path from a node S to every other node is computed most efficiently, in terms of *time complexity*, by [2007]
 - (A) Dijkstra's algorithm starting from S .
 - (B) Warshall's algorithm
 - (C) Performing a DFS starting from S .
 - (D) Performing a BFS starting from S .
8. A complete n -ary tree is a tree in which each node has n children or no children. Let I be the number of internal nodes and L be the number of leaves in a complete n -ary tree. If $L = 41$, and $I = 10$, what is the value of n ? [2007]
 - (A) 3
 - (B) 4
 - (C) 5
 - (D) 6
9. Consider the following C program segment where CellNode represents a node in a binary tree: [2007]


```
struct CellNode {
    struct CellNode *leftChild;
    int element;
    struct CellNode *rightChild;
};
```

```
int GetValue (struct CellNode *ptr) {
    int value = 0;
    if (ptr != NULL) {
        if ((ptr->leftChild == NULL) &&
            (ptr->rightChild == NULL))
            value = 1;
        else
            value = value + GetValue(ptr->leftChild)
                + GetValue(ptr->rightChild);
    }
    return (value);
}
```

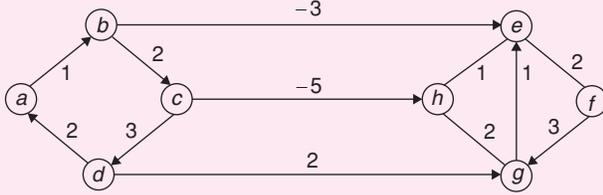
The value returned by GetValue when a pointer to the root of a binary tree is passed as its argument is:

- (A) The number of nodes in the tree
 - (B) The number of internal nodes in the tree
 - (C) The number of leaf nodes in the tree
 - (D) The height of the tree
10. Let w be the minimum weight among all edge weights in an undirected connected graph. Let e be a specific edge of weight w . Which of the following is FALSE? [2007]
 - (A) There is a minimum spanning tree containing e .
 - (B) If e is not in a minimum spanning tree T , then in the cycle formed by adding e to T , all edges have the same weight.
 - (C) Every minimum spanning tree has an edge of weight w .
 - (D) e is present in every minimum spanning tree.
 11. The Breadth first search algorithm has been implemented using the queue data structure. One possible order of visiting the nodes of the following graph is [2008]
 - (A) MNOPQR
 - (B) NQMPOR
 - (C) QMNPOR
 - (D) QMNPRO



- (A) MNOPQR
 - (B) NQMPOR
 - (C) QMNPOR
 - (D) QMNPRO
12. G is a graph on n vertices and $2n - 2$ edges. The edges of G can be partitioned into two edge-disjoint spanning trees. Which of the following is NOT true for G ? [2008]
 - (A) For every subset of k vertices, the induced subgraph has at most $2k - 2$ edges
 - (B) The minimum cut in G has at least two edges
 - (C) There are two edge-disjoint paths between every pair of vertices
 - (D) There are two vertex-disjoint paths between every pair of vertices

13.

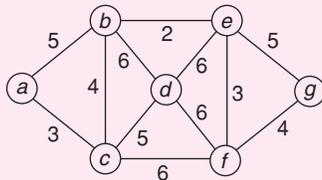


Dijkstra's single source shortest path algorithm when run from vertex a in the above graph, computes the correct shortest path distance to [2008]

- (A) Only vertex a
 - (B) Only vertices a, e, f, g, h
 - (C) Only vertices a, b, c, d
 - (D) all the vertices
14. You are given the post-order traversal, P , of a binary search tree on the n elements $1, 2, \dots, n$. You have to determine the unique binary search tree that has P as its post-order traversal. What is the time complexity of the most efficient algorithm for doing this? [2008]
- (A) $\Theta(\log n)$
 - (B) $\Theta(n)$
 - (C) $\Theta(n \log n)$
 - (D) None of the above, as the tree cannot be uniquely determined
15. Which of the following statement(s) is/are correct regarding Bellman–Ford shortest path algorithm? [2009]

- P. Always finds a negative weighted cycle, if one exists.
 - Q. Finds whether any negative weighted cycle is reachable from the source.
- (A) P only
 - (B) Q only
 - (C) Both P and Q
 - (D) Neither P nor Q

16. Consider the following graph: [2009]



Which one of the following is NOT the sequence of edges added to the minimum spanning tree using Kruskal's algorithm? [2009]

- (A) (b, e) (e, f) (a, c) (b, c) (f, g) (c, d)
- (B) (b, e) (e, f) (a, c) (f, g) (b, c) (c, d)
- (C) (b, e) (a, c) (e, f) (b, c) (f, g) (c, d)
- (D) (b, e) (e, f) (b, c) (a, c) (f, g) (c, d)

Common data for questions 17 and 18: Consider a binary max-heap implemented using an array.

17. Which one of the following array represents a binary max-heap? [2009]

- (A) $\{25, 12, 16, 13, 10, 8, 14\}$
- (B) $\{25, 14, 13, 16, 10, 8, 12\}$
- (C) $\{25, 14, 16, 13, 10, 8, 12\}$
- (D) $\{25, 14, 12, 13, 10, 8, 16\}$

18. What is the content of the array after two delete operations on the correct answer to the previous question? [2009]

- (A) $\{14, 13, 12, 10, 8\}$
- (B) $\{14, 12, 13, 8, 10\}$
- (C) $\{14, 13, 8, 12, 10\}$
- (D) $\{14, 13, 12, 8, 10\}$

Common data for questions 19 and 20: Consider a complete undirected graph with vertex set $\{0, 1, 2, 3, 4\}$. Entry W_{ij} in the matrix W below is the weight of the edge $\{i, j\}$.

$$W = \begin{pmatrix} 0 & 1 & 8 & 1 & 4 \\ 1 & 0 & 12 & 4 & 9 \\ 8 & 12 & 0 & 7 & 3 \\ 1 & 4 & 7 & 0 & 2 \\ 4 & 9 & 3 & 2 & 0 \end{pmatrix}$$

19. What is the minimum possible weight of a spanning tree T in this graph such that vertex 0 is a leaf node in the tree T ? [2010]

- (A) 7
- (B) 8
- (C) 9
- (D) 10

20. What is the minimum possible weight of a path P from vertex 1 to vertex 2 in this graph such that P contains at most 3 edges? [2010]

- (A) 7
- (B) 8
- (C) 9
- (D) 10

Common data for questions 21 and 22: A hash table of length 10 uses open addressing with hash function $h(k) = k \bmod 10$, and linear probing. After inserting 6 values into an empty hash table, the table is as shown below:

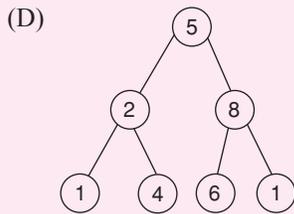
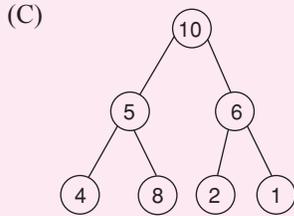
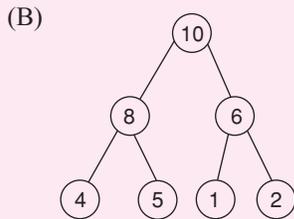
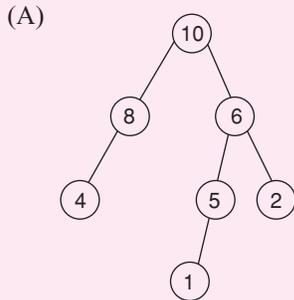
0	
1	
2	42
3	23
4	34
5	52
6	46
7	33
8	
9	

21. Which one of the following choices gives a possible order in which the key values could have been inserted in the table?

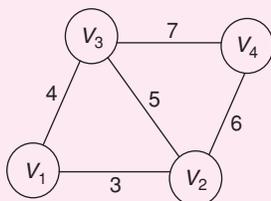
- (A) 46, 42, 34, 52, 23, 33
- (B) 34, 42, 23, 52, 33, 46
- (C) 46, 34, 42, 23, 52, 33
- (D) 42, 46, 33, 23, 34, 52

[2010]

22. How many different insertion sequences of the key values using the same hash function and linear probing will result in the hash table shown above? [2010]
 (A) 10 (B) 20
 (C) 30 (D) 40
23. A max-heap is a heap where the value of each parent is greater than or equal to the value of its children. Which of the following is a max-heap? [2011]

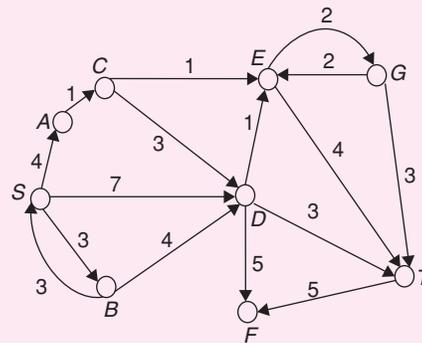


Common data for questions 24 and 25: An undirected graph $G(V, E)$ contains $n(n > 2)$ nodes named V_1, V_2, \dots, V_n . Two nodes V_i, V_j are connected if and only if $0 < |i - j| \leq 2$. Each edge (V_i, V_j) is assigned a weight $i + j$. A sample graph with $n = 4$ is shown below.



24. What will be the cost of the minimum spanning tree (MST) of such a graph with n nodes? [2011]
 (A) $\frac{1}{12}(11n^2 - 5n)$ (B) $n^2 - n + 1$
 (C) $6n - 11$ (D) $2n + 1$
25. The length of the path from V_5 to V_6 in the MST of previous question with $n = 10$ is [2011]
 (A) 11 (B) 25
 (C) 31 (D) 41

26. Consider the directed graph shown in the figure below. There are multiple shortest paths between vertices S and T . Which one will be reported by Dijkstra's shortest path algorithm? Assume that, in any iteration, the shortest path to a vertex v is updated only when a strictly shorter path to v is discovered. [2012]



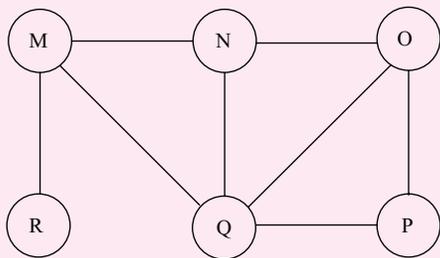
- (A) SDT (B) SBDT
 (C) SACDT (D) SACET

27. Let G be a weighted graph with edge weights greater than one and G^1 be the graph constructed by squaring the weights of edges in G . Let T and T^1 be the minimum spanning trees of G and G^1 , respectively, with total weights t and t^1 . Which of the following statements is **TRUE**? [2012]
 (A) $T^1 = T$ with total weight $t^1 = t^2$
 (B) $T^1 = T$ with total weight $t^1 < t^2$
 (C) $T^1 \neq T$ but total weight $t^1 = t^2$
 (D) None of the above

28. What is the time complexity of Bellman–Ford single-source shortest path algorithm on a complete graph of n vertices? [2013]
 (A) $\Theta(n^2)$ (B) $\Theta(n^2 \log n)$
 (C) $\Theta(n^3)$ (D) $\Theta(n^3 \log n)$

29. Consider the following operation along with Enqueue and Dequeue operations on queues, where k is a global parameter.

```
MultiDequeue(Q) {
    m = k
```

- (A) MNOPQR
- (B) NQMPOR
- (C) QMNROP
- (D) POQNMR

50. A message is made up entirely of characters from the set $X = \{P, Q, R, S, T\}$. The table of probabilities for each of the characters is shown below:

Character	Probability
P	0.22
Q	0.34
R	0.17
S	0.19
T	0.08
Total	1.00

If a message of 100 characters over X is encoded using Huffman coding, then the expected length of the encoded message in bits is _____. [2017]

51. Let G be a simple undirected graph. Let T_D be a depth first search tree of G . Let T_B be a breadth first search tree of G .

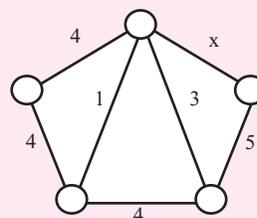
Consider the following statements.

- (I) No edge of G is a cross edge with respect to T_D . (A cross edge in G is between two nodes neither of which is an ancestor of the other in T_D .)
- (II) For every edge (u, v) of G , if u is at depth i and v is at depth j in T_B , then $|i - j| = 1$.

Which of the statements above must necessarily be true? [2018]

- (A) I only
- (B) II only
- (C) Both I and II
- (D) Neither I nor II

52. Consider the following undirected graph G :



Choose a value for x that will maximize the number of minimum weight spanning trees (MWSTs) of G . The number of MWSTs of G for this value of x is _____. [2018]

ANSWER KEYS

EXERCISES

Practice Problems 1

1. A 2. A 3. B 4. A 5. B 6. A 7. B 8. A 9. A 10. C
 11. B 12. A 13. A 14. B

Practice Problems 2

1. D 2. C 3. A 4. A 5. D 6. C 7. C 8. A 9. A 10. D
 11. C 12. C 13. A 14. C 15. D

Previous Years' Questions

1. A 2. B 3. C 4. D 5. D 6. A 7. D 8. C 9. C 10. B
 11. C 12. D 13. D 14. B 15. B 16. D 17. C 18. D 19. D 20. B
 21. C 22. C 23. B 24. B 25. C 26. D 27. 28. C 29. A 30. D
 31. C 32. C 33. 12 to 12 34. A 35. B 36. 6 to 6 37. 19 38. 110 39. 69
 40. A 41. C 42. 5 43. B 44. A 45. 7 46. B 47. 31 48. A 49. D
 50. 225 51. A 52. 4

Chapter 5

Dynamic Programming

LEARNING OBJECTIVES

- ☞ Dynamic programming
- ☞ Multi-stage graph
- ☞ All pairs shortest path problem
- ☞ Hashing methods
- ☞ Mid-square method
- ☞ Folding method
- ☞ Resolving collisions
- ☞ Matrix chain multiplication
- ☞ Longest common subsequence
- ☞ Optimal substructure of LCS
- ☞ NP-hard and NP-complete
- ☞ P-problem
- ☞ NP-problem
- ☞ P, NP, and Co-NP
- ☞ Cooks theorem
- ☞ Non-deterministic search

DYNAMIC PROGRAMMING

Dynamic programming is a method for solving complex problems by breaking them down into simpler sub problems. It is applicable to problems exhibiting the properties of overlapping sub problems which are only slightly smaller, when applicable; the method takes far less time than naive method.

- The key idea behind dynamic programming is to solve a given problem, we need to solve different parts of the problem (sub problems) then combine the solutions of the sub problems to reach an overall solution. Often, many of these sub problems are the same.
- The dynamic programming approach seeks to solve each sub problem only once, thus reducing the number of computations. This is especially useful when the number of repeating sub problems is exponentially large.
- There are two key attributes that a problem must have in order for dynamic programming to be applicable ‘optimal sub structure’ and ‘overlapping sub-problems’. However, when the overlapping problems are much smaller than the original problem, the strategy is called ‘divide-and-conquer’ rather than ‘dynamic programming’. This is why merge sort-quick sort are not classified as dynamic programming problems.

Dynamic programming is applied for:

- Multi stage graph
- All pairs shortestest path

Principle of Optimality

It states that whatever the initial state is, remaining decisions must be optimal with regard to the state following from the first decision.

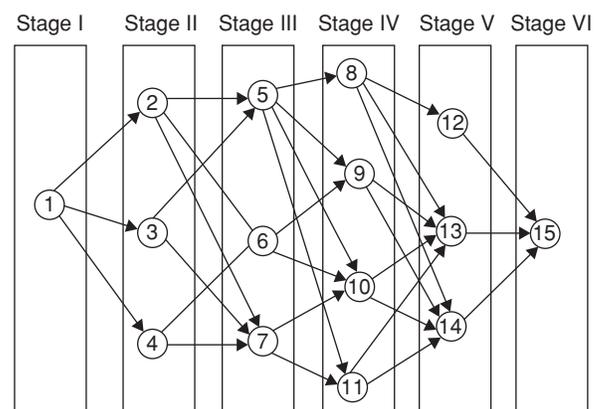
To solve a problem using dynamic programming strategy, it must observe the principle of optimality.

MULTI-STAGE GRAPH

A multi-stage graph is a graph

- $G = (V, E)$ with V partitioned into $K > = 2$ disjoint subsets such that if (a, b) is in E , then a is in V_p and b is in V_{i+1} for some sub sets in the partition;
- $|V_1| = |V_K| = 1$ the vertex S in V_1 is called the source; the vertex t is called the sink.
- G is usually assumed to be a weighted graph.
- The cost of a path from node V to node W is sum of the costs of edges in the path.
- The ‘multi-stage graph problem’ is to find the minimum cost path from S to t .

Example:



Costs of edges

- 1 – 2 → 10
- 1 – 3 → 20

1 – 4	→	30
2 – 5	→	10
2 – 6	→	20
2 – 7	→	30
3 – 5	→	40
3 – 7	→	50
4 – 6	→	40
4 – 7	→	30
5 – 8	→	10
5 – 9	→	20
5 – 10	→	30
5 – 11	→	40
6 – 9	→	20
6 – 10	→	30
7 – 10	→	30
7 – 11	→	20
8 – 12	→	10
8 – 13	→	20
8 – 14	→	30
9 – 13	→	20
9 – 14	→	10
10 – 13	→	10
10 – 14	→	20
11 – 13	→	10
11 – 14	→	30
12 – 15	→	20
13 – 15	→	10
14 – 15	→	30

Solution Using Backward Cost

Format: COST (Stage, node) = minimum cost of travelling to the node in stage from the source node (node 1)

Step I:

$$\text{Cost (I, 1)} = 0$$

Step II:

$$\text{Cost (II, 2)} = \text{cost (I, 1)} + \text{cost (1, 2)} = 0 + 10 = 10$$

$$\text{Cost (II, 3)} = \text{cost (I, 1)} + \text{cost (1, 3)} = 0 + 20 = 20$$

$$\text{Cost (II, 4)} = \text{cost (I, 1)} + \text{cost (1, 4)} = 0 + 30 = 30$$

Step III:

$$\begin{aligned} \text{Cost (III, 5)} &= \min \{ \text{cost (II, 2)} + \text{cost (2, 5)}, \\ &\quad \text{cost (II, 3)} + \text{cost (3, 5)}, \\ &\quad \text{cost (II, 4)} + \text{cost (4, 5)} \} \\ &= \min \{ 10 + 10, 20 + 40, 30 + \infty \} \\ &= 20 \rightarrow \text{Via path } 1 - 2 - 5 \end{aligned}$$

$$\begin{aligned} \text{Cost (III, 6)} &= \min \{ \text{cost (II, 2)} + \text{cost (2, 6)}, \\ &\quad \text{cost (II, 3)} + \text{cost (3, 6)}, \\ &\quad \text{cost (II, 4)} + \text{cost (4, 6)} \} \\ &= \min \{ 10 + 20, 20 + \infty, 30 + 40 \} \\ &= 30 \rightarrow \text{via the path } 1 - 2 - 6 \end{aligned}$$

$$\begin{aligned} \text{Cost (III, 7)} &= \min \{ \text{cost (II, 2)} + \text{cost (2, 7)}, \\ &\quad \text{Cost (II, 3)} + \text{cost (3, 7)}, \\ &\quad \text{Cost (II, 4)} + \text{cost (4, 7)} \} \\ &= \min \{ 10 + 30, 20 + 50, 30 + 30 \} \\ &= 40 \rightarrow \text{Via the path } 1 - 2 - 7 \end{aligned}$$

Step IV:

$$\begin{aligned} \text{Cost (IV, 8)} &= \min \{ \text{cost (III, 5)} + \text{cost (5, 8)}, \\ &\quad \text{Cost (III, 6)} + \text{cost (6, 8)}, \\ &\quad \text{Cost (III, 7)} + \text{cost (7, 8)} \} \\ &= \min \{ 20 + 10, 30 + \infty, 40 + \infty \} \\ &= 30 \rightarrow \text{Via path } 1 - 2 - 5 - 8 \end{aligned}$$

$$\begin{aligned} \text{Cost (IV, 9)} &= \min \{ \text{cost (III, 5)} + \text{cost (5, 9)}, \\ &\quad \text{Cost (III, 6)} + \text{cost (6, 9)}, \\ &\quad \text{Cost (III, 7)} + \text{cost (7, 9)} \} \\ &= \min \{ 20 + 20, 30 + 20, 40 + \infty \} \\ &= 40 \rightarrow \text{Via the path } 1 - 2 - 5 - 9 \end{aligned}$$

$$\begin{aligned} \text{Cost (IV, 10)} &= \min \{ \text{cost (III, 5)} + \text{cost (5, 10)}, \\ &\quad \text{Cost (III, 6)} + \text{cost (6, 10)}, \\ &\quad \text{Cost (III, 7)} + \text{cost (7, 10)} \} \\ &= \min \{ 20 + 30, 30 + 30, 40 + 30 \} \\ &= 50 \rightarrow \text{Via the path } 1 - 2 - 5 - 10 \end{aligned}$$

$$\begin{aligned} \text{Cost (IV, 11)} &= \min \{ \text{cost (III, 5)} + \text{cost (5, 11)}, \\ &\quad \text{Cost (III, 6)} + \text{cost (6, 11)}, \\ &\quad \text{Cost (III, 7)} + \text{cost (7, 11)} \} \\ &= \min \{ 20 + 40, 30 + \infty, 40 + 20 \} \\ &= 60 \rightarrow \text{Via the path } 1 - 2 - 5 - 11 \\ &\quad \text{or Via the path } 1 - 2 - 7 - 11 \end{aligned}$$

Step V:

$$\begin{aligned} \text{Cost (V, 12)} &= \min \{ \text{cost (IV, 8)} + \text{cost (8, 12)}, \\ &\quad \text{Cost (IV, 9)} + \text{cost (9, 12)}, \\ &\quad \text{Cost (IV, 10)} + \text{cost (10, 12)}, \\ &\quad \text{Cost (IV, 11)} + \text{cost (11, 12)} \} \\ &= \min \{ 30 + 10, 40 + \infty, 50 + \infty, 60 + \infty \} \\ &= 40 \rightarrow \text{Via the path } 1 - 2 - 5 - 8 - 12 \end{aligned}$$

$$\begin{aligned} \text{Cost (V, 13)} &= \min \{ \text{cost (IV, 8)} + \text{cost (8, 13)}, \\ &\quad \text{Cost (IV, 9)} + \text{cost (9, 13)}, \\ &\quad \text{Cost (IV, 10)} + \text{cost (10, 13)}, \\ &\quad \text{Cost (IV, 11)} + \text{cost (11, 13)} \} \\ &= \min \{ 30 + 20, 40 + 20, 50 + 10, 60 + 10 \} \\ &= 50 \rightarrow \text{Via the path } 1 - 2 - 5 - 8 - 13 \end{aligned}$$

$$\begin{aligned} \text{Cost (V, 14)} &= \min \{ \text{cost (IV, 8)} + \text{cost (8, 14)}, \\ &\quad \text{Cost (IV, 9)} + \text{cost (9, 14)}, \\ &\quad \text{Cost (IV, 10)} + \text{cost (10, 14)}, \\ &\quad \text{Cost (IV, 11)} + \text{cost (11, 14)} \} \\ &= \min \{ 30 + 30, 40 + 10, 50 + 20, 60 + 30 \} \\ &= 50 \rightarrow \text{Via the path } 1 - 2 - 5 - 9 - 14 \end{aligned}$$

Step VI:

$$\begin{aligned} \text{Cost (VI, 15)} &= \min \{ \text{cost (V, 12)} + \text{cost (12, 15)}, \\ &\quad \text{Cost (V, 13)} + \text{cost (13, 15)}, \\ &\quad \text{Cost (V, 14)} + \text{cost (14, 15)} \} \\ &= \min \{ 40 + 20, 50 + 10, 50 + 30 \} \\ &= 60 \rightarrow \text{Via the path } 1 - 2 - 5 - 8 - 13 - 15 \\ &\quad \text{(or) } 1 - 2 - 5 - 8 - 12 - 15 \end{aligned}$$

ALL PAIRS SHORTEST PATH PROBLEM (FLOYD-WARSHALL ALGORITHM)

A weighted graph is a collection of points (vertices) connected by lines (edges), where each edge has a weight (some real number) associated with it.

Example: A graph in the real world is a road map. Each location is a vertex and each road connecting locations is an edge. We can think of the distance travelled on a road from one location to another as the weight of that edge.

- The Floyd-Warshall algorithm determines the shortest path between all pairs of vertices in a graph.
- The vertices in a graph be numbered from 1 to n . Consider the subset $\{1, 2, \dots, K\}$ of these n vertices.
- Finding the shortest path from vertex i to vertex j that uses vertex in the set $\{1, 2, \dots, K\}$ only. There are two situations.
 1. K is an intermediate vertex on the shortest path.
 2. K is not an intermediate vertex on the shortest path.

In the first situation, we can break down our shortest path into two paths: i to K and then K to j . Note that all the vertices from i to K are from the set $\{1, 2, \dots, K-1\}$ and that all the intermediate vertices from K to j are from the set $\{1, 2, \dots, K-1\}$. Also in the second situation, we simply have that all intermediate vertices are from the set $\{1, 2, \dots, K-1\}$. Now define the function D for a weighted graph with the vertices $\{1, 2, \dots, n\}$ as follows.

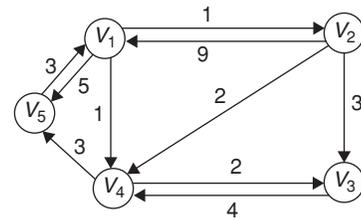
$D(i, j, K)$ = the shortest distance from vertex i to vertex j using the intermediate vertices. In the set $\{1, 2, \dots, K\}$

Using the above idea, we can recursively define the function D .

$$\begin{aligned} D(i, j, K) &= W(i, j) \text{ if } K = 0 \\ &= \min (D(i, j, K-1), D(i, K, K-1) + D(K, j, K-1)) \text{ if } K > 0 \end{aligned}$$

- The first line says that if we do not allow intermediate vertices, then the shortest path between two vertices is the weight of the edge that connects them. If no such weight exists, we usually define this shortest path to be of length infinity.
- The second line pertains to allowing intermediate vertices. It says that the minimum path from i to j through vertices $\{1, 2, \dots, K\}$ is either the minimum path from i to j through vertices $\{1, 2, \dots, K-1\}$ OR the sum of the minimum path from vertex i to K through $\{1, 2, \dots, K-1\}$ plus the minimum path from vertex K to j through $\{1, 2, \dots, K-1\}$. Since this is the case, we compute both and choose the smaller of these.

Example:



The weight matrix will be

	1	2	3	4	5
1	0	1	∞	1	5
2	9	0	3	2	∞
3	∞	∞	0	4	∞
4	∞	∞	2	0	3
5	3	∞	∞	∞	0

Let $D^{(K)} [i, j]$ = weight of a shortest path from v_i to v_j using only vertices from $\{v_1, v_2, \dots, v_k\}$ as intermediate vertices in the path.

- $D^{(0)} = W$
- $D^{(n)} = D$ which is the goal matrix.

How to compute $D^{(K)}$ from $D^{(K-1)}$?

Case I: A shortest path from v_i to v_j restricted to using only vertices from $\{v_1, v_2, \dots, v_K\}$ as intermediate vertices does not use V_K . Then $D^{(K)} [i, j] = D^{(K-1)} [i, j]$

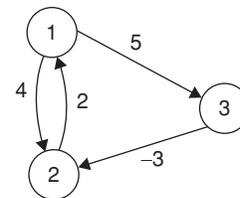
Case II: A shortest path from v_i to v_j restricted to using only vertices from $\{v_1, v_2, \dots, v_K\}$ as intermediate vertices does use V_K . Then $D^{(K)} [i, j] = D^{(K-1)} [i, K] + D^{(K-1)} [K, j]$

Since $D^{(K)} [i, j] = D^{(K-1)} [i, j]$

or $D^{(K)} [i, j] = D^{(K-1)} [i, K] + D^{(K-1)} [K, j]$

We conclude: $D^{(K)} [i, j] = \min \{ D^{(K-1)} [i, j], D^{(K-1)} [i, K] + D^{(K-1)} [K, j] \}$

Example: 1



$W = D^0 =$

	1	2	3
1	0	4	5
2	2	0	∞
3	∞	-3	0

$P =$

	1	2	3
1	0	0	0
2	0	0	0
3	0	0	0

$K = 1$, vertex 1 can be intermediate node

$$D^1 [2, 3] = \min (D^0[2, 3], D^0[2, 1] + D^0[1, 3])$$

$$= \min (\infty, 7)$$

$$= 7$$

$$D^1 [3, 2] = \min (D^0[3, 2], D^0[3, 1] + D^0[1, 2])$$

$$= \min (-3, \infty)$$

$$= -3$$

	1	2	3
1	0	4	5
2	2	0	7
3	∞	-3	0

	1	2	3
1	0	0	0
2	0	0	1
3	0	0	0

$K = 2$, vertices 1, 2 can be intermediate nodes,

$$D^2 [1, 3] = \min (D [1, 3], D [1, 2] + D [2, 3])$$

$$= \min (5, 4 + 7) = 5$$

$$D^2 [3, 1] = \min (D [3, 1], D [3, 2] + D [2, 1])$$

$$= \min (\infty, -3 + 2)$$

$$= -1$$

	1	2	3
1	0	4	5
2	2	0	7
3	-1	-3	0

	1	2	3
1	0	0	0
2	0	0	1
3	2	0	0

$K = 3$ vertices 1, 2, 3 can be intermediate

$$D^3 [1, 2] = \min (D^2[1, 2], D^2[1, 3] + D^2[3, 2])$$

$$= \min (4, 5 + (-3))$$

$$= 2$$

$$D^3 [2, 1] = \min (D^2[2, 1], D^2[2, 3] + D^2[3, 1])$$

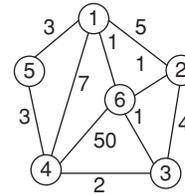
$$= \min (2, 7 + (-1))$$

$$= 2$$

	1	2	3
1	0	2	5
2	2	0	7
3	-1	-3	0

	1	2	3
1	0	3	0
2	0	0	1
3	2	0	0

Example 2:



The final distance matrix and P

	1	2	3	4	5	6
1	0	2(6)	2(6)	4(6)	3	1
2	2(6)	0	2(6)	4(6)	5(6)	1
3	2(6)	2(6)	0	2	5(4)	1
4	4(6)	4(6)	2	0	3	3(3)
5	3	5(4)	5(4)	3	0	4(1)
6	1	1	1	3(3)	4(1)	0

The values in parenthesis are the non-zero P values.

Table 1 Divide and conquer vs dynamic programming.

1. This design strategy divides the problem into sub problems, conquer the each sub problem recursively, finally combine all the sub problem solutions, for the original problem.	1. This design strategy chooses an optimal solution for the problem, by recursively defining the value of optimal solution, these values are computed in bottom up fashion or top down fashion.
2. each sub problem is solved recursively, and consumes more time at each sub problem	2. Each sub problem is solved only once and is stored in table
3. Sub problems are independent of each other e.g., Binary search	3. The sub problems are dependent e.g., Traveling sales person problem

Dynamic Programming vs Greedy Method

The main difference between greedy method (GM) and dynamic programming (DP) methodology is, DP considers all possible solutions for the given problem and picks the optimal one. Where as greedy, considers only one set of solutions to the problem.

The other difference between GM and DP is that, GM considers the choice, which is best at that step, which is done at each level of the sub problem. That is, it won't reconsider its choice. The choices reflect only present, won't consider the future choices, where as DP tries out all the best alternatives and finds the optimal solution. It implements principle of optimality. At each stage of the problem, it decides based on the previous decision made in the previous stage.

HASHING METHODS

Uniform Hash Function

If the keys, K , are integers randomly distributed in $[0, r]$ then hash function $H(K)$ is given as

$$H(K) = \left\lfloor \frac{mk}{r} \right\rfloor$$

$H(K)$ is a uniform hash function

Uniform hashing function should ensure

$$\sum_{K|h(K)=0} P(K) = \sum_{K|h(K)=1} P(K) = \dots = \sum_{K|h(K)=m} P(K) = \frac{1}{m}$$

$P(K)$ = probability that a key K , occurs that is the number of keys that map to each slot is equal.

Division method

Hashing an integer x is to divide x by M and then to use the remainder modulo M . This is called the division method of hashing. In this case the hash function is

$$h(x) = x \bmod M$$

Generally this approach is quite good for just about any value of M . However, in certain situations some extra care is needed in the selection of a suitable value for M . For example, it is often convenient to make M an even number. But this means that $h(x)$ is even if x is even, and $h(x)$ is odd if x is odd. If all possible keys are equiprobable, then this is not a problem. However, if say even keys are more likely than odd keys, the function $h(x) = x \bmod M$ will not spread the hashed values of those keys evenly.

- Let M be a power of two, i.e., $M = 2^k$ for some integer $k > 1$. In this case, the hash function $h(x) = x \bmod 2^k$ simply extracts the bottom k -bits of the binary representation of x . While this hash function is quite easy to compute, it is not a desirable function because it does not depend on all the bits in the binary representation of x .
- For these reasons M is often chosen to be a prime number. Suppose there is bias in the way the keys are created that makes it more likely for a key to be a multiple of some small constant, say two or three. Then making M a prime increases the likelihood that those keys are spread out evenly. Also if M is a prime number, the division of x by that prime number depends on all the bits of x , not just the bottom k -bits, for some small constant k .

Example: Hash table size = 10

Key value = 112

Hash function = $h(k) = k \bmod M$
 $= 112 \bmod 10 = 2$

Disadvantage: A potential disadvantage of the division method is due to the property that consecutive keys map to consecutive hash values.

$$h(i) = i$$

$$h(i + 1) = i + 1 \pmod{M}$$

$$h(i + 2) = i + 2 \pmod{M}$$

⋮
⋮
⋮

While this ensures that consecutive keys do not collide, it does not mean that consecutive array locations will be occupied. We will see that in certain implementations this can lead to degradation in performance.

Multiplication method

A variation on the middle-square method that alleviates its deficiencies is called, multiplication hashing method. Instead of multiplying the key x by itself, we multiply the key by a carefully chosen constant 'a' and then extract the middle k bits from the result. In this case, the hashing function is

$$h(x) = \left\lfloor \frac{M}{W} (ax \bmod W) \right\rfloor$$

if we want to avoid the problems that the middle-square method encounters with keys having a large number of leading (or) trailing zero's then we should choose an 'a' that has neither leading nor trailing zero's.

Furthermore, if we, choose an 'a' that is relatively prime to W , then there exists another number 'a' such that $aa' = 1 \pmod{W}$. Such a number has the nice property that if we take a key x , and multiply it by 'a' to get ax , we can recover the original key by multiplying the product again by a' , since $a \times a' = aa'x = 1x$.

The multiplication method for creating a hash function operates in two steps:

Step 1: Multiply the key K by a constant A in the range $0 < A < 1$ and extract the fractional part of KA .

Step 2: Multiply this value by M and take the floor of the result.

In short the hash function is

$$h(k) = \lfloor M \cdot (KA \bmod 1) \rfloor$$

Where $(KA \bmod 1)$ denotes the fractional part of KA , that is $KA - \lfloor KA \rfloor$

Example:

$$\begin{aligned} \text{Let } m = 10000, K = 123456 \text{ and } A = \frac{\sqrt{5}-1}{2} \\ = 0.618033 \end{aligned}$$

$$\begin{aligned} \text{Then } h(k) &= \lfloor 10000 \cdot (123456 \cdot 0.61803 \bmod 1) \rfloor \\ &= \lfloor 10000 \cdot (76300.00412 \bmod 1) \rfloor \\ &= \lfloor 10000 \cdot 0.00412 \rfloor = 41 \end{aligned}$$

Practical issues

- Easy to implement
 - On most machines multiplication is faster than division.
 - We can substitute one multiplication by shift operation.
 - We don't need to do floating-point operations.
- If successive keys have a large interval, $A = 0.6125423371$ can be recommended.

Mid-square method

A good hash function to use with integer key values is the mid-square method. The mid-square method squares the key value, and then takes out the middle 'r' bits of the result, giving a value in the range 0 to $2^r - 1$. This works well because most (or) all bits of the key value contribute to the result.

Example:

Consider records whose keys are 4-digit numbers in base 10. The goal is to hash these key values to a table of size 100(i.e., a range of 0 to 99).

This range is equivalent to two digits in base 10.

That is $r = 2$. If the input is the number 4567, squaring yields an 8-digit number, 20857489. The middle two digits of this result are 57. All digits of the original key value (equivalently, all bits when the number is viewed in binary) contribute to the middle two digits of the squared value. Thus, the result is not dominated by the distribution of the bottom or the top digit of the original key value. Of course, if the key values all tend to be small numbers, then their squares will only affect the low order digits of the hash value.

Example: To map the key 3121 into a hash table of size 1000, we square it $(3121)^2 = 9740641$ and extract 406 as the hash value.

Folding method

The folding method breaks up a key into precise segments that are added to form a hash value, and still another technique is to apply a multiplicative hash function to each segment individually before folding.

Algorithm $H(x) = (a + b + c) \bmod m$. Where a , b , and c represent the preconditioned key broken down into three parts, m is the table size, and mod stands for modulo. In other words: The sum of three parts of the pre conditioned key is divided by the table size. The remainder is the hash key.

Example:

Fold the key 123456789 into a hash table of ten spaces (0 through 9)

We are given $x = 123456789$ and the table size (i.e., $m = 10$)

Since we can break x into three parts any way, we will break it up evenly.

Thus $a = 123$, $b = 456$ and $c = 789$

$$H(x) = (a + b + c) \bmod M$$

$$H(123456789) = (123 + 456 + 789) \bmod 10 \\ = 1368 \bmod 10 = 8$$

123456789 are inserted into the table at address 8. The folding method is distribution independent.

Resolving collisions In collision resolution strategy algorithms and data structures are used to handle two hash keys that hash to the same hash keys. There are a number of collision resolution techniques, but the most popular are open addressing and chaining.

- Chaining: An array of linked list, Separate chaining
- Open Addressing: Array based implementation:
 - Linear probing (Linear Search)
 - Quadratic probing (non-linear search)
 - Double hashing (use two hash functions)

Separate chaining Every linked list has each element that collides to the similar slot. Insertion need to locate the accurate slot and appending to any end of the list in that slot wherever, deletion needs searching the list and removal.

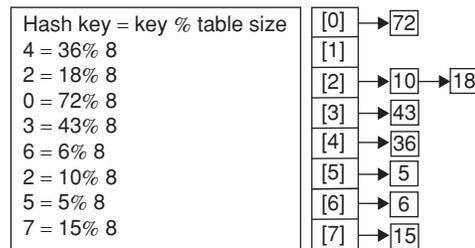


Figure 1 Separate chaining

Open addressing Open addressing hash tables are used to stock up the records straight inside the array. This approach is also known as closed hashing. This procedure is based on probing. Well known probe sequence include:

- Linear probing: In which the interval between probes is fixed often at 1.
- Quadratic probing: In which the interval between probes increases proportional to the hash value (the interval thus increasing linearly and the indices are described by a quadratic function).
- Double hashing: In which the interval between probes is computed by another hash function.

(i) Linear probing: Linear probing method is used for resolving hash collisions of values of hash functions by sequentially searching the hash table for a free location. The item will be stored in the next available slot in the table in linear probing. Also an assumption is made that the table is not already full.

This is implemented via a linear search for an empty slot, from the point of collision.

If the physical end of table is reached during the linear search, the search will again get start around to the beginning of the table and continue from there. The table is considered as full, if an empty slot is not found before reaching the point of collision.

[0]	72	Add the keys 10, 5 and 15 to the previous example. Hash key = key % table size 2 = 10% 8 5 = 5% 8 7 = 15% 8	[0]	72
[1]			[1]	5
[2]	18		[2]	18
[3]	43		[3]	43
[4]	36		[4]	36
[5]			[5]	10
[6]	6		[6]	6
[7]	7		[7]	7

Figure 2 Linear probing

Limitation: A problem with linear probe method is primary clustering. In primary clustering blocks of data may possibly be able to form collision. Several attempts may be required by any key that hashes into the cluster to resolve the collision.

(ii) **Quadratic probing:** To resolve the primary clustering problem, quadratic probing can be used. With quadratic probing, rather than always moving one spot, move i^2 spots from the point of collision where i is the number of attempts needed to resolve the collision.

[0]	49	89% 10 = 9 18% 10 = 8 49% 10 = 9 → 1 attempt needed → $1^2 = 1$ spot 58% 10 = 8 → 2 attempts needed → $2^2 = 4$ spot 69% 10 = 9 → 2 attempts needed → $2^2 = 4$ spot	[0]	69
[1]			[1]	
[2]	58		[2]	
[3]	69		[3]	58
[4]			[4]	
[5]			[5]	
[6]			[6]	49
[7]			[7]	
[8]	18	[8]	18	
[9]	89	[9]	89	

Limitation: Maximum half of the table can be used as substitute locations to resolve collisions. Once the table gets more than half full, its really hard to locate an unfilled spot. This new difficulty is recognized as secondary clustering because elements that hash to the same hash key will always probe the identical substitute cells.

(iii) **Double hashing:** Double hashing uses the idea of applying a second hash function to the key when a collision occurs, the result of the second hash function will be the number of positions from the point of collision to insert. There are some requirements for the second function:

1. It must never evaluate to zero
2. Must make sure that all cells can be probed.

A popular second hash function is:
Hash(key) = R - (Key mod R) where R is a prime number smaller than the size of the table.

Table size = 10 elements Hash1(key) = key % 10 Hash 2(key) = 7 - (key % 7) Insert keys: 89, 18, 49, 58 and 69 Hash key (89) = 89% 10 = 9 Hash key (18) = 18% 10 = 8 Hash key (49) = 49% 10 = 9 (collision) = (7 - (49% 7)) = (7 - (0)) = 7 positions from [9]		[0]	
		[1]	
		[2]	
		[3]	
		[4]	
		[5]	
		[6]	49
		[7]	
		[8]	18
		[9]	89

Figure 3 Double hashing

Insert keys = 58, 69 Hash key (58) = 58% 10 = 8 a collision! = (7 - (58% 7)) = (7 - 2) = 5 positions from [8] Hash key (69) = 69% 10 = 9 a collision! = (7 - (69 % 7)) = (7 - 6) = 1 position from [9]		[0]	69
		[1]	
		[2]	
		[3]	58
		[4]	
		[5]	
		[6]	49
		[7]	
		[8]	18
		[9]	89

Figure 4 Double hashing

MATRIX-CHAIN MULTIPLICATION

We are given a sequence of n matrices $m_1, m_2 \dots m_n$ to be multiplied. If the chain matrices is $\langle m_1, m_2, m_3, m_4 \rangle$, the product m_1, m_2, m_3, m_4 can be fully parenthesized in 5 distinct ways:

1. $(m_1 (m_2 (m_3 m_4)))$
2. $(m_1 ((m_2 m_3) m_4))$
3. $((m_1 m_2) (m_3 m_4))$
4. $((m_1 (m_2 m_3)) m_4)$
5. $((((m_1 m_2) m_3) m_4))$

The way we parenthesize a chain of matrices can have a dramatic impact on the cost of evaluating the product. We can multiply 2 matrices A and B only if they are compatible i.e., the number of columns of A must equal the number of rows of B . If A is a $(p \times q)$ matrix and B is a $(q \times r)$ matrix, the resulting matrix C is a $(p \times r)$ matrix. The time to compute C is the number of scalar multiplications, which is (pqr) .

Example: Consider the problem of a chain $\langle m_1, m_2, m_3 \rangle$ of three matrices. Suppose that the dimensions of the matrices are (6×8) , (8×14) , (14×20) respectively. Which parenthesization will give least number of multiplications?

Solution:

(i) $((m_1 m_2) m_3)$
 $[m_1]_{6 \times 8} \times [m_2]_{8 \times 14} = [m_1 m_2]_{6 \times 14}$
 Number of multiplications performed
 $= 6 \times 8 \times 14 = 672$
 $[m_1 m_2]_{6 \times 14} \times [m_3]_{14 \times 20} = ((m_1 m_2) m_3)_{6 \times 20}$
 Number of multiplications performed
 $= 6 \times 14 \times 20 = 1680$
 Total number of multiplications
 $= 672 + 1680 = 2352$

(ii) $(m_1 (m_2 m_3))$
 $[m_2]_{8 \times 14} \times [m_3]_{14 \times 20} = [m_2 m_3]_{8 \times 20}$
 Number of multiplications performed
 $= 8 \times 14 \times 20 = 2240$
 $[m_1]_{6 \times 8} \times [m_2 m_3]_{8 \times 20} = (m_1 (m_2 m_3))_{6 \times 20}$
 Number of multiplications performed
 $= 6 \times 8 \times 20 = 960$
 Total number of multiplications = $960 + 2240 = 3200$
 $\therefore ((m_1 m_2) m_3)$ gives least number of multiplications.

We need to define the cost of an optimal solution recursively in terms of the optimal solutions to sub problems. For Matrix-chain multiplication problem, we pick as our sub problem the problems of determining the minimum cost of a parenthesization of $A_i A_{i+1} \dots A_j$ for $1 \leq i \leq j \leq n$ let $m[i, j]$ be the minimum number of scalar multiplications needed to compute the matrix $A_i \dots A_j$; for the full problem, the cost of a cheapest way to compute $A_1 \dots A_n$ would be $m[1, n]$. We can define $m[i, j]$ recursively as follows:

$$m[i, j] = m[i, k] + m[k + 1, j] + P_{i-1} P_k P_j$$

If $i = j$, the problem is trivial. The chain consists of just one matrix $A_i \dots A_i = A_i$, so that no scalar multiplications are necessary to compute the product.

Minimum cost of parenthesizing the product $A_i A_{i+1} \dots A_j$ becomes

$$m[i, j] = \begin{cases} 0 & \text{if } i = j \\ \min\{m[i, k] + m[k + 1, j] + P_{i-1} P_k P_j\} & \text{if } i < j, i \leq k < j \end{cases}$$

The $m[i, j]$ values give the costs of optimal solutions to sub problems.

At this point, to write a recursive algorithm based on recurrence to compute the minimum cost $m[1, n]$ for multiplying $A_1 A_2 \dots A_n$. However, this algorithm takes exponential time, which is not better than the brute force method of checking each way of parenthesizing the product. The important observation we can make at this point is that we have relatively few sub problems, one problem for each choice of i and j satisfying $1 \leq i \leq j \leq n$ (or)

$\left(\frac{n}{2}\right) + n = \theta(n^2)$ in all. The property of overlapping sub problems is the second hallmark of the applicability of dynamic programming.

The first hall mark being optimal substructure.

Algorithm

1. $n \leftarrow \text{length}[p] - 1$
2. for $i \leftarrow 1$ to n
3. do $m[i, i] \leftarrow 0$
4. for $i \leftarrow 2$ to n
5. do for $i \leftarrow 1$ to $n - i + 1$
6. do $j \leftarrow i + i - 1$
7. $m[i, j] \leftarrow \infty$
8. for $k \leftarrow i$ to $j - 1$
9. do $q \leftarrow m[i, k] + m[k + 1, j] + P_{i-1} P_k P_j$
10. if $q < m[i, j]$
11. then $m[i, j] \leftarrow q$
12. $S[i, j] \leftarrow k$
13. return m and S

It first computes $m[i, j] \leftarrow 0$ for $i = 1, 2 \dots n$ (the minimum costs for chains of length 1). To compute $m[i, i + 1]$ for $i = 1, 2, \dots, n - 1$ (the minimum costs for chains of length $\lambda = 2$ and so on). At each step, the $m[i, j]$ cost computed depends only on table entries $m[i, k]$ and $m[k + 1, j]$ already computed. An entry $m[i, j]$ is computed using the products $P_{i-1} P_k P_j$ for $k = i, i + 1, \dots, j - 1$. A simple inspection of the nested loop structure of the above algorithm yields a running time of $O(n^3)$ for the algorithm.

LONGEST COMMON SUBSEQUENCE

A sub sequence of a given sequence is just the given sequence with 0 or more elements left out. Formally, given a sequence $x = \langle x_1, x_2 \dots x_m \rangle$, another sequence $z = \langle z_1, z_2 \dots z_k \rangle$ is a subsequence of x if there exists a strictly increasing sequence $\langle i_1, i_2 \dots i_k \rangle$ of indices of x such that for all $j = 1, 2 \dots k$, we have $x_{i_j} = z_j$

Example: $z = \langle B, C, D, B \rangle$ is a subsequence of $x = \langle A, B, C, B, D, A, B \rangle$ with corresponding index sequence $\langle 2, 3, 5, 7 \rangle$

Example: Given 2 sequences x and y , we say that a sequence z is a common sub sequence of x and y if z is a sub sequence of both x and y .

$$\begin{aligned} \text{If } x &= \langle A, B, C, B, D, A, B \rangle \\ y &= \langle B, D, C, A, B, A \rangle \end{aligned}$$

The sequence $\langle B, C, A \rangle$ is a common subsequence of both x and y .

The sequence $\langle B, C, A \rangle$ is not a longest common subsequence (LCS) of x and y since it has length '3' and the sequence $\langle B, C, B, A \rangle$, which is also common to both x and y , has length 4. The sequence $\langle B, C, B, A \rangle$ is an LCS of x and y , as is the sequence $\langle B, D, A, B \rangle$, since there is no common subsequence of length 5 or greater.

- In the longest-common-sub sequence problem, we are given 2 sequences $x = \langle x_1, x_2, x_3 \dots x_m \rangle$ and $y = \langle y_1, y_2 \dots y_n \rangle$ and wish to find a maximum length common subsequence of x and y .
- LCS problem can be solved efficiently using dynamic programming.
- A brute force approach to solve the LCS problem is to enumerate all subsequences of x and check each subsequence to see if it is also a subsequence of y , keeping track of the longest subsequence found. Each subsequence of x corresponds to a subset of the indices $\{1, 2 \dots m\}$ of x . There are 2^m subsequences of x , so this approach requires exponential time, making it impractical for long sequences.
- The classes of sub problems correspond to pairs of ‘pre fixes’ of 2 input sequences:

Given a sequence $x = \langle x_1, x_2 \dots x_m \rangle$, we define the i th prefix of x , for $i = 0, 1, \dots, m$, as

$$x_i = \langle x_1, x_2 \dots x_i \rangle$$

Example: If $x = \langle A, B, C, B, D, A, D \rangle$, then $x_4 = \langle A, B, C, B \rangle$ and x_0 is the empty sequence. LCS problem has an optimal sub-structure property.

Optimal Substructure of LCS

Let $x = \langle x_1, x_2 \dots x_m \rangle$ and $y = \langle y_1, y_2 \dots y_n \rangle$ be sequences and let $z = \langle z_1, z_2 \dots z_k \rangle$ be any LCS of x and y then

1. If $x_m = y_n$, then $z_k = x_m = y_n$ and z_{k-1} is an LCS of x_{m-1} and y_{n-1} .
2. If $x_m \neq y_n$, then $z_k \neq x_m$ implies that z is an LCS of x_{m-1} and y .
3. If $x_m \neq y_n$, then $z_k \neq y_n$ implies that z is an LCS of x and y_{n-1} .

NP-HARD AND NP-COMPLETE

A mathematical problem for which, even in theory, no shortcut or smart algorithm is possible that would lead to a simple or rapid solution. Instead the only way to find an optimal solution is a computationally intensive, exhaustive analysis in which all possible outcomes are tested. Examples of NP-hard problems include the travelling salesman problem.

P-problem

A problem is assigned to the P (polynomial time) class if there exists at least one algorithm to solve that problem, such that number of steps of the algorithm is bounded by a polynomial in n , where n is the length of the input.

NP-problem

A problem is assigned to the NP (non-deterministic polynomial time) class if it is solvable in polynomial time by a non-deterministic turing machine.

A P -problem (whose solution time is bounded by a polynomial) is always also NP . If a problem is known to be

NP , and a solution to the problem is somehow known, then demonstrating the correctness of the solution can always be reduced to a single P (polynomial time) verification. If P and NP are not equivalent then the solution of NP -problems requires (in the worst case) an exhaustive search.

A problem is said to be NP -hard, if an algorithm for solving it can be translated into one for solving any other NP -problem. It is much easier to show that a problem is NP than to show that it is NP -hard. A problem which is both NP and NP -hard is called an NP -complete problem.

P versus NP-problems

The P versus NP problem is the determination of whether all NP -problems are actually P -problems, if P and NP are not equivalent then the solution of NP -problem requires an exhaustive search, while if they are, then asymptotically faster algorithms may exist.

NP-complete problem

A problem which is both NP (verifiable in non-deterministic polynomial time) and NP -hard (any NP -problem can be translated into this problem). Examples of NP -hard problems include the Hamiltonian cycle and travelling sales man problems.

Example:

Circuit satisfiability is a good example of problem that we don't know how to solve in polynomial time. In this problem, the input is a Boolean circuit. A collection of and, or and not gates connected by wires. The input to the circuit is a set of m Boolean (true/false) values $x_1 \dots x_m$. The output is a single Boolean value. The circuit satisfiability problem asks, given a circuit, whether there is an input that makes the circuit output TRUE, or conversely, whether the circuit always outputs FALSE. Nobody knows how to solve this problem faster than just trying all 2^m possible inputs to the circuit but this requires exponential time.

P, NP, and Co-NP

- P is a set of yes/no problems that can be solved in polynomial time. Intuitively P is the set of problems that can be solved quickly.
- NP is the set of yes/no problems with the following property: If the answer is yes, then there is a proof of this fact that can be checked in polynomial time. Intuitively NP is the set of problems where we can verify a YES answer quickly if we have the solution in front of us.

Example: The circuit satisfiability problem is in NP .

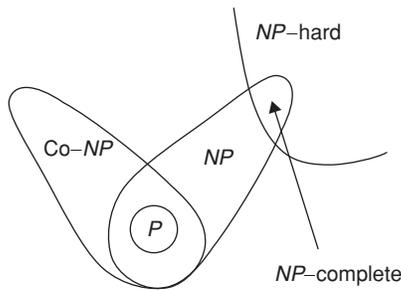
If the answer is yes, then any set of m input values that produces TRUE output is a proof of this fact, we can check the proof by evaluating the circuit in polynomial time.

- Co- NP is the exact opposite of NP . If the answer to a problem in co- NP is no, then there is a proof of this fact that can be checked in polynomial time.

- π is *NP*-hard \Rightarrow if π can be solved in polynomial time, then $P = NP$.

This is like saying that if we could solve one particular *NP*-hard problem quickly, then we could solve any problem whose solution is easy to understand, using the solution to that one special problem as a subroutine. *NP*-hard problems are atleast as hard as any problem in *NP*.

- Saying that a problem is *NP*-hard is like saying ‘If I own a dog, then it can speak fluent English’. You probably don’t know whether or not I own a dog, but you’re probably pretty sure that I don’t own a talking dog. Nobody has a mathematical proof that dogs can’t speak English. The fact that no one has ever heard a dog speak English is evidence as per the hundreds of examinations of dogs that lacked the proper mouth shape and brain power, but mere evidence is not a proof nevertheless, no sane person would believe me if I said I owned a dog that spoke fluent English. So the statement ‘If I own a dog then it can speak fluent English’ has a natural corollary: No one in their right mind should believe that I own a dog ! Likewise if a problem is *NP*-hard no one in their right mind should believe it can be solved in polynomial time.



Cooks Theorem

Cook’s theorem states that CNFSAT is *NP*-Complete

It means, if the problem is in *NP*, then the deterministic Turing machine can reduce the problem in polynomial time.

The inference that can be taken from these theorems is, if deterministic polynomial time algorithm exists for solving satisfiability, then to all problems present in *NP* can be solved in polynomial time.

Non-deterministic Search

Non-deterministic algorithms are faster, compared to deterministic ones. The computations are fast as it always chooses right step

The following functions are used to specify these algorithms

1. Choice (A), which chooses a random element from set A
2. Failure (A), specifies failure
3. Success (), Specifies success

The non-deterministic search is done as follows.

Let us consider an array $S[1 \dots n]$, $n \geq 1$ we need to get the indice of ‘ i ’ such that $S[i] = t$ (or) $i = 0$. The algorithm is given below.

Steps:

1. $i = \text{Choice}(1, n)$;
2. if $S[i] = t$, then
 - (i) Print (i);
 - (ii) Success ();
3. Print (0) failure
4. Stop.

If the search is successful it returns the indice of array ‘ S ’, otherwise it returns ‘0’, the time complexity is $\Omega(n)$.

EXERCISES

Practice Problems I

Directions for questions 1 to 15: Select the correct alternative from the given choices.

1. Hash the keys 12, 44, 13, 88, 23, 94, 11, 39, 20 using the hash function with chaining $(2k + 5) \bmod 11$, which of the following slots are empty?
 (A) 0, 1, 2, 3, 4 (B) 0, 2, 3, 4, 8, 10
 (C) 0, 1, 2, 4, 8, 10 (D) 0, 1, 2, 4, 8
2. Using linear probing on the list given in the above question with the same hash function, which slots are not occupied?
 (A) 3, 4 (B) 4, 5
 (C) 3, 6 (D) 4, 6
3. In hashing, key value 123456 is hashed to which address using multiplication method ($m = 10^4$)?
 (A) 40 (B) 41
 (C) 42 (D) 44

4. Insert element 14 into the given hash table with double hashing? $h_1(k) = k \bmod 13$, $h_2(k) = 1 + (k \bmod 11)$. The element will occupy, which slot?

0	
1	79
2	
3	
4	69
5	98
6	
7	72
8	
9	
10	
11	50
12	

- (A) 7th (B) 8th
(C) 2nd (D) 9th
5. Consider the below given keys:
257145368, 25842354, 12487654, 248645452. Find the hash values of keys using shift folding method?
(A) 770, 221, 153, 345 (B) 221, 770, 153, 345
(C) 760, 770, 153, 345 (D) 815, 770, 153, 345
6. Consider the following two problems on unidirected graphs.
 β : Given $G(V, E)$, does G have an independent set of size $|V|-4$?
 α : Given $G(V, E)$, does G have an independent set of size 5?
Which of the following is true?
(A) β is in P and α is in NP -Complete
(B) β is in NP -Complete and α is in P
(C) Both α and β are NP -Complete
(D) Both α and β are in P
7. Let S be an NP -complete problem and Q and R be two other problems not known to be in NP . Q is polynomial-time reducible to S and S is polynomial-time reducible to R . Which one of the following statements is true?
(A) R is NP -Complete (B) R is NP -Hard
(C) Q is NP -Complete (D) Q is NP -Hard
8. Let $FHAM_3$ be the problem of finding a Hamiltonian cycle in a graph $G = (V, E)$ with $|V|$ divisible by 3 and $DHAM_3$ be the problem of determining if a Hamiltonian cycle exists in such graphs. Which of the following is true?
(A) Both $FHAM_3$ and $DHAM_3$ are NP -hard
(B) $FHAM_3$ is NP -hard but $DHAM_3$ is not
(C) $DHAM_3$ is NP -hard but $FHAM_3$ is not
(D) Neither $FHAM_3$ nor $DHAM_3$ is NP -hard
9. Consider a hash table of size 7, with starting index '0' and a hash function $(3x + 4) \bmod 7$. Initially hash table is empty. The sequence 1, 3, 8, 10 is inserted into the table using closed hashing then what is the position of element 10?
(A) 1st (B) 2nd
(C) 6th (D) 0th
10. Place the given keys in the hash table of size 13, index from '0' by using open hashing, hash function is $h(k) \bmod 13$.
Keys: A, FOOL, HIS, AND
(hint : Add the positions of a word's letters in the alphabet, take $A \rightarrow 1, B \rightarrow 2, C \rightarrow 3, D \rightarrow 4 \dots Z \rightarrow 26$).
Which of the following shows the correct hash addresses of keys?
(A) $A - 1, FOOL - 10, HIS - 9, AND - 6$
(B) $A - 1, FOOL - 9, HIS - 10, AND - 6$
(C) $A - 0, FOOL - 6, HIS - 10, AND - 9$
(D) $A - 0, FOOL - 9, HIS - 9, AND - 6$

11. Consider the following input (322, 334, 471, 679, 989, 171, 173, 199) and the hash function is $x \bmod 10$ which statement is true?
I. 679, 989, 199 hash to the same value
II. 471, 171, hash to the same value
III. Each element hashes to a different value
IV. All the elements hash to the same value
(A) I Only (B) II Only
(C) I and II (D) III
12. For the input 30, 20, 56, 75, 31, 19 and hash function $h(k) = k \bmod 11$, what is the largest number of key comparisons in a successful search in the open hash table.
(A) 4 (B) 3
(C) 5 (D) 2
13. The keys 12, 18, 13, 2, 3, 23, 5 and 15 are inserted into an empty hash table of length 10 using open addressing with hash function, $h(k) = k \bmod 10$ and linear probing.
Which is the resultant hash table?
(A)

	0
	1
2	2
23	3
13	4
15	5
	6
	7
	8
	9

 (B)

	0
3	1
12	2
13	3
	4
15	5
	6
	7
	8
	9

(C)

	0
	1
12	2
13	3
2	4
3	5
23	6
5	7
18	8
15	9

 (D)

	0
	1
2	2
3	3
12	4
13	5
23	6
5	7
18	8
15	9
14. Which one of the following is correct?
(A) Finding shortest path in a graph is solvable in polynomial time.
(B) Finding longest path from a graph is solvable in polynomial time.
(C) Finding longest path from a graph is solvable in polynomial time, if edge weights are very small values.
(D) Both (A) and (B) are correct

15. In the following pair of problems

$$\frac{2 \text{ CNF Satisfiability}}{\text{I}} \text{ Vs } \frac{3 \text{ CNF Satisfiability}}{\text{II}}$$

(A) I is solvable in polynomial time, II is NP complete problem.

(B) II is solvable in polynomial time, I is NP complete problem.

(C) Both are solvable in polynomial time

(D) None can be solved in polynomial time.

Practice Problems 2

Directions for questions 1 to 15: Select the correct alternative from the given choices.

1. For NP-complete problems

- (A) Several polynomial time algorithms are available
- (B) No polynomial time algorithm is discovered yet
- (C) Polynomial time algorithms exist but not discovered
- (D) Polynomial time algorithms will not exist, hence cannot be discovered

2. In the division method for creating hash functions, we map a key k into one of m slots by taking the remainder of k divided by m . That is, the hash function is

- (A) $h(k) = m \bmod k$
- (B) $h(k) = m \bmod m/k$
- (C) $h(k) = k \bmod m$
- (D) $h(k) = mk \bmod k$

3. In the division method for creating hash function, which of the following hash table size is most appropriate?

- (A) 2
- (B) 7
- (C) 4
- (D) 8

4. Which of the following techniques are commonly used to compute the probe sequence required for open addressing?

- (A) Linear probing
- (B) Quadratic probing
- (C) Double hashing
- (D) All the above

5. Which of the following problems is not NP-hard?

- (A) Hamiltonian circuit problem
- (B) The 0/1 knapsack problem
- (C) The graph coloring problem
- (D) None of these

6. For problems x and y , y is NP-complete and x reduces to y in polynomial time. Which of the following is true?

- (A) If x can be solved in polynomial time, then so can y
- (B) x is NP-hard
- (C) x is NP-complete
- (D) x is in NP, but not necessarily NP-complete

7. If P_1 is NP-complete and there is a polynomial time reduction of P_1 to P_2 , then P_2 is

- (A) NP-complete
- (B) Not necessarily NP-complete
- (C) Cannot be NP-complete
- (D) None of these

8. A problem is in NP, and as hard as any problem in NP. The given problem is

- (A) NP hard
- (B) NP complete
- (C) NP
- (D) NP-hard \cap NP-complete

9. Which of the following is TRUE?

- (A) All NP-complete problems are NP-hard.
- (B) If an NP-hard problem can be solved in polynomial time, then all NP-complete problems can be solved in polynomial time.
- (C) NP-hard problems are not known to be NP-complete.
- (D) All the above

10. If a polynomial time algorithm makes polynomial number of calls to polynomial time subroutines, then the resulting algorithm runs in

- (A) Polynomial time
- (B) No-polynomial time
- (C) Exponential time
- (D) None of these

11. If a polynomial time algorithm makes atmost constant number of calls to polynomial time subroutines, then the resulting algorithm runs in

- (A) Polynomial time
- (B) No-polynomial time
- (C) Exponential time
- (D) None of these

12. When a record to be inserted maps to an already occupied slot is called

- (A) Hazard
- (B) Collision
- (C) Hashing
- (D) Chaining

13. Worst-case analysis of hashing occurs when

- (A) All the keys are distributed
- (B) Every key hash to the same slot
- (C) Key values with even number, hashes to slots with even number
- (D) Key values with odd number hashes to slots with odd number

14. Main difference between open hashing and closed hashing is

- (A) Closed hashing uses linked lists and open hashing does not.
- (B) Open hashing uses linked list and closed hashing does not
- (C) Open hashing uses tree data structure and closed uses linked list
- (D) None of the above

15. The worst case scenario in hashing occurs when

- (A) All keys are hashed to the same cell of the hash table
- (B) The size of hash table is bigger than the number of keys
- (C) The size of hash table is smaller than the number of keys
- (D) None of the above

PREVIOUS YEARS' QUESTIONS

1. Consider a hash table of size seven, with starting index zero, and a hash function $(3x + 4) \bmod 7$. Assuming the hash table is initially empty, which of the following is the contents of the table when the sequence 1, 3, 8, 10 is inserted into the table using closed hashing? Note that $-$ denotes an empty location in the table.

[2007]

- (A) 8, $-$, $-$, $-$, $-$, $-$, 10 (B) 1, 8, 10, $-$, $-$, $-$, 3
(C) 1, $-$, $-$, $-$, $-$, $-$, 3 (D) 1, 10, 8, $-$, $-$, $-$, 3

Common data for questions 2 and 3: Suppose the letters a, b, c, d, e, f have probabilities $\frac{1}{2}, \frac{1}{4}, \frac{1}{8}, \frac{1}{16}, \frac{1}{32}, \frac{1}{32}$, respectively.

2. Which of the following is the Huffman code for the letter a, b, c, d, e, f ? [2007]
(A) 0, 10, 110, 1110, 11110, 11111
(B) 11, 10, 011, 010, 001, 000
(C) 11, 10, 01, 001, 0001, 0000
(D) 110, 100, 010, 000, 001, 111
3. What is the average length of the correct answer to above question? [2007]

- (A) 3 (B) 2.1875
(C) 2.25 (D) 1.9375

4. The subset-sum problem is defined as follows: Given a set S of n positive integers and a positive integer W , determine whether there is a subset of S whose elements sum to W .

An algorithm Q solves this problem in $O(nW)$ time. Which of the following statements is false?

[2008]

- (A) Q solves the subset-sum problem in polynomial time when the input is encoded in unary
(B) Q solves the subset-sum problem in polynomial time when the input is encoded in binary
(C) The subset sum problem belongs to the class NP
(D) The subset sum problem is NP -hard
5. Let π_A be a problem that belongs to the class NP . Then which one of the following is TRUE? [2009]

- (A) There is no polynomial time algorithm for π_A .
(B) If π_A can be solved deterministically in polynomial time, then $P = NP$.
(C) If π_A is NP -hard, then it is NP -complete.
(D) π_A may be undecidable.

6. The keys 12, 18, 13, 2, 3, 23, 5 and 15 are inserted into an initially empty hash table of length 10 using open addressing with hash function $h(k) = k \bmod 10$

and linear probing. What is the resultant hash table? [2009]

(A)

0	
1	
2	12
3	23
4	
5	15
6	
7	
8	18
9	

(B)

0	
1	
2	12
3	13
4	
5	5
6	
7	
8	18
9	

(C)

0	
1	
2	12
3	13
4	2
5	3
6	23
7	5
8	18
9	15

(D)

0	
1	
2	12,2
3	13,3,23
4	
5	5,15
6	
7	
8	18
9	

Common data for questions 7 and 8: A sub-sequence of a given sequence is just the given sequence with some elements (possibly none or all) left out. We are given two sequences $X[m]$ and $Y[n]$ of lengths m and n , respectively, with indexes of X and Y starting from 0.

7. We wish to find the length of the longest common subsequence (LCS) of $X[m]$ and $Y[n]$ as $l(m, n)$, where an incomplete recursive definition for the function $l(i, j)$ to compute the length of the LCS of $X[m]$ and $Y[n]$ is given below:

$$l(i, j) = 0, \text{ if either } i = 0 \text{ or } j = 0$$

$$= \text{expr1}, \text{ if } i, j > 0 \text{ and } X[i - 1] = Y[j - 1]$$

$$= \text{expr2}, \text{ if } i, j > 0 \text{ and } X[i - 1] \neq Y[j - 1]$$

Which one of the following options is correct? [2009]

- (A) $\text{expr1} \equiv l(i - 1, j) + 1$
(B) $\text{expr1} \equiv l(i, j - 1)$
(C) $\text{expr2} \equiv \max(l(i - 1, j), l(i, j - 1))$
(D) $\text{expr2} \equiv \max(l(i - 1, j - 1), l(i, j))$

8. The values of $l(i, j)$ could be obtained by dynamic programming based on the correct recursive definition of $l(i, j)$ of the form given above, using an array $L[M, N]$, where $M = m + 1$ and $N = n + 1$, such that $L[i, j] = l(i, j)$.

Which one of the following statements would be TRUE regarding the dynamic programming solution for the recursive definition of $l(i, j)$? [2009]

- (A) All elements of L should be initialized to 0 for the values of $l(i, j)$ to be properly computed.

- (B) The values of $l(i, j)$ may be computed in a row major order or column major order of $L(M, N)$.
- (C) The values of $l(i, j)$ cannot be computed in either row major order or column major order of $L(M, N)$.
- (D) $L[p, q]$ needs to be computed before $L[r, s]$ if either $p < r$ or $q < s$.

9. The weight of a sequence a_0, a_1, \dots, a_{n-1} of real numbers is defined as $a_0 + a_1/2 + \dots + a_{n-1}/2^{n-1}$. A subsequence of a sequence is obtained by deleting some elements from the sequence, keeping the order of the remaining elements the same. Let X denote the maximum possible weight of a subsequence of a_0, a_1, \dots, a_{n-1} . Then X is equal to **[2010]**

- (A) $\max(Y, a_0 + Y)$
- (B) $\max(Y, a_0 + Y/2)$
- (C) $\max(Y, a_0 + 2Y)$
- (D) $a_0 + Y/2$

10. Four matrices M_1, M_2, M_3 and M_4 of dimensions $p \times q, q \times r, r \times s$ and $s \times t$ respectively, can be multiplied in several ways with different number of total scalar multiplications. For example when multiplied as $((M_1 \times M_2) \times (M_3 \times M_4))$, the total number of scalar multiplications is $pqr + rst + prt$. When multiplied $((M_1 \times M_2) \times M_3) \times M_4$ the total number of scalar multiplications is $pqr + prs + pst$.

If $p = 10, q = 100, r = 20, s = 5$ and $t = 80$, then the minimum number of scalar multiplications needed is **[2011]**

- (A) 248000
- (B) 44000
- (C) 19000
- (D) 25000

11. Assuming $P \neq NP$, which of the following is **TRUE?** **[2012]**

- (A) NP -complete = NP
- (B) NP -complete $\cap P = \emptyset$
- (C) NP -hard = NP
- (D) $P = NP$ -complete

12. Which of the following statements are TRUE?

- (i) The problem of determining whether there exists a cycle in an undirected graph is in P .
- (ii) The problem of determining whether there exists a cycle in an undirected graph is in NP .
- (iii) If a problem A is NP -Complete, there exists a non-deterministic polynomial time algorithm to solve A .

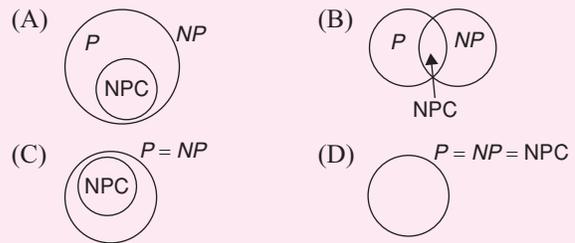
[2013]

- (A) 1, 2 and 3
- (B) 1 and 2 only
- (C) 2 and 3 only
- (D) 1 and 3 only

13. Suppose a polynomial time algorithm is discovered that correctly computes the largest clique in a given graph. In this scenario, which one of the following represents the correct Venn diagram of the complexity classes P, NP and NP -complete

(NPC)?

[2014]



14. Consider a hash table with 9 slots. The hash function is $h(K) = K \bmod 9$. The collisions are resolved by chaining. The following 9 keys are inserted in the order: 5, 28, 19, 15, 20, 33, 12, 17, 10. The maximum, minimum, and average chain lengths in the hash table, respectively, are **[2014]**

- (A) 3, 0 and 1
- (B) 3, 3 and 3
- (C) 4, 0 and 1
- (D) 3, 0 and 2

15. Consider two strings $A = \text{'qpqrr'}$ and $B = \text{'pqpqrqp'}$. Let x be the length of the longest common subsequence (not necessarily contiguous between A and B) and let y be the number of such longest common subsequences between A and B . then $x + 10y =$ **[2014]**

16. Suppose you want to move from 0 to 100 on the number line. In each step, you either move right by a unit distance or you take a *shortcut*. A shortcut is simply a pre-specified pair of integers i, j with $i < j$. Given a shortcut i, j if you are at position i on the number line, you may directly move to j . Suppose $T(k)$ denotes the smallest number of steps needed to move from k to 100. Suppose further that there is at most 1 shortcut involving any number, and in particular from 9 there is a shortcut to 15. Let y and z be such that $T(9) = 1 + \min(T(y), T(z))$. Then the value of the product yz is **[2014]**

17. Consider the decision problem 2CNFSAT defined as follows: **[2014]**

$\{\phi \mid \phi \text{ is a satisfiable propositional formula in CNF with at most two literals per clause}\}$

For example, $\phi = (x_1 \vee x_2) \wedge (x_1 \vee x_3) \wedge (x_2 \vee x_4)$ is a Boolean formula and it is in 2CNFSAT.

The decision problem 2CNFSAT is

- (A) NP -complete
- (B) Solvable in polynomial time by reduction to directed graph reachability.
- (C) Solvable in constant time since any input instance is satisfiable.
- (D) NP -hard, but not NP -complete

18. Consider a hash table with 100 slots. Collisions are resolved using chaining. Assuming simple uniform

hashing, what is the probability that the first 3 slots are unfilled after the first 3 insertions? [2014]

- (A) $(97 \times 97 \times 97)/100^3$
- (B) $(99 \times 98 \times 97)/100^3$
- (C) $(97 \times 96 \times 95)/100^3$
- (D) $(97 \times 96 \times 95)/(3! \times 100^3)$

19. Match the following [2014]

(P) prim's algorithm for minimum spanning tree	(i) Backtracking
(Q) Floyd-Warshall algorithm for all pairs shortest paths	(ii) Greedy method
(R) Mergesort	(iii) Dynamic programming
(S) Hamiltonian circuit	(iv) Divide and conquer

- (A) P-iii, Q-ii, R-iv, S-i
- (B) P-i, Q-ii, R-iv, S-iii
- (C) P-ii, Q-iii, R-iv, S-i
- (D) P-ii, Q-i, R-iii, S-iv

20. Given a hash table T with 25 slots that stores 2000 elements, the load factor ∞ for T is _____ [2015]

21. Language L_1 is polynomial time reducible to language L_2 . Language L_3 is polynomial time reducible to L_2 , which in turn is polynomial time reducible to language L_4 . Which of the following is/are true? [2015]

- (1) if $L_4 \in P$, then $L_2 \in P$
- (2) if $L_1 \in P$ or $L_3 \in P$, then $L_2 \in P$
- (3) $L_1 \in P$, if and only if $L_3 \in P$
- (4) if $L_4 \in P$, then $L_1 \in P$ and $L_3 \in P$

22. The Floyd - Warshall algorithm for all -pair shortest paths computation is based on [2016]

- (A) Greedy paradigm
- (B) Divide-and-Conquer paradigm
- (C) Dynamic Programming paradigm
- (D) Neither Greedy nor Divide-and-Conquer nor Dynamic Programming paradigm.

23. Let $A_1, A_2, A_3,$ and A_4 be four matrices of dimensions $10 \times 5, 5 \times 20, 20 \times 10,$ and $10 \times 5,$ respectively. The minimum number of scalar multiplications required to find the product $A_1 A_2 A_3 A_4$ using the basic matrix multiplication method is _____. [2016]

24. Consider the following table:

Algorithms	Design Paradigms
(P) Kruskal	(i) Divide and Conquer
(Q) Quicksort	(ii) Greedy
(R) Floyd-Warshall	(iii) Dynamic Programming

Match the algorithms to the design paradigms they are based on. [2017]

- (A) (P) \leftrightarrow (ii), (Q) \leftrightarrow (iii), (R) \leftrightarrow (i)
- (B) (P) \leftrightarrow (iii), (Q) \leftrightarrow (i), (R) \leftrightarrow (ii)
- (C) (P) \leftrightarrow (ii), (Q) \leftrightarrow (i), (R) \leftrightarrow (iii)
- (D) (P) \leftrightarrow (i), (Q) \leftrightarrow (ii), (R) \leftrightarrow (iii)

25. Assume that multiplying a matrix G_1 of dimension $p \times q$ with another matrix G_2 of dimension $q \times r$ requires pqr scalar multiplications. Computing the product of n matrices $G_1 G_2 G_3, \dots, G_n$ can be done by parenthesizing in different ways. Define $G_i G_{i+1}$ as an explicitly computed pair for a given paranthesization if they are directly multiplied. For example, in the matrix multiplication chain $G_1 G_2 G_3 G_4 G_5 G_6$ using parenthesization $(G_1(G_2 G_3))(G_4(G_5 G_6))$, $G_2 G_3$ and $G_5 G_6$ are the only explicitly computed pairs.

Consider a matrix multiplication chain $F_1 F_2 F_3 F_4 F_5$, where matrices $F_1, F_2, F_3, F_4,$ and F_5 are of dimensions $2 \times 25, 25 \times 3, 3 \times 16, 16 \times 1$ and $1 \times 1000,$ respectively. In the parenthesization of $F_1 F_2 F_3 F_4 F_5$ that minimizes the total number of scalar multiplications, the explicitly computed pairs is/are: [2018]

- (A) $F_1 F_2$ and $F_3 F_4$ only
- (B) $F_2 F_3$ only
- (C) $F_3 F_4$ only
- (D) $F_1 F_2$ and $F_4 F_5$ only

26. Consider the weights and values of items listed below. Note that there is only one unit of each item.

Item no.	Weight (in Kgs)	Value (in Rupees)
1	10	60
2	7	28
3	4	20
4	2	24

The task is to pick a subset of these items such that their total weight is no more than 11 kgs and their total value is maximized. Moreover, no item may be split. The total value of items picked by an optimal algorithm is denoted by V_{opt} . A greedy algorithm sorts the items by their value-to-weight ratios in descending order and packs them greedily, starting from the first item in the ordered list. The total value of items picked by the greedy algorithm is denoted by V_{greedy} .

The value of $V_{opt} - V_{greedy}$ is _____. [2018]

ANSWER KEYS**EXERCISES****Practice Problems 1**

1. B 2. A 3. B 4. D 5. A 6. C 7. C 8. A 9. B 10. B
11. C 12. B 13. C 14. A 15. A

Practice Problems 2

1. B 2. C 3. B 4. D 5. B 6. C 7. A 8. B 9. D 10. C
11. A 12. B 13. B 14. B 15. A

Previous Years' Questions

1. B 2. A 3. D 4. B 5. C 6. C 7. C 8. B 9. C 10. B
11. B 12. A 13. D 14. A 15. 34 16. 150 17. B 18. A 19. C 20. 80
21. C 22. C 23. 1500 24. C 25. C 26. 16

TEST

ALGORITHMS (PART 2)

Time: 45 min.

Directions for questions 1 to 30: Select the correct alternative from the given choices.

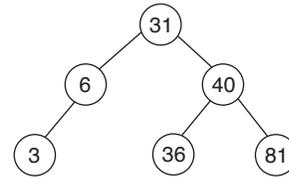
- The worst case running time of an algorithm means
 - The algorithm will never take any longer.
 - The algorithm will take less time than running time
 - The algorithm will run in a finite time
 - None of the above
- Analyzing an algorithm involves
 - Evaluating the complexity
 - Validating the Algorithm
 - Both A and B
 - None of the above
- $f(n) = O(g(n))$ is
 - $g(n)$ is asymptotic lower bound for $f(n)$
 - $g(n)$ is asymptotic tight bound for $f(n)$
 - $g(n)$ is asymptotic upper bound for $f(n)$
 - None of the above
- Which case yields the necessary information about an algorithm's behaviour on a random input?
 - Best-case
 - Worst-case
 - Average-case
 - Both A and C
- Algorithms that require an exponential number of operations are practical for solving.
 - Only problems of very small size
 - Problems of large size
 - Problems of any size
 - None of these
- Problems that can be solved in polynomial time are called
 - Tractable
 - Decidable
 - Solvable
 - Computable
- Problems that cannot be solved at all by any algorithm are known as
 - Tractable
 - Undecidable
 - Untractable
 - Unsolvable
- Which of the following problems is decidable but intractable?
 - Hamiltonian circuit
 - Traveling sales man
 - Knapsack problem
 - All the above
- Which method is used to solve recurrences?
 - Substitution method
 - Recursion-tree method
 - Master method
 - All the above
- Consider the following
 - Input
 - Output
 - Finiteness
 - Definiteness means clear and unambiguous
 - Effectiveness
 Which of the following is not a property of an algorithm?
 - (iv) only
 - (iv) and (v) only
 - (iii) and (iv) only
 - None of the above
 - Finiteness of an algorithm means
 - The steps of the algorithm should be finite
 - The algorithm should terminate after finite time
 - Algorithm must terminate after a finite number of steps
 - Algorithm should consume very less space
 - Asymptotic analysis on efficiency of algorithm means
 - The efficiency of the algorithm on a particular machine
 - How the running time of an algorithm increases as the size increases without bound
 - How efficiently the algorithm is applied to solve a problem without thinking of input size.
 - None of the above
 - What is the input size of a problem?
 - Number of variables used to solve the problem
 - Number of constants used to solve the problem
 - it is problem specific that is in case of graph it is number of edges and vertices and so on.
 - None of these
 - An algorithm must take input
 - An algorithm must give out put
 Which is true in the following options?
 - (i) Only
 - (ii) Only
 - (i) and (ii) Only
 - None of the above
 - As $n \rightarrow \infty$
 Which of the following is efficient?
 - (n^3)
 - (n^2)
 - (2^n)
 - (n^4)
 - Suppose

$$T_1(n) = O(f(n))$$

$$T_2(n) = O(f(n))$$
 which of the following is true,
 - $T_1(n) + T_2(n) = O(f(n))$
 - $\frac{T_1(n)}{T_2(n)} = O(1)$
 - $T_1(n) = O(T_2(n))$
 - None of these

17. The following program computes $n!$
 Find the complexity?
 Input: A non-negative integer
 Output: Value of $n!$
 If $n = 0$ return 1
 Else return $F(n - 1) \cdot n$
 (A) (n) (B) $(n \log n)$
 (C) (n^2) (D) (n^3)
18. Which of the following functions are often referred as 'exponential growth function'?
 (A) $2^n, \log n$ (B) $2^n, n!$
 (C) $n!, n \log n$ (D) $n!, \log n$
19. Consider the following code
`sort(a, n)`
`{`
`for i = 1 to n do`
`{`
`j = i;`
`for k = i + 1 to n do`
`if (a[k] < a[j]) then j = k;`
`t = a[i];`
`a[i] = a[j];`
`a[j] = t;`
`}`
`}`
`}`
 The above code implements which sorting?
 (A) Merge sort
 (B) selection sort
 (C) Insertion sort
 (D) Radix sort
20. Assume that the number of disks in a 'Towers of Hanoi problem' is ' n ', with '3' towers, Initially all disks are placed on tower 1, to get the largest disk to the bottom of 2nd tower, How many moves are required? ($n = 3$)
 (A) n
 (B) $(n - 1)$
 (C) $(n + 1)$
 (D) $2n$
21. Each new term in Fibonacci sequence is obtained by taking the sum of the two previous terms. The first term of the sequence is $f_0 = 0$, and the second term $f_1 = 1$. Which of the following gives Fibonacci sequence?
 (A) $f_n = f_{n+1} + f_{n-2}, n \geq 2$
 (B) $f_n = f_{n-1} + f_{n-2}, n \geq 2$
 (C) $f_n = f_{n-1} + f_{n+1}, n \geq 2$
 (D) All the above

22. Consider the binary search tree



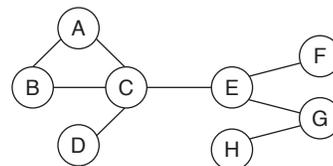
Delete node '31', what would be the parent node in the new binary search tree?

- (A) 36
 (B) 40
 (C) 81
 (D) 6
23. Consider the given array [4, 6, 7, 8, 21, 9, 3, 10, 13, 16, 31] after performing '1' delete max operation, on the max heap. What would be the sequence of elements in the array?
 (A) 9, 21, 13, 16, 3, 7, 10, 8, 4, 6
 (B) 21, 9, 13, 16, 7, 3, 10, 8, 4, 6
 (C) 21, 9, 13, 16, 3, 7, 10, 8, 4, 6
 (D) 21, 9, 13, 16, 7, 3, 10, 4, 8, 6
24. Consider the given Di-graph



How many strongly connected components does the above graph contain?

- (A) 1 (B) 2
 (C) 3 (D) many
25. Consider the given graph



Which of the following shows the adjacency matrix of the above graph?

$$(A) \begin{matrix} & A & B & C & D & E & F & G & H \\ \begin{matrix} A \\ B \\ C \\ D \\ E \\ F \\ G \\ H \end{matrix} & \begin{bmatrix} 0 & 1 & 1 & 0 & 0 & 0 & 0 & 0 \\ 1 & 0 & 1 & 0 & 0 & 0 & 0 & 0 \\ 1 & 1 & 0 & 1 & 1 & 0 & 0 & 0 \\ 0 & 0 & 1 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 1 & 0 & 0 & 1 & 1 & 0 \\ 0 & 0 & 0 & 0 & 1 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 1 & 0 & 0 & 1 \\ 0 & 0 & 0 & 0 & 0 & 0 & 1 & 0 \end{bmatrix} & \end{matrix}$$

(B)

	A	B	C	D	E	F	G	H
A	0	1	1	0	0	0	0	0
B	1	0	1	0	0	0	0	0
C	1	1	0	1	1	0	0	0
D	0	0	1	0	0	0	0	0
E	0	0	1	0	0	1	1	0
F	0	0	0	0	1	0	0	0
G	0	0	0	0	1	0	0	1
H	0	0	0	0	1	0	1	0

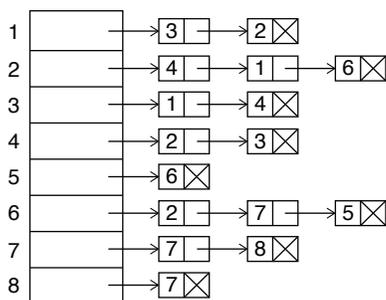
(C)

	A	B	C	D	E	F	G	H
A	0	1	1	0	0	0	0	0
B	1	0	1	0	0	0	0	0
C	1	1	0	1	1	0	0	0
D	0	0	1	0	0	0	0	0
E	0	0	1	0	0	1	1	0
F	0	0	0	0	1	0	0	0
G	0	0	0	0	0	1	0	1
H	0	0	0	0	0	0	1	0

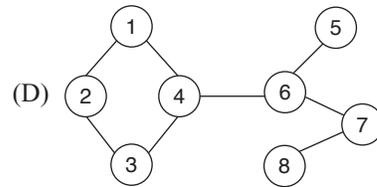
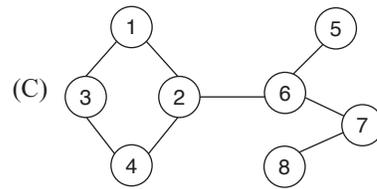
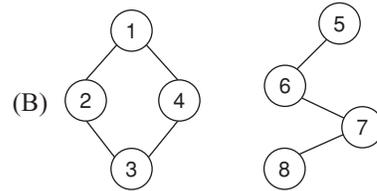
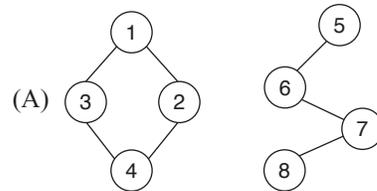
(D)

	A	B	C	D	E	F	G	H
A	0	1	1	0	0	0	0	0
B	1	0	1	0	0	0	0	0
C	1	1	0	1	1	0	0	0
D	0	0	1	0	0	0	0	0
E	0	0	1	0	0	1	1	0
F	0	0	0	0	1	0	0	0
G	0	0	0	0	1	0	0	1
H	0	0	0	0	1	1	0	1

26. Consider the given adjacency list



The above list is representation of which of the following graph?



27. Which of the following is FALSE?

- (A) In dynamic programming an optimal sequence of decisions is obtained by making explicit appeal to the principle of optimality
- (B) In greedy method only one decision sequence is generated.
- (C) In dynamic programming, many decision sequences may be generated.
- (D) In greedy method many decision sequences are generated.

28. Consider an array $a[n]$ of 'n' numbers that has ' $n/2$ ' distinct elements and ' $n/2$ ' copies of another element, to identify that repeated element, how many steps are required in the worst case?

- (A) $n/2$
- (B) $n/2 + 1$
- (C) $n/2 + 2$
- (D) n

29. Match the following, for a very large value of 'n'

- I. $36n^3 + 2n^2$
- II. $5n^2 - 6n$
- III. $n^{1.001} + n \log n$
- P. (n^2)
- Q. $\lfloor (n^3) \rfloor$
- R. $(n^{1.001})$

- (A) I - P, II - Q, III - R
- (B) I - Q, II - P, III - R
- (C) I - R, II - Q, III - P
- (D) I - R, II - P, III - R

3.154 | Unit 3 • Algorithms

30. Consider the following code

```
T(a, n)
{
  for i = 1 to n - 1 do
    for j = i + 1 to n do
      {
        t = a[i, j];
        a[i, j] = a[j, i];
        a[j, i] = t;
      }
}
```

The above code performs

- (A) Matrix multiplication
- (B) Matrix addition
- (C) Matrix transpose
- (D) Matrix chain multiplication

ANSWERS KEYS

- | | | | | | | | | | |
|-------|-------|-------|-------|-------|-------|-------|-------|-------|-------|
| 1. A | 2. C | 3. B | 4. C | 5. A | 6. A | 7. B | 8. D | 9. D | 10. D |
| 11. C | 12. D | 13. C | 14. B | 15. B | 16. A | 17. A | 18. B | 19. A | 20. C |
| 21. B | 22. A | 23. B | 24. B | 25. A | 26. C | 27. D | 28. C | 29. B | 30. C |

Databases

Chapter 1: ER Model and Relational Model	4.3
Chapter 2: Structured Query Language	4.21
Chapter 3: Normalization	4.49
Chapter 4: Transaction and Concurrency	4.65
Chapter 5: File Management	4.82

U

N

I

T

4

This page is intentionally left blank

Chapter 1

ER Model and Relational Model

LEARNING OBJECTIVES

- ☞ Data model
- ☞ Schemas
- ☞ Three-schema architecture
- ☞ ER model
- ☞ Types of attributes
- ☞ Mapping cardinality
- ☞ Complex attributes
- ☞ Entity types, entity sets and value sets
- ☞ Weak entity set
- ☞ Relational database
- ☞ NULL in tuples
- ☞ Inherent constraint
- ☞ Referential and entity integrity constraint

INTRODUCTION

A database is a collection of related data. By data, we mean facts that can be recorded and that have implicit meaning.

Example: Consider the names, telephone numbers and addresses of the people. We can record this data in an indexed address book and store it as Excel file on a hard drive using a personal computer. This is a collection of related data with an implicit meaning and hence is a database.

A database management system (DBMS) is a collection of programs that enables users to create and maintain a database. The DBMS is a general-purpose software system that facilitates the processes of defining, constructing, manipulating and sharing databases among various users and applications.

1. Defining the database involves specifying the data types, structures, and constraints for the data to be stored in the database.
2. Constructing the database is the process of storing the data itself on some storage medium that is controlled by the DBMS
3. Manipulating a database includes such functions as querying the database to retrieve specific data, updating the database to reflect changes.
4. Sharing a database allows multiple users and programs to access the database concurrently.
5. Fundamental characteristics of the database approach is that it provides some level of data abstraction by hiding details of data storage that are not needed by users.

Data Model

A data model is a collection of concepts that can be used to describe the structure of a database. It provides the necessary means to achieve abstraction.

SCHEMAS

In any data model, it is important to distinguish between the description of the database and the database itself. The description of a database is called the *database schema*, which is specified during database design and is not expected to change frequently.

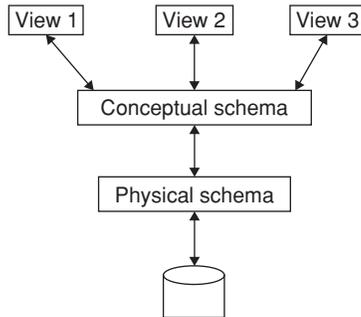
The actual data in a database may change frequently, for example the student database changes every time we add a student or enter a new grade for a student. The data in the database at a particular moment in time is called a *database state* or *snapshot*. It is also called the *current set of occurrences* or *instances in the database*.

The distinction between database schema and database state is very important. When we define a new database, we specify its database schema only to the DBMS. At this point, the corresponding database state is the empty state with no data. We get the initial state of the database when the database is first loaded with the initial data. The DBMS stores the description of the schema constructs and constraints, also called the *metadata* in the DBMS catalog so that DBMS software can refer to the schema whenever it needs. The schema is sometimes called the *intension*, and the database state an *extension* of the schema.

Three-schema Architecture

The goal of the three-schema architecture is to separate the user applications and the physical database.

Levels of Abstraction



1. The external or view level includes a number of external schemas or user views. Each external schema describes the part of the database that a particular user group is interested in and hides the rest of the database.
2. The conceptual level has a conceptual schema, which describes the structure of the whole database for a community of users. The conceptual schema hides the details of physical storage structures and concentrates on describing entities, data types, relationships, user operations, and constraints.
3. The internal level has an internal schema, which describes the physical storage structures of the database. It describes the complete details of data storage and access paths for the database.

ER MODEL

Entity relationship model is a popular high-level conceptual data model. This model and its variations are used for the conceptual design of database applications, and many database design tools employ its concepts. ER model describes data as entities, relationships and attributes. The basic object that the ER model represents is an entity.

Entity

It is an object that exists and is distinguishable from other objects

(or)

Entity is a “thing” in the real world with an independent existence.

(or)

An entity is something that has a distinct, separate existence, although it need not be a material existence. In particular, abstractions and legal functions are usually regarded

as entities. In general, there is no presumption that an entity is animate.

1. An object with a physical existence
Example: A particular person, car, house, employee.
2. An object with a conceptual existence
Example: A company, a job, a university course. Each entity has attributes, the particular properties that describe it.
Example: An employee entity can be described by employee’s name, age, address, salary and job.

Entity Set

Set of entities of same type that shares the same properties.

Example: All persons, all companies etc.

Example: Entity sets of customer and loan

Table 1 Customer Entity Set

Customer-id	Cust-name	Cust-street	City
C-143	John	MG Road	Sec.bad
C-174	Mary	SP Road	Hyd.bad
C-183	Tony	KD Road	Sec.bad
C-192	Satya	SG Road	Eluru

Table 2 Loan Entity Set

Loan-no	Amount
L-30	\$3000
L-31	\$4000
L-32	\$3500
L-33	\$4500
L-34	\$5000

An entity is represented by a set of attributes and by a descriptive properties possessed by all members of an entity set.

Types of Attributes

1. Simple versus composite
2. Single valued versus multivalued
3. Stored versus derived.

Composite attribute: Composite attributes can be divided into smaller subparts, which represent more basic attribute that has their own meaning (Figure 1).

Example: A common example is Address, it can be broken down into a number of subparts such as street address, city, postal code; street address is further broken down into Number, street Name and Apartment number.

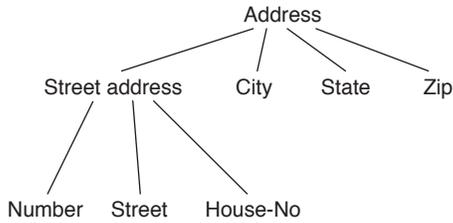


Figure 1 A hierarchy of composite attributes.

Street address is a composite attribute. Attributes that are not divisible are called simple (or) atomic attributes.

Single-valued versus multivalued attributes: Most attributes have a single value for a particular entity, such attributes are called *single-valued attribute*.

Example: Age is a single-valued attribute

Multivalued attributes: An attribute can have a set of values for the same entity.

Example: College degrees attribute for a person

Example: Name is also a multivalued attribute (Figure 2).

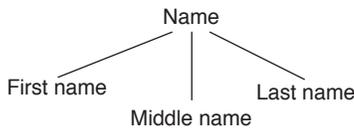


Figure 2 Multivalued attribute.

Stored versus derived attributes: Two (or) more attribute values are related.

Example: Age can be derived from a person’s date of birth. The age attribute is called derived attribute and is said to be derivable from the DOB attribute, which is called a stored attribute.

Domain: The set of permitted values for each attribute.

Example: A person’s age must be in the domain {0-130}

RELATIONSHIP SETS

A relationship is an association among several entities. Relationship sets that involve two entity sets are binary. Generally, most relationships in databases are binary. Relationship sets may involve more than two entity sets.

Example: Employee of a bank may have responsibilities at multiple braches, with different jobs at different branches, then there is a ternary relation between employee, job and branch.

Mapping Cardinality

For a binary relationship set, mapping cardinality must be:

1. One-to-one
2. One-to-many
3. Many-to-one
4. Many-to-many

One-to-one: An entity in *A* is associated with at most one entity in *B* and an entity in *B* is associated with at most one entity in *A* (Figure 3).

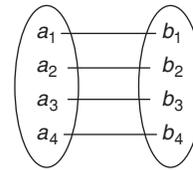


Figure 3 One-to-one relationship set.

One-to-many: An entity in *A* is associated with any number of entities in *B*. But an entity in *B* is associated with at most one entity in *A* (Figure 4).

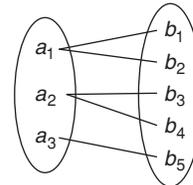


Figure 4 One-to-many relationship set.

Many-to-one: An entity in *A* is associated with at most one entity in *B*. But an entity in *B* can be associated with any number of entities in *A* (Figure 5).

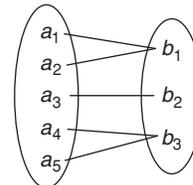


Figure 5 Many-to-one relationship set.

Many-to-many: An entity in *A* is associated with any number of entities in *B*. But an entity in *B* can be associated with any number of entities in *A* (Figure 6).

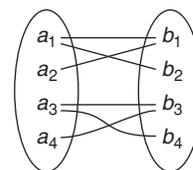


Figure 6 Many-to-many relationship set.

Example: One customer can have multiple accounts
Customer(c-Name) (Acc. no, Amount)

Table 3 Example of One-to-many Relationship Set

Arun	A-101	\$3000
Bunny	A-102	\$3500
Kate	A-103	\$2000
Mary	A-104	\$2500
John	A-105	\$4000

In Table 3, many-to-one relationship is not possible.

Complex Attributes

Composite and multivalued attributes can be nested in an arbitrary way. We can represent arbitrary nesting by grouping components of a composite attribute between parentheses () and separating the components with commas, and by displaying multivalued attributes between braces {}. Such attributes are called complex attributes.

Example: A person can have more than one residence and each residence can have multiple phones, an attribute AddressPhone for a person can be specified as shown below. {AddressPhone

({Phone (Areacode, phoneNumber)}, Address (StreetAddress (StreetNumber, streetName, ApartmentNumber), city, state, zip)))}

Entity types, entity sets and value sets

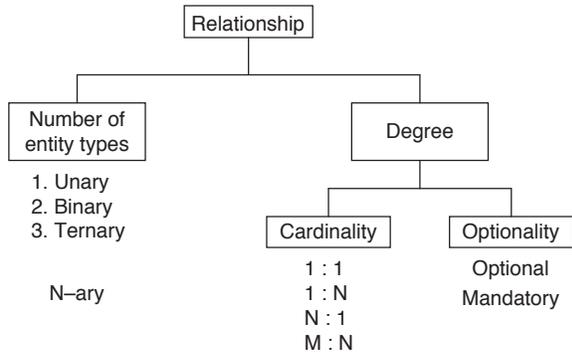
An entity type defines a collection of entities that have the same attributes. Each entity type in the database is described by its name and attribute. The following figure shows two entity types, named STUDENT and EMPLOYEE and a list of attributes for each

ENTITY	TYPE NAME	STUDENT	EMPLOYEE
	ATTRIBUTES:	R.No, Name, Grade	Name, Salary, Age
	ENTITY SET (EXTENSION)	<div style="border: 1px solid black; padding: 5px;"> <p>S₁. (86, Arun, A)</p> <p>S₂. (87, Pavan, B)</p> <p>S₃. (89, Karan, A)</p> </div>	<div style="border: 1px solid black; padding: 5px;"> <p>e₁. (Kamal, 20K, 42)</p> <p>e₂. (Bharat, 25K, 41)</p> <p>e₃. (Bhanu, 26K, 41)</p> </div>

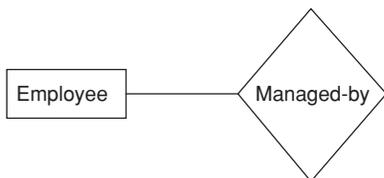
The collection of all entities of a particular entity type in the database at any point in time is called an *entity set*.

Types of relations

1. Unary relation.
2. Binary relation.
3. Ternary relation.
4. Quadrary relation.
5. N-ary relation



Unary relation If a relationship type is between entities in a single entity type then it is called a unary relationship type.



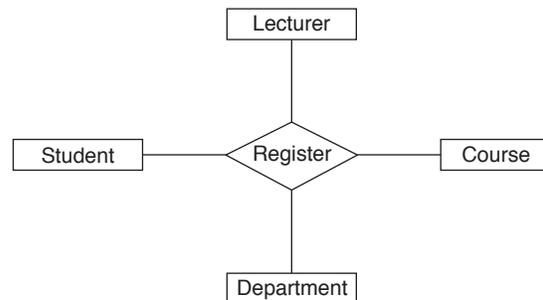
In employee entity, we will have all the employees including 'manager', this relation indicates, employees are managed by manager.

Binary relation If a relationship type is between entities in one type and entities in another type then it is called a *binary relation*, because two entity types are involved in the relation.



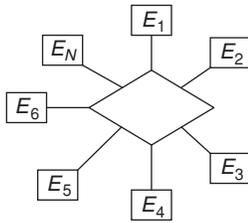
The above relation indicates that customers purchased product (or) products are purchased by customers.

Quadrary relation If a relationship type is among entities of four different types, then it is called *quadrary relation*.



In the above ER-diagram, any two entities can have a relation.

N-ary relation 'N' number of entities will participate in a relation, and each entity can have a relation with all the other entities.



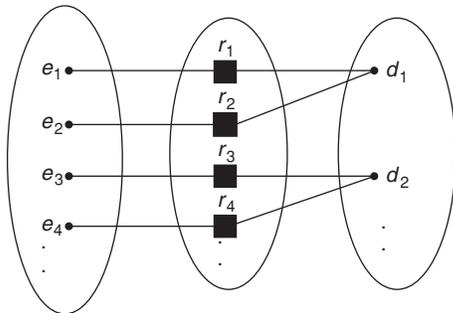
CARDINALITY RATIO AND PARTICIPATION CONSTRAINTS

The cardinality ratio for a binary relationship specifies the maximum number of relationship instances that an entity can participate in

1. The participation constraint specifies whether the existence of an entity depends on its being related to another entity via the relationship type. This constraint specifies the minimum number of relationship instances that each entity can participate in, and is some times called the *minimum cardinality constraint*.
2. There are two types of participation constraints:
 - a. Total participation
 - b. Partial participation

Example: If a company policy states that every employee must work for a department, then an employee entity can exist only if it participates (or) works for at least one department.

EMPLOYEE WORKS-FOR DEPARTMENT



Every entity in the EMPLOYEE set must be related to a DEPARTMENT entity via WORK-FOR. Total participation is also called *existence dependency*. If we do not expect every employee to manage a department, so the participation of EMPLOYEE in the relationship type is partial, means not necessarily all employees' entities are related to some department entity.

Cardinality ratio and participation constraints are taken together as the structural constraints of a relationship type. In ER diagrams, total participation is displayed as a double line connecting the participating entity type to the

relationship, whereas partial participation is represented by a single line.

ER Diagrams (Figure 7 and 8)

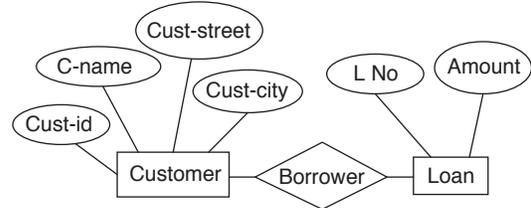


Figure 7 ER diagram.

Notations:

- Represents entity
- Represents relationship sets
- Represents links between attributes to entity sets and entity sets to relationship sets
- Represents attributes
- Represents multivalued attributes
- Represents primary key attribute

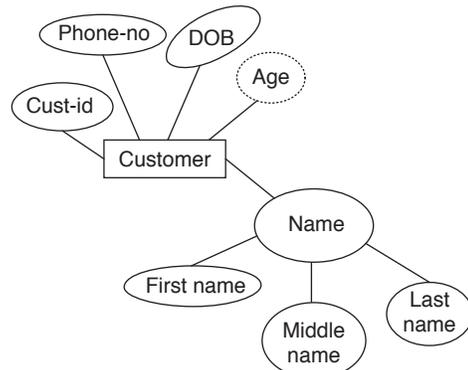
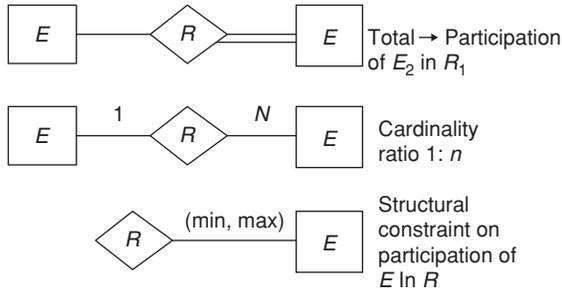


Figure 8 ER diagram.

- Derived attribute
- Multivalued attribute
- Weak entity
- Weak entity set relation
- Composite attribute



Cardinality Constraints

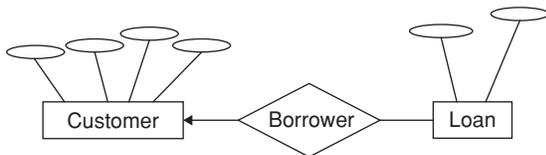
One-to-one

Each entity of one entity set is related to at most one entity of the other set. Only one matching record exists between two tables.

Example: Assume each owner is allowed to have only one dog and each dog must belong to one owner. The own relationship between dog and owner is one-to-one. One-to-one relationships can often combine the data into one table.

Examples:

1. One birdfeeder is located in one place in the yard.
2. One state has one governor.
3. One yard has one address.
4. One patient has one phone number.
2. One student has one ID.



One customer is associated to one loan via borrower

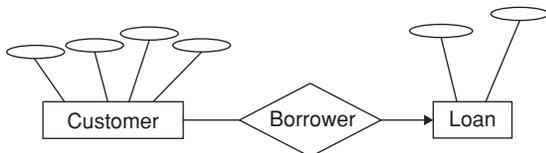
One-to-many

Examples:

1. One birdfeeder is visited by many birds.
2. One student can have many degrees.
3. One Book can be written by many authors.
4. One yard contains many bird feeders.
5. One patient has many prescriptions.

In the one-to-many relationship, a loan is associated with one customer via borrower.

Many-to-one

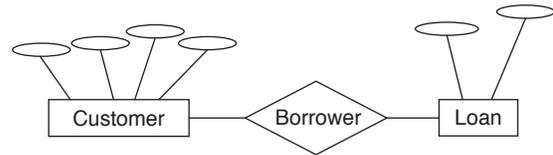


A customer is associated with at most one loan via borrower.

Many-to-many

Examples:

1. Many students are taught by many teachers.
2. Many patients are treated by many doctors.
3. Many medications are taken by many patients.
4. Many customers buy many products.
5. Many books are written by many authors.



Customer is associated with several loans and loan is associated with several customers.

ER Diagram with a Ternary Relationship

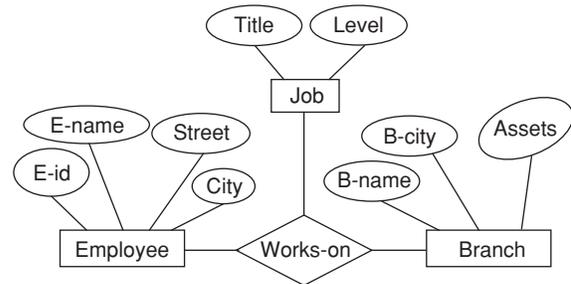


Figure 9 ER diagram with a ternary relationship.

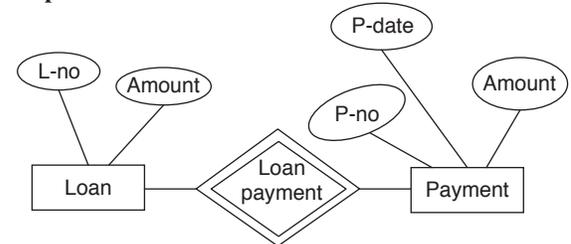
Weak Entity Set

An entity set that does not have a primary key is called *weak entity set*.

Weak entity set is represented by \rightarrow

Underline the primary key of a weak entity with a dashed line.

Example:



RELATIONAL DATABASE

Relational Model

The relational model represents the database as a collection of relations. When a relation is thought of as a table of values, each row in the table represents a collection of related data values. Each row in the table represents a fact that typically corresponds to a real-world entity. The table name and

column names are used to help in interpreting the meaning of the values in each row.

In formal relational model terminology, a row is called a *tuple*, a column header is called an *attribute*, and the table is called a *relation*.

Domain

A domain is a set of atomic values. A common method of specifying a domain is to specify a data type from which the data values forming the domain are drawn. It is also useful to specify a name for the domain, to help in interpreting its values

Example:

1. *Set of telephone numbers*: The set of valid numbers in a particular country.
2. *Employee id numbers*: The set of valid employee numbers in a company.
3. *Names*: The set of character strings that represent names of persons.
4. *Grade-point average*: Possible values of computed grade point averages, each must be real (floating point) number between 0 and 4
5. *Research department names*: The set of research department names in a specialization, such as computer science, chemistry and applied mathematics.
6. *Research department codes*: The set of Research department codes, such as CS, CHE, AM.

The preceding is called *logical definitions of domains*. The data type for research department names is set of character strings that represent valid department names. A domain is thus given a name, data type, and format. Additional information for interpreting the values of a domain can also be given, for example a numeric domain such as person weights should have the units of measurements, such as kilograms or pounds.

1. The relational model is often described as having the following three aspects:
 - *Structural aspect*: The data in the database is perceived by the user as tables.
 - *Integrity aspect*: Those tables has to satisfy certain integrity constraints.

- *Manipulative aspect*: The operators available to the user for manipulating those tables, for purposes of data retrieval, these operators derive tables form tables, the most important operators are ‘SELECT’, ‘PROJECT’ and JOIN.

Relation Schema

A relation schema ‘ R ’ denoted by $R(A_1, A_2, \dots, A_n)$ is made up of a relation name R and a list of attributes A_1, A_2, \dots, A_n . Each Attribute A_i is the name of role played by some domain D in the relation schema R . D is called the domain of A_i and is denoted by $\text{dom}(A_i)$.

A relation schema is used to describe a relation and R is called the name of this relation. The degree of a relation is the number of attributes ‘ n ’ of its relation schema.

Example:

A relation schema of degree ‘7’, which describes an employee is given below:

EMPLOYEE (Name, EId, HomePhone, Address, Office phone, Age, Salary)

Using the data type of each attribute, the definition is written as:

EMPLOYEE (Name: String, EId: INT, Homephone: INT, Address: String, OfficePhone: String, Age: Real, Salary: INT)

For this relation schema, EMPLOYEE is the name of the relation, which has ‘7’ attributes.

2. A relation ‘ r ’ of the relation schema $R(A_1, A_2 \dots A_n)$, also denoted by $r(R)$, is a set of n -tuples. $r = \{t_1, t_2 \dots t_m\}$. Each n -tuple ‘ t ’ is an ordered list of n values $t = \langle V_1, V_2 \dots V_n \rangle$, where each value V_i , $1 \leq i \dots n$, is an element of $\text{dom}(A_i)$ or is a special null value.
3. The i th value in tuple t , which corresponds to the attribute A_i , is referred to as $t[A_i]$.
4. The following figure shows EMPLOYEE relation. Each type in a relation represents a particular employee entity. We display the relation as a table where each tuple is shown as a row and each attribute corresponds to a column header, indicating a role or interpretation of the values in that column. Null values represent attributes whose values are unknown or do not exist for some individual EMPLOYEE tuple.

Employee						
Name	EId	Home Phone	Address	Office Phone	Age	Salary
Mahesh	30-01	870-223366	Warangal	NULL	35	40 k
Ramesh	30-02	040-226633	Hyderabad	NULL	36	40 k
Suresh	30-03	040-663322	Kolkata	040-331123	35	42 k
Dinesh	30-04	040-772299	Bangalore	040-321643	36	40 k

Fig: The attributes and tuples of a relation EMPLOYEE. The earlier definition of a relation can be restated as follows:

A relation $r(R)$ is a mathematical relation of degree ‘ n ’ on the domains $\text{dom}(A_1), \dots, \text{dom}(A_n)$, which is a subset of the Cartesian product of the domains that define R :

$$r(R) \subseteq (\text{dom}(A_1) \times \text{dom}(A_2) \dots \times \text{dom}(A_n))$$

Characteristics of Relations

There are certain characteristics that make a relation different from a file or a table.

Ordering of tuples in a relation: A relation is defined as a set of tuples. Tuples in a relation do not have any particular order. In a file, records are stored on disk in order. This ordering indicates first, second, i^{th} , and last records in the file. Similarly, when we display a relation as a table, the rows are displayed in a certain order.

Tuple ordering is not part of a relation definition, because a relation attempts to represent facts at a logical or abstract level. Many logical orders can be specified on a relation. Tuples in the EMPLOYEE relation could be logically ordered by values of Name or by EID or by Age or by some other attribute. When a relation is implemented as a file or displayed as a table, a particular ordering may be specified on the records of the file or the rows of the table.

NULL IN TUPLES

Each value in a tuple is an atomic value; that is, it is not divisible into components. Hence, composite and multi-valued attributes are not allowed. This model is sometimes called the *flat relational model*. Multivalued attributes must be represented by separate relations, and composite attributes are represented only by their simple component attributes in the basic relational model.

NULLS are used to represent the values of attributes that may be unknown or may not apply to tuple. For example, some student tuples have null in their office phones because they do not have an office. In this case, the meaning behind NULL is not applicable. If a student has a NULL for home phone, it means either he/she does not have a home phone or he/she has one but we do not know it, in this case the meaning of NULL is ‘Unknown’.

RELATIONAL MODEL CONSTRAINTS

In a relational database, there will be many relations, and the tuples in those relations are related in various ways. There are many restrictions or constraints on the actual values in a database state.

Constraints on database can generally be divided into three main categories as follows:

1. Constraints that are inherent in the data model, we call them *inherent model-based constraints*.
2. Constraints that can be directly expressed in the schemas of the data model, by specifying them in the DDL (Data Definition Language). We call these *schema-based constraints*.

3. Constraints that cannot be directly expressed in the schemas of the data model, and they must be expressed and enforced by the application programs are *application-based constraints*.

Inherent Constraint

The constraint that a relation cannot have duplicate tuples is an *inherent constraint*. Another important category of constraints is data dependencies, which include ‘functional dependencies’ and ‘multivalued dependencies’. They are used mainly for testing the ‘goodness’ of the design of a relational database and are utilized in a process called normalization.

Schema-based Constraints

1. Domain constraints
2. Key constraints
3. Constraints on nulls (Not null constraint)
4. Entity integrity constraints
5. Referential integrity constraints
6. Unique constraint
7. Check constraint

Domain constraints

Domain constraints specify that within each tuple, the value of each attribute ‘ X ’ must be an atomic value from the domain $\text{dom}(X)$.

The data types associated with domains include standard numeric data types for integers

1. Short integer
2. Integer
3. Long integer
4. Real numbers
 - Float
 - Double-precision float
5. Characters
6. Booleans
7. Fixed-length strings
8. Variable-length strings
9. Date, time, time stamp
10. Money data types

Key constraints

A relation is a set of tuples. All elements of a set are distinct; hence, all tuples in a relation must also be distinct. This means no two tuples can have the same combination of values for all their attributes.

1. There are other subsets of attributes of a relation schema R with the property that no two tuples in any relation state ‘ r ’ of R should have the same combination of values of these attributes.

Suppose that we denote one subset of attributes by ‘SK’, then for two distinct tuples t_1 and t_2 in a relation state ‘ r ’ of R , we have the following constraint:

$$t_1[\text{SK}] = t_2[\text{SK}]$$

- Any such set of attributes SK is called a *super key* of the relation schema R .

SK specifies a uniqueness constraint that no two distinct tuples in any state r or R can have the same value for SK.

- Every relation has at least one default super key, the set of all its attributes. A key, ' K ' of a relation schema R is a super key of R with the additional property that removing any attributes ' X ' from K leaves a set of attributes K' that is not a super key of R any more.

A key satisfies the following two constraints:

- Two distinct tuples in any state of the relation cannot have identical values for all the attributes in the key.
- A super key from which we cannot remove any attributes and still have the uniqueness constraints mentioned in above condition is known as a *minimal super key*.

The first condition applies to both keys and super keys. The second condition is required only for keys.

Example: Consider the employee relation in Page no. 9. The attribute set {EId} is a key of employee because no two employee tuples can have the same value for EId.

Any set of attributes that include EId will form a super key.

- {EId, Homephone, Name}
- {EId, Age, Salary}
- {Name, EId, Address}

However, the super key {EId, Name, Age} is not a key of EMPLOYEE, because removing Name or age or both from the set leaves us with a super key. Any super key formed from a single attributes is also a key. A key with multiple attributes must require all its attributes to have the uniqueness property.

A relation schema may have more than one key. In that case, each of the keys is called a *candidate key*.

Example: Employee relation has three candidate keys. {Name, EId, Homephone}

One of the candidate keys is chosen as primary key of the relation.

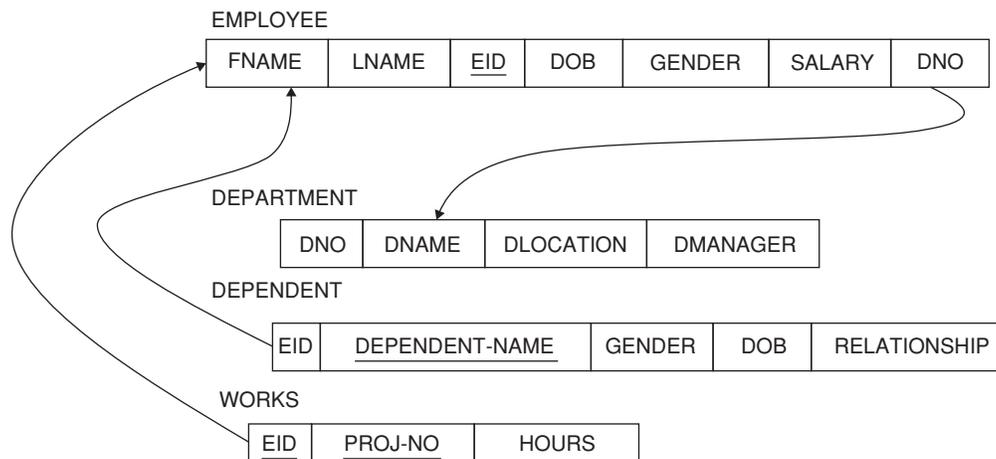
Another constraint on attributes specifies whether null values are permitted in tuples or not. If we want some tuples to have a valid (or) non-null value, we need to use NOT NULL constraint on that attribute.

Referential and entity integrity constraint

The entity Integrity constraint states that no primary key value can be null. If we have NULL values in the primary key column, we cannot identify some tuples in a relation.

- Key constraints and entity Integrity constraints are specified on individual relations
- Referential integrity constraint is specified between relations and used to maintain the consistency among tuples in the two relations.
- Referential Integrity constraints states that a tuple in one relation that refers to another relation must refer to an existing tuple in that relation.
- To understand the concept of Referential Integrity, first we have to understand the concept of FOREIGN KEY.
- Suppose we have two relations R_1 and R_2 . A set of attributes FK in relation schema R_1 is a foreign key of R_1 that references relation R_2 if it satisfies the following two rules:
 - The attributes in FK have the same domains as the primary key attributes PK of R_2 , FK will have to refer to PK.
 - A value of FK in a tuple t_1 of the current state $r_1(R_1)$ either occurs as a value of PK for some tuple t_2 in the current state $r_2(R_2)$ or is null. We have $t_1[FK] = t_2[PK]$, and we say that the tuple t_1 references to the tuple t_2 . In this definition, R_1 is called the *referencing relation* and R_2 is the *referenced relation*.
- A *foreign key* can refer to its own relation. We can diagrammatically display referential integrity constraints by drawing a directed Arc from each foreign key to the relation it references. The arrow head may point to the primary key of the referenced relation.

Example:



7. *Referential integrity rule*: The database must not contain any unmatched foreign key values.

If 'B' References 'A', then A must exist.

NOT NULL constraint

1. NOT NULL constraint restricts a column from having a NULL value. NOT NULL constraint can be applied to any column in a table.
2. We cannot give NULL values under that column
3. NOT NULL Constraint enforces a column to contain a proper value.
4. This constraint cannot be defined at table level.

Example: CREATE TABLE student(RNo:INT Name:varchar(70)NOTNULL age:INT)
Suppose a row is inserted into the following table,

Insert into student values <11, NULL, 20>

In the schema, we enforced NOT NULL constraint on Name column, means Name cannot have NULL value, when the above insert command is executed, the system gives, NOT NULL constraint violation.

UNIQUE constraint

The column on which UNIQUE constraint is enforced should not have any duplicate values.

1. UNIQUE constraint can be enforced on any column except the primary key column.
2. By default primary key column will not accept any duplicate values that are handled by key constraint.
3. UNIQUE constraint can be applied at column level or table level.

Example: CREATE TABLE student (RNo:INT Name:varchar(60) Grade:CHAR(1))
Assume that the table contains following tuples

Student		
R no.	Name	Grade
11	Sita	B
12	Anu	A
13	Bala	A

Suppose the following tuple is inserted into the student table.

1. Insert into student values <14, 'Anu', 'B'>
2. UNIQUE constraint is enforced on Name column, in the student table we have 'Anu', and again the new tuple contains name 'Anu', this Insert command violates the UNIQUE constraint.

CHECK constraint

This constraint is used to restrict a value of a column between a range.

1. When a value is inserted into particular column, before storing that, a check will be performed to see whether the values lie within the specified range.
2. If the value entered is out of range, it will not accept and violation happens.

3. It is like checking a condition before saving data into a column.

Example: Create table student (RNo:INT CHECK (Rno>0)

Name:varchar(60)

Dept:varchar(4))

Suppose the following tuple is inserted into student table.

1. Insert into student values <-4, 'Bhanu', 'CS'>
2. CHECK constraint is enforced on RNo column, the RNo should be greater than '0', but '-4' is given.
3. CHECK constraint is violated.

Creating table from a table: A view is called a *derived table* (or) *virtual table*, because the data stored in views is taken from already existing tables.

A view can also be defined as a logical subset of data from one or more tables.

Syntax:

```
CREATE view view-name As
SELECT column-names
FROM table-name
WHERE condition
```

Example: Consider the following table "sales".

Sales			
Order-Id	Order Name	Previous Balance	Customer
21	Order 3	3000	Ana
22	Order 4	1000	Adam
23	Order 5	3000	Brat
24	Order 6	2000	John
25	Order 7	2000	Ana
26	Order 8	4000	Ana

Query to create a view:

```
CREATE view sales-view As
```

```
SELECT *
FROM Sales
WHERE customer = 'Ana'
```

1. The data fetched from select statement will be stored in an object called 'sales-view'.
2. To display the contents stored in view, execute the following statement.

```
SELECT *
FROM Sales-view
```

Removal of specific rows:

Consider the following SQL query:

```
Delete *
FROM Sales
```

The above query will delete all the tuples from sales.

To remove specific rows, we have to specify the condition in WHERE clause.

Consider the table "sales" given in the above example.

Remove the rows from sales table whose previous balance is 3000.

SQL query:

```
Delete      *
FROM        Sales
WHERE       Previous-balance = 3000
```

Output:

Order Id	Order Name	Previous Balance	Customer
22	Order 4	1000	Adam
24	Order 6	2000	John
25	Order 7	2000	Ana
26	Order 8	4000	Ana

Referential actions: Referential Integrity can be violated if the value of any foreign key refers to a tuple that does not exist in the referenced relation.

When certain violations occur, we need to perform some alternate action. Those actions are as follows:

1. ON DELETE CASCADE
2. ON UPDATE CASCADE
3. ON DELETE SET NULL
4. ON DELETE SET DEFAULT

Example: Consider the given database:

SUPPLIERS:

Supplier Number	Supplier Name	Status	City
SN1	Suma	30	Hyderabad
SN2	Hari	20	Chennai
SN3	Anu	10	Hyderabad
SN4	Mahesh	20	Bombay
SN5	Kamal	30	Delhi

PARTS:

Part number	Part name	Colour	Weight	City
PN1	X	Red	13.0	Chennai
PN2	Y	Green	13.5	Bombay
PN3	X	Yellow	13.2	Hyderabad
PN4	Y	Green	14.1	Calcutta
PN5	Z	Red	14.3	Hyderabad
PN6	Z	Blue	14.2	Bombay

PROJECT:

Project Number	Project Name	City
PJ1	Display	Chennai
PJ2	OCR	Bombay
PJ3	RAID	Chennai
PJ4	SORTER	Hyderabad
PJ5	EDS	Chennai
PJ6	Tape	Bombay
PJ7	Console	Hyderabad

SHIPMENTS:

Supplier Number	Part Number	Project Number	Quantity
SN1	PN1	PJ1	300
SN1	PN1	PJ4	400
SN2	PN3	PJ1	350
SN2	PN3	PJ2	450
SN2	PN3	PJ3	640
SN2	PN3	PJ4	320
SN2	PN3	PJ5	330
SN2	PN3	PJ6	520
SN2	PN3	PJ7	480
SN2	PN5	PJ2	460
SN3	PN3	PJ1	440
SN3	PN4	PJ2	410
SN4	PN6	PJ3	310
SN4	PN6	PJ7	320
SN5	PN2	PJ2	340
SN5	PN2	PJ4	350
SN5	PN5	PJ5	360
SN5	PN5	PJ7	370
SN5	PN6	PJ2	380
SN5	PN1	PJ4	420
SN5	PN3	PJ4	440
SN5	PN4	PJ4	450
SN5	PN5	PJ4	400
SN5	PN6	PJ4	410

Consider the following statement:

```
DELETE FROM SUPPLIER
WHERE SUPPLIER – NUMBER = ‘SN1’
```

It deletes the supplier tuple for supplier ‘SN1’. The database has some other tables which have ‘SN1’ tuple (Shipments table). The application does not delete those suppliers, then it will find a violation, and an exception will be raised.

An alternate approach is possible, one that might be preferable in some cases, and that is for the system to perform an appropriate ‘compensating action’ that will guarantee that the overall result does still satisfy the constraint. In the example, the compensating action would be for the system to delete the shipments for supplier SN1 “automatically”.

We can achieve this effect by extending the foreign key as indicated below:

```
CREATE TABLE SHIPMENT{.....}.....
FOREIGN KEY {SUPPLIER – NUMBER}
REFERENCES
SUPPLIER ON DELETE CASCADE
```

The specification ON DELETE CASCADE defines a delete rule for this particular foreign key, and the specification CASCADE is the referential action for that delete rule. The meaning of these specifications is that a DELETE operation on the suppliers relvar will ‘Cascade’ to delete matching tuples (if any) in the shipments relvar as well.

Same procedure is applied for all the referential actions.

TRIGGERS

Triggers are precompiled procedures that are stored along with the database and invoked automatically whenever some specified event occurs.

Suppose we have a view called HYDERABAD - SUPPLIER defined as follows:

```
CREATE VIEW HYDERABAD-SUPPLIER
AS SELECT SUPPLIER - NUMBER, SUPPLIER-
NAME, STATUS
FROM SUPPLIER
WHERE CITY = ‘HYDERABAD’,
```

Normally, if the user tries to insert a row into this view, SQL will actually insert a row into the underlying base table SUPPLIERS with CITY value whatever the default is for the CITY column. Assuming that default is not Hyderabad, the net effect is that the new row will not appear in the view; therefore, let us create a triggered procedure as follows:

```
CREATE TRIGGER HYDERABAD -
SUPPLIER - INSERT
INSTEAD OF INSERT ON HYDERABAD
- SUPPLIER
REFERENCING NEW ROW AS R
FOR EACH ROW
INSERT INTO SUPPLIERS (SUPPLIER -
NUMBER, SUPPLIER - NAME, STATUS, CITY)
VALUES (R. SUPPLIER - NUMBER, R.
SUPPLIER - NAME, R. STATUS, ‘HYDERABAD’);
```

Inserting a row into the view will now cause a row to be inserted into the underlying base table with CITY value equal to Hyderabad inserted of the default value.

In general, CREATE TRIGGER specifies, among other things, an event, a condition, and an action.

The event is an operation on the database (“INSERT ON HYDERABAD - SUPPLIER” in the example)

1. The “condition” is a Boolean expression that has to evaluate to TRUE in order for the action to be executed.
2. The ‘action’ is the triggered procedure (“INSERT INTO SUPPLIERS ...”)
3. The event and condition together are sometimes called the *triggering event*. The combination of all three (event, condition, and action) is usually called a trigger.
4. Possible events include INSERT, DELETE, UPDATE, reaching end-of-transaction (COMMIT) reaching a specified time of day, exceeding a specified elapsed time, violating a specified constraint, etc.
5. A database that has associated triggers is sometimes called an active database.

Base Table Constraints

SQL-base table constraints are specified on either CREATE TABLE or ALTER TABLE. Each such constraint is a candidate key constraint, a foreign key constraint, or a CHECK constraint.

Candidate keys: An SQL candidate key definition takes one of the following two forms:

```
PRIMARY KEY (< column name comma list>)
UNIQUE (< column name comma list>)
```

The following example illustrates base table constraints of all three kinds:

```
CREATE TABLE SHIPMENTS
(SUPPLIER-NUMBER. SUPPLIER-NUMBER
NOT NULL, PART - NUMBER PART-NUMBER
NOT NULL, QUANTITY NOT NULL
PRIMARY KEY (SUPPLIER - NUMBER,
PART NUMBER)
FOREIGN KEY (SUPPLIER-NUMBER)
REFERENCES SUPPLIERS
ON DELETE CASCADE
ON UPDATE CASCADE,
FOREIGN KEY (PART-NUMBER)
REFERENCES PARTS
ON DELETE CASCADE
ON UPDATE CASCADE
CHECK(QUANTITY ≤ QUANTITY (0) AND
QUANTITY ≤ QUANTITY (1000));
```

A check constraint of the form CHECK (< column name > IS NOT NULL) can be replaced by a simple NOT NULL specification in the definition of the column.

EXERCISES

Practice Problems I

Directions for questions 1 to 20: Select the correct alternative from the given choices.

1. Consider the following two tables T_1 and T_2 . Show the output for the following operations:

Table T_1

P	Q	R
11	a	6
16	b	9
26	a	7

Table T_2

A	B	C
11	b	7
26	c	4
11	b	6

What is the number of tuples present in the result of given algebraic expressions?

- (i) $T_1 \bowtie_{T_1.P = T_2.A} T_2$
- (A) 2 (B) 3
(C) 4 (D) 5
- (ii) $T_1 \bowtie_{T_1.Q = T_2.B} T_2$
- (A) 2 (B) 3
(C) 4 (D) 5
- (iii) $T_1 \bowtie_{(T_1.P = T_2.A \text{ AND } T_1.R = T_2.C)} T_2$
- (A) 1 (B) 2
(C) 3 (D) 4
2. Suppose $R_1(A, B)$ and $R_2(C, D)$ are two relation schemas. Let R_1 and R_2 be the corresponding relation instances. B is a foreign key that refers to C in R_2 . If data in R_1 and R_2 satisfy referential integrity constraints, which of the following is true?
- (A) $\pi_B(R_1) - \pi_C(R_2) = \phi$
(B) $\pi_C(R_2) - \pi_B(R_1) = \phi$
(C) $\pi_B(R_1) - \pi_C(R_2) \neq \phi$
(D) Both A and B
3. Consider the following relations:
 A, B and C

A		
Id	Name	Age
12	Arun	60
15	Shreya	24
99	Rohit	11

B		
Id	Name	Age
15	Shreya	24
25	Hari	40
98	Rohit	20
99	Rohit	11

C		
Id	Phone	Area
10	2200	02
99	2100	01

How many tuples does the result of the following relational algebra expression contain? Assume that the scheme of $(A \cup B)$ is the same as that of A .

$$(A \cup B) \bowtie_{A.Id > 40 \vee C.Id < 15} C$$

- (A) 6 (B) 7
(C) 8 (D) 9

4. Consider the relations A, B and C given in Question 3. How many tuples does the result of the following SQL query contain?

```
SELECT  A.Id
FROM    A
WHERE   A.Age >
ALL    (SELECT B.Age
FROM    B
WHERE   B.Name = 'Arun')?
```

(A) 0 (B) 1
(C) 2 (D) 3

5. Consider a database table T containing two columns X and Y each of type integer. After the creation of the table, one record ($X = 1, Y = 1$) is inserted in the table. Let MX and MY denote the respective maximum value of X and Y among all records in the table at any point in time. Using MX and MY , new records are inserted in the table 128 times with X and Y values being $MX + 1, 2 * MY + 1$, respectively. It may be noted that each time after the insertion, values of MX and MY change. What will be the output of the following SQL query after the steps mentioned above are carried out?

```
SELECT Y FROM T WHERE X = 7;?
```

(A) 15 (B) 31
(C) 63 (D) 127

6. Database table by name loan records is given below:

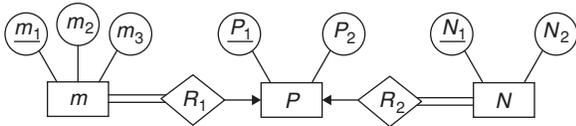
Borrower	Bank Manager	Loan Amount
Ramesh	Sunderajan	10000
Suresh	Ramgopal	5000
Mahesh	Sunderajan	7000

What is the output of the following SQL Query
SELECT count (*)

```
FROM((SELECT Borrower, Bank-manager
FROM Loan-Records)AS S
NATURAL JOIN
(SELECT Bank-manager, Loan-Amount
FROM Loan-Records) AS T);
```

(A) 3 (B) 4
(C) 5 (D) 6

7. Consider the following ER diagram:



What is the minimum number of tables needed to represent M, N, P, R_1, R_2 ?

- (A) 2 (B) 3
(C) 4 (D) 5
8. Let E_1 and E_2 be two entities in an ER diagram with simple single-valued attributes. R_1 and R_2 are two relationships between E_1 and E_2 , where R_1 is one-to-many and R_2 is many-to-many. R_1 and R_2 do not have any attributes of their own. What is the minimum number of tables required to represent this situation in the relational model?
- (A) 2 (B) 3
(C) 4 (D) 5
9. The following table has two attributes A and C where A is the primary key and C is the foreign key referencing A with ON DELETE CASCADE.

A	C
2	4
3	4
4	3
5	2
7	2
9	5
6	4

What is the set of all tuples that must be additionally deleted to preserve referential integrity when the tuple (2, 4) is deleted?

- (A) (5, 2), (7, 2) (B) (5, 2), (7, 2), (9, 5)
(C) (5, 2), (9, 5) (D) (2, 4), (7, 2)

10. Consider the following SQL query

```
SELECT DISTINCT a1, a2, a3 ... an
FROM R1, R2 ... Rm
WHERE P
```

For any arbitrary predicate P , this query is equivalent to which relational algebra expression?

- (A) $\Pi_{a_1, a_2, \dots, a_n} (\sigma_P R_1 \times R_2 \times \dots \times R_m)$
(B) $\sigma_{a_1, a_2, \dots, a_n} (\sigma_P R_1 \times R_2 \times \dots \times R_m)$
(C) $\sigma_{a_1, a_2, \dots, a_n} (\prod_P R_1 \times R_2 \times \dots \times R_m)$
(D) $\Pi_{R_1, R_2, \dots, R_m} (\sigma_P a_1 \times a_2 \times \dots \times a_n)$

11. Consider the following relation schema pertaining to a student's database:

Student (Rollno, name, address)

Enroll (Rollno, courseno, coursenam)

Where the primary keys are shown underlined. The number of tuples in the student and Enroll tables are 120 and 6, respectively. What are the maximum and minimum number of tuples that can be present in

(student * Enroll)

Where '*' denotes natural join?

- (A) 6, 6 (B) 6, 120
(C) 120, 6 (D) 120, 120
12. A relational schema for a train reservation database is given below

Table 4 Passenger

Pid	P Name	Age
0	'Sachin'	65
1	'Rahul'	66
	'Sourav'	67
3	'Anil'	69

Table 5 Reservation

Pid	Class	Tid
0	AC	8200
1	AC	8201
2	SC	8201
5	AC	8203
1	SC	8204
3	AC	8202

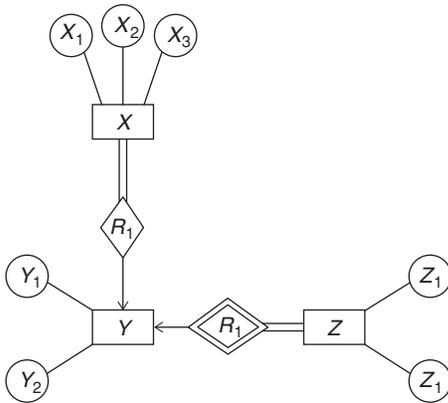
What pid's are returned by the following SQL query for the above instance of the tables?

```
SELECT pid
FROM Reservation
WHERE class = 'AC' AND
      EXISTS (SELECT *
              FROM passenger
              WHERE age > 65 AND
                    Passenger.pid = Reservation.pid)
```

- (A) 0, 1
- (B) 1, 3
- (C) 1, 5
- (D) 0, 3

13. Given {customer} is a candidate key, [customer name, customer street] is another candidate key then
- (A) {customer id, customer name} is also a candidate key.
 - (B) {customer id, customer street} is also a candidate key.
 - (C) {customer id, customer name, customer street} is also a candidate key.
 - (D) None

Common data for questions 14 and 15: Consider the following diagram,



14. The minimum number of tables needed to represent X, Y, Z, R₁, R₂ is
- (A) 2
 - (B) 3
 - (C) 4
 - (D) 5
15. Which of the following is a correct attribute set for one of the tables for the correct answer to the above questions?
- (A) {X₁, X₂, X₃, Y₁}
 - (B) {X₁, Y₁, Z₁, Z₂}
 - (C) {X₁, Y₁, Z₁}
 - (D) {M₁, Y₁}
16. UPDATE account SET
DA = basic * .2,
GROSS = basic * 1.3,
Where basic > 2000;

- (A) The above query displays DA and gross for all those employees whose basic is ≥ 2000
- (B) The above query displays DA and gross for all employees whose basic is less than 2000
- (C) The above query displays DA as well as gross for all those employees whose basic is >2000
- (D) Above all

17. Which of the following query transformations is correct? R₁ and R₂ are relations C₁, C₂ are selection conditions and A₁ and A₂ are attributes of R₁
- (A) $\sigma_{C_1}(\sigma_{C_1}(R_1)) \rightarrow \sigma_{C_2}(\sigma_{C_2}(R_1))$
 - (B) $\sigma_{C_1}(\sigma_{A_1}(R_1)) \rightarrow \sigma_{A_1}(\sigma_{C_1}(R_1))$
 - (C) $\pi_{A_2}(\pi_{A_1}(R_1)) \rightarrow \pi_{A_1}(\pi_{A_2}(R_1))$
 - (D) All the above

18. Consider the following query select distinct a₁, a₂, ... a_n from r₁, r₂ ... r_m where P for an arbitrary predicate P, this query is equivalent to which of the following relational algebra expressions:
- (A) $\pi_{a_1 \dots a_n} \sigma_P(r_1 \times r_2 \times \dots \times r_m)$
 - (B) $\pi_{a_1 \dots a_n} \sigma_P(r_1 \times r_2 \times r_3 \times \dots \times r_m)$
 - (C) $\sigma_{a_1 \dots a_n} \pi_P(r_1 \times r_2 \times \dots \times r_m)$
 - (D) $\sigma_{a_1 \dots a_n} \pi_P(r_1 \times r_2 \times \dots \times r_m)$

19. The relational algebra expression equivalent to the following tuple calculus expression
 $\{a\} a \in r \wedge (a[A] = 10 \wedge a[B] = 20)$ is
- (A) $\sigma_{(A=10 \wedge B=20)} r$
 - (B) $\sigma_{(A=10)}(r) \cup \sigma_{(B=20)}(r)$
 - (C) $\sigma_{(A=10)}(r) \cap \sigma_{(B=20)}(r)$
 - (D) $\sigma_{(A=10)}(r) - \sigma_{(B=20)}(r)$
20. Which of the following is/are wrong?
- (A) An SQL query automatically eliminates duplicates.
 - (B) An SQL query will not work if there are no indexes on the relations.
 - (C) SQL permits attribute names to be repeated in the same relation
 - (D) All the above

Practice Problems 2

Directions for questions 1 to 20: Select the correct alternative from the given choices.

1. If ABCDE are the attributes of a table and ABCD is a super key and ABC is also super key then
- (A) A B C must be candidate key
 - (B) A B C cannot be super key
 - (C) A B C cannot be candidate key
 - (D) A B C may be candidate key

2. The example of derived attribute is
- (A) Name if age is given as other attribute
 - (B) Age if date_of_birth is given as other attribute
 - (C) Both (A) and (B)
 - (D) None
3. The weak entity set is represented by
- (A) box
 - (B) ellipse
 - (C) diamond
 - (D) double outlined box

4. In entity relationship diagram double lines indicate
 - (A) Cardinality
 - (B) Relationship
 - (C) Partial participation
 - (D) Total participation
5. An edge between an entity set and a binary relationship set can have an associated minimum and maximum cardinality, shown in the form $1 \dots h$ where 1 is the minimum and h is the maximum cardinality A minimum value 1 indicates:
 - (A) total participation
 - (B) partial participation
 - (C) double participation
 - (D) no participation
6. Let R be a relation schema. If we say that a subset k of R is a super key for R , we are restricting R , we are restricting consideration to relations $r(R)$ in which no two distinct tuples have the same value on all attributes in K . That is if t_1 and t_2 are in r and $t_1 \uparrow t_2$
 - (A) $t_1[k] = 2t_2[k]$
 - (B) $t_2[k] = 2t_1[k]$
 - (C) $t_1[k] = t_2[k]$
 - (D) $t_1[k] \uparrow t_2[k]$
7. Which one is correct?
 - (A) Primary key \subset Super key \subset Candidate key
 - (B) Candidate key \subset Super key \subset Primary key
 - (C) Primary key \subset Candidate key \subset Super key
 - (D) Super key \subset Primary key \subset Candidate key
8. If we have relations $r_1(R_1)$ and $r_2(R_2)$, then $r_1 \cap r_2$ is a relation whose schema is the
 - (A) concatenation
 - (B) union
 - (C) intersection
 - (D) None

9. Match the following:

I	Empid	1	Multivalued
II	Name	2	Derived
III	Age	3	Composite
IV	Contact No.	4	Simple

- (A) I – 4, II – 3, III – 2, IV – 1
 - (B) I – 3, II – 2, III – 4, IV – 1
 - (C) I – 2, II – 1, III – 4, IV – 3
 - (D) I – 1, II – 3, III – 2, IV – 4
10. Match the following:
- | | | | |
|-----|----------------------|---|-----------------------|
| I | Double-lined ellipse | 1 | Multivalued attribute |
| II | Double line | 2 | Total participation |
| III | Double-lined box | 3 | Weak entity set |
| IV | Dashed ellipse | 4 | Derived attribute |
- (A) I – 1, II – 2, III – 3, IV – 4
 - (B) I – 2, II – 3, III – 4, IV – 1
 - (C) I – 3, II – 4, III – 2, IV – 2
 - (D) I – 4, II – 3, III – 2, IV – 1
11. The natural join is a
- (A) binary operation that allows us to combine certain selections and a Cartesian product into one operation
 - (B) unary operations that allows only Cartesian product

- (C) query which involves a Cartesian product and a projection
- (D) None

12. The number of entities participating in the relationship is known as
 - (A) maximum cardinality
 - (B) composite identifiers
 - (C) degree
 - (D) None
13. A minimum cardinality of 0 specifies
 - (A) non-participation
 - (B) partial participation
 - (C) total participation
 - (D) zero participation
14. What is not true about weak entity?
 - (A) They do not have key attributes.
 - (B) They are the examples of existence dependency.
 - (C) Every existence dependency results in a weak entity
 - (D) Weak entity will have always discriminator attributes
15. Which one is the fundamental operation in the relational algebra?
 - (A) Natural join
 - (B) Division
 - (C) Set intersection
 - (D) Cartesian product
16. For the given tables

A		B
X	Y	Y
a_1	b_1	b_1
a_2	b_1	b_2
a_1	b_2	
a_2	b_2	

$A \div B$ will return

- (A) a_1, a_2
 - (B) a_1
 - (C) a_2
 - (D) None
17. The number of tuples selected in the above answer is
- (A) 2
 - (B) 1
 - (C) 0
 - (D) 4

Common data for questions 18 and 19: Consider the following schema of a relational database.

Emp (empno, name, add)

Project (Pno, Prame)

Work on (empno, Pno)

Part (partno, Pname, qty, size)

Use (empno, pno, partno, no)

18. ((name(emp) ((name(emp) \bowtie workon) displays
- (i) The names of the employees who are not working in any project
 - (ii) The names of the employees who were working in every project.
- (A) Only (i)
 - (B) Only (ii)
 - (C) Both (A) and (B)
 - (D) None

19. List the partno and names of the parts used in both the projects DBMS & MIS:
- (A) $\sigma_{partno, pname, (part \bowtie (\pi_{partno} (\sigma_{pname = "DBMS"}(project) \bowtie use)) \cap \pi_{partno} (\sigma_{pname = "MIS"}(project) \bowtie use))}$
 - (B) $(partno, pname(part \bowtie ((partno ((pname = "DBMS"(project) \bowtie use) ((partno(\sigma_{pname = "MIS"}(project) \bowtie use))))$
 - (C) $\pi_{partno, pname}(part \bowtie (\sigma_{partno} (\sigma_{pname = "DBMS"}(project) \bowtie use) \cap (partno ((pname = "MIS"(project) \bowtie use))))$
 - (D) None

20. The following query shows. SELECT job-status, sum (basic – salary) AVG (basic – salary) from employees group by job – status.
- (i) It shows job status, sum, AVG of all data
 - (ii) It shows job status, Sum, AVG, with group by clause in use.
- (A) only (i)
 - (B) only (ii)
 - (C) both (A) and (B)
 - (D) None

PREVIOUS YEARS' QUESTIONS

1. Let E_1 and E_2 be two entities in an E/R diagram, with simple single-valued attributes. R_1 and R_2 are two relationships between E_1 and E_2 , where R_1 is one-to-many and R_2 is many-to-many. R_1 and R_2 do not have any attributes of their own. What is the minimum number of tables required to represent this situation in the relational model? [2005]
- (A) 2
 - (B) 3
 - (C) 4
 - (D) 5

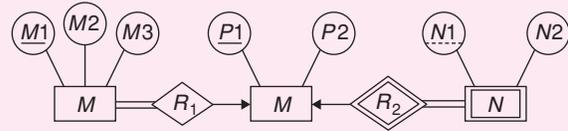
2. The following table has two attributes A and C where A is the primary key and C is the foreign key referenc- ing A with on-delete cascade.

A	C
2	4
3	4
4	3
5	2
7	2
9	5
6	4

The set of all tuples that must be additionally deleted to preserve referential integrity when the tuple (2, 4) is deleted is: [2005]

- (A) (3, 4) and (6, 4)
 - (B) (5, 2) and (7, 2)
 - (C) (5, 2), (7, 2) and (9, 5)
 - (D) (3, 4), (4, 3) and (6, 4)
3. Which of the following tuple relational calculus expression(s) is/are equivalent to $\forall t \in r (P(t))$?
- I. $\neg \exists t \in r (P(t))$
 - II. $\exists t \notin r (P(t))$
 - III. $\neg \exists t \in r (\neg P(t))$
 - IV. $\exists t \in r (\neg P(t))$ [2008]
- (A) I only
 - (B) II only
 - (C) III only
 - (D) III and IV only

Common data for questions 4 and 5: Consider the following ER diagram:



4. The minimum number of tables needed to represent M, N, P, R_1, R_2 is? [2008]
- (A) 2
 - (B) 3
 - (C) 4
 - (D) 5
5. Which of the following is a correct attribute set for one of the tables for the correct answer to the above question? [2008]
- (A) $\{M1, M2, M3, P1\}$
 - (B) $\{M1, P1, N1, N2\}$
 - (C) $\{M1, P1, N1\}$
 - (D) $\{M1, P1\}$
6. Consider a relational table with a single record for each registered student with the following attributes.
1. *Registration_Num*: Unique registration number of each registered student
 2. *UID*: Unique identity number, unique at the national level for each citizen
 3. *BankAccount_Num*: Unique account number at the bank. A student can have multiple accounts or joint accounts. This attribute stores the primary account number
 4. *Name*: Name of the student
 5. *Hostel_Room*: Room number of the hostel
- Which of the following options is incorrect? [2011]
- (A) *BankAccount_Num* is a candidate key
 - (B) *Registration_Num* can be a primary key
 - (C) *UID* is a candidate key if all students are from the same country
 - (D) If S is a super key such that $S \cap UID$ is NULL then $S \cup UID$ is also super key.
7. Given the basic ER and relational models, which of the following is incorrect? [2012]
- (A) An attribute of an entity can have more than one value
 - (B) An attribute of an entity can be composite
 - (C) In a row of a relational table, an attribute can have more than one value
 - (D) In a row of a relational table, an attribute can have exactly one value or a NULL value

8. An ER model of a database consists of entity types A and B. These are connected by a relationship R which does not have its own attribute. Under which one of the following conditions, can the relational table for R be merged with that of A? [2017]
- (A) Relationship R is one-to-many and the participation of A in R is total.
 (B) Relationship R is one-to-many and the participation of A in R is partial.
 (C) Relationship R is many-to-one and the participation of A in R is total.
 (D) Relationship R is many-to-one and the participation of A in R is partial.
9. In an Entity-Relationship (ER) model, suppose R is a many-to-one relationship from entity set E1 to entity

set E2. Assume that E1 and E2 participate totally in R and that the cardinality of E1 is greater than the cardinality of E2.

Which one of the following is true about R? [2018]

- (A) Every entity in E1 is associated with exactly one entity in E2.
 (B) Some entity in E1 is associated with more than one entity in E2.
 (C) Every entity in E2 is associated with exactly one entity in E1.
 (D) Every entity in E2 is associated with at most one entity in E1.

ANSWER KEYS

EXERCISES

Practice Problems 1

1. (i) B (ii) A (iii) A 2. A 3. B 4. D 5. D 6. D 7. B 8. B
 9. B 10. A 11. A 12. B 13. D 14. A 15. A 16. C 17. B 18. B
 19. C 20. D

Practice Problems 2

1. D 2. B 3. C 4. A 5. A 6. D 7. C 8. A 9. A 10. A
 11. A 12. C 13. C 14. A 15. D 16. A 17. A 18. A 19. C 20. B

Previous Years' Questions

1. B 2. C 3. D 4. B 5. A 6. A 7. C 8. C 9. A

Chapter 2

Structured Query Language

LEARNING OBJECTIVES

- Relational algebra
- Select operator
- Project operator
- Set operators
- Union compatible relations
- Union operation
- Aggregate operators
- Correlated nested queries
- Relational calculus
- Tuple relational calculus
- Tuple relational calculus
- DML
- Super key
- SQL commands

RELATIONAL ALGEBRA

- A set of operators (unary or binary) that take relation instances as arguments and return new relations.
- Gives a procedural method of specifying a retrieval query
- Forms the core component of a relational query engine
- SQL queries are internally translated into RA expressions
- Provides a framework for query optimization

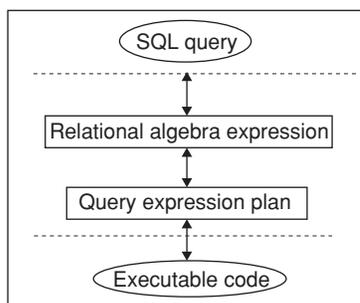
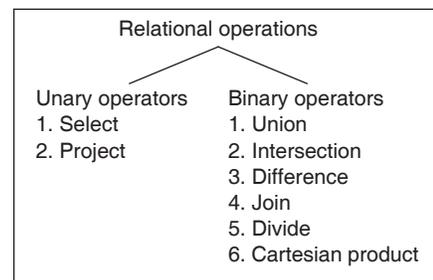


Figure 1 Role of relational algebra in DBMS:

Relational Operations

A collection of simple 'low-level' operations used to manipulate relations.

- It provides a procedural way to query a database.
- Input is one (or) more relations.
- Output is one relation.



Select Operator (σ)

Select operator is an unary operator. It can be used to select those tuples of a relation that satisfy a given condition.

Notation: $\sigma_{\theta}(r)$
 σ : Select operator (read as sigma)
 θ : Selection condition
 r : Relation name

Result is a relation with the same scheme as r consisting of the tuples in r that satisfy condition θ

Syntax: $\sigma_{\text{condition}}(\text{relation})$

Example:

Table 2.1 Person

Id	Name	Address	Hobby
112	John	12, SP Road	Stamp collection
113	John	12, SP Road	Coin collection
114	Mary	16, SP Road	Painting
115	Brat	18, GP Road	Stamp collection

$\sigma_{\text{Hobby} = \text{'stamp. Collection'}}(\text{person})$

The above given statement displays all tuples (or) records with hobby 'stamp collection'.

Output:

Id	Name	Address	Hobby
112	John	12, SP Road	Stamp collection
115	Brat	18, GP Road	Stamp collection

Selection condition can use following operators:

<, ≤, >, ≥, =, ≠

1. <attribute> operator <attribute>
2. <attribute> operator <constant>

Example: Salary ≥ 1000

3. <Condition> AND/OR <condition>

Example: (Experience > 3) AND (Age < 58)

4. NOT <condition>

Selection operation examples:

1. $\sigma_{\text{Id} > 112 \text{ OR } \text{Hobby} = \text{'paint'}}(\text{person})$

It displays the tuples whose ID > 112 or Hobby is paint

2. $\sigma_{\text{Id} > 112 \text{ AND } \text{Id} < 115}(\text{person})$

It displays tuples whose ID is greater than 112 and less than 115

3. $\sigma_{\text{NOT } (\text{hobby} = \text{'paint'})}(\text{person})$

It displays tuples whose hobby is not paint

4. $\sigma_{\text{Hobby} \neq \text{'paint'}}(\text{person})$

It displays tuples whose hobby is not paint, displays all tuples other than hobby paint.

Selection operator: Produces table containing subset of rows of argument table which satisfies condition.

Project Operator (π)

The project operator is unary operator. It can be used to keep only the required attributes of a relation instance and throw away others.

Notation: $\pi_{A_1, A_2, \dots, A_k}(r)$ Where A_1, A_2, \dots, A_k is a list L of desired attributes in the scheme of r .

Result = $\{(V_1, V_2, \dots, V_k) / V_i \in \text{DOM}(A_i), 1 \leq i \leq k \text{ and there is some tuple } t \text{ in } r, \text{ such that } t.A_1 = v_1, t.A_2 = v_2, \dots, t.A_k = v_k\}$
 $\pi_{\text{Attribute List}}(\text{Relation})$

Take table 2.1 as reference.

1. $\pi_{\text{Name}}(\text{person})$

Output:

Name
John
Mary
Bart

In the output table, John name has appeared once, project operation eliminated duplicates.

2. $\pi_{\text{Name, address}}(\text{person})$

Output:

Name	Address
John	12, SP Road
Mary	16, SP Road
Bart	18, GP Road

Expressions:

$\pi_{\text{id}^{\text{name}}}(\sigma_{\text{hobby} = \text{'stamp collection'} \text{ OR } \text{Hobby} = \text{'coin collection'}}(\text{person}))$

Output:

Id	Name
112	John
115	Bart

The above given relational algebra expression gives Ids, names of a person whose hobby is either stamp collection (or) coin collection.

Set Operators

Union (\cup), Intersection (\cap), set difference ($-$) are called *set operators*. Result of combining two relations with a set operator is a relation \Rightarrow all its elements must be tuples having same structure. Hence scope of set operations is limited to union compatible relations.

Union Compatible Relations

Two relations are union compatible if

1. Both have same number of columns
2. Names of attributes are same
3. Corresponding fields have same type
4. Attributes with the same name in both relations have same domain.
5. Union compatible relations can be combined using Union, Intersection, and set difference.

Example:

Consider the given tables.

Person (SSN, Name, Address, Hobby)

Professor (Id, Name, office, phone)

person and professor tables are not union compatible.

Union

The result of union will be a set consisting of all tuples appearing in either or both of the given relations. Relations cannot contain a mixture of different kinds of tuples, they must be 'tuple – homogeneous'. The union in the relational algebra is not the completely general mathematical union; rather, it is a special kind of union, in which we require the two input relations to be of the same type.

The general definition of relational union operator:

Given are two relations 'a' and 'b' of the same type. The union of those two relations, a union b, is a relation of the same type, with body consisting of all tuples 't' such that 't' appears in a or b or both.

* Union operation eliminates duplicates.

Here is a different but equivalent definition:

Given are two relations 'a' and 'b' of the same type. The union of those two relations, a union b, is a relation of the same type, with body consisting of all tuples t such that t is equal to (i.e., is a duplicate of) some tuple in a or b or both.

Union Operation (U)

When union operation is applied on two tables it gives all the tuples in both without Repetition.

Example:

Table 2 Result of union operation

	Roll no.	Name	Semester	Percentage
R	22	Arun	7	45%
	31	Bindu	6	55%
	58	Sita	5	35%
S	28	Suresh	4	65%
	31	Bindu	6	55%
	44	Pinky	4	75%
	58	Sita	5	35%
$R \cup S$	22	Arun	7	45%
	31	Bindu	6	55%
	58	Sita	5	35%
	44	Pinky	4	75%
	28	Sita	5	35%

Intersection

Like union, Intersection operator requires its operands to be of the same type. Given are two relations a and b of the same type, then, the intersection of those two relations, 'a' INTERSECT 'b', is a relation of the same type, with body consisting of all tuples t such that t appears in both 'a' and 'b'.

Intersection operation returns tuples which are common to both tables

Table 3 Result of intersection operation

$R \cap S =$	Roll no.	Name	Semester	Percentage
	31	Bindu	6	55%
	58	Sita	5	35%

Difference

Like union and intersection, the relational difference operator also requires its operands to be of the same type. Given are two relations 'a' and 'b' of the same type, Then, the difference between those two relations, 'a' MINUS 'b' (in that order), is a relation of the same type, with body consisting of all types t such that t appears in a and not b.

1. MINUS has a directionality to it, just as subtraction does in ordinary arithmetic (e.g., '6 - 3' and '3 - 6' are not the same thing)
2. Redundant duplicate rows are always eliminated from the result of UNION, INTERSECTION, EXCEPT operations.
3. SQL also provides the qualified variants UNION ALL, INTERSECT ALL and EXCEPT ALL, where duplicates are retained

Set difference operation returns the tuples in the first table which are not matching with the tuples of other table.

Table 4 Result of $R - S$

$R - S =$	Roll no.	Name	Semester	Percentage
	22	Arun	7	45%

Table 5 Result of $S - R$

$S - R =$	Roll no.	Name	Semester	Percentage
	28	Suresh	4	65%
	44	Pinky	4	75%

* $R - S \neq S - R$ (both are different)

Example:

A			
Supplier number	Supplier name	Status	City
SN1	MAHESH	40	HYDERABAD
SN3	SURESH	40	HYDERABAD

B			
Supplier number	Supplier number	Status	City
SN3	SURESH	40	HYDERABAD
SN4	RAMESH	30	CHENNAI

UNION ($A \cup B$)			
Supplier number	Supplier name	Status	City
SN1	MAHESH	40	HYDERABAD
SN3	SURESH	40	HYDERABAD
SN4	RAMESH	30	CHENNAI

INTERSECTION ($A \cap B$)			
Supplier number	Supplier name	Status	City
SN3	SURESH	40	HYDERABAD

DIFFERENCE (A – B)			
Supplier name	Supplier name	Status	City
SN1	MAHESH	40	HYDERABAD

DIFFERENCE (B – A)			
Supplier name	Supplier name	Status	City
SN4	RAMESH	30	CHENNAI

Cartesian Product

The Cartesian product of two sets is the set of all ordered pairs such that in each pair, the first element comes from the first set and the second element comes from second set.

The result consists of all the attributes from both of the two input headings. We define the Cartesian product of two relations ‘a’ and ‘b’, as

‘a’ times ‘b’, where a and b have no common attribute names (If we need to construct the Cartesian product of two relations that do have any such common attribute names, therefore, we must use the RENAME operator first to rename attributes appropriately).

The Cartesian product operation is also known as CROSS PRODUCT. This is also a binary set operation, but the relations on which it is applied need not to be union compatible. This operation is used to combine tuples from two relations in a combinational fashion.

Example:

	A	B
R	X ₁	X ₂
	X ₃	X ₄

	C	D
S	Y ₁	Y ₂
	Y ₃	Y ₄

$R \times S =$	A	B	C	D
	X ₁	X ₂	Y ₁	Y ₂
	X ₁	X ₂	Y ₃	Y ₄
	X ₃	X ₄	Y ₁	Y ₂
	X ₃	X ₄	Y ₃	Y ₄

Example: Transcript (StuId, coursecode, semester, grade)
Teaching (ProfId, coursecode, semester)

$$\pi_{\text{stuId, coursecode}}(\text{Transcript}) \times \pi_{\text{profId, coursecode}}(\text{Teaching})$$

The above expression returns

Table 6 Result of cross product

Stu Id	Course code	Prof Id	Course code
...

Aggregate Operators

SQL Supports the usual aggregate operators COUNT, SUM, AVG, MAX, MIN, EVERY and ANY, but there are a few SQL-specific points.

1. The argument can optionally be preceded by the keyword DISTINCT, for example SUM (DISTINCT column -name) to indicate that duplicates are to be eliminated before the aggregation is done. For MAX, MIN, EVERY and ANY, however, DISTINCT has no effect and should not be specified.
2. The operator COUNT (*), for this DISTINCT is not allowed, and is provided to count all rows in a table without any duplicate elimination.
3. Any NULLS in the argument column are eliminated before the aggregation is done, regardless of whether DISTINCT is specified, except in the case of COUNT (*), where nulls behave as if they were values.
4. After NULLS if any have been eliminated, if what is left is an empty set, COUNT returns zero. The other operators return NULL.

AVG, MIN, MAX, SUM, COUNT

These functions operate on the multiset of values of column of a relation and returns a value.

1. Find the average account balance at the Perryridge branch.

Solution: SELECT AVG (balance) FROM account WHERE branch.name = ‘perryridge’

2. Find the number of tuples in customer relation.

Solution: SELECT count (*) FROM customer

3. Find the number of depositors for each branch.

Solution: SELECT branch.name, COUNT (distinct customer-name) FROM depositor, account WHERE depositor. account-no = account account-no GROUPBY branch.name.

Nested Queries

Some queries require that existing values in the database be fetched and then used in a comparison condition.

Such queries can be conveniently formulated by using nested queries, which are complete SELECT – FROM –WHERE blocks within the WHERE clause of another query. The other query is called the *outer query*.

In-Comparison Operator

The comparison operator IN, which compares a value ‘v’ with a set of values ‘V’ and evaluates to TRUE if ‘v’ is one of the elements in V.

Example: Consider the given database scheme and the statement:

EMPLOYEE							
FNAME	INITIAL	LNAME	ENO	DOB	ADDRESS	SALARY	DNO

DEPARTMENT		
D NAME	DNO	MANAGER-NO

DEPARTMENT-LOCATIONS	
DNO	D-LOCATION

PROJECT			
PNAME	PNO	P-LOCATION	DNO

WORKS-ON		
ENO	PNO	HOURS

Example: Select distinct PNO from project where PNO IN (select PNO from project, department, employee where P.DNO = D.DNO AND MANAGER.NO = ENO AND LNAME = 'RAMYA')

The first query selects. The project numbers that have a 'Ramya' involved as manager, while the second selects the project numbers of projects that have a 'Ramya' involved as worker.

If a nested returns a single value, in such cases, it is permissible to use = instead of IN for the comparison operator.

In general, the nested query will return a table, Which is a set of multiset of tuples.

* SQL allows the use of tuples of values in comparisons by placing them within parentheses.

Example: SELECT DISTINCT ENO FROM WORKS - ON WHERE (PNO, HOURS) IN (SELECT PNO, HOURS FROM WORKS - ON WHERE ENO = 929).

This query will select the employee numbers of all employees who work on the same (PROJECT, HOURS) combination on some project a particular employee whose ENO = '929' works on. In this example, the IN operator compares the subtuple of values in parentheses (PNO, HOURS) for each tuple in works on with the set of union-compatible tuples produced by the nested query.

Correlated nested queries: Nested queries can be evaluated by executing the sub query (or) Inner query once and substituting the resulting value (or) values into the WHERE clause of the outer query.

1. In co-related nested queries, the inner query depends on the outer query for its value.
2. Sub-query is executed repeatedly, once for each row that is selected by the outer query.
3. A correlated subquery is a sub query that contains a reference to a table that also appears in the outer query.

Example: Consider the following correlated nested query:

```
SELECT *
FROM table1
WHERE col1 ≥ ALL
(SELECT col1 FROM table2 WHERE table2. col2 =
table1. col2)
```

1. The subquery contains reference to a column of table1, even though the sub-queries FROM clause does not mention a table table1
2. SQL has to check outside the sub-query and find Table 1 in the outer query
3. Suppose that Table 1 contains a row where col1 = 3 and col2 = 4 and Table 2 contains a row where col1 = 5 and col2 = 4
4. The expression

WHERE col1 ≥ All (SELECT col1 FROM table2

3 ≥ 5 (false)

(WHERE condition TRUE) Table1. col2 = table2. col2 4 = 4

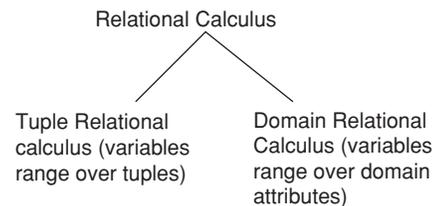
So the expression as a whole is FALSE.

5. It is evaluated from outside to inside

Relational Calculus

Relational calculus can define the information to be retrieved

1. In this, there is no specific series of operations.
2. Relational algebra defines the sequence of operations.
3. Relational calculus is closer to how users would formulate queries, in terms of information requirements, rather than in terms of operations.



4. Relational calculus is based on predicate logic, gives the usual quantifiers to construct complex queries.

Tuple Relational Calculus

Example: Employee

E Id	F Name	L Name	S alary
201	John	James	3000
202	Brat	Frank	2000
203	Mary	Jennifer	3000
204	Adam	Borg	2000
205	Smith	Joyce	1000

Query 1: Display the employees whose salary is above 2000.
 $\{E | \exists E \in \text{Employee}(E.\text{salary} > 2000)\}$

Output:

E Id	F Name	L Name	Salary
201	John	James	3000
203	Mary	Jennifer	3000

Query 2: Display the employee Ids whose salary is above 1000 and below 3000.

$\{P | \exists E \in \text{Employee} ((E.\text{salary} > 1000 \wedge E.\text{salary} < 3000) \wedge P.\text{EId} = E.\text{EId})\}$

P is a table, in which EIds are stored, from tuples which satisfies the given condition.

TUPLE RELATIONAL CALCULUS

A non-procedural language, where each query is of the form $\{t | p(t)\}$. It is a set of all tuples t such that predicate p is true for t , t is a tuple variable, $t[A]$ denotes the value of tuple ' t ' on attribute A .

$$\{T | p(T)\}$$

T is a tuple and $P(T)$ denotes a formula in which tuple variable T appears.

- $\forall \times (P(X))$

\forall is called the universal or 'for all' quantifier because every tuple in 'the universe of' tuples must make F true to make the quantified formula true.

Only true if $p(X)$ is true for every X in the universe.

Example: $\forall \times (x.\text{color} = \text{'Red'})$
 means everything that exists is red.

Example: $\forall \times ((x \in \text{Boats}) \Rightarrow (X.\text{color} = \text{'Red'}))$
 ' \Rightarrow ' is a logical implication. $a \Rightarrow b$ means that if a is true, b must be true

(or)

$\exists \times \in \text{Boats} (X.\text{color} = \text{'Red'})$

For every ' x ' in the boats relation, the color must be red.

- $\exists \times (P(X))$

\exists is called the *existential* or '*there exists*' quantifier because any tuple that exists in 'the universe of' tuples may take F true, to make the quantified formula true.

Example: $\exists \times ((X \in \text{Boats}) \wedge X.\text{color} = \text{'Red'})$

There exists a tuple X in the boats relation whose color is red.

(or)

$\exists \times \in \text{Boats} (X.\text{color} = \text{'Red'})$

Examples:

- Find all sailors with rating above 8.

Sid	Sname	Rating	Age
28	Yuppy	9	35
35	Rubber	8	55
44	grove	5	35
58	rusty	10	35

Sailors:

Solution: $\{s | s \in \text{sailors} \wedge s.\text{rating} > 8\}$

Output

Sid	Sname	Rating	Age
28	yuppy	9	35
58	rusty	10	35

$R =$

- Find names and ages of sailors with rating > 8 .

Solution: $\{R | \exists S \in \text{sailors} (s.\text{rating} > 8 \wedge R.\text{sname} = s.\text{sname} \wedge R.\text{age} = s.\text{age})\}$

Output:

sname	age
yuppy	35
rusty	35

Join Operation in Tuple Relational Calculus

Examples:

- Find sailors rated > 7 who have reserved boat = 103.

Solution: $\{S | S \in \text{sailors} \wedge s.\text{rating} > 7 \wedge \exists R (R \in \text{reserves} \wedge R.\text{sid} = s.\text{sid} \wedge R.\text{bid} = 103)\}$

- Find sailors rated > 7 who have reserved a red boat.

Solution: $\{S | S \in \text{sailors} \wedge s.\text{rating} > 7 \wedge \exists R (R \in \text{reserves} \wedge R.\text{sid} = s.\text{sid} \wedge \exists B (B \in \text{Boats} \wedge B.\text{bid} = R.\text{bid} \wedge B.\text{color} = \text{'Red'}))\}$

Division Operation in Tuple Relational Calculus

Examples

- Find sailors who have reserved all boats.

Solution: $\{S | S \in \text{sailors} \wedge \forall B \in \text{boats} (\exists R \in \text{reserves} (s.\text{sid} = R.\text{sid} \wedge B.\text{bid} = R.\text{bid}))\}$

Domain Relational Calculus

- Tuple relational and domain relational are semantically similar.
- In TRC, tuples share an equal status as variables, and field referencing can be used to select tuple parts.

3. In *DRC* formed variables are explicit.
4. *DRC* query has the following form.
 $\{ \langle x_1, x_2, \dots, x_n \rangle / P(\langle x_1, x_2, \dots, x_n \rangle) \}$
 Result included all tuples $\langle x_1, x_2, \dots, x_n \rangle$
 That make the formula $p(\langle x_1, x_2, \dots, x_n \rangle)$ true.
5. Formula given in *DRC* is recursively defined. First start with simple atomic formula and expand the formulas by using the logical connectives.
6. A variable that is not bound is free.
7. The variable X_1, X_2, \dots, X_n that appear in the left side of '/' must be the only free variable in the formula $p(\dots)$.

Example: Consider the employee table given in the above Example.

The use of quantifiers $\exists x$ and $\forall x$ in a formula is said to bind x

Query 1: Display the Employees whose salary is above 2000?

$\{ \langle I, F, L, S \rangle / \langle I, F, L, S \rangle \in \text{Employee} \wedge S > 2000 \}$

Query 2: Display the Elds of Employees, whose salary is above 1000 and below 3000?

$\{ \langle I \rangle / \exists F, L, S (\langle I, F, L, S \rangle \in \text{Employee} \wedge (S > 1000 \wedge S < 3000)) \}$

SQL (STRUCTURED QUERY LANGUAGE)

When a user wants to get some information from a database file, he/she can issue a query. A query is a user-request to retrieve data (or) information with a certain condition. SQL is a query language that allows user to specify the conditions (instead of algorithms)

Concept of SQL

The user specifies a certain condition. The program will go through all the records in the database file and select those records that satisfy the condition. The result of the query will be stored in the form of a table.

Features of SQL

1. SQL is a language of database. It includes database creation, deletion, fetching rows and modifying rows.
2. SQL is a structured query language for storing, manipulating and retrieving data stored in relational database.
3. It allows users to describe the data.
4. It allows users to create and drop database and tables.
5. It allows users to create view, functions in a database.
6. Allows users to set permissions on tables and views.
7. The standard SQL commands to interact with relational database are CREATE, SELECT, UPDATE, INSERT, DROP and DELETE.
8. The commands can be classified as follows:
 - *Data query language:* SELECT – It retrieves particular rows which satisfies the given condition.
 - *Data definition language:* CREATE, ALTER, DROP.
 - *Data manipulation language:* INSERT, UPDATE, DELETE

Features

1. Strong data protection
2. Robust transactional support
3. High performance
4. High availability
5. Security and flexibility to run anything
6. Easy to manage
7. User friendly

General Structure

SELECT ... FROM ... WHERE

SQL is divided into two languages

1. DML (data manipulation language)
 - SELECT: Extracts data from a database table.
 - UPDATE: Updates data in a database table.
 - DELETE: Deletes data from a database table.
 - INSERT INTO: Inserts new data into database table.
2. DDL (data definition language)
 - CREATE TABLE - creates a new database table.
 - ALTER TABLE - Alters a database table.
 - DROP TABLE - deletes a database table.
 - CREATE INDEX - Creates an index (search key).
 - DROP INDEX - Deletes an index.
 - RENAME – Changes the name of the table.

Types of keys:

1. Candidate key
2. Primary key
3. Super key
4. Foreign key
5. Composite primary key

In relational database, 'keys' play a major role. Keys are used to establish and identify relation between relations (or) tables.

Keys are used to ensure that each record within a table can be uniquely identified by combining one or more fields (or) column headers within a table.

Candidate key: A candidate key is a column or set of columns in a table that contains unique values, with these we can uniquely identify any database record without referring to any other columns data.

Each table may have one or more candidate keys, among the available candidate keys, one key is preserved for primary key.

A candidate key is a subset of a super key.

Example: Student

StudentId	First name	Last name	Course Id
CS00345	Jim	Black	C2
CS00254	Carry	Norris	C1
CS00349	Peter	Murray	C1
CS00196	John	Mc Cloud	C3
CS00489	Brat	Holland	C4
CS00553	Mary	Smith	C5

In the above table, we have studentId that uniquely identifies the students in a student table. This would be a candidate key.

In the same table, we have student’s first name and last name, which are also candidate keys.

1. If we combine first name and last name then also it becomes a candidate key.
2. Candidate key must (have)
 - Unique values
 - No null values
 - Minimum number of fields to ensure uniqueness.
 - Uniquely identify each record in the table.
3. The candidate keys which are not selected for primary key are known as secondary keys or alternative keys.

Primary key: A primary key is a candidate key that is most suitable (or) appropriate to become main key of the table.

1. It is a special relational database table column ((or) combination of columns)
2. Primary key main features are
 - It must contain a unique value for each row of data.
 - It cannot contain null values.

Example: We can choose primary key as studentId which is mentioned in the table given in above example.

Composite primary key: A key that consists of two or more attributes that uniquely identify an entity is called *composite key* or *composite primary key*.

Example: Customer

Cust-Id	Order-Id	Sale-details
C1	O-2	Sold
C1	O-3	Sold
C2	O-2	Sold
C2	O-3	Sold

Composite primary key is {cust-Id, order-Id}

Super key: A super key is a combination of attributes that can be uniquely used to identify a database record. A table can have any number of super keys.

1. Candidate key is a special subset of super keys.

Example: Customer

Customer name	Customer Id	SSN	Address	DOB
---------------	-------------	-----	---------	-----

Assume that we can guarantee uniqueness only for SSN field, then the following are some of the super keys possible.

1. {Name, SSN, DOB}
2. {ID, Name, SSN}

In a set of attributes, there must be at least one key (could be primary key or candidate key)

Foreign key: A foreign key is a column or group of columns in a relational database table that provides connectivity between data in two tables.

1. The majority of tables in a relational database system adhere to the concept of foreign key.
2. In complex databases, data must be added across multiple tables, thus the link or connectivity has to be maintained among the tables.
3. The concept of Referential Integrity constraint is derived from Foreign key.

Example: Emp

EId	EName	Dept – No
-----	-------	-----------

Dept

Dept-No	DName
---------	-------

In the above specified tables, Dept-No is common to both the tables, In Dept table it is called as primary key and in Emp table it is called as *foreign key*.

These two tables are connected with the help of ‘Dept-No’ field

1. For any column acting as a foreign key, a corresponding value should exist in the link (or) connecting table.
2. While inserting data and removing data from the foreign key column, a small incorrect insertion or deletion destroys the relationship between the two tables.

SQL Commands

SELECT statement

The most commonly used SQL command is SELECT statement. The SQL SELECT statement is used to query or retrieve data from a table in the database. A query may retrieve information from specified columns or from all of the columns in the table. To create a simple SQL SELECT statement, you must specify the column(s) names and the table name.

Syntax: SELECT column-name (s) from table name

Example: Persons

Lastname	Firstname	Address	City
Hansen	Ola	SpRoad,-20	Hyd
Svendson	Tove	GPRoad,-18	Secbad
Petterson	Kari	RpRoad,-19	Delhi

1. SELECT lastname FROM persons

Output:

Lastname
Hansen
Svendson
Petterson

2. SELECT lastname, firstname FROM persons

Output:

Lastname	Firstname
Hansen	Ola
Svendson	Tove
Petterson	Kari

DISTINCT statement

Returns distinct values. It eliminates duplicate values.

Syntax: Select DISTINCT column_name (s) from table-name

Example: Orders

Company	Order.No
IBM	3412
DELL	5614
WIPRO	4412
DELL	4413

1. SELECT company FROM orders

Output:

Company
IBM
DELL
WIPRO
DELL

2. SELECT DISTINCT company FROM orders

Company
IBM
DELL
WIPRO

WHERE statement

The WHERE clause is used when you want to retrieve specific information from a table excluding other irrelevant data. By using WHERE clause, we can restrict the data that is retrieved. The condition provided in the WHERE clause filters the rows retrieved from the table and gives only those rows which were expected. WHERE clause can be used along with SELECT, DELETE, UPDATE statements.

The WHERE clause is used to specify a selection condition. All conditions are specified in this clause.

Syntax: SELECT column FROM table WHERE column operator value.

Operates used in where clause:

- =
- <> (not equal) (or) !=
- >
- <
- >=
- <=

BETWEEN - Between an inclusive range.

LIKE - Search for a pattern

Example: Persons

Lastname	Firstname	Address	City	Year
Hansen	Ola	SPRoad, 16	Hyd	1956
Svendson	Tiva	GPRoad, 18	Sec	1977
Smith	Ole	RPRoad, 19	Hyd	1986
Petterson	Kari	SPRoad, 17	Sec	1985

1. SELECT * FROM persons

Output: It displays the entire table

2. SELECT * FROM persons WHERE city = 'Hyd'

Output:

Lastname	Firstname	Address	City	Year
Hansen	Ola	SPRoad, 16	Hyd	1956
Smith	Ole	RPRoad, 19	Hyd	1986

LIKE condition

The LIKE operator is used to list all rows in a table whose column values match a specified pattern. It is useful when you want to search rows to match a specific pattern, or when you do not know the entire value. For this purpose, we use a wildcard character '%'.

The LIKE condition is used to specify a search for a pattern in a column.

A '%' sign can be used to define wildcards (missing letters in the pattern) both before and after the pattern.

Syntax: SELECT column FROM table WHERE column LIKE pattern

1. SELECT * FROM persons WHERE Firstname LIKE 'O%'

Solution: SQL statement will return persons with first names that start with a letter 'O'

Output:

Lastname	Firstname	Address	City	Year
Hansen	Ola	SPRoad, 16	Hyd	1956
Smith	Ole	RPRoad, 19	Hyd	1986

2. SELECT * FROM persons WHERE Firstname LIKE '%a'

Solution: SQL statement will return persons whose first name ends with letter 'a'.

Output:

Last name	First name	Address	City	Year
Hansen	Ola	SPRoad, 16	Hyd	1956
Svendson	Tiva	GPRoad, 18	Sec. bad	1977

3. `SELECT * FROM persons WHERE firstname LIKE 'la%'`

Solution: SQL statement returns persons whose firstname contains 'la'. The word sequence 'la' may come at any place in the word.

Output:

Last name	First Name	Address	City	Year
Hansen	Ola	SPRoad, 16	Hyd	1956

String operations

- '%idge%' matches 'Rockridge', 'Ridgeway', 'Perryridge'.
- '_____' matches a string of three characters.
- '____%' matches a string of at least three characters.

INSERT INTO statement

This statement is used to insert new rows into a table. While inserting a row, if you are adding values for all the columns of the table you need not specify the column(s) name in the SQL query. But you need to make sure the order of the values is in the same order as the columns in the table. When adding a row, only the characters or data values should be enclosed with single quotes and ensure the data type of the value and the column matches. One can specify the columns for which you want to insert data

Syntax: `INSERT INTO table-name (column1, column2 ...) VALUES (value 1, value2 ...)`

- `INSERT INTO persons VALUES ('Hetland', 'Camilla', 'HPRoad 20', 'Hyd')`

Output:

Last name	First Name	Address	City
Hansen	Ola	S.P Road 16	Hyd
Svesdon	Tiva	GP Road 18	Secbad
Smith	Ole	RP Road 19	Hyd
Petterson	Kari	SP Road 17	Secbad
Hetlan	Camilla	HPRoad, 20	Hyd

- Insert data into specified columns
`INSERT INTO persons (Lastname, Address) VALUES ('Rasmussen', 'street 67')`

Output:

Last name	First Name	Address	City
Hansen	Ola	SP Road 16	Hyd
Svesdon	Tiva	GP Road 18	Secbad
Smith	Ole	RP Road 19	Hyd
Petterson	Kari	SP Road 17	Secbad
Hetlan	Camilla	HP Road 20	Hyd
Rasmussen		Street 67	

UPDATE

The update statement is used to modify the data in a table.

Syntax: `UPDATE table_name
SET Column_name = new_value
WHERE column_name = some_value.`

- Add a first name (Nine) to the person whose last name is 'Rasmussen'?

Solution: `UPDATE person SET Firstname = 'Nine'
WHERE Lastname = 'Rasmussen'`

- Change the address and add the name of the city as Hyd of a person with last name Rasmussen?

Solution: `UPDATE person
SET Address = 'street 12',
city = 'Hyd'
WHERE Lastname = 'Rasmussen'`

DELETE statement

The DELETE statement is used to delete rows from a table. The WHERE clause in the SQL delete command is optional, and it identifies the rows in the column that gets deleted. If you do not include the WHERE clause, all the rows in the table will be deleted.

Syntax: `DELETE FROM table_name
WHERE column_name = some_value`

- Delete all rows?

Solution: `DELETE * FROM table_name`

Cartesian product

The Cartesian product of two sets is the set of all ordered pairs of elements such that the first element in each pair belongs to the first set and the second element in each pair belongs to the second set. It is denoted by $\text{cross}(X)$.

For example, given two sets:

$$S1 = \{1, 2, 3\} \text{ and } S2 = \{4, 5, 6\}$$

The Cartesian product $S1 \times S2$ is the set

$$\{(1, 4), (1, 5), (1, 6), (2, 4), (2, 5), (2, 6), (3, 4), (3, 5), (3, 6)\}$$

Example:

Female		Male	
Name	Job	Name	Job
Komal	Clerk	Rohit	Clerk
Ankita	Sales	Raju	Sales

Assume that the tables refer to male and female staff, respectively. Now, in order to obtain all possible inter-staff marriages, the cartesian product can be taken.

Male-Female

Female name	Female job	Male name	Male job
Komal	Clerk	Rohit	Clerk
Komal	Clerk	Raju	Sales
Ankita	Sales	Rohit	Clerk
Ankita	Sales	Raju	Sales

Examples:

1. Find the Cartesian product of borrower and loan?

Solution: SELECT * FROM borrower, loan

2. Find the name, loan-no, and loan amount of all customers having a loan at the Perryridge branch?

Solution: SELECT customer_name, borrower. loan_number, amount FROM borrower, loan WHERE borrower. loan-no = Loan.loan_no AND branch_name = 'perryridge'

3. Find all loan numbers for loans made at the perryridge branch with loan amount greater than 1200?

Solution: SELECT loan-no FROM loan WHERE branch.name = 'perryridge' AND amount > 1200

Comparison operator

Relation algebra includes six comparison operators (=, <>, <, >, <=, >=). These are proposition forming operators on terms. For example, $x <> 0$ asserts that x is not equal to 0. It also includes three logical operators (AND, OR, NOT). These are proposition forming operators on propositions.

Example: $x > 0$ and $x < 8$

Comparison results can be combined using the logical connections AND, OR NOT

1. Find the loan-no of those loans with amounts between 90,000 and 1,00,000?

Solution: SELECT loan-no FORM loan WHERE amount BETWEEN 90,000 AND 1,00,000. SQL allows renewing relations and attributes using 'AS' clause

2. Find the name, loan-no and loan amount of all customers, rename the column name loan-no as loan.id?

Solution: SELECT customer.name, borrower.loan no AS loan.id, amount FROM borrower, loan, WHERE borrower. loan-no = loan.loan-no

Ordering of Tuples

It lists the tuples in alphabetical order.

Example: List in alphabetic order, the names of all customers having a loan in Perryridge branch?

Solution: SELECT customer-name FROM borrower WHERE branch.name = 'perryridge' ORDERBY customer-name

We may specify 'desc' for descending order (or) 'asc' for ascending order. - 'asc' is default.

Example: ORDERBY customer-name desc.

Join (⋈)

SQL Join is used to get data from two (or) more tables, which appear as single table after joining.

1. Join is used for combining columns from two or more tables by using values common to both tables.
2. Self Join: A table can also join to itself is known as self join. Types of JOIN

1. INNER JOIN
2. OUTER JOIN
 - (i) LEFT OUTER JOIN
 - (ii) RIGHT OUTER JOIN
 - (iii) FULL OUTER JOIN

1. INNER JOIN (or) EQUI JOIN

It is a simple JOIN in which result is based on matching tuple, depending on the equality condition specified in the query.

Syntax: SELECT Column-names FROM table name1 INNER JOIN table name 2 WHERE table name 1. Column name = table name 2.column – name.

Example: Class

SID	Name
11	Ana
12	Bala
13	Sudha
14	adam

Info

SID	City
11	Bangalore
12	Delhi
13	Hyderabad

SELECT *
FROM Class INNER JOIN Info
WHERE Class.SID = Info.SID

Result:

SID	Name	SID	City
11	Ana	11	Banglore
12	Bala	12	Delhi
13	Sudha	13	Hyderabad

NATURAL JOIN:

NATURAL JOIN is a type of INNER JOIN which is based on column having same name and same data type present in two tables on which join is performed.

Syntax: SELECT *
FROM table-name1 NATURAL JOIN table-name 2

Example: Consider the tables class and Info, and the following Query

SELECT *
FROM class NATURAL JOIN Info

Result:

SID	Name	City
11	Ana	Bangalore
12	Bala	Delhi
13	Sudha	Hyderabad

Both tables being joined have SID column (same name and same data type), the tuples for which value of SID matches in both the tables, appear in the result.

Dangling tuple: When NATURAL JOIN is performed on two tables, there would be some missing tuples in the result of NATURAL JOIN

Those missing tuples are called *Dangling tuples*. In the above example, the number of dangling tuples is 1 that is

14	Adam
----	------

OUTER JOIN: Outer Join is based on both matched and unmatched data.

LEFT OUTER JOIN: Left outer Join returns the tuples available in the left side table with the matched data of 2 tables and null for the right tables column.

Example: Consider the table's class and Info
SELECT *
FROM class LEFT OUTER JOIN Info
ON(class.SID = Info. SID)

Result:

SID	Name	City
11	Ana	Banglore
12	Bala	Delhi
13	Sudha	Hyderabad
14	adam	NULL

RIGHT OUTER JOIN: RIGHT OUTER JOIN returns the tuples available in the Right side table with the matched data of 2 tables and NULL for the left table's column.

Example: Class 1

SID	Name
16	Arun
17	Kamal

Info 1

SID	City
16	Chennai
17	Noida

Query:
SELECT *
FROM Class1 RIGHT OUTER JOIN Info1
ON(class1.SID = Info1.SID)

Result:

SID	Name	City
16	Arun	Chennai
18	NULL	Noida

FULL OUTER JOIN: The full outer Join returns the tuples with the matched data of two tables, remaining rows of both left table and Right table are also included.

Example: Consider the tables class 1 and Info1

Query:
SELECT *
FROM class1 FULL OUTER JOIN Info1
ON(class1.SID = Info1.SID)

Result:

SID	Name	City
16	Arun	Chennai
17	Kamal	NULL
18	NULL	Noida

ALTER command: ALTER command is used for altering the table structure

1. It is used to add a new column to existing table.
2. To rename existing column.
3. ALTER is used to drop a column.
4. It is used to change data type of any column or modify its size.

Add new column: By using alter command, we can add a new column to the table.

Syntax: ALTER table table-name ADD(column-name data type).

Example: Consider a student table.

SID	S Name	Grade
-----	--------	-------

Add a new column called address

ALTER table student ADD (address char);

Example: Add multiple columns, parent-name, course-Name, date-of-birth to student table.

ALTER table student ADD (parent-name varchar(60), course-Name varchar(20), date-of-birth date);

Example: Change the data type of column address to varchar?

ALTER table student modify(address varchar(30))

Example: Rename a column address to Location

```
ALTER table student rename address to Location
```

TRUNCATE command: Truncate command removes all tuples from a table, this command will not destroy the tables structure.

Syntax: Truncate table table-name

DROP Command: DROP query removes a table completely from database. This command will destroy the table structure.

Syntax: Drop table table-name

Rename: This command is used to rename a table.

Syntax: Rename table old-table-name to new-table-name.

Example: Rename table Employee to New-Employee.

DROP a column: Alter command can be combined with DROP command to remove columns from a table.

Syntax: alter table table-name DROP(column-name)

Example: Alter table student DROP (grade)

EXERCISES

Practice Problems I

Directions for questions 1 to 20: Select the correct alternative from the given choices.

1. Consider the given table called *Persons*

P-Id	Lastname	Firstname	Address	City
1	Hansen	ola	Timoteivn -10	Sandnes
2	Svendson	Tove	Brazil-50	Sandnes
3	Petterson	Kari	Storgt-20	Stavanger
4	Joseph	ole	Brazil-20	Sandnes

Write a query to select the persons with first name 'Tove' and last name 'Svendson'?

- (A) SELECT *
FROM Persons
WHERE first-name='tove'
AND last-name='svendson'
- (B) SELECT *
FROM Persons
WHERE first-name='tove'
OR last-name='svendson'
- (C) SELECT first-name
FROM Persons
WHERE first-name='tove'
AND last-name='svendson'
- (D) SELECT last-name
FROM Persons
WHERE first-name='tove'
AND last-name='svendson'
2. Write a query to select only the persons with last name 'Svendson' and the first name equal to 'Tove' or 'ola'?
- (A) SELECT *
FROM Persons
WHERE last-name='svendson'
AND first-name='tove'
- (B) SELECT *
FROM Persons
WHERE last-name='svendson'
AND (first-name='tove' OR first-name='ola')

(C) SELECT *
FROM Persons
WHERE last-name='svendson'
AND (first-name='tove' AND first-name='ola')

(D) SELECT *
FROM Persons
WHERE last-name='svendson'
OR (first-name='tove' AND first-name='ola')

3. Write an SQL statement to add a new row, but only in specified columns, for the persons table add data into columns 'P-Id', 'Last name' and the 'First name' with values (5, Teja, Jakob)?

- (A) INSERT INTO Persons VALUES(5,'teja','jakob')
- (B) INSERT INTO Persons VALUES(5,teja,jakob)
- (C) INSERT INTO Persons (P-Id, last-name, first-name) VALUES(5,'teja','jakob')
- (D) INSERT INTO Persons(P-Id, last-name, first-name) VALUES(5,teja,jakob)

4. Write an SQL statement:

(i) To select the persons living in a city that starts with 'S' from the 'Persons' table?

(A) SELECT *
FROM Persons
WHERE city LIKE 's__'.

(B) SELECT *
FROM Persons
WHERE city LIKE 's%'.

(C) SELECT *
FROM Persons
WHERE city LIKE '%s'.

(D) SELECT *
FROM Persons
WHERE city LIKE '_s%'.

(ii) To select the persons living in a city that contains the pattern 'tav' from 'Persons' table?

(A) SELECT *
FROM Persons
WHERE city LIKE '_tav_'.

(B) SELECT *
FROM Persons
WHERE city LIKE '_tav%'.

- (C) SELECT *
FROM Persons
WHERE city LIKE '%tav_'
- (D) SELECT *
FROM Persons
WHERE city LIKE '%tav%'
- (iii) To select the persons whose last name starts with 'b' or 's' or 'p' ?
 - (A) SELECT *
FROM Persons
WHERE last-name LIKE 'b-s-p'
 - (B) SELECT *
FROM Persons
WHERE last-name LIKE 'b%s%p'
 - (C) SELECT *
FROM Persons
WHERE last-name LIKE 'b%s%p%'
 - (D) SELECT *
FROM Persons
WHERE last-name LIKE '[bsp]%'

5. Consider the given table called 'Persons'

P-Id	Last-name	First-name	Address	City
1	Hansen	ola	Timoteivn-10	Sandnes
2	Svendson	Tove	Brazil-50	Sandnes
3	Petterson	Kari	Storgt-20	Stavanger

and the 'Orders' table

O-Id	Order No	P-Id
11	77895	3
12	44678	3
13	22456	1
14	24562	1
15	34764	5

perform NATURAL JOIN operation on both the tables and what is are the O_Id's displayed in the result?

- (A) 11, 12, 13 (B) 11, 13, 14
 - (C) 11, 12, 13, 14 (D) 12, 13, 14
6. Write an SQL to perform FULL JOIN operation on both 'Person' and 'Orders' tables and What is the number of tuples in the Result?
- (A) 4 (B) 5
 - (C) 6 (D) 7
7. Consider the given table 'Result'.

Student Name	Marks
A	55
B	90
C	40
D	80
E	85
F	95
G	82

- (i) Find out the students who have scored more than 80 marks, and display them in descending order according to their marks?
 - (A) SELECT student-name,marks
FROM Result
WHERE marks > 80
ORDERBY marks DESC
 - (B) SELECT *
FROM Result
WHERE marks > 80
ORDERBY marks DESC
 - (C) SELECT student-name,marks
FROM Result
WHERE marks > 80
ORDERBY marks
 - (D) (A) and (B)
 - (ii) From the above table, find out the top-most three students.
 - (A) SELECT student-name
FROM Result
ORDERBY marks DESC > 3
 - (B) SELECT student-name
FROM Result
ORDERBY marks DESC = 3
 - (C) SELECT student-name
FROM Result
ORDERBY marks DESC limit 3
 - (D) None of these
8. From the table 'Results', Identify the suitable SQL expression?
- (i) Find out the student Who stood 2nd?
 - (A) SELECT student-name
FROM Result
ORDERBY marks DESC limit 2
 - (B) SELECT student-name
FROM Result
ORDERBY marks DESC limit 1,1
 - (C) SELECT student-name
FROM Result
ORDERBY marks DESC limit 1,2
 - (D) SELECT student-name
FROM Result
ORDERBY marks DESC limit 2,1
 - (ii) Find out how many students scored >= 80.
 - (A) SELECT COUNT(*)
FROM Result
WHERE marks >= 80
 - (B) SELECT COUNT
FROM Result
WHERE marks >= 80
 - (C) SELECT SUM(*)
FROM Result
WHERE marks >= 80
 - (D) SELECT SUM
FROM Result
WHERE marks >= 80

9. Consider the given tables:

Customer

Customer name	Customer street	Customer city
Sonam	Mirpurroad	Dhaka
Sonam	Aga KhaRoad	Bogra
Anusha	XYZRoad	Kanchi
Nandy	MirpurRoad	Dhaka

Account

Account number	Customer name	Balance
A-101	Anusha	1000
A-102	Anusha	1500
A-103	Sonam	2000
A-104	Nandy	2500

From the customer table, find out the names of all the customers who live in either Dhaka or Bogra?

- (A) SELECT customer-name
FROM customer
WHERE customer-city='dhaka' OR
customer-city='bogra'
- (B) SELECT customer-name
FROM customer
WHERE customer-city=dhaka OR
customer-city='bogra'
- (C) SELECT customer-name
FROM customer
WHERE customer-city='dhaka' AND
customer-city='bogra'
- (D) SELECT customer-name
FROM customer
WHERE customer-city='dhaka' EXIST
customer-city='bogra'

10. Consider the given tables

Loan

Loan Number	Branch Name	Amount
L-101	Dhaka	1000
L-103	Khulna	2000

Borrower:

Customer name	Loan number
Sonam	L-101
Nandy	L-103
Anusha	L-103

(i) What are the number of tuples present in the result of cross product of the above two tables?

- (A) 4 (B) 5
(C) 6 (D) 7

(ii) Find the loan-numbers from loan table where branch-name is Dhaka?

- (A) SELECT loan-number
FROM loan
WHERE branch-name='dhaka'
- (B) SELECT loan-number
FROM branch-name='dhaka'
- (C) SELECT loan-number
FROM Loan \times Borrower
- (D) Both (A) and (C)

11. (i) Find all customers who have only accounts but no loans.

- (A) SELECT customer-name
FROM depositor LEFT OUTER JOIN Borrower ON
Depositor.customer-name=Borrower.customer-name
WHERE loan-number IS NULL
- (B) SELECT customer-name
FROM depositor LEFT OUTER JOIN Borrower ON
Depositor.customer-name = Borrower.customer-name
WHERE loan-number=NULL
- (C) SELECT customer-name
FROM depositor RIGHT OUTER JOIN Borrower ON
Depositor.customer-name=Borrower.customer-name
WHERE loan-number IS NULL
- (D) SELECT customer-name
FROM depositor RIGHT OUTER JOIN Borrower ON
Depositor.customer-name=Borrower.customer-name
WHERE loan-number=NULL

(ii) Find the names of all customers who have either an account or loan but not both.

Borrower

Customer name	Loan no.
Sonam	L-101
Sonam	L-102
Anusha	L-103

Depositor

Customer name	Account no.
Anusha	A-102
Sonam	A-103
Nandy	A-104

- (A) SELECT customer name
FROM depositor FULL OUTER JOIN
Borrower ON
Depositor.customer-name=Borrower.customer-
name
WHERE loan-number IS NULL OR Account-
number=NULL
- (B) SELECT customer-name
FROM depositor FULL OUTER JOIN
Borrower ON
Depositor.customer-name = Borrower.customer-
name
WHERE loan-number IS NULL OR Account-
number IS NULL
- (C) SELECT customer-name
FROM depositor FULL OUTER JOIN
Borrower ON
Depositor.customer-name = Borrower.customer-
name
WHERE loan-number = NULL OR Account-num-
ber = NULL
- (D) SELECT customer-name
FROM depositor FULL OUTER JOIN Borrower
ON
Depositor.customer-name = Borrower.customer-
name
WHERE loan-number=NULL OR Account-num-
ber IS NULL

12. Consider the following 'employee' table

Employee name	Branch name	Branch city	Salary
A	DU	Dhaka	1000
B	DU	Dhaka	2000
C	BUET	Dhaka	3000
D	KUET	Khulna	4000
E	KU	Khulna	5000
F	RU	Rajshahi	6000

- (i) Find the distinct number of branches appearing in the employee relation.
 - (A) SELECT COUNT(branch-name)
FROM Employee
 - (B) SELECT COUNT(DISTINCT branch-name)
FROM Employee
 - (C) SELECT DISTINCT COUNT(branch-name)
FROM Employee
 - (D) SELECT COUNT(*)
FROM Employee
- (ii) Find the total salary of all employees at each branch of the bank.
 - (A) SELECT branch-name, SUM(salary)
FROM Employee
GROUP BY Branch-city
 - (B) SELECT branch-name, SUM(salary)
FROM Employee
GROUP BY Branch-name

- (C) SELECT SUM(salary)
FROM Employee
GROUP BY Branch-name
- (D) SELECT branch-name, SUM(salary)
FROM Employee
- (iii) Find branch city, branch name Wise total salary, average salary and also number of employees.
 - (A) SELECT branch-city, branch-name,
SUM (salary), AVG(salary),
COUNT (Employee-name)
FROM Employee
GROUP BY branch-city, branch-name
 - (B) SELECT branch-city, branch-name,
SUM (salary), AVG (salary),
COUNT (Employee-name)
FROM Employee
GROUP BY branch-city
 - (C) SELECT branch-city, branch-name,
SUM (salary), AVG (salary), COUNT (Em-
ployee-name)
FROM Employee
GROUP BY branch-name
 - (D) SELECT branch-name, SUM (salary),
AVG (salary), COUNT (Employee-name)
FROM Employee
GROUP BY branch-city, branch-name

Common data for questions 13 to 15: Consider the SHIPMENTS relation and write the SQL statements for the below

SUPPLIERS

Supplier number	Supplier name	Status	City
SN1	Suma	30	Hyderabad
SN2	Hari	20	Chennai
SN3	Anu	10	Hyderabad
SN4	Mahesh	20	Bombay
SN5	Kamal	30	Delhi

PARTS

Part number	Part name	Color	Weight	City
PN1	X	Red	13.0	Chennai
PN2	Y	Green	13.5	Bombay
PN3	X	Yellow	13.2	Hyderabad
PN4	Y	Green	14.1	Calcutta
PN5	Z	Red	14.3	Hyderabad
PN6	Z	Blue	14.2	Bombay

PROJECT

Project number	Project name	City
PJ1	Display	Chennai
PJ2	OCR	Bombay
PJ3	RAID	Chennai
PJ4	SORTER	Hyderabad
PJ5	EDS	Chennai
PJ6	Tape	Bombay
PJ7	Console	Hyderabad

SHIPMENTS

Supplier number	Part number	Project number	Quantity
SN1	PN1	PJ1	300
SN1	PN1	PJ4	400
SN2	PN3	PJ1	350
SN2	PN3	PJ2	450
SN2	PN3	PJ3	640
SN2	PN3	PJ4	320
SN2	PN3	PJ5	330
SN2	PN3	PJ6	520
SN2	PN3	PJ7	480
SN2	PN5	PJ2	460
SN3	PN3	PJ1	440
SN3	PN4	PJ2	410
SN4	PN6	PJ3	310
SN4	PN6	PJ7	320
SN5	PN2	PJ2	340
SN5	PN2	PJ4	350
SN5	PN5	PJ5	360
SN5	PN5	PJ7	370
SN5	PN6	PJ2	380
SN5	PN1	PJ4	420
SN5	PN3	PJ4	440
SN5	PN4	PJ4	450
SN5	PN5	PJ4	400
SN5	PN6	PJ4	410

13. (i) For each part supplied, get the part number and the total shipment quantity?
- SELECT shipments.part-number, SUM(shipments.quantity)
FROM Shipments
GROUP BY shipments.part-number
 - SELECT SUM(shipments.quantity)
FROM Shipments
GROUP BY shipments.part-number
 - SELECT shipments.part-number, SUM(shipments.quantity)
FROM Shipments
GROUP BY shipments.quantity
 - SELECT shipments.part-number, SUM(shipments.part-number)
FROM Shipments
GROUP BY shipments.part-number
- (ii) Get part numbers for parts supplied by more than two suppliers?
- SELECT shipments.part-number
FROM Shipments
GROUP BY shipments.part-number
HAVING COUNT(shipments.supplier-number) > 2
 - SELECT shipments.part-number
FROM Shipments
GROUP BY shipments.part-number
HAVING COUNT(shipments.supplier-number) >= 2
 - SELECT shipments.part-number
FROM Shipments
GROUP BY shipments.part-number > 2
 - SELECT shipments.part-number, COUNT(shipments.supplier-number) > 2
FROM Shipments
GROUP BY shipments.part-number
- (iii) Get supplier names for suppliers who supply part PN3?
- SELECT DISTINCT suppliers.supplier-name
FROM Supplier
WHERE suppliers.supplier-number IN (SELECT Shipments.supplier-number
FROM Shipments
WHERE Shipments.part-number='PN3')
 - SELECT DISTINCT suppliers.supplier-name
FROM Supplier
WHERE suppliers.supplier-number NOT IN(SELECT Shipments.supplier-number
FROM Shipments
WHERE Shipments.part-number='PN3')
 - SELECT DISTINCT suppliers.supplier-name
FROM Supplier
WHERE suppliers.supplier-number EXCEPT (SELECT Shipments.supplier-number
FROM Shipments
WHERE Shipments.part-number='PN3')
 - SELECT DISTINCT suppliers, supplier-name
FROM Supplier
WHERE suppliers.supplier-number
UNION
SELECT Shipments.supplier-number
FROM Shipments
WHERE Shipments.part-number='PN3'
14. (i) Get supplier names for suppliers who supply at least one blue part.
- SELECT DISTINCT suppliers.supplier-name
FROM Suppliers
WHERE suppliers.supplier-number
IN (SELECT Shipments.supplier-number
FROM Shipments
WHERE Shipments.part-number
IN (SELECT Parts.part-number
FROM Parts
WHERE Parts.color='Blue'))
 - SELECT DISTINCT suppliers.supplier-name
FROM Suppliers
WHERE suppliers.supplier-number
IN (SELECT Shipments.supplier-number
FROM Shipments

- WHERE Shipments.part-number NOT
IN(SELECT Parts.part-number
FROM Parts
WHERE Parts.color='Blue'))
- (C) SELECT DISTINCT suppliers.supplier-name
FROM Suppliers
WHERE suppliers.supplier-number NOT
IN(SELECT Shipments.supplier-number
FROM Shipments
WHERE Shipments.part-number
IN (SELECT Parts.part-number
FROM Parts
WHERE Parts.color='Blue'))
- (D) SELECT DISTINCT suppliers.supplier-name
FROM Suppliers
WHERE suppliers.supplier-number
IN (SELECT Shipments.supplier-name
FROM Shipments
WHERE Shipments.part-number
IN (SELECT Parts.part-number
FROM Parts
WHERE Parts.color='Blue'))
- (ii) Get supplier numbers for suppliers with status less than the current maximum status in the suppliers table:
- (A) SELECT Suppliers.supplier-number
FROM suppliers
WHERE Suppliers.status < (SELECT MAX
(Suppliers.status)
FROM Suppliers)
- (B) SELECT Suppliers.supplier-number
FROM suppliers
WHERE Suppliers.status<=(SELECT MAX
(Suppliers.status)
FROM Suppliers)
- (C) SELECT Suppliers.supplier-number,
MAX (Suppliers.status)
FROM suppliers
WHERE Suppliers.status
- (D) SELECT Suppliers.supplier-number
FROM suppliers
WHERE Suppliers.status=MAX(Suppliers.
status)
- (iii) Get supplier names for suppliers who supply part PN2?
- (A) SELECT DISTINCT suppliers.supplier-name
FROM Suppliers
WHERE EXIST(SELECT *
FROM Shipments
WHERE Shipments.supplier-number = sup-
pliers.supplier-number
AND
Shipments.part-number='PN2')
- (B) SELECT DISTINCT suppliers.supplier-name
FROM Suppliers
WHERE NOT EXIST(SELECT *
FROM Shipments
WHERE Shipments.supplier-number = sup-
pliers.supplier-number
AND
Shipments.part-number='PN2')
- (C) SELECT DISTINCT suppliers.supplier-name
FROM Suppliers
WHERE EXIST(SELECT *
FROM Shipments
WHERE Shipments.supplier-number = sup-
pliers.supplier-number
OR
Shipments.part-number='PN2')
- (D) SELECT DISTINCT suppliers.supplier-name
FROM Suppliers
WHERE EXIST(SELECT *
FROM Shipments
WHERE Shipments.supplier-number = sup-
pliers.supplier-number
UNION
Shipments.part-number='PN2')
15. (i) Get supplier names for suppliers who do not supply part PN2.
- (A) SELECT DISTINCT suppliers.supplier-name
FROM Suppliers
WHERE NOT EXIST(SELECT *
FROM Shipments
WHERE Shipments.supplier-number = sup-
pliers.supplier-number
AND
Shipments.part-number='PN2')
- (B) SELECT DISTINCT suppliers.supplier-name
FROM Suppliers
WHERE EXIST(SELECT *
FROM Shipments
WHERE Shipments.supplier-number = sup-
pliers.supplier-number
AND
Shipments.part-number='PN2')
- (C) SELECT DISTINCT suppliers.supplier-name
FROM Suppliers
WHERE EXCEPT(SELECT *
FROM Shipments
WHERE Shipments.supplier-number = sup-
pliers.supplier-number
AND
Shipments.part-number='PN2')
- (D) SELECT DISTINCT suppliers.supplier-name
FROM Suppliers
WHERE NOT EXIST(SELECT *
FROM Shipments

- WHERE Shipments.supplier-number = suppliers.supplier-number
OR
Shipments.part-number='PN2')
- (ii) Get supplier names for suppliers who supply all parts.
- (A) SELECT DISTINCT suppliers.supplier-name
FROM Suppliers
WHERE NOT EXIST(SELECT *
FROM Part
WHERE NOT EXIST(SELECT * FROM Shipments
WHERE Shipments.supplier-number = suppliers.supplier-number
AND
Shipments.part-number=Parts.part-number))
- (B) SELECT DISTINCT suppliers.supplier-name
FROM Suppliers
WHERE EXIST(SELECT *
FROM Part
WHERE NOT EXIST(SELECT * FROM Shipments
WHERE Shipments.supplier-number = suppliers.supplier-number
AND
Shipments.part-number=Parts.part-number))
- (C) SELECT DISTINCT suppliers.supplier-name
FROM Suppliers
WHERE NOT EXIST(SELECT *
FROM Part
WHERE EXIST(SELECT * FROM Shipments
WHERE Shipments.supplier-number = suppliers.supplier-number
AND
Shipments.part-number=Parts.part-number))
- (D) SELECT DISTINCT suppliers.supplier-name
FROM Suppliers
WHERE EXIST(SELECT *
FROM Part
WHERE EXIST(SELECT * FROM Shipments
WHERE Shipments.supplier-number = suppliers.supplier-number
AND
Shipments.part-number=Parts.part-number))
- (iii) Get part numbers for parts that either weigh more than-16 pounds or are supplied by supplier SN3, or both?
- (A) SELECT parts.part-number
FROM parts
WHERE Parts.weight>18
UNION
SELECT Shipments.part-number
FROM shipments
WHERE Shipments.supplier-number='SN2'

- (B) SELECT parts.part-number
FROM parts
WHERE Parts.weight>18
UNION
SELECT Shipments.supplier-name
FROM shipments
WHERE Shipments.supplier-number='SN2'
- (C) SELECT parts.part-number
FROM parts
WHERE Parts.weight>18
UNION
SELECT Shipments.part-number,Shipments.supplier-name
FROM shipments
WHERE Shipments.supplier-number='SN2'
- (D) SELECT parts.part-Number, parts.color
FROM parts
WHERE Parts.weight>18
UNION
SELECT Shipments.part-number
FROM shipments
WHERE Shipments.supplier-number='SN2'

Common data for questions 16 and 17: Consider the following relation: Teach

Name	Address	course
Zohar	40B,east city	MD
Nisha	16/2, hyd	BDS
Zohar	40B, East city	MS
Ravi	New York	MBA

16. The teacher with name Zohar teaching the course MS?
- (A) $\sigma_{Name = 'Zohar'}(Teach) = MS$.
 (B) $\pi_{Name = 'Zohar'}(Teach) = MS$.
 (C) $\sigma_{Name = 'Zohar' \text{ and } course = 'MS'}(Teach)$.
 (D) $\pi_{Name = 'Zohar' \text{ and } course = 'MS'}(Teach)$.
17. Select the names of courses taught by Zohar?
- (A) $\pi_{course}(\sigma_{Name = 'Zohar'}(Teach))$
 (B) $\sigma_{course}(\pi_{Name = 'Zohar'}(Teach))$
 (C) $\pi_{course}(\sigma_{Name = 'MD'}(Teach))$
 (D) None
18. Consider the join of a relation A with a relation B . If A has m tuples and B has n tuples. Then the maximum and minimum sizes of the join respectively are.
- (A) mn and $m + n$ (B) $m + n$ and $(m - n)$
 (C) mn and m (D) mn and 0
19. Match the following:

I	Set intersection	1	$R \times S$
II	Natural join	2	$r - (r - s)$
III	Division	3	\leftarrow
IV	Assignment	4	$\pi_{R-S}(r) - \pi_{R-S}$ $(\pi_{R-S}(r) \times s)$ $-\pi_{R-S}, S(r)$

- (A) I – 2, II – 1, III – 4, IV – 3
- (B) I – 3, II – 4, III – 2, IV – 1
- (C) I – 1, II – 2, III – 3, IV – 4
- (D) I – 2, II – 3, III – 4, IV – 1

20. Which one is correct for division operations for relation r and s

- (A) $r \div s$
- (B) $\pi_{R-S}(r) - \pi_{R-S}((\pi_{R-S}(r) \times s) - \pi_{R-S}, s(r))$
- (C) Temp 1 $\leftarrow \pi_{R-S}(r)$
Temp 2 $\leftarrow \pi_{R-S}(\text{temp1} \times s) - \pi_{R-S}, s(r)$
result = temp 1 – temp 2
- (D) All the above

Practice Problems 2

Directions for questions 1 to 20: Select the correct alternative from the given choices.

1. The correct order of SQL expression is
 - (A) Select, group by, where, having
 - (B) Select, where, group by, having
 - (C) Select, group by, having, where
 - (D) Select, having, where, group by
2. Which one is not a query language?
 - (A) SQL
 - (B) QBE
 - (C) Data log
 - (D) MySQL
3. Like ‘ $a b \% c d$ ’ escape ‘\’ matches all the strings
 - (A) Ending with $a b c d$
 - (B) Beginning with $a b c d$
 - (C) Beginning with $a b c d$
 - (D) Beginning with $a b \% c d$
4. ‘ $_ _ \%$ ’ matches any string of
 - (A) At least three characters
 - (B) At most three characters
 - (C) Exactly three characters
 - (D) exactly three characters ending with %
5. Which of the following are set operations?
 - (i) Union
 - (ii) Intersection
 - (iii) Set Difference
 - (iv) Cartesian Product
 - (A) (i), (ii), (iii)
 - (B) (i), (iii), (iv)
 - (C) (i), (iii), (ii), (iv)
 - (D) (i), (ii), (iv)
6. What is the purpose of project operation?
 - (A) It selects certain columns
 - (B) It selects certain rows
 - (C) It selects certain strings
 - (D) It selects certain integers

Common data for questions 7 and 8: Person

Id	Name	Age	Hobby
11	Anu	21	Stamp Collection
22	Kamal	32	Painting
33	Ravi	24	Dancing
44	Ram	22	Singing

7. Select the persons whose hobby is either painting (or) singing.
 - (A) $\sigma_{\text{Hobby} = \text{'painting'} \text{ OR } \text{Hobby} = \text{'singing'}}(\text{person})$
 - (B) $\sigma_{\text{Hobby} = \text{'painting'}, \text{'singing'}}(\text{person})$
 - (C) $\sigma_{\text{Hobby} = \text{'painting'} \text{ OR } \text{'singing'}}(\text{person})$
 - (D) All are correct
8. Select the persons whose age is above 21 and below 32:
 - (A) $\sigma_{\text{age} > 21 \text{ AND } \text{age} < 32}(\text{person})$
 - (B) $\sigma_{21 < \text{age} < 32}(\text{person})$
 - (C) $\sigma_{\text{age} > 21 \text{ OR } \text{age} < 32}(\text{person})$
 - (D) $\sigma_{\text{age} < 21 \text{ AND } \text{age} > 32}(\text{person})$

Common data for questions 9 and 10: Consider the following relation: Teach

Name	course	Rating	Age
Zohar	MD	7	35
Nisha	BDS	8	27
Zohar	MS	7	34
Ravi	MBA	9	33

9. Select the teachers whose rating is above 7 and whose age is less than 32?
 - (A) $S_{\text{Rating} > 7 \text{ AND } \text{Age} < 32}(\text{Teach})$
 - (B) $S_{\text{Rating} \geq 7 \text{ AND } \text{Age} < 32}(\text{Teach})$
 - (C) $S_{\text{Rating} > 7 \text{ AND } < 32}(\text{Teach})$
 - (D) Both (A) and (B)
10. Select the courses with rating above 7?
 - (A) $\pi_{\text{course}}(\sigma_{\text{rating} > 7}(\text{Teach}))$
 - (B) $\sigma_{\text{course}}(\pi_{\text{rating} > 7}(\text{Teach}))$
 - (C) $\pi_{\text{name}, \text{course}}(\sigma_{\text{rating} > 7}(\text{Teach}))$
 - (D) None

Common data for questions 11 and 12: Consider the following schema of a relational database employee (empno, ename, eadd) project (pno, pname) Work-on (empno, pno) Part(partno, partname, qty-on-hand, size) Use (empno, pno, partno, number)

11. Display the names of the employees who are working on a project named ‘VB’.
 - (A) $\sigma_{\text{name}}(\text{employee} \bowtie (\sigma_{\text{pname} = \text{'VB'}}(\text{project}) \bowtie \text{worked on}))$
 - (B) $\sigma_{\text{name}}(\text{employee} \bowtie (\pi_{\text{pname} = \text{'VB'}}(\text{project}) \bowtie \text{work on}))$
 - (C) $\pi_{\text{name}}(\text{employee} \bowtie (\sigma_{\text{pname} = \text{'VB'}}(\text{project}) \bowtie \text{work on}))$
 - (D) $\pi_{\text{name}}(\text{employee} \bowtie (\pi_{\text{pname} = \text{'VB'}}(\text{project}) \bowtie \text{work on}))$

- (C) Courses in which only male students are enrolled.
- (D) None of the above

5. Consider the relation **employee** (name, sex, supervisorName) with *name* as the key. *supervisorName* gives the name of the supervisor of the employee under consideration. What does the following Tuple Relational Calculus query produce?

$e \cdot \text{name} \mid \text{employee}(e) \wedge$
 $(\forall x)[\neg \text{employee}(x) \vee x \cdot \text{supervisor Name} \neq e \cdot \text{name} \vee$
 $x \cdot \text{sex} = \text{"male"}]$ [2007]

- (A) Names of employees with a male supervisor.
 - (B) Names of employees with no immediate male subordinates.
 - (C) Names of employees with no immediate female subordinates.
 - (D) Names of employees with a female supervisor.
6. Consider the table **employee** (empId, name, department, salary) and the two queries Q_1 , Q_2 below. Assuming that department 5 has more than one employee, and we want to find the employees who get higher salary than anyone in the department 5, which one of the statements is **TRUE** for any arbitrary employee table?

Q_1 : SELECT e.empId
 FROM employee e
 WHERE not exists
 (Select * From employee s where s.department = '5'
 and s.salary >=e.salary)
 Q_2 : SELECT e.empId
 FROM employee e
 WHERE e.salary > Any
 (Select distinct salary From employee s Where
 s.department = '5') [2007]

- (A) Q_1 is the correct query
- (B) Q_2 is the correct query
- (C) Both Q_1 and Q_2 produce the same answer.
- (D) Neither Q_1 nor Q_2 is the correct query

7. Let R and S be two relations with the following schema
 $R(\underline{P}, Q, R1, R2, R3)$
 $S(\underline{P}, Q, S1, S2)$
 Where $\{P, Q\}$ is the key for both schemas. Which of the following queries are equivalent?

- I. $\Pi_p(R \bowtie S)$
 - II. $\Pi_p(R) \bowtie \Pi_p(S)$
 - III. $\Pi_p(\Pi_{p,q}(R) \cap \Pi_{p,q}(S))$
 - IV. $\Pi_p(\Pi_{p,q}(R) - (\Pi_{p,q}(R) - (\Pi_{p,q}(S))))$ [2008]
- (A) Only I and II
 - (B) Only I and III
 - (C) Only I, II and III
 - (D) Only I, III and IV

8. Let R and S be relational schemes such that $R = \{a,b,c\}$ and $S = \{c\}$. Now consider the following queries on the database:

- I. $\pi_{R-S}(r) - \pi_{R-S}(\pi_{R-S}(r) \times_S - \pi_{R-S,S}(r))$
- II. $\{t \mid t \in \pi_{R-S}(r) \wedge \forall u \in S(\exists v \in r(u = v[S] \wedge t = v[R-S]))\}$
- III. $\{t \mid t \in \pi_{R-S}(r) \wedge \forall v \in r(\exists u \in S(u = v[S] \wedge t = v[R-S]))\}$
- IV. SELECT R.a, R.b
 FROM R, S
 WHERE R.c = S.c

Which of the above queries are equivalent? [2009]
 (A) I and II (B) I and III
 (C) II and IV (D) III and IV

Common data for questions 9 and 10: Consider the following relational schema: Suppliers (sid: integer, sname: string, city: string, street: string) Parts(pid: integer, pname: string, color: string) Catalog (sid: integer, pid: integer, cost: real)

9. Consider the following relational query on the above database:
 SELECT S.sname
 FROM Suppliers S
 WHERE S.sid NOT IN (SELECT C.sid
 FROM Catalog C
 WHERE C.pid NOT IN (SELECT P.pid FROM Parts P
 WHERE P.color <> 'blue'))

- Assume that relations corresponding to the above schema are not empty. Which one of the following is the correct interpretation of the above query? [2009]
- (A) Find the names of all suppliers who have supplied a non-blue part.
 - (B) Find the names of all suppliers who have not supplied a non-blue part.
 - (C) Find the names of all suppliers who have supplied only blue parts.
 - (D) Find the names of all suppliers who have not supplied only blue parts.

10. A relational schema for a train reservation database is given below

Passenger (pid, pname, age)
 Reservation (pid, class, tid)
 Table :Passenger
 Table :Reservation

Pid	pname	Age	Pid	class	tid
0	'Sachin'	65	0	'AC'	8200
1	'Rahul'	66	1	'AC'	8201
2	'Sourav'	67	2	'SC'	8201

3	'Anil'	69	5	'AC'	8203
			1	'SC'	8204
			3	'AC'	8202

What pids are returned by the following SQL query for the above instance of the tables?

```
SELECT pid
FROM Reservation
WHERE class = 'AC' AND
EXISTS (SELECT *
FROM Passenger
WHERE age > 65 AND
Passenger.pid = Reservation.pid) [2010]
```

- (A) 1, 0
- (B) 1, 2
- (C) 1, 3
- (D) 1, 5

11. Consider a relational table r with sufficient number of records, having attributes A_1, A_2, \dots, A_n and let $1 \leq p \leq n$. Two queries $Q1$ and $Q2$ are given below.

$Q1: \pi_{A_1 \dots A_n}(\sigma_{A_p=c}(r))$ where c is a constant.

$Q2: \pi_{A_1 \dots A_n}(\sigma_{c_1 \leq A_p \leq c_2}(r))$ where c_1 and c_2 are constants.

The database can be configured to do ordered indexing on A_p or hashing on A_p . Which of the following statements is TRUE? [2011]

- (A) Ordered indexing will always outperform hashing for both queries
- (B) Hashing will always outperform ordered indexing for both queries.
- (C) Hashing will outperform ordered indexing on $Q1$, but not on $Q2$.
- (D) Hashing will outperform ordered indexing on $Q2$, but not on $Q1$.

12. Database table by name Loan_Records is given below.

Borrower	Bank manager	Loan amount
Ramesh	Sunderajan	10000.00
Suresh	Ramgopal	5000.00
Mahesh	Sunderajan	7000.00

What is the output of the following SQL query?

```
SELECT count ( * )
FROM (Select Borrower, Bank_Manager FROM
Loan Records) AS S
NATURAL JOIN
(SELECT Bank_Manager, Loan_Amount FROM
Loan_Records) AS T; [2011]
```

- (A) 3
- (B) 9
- (C) 5
- (D) 6

13. Consider a database table T containing two columns X and Y each of type *integer*. After the creation of the table, one record ($X = 1, Y = 1$) is inserted in the table.

Let MX and MY denote the respective maximum values of X and Y among all records in the table at any point in time. Using MX and MY , new records are inserted in the table 128 times with X and Y values being $MX + 1, 2 * MY + 1$ respectively. It may be noted that each time after the insertion, values of MX and MY change. What will be the output of the following SQL query after the steps mentioned above are carried out?

```
SELECT Y FROM T WHERE X = 7; [2011]
```

- (A) 127
- (B) 255
- (C) 129
- (D) 257

14. Which of the following statements are true about an SQL query?

- P: An SQL query can contain a HAVING clause even if it does not have a GROUP BY clause
- Q: An SQL query can contain a HAVING clause only if it has a GROUP BY clause
- R: All attributes used in the GROUP BY clause must appear in the SELECT clause
- S: Not all attributes used in the GROUP BY clause need to appear in the SELECT clause [2012]

- (A) P and R
- (B) P and S
- (C) Q and R
- (D) Q and S

15. Suppose $R_1(\underline{A}, B)$ and $R_2(\underline{C}, D)$ are two relation schemas. Let r_1 and r_2 be the corresponding relation instances. B is a foreign key that refers to C in R_2 . If data in r_1 and r_2 satisfy referential integrity constraints, which of the following is always true? [2012]

- (A) $\Pi_B(r_1) - \Pi_C(r_2) = \emptyset$
- (B) $\Pi_C(r_2) - \Pi_B(r_1) = \emptyset$
- (C) $\Pi_B(r_1) = \Pi_C(r_2)$
- (D) $\Pi_B(r_1) - \Pi_C(r_2) \neq \emptyset$

Common data for questions 16 and 17: Consider the following relations A, B and C :

(A)

Id	Name	Age
12	Arun	60
15	Shreya	24
99	Rohit	11

(B)

Id	Name	Age
15	Shreya	24
25	Hari	40
98	Rohit	20
99	Rohit	11

(C)

Id	Phone	Area
10	2200	02
99	2100	01

16. How many tuples does the result of the following SQL query contain?
 SELECT A.Id
 FROM A
 WHERE A. Age > ALL (SELECT B. Age
 FROM B
 WHERE B. Name = 'Arun') [2012]
- (A) 4 (B) 3
 (C) 0 (D) 1

17. How many tuples does the result of the following relational algebra expression contain? Assume that the schema of $A \cup B$ is the same as that of A .
 $(A \cup B) \bowtie_{A.Id > 40 \vee C.Id < 15} C$ [2012]
- (A) 7 (B) 4
 (C) 5 (D) 9

18. Consider the following relational schema. Students (rollno: integer, sname: string) Courses (courseno: integer, cname: string) Registration(rollno:integer,co urserno: integer, percent: real)

Which of the following queries are equivalent to this query in English?

'Find the distinct names of all students who score more than 90% in the course numbered 107'

- (I) SELECT DISTINCT S.sname FROM Students as S, Registration as R WHERE R.rollno=S.rollno AND R.courseno=107 AND R.percent>90
- (II) $\pi_{\text{sname}}(\sigma_{\text{courseno}=107 \wedge \text{percent}>90} \text{Registration} \times \text{Students})$
- (III) $\{T | \exists S \in \text{Students}, \exists R \in \text{Registration} (S.\text{rollno}=R.\text{rollno} \wedge R.\text{courseno}=107 \wedge R.\text{percent}>90 \wedge T.\text{sname}=S.\text{sname})\}$
- (IV) $\{ \langle S_N \rangle | \exists S_R \exists R_p (\langle S_R, S_N \rangle \in \text{Students} \wedge \langle S_R, 107, R_p \rangle \in \text{Registration} \wedge R_p > 90) \}$ [2013]
- (A) I, II, III and IV (B) I, II and III only
 (C) I, II and IV only (D) II, III and IV only

19. Given the following statements:

S_1 : A foreign key declaration can always be replaced by an equivalent check assertion in SQL.

S_2 : Given the table $R(a, b, c)$ where a and b together form the primary key, the following is a valid table definition.

```
CREATE TABLE S (
a INTEGER
d INTEGER,
e INTEGER,
PRIMARY KEY (d),
FOREIGN KEY (a) references R)
```

Which one of the following statements is CORRECT? [2014]

- (A) S_1 is TRUE and S_2 is FALSE
 (B) Both S_1 and S_2 are TRUE
 (C) S_1 is FALSE and S_2 is TRUE
 (D) Both S_1 and S_2 are FALSE

20. Given the following schema:

Employees (emp-id, first-name, last-name, hire-date, dept-id, salary)

Departments (dept-id, dept-name, manager-id, location-id)

you want to display the last names and hire dates of all latest hires in their respective departments in the location ID 1700. You issue the following query:

```
SQL > SELECT last-name, hire-date
FROM employees
WHERE (dept-id, hire-date) IN
(SELECT dept-id, MAX (hire-date)
FROM employees JOIN departments USING
(dept-id)
WHERE location-id = 1700
GROUP BY dept-id);
```

What is the outcome? [2014]

- (A) It executes but does not give the correct result.
 (B) It executes and gives the correct result.
 (C) It generates an error because of pair wise comparison.
 (D) It generates an error because the GROUP BY clause cannot be used with table joins in a subquery.

21. Given an instance of the STUDENTS relation as shown below:

Student ID	Student		Student	
	Name	Student Email	Age	CPI
2345	Shankar	shaker @ math	X	9.4
1287	Swati	swati @ ee	19	9.5
7853	Shankar	shankar @ cse	19	9.4
9876	Swati	swati @ mech	18	9.3
8765	Ganesh	ganesh@ civil	19	8.7

For (StudentName, StudentAge) to be a key for this instance, the value X should NOT be equal to _____.

[2014]

22. Consider a join (relation algebra) between relations $r(R)$ and $s(S)$ using the nested loop method. There are three buffers each of size equal to disk block size, out of which one buffer is reserved for intermediate results. Assuming $\text{size } r(R) < \text{size } s(S)$, the join will have fewer number of disk block accesses if [2014]
- (A) Relation $r(R)$ is in the outer loop
 (B) Relation $s(S)$ is in the outer loop
 (C) Join selection factor between $r(R)$ and $s(S)$ is more than 0.5
 (D) Join selection factor between $r(R)$ and $s(S)$ is less than 0.5
23. SQL allows duplicate tuples in relations, and correspondingly defines the multiplicity of tuples in the result of joins. Which one of the following queries always gives the same answer as the nested query shown below:

Select * from R where a in (select S. a from S)[2014]

- (A) Select $R.a$ from R, S where $R.a = S.a$
- (B) Select distinct $R.a$ from R, S where $R.a = S.a$
- (C) Select $R.a$ from R, S (select distinct a from S) as $S1$ where $R.a = S1.a$
- (D) Select $R.a$ from R, S where $R.a = S.a$ and is unique R

24. What is the optimized version of the relation algebra expression $\pi_{A_1}(\pi_{A_2}(\sigma_{F_1}(\sigma_{F_2}(r))))$, where A_1, A_2 are sets of attributes in r with $A_1 \subset A_2$ and F_1, F_2 are Boolean expressions based on the attributes in r ? [2014]

- (A) $\pi_{A_1}(\sigma_{(F_1 \wedge F_2)}(r))$
- (B) $\pi_{A_1}(\sigma_{(F_1 \vee F_2)}(r))$
- (C) $\pi_{A_2}(\sigma_{(F_1 \wedge F_2)}(r))$
- (D) $\pi_{A_2}(\sigma_{(F_1 \vee F_2)}(r))$

25. Consider the relational schema given below, where **empId** of the relation **dependent** is a foreign key referring to **empId** of the relation **employee**. Assume that every employee has at least one associated dependent in the **dependent** relation.

Consider the following relational algebra query:

employee (empId, empName, empAge)
 dependent (depId, eId, depName, depAge)

$\pi_{\text{empId}}(\text{employee}) - \pi_{\text{empId}}(\text{employee} \bowtie_{(\text{empId} = \text{eId}) \wedge (\text{empAge} \leq \text{depAge})} \text{dependent})$

The above query evaluates to the set of empIds of employees whose age is greater than that of [2014]

- (A) some dependent.
- (B) all dependents.
- (C) some of his/her dependents.
- (D) all of his/her dependents.

26. Consider the following relational schema:
 employee (empId, empName, empDept)
 customer (custId, custName, salesRepId, rating)
 salesRepId is a foreign key referring to empId of the employee relation. Assume that each employee makes a sale to at least one customer. What does the following query return?

```
SELECT empName
FROM employee E
WHERE NOT EXISTS
(SELECT custId
FROM customer C
WHERE C.salesRepId = E.empId
AND C.Rating <> 'GOOD');
```

[2014]

- (A) Names of all the employees with at least one of their customers having a 'GOOD' rating.
- (B) Names of all the employees with at most one of their customers having a 'GOOD' rating.

- (C) Names of all the employees with none of their customers having a 'GOOD' rating.
- (D) Names of all the employees with all their customers having a 'GOOD' rating.

27. SELECT operation in SQL is equivalent to [2015]

- (A) The selection operation in relational algebra
- (B) The selection operation in relational algebra, except that SELECT in SQL retains duplicates.
- (C) The projection operation in relational algebra.
- (D) The projection operation in relational algebra, except that SELECT in SQL retains duplicates.

28. Consider the following relations:

Student

Roll No	Student Name
1	Raj
2	Rohit
3	Raj

Performance

Roll No	Course	Marks
1	Math	80
1	English	70
2	Math	75
3	English	80
2	Physics	65
3	Math	80

Consider the following SQL query.

```
SELECT S.Student_Name, sum (P.Marks)
FROM Student S, Performance P
WHERE S.Roll_No = P.Roll_No
GROUP BY S.Student_Name
```

The number of rows that will be returned by the SQL query is _____ [2015]

29. Consider two relations $R_1(A, B)$ with the tuples (1, 5), (3, 7) and $R_2(A, C) = (1, 7), (4, 9)$. Assume that $R(A, B, C)$ is the full natural outer join of R_1 and R_2 . Consider the following tuples of the form (A, B, C) : $a = (1, 5, \text{null})$, $b = (1, \text{null}, 7)$, $c = (3, \text{null}, 9)$, $d = (4, 7, \text{null})$, $e = (1, 5, 7)$, $f = (3, 7, \text{null})$, $g = (4, \text{null}, 9)$. Which one of the following statements is correct? [2015]

- (A) R contains a, b, e, f, g but not c, d .
- (B) R contains all of a, b, c, d, e, f, g .
- (C) R contains e, f, g but not a, b .
- (D) R contains e but not f, g .

30. Consider the following relation

Cinema (theater, address, capacity)

Which of the following options will be needed at the end of the SQL query

SELECT P_1 .address
FROM Cinema P_1

such that it always finds the addresses of theaters with maximum capacity? [2015]

- (A) WHERE P_1 .capacity >= All (select P_2 . Capacity from Cinema P_2)
- (B) WHERE P_1 .capacity >= Any (select P_2 . Capacity from Cinema P_2)
- (C) WHERE P_1 .capacity > All (select max(P_2 . capacity) from Cinema P_2)

(D) WHERE P_1 .capacity > Any (select max(P_2 . capacity) from Cinema P_2)

31. Which of the following is NOT a superkey in a relational schema with attributes V, W, X, Y, Z and primary key VY ? [2016]

- (A) $VXYZ$
- (B) $VWXZ$
- (C) $VWXY$
- (D) $VWXYZ$

32. Consider a database that has the relation schema EMP (EmpId, EmpName and DeptName). An instance of the schema EMP and a SQL query on it are given below.

EMP		
EmpId	EmpName	DeptName
1	XYA	AA
2	XYB	AA
3	XYC	AA
4	XYD	AA
5	XYE	AB
6	XYF	AB
7	XYG	AB
8	XYH	AC
9	XYI	AC
10	XYJ	AC
11	XYK	AD
12	XYL	AD
13	XYM	AE

```
SELECTIVE AVG (EC.Num)
FROM EC
WHERE (DeptName, Num) IN
      (SELECT DeptName, COUNT (EmpId) AS
        EC(DeptName, Num)
       FROM EMP
        GROUP BY DeptName)
```

The output of executing the SQL query is _____.

[2017]

33. Consider a database that has the relation schemas EMP(EmpId, EmpName, DeptId), and DEPT(DeptName, DeptId), Note that the DeptId can be permitted to be NULL in the relation EMP. Consider the following queries on the database expressed in tuple relational calculus.

- (I) $\{t \mid \exists u \in \text{EMP}(t[\text{EmpName}] = u[\text{EmpName}] \wedge \forall v \in \text{DEPT}(t[\text{DeptId}] \neq v[\text{DeptId}]))\}$
- (II) $\{t \mid \exists u \in \text{EMP}(t[\text{EmpName}] = u[\text{EmpName}] \wedge \exists v \in \text{DEPT}(t[\text{DeptId}] \neq v[\text{DeptId}]))\}$
- (III) $\{t \mid \exists u \in \text{EMP}(t[\text{EmpName}] = u[\text{EmpName}] \wedge \exists v \in \text{DEPT}(t[\text{DeptId}] = v[\text{DeptId}]))\}$

Which of the above queries are safe? [2017]

- (A) (I) and (II) only
- (B) (I) and (III) only
- (C) (II) and (III) only
- (D) (I), (II) and (III)

34. Consider a database that has the relation schema CR (studentName, CourseName). An instance of the schema CR is as given below.

CR	
StudentName	CourseName
SA	CA
SA	CB
SA	CC
SB	CB
SB	CC
SC	CA
SC	CB
SC	CC
SD	CA
SD	CB
SD	CC
SD	CD
SE	CD
SE	CA
SE	CB
SF	CA
SF	CB
SF	CC

The following query is made on the database.

$T1 \leftarrow \pi_{CourseName}(\sigma_{StudentName='SA'}(CR))$

$T2 \leftarrow CR \div T1$

The number of rows in $T2$ is _____. [2017]

35. Consider the following database table named *top_scorer*:

top_scorer		
player	country	goals
Klose	Germany	16
Ronaldo	Brazil	15
G Müller	Germany	14
Fontaine	France	13
Pelé	Brazil	12
Klinsmann	Germany	11
Kocsis	Hungary	11
Batistuta	Argentina	10
Cubillas	Peru	10
Lato	Poland	10
Lineker	England	10
T Muller	Germany	10
Rahn	Germany	10

Consider the following SQL query:

```
SELECT ta.player FROM top_scorer AS ta
WHERE ta.goals > ALL (SELECT tb.goals
FROM top_scorer AS tb
WHERE tb.country = 'Spain')
AND ta.goals > ANY (SELECT tc.goals
FROM top_scorer AS tc
WHERE tc.country = 'Germany')
```

The number of tuples returned by the above SQL query is _____. [2017]

36. Consider the following two tables and four queries in SQL.

Book (*isbn*, *bname*), Stock (*isbn*, *copies*)

Query 1: SELECT B.isbn, S.copies
FROM Book B INNER JOIN Stock S
ON B.isbn = S.isbn;

Query 2: SELECT B.isbn, S.copies
FROM Book B LEFT OUTER
JOIN Stock S
ON B.isbn = S.isbn;

Query 3: SELECT B.isbn, S.copies
FROM Book B RIGHT OUTER
JOIN Stock S
ON B.isbn = S.isbn;

Query 4: SELECT B.isbn, S.copies
FROM Book B FULL OUTER
JOIN Stock S
ON B.isbn = S.isbn;

Which one of the queries above is certain to have an output that is a superset of the outputs of the other three queries? [2018]

- (A) Query 1 (B) Query 2
(C) Query 3 (D) Query 4

37. Consider the relations $r(A, B)$ and $s(B, C)$, where $s \cdot B$ is a primary key and $r \cdot B$ is a foreign key referencing $s \cdot B$. Consider the query

$Q: r \bowtie (\sigma_{B<5}(S))$

Let LOJ denote the natural left outer-join operation. Assume that r and s contain no null values.

Which one of the following queries is NOT equivalent to Q ? [2018]

- (A) $\sigma_{B<5}(r \bowtie s)$ (B) $\sigma_{B<5}(r \text{ LOJ } s)$
(C) $r \text{ LOJ } (\sigma_{B<5}(s))$ (D) $\sigma_{B<5}(r) \text{ LOJ } s$

ANSWER KEYS

EXERCISES

Practice Problems 1

1. A 2. B 3. C 4. (i) B (ii) D (iii) D 5. C 6. C 7. (i) A (ii) C
8. (i) B (ii) A 9. A 10. A 11. (i) A (ii) B 12. (i) B (ii) B (iii) A
13. (i) A (ii) A (iii) A 14. (i) A (ii) A (iii) A 15. (i) A (ii) A (iii) A 16. C
17. A 18. D 19. A 20. D

Practice Problems 2

1. B 2. D 3. D 4. A 5. C 6. A 7. A 8. A 9. A 10. A
11. C 12. C 13. C 14. C 15. C 16. B 17. C 18. D 19. B 20. A

Previous Years' Questions

1. C 2. A 3. C 4. B 5. C 6. B 7. D 8. A 9. A 10. C
11. C 12. C 13. A 14. C 15. A 16. B 17. A 18. A 19. D 20. B
21. 19 22. A 23. C 24. A 25. D 26. D 27. D 28. 2 29. C 30. A
31. B 32. 2.6 33. D 34. 4 35. 7 36. D 37. C

Normalization

LEARNING OBJECTIVES

- ☞ Normalization
- ☞ Anomalies
- ☞ First normal form
- ☞ Functional dependency
- ☞ Inference rules
- ☞ Second normal form
- ☞ Third normal form
- ☞ Higher normal forms (Boyce-Codd normal form)
- ☞ Fifth normal form
- ☞ Courses

NORMALIZATION

Database design theory includes design standards called *normal forms*. The process of making data and tables match these standards is called *normalizing data* or *data normalization*. By normalizing data, we eliminate redundant information and organize table to make it easier to manage the data and make future changes to the table and database structure. This process removes the insertion, deletion, and modification anomalies. In normalizing your data, we usually divide large tables into smaller, easier to maintain tables. We can then use the technique of adding foreign keys to enable connections between the tables.

Data normalization is part of the database design process and is neither specific nor unique to any particular RDBMS. These are in order, such as first, second, third, Boyce-Codd, fourth, and fifth normal forms. Each normal form represents an increasingly stringent set of rules; that is, each normal form assumes that the

requirements of the preceding forms have been met. Many relational database designers feel that, if their tables are in third normal form, most common design problems have been addressed. However, the higher-level normal forms can be of use and are included here.

Database normalization is the process of removing redundant data from tables to improve storage efficiency, data integrity and scalability.

1. In the relational model, methods exists for quantifying how efficient a database is, these classifications are called q' .
2. Normalization generally involves splitting existing tables into multiple ones, which must be rejoined (or) linked each time a query is issued.
3. Edgar F. Codd originally established three normal forms: 1NF, 2NF, 3NF. There are others also, but 3NF is widely considered to be sufficient for most applications, most tables when reaching 3NF are also in BCNF (Boyce-Codd normal form).

Table 1

Title	Author 1	Author 2	I SBN	Subject	Pages	Publisher
Database system concepts	Abraham Silber schatz	Henry F. Korth	0072958863	My SQL, computers	1160	McGraw-Hill
OS concepts	Abraham Silberschatz	Henry F. Korth	0471694665	Computers	990	McGraw-Hill

Problems:

1. This table is not very efficient with storage.
2. This design doesn't protect data integrity.
3. This table doesn't scale well.

Anomalies

An anomaly is a variation that differs in some way from what is said to be normal, with respect to maintaining a database.

1. The basic operations performed on Databases are Record insertion, Record updation, Record deletion.

2. It is desirable for these operations to be straight forward and efficient.
3. When relations are not fully normalized they exhibit anomalies.
4. The design goal of database is too easily to understand and to maintain.
5. Anomalies are problems that occur in un-normalized databases where all the data is stored in one table.

Types of anomalies

There are three types of anomalies that can arise in the database because of redundancy as follows:

1. Insertion anomaly
2. Deletion anomaly
3. Updation anomaly

Insertion anomaly An insertion anomaly occurs when particular attributes cannot be inserted into the database without the presence of other attributes.

Example: Consider the following table: Sales

Sales-Rep-Id	Name	Hire-Date	Client
1	Ana	1/1/2015	Madison
2	Sudha	2/4/2014	Peterson
3	Joey	3/2/2014	John
* New			

Insertion anomaly occurs in the above table which stores records for a company’s sales representatives and the clients for whom they are responsible.

1. It is not possible to add records for newly hired Sales representatives until they have been assigned to one or more clients.
2. If we insert a record for newly hired, client column will be NULL, which is a required field for the table.
3. It is not possible to record newly hired in the table during training.

Deletion anomaly Deletion anomaly occurs when some particular attributes are lost because of the deletion of other attributes.

Example: Consider the following table ‘course’.

S No	C No	S Name	Course
S41	C9201	John	Sales
S42	C9401	Brat	Finance
S40	C9201	Amit	Sales
S43	C9608	Arun	Accounts

Execute the following SQL query:

```
Delete *
From course
Where S No = S43
```

If we delete a tuple where SNo = S43, he is the only (or) last student in the accounts department, we will lose data about student ‘S43, Arun’ as well as data about Accounts course that is ‘C9608, Accounts’.

Updation anomaly An updation anomaly occurs when one or more instances of duplicated data are updated but not all.

Example: Consider the ‘course’ table given in the above example.

If we want to update course – No (Cno) of sales C9201 to C8686, in the course table.

1. It might happen that, the tuple with S No = S41 updated its CNo to C8686, but not the tuple with SNo = S43.
2. Inconsistency occurs in the table, because for the same course sales we have 2 different course Numbers.

Determining keys For a table ‘R’, its schema R consists of all attributes of R, we say X is a key to R if $X \rightarrow R$ means

X determines R
R is dependent upon X
If you know x then you know R

Example: Consider a relation schema R(ABCDE) and the functional dependencies:

- AC → D
- B → E
- DA → B

The closure of AC determines all the attributes present in Relation R, so the key for R is ‘AC’.

$$\begin{aligned}
 AC^+ &= \{AC\} && \text{(self determination)} \\
 &\{ACD\} && (AC \rightarrow D) \\
 &\{ACDB\} && (DA \rightarrow B) \\
 &\{ACDBE\} && (B \rightarrow E) \\
 \therefore \text{key} &= AC
 \end{aligned}$$

Any attribute which does not appear on the right-hand-side of a given functional dependency appears in any one of the candidate keys.

1. From the above example, neither A (or) C appears in the right hand side of any functional dependency.

FIRST NORMAL FORM (1NF)

In Table 1, we have two violations of 1NF such as:

1. More than one author field and
2. Subject field contains more than one piece of information with more than one value in a single field; thus, it would be very difficult to search for all books on a given subject.

Table 2 1NF table

Title	Author	ISBN	Subject	Pages	Publisher
Database system concept	Abraham Silbers Chatt	0072958863	My SQL	1160	McGraw-Hill
Database system concept	Henry K. Forth	0072958863	Computers	1160	McGraw-Hill
OS concepts	Henry K. Forth	0471694665	Computers	990	McGraw-Hill
OD concepts	Abraham Silber Schatz	0471694665	Computers	990	McGraw-Hill

In Table 2, we have two rows for a single book. Additionally, we would be violating the second NF. A better solution to the problem would be to separate the data into separate tables—an author table and a subject table to store our information, removing that information from the book table.

Table 3 Subject table

Subject-ID	Subject
1	My SQL
2	computers

Table 4 Author table

Author-ID	Last Name	First name
1	Silberschatz	Abraham
2	Korth	Henry

Table 5 Book table

ISBN	Title	Pages	Publisher
0072958863	Database System Concepts	1160	McGraw-Hill
0471694665	OS concepts	990	McGraw-Hill

Each table has a primary key, used for joining tables together when querying the data.

A table is in first normal form (1NF) if there are no repeating groups. A repeating group is a set of logically related fields or values that occur multiple times in one record. The sample tables below do not comply with first normal form. Look for fields that contain too much data and repeating group of fields.

EMPLOYEES_PROJECTS_TIME

A table with fields containing too much data.

Employee ID	Name	Project	Time
EN1-26	Sean O'Brien	30-452-T3, 30-457-T3, 32-244-T3	0.25, 0.40, 0.30
EN1-33	Amy Guya	30-452-T3, 30-382-TC, 32-244-T3	0.05, 0.35, 0.60
EN1-35	Steven Baranco	30-452-T3, 31-238-TC	0.15, 0.80
EN1-36	Elizabeth Roslyn	35-152-TC	0.90
EN1-38	Carol Schaaf	36-272-TC	0.75
EN1-40	Alexandra Wing	31-238-TC, 31-241-TC	0.20, 0.70

The example above is also related to another design issue, namely, that each field should hold the smallest meaningful value and that there should not be multiple values in a single field.

Why is this table design a problem?

There would be no way to sort by last names or to know which allocation of time belonged to which project.

EMPLOYEES_PROJECTS_TIME

Table 5 A table with repeating groups of fields.

Emp ID	Last Name	First Name	Project1	Time1	Project2	Time2	Project3	Time3
EN1-26	O'Brien	Sean	30-452-T3	0.25	30-457-T3	0.40	32-244-T3	0.30
EN1-33	Guya	Amy	30-452-T3	0.05	30-382-TC	0.35	32-244-T3	0.60
EN1-35	Baranco	Steven	30-452-T3	0.15	31-238-TC	0.80		
EN1-36	Roslyn	Elizabeth	35-152-TC	0.90				
EN1-38	Schaaf	Carol	36-272-TC	0.75				
EN1-40	Wing	Alexandra	31-238-TC	0.20	31-241-TC	0.70		

If an employee was assigned to a fourth project, you would have to add two new fields to the table. Also, it would be very difficult to total the amount of time devoted to a particular project.

The design problems addressed are very common, particularly among new designers who are accustomed to tracking data in a spreadsheet. Often, when building a spreadsheet, we arrange the data horizontally, laying it out across the spreadsheet. When designing tables, we have to think more vertically. Similar data belongs in the same column or field with a single value in each row.

Now we will take the table you saw above and redesign it so it will comply with first normal form.

Look at the repeating groups of data. Identify tables and fields that will hold this data without the repeating groups. Think vertically and remember that similar data belongs in the same field.

Enter the sample data from the table to make sure you don't have repeating groups. If necessary, include foreign key field(s) to connect the tables.

EMPLOYEES

EmployeeID	Last Name	First Name
EN1-26	O'Brien	Sean
EN1-33	Guya	Amy
EN1-35	Baranco	Steven
EN1-36	Roslyn	Elizabeth
EN1-38	Schaaf	Carol
EN1-40	Wing	Alexandra

PROJECTS_EMPLOYEES_TIME

Project Num	EmployeeID	Time
30-328-TC	EN1-33	0.35
30-452-T3	EN1-26	0.25
30-452-T3	EN1-33	0.05
30-452-T3	EN1-35	0.15
31-238-TC	EN1-35	0.80
30-457-T3	EN1-26	0.40
31-238-TC	EN1-40	0.20
31-241-TC	EN1-40	0.70
32-244-T3	EN1-33	0.60
35-152-TC	EN1-36	0.90
36-272-TC	EN1-38	0.75

Mark the primary key field(s) and foreign keys in each table. Shown below with * indicating the Primary key.

EMPLOYEES

EmployeeID	Last Name	First Name
EN1-26	O'Brien	Sean
EN1-33	Guya	Amy
EN1-35	Baranco	Steven
EN1-36	Roslyn	Elizabeth
EN1-38	Schaaf	Carol
EN1-40	Wing	Alexandra

PROJECTS_EMPLOYEES_TIME

Project Num	EmployeeID	Time
30-328-TC	EN1-33	0.35
30-452-T3	EN1-26	0.25
30-452-T3	EN1-33	0.05
30-452-T3	EN1-35	0.15
31-238-TC	EN1-35	0.80
30-457-T3	EN1-26	0.40
31-238-TC	EN1-40	0.20
31-241-TC	EN1-40	0.70
32-244-T3	EN1-33	0.60
35-152-TC	EN1-36	0.90
36-272-TC	EN1-38	0.75

If an employee was assigned to an additional project, it would involve merely adding a new record. Also, it would be much easier to search for a particular project number as they are all held in a single column.

Functional Dependency

A functional dependency is a relationship between fields so that the value in Field *A* determines the value in Field *B*, and there can be only one value in Field *B*. In that case, Field *B* is functionally dependent on Field *A*. Consider the following sample table:

Airport	City
National	Washington, DC
JFK	New York
LaGuardia	New York
Logan	Boston
Dulles	Washington, DC

Each airport name is unique and each airport can be in only one city. Therefore, City is functionally dependent on Airport. The value in the Airport field determines what the value will be in the City field (making Airport the determinant field) and there can be only one value in the City field. This does not need to work in the reverse. As shown in the table, a city can have more than one airport, so Airport is not functionally

dependent on City; the value in City does not necessarily determine what the value in Airport will be.

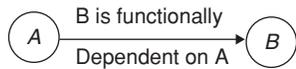
You will sometimes see a functional dependency written in this format:

Determinant field(s) → Functionally dependent field as in:

Airport → City

Functional dependency describes the relationship between attributes in a relation.

Example: If *A* and *B* are attributes of relation *R*, and *B* is functionally dependent on *A* ($A \rightarrow B$) if each value of *A* is associated with one value of *B*.



Determinant refers to the attributed (or) group attributes on the left-hand side of the arrow of a functional dependency.

Inference Rules

The following inference rules IR 1 through IR 6 form a complete set for inferring functional and multi-valued dependencies from a given set of dependencies

Assume that all attributes are included in a ‘universal’ relation schema $R = \{A_1, A_2, \dots, A_N\}$ and that *X*, *Y*, *Z* and *W* are subsets of *R*.

- IR 1 (reflexive rule): if $X \supseteq Y$, then $X \rightarrow Y$
- IR 2 (Augmentation rule): $\{X \rightarrow Y\} = XZ \rightarrow YZ$
- IR 3 (transitive rule): $\{X \rightarrow Y, Y \rightarrow Z\} = X \rightarrow Z$
- IR 4 (complementation rule): $\{X \rightarrow \rightarrow Y\} = \{X \rightarrow \rightarrow (R - (X \cup Y))\}$
- IR5 (augmentation rule for MVD's): if $X \rightarrow \rightarrow y$ and $W \rightarrow Z$ then $WX \supseteq YZ$
- IR6 (transitive rule for MVD's): $\{X \rightarrow \rightarrow Y, Y \rightarrow \rightarrow Z\} = X \rightarrow \rightarrow (Z - Y)$

SECOND NORMAL FORM

A table is said to be in second normal form if it is in first normal form and each non-key field is functionally dependent on the entire primary key.

Look for values that occur multiple times in a non-key field. This tells us that we have too many fields in a single table.

Example: In the example below, see all the repeating values in the name and Project Title fields. This is an inefficient way to store and maintain data. In a well-designed database, the only data that is duplicated is in key fields used to connect tables. The presumption is that the data in key fields will rarely change, while the data in non-key fields may change frequently.

A table with a multifield primary key and repeating data in non-key fields

EmployeeID	Last Name	First Name	Project Number	Project Title
EN1-26	O'Brien	Sean	30-452-T3	STAR manual
EN1-26	O'Brien	Sean	30-457-T3	ISO procedures
EN1-26	O'Brien	Sean	31-124-T3	Employee handbook
EN1-33	Guya	Amy	30-452-T3	STAR manual
EN1-33	Guya	Amy	30-482-TC	Web Site
EN1-33	Guya	Amy	31-241-TC	New catalogue
EN1-35	Baranco	Steven	30-452-T3	STAR manual
EN1-35	Baranco	Steven	31-238-TC	STAR prototype
EN1-36	Roslyn	Elizabeth	35-152-TC	STAR pricing
EN1-38	Schaaf	Carol	36-272-TC	Order system
EN1-40	Wing	Alexandra	31-238-TC	STAR prototype
EN1-40	Wing	Alexandra	31-241-TC	New catalogue

If a ProjectTitle changed, we would have to edit it in several records. And what would happen in this table if the EmployeeID was part of the primary key and we wanted to add a new ProjectNum and ProjectTitle even though no employees had yet been assigned?

The primary key cannot contain a null value so you couldn't add the new project. Additionally, if a project ended and you wanted to delete it, you would have to delete the individual values because, if we deleted the records

containing the titles and an employee was assigned to only that project, you would also delete that employee's record, something that we may not want to do.

In the above example, the asterisks indicate the fields that make up the primary key of this table as it now stands. A multifield primary key is necessary because neither the EmployeeID nor the ProjectNum fields contain unique values.

The reason there are repeated values in LastName, FirstName, and ProjectTitle is that these fields are dependent

on only part of the primary key. The value in EmployeeID determines what the value in LastName will be, but the value in ProjectNum has nothing to do with it. Similarly, the value in ProjectNum determines the value in ProjectTitle, but EmployeeID does not. These non-key fields relate to only part of the primary key. They are not functionally dependent on the entire primary key.

The solution to this lies in breaking the table into smaller tables that do meet second normal form. You will find that more tables are the solution to most problems encountered during data normalisation.

EMPLOYEES

EmployeeID	Last Name	First Name
EN1-26	O'Brien	Sean
EN1-33	Guya	Amy
EN1-35	Baranco	Steven
EN1-36	Roslyn	Elizabeth
EN1-38	Schaaf	Carol
EN1-40	Wing	Alexandra

EMPLOYEES_PROJECTS

EmployeeID	Project Num
EN1-26	30-452-T3
EN1-26	30-457-T3
EN1-26	31-124-T3
EN1-33	30-328-TC
EN1-33	30-452-T3
EN1-33	32-244-T3
EN1-35	30-452-T3
EN1-35	31-238-TC
EN1-36	35-152-TC
EN1-38	36-272-TC
EN1-40	31-238-TC
EN1-40	31-241-TC

Now we'll take the table above and design new tables that will eliminate the repeated data in the non-key fields.

1. To decide what fields belong together in a table, think about which field determines the values in other fields. Create a table for those fields and enter the sample data.
2. Think about what the primary key for each table would be and about the relationship between the tables. If necessary, add foreign keys or a junction table.
3. Mark the primary key for each table and make sure that you don't have repeating data in non-key fields.

PROJECTS

Project Num	Project Title
30-452-T3	STAR manual
30-457-T3	ISO procedures
30-482-TC	Web site
31-124-T3	Employee handbook
31-238-TC	STAR prototype
31-238-TC	New catalog
35-152-TC	STAR pricing
36-272-TC	Order system

Examine the tables to make sure there are no repeating values in non-key fields and that the value in each non-key field is determined by the value(s) in the key field(s). This removes the modification anomaly of having the repeated values.

THIRD NORMAL FORM

A table is said to be in third normal form if it is in second normal form (2NF) and there are no transitive dependencies.

A transitive dependency is a type of functional dependency in which the value in a non-key field is determined by the value in another non-key field and that field is not a candidate key. Again, look for repeated values in a non-key field as in the following example.

A table with a single field primary key and repeating values in non-key fields.

Project Num	Project Title	Project Mgr	Phone
30-452-T3	STAR manual	Garrison	2756
30-457-T3	ISO procedures	Jacanda	2954
30-482-TC	Web site	Friedman	2846
31-124-T3	Employee handbook	Jones	3102
31-238-TC	STAR prototype	Garrison	2756
31-241-TC	New catalog	Jones	3102
35-152-TC	STAR pricing	Vance	3022
36-272-TC	Order system	Jacanda	2954

The phone number is repeated each time a manager's name is repeated. It is dependent on the manager, which is dependent on the project number (a transitive dependency).

The Project Manager field is not a candidate key, because the same person manages more than one project. Again, the solution is to remove the field with repeating data to a separate table.

Take the above table and create new tables to fix the problem.

1. Think about which fields belong together and create new tables to hold them.
2. Enter the sample data and check for unnecessarily (not part of primary key) repeated values.
3. Identify the primary key for each table and, if necessary, add foreign keys.

PROJECTS

Project Num	Project Title	Project Mgr
30-452-T3	STAR manual	Garrison
30-457-T3	ISO procedures	Jacanda
30-482-TC	Web site	Friedman
31-124-T3	Employee handbook	Jones
31-238-TC	STAR prototype	Garrison
31-241-TC	New catalog	Jones
35-152-TC	STAR pricing	Vance
36-272-TC	Order system	Jacanda

MANAGERS

Project Manager	Phone
Friedman	2846
Garrison	2756
Jacanda	2954
Jones	3102
Vance	3022

Reexamine your tables to make sure there are no unnecessarily repeating values in non-key fields and that the value in each non-key field is determined by the value(s) in the key field(s). In most cases, 3NF should be sufficient to ensure that your database is properly normalised.

HIGHER NORMAL FORMS (BOYCE–CODD NORMAL FORM)

A table is in third normal form (3NF), and all determinants are candidate keys.

Boyce–Codd normal form (BCNF) can be thought of as a 'new' third normal form. It was introduced to cover situations that the 'old' third normal form did not address. The mean of a determinant (determines the value in another field) and candidate keys (qualify for designation as primary

key). This normal form applies to situations where you have overlapping candidate keys.

If a table has no non-key fields, it is automatically in BCNF (Figure 1). Look for potential problems in updating existing data (modification anomaly) and in entering new data (insertion anomaly).

Imagine that we were designing a table for a college to hold information about courses, students, and teaching assistants. We have the following business rules:

1. Each course can have many students.
2. Each student can take many courses.
3. Each course can have multiple teaching assistants (TAs).
4. Each TA is associated with only one course.
5. For each course, each student has one TA.

Some sample data:

COURSES_STUDENTS_TA's

CourseNum	Student	TA
ENG101	Jones	Clark
ENG101	Grayson	Chen
ENG101	Samara	Chen
MAT350	Grayson	Powers
MAT350	Jones	O'Shea
MAT350	Berg	Powers

To uniquely identify each record, we could choose CourseNum + Student as a primary key. This would satisfy third normal form also because the combination of CourseNum and Student determines the value in TA. Another candidate key would be Student + TA. In this case, you have overlapping candidate keys (Student is in both). The second choice, however, would not comply with third normal form, because the CourseNum is not determined by the combination of Student and TA; it only depends on the value in TA. This is the situation that Boyce–Codd normal form addresses; the combination of Student + TA could not be considered to be a candidate key.

If we wanted to assign a TA to a course before any students enrolled, we couldn't because Student is part of the primary key. Also, if the name of a TA changed, would have to update it in multiple records. If assume have just these fields, this data would be better stored in three tables: one with CourseNum and Student, another with Student and TA, and third with CourseNum and TA.

COURSES

Course Num	Student
ENG101	Jones
ENG101	Grayson
ENG101	Samara
MAT350	Grayson
MAT350	Jones
MAT350	Berg

STUDENTS

Student	TA
Jones	Clark
Grayson	Chen
Samara	Chen
Grayson	Powers
Jones	O'Shea
Berg	Powers

TA's

FOURTH NORMAL FORM

A table is in Boyce-Codd normal form (BCNF) and there are no multi-valued dependencies.

A *multi-valued dependency* occurs when, for each value in field *A*, there is a set of values for field *B* and a set of values for field *C* but fields *B* and *C* are not related.

Look for repeated or null values in non-key fields. A multi-valued dependency occurs when the table contains fields that are not logically related. An often used example is the following table:

MOVIES

Movie	Star	Producer
Once Upon a Time	Julie Garland	Alfred Brown
Once Upon a Time	Mickey Rooney	Alfred Brown
Once Upon a Time	Julie Garland	Muriel Humphreys
Once Upon a Time	Mickey Rooney	Muriel Humphreys
Moonlight	Humphrey Bogart	Alfred Brown
Moonlight	Julie Garland	Alfred Brown

A movie can have more than one star and more than one producer. A star can be in more than one movie. A producer can produce more than one movie. The primary key would have to include all three fields, and so this table would be in BCNF. But you have unnecessarily repeated values, with the

PROJECTS_EQUIPMENT

Dept Code	Project Num	Project Mgr ID	Equipment	Property ID
IS	36-272-TC	EN1-15	CD-ROM	657
IS			VGA desktop monitor	305
AC	35-152-TC	EN1-15		
AC			Dot-matrix printer	358
AC			Calculator with tape	239
TW	30-452-T3	EN1-10	486 PC	275
TW	30-457-T3	EN1-15		
TW	31-124-T3	EN1-15	Laser printer	109
TW	31-238-TC	EN1-15	Handheld scanner	479
RI			Fax machine	775
MK			Laser printer	858
MK			Answering machine	187
TW	31-241-TC	EN1-15	Standard 19200 bps modem	386
SL			486 Laptop PC	772
SL			Electronic notebook	458

*CourseNum	*TA
ENG101	Clark
ENG101	Chen
MAT350	O'Shea
MAT350	Powers

Figure 1 Tables that comply with BCNF.

data maintenance problems that causes and you would have trouble with deletion anomalies.

The Star and the Producer really aren't logically related. The Movie determines the Star and the Movie determines the Producer. The answer is to have a separate table for each of those logical relationships: one holding Movie and Star and the other with Movie and Producer, as shown below:

STARS

*Movie	*Star
Once Upon a Time	Julie Garland
Once Upon a Time	Mickey Rooney
Moonlight	Humphrey Bogart
Moonlight	Julie Garland

PRODUCERS

*Movie	*Producer
Once Upon a Time	Alfred Brown
Once Upon a Time	Muriel Humphreys
Moonlight	Alfred Brown

Above, showing tables that comply with 4NF

Below is another example of a common design error, and it's easily spotted by all the missing or blank values.

A table with many null values (Note: It also does not comply with 3NF and BCNF).

It is the same problem here because not all of the data is logically related. As usual, the answer is more tables: one to hold the information on the equipment assigned to departments (with PropertyID as the primary key) and another with projects and departments. We would now the business rules to know whether a project might involve more than one department or manager and be able to figure out the primary key. Assuming a project can have only one manager and be associated with only one department, the tables would be as follows:

EQUIPMENT

*Property ID	Equipment	DeptCode
657	CD-ROM	IS
305	VGA desktop monitor	IS
358	Dot-matrix printer	AC
239	Calculator with tape	AC
275	486 PC	TW
109	Laser printer	TW
479	Handheld scanner	TW
775	Fax machine	RI
858	Laser printer	MK
187	Answering machine	MK
386	Standard 19200 bps modem	TW
772	486 Laptop PC	SL
458	Electronic notebook	SL

PROJECTS_EQUIPMENT

Project Num	Project Mgr ID	Dept Code
36-272-TC	EN1-15	IS
35-152-TC	EN1-15	AC
30-452-T3	EN1-10	TW
30-457-T3	EN1-15	TW
31-124-T3	EN1-15	TW
31-238-TC	EN1-15	TW
31-241-TC	EN1-15	TW

Figure 2 Tables that eliminate the null values and comply with 4NF.

FIFTH NORMAL FORM

A table is in fourth normal form (4 NF) and there are no cyclic dependencies.

A *cyclic dependency* can occur only when you have a multifiend primary key consisting of three or more fields. For example, let's say your primary key consists of fields *A*, *B*, and *C*. A cyclic dependency would arise if the values in those fields were related in pairs of *A* and *B*, *B* and *C*, and *A* and *C*.

Fifth normal form is also called *projection-join normal form*. A *projection* is a new table holding a subset of fields from an original table. When properly formed projections are joined, they must result in the same set of data that was contained in the original table.

Look for the number of records that will have to be added or maintained

Following is some sample data about buyers, the products they buy, and the companies they buy from.

BUYING

Buyer	Product	Company
Chris	Jeans	Levi
Chris	Jeans	Wrangler
Chris	Shirts	Levi
Lori	Jeans	Levi

Figure 3 A table with cyclic dependencies.

The primary key consists of all three fields. One data maintenance problem that occurs is that you need to add a record for every buyer who buys a product for every company that makes that product or they can't buy from them. That may not appear to be a big deal in this sample of two buyers, two products, and two companies ($2 \times 2 \times 2 = 8$ total records). But what if we went to 20 buyers, 50 products, and 100 companies ($20 \times 50 \times 100 = 100,000$ potential records)? It quickly gets out of hand and becomes impossible to maintain.

We might solve this by dividing this into the following two tables:

BUYERS

Buyer	Product
Chris	jeans
Chris	shirts
Lori	jeans

PRODUCTS

Product	Company
jeans	Wrangler
jeans	Levi
shirts	Levi

However, if you joined the two tables above on the Product field, it would produce a record not part of the original data set (it would say that Lori buys jeans from Wrangler). This is where the projection-join concept comes in.

The correct solution would be three tables:

BUYERS

*Buyer	*Product
Chris	jeans
Chris	shirts
Lori	jeans

PRODUCTS

*Product	*Company
jeans	Wrangler
jeans	Levi
shirts	Levi

COMPANIES

*Buyer	*Company
Chris	Levi
Chris	Wrangler
Lori	Levi

Figure 4 Tables that comply with 5NF.

When the first two tables are joined by Product and the result joined to the third table by Buyer and Company, the result is the original set of data.

EXERCISES

Practice Problems I

Directions for questions 1 to 20: Select the correct alternative from the given choices.

1. Consider the given functional dependencies

$A \rightarrow B$
 $BC \rightarrow DE$
 $AEF \rightarrow G$

Which of the following is true?

- (A) Functional dependency $ACF \rightarrow DG$ implied by the set
- (B) Functional dependency $ACF \rightarrow DG$ cannot be implied by the set
- (C) Functional dependency $AB \rightarrow G$ implied by the set
- (D) Both (B) and (C)

2. Consider the given relation

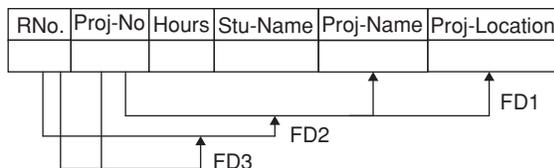
DNAME	DNO	MGRNO	LOCATION
RESEARCH	5	333	{BANGLORE,DELHI, HYDERABAD}
ADMINISTRATION	4	987	{CHENNAI}
EXECUTIVES	1	885	{HYDERABAD}

Department
 The given relation is

- (A) is not in 1NF
- (B) in 1NF
- (C) in 2NF
- (D) in 3NF

3. Consider the given Relational scheme

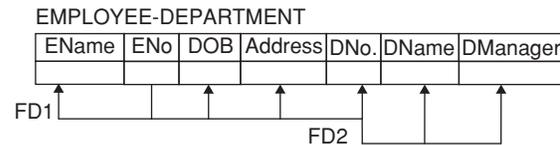
Student-project



Which functional dependencies are violating 2NF property?

- (A) FD1
- (B) FD2
- (C) FD3
- (D) Both A and B

4. Consider the given relation



Which functional dependencies are violating 3NF?

- (A) FD1
- (B) FD2
- (C) Both
- (D) None of these

5. Consider the given relation $R(A, B, C, D)$ and functional dependencies:

$FD = (AB \rightarrow C$
 $C \rightarrow B$
 $C \rightarrow D)$

Determine the key, prime attributes and non-prime attributes.

- (A) $\{A\}, \{AB\}, \{CDE\}$
- (B) $\{AB, AC\}, \{ABC\}, \{D\}$
- (C) $\{AB, BC\}, \{ABC\}, \{D\}$
- (D) $\{AB, AC\}, \{AB\}, \{D\}$

6. Consider the given relation and functional dependencies

$R(ABCDE)$
 $FD = (ABD \rightarrow C$
 $BC \rightarrow D$
 $CD \rightarrow E)$

Determine the key, prime attributes, non-prime attributes and the normal form of the relation?

- (A) $\{AB, AD\}, \{ABCD\}, \{E\}$
- (B) $\{ABC, ABD\}, \{ABCD\}, \{E\}$
- (C) $\{AB, AD\}, \{ABC\}, \{DE\}$
- (D) $\{ABC, ABD\}, \{AB\}, \{CDE\}$

7. Consider the given relation and functional dependencies

$R(ABC)$
 $FD = (AB \rightarrow C$
 $C \rightarrow A)$

The relation is in which normal form?

- (A) 1NF
- (B) 2NF
- (C) 3NF
- (D) BCNF

8. Consider the given relation and its functional dependencies:

$R(ABCDE)$
 $FD = (AB \rightarrow C$
 $C \rightarrow E$
 $B \rightarrow D$
 $E \rightarrow A)$

The relation is further decomposed into two relations:

$R_1(BCD), R_2(ACE)$

- (A) Decomposition is lossy and dependency preserving
- (B) Decomposition is lossless and dependency preserving
- (C) Decomposition is lossy and not dependency preserving
- (D) Decomposition is lossless and not dependency preserving

9. Consider the following relational instance:

X	Y	Z
1	4	2
1	5	3
1	6	3
3	2	2

Which of the following functional dependencies are satisfied by the instance?

- (A) $xy \rightarrow z$ and $z \rightarrow y$
- (B) $yz \rightarrow x$ and $y \rightarrow z$
- (C) $yz \rightarrow x$ and $x \rightarrow z$
- (D) $xz \rightarrow y$ and $y \rightarrow x$

10. Consider the following functional dependencies:

$DOB \rightarrow Age$
 $Age \rightarrow Eligibility$
 $Name \rightarrow RNo$
 $RNo \rightarrow Name$
 $CourseNo \rightarrow CourseName$
 $CourseNo \rightarrow Instructor$
 $(RNo, CNo) \rightarrow Grade$

The relation (RNo, Name, DOB, Age) is in which normal form?

- (A) 1NF
- (B) 2NF
- (C) 3NF
- (D) BCNF

11. Consider the given functional dependencies:

$AB \rightarrow CD$

$AF \rightarrow D$
 $DE \rightarrow F$
 $C \rightarrow G$
 $F \rightarrow E$
 $G \rightarrow A$

Which of the following is false?

- (A) $\{CF\}^+ = \{ACDEFG\}$
- (B) $\{BG\}^+ = \{ABCDG\}$
- (C) $\{AF\}^+ = \{ACDEFG\}$
- (D) $\{AB\}^+ = \{ABCDG\}$

12. What should be the key to make the given relation to be in BCNF? The dependencies for the following, 'Grades' relation are GRADES (student-Id, course#, semester#, Grade) student-Id, course#, semester# \rightarrow Grade

- (A) student-Id
- (B) course#
- (C) semester#
- (D) student-Id, course#, semester #

13. What normal form is the following relation in?

STORE_ITEM (SKU, promotionID, vendor, style, price)

$SKU, promotionID \rightarrow vendor, style, price$

$SKU \rightarrow vendor, style$

- (A) 1NF
- (B) 2NF
- (C) 3NF
- (D) 4NF

14. What normal form is the following relation in?

Only H, I can act as the key
 STUFF (H, I, J, K, L, M, N, O)

$H, I \rightarrow J, K, L$

$J \rightarrow M$

$K \rightarrow N$

$L \rightarrow O$?

- (A) 1NF
- (B) 2NF
- (C) 3NF
- (D) BCNF

15. What normal form the following relation is in?

STUFF2(D, O, N, T, C, R, Y)

$D, O \rightarrow N, T, C, R, Y$

$C, R \rightarrow D$

$D \rightarrow N$?

- (A) 1NF
- (B) 2NF
- (C) 3NF
- (D) BCNF

16. The given table is in the BCNF form, convert it to the 4th normal form.

Employee	Skill	Language
Jones	Electrical	French
Jones	Electrical	German
Jones	Mechanical	French
Jones	Mechanical	German
Smith	Plumbing	Spanish

(A)

Employee	Skill

(B)

Employee	Language

(C)

Skill	Language

(D) Both A and B

17. For a database relation $x(a, b, c, d)$, where all the domains of a, b, c, d , include only atomic values, only the following FDs and those that can be inferred from them hold.

$a \rightarrow b, c \rightarrow d$

the relation is

- (A) In 1st NF but not in 2nd NF
- (B) In 2nd NF but not in 3rd NF
- (C) In 2nd NF
- (D) In 3rd NF

18. Which of the following FDs are satisfied by the instance from the below relation:

A	B	C
2	8	4
2	10	6
2	12	6
6	4	4

- (A) $AB \rightarrow C$ and $C \rightarrow B$
- (B) $BC \rightarrow A$ and $B \rightarrow C$
- (C) $BC \rightarrow A$ and $A \rightarrow C$
- (D) $AC \rightarrow B$ and $B \rightarrow A$

19. Consider the following database:

Course # \rightarrow Title

Course # time \rightarrow location

Emp – ID \rightarrow T –Name salary
is in

- (A) 3NF
- (B) 2NF
- (C) 1NF
- (D) BCNF

20. Consider the following schema

$A = (w, x, y, z)$ and the dependencies are

$W \rightarrow X, X \rightarrow Y, Y \rightarrow Z$, and $Z \rightarrow W$

Let $A = (A_1$ and $A_2)$ be a decomposition such that $A_1 \cap A_2 = \phi$

The decomposition is

- (A) In 1NF and in 2NF
- (B) In 2NF and not in 3NF
- (C) In 2NF and in 3NF
- (D) Not in 2NF and in 3NF

Practice Problems 2

Directions for questions 1 to 20: Select the correct alternative from the given choices.

1. Integrity constraints ensures that changes made to the database by authorized users do not result in
 - (A) Loss of FDs
 - (B) Loss of keys
 - (C) Loss of tables
 - (D) Loss of data consistency
2. Relation $R = (\overline{A}, B, C, D)$ with AB as primary key. Choose one FD such that R should be in 1NF but not in 2NF
 - (A) $AB \rightarrow C$
 - (B) $AB \rightarrow D$
 - (C) $A \rightarrow D$
 - (D) $AB \rightarrow CD$
3. A normalized relation (1NF) can be retrieved from unnormalized relation by removing
 - (A) repeating groups
 - (B) duplicate tuples
 - (C) transitive dependency
 - (D) primary key

4. A relation will be in 2NF, if we
 - (A) remove repeating groups
 - (B) remove partial dependency
 - (C) remove transitive dependency
 - (D) have overlapping candidate key
5. Relation $R = (A, B, C, D)$ with AB as primary key, choose the FD so that R should be in 2NF but not in 3NF.
 - (A) $D \rightarrow C$
 - (B) $AB \rightarrow C$
 - (C) $AB \rightarrow D$
 - (D) $A \rightarrow B$
6. If a relation is in 2NF, then it can be in 3NF by removing
 - (A) repeating groups
 - (B) partial dependencies
 - (C) transitive dependencies
 - (D) overlapping dependencies
7. BCNF can be achieved from 3NF by removing
 - (A) repeating groups
 - (B) partial dependencies
 - (C) transitive dependencies
 - (D) overlapping dependencies

8. Which one of the following is not possible?
 (A) Relation is in BCNF but not in 4NF
 (B) Relation is in 3NF but not in BCNF
 (C) Relation is in 2NF but not in 3NF
 (D) Relation is in 3NF but not in 2NF

Common data for questions 9 and 10: Let R be a relation schema $R(A, B, C, D)$;

$F = \{AB \rightarrow CD; C \rightarrow A\}$ F is the set of functional dependencies

9. How many prime attributes are there?
 (A) 1 (B) 2
 (C) 3 (D) 4
10. The highest normal form of the above relation is
 (A) 1NF (B) 2NF
 (C) 3NF (D) 4NF

Linked answer questions

11. For a given relation schema $R = \{A, B, C, D, E\}$

$A \rightarrow BC$

$CD \rightarrow E$

$B \rightarrow D$

$E \rightarrow A$

Which of the following is not a candidate key?

- (A) A (B) B
 (C) E (D) BC

12. For the above answer, what is the closure?

- (A) BD (B) ABC
 (C) $ABCDE$ (D) BC

13. Consider the following functional dependencies:

$A \rightarrow B$

$C \rightarrow D$

$B \rightarrow E$

$F \rightarrow A$

The relation (A, B, C, D) is

- (A) in second normal form, but not in third normal form
 (B) in third normal form, but not in $BCNF$
 (C) in $BCNF$
 (D) None of the above

Common data for questions 14 and 15:

$R = (A, B, C, D, E, F)$

FDs = $A \rightarrow B$

$C \rightarrow DF$

$AC \rightarrow E$

$D \rightarrow F$

14. Determine the key from the given FDs:

- (A) AB (B) AC
 (C) ACB (D) ACD

15. Decompose the FDs into 2NF

- (A) $R_1(AB) R_2(CDF) R_3(ACE)$
 (B) $R_1(AB) R_2(CDEF)$
 (C) $R_1(ABC) R_2(CDF)$
 (D) $R_1(AB) R_2(CD) R_3(EF)$

16. For a database relation $x(a, b, c, d)$, where all the domains of a, b, c, d , include only atomic values, only the following FDs and those that can be inferred from them hold.

$a \rightarrow b, c \rightarrow d$

The relation is decomposed into $R_1(ab), R_2(cd)$. Which of the following is true, The decomposition

- (A) is dependency preserving
 (B) is not dependency preserving
 (C) is loss less
 (D) Both A and C

17. Which of the following FDs are satisfied by the instance from the below relation?

A	B	C
4	12	8
4	14	10
4	16	10
10	8	8

- (A) $AB \rightarrow C$ and $C \rightarrow B$
 (B) $BC \rightarrow A$ and $B \rightarrow C$
 (C) $BC \rightarrow A$ and $A \rightarrow C$
 (D) $AC \rightarrow B$ and $B \rightarrow A$

18. Indicate which of the following statements are false: 'A relational database, which is in 3NF still have undesirable data redundancy because there may exist.

- (A) Below all
 (B) Non trivial FDs involving prime attributes on the right side.
 (C) Non-trivial FDs involving prime attributes on the left side
 (D) Non-trivial FDs involving only prime attributes

19. Consider the following database:

SOFTWARE (software-vendor, product, Release-date, systemReq, warranty)

FD: (software-vendor, product, Releasedate) \rightarrow system Req, price, Warranty.

Which of the following are non-prime attributes?

- (A) SystemReq
 (B) Price
 (C) Warranty
 (D) All the above

20. Consider a relation schema $R(A, B, C, D, E, X, Y)$ with the following FDs

$F = \{0 \rightarrow A, XD \rightarrow C, DA \rightarrow B, A \rightarrow X, XE \rightarrow B, E \rightarrow A, B \rightarrow D, DA \rightarrow B, EB \rightarrow C, AB \rightarrow C, Y \rightarrow B, C \rightarrow B\}$ is in

- (A) 2NF
 (B) 3NF
 (C) 4NF
 (D) BCNF

PREVIOUS YEARS' QUESTIONS

1. Which one of the following statements is false? [2007]
- (A) Any relation with two attributes is in BCNF
 (B) A relation in which every key has only one attribute is in 2NF
 (C) A prime attribute can be transitively dependent on a key in a 3NF relation.
 (D) A prime attribute can be transitively dependent on a key in a BCNF relation.
2. Consider the following relational schemas for a library database:
 Book (Title, Author, Catalog_no, Publisher, Year, Price)
 Collection (Title, Author, Catalog_no)
 with the following functional dependencies:
 I. Title Author \rightarrow Catalog_no
 II. Catalog_no \rightarrow Title Author Publisher Year
 III. Publisher Title Year \rightarrow Price
 Assume {Author, Title} is the key for both schemas. Which of the following statements is true? [2008]
- (A) Both Book and Collection are in BCNF
 (B) Both Book and Collection are in 3NF only
 (C) Book is in 2NF and Collection is in 3NF
 (D) Both Book and Collection are in 2NF only
3. The following functional dependencies hold for relations $R(A, B, C)$ and $S(B, D, E)$
 $B \rightarrow A$,
 $A \rightarrow C$
 The relation R contains 200 tuples and the relation S contains 100 tuples. What is the maximum number of tuples possible in the natural join $R \bowtie S$? [2010]
- (A) 100 (B) 200
 (C) 300 (D) 2000
4. Which of the following is true? [2012]
- (A) Every relation in 3NF is also in BCNF
 (B) A relation R is in 3NF if every non-prime attribute of R is fully functionally dependent on every key of R
 (C) Every relation in BCNF is also in 3NF
 (D) No relation can be in both BCNF and 3NF
- Common data questions 5 and 6:** Relation R has eight attributes ABCDEFGH, Fields of R contain only atomic values.
 $F = \{CH \rightarrow G, A \rightarrow BC, B \rightarrow CFH, E \rightarrow A, F \rightarrow EG\}$
 is a set of functional dependencies (FDs) so that F^+ is exactly the set of FDs that hold for R .
5. How many candidate keys does the relation R have? [2013]
- (A) 3 (B) 4
 (C) 5 (D) 6
6. The relation R is [2013]
- (A) in 1NF, but not in 2NF
 (B) in 2NF, but not in 3NF
 (C) in 3NF, but not in BCNF
 (D) in BCNF
7. Assume that in the suppliers relation above, each supplier and each street within a city has a unique name, and (sname, city) forms a candidate key. No other functional dependencies are implied other than those implied by primary and candidate keys. Which one of the following is true about the above schema? [2009]
- (A) The schema is in BCNF
 (B) The schema is in 3NF but not in BCNF
 (C) The schema is in 2NF but not in 3NF
 (D) The schema is not in 2NF
8. Consider the relation scheme $R = (E, F, G, H, I, J, K, L, M, N)$ and the set of functional dependencies $\{\{E, F\} \rightarrow \{G\}, \{F\} \rightarrow \{I, J\}, \{E, H\} \rightarrow \{K, L\}, \{K\} \rightarrow \{M\}, \{L\} \rightarrow \{N\}\}$ on R . What is the key for R ? [2014]
- (A) $\{E, F\}$ (B) $\{E, F, H\}$
 (C) $\{E, F, H, K, L\}$ (D) $\{E\}$
9. Given the following two statements:
 S_1 : Every table with two single-valued attributes is in 1NF, 2NF, 3NF and BCNF
 S_2 : $AB \rightarrow C, D \rightarrow E, E \rightarrow C$ is a minimal cover for the set of functional dependencies $AB \rightarrow C, D \rightarrow E, AB \rightarrow E, E \rightarrow C$
 Which one of the following is correct? [2014]
- (A) S_1 is true and S_2 is false
 (B) Both S_1 and S_2 are true
 (C) S_1 is false and S_2 is true
 (D) Both S_1 and S_2 are false
10. The maximum number of super-keys for the relation schema $R(E, F, G, H)$ with E as the key is _____. [2014]
11. A *prime attribute* of a relation scheme R is an attribute that appears [2014]
- (A) in all candidate keys of R
 (B) in some candidate key of R
 (C) in a foreign key of R
 (D) only in the primary key of R
12. Consider an entity-Relationship (ER) model in which entity sets E_1 and E_2 are connected by an m:n relationship R_{12} . E_1 and E_3 are connected by a 1:n (1 on the side of E_1 and n on the side of E_3) relationship R_{13} . E_1 has two single-valued attributes a_{11} and a_{12} of which a_{11} is the key attribute. E_2 has two single-valued

attributes a_{21} and a_{22} of which a_{21} is the key attribute. E_3 has two single-valued attributes a_{31} and a_{32} of which a_{31} is the key attribute. The relationships do not have any attributes.

If a relational model is derived from the above ER model, then the minimum number of relations that would be generated if all the relations are in 3 NF is _____.

13. Consider the relation $X(P, Q, R, S, T, U)$ with the following set of functional dependencies

$$F = \{ \{P, R\} \rightarrow \{S, T\} \{P, S, U\} \rightarrow \{Q, R\} \}$$

Which of the following is the trivial functional dependency in F^+ , where F^+ is closure of F ? [2015]

- (A) $\{P, R\} \rightarrow \{S, T\}$
 - (B) $\{P, R\} \rightarrow \{R, T\}$
 - (C) $\{P, S\} \rightarrow \{S\}$
 - (D) $\{P, S, U\} \rightarrow \{Q\}$
14. A database of research articles in a journal uses the following schema. [2016]

(VOLUME, NUMBER, STARTPAGE, ENDPAGE, TITLE, YEAR, PRICE)

The primary key is (VOLUME, NUMBER, STARTPAGE, ENDPAGE) and the following functional dependencies exist in the schema.

(VOLUME, NUMBER, STARTPAGE, ENDPAGE) \rightarrow TITLE

(VOLUME, NUMBER) \rightarrow YEAR

(VOLUME, NUMBER, STARTPAGE, ENDPAGE) \rightarrow PRICE

The database is redesigned to use the following schemas.

(VOLUME, NUMBER, STARTPAGE, ENDPAGE, TITLE, PRICE)

(VOLUME, NUMBER, YEAR)

Which is the weakest normal form that the new database satisfies, but the old one does not?

- (A) 1NF
 - (B) 2NF
 - (C) 3NF
 - (D) BCNF
15. Consider the following database table water_schemes: [2016]

water_schemes		
scheme_no	District name	Capacity
1	Ajmer	20
1	Bikaner	10
2	Bikaner	10

3	Bikaner	20
1	Churu	20
2	Churu	20
1	Dungargarh	10

The number of tuples returned by the following SQL query is _____.

```
with total (name, capacity) as
select district_name, sum(capacity)
from water_schemes
group by district_name
with total_avg (capacity) as
select avg(capacity)
from total
select name
from total, total_avg
where total.capacity >= total_avg.capacity
```

16. The following functional dependencies hold true for the relational schema $R \{V, W, X, Y, Z\}$:

$V \rightarrow W$

$VW \rightarrow X$

$Y \rightarrow VX$

$Y \rightarrow Z$

Which of the following is irreducible equivalent for this set of set of functional dependencies? [2017]

- (A) $V \rightarrow W$
- (B) $V \rightarrow W$
- (C) $V \rightarrow W$
- (D) $V \rightarrow W$

17. Consider the following tables T1 and T2.

T1		T2	
P	Q	R	S
2	2	2	2
3	8	8	3
7	3	3	2
5	8	9	7
6	9	5	7
8	5	7	2
9	8		

In table T1, **P** is the primary key and **Q** is the foreign key referencing **R** in table T2 with on-delete cascade and on-update cascade. In table T2, **R** is the primary key and **S** is the foreign key referencing **P** in table T1

with on-delete set NULL and on-update cascade. In order to delete record $\langle 3, 8 \rangle$ from table T1, the number of additional records that need to be deleted from table T1 is _____. [2017]

18. Consider the following four relational schemas. For each schema, all non-trivial functional dependencies are listed. The underlined attributes are the respective primary keys.

Schema I:

Registration (rollno, courses)

Field 'courses' is a set-valued attribute containing the set of courses a student has registered for.

Non-trivial functional dependency:

Rollno \rightarrow courses

Schema II:

Registration (rollno, courseid, email)

Non-trivial functional dependencies:

Rollno, courseid \rightarrow email

email \rightarrow rollno

Schema III:

Registration (rollno, courseid, marks, grade)

Non-trivial functional dependencies:

Rollno, courseid \rightarrow marks, grade

Marks \rightarrow grade

Schema IV:

Registration (rollno, courseid, credit)

Non-trivial functional dependencies:

Rollno, courseid \rightarrow credit

Courseid \rightarrow credit

Which one of the relational schemas above is in 3NF but not in BCNF? [2018]

- (A) Schema I
(B) Schema II
(C) Schema III
(D) Schema IV

ANSWER KEYS

EXERCISES

Practice Problems 1

1. A 2. A 3. D 4. B 5. B 6. B 7. C 8. D 9. B 10. A
11. C 12. D 13. A 14. B 15. A 16. D 17. A 18. B 19. C 20. C

Practice Problems 2

1. D 2. C 3. A 4. B 5. A 6. C 7. D 8. D 9. C 10. C
11. B 12. A 13. D 14. B 15. A 16. A 17. C 18. C 19. D 20. B

Previous Years' Questions

1. D 2. C 3. A 4. C 5. B 6. A 7. B 8. B 9. A 10. 8
11. B 12. 4 13. C 14. B 15. 2 16. A 17. 0 18. B

Chapter 4

Transaction and Concurrency

LEARNING OBJECTIVES

- ☞ Transactions and concurrency control
- ☞ Transaction
- ☞ Transaction properties
- ☞ Uncommitted data
- ☞ Transaction processing systems
- ☞ Concurrency control with locking methods
- ☞ Two-phase locking to ensure serializability
- ☞ Concurrency control with time stamping methods
- ☞ Concurrency control with optimistic methods
- ☞ Recoverability
- ☞ Equivalence of schedules
- ☞ Testing for conflict serializability

INTRODUCTION

A transaction is a logical unit of work. It begins, with the execution of a `BEGIN TRANSACTION` operation, and ends with the execution of a `COMMIT` or `ROLLBACK` operation. The logical unit of work that is, a transaction does not necessarily involve just a single database operation. Rather, it involves a sequence of several such operations as follows:

1. Database updates are kept in buffers in main memory and not physically written to disk until the transaction commits. That way, if the transaction terminates unsuccessfully, there will be no need to undo any disk updates.
2. Database updates are physically written to disk as part of the process of honouring the transaction's `COMMIT` request. That way if the system subsequently crashes, we can be sure that there will be no need to redo any disk updates.

Transactions and Concurrency Control

Database transactions reflect real-world transactions that are triggered by events, such as buying a product, registering for a course, or making a deposit in your checking account. Transactions are likely to contain many parts, for example, a sales transaction consists of at least two parts.

`UPDATE` inventory by subtracting number of units sold from the `PRODUCT` table's available quantity on hand and `UPDATE` the `ACCOUNTS RECEIVABLE` table in order to bill the `CUSTOMER`. All parts of a transaction must be completed to prevent data integrity problems. Therefore, executing and managing transactions are important database system activities.

Concurrency control is the management of concurrent transactions execution. When many users are able to access the database, the number of concurrent transactions tends to grow rapidly; as a result, concurrency control is especially important in multiuser database environments.

TRANSACTION

A transaction is a logical unit of work that must be either entirely completed or aborted, no intermediate states are acceptable, that is, multicomponent transactions like the previously mentioned sale, must not be partially completed. If you read from and/or write to (update) the database, you create a transaction. Another example is using `SELECT`, to generate a list of table contents. Many real-world database transactions are formed by two or more database requests. A database request is the equivalent of a single SQL statement in an application program or transaction. Each database request generates several input/output operations. A transaction that changes the contents of a database must alter the database from one consistent state to another. A consistent database state is one in which all data integrity constraints are satisfied.

Example:

1. Checking an account balance:

```
SELECT ACC_NUM, ACC_BALANCE
FROM CHECKACC
WHERE ACC_NUM = '0908110638';
```

Even though we did not make any changes to the `CHECKACC` table, the SQL code represents a transaction, because we accessed the database.

2. Registering a credit sale of 100 units of product X to customer Y in the amount of \$500.00 first, product X 's quantity on hand (QOH) needs to be reduced by 100.

```
UPDATE PRODUCT
SET PROD_QOH = PROD_QOH_100
WHERE PROD_CODE = 'x';
Then, $500 needs to be added to customer Y's
accounts receivable
UPDATE ACCT_RECEIVABLE
SET ACCT_RECEIVABLE = ACCT_BALANCE +
500
WHERE ACCT_NUM = 'Y';
```

In Example 2, both the SQL, transactions must be completed in order to represent the real-world sales transaction. If both transactions are not completely executed, the transaction yields an inconsistent database.

If a transaction yields an inconsistent database, the DBMS must be able to recover the database to a previous consistent state.

Transaction Properties

All transactions must display atomicity, consistency, isolation and durability. These are known as ACID properties of transactions.

Atomicity

It requires that all operations of a transaction be completed; if not, the transaction is aborted. Therefore, a transaction is treated as a single, logical unit of work.

Consistency

It describes the result of the concurrent execution of several transactions. The concurrent transactions are treated as though they were executed in serial order. This property is important in multiuser and distributed database, where several transactions are likely to be executed concurrently.

Isolation

It means that the data used during the execution of a transaction cannot be used by a second transaction until the first one is completed. Therefore, if a transaction T_1 is being executed and is using the data item X_1 , that data item cannot be accessed by any other transaction ($T_2 \dots T_n$) until T_1 ends. This property is particularly useful in multiuser database environment, because several different users can access and update the database at the same time.

Durability

It indicates the permanence of the database's consistent state. When a transaction is completed, the database reaches a consistent state, and that state cannot be lost, even in the event of the system's failure.

Transaction Management with SQL

The ISO standard defines a transaction model based on two SQL statements: COMMIT and ROLLBACK. The standard specifies that an SQL transaction automatically begins with a transaction-initiating SQL statement executed by a user or program (e.g., SELECT, INSERT, UPDATE). Changes made by a transaction are not visible to other concurrently executing transactions until the transaction completes. When a transaction sequence is initiated, it must continue through all succeeding SQL, statements until one of the following four events occur:

1. A COMMIT statement ends the transaction successfully, making the database changes permanent. A new transaction starts after COMMIT with the next transaction initiating statement.
2. For programmatic SQL, successful program termination ends the final transaction successfully, even if a commit statement has not been executed (equivalent to COMMIT)
3. For programmatic SQL, abnormal program termination aborts the transaction (equivalent to ROLLBACK)

Example:

```
UPDATE PRODUCT
SET PROD_QOH = PROD_QOH_100
WHERE PROD_CODE = '345TYX';
UPDATE ACCREC
SET AR_BALANCE = AR_BALANCE + 3500
WHERE AR_NUM = '60120010';
COMMIT;
```

CONCURRENCY CONTROL

The coordination of simultaneous execution of transactions in a multiprocessing database system is known as *concurrency control*. The objective of concurrency control is to ensure the serializability of transaction in a multiuser database environment. Concurrency is important, because the simultaneous execution of transactions over a shared database can create several data integrity and consistency problems. Three main problems are lost updates, uncommitted data and inconsistent retrievals.

Lost Updates

Consider the following two concurrent transactions where PROD_QOH represents a particular PRODUCT's quantity on hand. (PROD_QOH is an attribute in the Product table)

Assume the current PROD_QOH value for the product concerned is 35.

Table 1

Transaction	Computation
T_1 : purchase 100 units	PROD_QOH = PROD_QOH + 100
T_2 : sell 30 units	PROD_QOH = PROD_QOH - 30

Table 1 shows the serial execution of these transactions under normal circumstances, yielding the correct answer: PROD_QOH = 105. But suppose that a transaction is able to read a product's PROD_QOH value from the table before a previous transaction (using the same product) has been committed. The sequence depicted in Table 2 shows how the cost update problem can arise. Note that the first transaction (T_1) has not yet been committed when the second transaction (T_2) is executed. Therefore T_2 still operates on the value 135 to disk which is promptly over written by T_2 . As a result, the addition of 100 units is "lost" during the process.

Uncommitted Data

Data are not committed when two transactions, T_1 and T_2 are executed concurrently and the first transaction (T_1) is rolled back after the second transaction (T_2) has already accessed the uncommitted data, thus violating the isolation property of transaction. Consider the same transactions from T_1 and T_2 , from above. However, this time T_1 is rolled back to eliminate the addition of the 100 units. Because T_2 subtracts 30 from the original 35 units, the correct answer should be 5.

Table 2

	Computation
T_1 : purchase 100 units	PROD_QOH = PROD_QOH + 100 (Rolled Back)
T_2 : sell 30 units	PROD_QOH = PROD_QOH - 30

Table 2 shows how, under normal circumstances, the serial execution of these transactions yield the correct answer. The uncommitted data problem can arise when the ROLLBACK is completed after T_2 has begun its execution.

Inconsistent Retrievals

Inconsistent retrievals occur when a transaction calculates some summary (aggregate) functions over a set of data, while other transactions are updating the data. The problem is that the transaction might read some data before they are changed and other data after they are changed, thereby yielding inconsistent results.

Example:

1. T_1 calculates the total PROD_QOH of the products stored in the PRODUCT table

2. At the same time, T_2 updates the PROD_QOH for two of the PRODUCT table's products (T_2 represents the correction of a typing error: the user added 30 units to product 345TYX's PROD_QOH but meant to add the 30 units to 125TYZ's PROD_QOH to correct the problem, the user subtracts 30 from product 345TYX's PROD_QOH and adds 30 to product 125TYZ's PROD_QOH).

The computed answer 485 is obviously wrong, because we know the correct answer to be 455.

TRANSACTION PROCESSING SYSTEMS

Transaction processing systems are systems with large databases and hundreds of concurrent users that are executing database transactions. For example, banking, credit card processing, stock markets, supermarket checkout, etc.

They require high availability and fast response time for hundreds of concurrent users.

1. A transaction includes one or more database access operations. These can include insertion, deletion, modification, or retrieval operations.
2. *Basic operations*: The basic database access operations that a transaction can include are as follows:
 - read_item (X): Reads a database item named X into a program variable.
 - Write_item (X): Writes the value of program variable X into the database item named X

Executing a read_item (X): Command includes the following steps:

1. Find the address of the disk block that contains item X
2. Copy that disk block into a buffer in main memory (if that disk block is not in main memory buffer).
3. Copy item X from the buffer to the program variable named x

Executing a write_item (X) command includes the following steps:

1. Find the address of the disk block that contains item X
2. Copy that disk block into a buffer in main memory (if that disk block is not in main memory buffer)
3. Copy item X from the program variable named X into its correct location in the buffer
4. Store the updated block from the buffer back to disk.

Step 4 actually updates the database on disk.

The decision about when to store back a modified disk block that is in a main memory buffer is handled by the recovery manager of the DBMS in cooperation with the underlying operating system.

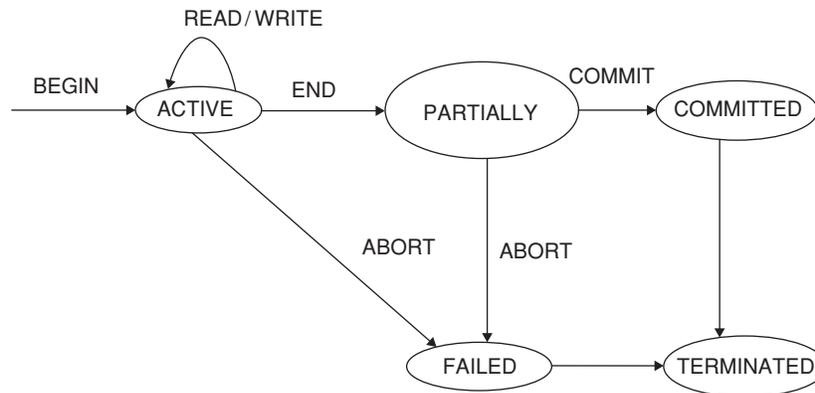


Figure 1 Transactions execution state transition diagram.

For the purpose of recovery, the system needs to keep track of when the transaction starts, terminates, and commits or aborts. Hence, the recovery manager keeps track of the following operations:

1. **BEGIN_TRANSACTION:** It shows the beginning of Execution of a transaction.
2. **READ/WRITE:** These specify read or write operations on the database items.
3. **END_TRANSACTION:** This specifies that READ and WRITE operations have ended and marks the end of transaction execution.
4. **COMMIT_TRANSACTION:** This shows a successful end of the transaction so that any changes executed by the transaction can be safely committed to the database and will not be undone.
5. **ROLL BACK OR ABORT:** This shows that the transaction has ended unsuccessfully, so that any changes or effects that the transaction may have applied to the database must be undone.
6. **ACTIVE STATE:** A transaction goes into an active state immediately after it starts execution where it can issue READ and WRITE operations.
7. **PARTIALLY COMMITTED:** When the transaction ends, it moves to the partially committed state
8. **COMMIT:** A transaction reaches its commit point when all its operations that access the database have been executed successfully, and the effect of all the transaction operations on the database have been recorded in the log. Beyond the commit point, the transaction is said to be committed, and its effect is assumed to be permanently recorded in the database.
9. **FAILED_STATE:** A transaction can go to the failed state if the transaction is aborted during its active state. The transaction may then have to be rolled back to undo the effect of its WRITE operations on the database.
10. **TERMINATED:** The terminated state corresponds to the transactions leaving the system.

CONCURRENCY CONTROL WITH LOCKING METHODS

A lock guarantees exclusive use of a data item to a transaction. In general, if transaction T_1 holds a lock on a data item (e.g., an employee's salary) then transaction T_2 does not have access to that data item. A transaction acquires a lock prior to data access; the lock is released (unlocked) when the transaction is completed, so that another transaction can lock the data item for its exclusive use. All lock information is managed by a lock manager, which is responsible for assigning and policing the locks used by the transactions.

Lock Granularity

Lock granularity indicates the level of lock use. Locking can take place at the following levels: database level, table level, page level, row level and field (or attribute) level.

Database Level

In a database-level lock, the entire database is locked, thus preventing the use of any tables in the database by transaction T_2 while transaction T_1 is being executed. Transaction T_1 and T_2 cannot access the database concurrently, even if they use different tables. This level of locking is suitable for batch processes, but it is not unsuitable for online multiuser DBMSs.

Table Level

In a table-level lock, the entire table is locked, preventing access to any row by transaction T_2 while transaction T_1 is using Table 2 transactions can access the same database, as long as they access different tables. Transactions T_1 and T_2 cannot access the same table even if they try to use different rows, T_2 must wait until T_1 unlocks the table.

Page level

In a page level lock, the DBMS will lock an entire disk page (a disk page or page is the equivalent of a disk block, which can

be described as a (referenced) section of a disk). Transactions T_1 and T_2 access the same table while locking different disk pages. If T_2 requires the use of a row located on a page that is locked by T_1 , T_2 must wait until the page is unlocked by T_1 .

Row level

The row-level lock is much less restrictive than the locks discussed earlier. The DBMS allows concurrent transactions to access different rows of the same table, even if the rows are located on the same page. A lock exists for each row in each table of the database.

Field level

The field-level lock allows concurrent transactions to access same row as long as they require the use of different fields (attributes) within the row. Although, field-level locking clearly yields the most flexible multi user data access, it requires a high level of computer overhead.

Lock Types

1. Binary locks
2. Shared/Exclusive locks

Binary Locks

A binary lock has only two states: locked (1) or unlocked (0). If an object, that is, a database, table, page, or row is locked by a transaction, no other transaction can use that object. If an object is unlocked, any transaction can lock the object for its use. As a rule, a transaction must unlock the object after its termination. Every database operation requires that the affected object be locked. Therefore, every transaction requires a lock and unlock operation for each data item that is accessed. Such operations are automatically scheduled by the DBMS, the user need not be concerned about locking or unlocking data items. Binary locks are now considered too restrictive to yield optimal concurrency conditions. For example if two transactions want to read the same database object, the DBMS will not allow this to happen, even though neither transaction updates the database (and therefore, no concurrency problems can occur) concurrency conflicts occur only when two transactions execute concurrently and one of them updates the database.

Shared/Exclusive locks

The terms “shared” and “exclusive” indicate the nature of the lock. The following table comparatively explains both locks.

Exclusive Locks	Shared Locks
An <i>exclusive lock</i> exists when access is specifically reserved for the transaction that locked the object.	A <i>shared lock</i> exists when concurrent transactions are granted READ access on the basis of a common lock.
The exclusive lock must be used when the potential for conflict exists.	A shared lock produces no conflict as long as the concurrent transactions are read only.
(An exclusive lock is issued when a transaction wants to write (update) a data item and no locks are currently held on that data item by any other transaction.	A shared lock is issued when a transaction wants to read data from the database and no exclusive lock is held on that data item.

Using the shared/exclusive locking concept, a lock can have three states: unlocked, shared (READ) and exclusive (WRITE). 2 READ transactions can be safely executed and shared locks allow several READ transactions to concurrently read the same data item. For example, if transaction T_1 has a shared lock on data item X , and transaction T_2 wants to read data item X , T_2 may also obtain a shared lock on data item X .

If transaction T_2 updates data item X , then an exclusive lock is required by T_2 over data item X . The exclusive lock is granted if and only if no other locks are held on the data item. Therefore, if a shared or exclusive lock is already held on data item X by transaction T_1 , an exclusive lock cannot be granted to transaction T_2 .

Potential problems with locks

Although locks prevent serious data inconsistencies, their use may lead to two major problems:

1. The resulting transaction schedule may not be serializable.
2. The schedule may create deadlocks. Database *deadlocks* are the equivalent of a traffic gridlock in

a big city and are caused when *two transactions wait for each other to unlock data*.

Both problems can be solved. Serializability is guaranteed through a locking protocol known as two-phase locking and deadlocks can be eliminated by using deadlock detection, and prevention techniques. We shall examine these techniques next.

Two-phase Locking to Ensure Serializability

The *two-phase locking protocol* defines how transactions acquire and relinquish locks. It guarantees serializability, but it does NOT prevent deadlocks. The two phases are as follows:

1. A *growing phase*, in which a transaction acquires all the required locks without unlocking any data. Once all locks have been acquired, the transaction is in its locked point.
2. A *shrinking phase*, in which a transaction releases all locks and cannot obtain any new lock.

The two-phase locking protocol is governed by the following rules”

- Two transactions cannot have conflicting locks.
- No unlock operation can precede a lock operation in the same transaction.
- No data are affected until all locks are obtained, that is, until the transaction is in its locked point.

DEADLOCKS

Deadlocks exist when two transactions, T_1 and T_2 , exist in the following mode:

- T_1 would like to access data item X and then data item Y . (So far, T_1 has locked data item X and T_1 is in progress, it will eventually require to lock data item Y .)
- T_2 needs to access data items X and Y , to begin. (So far, T_2 has locked data item Y .)

If T_1 has not unlocked data item X , T_2 cannot begin; if T_2 has not unlocked data item Y , T_1 cannot continue. Consequently, T_1 and T_2 wait indefinitely, each waiting for the other to unlock the required data item. Such in a real-world DBMS, many transactions can be executed simultaneously, thereby increasing the probability of generating deadlocks. Note that deadlocks are possible only if one of the transactions wants to obtain an exclusive lock on a data item; no deadlock condition can exist among shared locks.

Three basic techniques exist to control deadlocks:

- Deadlock prevention:** A transaction requesting a new lock is aborted if there is a possibility that a deadlock can occur. If the transaction is aborted, all the changes made by this transaction are ROLLED BACK, and all locks obtained by the transaction are released. The transaction is then rescheduled for execution. Deadlock prevention works because it avoids the conditions that lead to deadlocking.
- Deadlock detection:** The DBMS periodically tests the database for deadlocks. If a deadlock is found, one of the transactions (the “victim”) is aborted (ROLLED BACK and restarted), and the other transaction continues.
- Deadlock avoidance:** The transaction must obtain all the locks it needs before it can be executed.

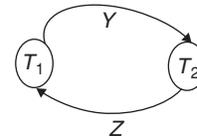
The best deadlock control method depends on the database environment. For example, if the probability of deadlocks is low, deadlock detection is recommended. However, if the probability of deadlocks is high, deadlock prevention is recommended. If response time is not high on the system priority list, deadlock avoidance might be employed.

Deadlock occurs when each transaction T in a set of two or more transactions is waiting for some item that is locked by some other transaction T^i in the set. Each transaction in the set is on a waiting queue, waiting for one of the other transactions in the set to release the lock on an item

Example:

T_1	T_2
Write lock (z)	
	Read lock (z)
	Read lock (y)
Write lock (y)	

Transaction T_1 is waiting for Y which is locked by Transaction T_2 and transaction T_2 is waiting for z which is locked by transaction T_1 . The below graph is called *wait for graph*.



Deadlock Prevention

There are number of deadlock prevention schemes that make a decision about what to do with a transaction involved in a possible deadlock situation:

- Should it be blocked and made to wait
- Should it be aborted
- Should the transaction pre-empt and abort another transaction.

CONCURRENCY CONTROL WITH TIME STAMPING METHODS

The *time stamping* approach to scheduling concurrent transactions assigns a global unique time stamp to each transaction. The time stamp value produces an explicit order in which transactions are submitted to the DBMS. Time stamps must have two properties: uniqueness and monotonicity. *Uniqueness* ensures that no equal time stamp values can exist, and *monotonicity* ensures that time increases.

All database operations (READ and WRITE) within the same transaction must have the same time stamp. The DBMS executes conflicting operations in time stamp order, thereby ensuring serializability of the transactions. If two transactions conflict, one often is stopped, rescheduled, and assigned a new time stamp value.

The concept of transaction time stamp $TS(T)$, which is a unique identifier assigned to each transaction. The time stamps are based on the order in which transactions start. If transaction T_1 starts before transaction T_2 , then $TS(T_1) < TS(T_2)$

- Older transaction will have the smaller time stamp value
- Two schemes that prevent deadlock are
 - Wait-die
 - Wound-wait

Suppose that transaction T_k tries to lock an item x but is not able to because x is locked by some other transaction T_L with a conflicting lock.

The rules followed by these schemes are as follows:

1. Wait-die: If $T_s(T_k) < T_s(T_L)$, then (T_k older than T_L) T_k is allowed to wait; otherwise (T_k younger than T_L) abort T_k and restart it later with the same time stamp
2. Wound-wait: If $T_s(T_k) < T_s(T_L)$ then (T_k older than T_L) abort T_L and restart it later with the same time stamp; otherwise (T_k younger than T_L) T_k is allowed to wait

Both schemes end up aborting the younger of the two transactions that may be involved in a deadlock

Concurrency Control with Optimistic Methods

Optimistic methods are based on the assumption that the majority of the database operations do not conflict. A transaction is executed without restrictions until it is committed. Each transaction moves through *two* or *three* phases:

Read phase: The transaction reads the database, executes the needed computations, and makes the updates to a private copy of the database values.

Validation phase: The transaction is validated to assure that the changes made will not affect the integrity and consistency of the database.

If the validation test is positive, transaction goes to the Write Phase.

If the validation test is negative, transaction is restarted, and changes are discarded.

Write phase: The changes are permanently applied to the database.

SERIALIZABILITY

Serializability is accepted as ‘criterion for correctness’ for the interleaved execution of a set of transactions, such an execution is considered to be correct if and only if it is serializable.

1. A set of transactions is serializable if and only if it is equivalent to some serial execution of the same transactions
2. A serial execution is one in which the transactions are run one at a time in some sequence

Schedule: Given a set of transactions, any execution of those transactions interleaved or otherwise is called a *schedule*.

1. Executing the transactions one at a time, with no interleaving constitutes a serial schedule. A schedule that is not serial is an interleaved schedule (or) non-serial schedule.
2. Two schedules are said to be equivalent if and only if they are guaranteed to produce the same result as each other. Thus, a schedule is serializable, and correct, if and only if it is equivalent to some serial schedule.

Two-phase Locking Theorem

If all transactions obey the two phase locking protocol, then all possible interleaved schedule are serializable.

1. Before operating on any object (it could be a database tuple), a transaction must acquire a lock on the object
2. After releasing a lock, a transaction must never go on to acquire any more locks.

A transaction that obeys this protocol thus has two phases: a lock acquisition or “growing phase and a lock releasing or “shrinking” phase

Let ‘ T ’ be an interleaved schedule involving some set of transactions $T_1, T_2, T_3, \dots, T_n$.

If ‘ T ’ is serializable, then there exists at least one serial schedule ‘ S ’ involving T_1, T_2, \dots, T_n such that ‘ T ’ is equivalent to ‘ S ’ is said to be a serialization of ‘ T ’

Let T_i and T_j be any two distinct transactions in the set $T_1, T_2, T_3, \dots, T_n$. Let T_i precede T_j in the serialization ‘ S ’. In the interleaved schedule I , then the effect must be as if T_i really did execute before T_j . In other words, if A and B are any two transactions involved in some serializable schedule, then either A logically precedes B or B logically precedes A in that schedule, that is, either B can see A ’s output or A can see B ’s. If the effect is not as if either A ran before B or B ran before A , then the schedule is not serializable and not correct.

1. A schedule ‘ S ’ of ‘ n ’ transactions T_1, T_2, \dots, T_n is an ordering of the operations of the transactions subject to the constraint that, for each transaction T_i that participates in ‘ S ’, the same order in which they occur in T_i .
2. For the purpose of recovery and concurrency control, we are mainly interested in the ‘read item’ and ‘write item’ operations of the transactions, as well as the COMMIT and ABORT operations. A shorthand notation for describing a schedule uses the symbols, ‘ R ’, ‘ W ’, ‘ C ’ and ‘ A ’ for the operations read item, write item, commit, and abort respectively, and appends as subscript the transaction-id (transaction number) to each operation in the schedule

Example: The schedule of the given set of transactions can be written as follows:

T_1	T_2
Read item (x); $X = X - N$;	Read item (x); $X = X + M$;
Write item (x)	Write item (x)
Read item (y)	Write item (y)
$Y = Y + N$	Write item (y)
Write item (y)	Commit

Schedule:

$S: R_1(X); R_2(X); W_1(X); R_1(Y); W_2(X); W_2(Y); W_1(Y); C_2$

Conflicts: Two operations in a schedule are said to have conflict if they satisfy all three conditions, if

1. they belong to different transactions
2. they access the same data item
3. at least one of the operations is a write item

For example, In the schedule ‘S’ given above, the operations $r_1(x)$ and $w_2(x)$ conflict, as do

The operations $r_2(x)$ and $w_1(x)$ and the operations $w_1(x)$ and $w_2(x)$. However the operations $r_1(x)$ and $r_2(x)$ do not conflict, since they are both read operations;

The operations $w_1(x)$ and $w_2(y)$ do not conflict because they operate on distinct data items x and y . The operations $r_1(x)$ and $w_1(x)$ do not conflict, because they belong to the same transaction.

Complete schedule: A schedule S of ‘ n ’ transactions $T_1, T_2, T_3 \dots T_n$ is said to be a complete schedule if the following conditions hold.

1. The operations in ‘S’ are exactly those operations in $T_1, T_2, \dots T_n$ including a commit or abort operation as the last operation for each transaction in the schedule.
2. For any pair of operations from the same transaction T_i , their order of appearance in ‘S’ is the same as their order of appearance in T_i
3. For any two conflicting operations, one of the two must occur before the other in the schedule.

The preceding condition (3) allows for two non-conflicting operations to occur in the schedule without defining which occurs first, thus leading to the definition of a schedule as a partial order of the operations in the ‘ n ’ transactions.

It is difficult to encounter complete schedules in a transaction processing system, because new transactions are continually being submitted to the system. Hence, it is useful to define the concept of the ‘committed projection $C(S)$ ’ of schedule S ; which include only the operations in S that belong to committed transactions, that is, transaction T_i whose commit operation is C_i .

RECOVERABILITY

Recoverability ensures that once a transaction T is committed, it should never be necessary to roll back T . The schedules that theoretically meet this criterion are called *recoverable schedules* and those that do not are called *non-recoverable*, and hence should not be permitted

A schedule ‘S’ is recoverable if no transaction T in ‘S’ commits until all transactions T^1 that have written an item that T reads have committed

A transaction T reads from transaction T^1 in a schedule S if some item x is first written by T^1 and later read by T . In addition, T^1 should not have been aborted before T reads item x , and there should be no transactions that write x after T^1 writes it and before T reads it (unless those transactions, if any, have aborted before T reads x)

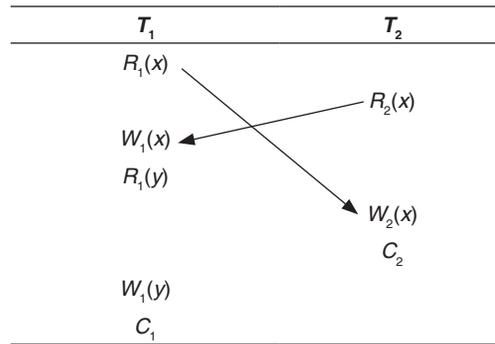
Example:

Consider the given schedule, check whether it is recoverable or not:

$S: R_1(X); R_2(X); W_1(X); R_1(Y); W_2(X); C_2, W_1(Y); C_1?$

Solution:

The given schedule can also be represented as follows:



There are two *WR* conflicts, if the schedule consists of *RW* conflict, then we may say that the schedule is not recoverable (if the transaction which is performing read operation commits first)

Cascadeless Schedule

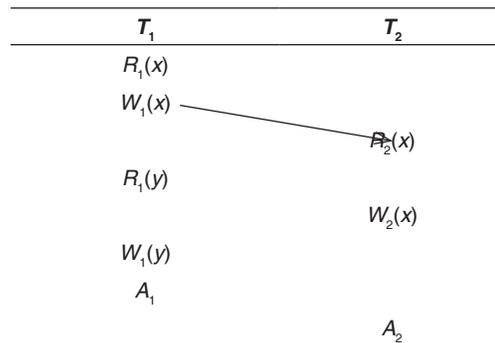
In a recoverable schedule, no committed transaction ever needs to be rolled back. It is possible for a phenomenon known as cascading rollback to occur, when an uncommitted transaction has to be rolled back because it read an item from a transaction that failed.

This is illustrated in the following schedule:

Example:

$S: R_1(X); W_1(X); R_2(X); R_1(Y); W_2(X); W_1(Y); A_1, A_2$

The above schedule is represented as follows:



Transaction T_2 has to be rolled back because it reads item x from T_1 , and T_1 is then aborted, because cascading rollback

can be quite time consuming since numerous transactions can be rolled back. It is important to characterize the schedules where this phenomenon is guaranteed not to occur.

A schedule is said to be cascadeless if every transaction in the schedule reads only items that were written by committed transaction. In this case, all items read will not be discarded, so no cascading rollback will occur.

Strict Schedule

A schedule is called *strict schedule*, in which transactions can neither read nor write an item x until the last transaction that wrote x has committed or aborted

1. All strict schedules are cascadeless
2. All cascadeless schedules are recoverable

EQUIVALENCE OF SCHEDULES

There are several ways to define equivalence of schedules as follows:

1. Result equivalent
2. Conflict equivalent
3. View equivalent

Result Equivalent

Two schedules are called *result equivalent* if they produce the same final state of the database. However, two different schedules may accidentally produce the same final state.

Example: Check whether the two schedules are result equivalent or not:

S_1	S_2
Read item (x);	Read item (x)
$X = X + 20$;	$X = X * 1.1$;
Write item (x);	Write item (x);

Solution:

Schedules S_1 and S_2 will produce the same final database state if they execute on a database with an initial value of $x = 200$; but for other initial values of x , the schedules are not result equivalent

For two schedules to be equivalent, the operations applied to each data item affected by the schedules should be applied to that item in both schedules in the same order. The other two definitions of equivalence of schedules generally used are conflict equivalence and view equivalence

Conflict Equivalence

Two schedules are said to be conflict equivalent if the order of any two conflicting operations is the same in both schedules. If two conflicting operations are applied in different orders in two schedules, the effect can be different on the database or on other transactions in the schedule, and hence the schedules are not conflict equivalent.

Example:

S_1	S_3
$r_1(x)$	$w_1(x)$
$w_2(x)$	$w_2(x)$
S_2	S_4
$w_2(x)$	$w_2(x)$
$r_1(x)$	$w_1(x)$

The value read by $r_1(x)$ can be different in the two schedules. Similarly, if two write operations occur in the order $w_1(x), w_2(x)$ in s_3 , and in the reverse order $w_2(x), w_1(x)$ in s_4 , the next $r(x)$ operation in the two schedules will read potentially different values.

Testing for Conflict Serializability

The following algorithm can be used to test a schedule for conflict serializability. The algorithm takes read item and write item operations in a schedule to construct a precedence graph or serialization graph, which is a directed graph $G(N, E)$ here N is a set of Nodes $N = \{T_1, T_2, \dots, T_n\}$ and E is a set of directed edges $E = \{e_1, e_2, \dots, E_m\}$ There is one node in the graph for each transaction. T_i in the schedule. Each edge e_i in the graph is of the form $(T_j \rightarrow T_k), 1 \leq j < n, 1 \leq k < n$,

Where T_j is the starting node of e_i and T_k is the ending node of e_i .

Edge is created if one of the operations in T_j appears in the schedule before some conflicting operation in T_k

Algorithm

1. For each transaction T_i participating in schedule S , create a node labelled T_i in the precedence graph
2. For each case in S where T_j executes a read item (x) after T_i executes a write item (x), create an edge $(T_i \rightarrow T_j)$ in the graph
3. For each case in S where T_j executes a write item (x) after T_i executes a read item (x), create an edge $(T_i \rightarrow T_j)$ in the graph
4. For each case in S where T_j executes a write item (x) after T_i executes a write item (x), create an edge $(T_i \rightarrow T_j)$ in the precedence graph
5. The schedule S is serializable, if the precedence graph contains no cycles

A cycle in a directed graph is a sequence of edges $C = ((T_i \rightarrow T_k), (T_k \rightarrow T_p), \dots, (T_i \rightarrow T_j))$

With the property that the starting node of each edge, except the first edge is the same as the ending node of the previous edge, and the starting node of the first edge is the same as the ending node of the last edge.

Example: Check whether the given schedule is conflict serializable or not by drawing precedence graph:

T_1	T_2	T_3
$R_1(x)$		
	$W_2(x)$	
		$R_3(x)$
$W_1(x)$		
		$W_3(x)$
$R_1(x)$		

Solution:

First identify the conflicts:

$(T_1 \rightarrow T_2)$ WR conflict

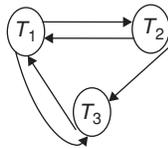
$(T_2 \rightarrow T_1)$ WW conflict

$(T_2 \rightarrow T_3)$ RW conflict

$(T_3 \rightarrow T_1)$ WR conflict

$(T_1 \rightarrow T_3)$ WW conflict

Take transactions as nodes in the precedence graph:



The precedence graph has cycle, which says that the schedule is not serializable.

View Equivalence and View Serializability

View equivalence is less restrictive compared to conflict equivalence. Two schedules S and S' are said to be view equivalent if the following three conditions hold:

1. The same set of transactions participate in S and S' , and S and S' include the same operations of those transactions
2. For any operation $r_i(x)$ of T_i in S , if the value of x read by the operation has been written by an operation $w_j(x)$ of T_j , the same condition must hold for the value of x read by operation $r_i(x)$ of T_i in S'
3. If the operation $w_k(y)$ of T_k is the last operation to write Y in S , then $W_k(y)$ of T_k must also be the last operation to write item Y in S'

The idea behind view equivalence is that as long as each read operation of a transaction reads the result of the same write operation in both schedules, the write operations of each transaction must produce the same result. Hence the read operation is said to see the same view in both schedules:

1. A schedule S is said to be view serializable if it is view equivalent to a serial schedule.
2. All conflict serializable schedules are view serializable, but vice versa is not true.

EXERCISES

Practice Problem 1

Directions for questions 1 to 20: Select the correct alternative from the given choices.

1. Consider the given schedules S_1 and S_2
 $S_1: r_1(x), r_1(y), r_2(x), r_2(y), w_2(y), w_1(x)$
 $S_2: r_1(x), r_2(x), r_2(y), w_2(y), r_1(y), w_1(x)$
 Which schedule is conflict serializable?
 (A) S_1 (B) S_2
 (C) S_1 and S_2 (D) None of these
2. Consider the given schedule with three transactions T_1, T_2 and T_3 :

T_1	T_2	T_3
$r_1(x)$		
	$r_2(y)$	
		$r_3(y)$
	$w_2(y)$	
$w_1(x)$		
		$w_3(x)$
	$r_2(x)$	
	$w_2(x)$	

Which of the following is correct serialization?

- (A) $T_2 \rightarrow T_1 \rightarrow T_3$
- (B) $T_1 \rightarrow T_3 \rightarrow T_2$
- (C) $T_3 \rightarrow T_1 \rightarrow T_2$
- (D) None of these

3. Consider the three data items D_1, D_2 and D_3 and the following execution of schedules of transactions T_1, T_2 and T_3 :

T_1	T_2	T_3
	$R(D_2)$	
	$R(D_2)$	
	$W(D_2)$	
		$R(D_2)$
		$R(D_3)$
$R(D_1)$		
$W(D_1)$		
		$W(D_2)$
		$W(D_3)$
	$R(D_1)$	
$R(D_2)$		
$W(D_2)$		
	$W(D_1)$	

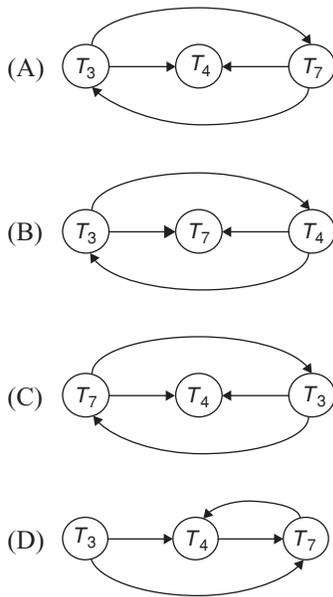
Which of the following is true?

- (A) The schedule is conflict serializable
- (B) The schedule is not conflict serializable
- (C) The schedule has deadlock
- (D) Both (A) and (C)

4. Consider the given schedule

T_3	T_4	T_7
$R(Q)$		
	$W(Q)$	
$W(Q)$		
		$R(Q)$
		$W(Q)$

Which of the following is the correct precedence graph for the above schedule?



5. Consider two Transactions T_1 and T_2 and four schedules: S_1, S_2, S_3 and S_4 of T_1 and T_2 :

T_1 : $r_1(x), w_1(x), w_1(y)$

T_2 : $r_2(x), r_2(y), w_2(y)$

S_1 : $r_1(x), r_2(x), r_2(y), w_1(x), w_1(y), w_2(y)$

S_2 : $r_1(x), r_2(x), r_2(y), w_1(x), w_2(y), w_1(y)$

S_3 : $r_1(x), w_1(x), r_2(x), w_1(y), r_2(y), w_2(y)$

S_4 : $r_2(x), r_2(y), r_1(x), w_1(x), w_1(y), w_2(y)$

Which schedules are conflict serializable in the given schedules?

- (A) S_1 and S_2
- (B) S_1 and S_3
- (C) S_2 and S_3
- (D) S_1 and S_4

6. Consider the following transactions with data items P and Q initialized to '0':

T_1 : $read(P)$

$Read(Q)$

if $p = 0$ then $Q = Q + 1$

$Write(Q)$

T_2 : $read(Q)$

$Read(P)$

if $Q = 0$ then $p = p + 1$

$Write(P)$

Any non-serial interleaving of T_1 and T_2 for concurrent execution leads to

- (A) a serializable schedule
- (B) a schedule that is not conflict serializable
- (C) a conflict serializable schedule
- (D) a schedule for which a precedence graph cannot be drawn

7. Consider the concurrent execution of two transactions T_1 and T_2 , if the initial values of x, y, M and N are 200, 100, 10, 20 respectively. What are the final values of x and y ?

T_1	T_2
$read-item(x)$	
$x = x - N$	
	$read-item(x)$
	$x = x + M$
$Write-item(x)$	
$read-item(y)$	
	$Write-item(x)$
$y = y + N$	
$Write-item(y)$	

- (A) 220, 110
- (B) 210, 120
- (C) 220, 120
- (D) 210, 110

8. For the above data, if the transactions are executed in serial manner, what would be the values of X and Y at the end of the serial execution of T_1 and T_2 ?

T_1	T_2
$Read-item(x)$	
$X = X - N$	
$Write-item(x)$	
$Read-item(y)$	
$Y = Y + N$	
$Write-item(y)$	
	$Read-item(x)$
	$X = X + M$
	$Write-item(x)$

- (A) 190, 120
- (B) 180, 120
- (C) 190, 110
- (D) 180, 110

9. Consider the given two transactions T_1 and T_2 :

T_1 : $r_1(x), w_1(x), r_1(y)$

$T_2: r_2(x), r_2(y), w_2(x), w_2(y)$

Which of the following schedules are complete schedules?

- (A) $r_1(x), r_2(x), w_1(x), r_1(y), r_2(y), w_2(x), w_2(y)$
- (B) $r_2(x), r_1(x), r_2(y), w_1(x), w_2(x), r_1(y), w_2(y)$
- (C) $r_1(x), r_1(y), r_2(x), r_2(y), w_1(x), w_2(x), w_2(y)$
- (D) All the above

10. Consider the given schedule with data-locks on data-items, check whether it has dead-lock or not. The locks are shared-lock(S) and Exclusive-lock(X). Shared-lock is also called Read-lock, Exclusive-lock is also called Write-lock. Read and Write operations are denoted by R and W, respectively.

T_1	T_2	T_3	T_4
S(A)			
R(A)			
	X(B)		
	W(B)		
S(B)		S(C)	
		R(C)	
	X(C)		
			X(B)
		X(A)	

Which of the following is incorrect?

- (A) $T_1 \rightarrow T_2$
- (B) $T_3 \rightarrow T_1$
- (C) $T_2 \rightarrow T_3$
- (D) $T_4 \rightarrow T_3$

11. Consider the three transactions T_1, T_2 and T_3 and the schedule S_1 as given below. Draw the serializability (precedence) graph for S_1 , and state whether the schedule is serializable or not. If a schedule is serializable, which one of the following is equivalent serial schedule?

$T_1: r_1(x), r_1(z), w_1(x)$

$T_2: r_2(z), r_2(y), w_2(z), w_2(y)$

$T_3: r_3(x), r_3(y), w_3(y)$

$S_1: r_1(x), r_2(z), r_1(z), r_3(x), r_3(y), w_1(x), w_3(y), r_2(y), w_2(z), w_2(y)$?

- (A) $r_3(x), r_3(y), w_3(y), r_1(x), r_1(z), w_1(x), r_2(z), r_2(y), w_2(z), w_2(y)$
- (B) $r_1(x), r_1(z), w_1(x), r_2(z), r_2(y), w_2(z), w_2(y), r_3(x), r_3(y), w_3(y)$
- (C) $r_2(z), r_2(y), w_2(z), w_2(y), r_3(x), r_3(y), w_3(y), r_1(x), r_1(z), w_1(x)$
- (D) $r_2(z), r_2(y), w_2(z), w_2(y), r_1(x), r_1(z), w_1(x), r_3(x), r_3(y), w_3(y)$

12. Consider the data given in the above question. Draw the precedence graph for S_2 and state whether each schedule is serializable or not. If a schedule is serializable, which of the following is equivalent serial schedule?

$S_2: r_1(x), r_2(z), r_3(x), r_1(z), r_2(y), r_3(y), w_1(x), w_2(z), w_3(y), w_2(y)$

- (A) $r_3(x), r_3(y), w_3(y), r_1(x), r_1(z), w_1(x), r_2(z), r_2(y), w_2(z), w_2(y)$
- (B) $r_1(x), r_1(z), w_1(x), r_2(z), r_2(y), w_2(z), w_2(y), r_3(x), r_3(y), w_3(y)$
- (C) $r_2(z), r_2(y), w_2(z), w_2(y), r_3(x), r_3(y), w_3(y), r_1(x), r_1(z), w_1(x)$
- (D) $r_2(z), r_2(y), w_2(z), w_2(y), r_1(x), r_1(z), w_1(x), r_3(x), r_3(y), w_3(y)$

13. Consider schedule S_3 , which is a combination of transactions T_1, T_2 and T_3 from Q. No.11.

$S_3: r_1(x), r_2(z), r_1(z), r_3(x), r_3(y), w_1(x), c_1, w_3(y), c_3, r_2(y), w_2(z), w_2(y), c_2$?

Which of the following is true?

- (A) Recoverable and conflict serializable
- (B) Recoverable but not conflict serializable
- (C) Conflict serializable but not Recoverable
- (D) Not recoverable and not conflict serializable

14. Consider the given schedule:

$S_4: r_1(x), r_2(z), r_1(z), r_3(x), r_3(y), w_1(x), w_3(y), r_2(y), w_2(z), w_2(y), c_1, c_2, c_3$;

Which of the following is true?

- (A) Recoverable and conflict serializable
- (B) Recoverable but not conflict serializable
- (C) Conflict serializable but not Recoverable
- (D) Not recoverable and not conflict serializable

15. Which of the following is correct for the below compatibility matrix?

Mode of Locks Currently held by other transactions		Shared-Lock	Exclusive-Lock
S			
X			

S – shared - Lock, X – Exclusive - Lock

- (A)

	S	X
S	No	No
X	Yes	No

- (B)

	S	X
S	Yes	No
X	No	No

- (C)

	S	X
S	Yes	Yes
X	No	No

- (D)

	S	X
S	No	Yes
X	No	No

16. Consider the following schedule with locking:

T_1	T_2
Lock - X(A)	
R(A)	
W(A)	
	Lock - X(B)
	R(B)
	W(B)
	Lock - X(A)
	Lock - X(B)

Which of the following is true?

- (A) schedule is in Dead-Lock state
- (B) schedule is conflict serializable
- (C) schedule is not conflict serializable
- (D) Both A and B

17. Consider the given set of transactions:

T_1	T_2
	SELECT AVG (balance) FROM Account
INSERT INTO Account VALUES (487, 2000); COMMIT	
	SELECT AVG (balance) FROM Account COMMIT

The above problem is a case of

- (A) READ UNCOMMITTED
- (B) RAD COMMITTED
- (C) REPEATABLE READ
- (D) DIRTY READ

18. Consider the given set of transactions

T_1	T_2
UPDATE ACCOUNT SET balance = balance - 1000 WHERE number = 586;	
	SELECT AVG (balance)

ROLL BACK	FROM Account COMMIT
-----------	------------------------

The above problem is a case of

- (A) READ UNCOMMITTED
- (B) READ COMMITTED
- (C) DIRTY READ
- (D) BOTH A and C

19. Consider the following set of transactions

T_1	T_2
	SELECT AVG (balance) FROM Account
UPDATE Account SET balance = balance - 4000 WHERE number = 586; COMMIT	
	SELECT AVG (balance) FROM Account COMMIT

The above problem is a case of

- (A) READ UNCOMMITTED
- (B) READ COMMITTED
- (C) REPEATABLE READ
- (D) DIRTY READ

20. Consider the following schedule with locks on data items:

T_1	T_2	T_3
X(A)		
	X(A)	
		X(A)
S(B)		
	S(B)	

Which of the following is incorrect?

- (A) $T_2 \rightarrow T_1$
- (B) $T_3 \rightarrow T_2$
- (C) $T_3 \rightarrow T_1$
- (D) $T_1 \rightarrow T_3$

Practice Problem 2

Directions for questions 1 to 20: Select the correct alternative from the given choices.

1. Which of the following is false with respect to B^+ - trees of order p ?
 - (A) Each internal node has at most p tree pointers.

(B) Each leaf node has at most $\left\lceil \left(\frac{p}{2} \right) \right\rceil$ values.

(C) Each internal node, except the root, has at least

$\left\lceil \left(\frac{p}{2} \right) \right\rceil$ tree pointers.

(D) All leaf nodes are at same level.

2. Consider below transactions:

T_1	T_2
Read – item(X); $X := X - N;$	
	Read – item (X); $X := X + M;$
Write – item(X); Read – item (Y);	
	Write – item (X);
$Y := Y + N;$ Write – item (Y);	

Which of the following problem will occur during the concurrent execution of the above transactions?

- (A) Lost update problem because of incorrect X .
 - (B) Lost update problem because of incorrect Y .
 - (C) Dirty read problem because of incorrect X .
 - (D) Dirty read problem because of incorrect Y .
3. Consider the scheduled:
- $S: r_1(X); r_2(X); w_1(X); r_1(Y); w_2(X); C_2; w_1(Y); C_1;$
- This schedule is
- (A) Recoverable
 - (B) Non-recoverable
 - (C) Strict schedule
 - (D) Both (A) and (C)

4. Consider below schedule:

T_1	T_2
Read – item(X); $X := X - N;$	
	Read – item (X); $X := X + M;$
Write – item(X); Read – item (Y);	
	Write – item (X);
$Y := Y + N;$ Write – item (Y);	

This schedule is

- (A) Serializable
 - (B) Not serializable
 - (C) Under dead lock
 - (D) Both (B) and (C)
5. Let, current number of file records = r
 maximum number of records = bfr
 current number of file buckets = N
 Then what will be the file load factor?
- (A) $\frac{r}{(bfr * N)}$
 - (B) $r + (bfr * N)$
 - (C) $r * (bfr * N)$
 - (D) $r * (bfr + N)$

6. Match the following:

LIST I		LIST II	
1. Primary index	A. Ordered key field		
2. Clustering index	B. Non-ordered field		
3. Secondary index	C. Ordered non-key field		

- (A) 1 – A, 2 – B, 3 – C
 - (B) 1 – A, 2 – C, 3 – B
 - (C) 1 – C, 2 – B, 3 – A
 - (D) 1 – C, 2 – A, 3 – B
7. Consider a file with 30,000 fixed length records of size 100 bytes stored on a disk with block size 1024 bytes. Suppose that a secondary index on a non-ordering key field is constructed with key field size 9 bytes and block pointer 6 bytes. What will be the number of blocks needed for the index?
- (A) 68
 - (B) 442
 - (C) 1500
 - (D) 3000
8. Match the following:

Index type	Number of Index entries
1. Primary Index	A. Blocks in data file
2. Clustering index	B. Record in data file
3. Secondary key index	C. Distinct index filed values

- (A) 1 – A, 2 – B, 3 – C
 - (B) 1 – A, 2 – C, 3 – B
 - (C) 1 – C, 2 – B, 3 – A
 - (D) 1 – C, 2 – A, 3 – B
9. Which of the following is true with respect to B – Tree of order p ?
- (A) Each node has at most p tree pointers.
 - (B) Each node, except the root and leaf nodes, has at least $\left\lceil \frac{p}{2} \right\rceil$ tree pointers.
 - (C) All leaf nodes are at the same level.
 - (D) All of these.
10. What is the amount of unused space in allocation of unspanned fixed records of size R on a block of size B bytes?
- (A) $B - R$
 - (B) $B - \left\lfloor \frac{B}{R} \right\rfloor$
 - (C) $B - \left(\left\lfloor \frac{B}{R} \right\rfloor * R \right)$
 - (D) $\left(\left\lfloor \frac{B}{R} \right\rfloor * R \right) - B$
11. What is the average time required to access a record in a file consisting of b blocks using unordered heap linear search?
- (A) b
 - (B) $b/2$
 - (C) \log_2^b
 - (D) b^2
12. Consider a file of fixed length records of size R bytes. If the block size is B bytes, then the blocking factor will be

- (A) $B \times R$ records (B) $\left\lfloor \frac{B}{R} \right\rfloor$ records
 (C) $\left\lceil \frac{B}{R} \right\rceil$ records (D) $B + R$ records

13. Consider the following relation instance:

P	Q	R
1	4	2
1	5	3
1	6	3
3	2	2

Which of the following FDs are satisfied by the instance?

- (A) $PQ \rightarrow R$ and $R \rightarrow Q$ (B) $QR \rightarrow P$ and $Q \rightarrow R$
 (C) $QR \rightarrow P$ and $P \rightarrow R$ (D) $PR \rightarrow Q$ and $Q \rightarrow P$
14. Consider an ordered file with 30,000 records stored on a disk with block size of 1024 bytes. The records are of fixed size and are of unspanned, with record length 100 bytes. What is the number of accesses required to access a data file using binary search?
 (A) 10 (B) 12
 (C) 1500 (D) 3000
15. What is the blocking factor for an index if the ordering key field size is 9 bytes and block pointer is 6 bytes long, and the disk block size is 1024 bytes?

- (A) 114 (B) 171
 (C) 341 (D) 68
16. For a set of n transactions, there exist _____ different valid serial schedules
 (A) n (B) n^2
 (C) $n/2$ (D) $n!$
17. The number of possible schedules for a set of n transactions is
 (A) lesser than $n!$ (B) much larger than $n!$
 (C) $n!$ (D) None
18. Which one of the following is conflict operation?
 (A) Reads and writes from the same transaction
 (B) Reads and writes from different transaction
 (C) Reads and writes from different transactions on different data items.
 (D) Reads and writes from different transaction on same data.
19. The following schedule $S: r_3(x), r_2(x), w_3(x), r_1(x), w_1(x)$ is conflict equivalent to serial schedule
 (A) $T_1 \rightarrow T_3 \rightarrow T_2$ (B) $T_2 \rightarrow T_1 \rightarrow T_3$
 (C) $T_1 \rightarrow T_2 \rightarrow T_3$ (D) None
20. The following schedule $S: R_1(x), R_2(x), W_1(x), W_2(x)$ is
 (A) Conflict serializable (B) View serializable
 (C) Both (D) None

PREVIOUS YEARS' QUESTIONS

1. Consider the following four schedules due to three transactions (indicated by the subscript) using *read* and *write* on a data item x , denoted by $r(x)$ and $w(x)$, respectively. Which one of the them is conflict serializable? [2014]
 (A) $r_1(x); r_2(x); w_1(x); r_3(x); w_2(x)$
 (B) $r_2(x); r_1(x); w_2(x); r_3(x); w_1(x)$
 (C) $r_3(x); r_2(x); r_1(x); w_2(x); w_1(x)$
 (D) $r_2(x); w_2(x); r_3(x); r_1(x); w_1(x)$
2. Consider the following schedule S of transactions T_1, T_2, T_3, T_4 :

T_1	T_2	T_3	T_4
	Reads (X)		
		Writes (X) Commit	
Writes (X) Commit			
	Writes (Y) Reads (Z) Commit		
			Reads (X) Reads (Y) Commit

Which one of the following statements is correct? [2014]

- (A) S is conflict-serializable but not recoverable
 (B) S is not conflict-serializable but is recoverable
 (C) S is both conflict-serializable and recoverable
 (D) S is neither conflict-serializable nor it is recoverable
3. Consider the following transaction involving two bank accounts x and y .
 $read(x); x := x - 50; write(x); read(y); y := y + 50; write(y)$
 The constraint that the sum of the accounts x and y should remain constant is that of [2015]
 (A) Atomicity (B) Consistency
 (C) Isolation (D) Durability
4. Consider a simple checkpointing protocol and the following set of operations in the log.
 (start, T_4); (write, $T_4, y, 2, 3$); (start, T_1); (commit, T_4);
 (write, $T_1, z, 5, 7$);
 (checkpoint);
 (start, T_2); (write, $T_2, x, 1, 9$); (commit, T_2); (start, T_3);
 (write, $T_3, z, 7, 2$);
 If a crash happens now and the system tries to recover using both undo and redo operations. What are the contents of the undo list and the redo list? [2015]

- (A) Undo: T_3, T_1 ; Redo: T_2
- (B) Undo: T_3, T_1 ; Redo: T_2, T_4
- (C) Undo: none; Redo: T_2, T_4, T_3, T_1
- (D) Undo: T_3, T_1, T_4 ; Redo: T_2

5. Consider the following partial schedule S involving two transactions T_1 and T_2 . Only the read and the write operations have been shown. The read operation on data item P is denoted by $read(P)$ and the write operation on data item P is denoted by $write(P)$

Time Instance	Transaction – id	
	T_1	T_2
1	read(A)	
2	write(A)	
3		read(C)
4		write(C)
5		read(B)
6		write(B)
7		read(A)
8		commit
9	read(B)	

Schedule S

Suppose that the transaction T_1 fails immediately after time instance 9. Which one of the following statements is correct? [2015]

- (A) T_2 must be aborted and then both T_1 and T_2 must be re-started to ensure transaction atomicity.
 - (B) Schedule S is non-recoverable and cannot ensure transaction atomicity.
 - (C) Only T_2 must be aborted and then re-started to ensure transaction atomicity.
 - (D) Schedule S is recoverable and can ensure atomicity and nothing else needs to be done.
6. Which one of the following is **NOT** a part of the ACID properties of database transactions? [2016]
- (A) Atomicity
 - (B) Consistency
 - (C) Isolation
 - (D) Deadlock - freedom
7. Consider the following two phase locking protocol. Suppose a transaction T accesses (for read or write operations), a certain set of objects $\{O_1, \dots, O_k\}$. This is done in the following manner: [2016]
- Step 1. T acquires exclusive locks to O_1, \dots, O_k in increasing order of their addresses.
- Step 2. The required operations are performed.
- Step 3. All locks are released.
- This protocol will

- (A) guarantee serializability and deadlock-freedom.
- (B) guarantee neither serializability nor deadlock-freedom.
- (C) guarantee serializability but not deadlock-freedom.
- (D) guarantee deadlock-freedom but not serializability.

8. Suppose a database schedule S involves transactions T_1, \dots, T_n . Construct the precedence graph of S with vertices representing the transactions and edges representing the conflicts. If S is serializable, which one of the following orderings of the vertices of the precedence graph is guaranteed to yield a serial schedule? [2016]

- (A) Topological order
- (B) Depth - first order
- (C) Breadth - first order
- (D) Ascending order of transaction indices

9. Consider the following database schedule with two transactions T_1 and T_2 .

$$S = r_2(X); r_1(X); r_2(Y); w_1(X); r_1(Y); w_2(X); a_1; a_2$$

Where $r_i(Z)$ denotes a *read* operation by transaction T_i on a variable Z , $w_i(Z)$ denotes a *write* operation by T_i on a variable Z and a_i denotes an *abort* by transaction T_i .

Which one of the following statements about the above schedule is **TRUE**? [2016]

- (A) S is non - recoverable
- (B) S is recoverable, but has a cascading abort
- (C) S does not have a cascading abort
- (D) S is strict.

10. In a database system, unique timestamps are assigned to each transaction using Lamport's logical clock. Let $TS(T_1)$ and $TS(T_2)$ be the timestamps of transactions T_1 and T_2 respectively. Besides, T_1 holds a lock on the resource R, and T_2 has requested a conflicting lock on the same resource R. The following algorithm is used to prevent deadlocks in the database system assuming that a killed transaction is restarted with the same timestamp.

```

if  $TS(T_2) < TS(T_1)$  then
     $T_1$  is killed
else  $T_2$  waits.
    
```

Assume any transaction that is not killed terminates eventually. Which of the following is TRUE about the database system that uses the above algorithm to prevent deadlocks? [2017]

- (A) The database system is both deadlock-free and starvation-free.
- (B) The database system is deadlock-free, but not starvation-free.

- (C) The database system is starvation-free, but not deadlock-free.
 (D) The database system is neither deadlock-free nor starvation-free.

11. Two transactions T_1 and T_2 are given as

$$T_1 : r_1(X)w_1(X)r_1(Y)w_1(Y)$$

$$T_2 : r_2(Y)w_2(Y)r_2(Z)w_2(Z)$$

where $r_i(V)$ denotes a *read* operation by transaction T_i on a variable V and $w_i(V)$ denotes a *write* operation by transaction T_i on a variable V . The total number of conflict serializable schedules that can be formed by T_1 and T_2 is _____. [2017]

ANSWER KEYS

Practice Problem I

1. B 2. B 3. B 4. B 5. C 6. B 7. B 8. A 9. A 10. D
 11. A 12. A 13. A 14. C 15. B 16. D 17. C 18. D 19. B 20. D

Practice Problem I

1. B 2. A 3. A 4. D 5. A 6. B 7. B 8. B 9. D 10. C
 11. B 12. B 13. B 14. B 15. D 16. D 17. C 18. D 19. A 20. D

Previous Years' Questions

1. D 2. C 3. B 4. A 5. B 6. D 7. A 8. A 9. C 10. A
 11. 54

Chapter 5

File Management

LEARNING OBJECTIVES

- Files
- Memory hierarchies
- Description of disk devices
- File records
- Sorted files
- Hashing techniques
- Extendible hashing
- Index update
- Clustering index
- B-Trees
- B+Trees
- Over flow in internal node

FILES

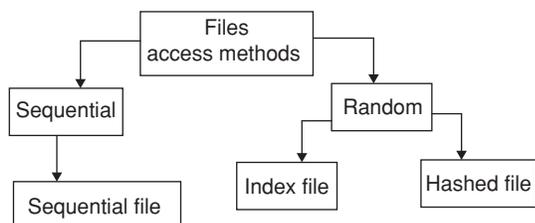
Databases are stored on magnetic disks as files of records. Computer storage media form a storage hierarchy that includes two main categories.

Primary storage This category includes storage media that can be operated on, directly by CPU, such as the computer main memory and cache memory. Primary storage provides fast access but is of limited storage capacity.

Secondary storage This category includes magnetic disks, optical disks, and tapes. These devices usually have a larger capacity, less cost, and slower access to data. Data in secondary storage cannot be processed directly by the CPU, it must be copied into primary storage.

File Structure

Taxonomy of file structure



Sequential file A sequential file is one in which records can only be accessed sequentially, one after another from beginning to end. Records are stored contiguously on the storage device.

Index files These files are used to access a record in the file. The entire index file is loaded into main memory data and indexes are stored in the same file. The term 'index file' is used as a synonym for the term 'database file'. The index file contains parameters that specify the name and location of file used to store DB.

Indexing Indexing mechanism is used to speed up access to desired data. An index file consists of records (called *index entries*) of the form.

Search-key	Pointer
------------	---------

Index files are typically much smaller than the original file.

Ordered indices In ordered index, index entries are stored, sorted on the search-key value.

Example: Author catalogue in library.

MEMORY HIERARCHIES

At the primary storage level, the memory hierarchy includes cache memory which is a static RAM.

The next level of primary storage is DRAM (dynamic RAM) which provides the main work area for the CPU for keeping programs and data and is called the *main memory*.

At the secondary storage level, the hierarchy includes magnetic disks, as well as mass storage in the form of CD-ROM (compact disk read-only memory) and tapes. Programs reside in DRAM and large permanent databases reside on secondary storage.

Another form of memory, *flash memory*, is non-volatile. Flash memories are high-density, high-performance memories using EEPROM (electrically erasable programmable read-only

memory) technology. The advantage of flash memory is the fast access speed, the disadvantage is that an entire block must be erased and written over at a time. Finally, magnetic tapes are used for archiving and backup storage of data.

DESCRIPTION OF DISK DEVICES

Magnetic disks are used for storing large amounts of data. The capacity of a disk is the number of bytes it can store. A disk is single sided if it stores information on only one of its surfaces and double sided if both surfaces are used. To increase storage capacity, disks are assembled into a disk pack, which may include many disks and hence many surfaces. Information is stored on a disk surface in concentric circles with small width, each having a distinct diameter. Each circle is called a *track*. For disk packs, the tracks with the same diameter on the various surfaces are called a *cylinder* because of the shape they would form if connected in space.

A track usually contains a large amount of information; it is divided into smaller blocks (or sectors). The division of track into equal-sized disk blocks (or pages) is set by the operating system during disk formatting.

Blocks are separated by fixed-size inter-block gaps, which include specially coded control information written during disk initialization. This information is used to determine which block on the track follows each inter block gap. Transfer of data between main memory and disk takes place in units of disk blocks. The hardware address of a block is the combination of a cylinder number, track number and block number is supplied to the disk I/O hardware.

The actual hardware mechanism that reads or writes a block is the disk read/write head, which is part of a system called a *disk drive*. A disk is mounted in the disk drive, which includes a motor that rotates the disk. To transfer a disk block, given its address, the disk controller must first mechanically position the read/write head on the correct track. The time required to do this is called the *seek time*. There is another delay called *rotational delay* or *latency*; the beginning of the desired block rotates into position under the read/write head. It depends on the RPM of the disk. Finally, some additional time is needed to transfer the data, which is called *block-transfer time*. Hence, the total time needed to locate and transfer an arbitrary block, given its address is the sum of the seek time, rotational delay and block transfer time. The seek time and rotational delay are usually much larger than the block transfer time.

FILE RECORDS

Data is usually stored in the form of *records*. Each record consists of a collection of related data values or items where each value is of one or more bytes and corresponds to a particular field of the record. Records describe entities and their attributes.

Record type A collection of field names and their corresponding data types constitutes a record type (or) record format.

A file is a sequence of records. If every record in the file has exactly the same size (in bytes), the file is said to be made up of fixed-length records. If different records in the file have different sizes, the file is said to be made up of variable-length records.

Spanned Versus Unspanned Records

The records of a file must be allocated to disk blocks because a block is the unit of data transfer between disk and memory. When the block size is larger than the record size, each block will contain numerous records, although some files may have unusually large records that cannot fit in one block.

Suppose that the block size is B bytes. For a file of fixed-length records of size R bytes, with $B \geq R$, we can fit

$$bfr = \lfloor B/R \rfloor \text{ records per block}$$

The value bfr is called the *blocking factor* for the file. Some times R may not divide B exactly, so we have some unused space in each block equal to $B - (bfr * R)$ bytes. To utilize this unused space, we can store part of a record on one block and the rest on another. A pointer at the end of the first block points to the block containing the remainder of the record in case it is not the next consecutive block on disk. This organization is called *spanned*, because records can span more than one block. Whenever a record is larger than a block, we must use a spanned organization. If records are not allowed to cross block boundaries, the organization is called *unspanned*. This is used with fixed-length records having $B > R$, because it makes each record start at a known location in the block. For variable-length records, either a spanned or an unspanned organization can be used.

For variable-length records using spanned organization, each block may store a different number of records. In this case, the blocking factor bfr represents the average number of records per block for the file. We can use bfr to calculate the number of blocks ' b ' needed for a file of ' r ' records.

$$b = \lceil (r/bfr) \rceil \text{ blocks}$$

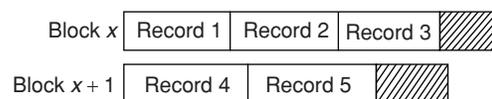


Figure 1 Unspanned records

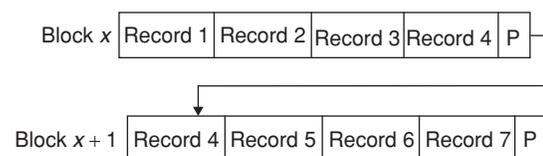


Figure 2 Spanned

There are several standard techniques for allocating the blocks of a file on disk. In contiguous allocation, the file blocks are allocated to consecutive disk blocks. In linked

allocation, each file block contains a pointer to the next file block. A combination of the two allocates clusters of consecutive disk blocks, and the clusters are linked. Clusters are sometimes called file *segments* (or) *extents*. Another possibility is to use indexed allocation, where one or more index blocks contain pointers to the actual file blocks.

SORTED FILES (ORDERED RECORDS)

We can physically order the records of a file on disk based on the values of the one of their fields called the *ordering field*. This leads to an ordered or sequential file. If the ordering field is also a key field of the file, a field guaranteed to have a unique value in each record, then the field is called the *ordering key for the file*.

Advantages

1. Reading the records in order of the ordering key values becomes extremely efficient, because no sorting is required.
2. Finding the next record from the current one in order of the ordering key usually requires no additional block access, because the next record is in the same block as the current one.
3. Using a search condition based on the value of an ordering key field results in faster access when the binary search technique is used. This constitutes an improvement over linear searches, although it is not often used for disk files.

A binary search for disk files can be done on the blocks rather than on the records. Suppose that a file has ‘*b*’ blocks numbered 1, 2, ..., *b*, the records are ordered by ascending value of their ordering key field and we are searching for a record whose ordering key field value is *K*. Assuming that disk addresses of the file blocks are available in the file header, the binary search usually accesses $\log_2^{(b)}$ blocks, whether the record is found (or) not, an improvement over linear searches, where, on the average, (*b*/2) blocks are accessed when the record is found and ‘*b*’ blocks are accessed when the record is not found.

Type of Organization	Access Method	Average Time to Access a Specific Record
Heap (unordered)	Sequential scan (linear search)	<i>b</i> /2
Ordered	Ordered scan	<i>b</i> /2
Ordered	Binary search	$\log_2 b$

Ordered files are rarely used in database applications unless an additional access path, called a *primary index*, is used; this results in an indexed sequential file. This further improves the random access time on the ordering key field.

HASHING TECHNIQUES

The other type of primary file organization is based on hashing, which provides very fast access to records on certain search conditions. This organization is usually called a *hash file*.

The search condition must be an equality condition on a single field, called the *hash field* of the file. If the hash is also a key field of the file, in which case it is called the *hash key*.

The idea behind hashing is to provide a function ‘*h*’, called a hash function or randomizing function, which is applied to the hash field value of a record and yields the address of the disk block in which the record is stored. We need only a single-block access to retrieve that record.

Example:

	NAME	RNO	CLASS	GRADE
0				
1				
2				
3				
⋮				
⋮				
⋮				
<i>m</i> -2				
<i>m</i> -1				

Internal Hashing

For internal files, hashing is implemented as a hash table through the use of an array of records. Suppose that the array index range is from 0 to *M* - 1, then we have *M* slots whose addresses corresponds to the array indexes. We choose a hash function that transforms the hash field value into an integer between 0 and *M* - 1. One common hash function is the $h(K) = K \text{ mod } M$ function, which returns the remainder of an integer hash field value *K* after division by *M*; this value is then used for the record address.

Non-integer hash field values can be transformed into integers before the mod function is applied. For character strings, the numeric (ASCII) codes associated with characters can be used in the transformation.

A collision occurs when the hash field value of a record that is being inserted hashes to an address that already contains a different record. In this situation, we must insert the new record in some other position, since its hash address is occupied. The process of finding another position is called *collision resolution*. There are different methods for collision resolution as follows:

Open addressing Proceeding from the occupied position specified by the hash address, the program checks the subsequent positions in order until an unused (empty) position is found.

Chaining For this method, various overflow locations are kept, usually by extending the array with a number of overflow positions. In addition, a pointer field is added to each record location. A collision is resolved by placing the new record in an unused overflow location and setting the pointer of the occupied hash address location to the address of that overflow location. A linked list of overflow records for each hash address is thus maintained.

Multiple hashing The program applies a second hash function if the first results in a collision. If another collision results, the program uses open addressing or applies a third hash function and then uses open addressing if necessary.

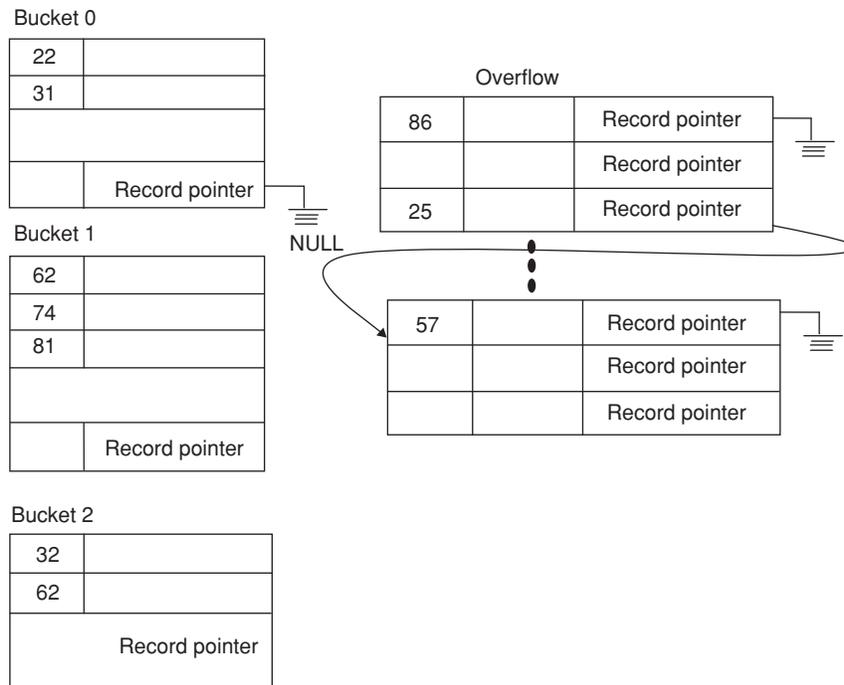
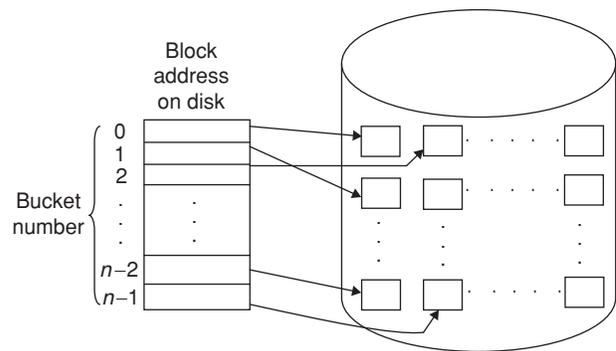
If we expect to have ‘*r*’ records to store in the table, we should choose *M* locations for the address space such that (*r*/*M*) is between 0.7 and 0.9. It may also be useful to choose a prime number for *M*, since it has been demonstrated that this distributes the hash addresses better over the address space when the ‘mod’ hashing function is used. Other hash functions may require *M* to be a power of 2.

External Hashing

Hashing for disk files is called *external hashing*. To suit the characteristics of disk storage, the target address space is made of buckets, each of which holds multiple records. A bucket is either one disk block or a cluster of contiguous

blocks. The hashing function maps a key into a relative bucket number, rather than assigning an absolute block address to the bucket. A table maintained in the file header converts the bucket number into the corresponding disk block address.

The collision problem is less severe with buckets, because as many records as will fit in a bucket can hash to the same bucket without causing problems. If the capacity of bucket exceeds, we can use a variation of chaining in which a pointer is maintained in each bucket to a linked list of overflow records for the bucket. The pointers in the linked list should be record pointers, which include both a block address and a relative record position within the block.



The hash function is $h(k) = k \text{ mod } 10$ and the hashing scheme described above is called *static hashing* because a fixed number of buckets *M* is allocated. It can be a drawback for dynamic files. Suppose that we allocate *M* buckets for the address space and let ‘*m*’ be the maximum number of records that can fit in one bucket, then at most (*m* * *M*)

records will fit in the allocated space. If the number of records turns out to be substantially fewer than (*m* * *M*), we are left with a lot of unused space.

If the number of records increases to substantially more than (*m* * *M*), numerous collisions will result and retrieval will be slowed down because of the long lists of overflow

records. In either case, we may have to change the number of blocks M allocated and then use a new hashing function (based on the new value of M) to redistribute the records. These organizations can be quite time consuming for large files. Newer dynamic file organizations based on hashing allows the number of buckets to vary dynamically with only localized reorganization.

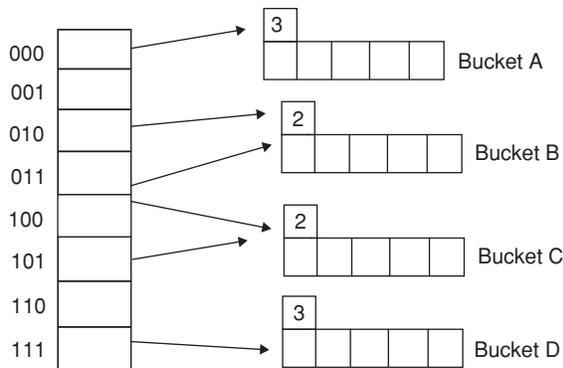
Hashing Techniques with Dynamic File Expansion

The disadvantage of static hashing is that the hash address space is fixed. Hence, it is difficult to expand or shrink the file dynamically. The first scheme is extendible hashing. It stores an access structure in addition to the file hence it is similar to indexing. The main difference is that the access structure is based on the values that result after application of the hash function to the search field. In indexing, the access structure based on the values of the search field itself. The second technique, called *linear hashing*, does not require additional access structure.

These hashing schemes take advantage of the fact that the result of applying a hashing function is a non-negative integer and hence can be represented as a binary number. The access structure is built on the binary representation of the hashing function result, which is a string of bits. We call this the *hash value of a record*. Records are distributed among buckets based on the values of the leading bits in their hash values.

Extendible Hashing

In extendible hashing, a type of directory, an array of 2^d bucket addresses is maintained, where d is called the *global depth of the directory*. The integer value corresponding to the first (high-order) d bits of a hash value is used as an index to the array to determine a directory entry, and the address in that determines the bucket in which the corresponding records are stored. Several directory locations with the same first d' bits for their hash values may contain the same bucket address if all the records that hash to these locations fit in a single bucket. A local depth d' is stored with each bucket specifies the number of bits on which the bucket contents are based.



The value of d can be increased and decreased by one at a time, thus doubling or halving the number of entries in the directory array. Doubling is needed if a bucket, whose local depth d' is equal to the global depth d , overflows. Halving occurs if $d > d'$ for all the buckets after some locations occur. Most record retrievals require two block accesses: one to the directory and the other to the bucket.

The main advantage of extendible hashing is the performance of the file does not degrade as the file grows, as opposed to static external hashing where collisions increases and the corresponding chaining causes additional accesses. No space is allowed in extendible hashing for future growth, but additional buckets can be allocated dynamically as needed. The space overhead for the directory table is negligible.

Another advantage is that splitting causes minor reorganization in most cases, since only the records in one bucket are redistributed to the two new buckets. The only time a reorganization is more expensive is when the directory has to be doubled (or) halved.

A disadvantage is that the directory must be searched before accessing the buckets themselves, resulting in two block accesses instead of one in static hashing.

INDEXING

Indexes are auxiliary access structures, which are used to speed up the retrieval of records in response to certain search conditions. The index structure typically provides secondary access paths, which provide alternative ways of accessing the records without affecting the physical placement of records on disk. They enable efficient access to records based on the indexing fields that are used to construct the index.

Any field of the file can be used to create an index and multiple indexes on different fields can be constructed on the same file. To find a record or records in the file based on a certain selection criterion on an indexing field, one has to initially access the index, which points to one or more blocks in the file where the required records are located. The most prevalent types of indexes are based on ordered files (single-level indexes) and tree data structures (multilevel indexes, B^+ trees).

Dense Index files: Index record appears for every search-key value in the file.

Brighton		A-217	Brighton	750
Downtown		A-101	Downtown	600
Mianus		A-110	Downtown	300
Perryridge		A-215	Mianus	400
		A102	Perryridge	800

Figure 3 Dense index file.

Sparse index files These files contain index records for only some search-key values. Applicable when records are sequentially ordered on search key.

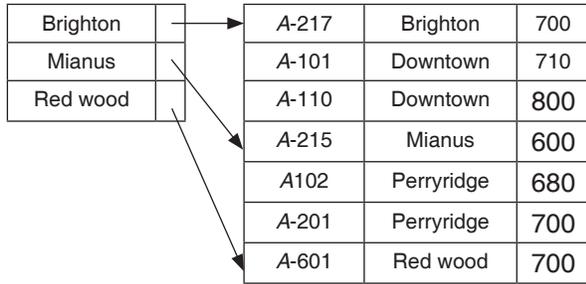


Figure 4 Sparse index file.

Compared to dense index, sparse index takes less space and less maintenance overhead for insertions and deletions. It is slower than dense index for locating records.

Index Update

Record deletion If delete key was the only record in the file with its particular search-key value, the search key is deleted from the index also.

In dense index, delete the search key.

In sparse index, if deleted key value exists in the index, the value is replaced by next search-key value in the file. If the next search-key value already has an index entry, the entry is deleted instead of being replaced.

Record insertion In dense index, if the search-key value doesn't appear in the index insert it.

If index stores an entry for each block of the file, no change needs to be made to the index unless a new block is created. If a new block is created, the first search-key value appearing in the new block is inserted into the index.

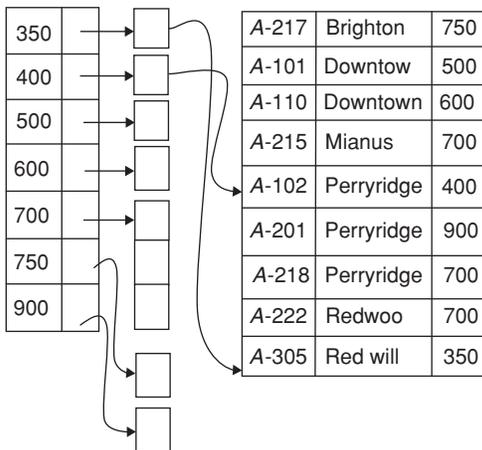


Figure 5 Secondary index.

Secondary index example:

1. Index record points to a bucket that contains pointers to all the actual records with that particular search – key value
2. secondary index have to be dense

Single-level Ordered Indexes

A file with a given record structure consisting of several fields (or attributes), an index access structure is usually defined on a single field of a file is called an *indexing field* or *indexing attribute*. The index typically stores each value of the index field along with a list of pointers to all disk blocks that contain records with that field value. The values in the index are ordered so that we can do a binary search on the index.

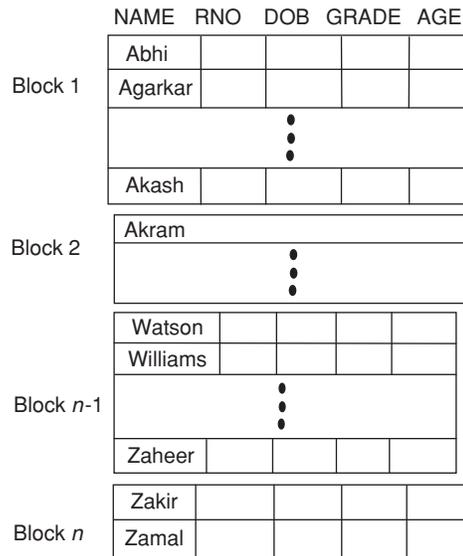
The index file is much smaller than the data file, so searching the index using a binary search is reasonably efficient. Multilevel indexing does away the need for a binary search at the expense of creating indexes to the index itself.

Types of Ordered Indexes

1. Primary index
2. Clustering index
3. Secondary index

Primary index A primary index is an ordered file whose records are of fixed length with two fields. The first field is of the same data type as the ordering key field called the *primary key of the data file*, and the second field is a pointer to a disk block (block address). There is one index entry (index record) in the index file for each block in the data file. Each index entry has the value of the primary key field for the record in a block and a pointer to that block as its two field values. The two field values of index entry i is $\langle k(i), p(i) \rangle$.

Example:



To create a primary index on the ordered file shown in the above figure, we use the NAME field as primary key, because that the ordering key field on the file (assuming that each value of NAME is unique). Each entry in the index has a NAME value and a pointer. Some sample index entries are as follows:

< $k(1) = (\text{Abhi}), p(1) = \text{address of block 1}$ >
 < $k(2) = (\text{Akram}), p(2) = \text{address of block 2}$ >
 < $k(3) = (\text{Brat}), p(3) = \text{address of block 3}$ >

The below figure illustrates this primary index. The total number of entries in the index is the same as the number of disk blocks in the ordered data file. The first record in each block of the data file is called the *anchor record of the block (or) block anchor*.

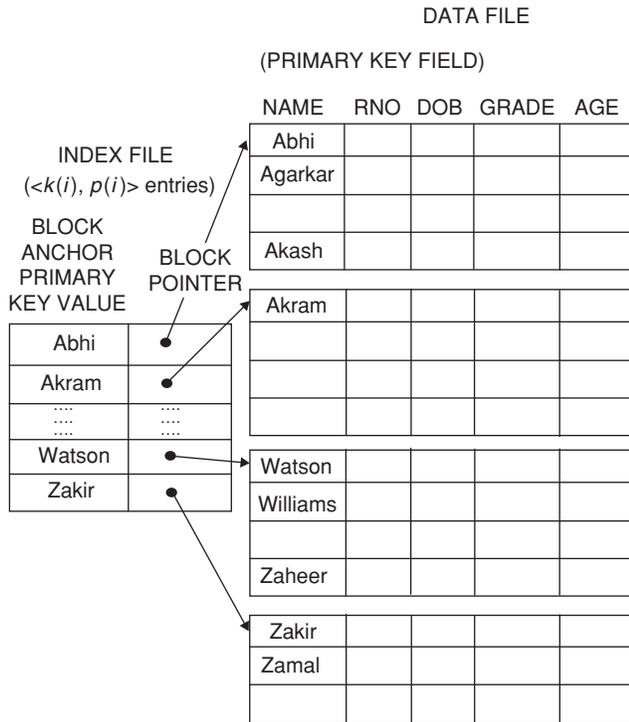


Figure 6 Primary index on the ordering key field of the file.

Indexes can also be characterized as dense or sparse. A dense index has an index entry for every search-key value (every record) in the data file. A sparse (non-dense) index has index entries only for some of the search values. A primary index is non-dense (sparse) index, since it includes an entry for each disk block of the data file and the keys of its anchor record rather than for every search value.

The index file for primary index needs fewer blocks than does the data file, for two reasons as follows:

1. There are fewer index entries than there are records in the data file.
2. Each index entry is typically smaller in size than a data record because it has only two fields. So more index entries than data records can fit in one block.

A binary search on the index file requires fewer block accesses than a binary search on the data file. The binary search for an ordered data file required \log_2^b block accesses. But if the primary index file contains b_i blocks, then to locate a record with a search-key value requires a binary search of that index and access to the block containing that record, a total of $\log_2^b b_i + 1$ accesses.

A record whose primary key value is k lies in the block whose address is $p(i)$, where $k(i) \leq k \leq k(i + 1)$. The i th block in the data file contains all such records because of the physical ordering of the file records on the primary key field. To retrieve a record, given the value k of its primary key field, we do a binary search on the index file to find the appropriate index entry i , and then retrieve the data file block whose address is $p(i)$.

The following example illustrates the saving in block accesses that is attainable when a primary index is used to search for a record.

Example: Suppose that we have an ordered file with $r = 24,000$ records stored on a disk with block size $B = 512$ bytes. File records are of fixed size and are unspanned, with record length $R = 120$ bytes.

The blocking factor for the file would be $bfr = \lfloor B/R \rfloor$

$$= \left\lfloor \frac{512}{120} \right\rfloor = \lfloor 4.26 \rfloor = 4 \text{ records per block}$$

The number of blocks needed for the file is

$$b = \left\lceil \left(\frac{r}{bfr} \right) \right\rceil = \left\lceil \frac{24,000}{42} \right\rceil = 6000 \text{ blocks}$$

A binary search on the data file would need

$$\lceil \log_2^b = \log_2^{6000} \rceil = 13 \text{ block accesses}$$

Example: For the above data, suppose that the ordering key field of the file is $V = 7$ bytes long, a block pointer, $P = 5$ bytes long, and we have constructed a primary index for the file.

The size of each index entry is $R_i = (7 + 5) = 12$ bytes, so the blocking factor for the index is $bfr_i = \lfloor (B/R_i) \rfloor$

$$\lfloor 512/12 \rfloor = \lfloor 42.66 \rfloor = 42 \text{ entries per block.}$$

The total number of index entries r_i is equal to number of blocks in the data file, which is 6000. The number of index blocks is hence

$$b_i = \lceil (r_i/bfr_i) \rceil = \lceil 6000/42 \rceil = 142 \text{ blocks}$$

To perform a binary search on the index file would need $\lceil \log_2^b b_i \rceil = \lceil \log_2 142 \rceil = 8$ block accesses. To search for a record using the index, we need one additional block access to the data file for a total of '9' block accesses.

Disadvantage: A major problem with a primary index is insertion and deletion of records. If we attempt to insert a record in its correct position in the data file, we have to not only move records to make space for the new record but also change some index entries.

Clustering Index If records of a file are physically ordered on a non-key field, which does not have a distinct value for each record, that field is called the *clustering field*. We can create a different type of index called *clustering index* to

speed up the retrieval of records that have the same value for the clustering field.

A clustering index is also an ordered file with two fields, the first field is of the same type as the clustering field of the data file, and the second field is a block pointer.

There is one entry in the clustering index for each distinct value of the clustering field, containing the value and a pointer to the first block in the data file that has a record with that value for its clustering field.

The record insertion and deletion still cause problems, because the data records are physically ordered. To alleviate the problem of insertion, reserve a whole block (or a cluster of contiguous blocks) for each value of the clustering field, all records with that value are placed in the block (or block cluster). A clustering index is an example of a non-dense index, because it has an entry for every distinct value of the indexing field which is a non-key.

Secondary Index A secondary index provides a secondary means of accessing a file for which some primary access already exists. The secondary index may be on a field which is a candidate key and has a unique value in every record, or a non-key with duplicate values. The index is an ordered file with two fields. The second field is either a block pointer or a record pointer. There can be many secondary indexes for the same file.

First consider a secondary index access structure on a key field that has a distinct value for every record such a field is some times called a *secondary key*.

The records of the data file are not physically ordered by values of the secondary key field, we cannot use block anchors. That is why an index entry is created for each record in the data file, rather than for each block, as in the case of a primary index.

B⁻ Trees

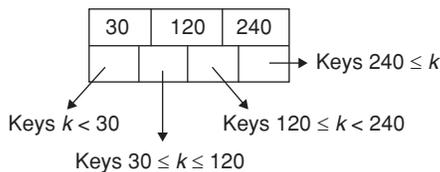
1. A commonly used index structure
2. Non-sequential, ‘balanced’
3. Adapts well to insertions and deletions
4. Consists of blocks holding at most n keys and $n + 1$ pointers.
5. We consider a variation actually called a B⁺ tree

B⁺ Trees

B⁺ trees are a variant of B⁻ trees. In B⁺ trees data stored only in leaves, leaves form a sorted linked list.

Parameter – n

Branching factor – $n + 1$



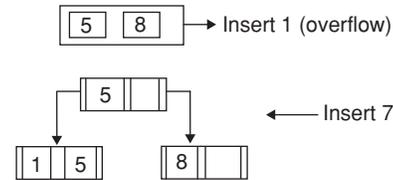
Each node (except root) has at least $n/2$ keys. B⁻ tree stands for balanced tree. All the paths through a B⁻ tree from root

to different leaf nodes are of the same length (balanced path length). All leaf nodes are at the same depth level.

This ensures that number of disk accesses required for all the searches are same. The lesser the depth (level) of an index tree, the faster the search.

Insertion into B⁺ tree

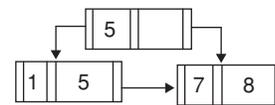
Given nodes 8 5 1 7 3 12 Initially start with root node (has no children)



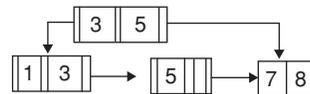
Overflow in Leaf Node

Split the leaf node First, $j = \text{ceiling}((p_{\text{leaf}} + 1)/2)$ entries are kept in the original node and the remaining moved to the new leaf.

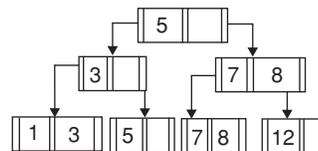
1. Create a new internal node, and j th index value is replicated in the parent internal node.
2. A pointer is added to the newly formed leaf node.



Insert 3 → overflow



Insert 12 (overflow, split propagates, new level)



Overflow in Internal Node

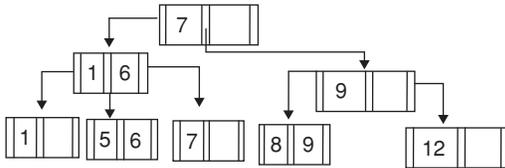
Split the internal node, the entries up to P_j where $j = \text{floor}((p + 1)/2)$ are kept in the original node and remaining moved to the new internal node

1. Create a new internal node and the j th index value is moved to the parent internal node (without replication)
2. Pointers are added to the newly formed nodes.

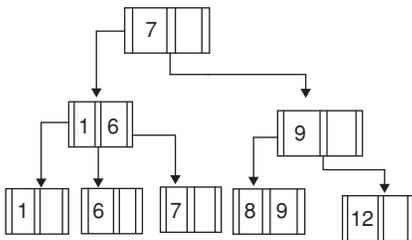
- B⁺ tree ensures some space always left in nodes for new entries. Also makes sure all nodes are at least half full.

Deletion in B⁺ Trees

Delete 5,12,9 from the below B⁺ tree:

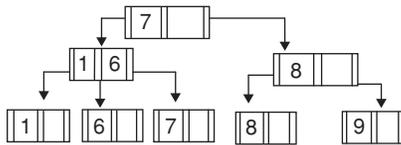


Delete 5:

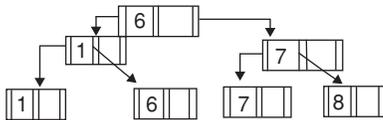


Delete 12:

Under flow has occurred, so redistribute.



Delete 9: Underflow (merge with left) redistribute.



Advantages

- B⁻ Trees and B⁺ trees: B⁻ tree is a data structure used for external memory.
- B⁻ trees are better than binary search trees if data is stored in external memory.

- Each node in a tree should correspond to a block of data.
- Each node can store many data items and has many successors.
- The B⁻ tree has fewer levels but search for an item takes more comparisons at each level.
- If a B⁻ tree has order 'd', then each node (except root) has at least d/2 children, then the depth of the tree is at most $\log_{d/2}(\text{size}) + 1$.
- In the worst case, we need (d - 1) comparisons in each node (using linear search)
- Fewer disk accesses are required compared to binary Tree.
- The usual data structure for an index is the B⁺ tree.
- Every modern DBMS contains some variant of B⁻ trees in addition with other index structures depending on the application.
- B⁻ trees and B⁺ trees are one and the same. They differ from B⁻ trees in having all data in the leaf blocks.
- Compared to binary trees, B⁻ trees will have higher branching factor.
- Binary trees can degenerate to a linear list, B⁻ trees are balanced, so this is not possible.
- In B⁺ tree, the values in inner nodes are repeated in the leaf nodes.
- The height of the tree might decrease, because the data pointer is needed only in the leaf nodes, we can also get a sorted sequence.
- In B⁻ trees, all leaves have the same distance from root hence B⁻ trees are balanced. This ensures that the chain of links followed to access a leaf node is never too long.
- The time complexity of search operation in B⁻ tree (tree height) is $O(\log n)$, where 'n' is the number of entries.
- Advantage of B⁺ tree automatically reorganizes itself with small and local changes while doing insertions and deletions, reorganization of entire file is not required to maintain performance.
- Disadvantage of B⁺ tree, extra Insertion and deletion overhead, space overhead.
- B⁺ trees can be used as dynamic multilevel Indexes.

EXERCISES

Practice Problems I

Directions for questions 1 to 20: Select the correct alternative from the given choices.

- Consider the following specifications of a disk. Block size of a disk is 512 bytes, inter-block gap size is 128 bytes. Number of blocks per track is 20 and number of tracks per surface is 400.
 - What is the capacity of disk including Inter block gap?

- (A) 124000 (B) 1260000
(C) 5120000 (D) 512000

- (ii) What is the capacity of disk excluding Inter block gap?

- (A) 25400 (B) 25600
(C) 25800 (D) 25900

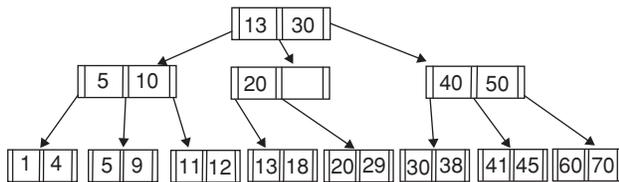
- Consider the following specifications of a disk. Block size of disk is 512 bytes, 2000 tracks per surface, 50 sectors per track and 5 double sided platters.

- (i) What is the capacity of track in bytes?
 (A) 4096000 (B) 4086000
 (C) 4076000 (D) 4066000
- (ii) What is the capacity of surface in bytes?
 (A) 25600000 (B) 512000
 (C) 5120000 (D) 51200000
- (iii) What is the capacity of disk in bytes?
 (A) 512×10^4 (B) 512×10^5
 (C) 512×10^6 (D) 512×10^7
- (iv) How many cylinders does it have?
 (A) 512 (B) 1000
 (C) 2000 (D) 2048
- (v) Identify the In valid disk block size from below:
 (A) 2048 (B) 51200
 (C) 4098 (D) 4096
3. What is the order of internal node of B⁺ tree suppose that a child pointer takes 6 bytes, the search field value takes 14 bytes and the block size is 512 bytes?
 (A) 23 (B) 24
 (C) 25 (D) 26
4. The order of a leaf node in a B⁺ tree is the maximum number of (value, data, record pointer) pairs it can hold. Given that block size is 1 k bytes (1024 bytes), data record pointer is 7 bytes long, the value field is '9' bytes long and block pointer is 6 bytes.
 (A) 63 (B) 64
 (C) 65 (D) 66
5. The following key values are inserted into a B⁺ tree in which order of the internal nodes is 3, and that of the leaf nodes is 2, in the sequence given below. The order of internal nodes is the maximum number of tree pointers in each node, and the order of leaf nodes is the maximum number of data items that can be stored in it. The B⁺ tree is initially empty. 10, 3, 6, 8, 4, 2, 1. What is the maximum number of times leaf nodes would get split up as a result of these insertions?
 (A) 3 (B) 4
 (C) 5 (D) 6
6. For the same key values given in the above question, suppose the key values are inserted into a B⁻ tree in which order of the internal nodes is 3 and that of leaf nodes is 2. The order of internal nodes is the maximum number of tree pointers in each node and the order of leaf nodes is the maximum number of data items that can be stored in it. The B⁻ tree is initially empty. What is the maximum number of times leaf nodes would get split up as a result of these insertions?
 (A) 1 (B) 2
 (C) 3 (D) 4
7. Suppose that we have an ordered file with 45,000 records stored on a disk with block size 2048 bytes. File records are of fixed size and are unspanned with record length 120 bytes.
- (i) What is the blocking factor?
 (A) 16 (B) 17
 (C) 18 (D) 19
- (ii) What is the number of blocks needed for the file?
 (A) 2642 (B) 2644
 (C) 2646 (D) 2648
- (iii) How many block accesses are required to search for a particular data file using binary search?
 (A) 10 (B) 11
 (C) 12 (D) 13
8. Suppose that the ordering key field of the file is 12 bytes long, a block pointer is 8 bytes long, and we have constructed a primary index for the file. Consider the file specifications given in the above questions.
- (i) What is the size of each index entry?
 (A) 16 (B) 18
 (C) 20 (D) 22
- (ii) What is the blocking factor for the index?
 (A) 101 (B) 102
 (C) 103 (D) 104
- (iii) What is the total number of index entries?
 (A) 2642 (B) 2644
 (C) 2646 (D) 2648
- (iv) What is the number of index blocks?
 (A) 22 (B) 24
 (C) 26 (D) 28
- (v) How many block accesses are required, if binary search is used?
 (A) 3 (B) 4
 (C) 5 (D) 6
9. For the file specifications given in Q. No. 7, if we construct secondary index on a non-ordering key field of the file that is 12 bytes long, a block-pointer of size 8 bytes, each index entry is 20 bytes long and the blocking factor is 102 entries per block.
- (i) What is the total number of index blocks?
 (A) 422 (B) 424
 (C) 442 (D) 444
- (ii) How many block accesses are required to access the secondary index using binary search?
 (A) 6 (B) 7
 (C) 8 (D) 9
10. For the file specifications given in Q. No. 8, if we construct a multilevel index, number of 1st-level blocks are 442, blocking factor is 102, each index entry is 20 bytes long.
- (i) What is the number of 2nd-level blocks?
 (A) 4 (B) 5
 (C) 6 (D) 7
- (ii) What is the number of 3rd-level blocks?
 (A) 0 (B) 1
 (C) 2 (D) 3

11. Construct a B⁺ tree for (1,4,7,10,17,21,31) with $n = 4$, which nodes will appear two times in a final B⁺ tree?
 (A) 17,7,20 (B) 17,7,20,25
 (C) 17,20,25 (D) 7,17,25

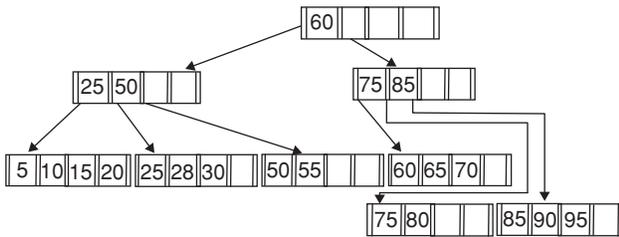
12. Suppose the hash function is $h(x) = x \bmod 8$ and each bucket can hold at most two records. The extendable hash structure after inserting 1, 4, 5, 7, 8, 2, 20, what is the local depth of '4'?
 (A) 0 (B) 1
 (C) 2 (D) 3

13. Consider the given B⁺ tree, insert 19 into the tree, what would be the new element in level 2?



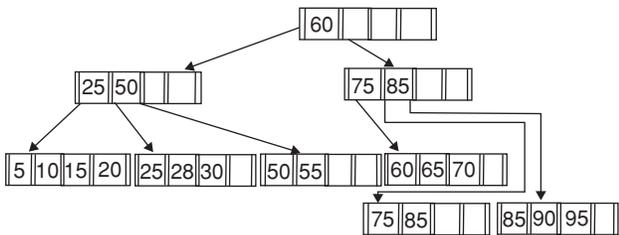
- (A) 13 (B) 18
 (C) 20 (D) 29

14. Consider the given B⁺ tree, delete 70 and 25 from the tree, what are the elements present in level 2? (\therefore root is at level 1)



- (A) 25, 50, 75 (B) 25, 50, 75, 85
 (C) 28, 50, 75, 85 (D) 28, 50, 65, 75

15. Delete 60 from the above given tree (Q. No. 14). After deletion, what is the total number of nodes present in the tree?



- (A) 5 (B) 6
 (C) 7 (D) 8

16. What will be the number of index records/block?
 (A) 68 (B) 65
 (C) 69 (D) None

17. What will be the number of index blocks?
 (A) 442 (B) 440
 (C) 400 (D) None

18. Consider the following:

Block size = 1025 bytes

Record length in data file = 100 bytes

Total number of records = 30000

Search key = 9 bytes

Pointer = 6 bytes

What is the number of index blocks?

- (A) 44 (B) 45
 (C) 46 (D) None

19. Which of the following is maximum search time t_{\max} in B⁻ trees?

$$(A) t_{\max} = a \log_2 \left(\frac{N}{2} \right) \left[\frac{a+d}{\log_2 m} + \frac{bm}{\log_2 m} + c \right]$$

$$(B) t_{\max} = a \log_2 \left(\frac{N}{2} \right) \left[\frac{a+d}{\log_2 m} + \frac{bm}{\log_2 m} + c \right]$$

$$(C) t_{\max} = a \log_2 N \left[\frac{a+d}{\log_2 m} + \frac{bm}{\log_2 m} + c \right]$$

$$(D) t_{\max} = a \log_2 (N) \left[\frac{a}{\log_2 m} + \frac{bm}{\log_2 m} + c \right]$$

20. Consider a B⁺ tree. A child pointer takes 3 bytes, the search field value takes 7 bytes, and the block size is 256 bytes. What is the order of the internal node?

- (A) 63 (B) 64
 (C) 65 (D) 66

Practice Problems 2

Directions for questions 1 to 20: Select the correct alternative from the given choices.

- Which of the following is true?
 - Every conflict serializable is view serializable
 - Every view serializable is conflict serializable
 - Both A and B
 - A schedule can be either only conflict serializable or only view-serializable.
- Which one is the 2-phase locking rule?
 - Two transactions cannot have conflicting locks
 - No unlock operation can precede a lock operation in the same transaction.
 - No data is/are affected until all locks are obtained and until the transaction is in its locked point.
 - All of the above
- If Transaction T_i has obtained an exclusive mode lock on item Q , then
 - T_i can read Q
 - T_i can write Q
 - T_i can read and write
 - Neither read nor write
- Phantom phenomenon is
 - A transaction retrieves a collection of objects but sees same result.
 - A transaction retrieves a collection of objects but sees different results.
 - Transaction T_1 waits for T_2 and T_2 waits for T_1
 - This problem arises when the transaction has not locked all the objects.
- We can avoid the starvation of transactions by granting locks by following manner:
When a transaction T_i requests a lock on a data item Q in a particular mode M , the concurrency control manager grants the lock provided that
 - There is no other transaction holding a lock on Q in a mode that conflicts with M .
 - There is no other transaction that is waiting for a lock on Q ,
 - (A) and (B)
 - None
- Which one is correct?
 - Upgrading can take place only in shrinking phase
 - Upgrading can take place only in growing phase.
 - Downgrading can take place only in growing phase
 - (A) and (C) both
- A simple but widely used scheme automatically generates the appropriate lock and unlock instructions for a transaction, on the basis of read and write requests from the transaction:
 - When a transaction T_i issues a read (Q) operation, the system issues a lock $s(Q)$ instruction followed by the read instruction.
 - When T_i issues a write Q operation, the system checks to see whether T_i already holds a shared lock on Q . If it does, then the system issues an upgrade Q instruction followed by the write Q instruction, otherwise the system issues a lock $-X(Q)$ instruction, followed by the write Q instruction.
 - All locks obtained by a transaction are unlocked after that transaction commits or aborts.
 - All of the above

8. Which one is correct?

- A lock manager can be implemented as a process that receives messages from transactions and sends messages in reply.
- It uses linked list of records.
- It uses hash table called lock table.
- All of the above

Common data questions 9 and 10: Transaction T_1 has 5 instructions. Transaction T_2 has 3 instructions.

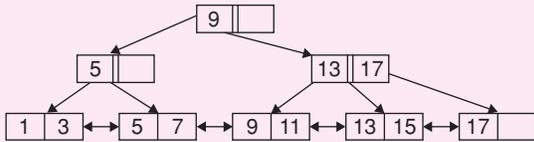
- The number of non-serial transactions will be
 - 15
 - 8
 - 2
 - 56
- The number of serial transaction schedules will be
 - 15
 - 8
 - 2
 - 56
- In a heap file system, which of the following function finds 'average number of blocks to be read'?
 - $\frac{i}{n} = \frac{1}{2}(1+n) = \frac{n}{2}$
 - $\sum_{i=1}^n \frac{i}{n} = \frac{1}{2}(1+n) = \frac{n}{2}$
 - $\sum_{i=0}^{n-1} \frac{i}{n} = \frac{1}{2}(1+n) = \frac{n}{2}$
 - All of the above
- What is the disadvantage in one directory per user?
 - Different applications can be divided into separate groups.
 - Different applications cannot be divided into separate groups
 - All files are in a single group
 - All of the above
- What are the possible violations if an application program uses isolation-level 'Read uncommitted'?
 - Dirty read problem
 - Non-repeatable read problem
 - Phantom phenomenon
 - All of the above

14. The two-phase locking protocol
 (A) ensures serializability
 (B) issues locks in two phases
 (C) unlocks in two phases
 (D) All of the above
15. The point in the schedule where the transaction has obtained its final lock (the end of its growing phase) is called the
 (A) block point (B) critical section
 (C) growing point (D) lock point
16. Which of the following is not a problem of file management system?
 (A) Data redundancy
 (B) Lack of data independence
 (C) Program dependence
 (D) All of the above
17. Which of the following is/are true about master list of an index file?
 (i) Is sorted in ascending order
 (ii) A number is assigned to each record.
 (A) Only (i) (B) Only (ii)
 (C) Both (i) and (ii) (D) None of the above
18. To have a file, holding a list is necessary to
 (i) Identify the records in the list
 (ii) Identify the name, and type of the fields of each record.
 (iii) Decide which fields will be used as sort of index keys.
 (A) Only (i) and (ii)
 (B) Only (i) and (iii)
 (C) Only (ii) and (iii)
 (D) All of the above
19. Two files may be joined into a third file, if the following is true:
 (A) if they have row in common
 (B) if they have a field in common
 (C) Both (A) and (B)
 (D) None
20. The minimum number of record movements required to merge four files w (with 10 records), x (with 20 records), y (with 15 records) and z (with 5 records) is:
 (A) 50 (B) 40
 (C) 30 (D) 35

PREVIOUS YEARS' QUESTIONS

1. A clustering index is defined on the fields which are of type **[2008]**
 (A) non-key and ordering
 (B) non-key and non-ordering
 (C) key and ordering
 (D) key and non-ordering
2. A B-tree of order 4 is built from scratch by 10 successive insertions. What is the maximum number of node splitting operations that may take place? **[2008]**
 (A) 3 (B) 4
 (C) 5 (D) 6
3. Consider a file of 16384 records. Each record is 32 bytes long and its key field is of size 6 bytes. The file is ordered on a non-key field, and the file organization is unspanned. The file is stored in a file system with block size 1024 bytes, and the size of a block pointer is 10 bytes. If the secondary index is built on the key field of the file, and a multilevel index scheme is used to store the secondary index, the number of first-level and second-level blocks in the multilevel index are respectively **[2008]**
 (A) 8 and 0 (B) 128 and 6
 (C) 256 and 4 (D) 512 and 5
4. The following key values are inserted into a B⁺ tree in which order of the internal node s is 3, and that of the leaf nodes is 2, in the sequence given below. The order of internal nodes is the maximum number of tree pointers in each node, and the order of leaf nodes is the maximum number of data items that can be stored in it. The B⁺ tree is initially empty.
 10, 3, 6, 8, 4, 2, 1
 The maximum number of times leaf nodes would get split up as a result of these insertions is **[2009]**
 (A) 2 (B) 3
 (C) 4 (D) 5
5. Consider a B⁺ tree in which the maximum number of keys in a node is 5. What is the minimum number of keys in any non-root node? **[2010]**
 (A) 1 (B) 2
 (C) 3 (D) 4
6. An index is clustered, if **[2013]**
 (A) it is on a set of fields that form a candidate key.
 (B) it is on a set of fields that include the primary key.
 (C) the data records of the file are organized in the same order as the data entries of the index.
 (D) the data records of the file are organized not in the same order as the data entries of the index.
7. A file is organized so that the ordering of data records is the same as or close to the ordering of data entries in some index. Then that index is called **[2015]**
 (A) Dense (B) Sparse
 (C) Clustered (D) Unclustered

8. With reference to the B+ tree index of order 1 shown below, the minimum number of nodes (including the Root node) that must be fetched in order to satisfy the following query: “Get all records with a search key greater than or equal to 7 and less than 15” is _____ [2015]



9. Consider a B+ tree in which the search key is 12 bytes long, block size is 1024 bytes, record pointer is 10 bytes long and block pointer is 8 bytes long. The maximum number of keys that can be accommodated in

each non-leaf node of the tree is _____. [2015]

10. B+ Trees are considered **BALANCED** because _____ [2016]
- (A) The lengths of the paths from the root to all leaf nodes are all equal.
 - (B) The lengths of the paths from the root to all leaf nodes differ from each other by at most 1.
 - (C) The number of children of any two non - leaf sibling nodes differ by at most 1.
 - (D) The number of records in any two leaf nodes differ by at most 1.
11. In a B+ tree, if the search-key value is 8 bytes long, the block size is 512 bytes and the block pointer size is 2 bytes, then the maximum order of the B+ tree is _____. [2017]

ANSWER KEYS

EXERCISES

Practice Problems 1

1. (i) C (ii) A 2. (i) B (ii) D (iii) C (iv) C (v) C 3. C 4. A 5. C 6. B
 7. (i) B (ii) D (iii) C 8. (i) C (ii) B (iii) D (iv) C (v) C 9. (i) C (ii) D
 10. (i) B (ii) B 11. B 12. D 13. B 14. C 15. B 16. A 17. A 18. B
 19. A 20. C

Practice Problems 2

1. A 2. D 3. C 4. B 5. C 6. B 7. D 8. D 9. D 10. C
 11. B 12. B 13. D 14. D 15. D 16. D 17. B 18. D 19. B 20. B

Previous Years' Questions

1. A 2. C 3. C 4. C 5. B 6. C 7. C 8. 5 9. 50 10. A
 10. 52

TEST

DATABASES

Time: 60 min.

- Which of the following are used in DBMS files?
 (i) Data dictionary (ii) DML
 (iii) Query language (iv) Transaction log
 (A) (i) and (ii) (B) (ii) and (iii)
 (C) (iii) and (iv) (D) (i) and (iv)
- Which among the following is not a problem of file management system?
 (A) Data redundancy
 (B) Lack of data independence
 (C) Program independence
 (D) None of these
- A transparent DBMS
 (A) cannot hide sensitive information from users
 (B) keeps its logical structure hidden from users
 (C) keeps its physical structure hidden from users
 (D) All of the above
- If the field size is too small, for the longest piece of data to be entered,
 (A) database program will be frozen
 (B) field will automatically expand
 (C) part of the data will be cut off
 (D) All of the above
- Which of the following functional dependencies are satisfied by the instance from the below relation?

A	B	C
1	7	3
1	9	5
1	11	5
5	3	3

- $AB \rightarrow C$ and $C \rightarrow B$
 - $BC \rightarrow A$ and $B \rightarrow C$
 - $BC \rightarrow A$ and $A \rightarrow C$
 - $AC \rightarrow B$ and $B \rightarrow A$
- Let E_1 and E_2 be two entities in an E/R diagram with single-valued attributes, R_1 and R_2 are two relationships between E_1 and E_2 , R_1 is one to many R_2 is many-to-one. R_1 and R_2 do not have any attributes of their own. What is the minimum number of tables required to represent this situation in the relation model?
 (A) 2 (B) 3
 (C) 4 (D) 5
 - Which of the following is true about DBMS?
 (i) Low-level DMLs are record-at-a-time
 (ii) High-level DMLs are set oriented or set-at-a-time
 (iii) Query in high-level DML specify which data to retrieve rather than how.
 (iv) When used as standalone, DML is called 'host language'

- (i) only (B) (i) and (iii)
 (C) (i), (ii) and (iii) (D) (iii) and (iv)
- In which of the following, the structure of data files is stored?
 (A) Metadata (B) Database catalog
 (C) Database schema (D) Data model
 - A schedule is a collection of
 (A) Data models (B) Transactions
 (C) Schemas (D) Tables
 - Select from the following which matches the term 'Impedance mismatch problem':
 (A) In compatibility of storage and data structure
 (B) Mismatch in user authentication
 (C) File structure mismatching
 (D) None of these
 - Which of the following is not a/an integrity constraint?
 (A) Entity integrity
 (B) Candidate key constraint
 (C) Business rules
 (D) None of the above
 - Select from the following which is concerned with 'Query Optimizer':
 (A) Extracts DML commands from an application program in a high-level language
 (B) Parsing and analyzing interactive query
 (C) Rearrangement and reordering of operations and elimination of redundancies
 (D) Performance monitoring
 - Which of the following does not belong to database model?
 (A) Relational Model (B) Distributed Model
 (C) Hierarchical Model (D) Network Model
 - What is the correct sequence of database design process?
 (i) Create conceptual schema
 (ii) Data model mapping
 (iii) Requirement collection and analysis
 (iv) Physical design
 (A) iii \rightarrow i \rightarrow ii \rightarrow iv
 (B) iii \rightarrow ii \rightarrow i \rightarrow iv
 (C) i \rightarrow ii \rightarrow iii \rightarrow iv
 (D) i \rightarrow iii \rightarrow ii \rightarrow iv
 - Consider the following schema definitions
 Employee {Name, SSN, Address, DNo}
 Department {DName, DNumber, Manager, SSN}
 Which among the following expressions represent the query $\Pi_{\text{name, address}}(\sigma_{\text{Dname} = \text{'Res'} \wedge \text{DNumber} = \text{DNo}}(\text{Department} \bowtie \text{Employee}))$?

- (A) Retrieve the name and address of employees who work for the project no 'Dno'
- (B) Retrieve the name and address of all employees who control the 'Res' department.
- (C) Retrieve the name and address of all employees who work for the 'Res' department.
- (D) None of these

16. Select from the following which closely resembles the concept 'Degree of a relationship':

- (A) Number of entities participating in a relation
- (B) Number of entity types participating in a relation
- (C) Number of strong entity types in a relation
- (D) Number of weak entity types in a relation

17. Consider the following statements in a database:

- (i) No primary key value can be NULL
- (ii) A tuple in one relation which refers to another relation must refer to an existing tuple in that relation
- (iii) The value of x determines the value of y in all states of a relation, where x and y are two attributes of the relation Which of the following combinations matches the given statements in order?

- (A) Referential integrity, functional dependency, entity integrity.
- (B) Functional dependency, entity integrity, referential integrity
- (C) Entity integrity, functional dependency, referential integrity.
- (D) Entity integrity, referential integrity, functional dependency

18. Consider the following relation schemas:

Works (emp_name, comp_name, salary)

Livesin (emp_name, street, city)

Location (comp_name, city)

Manager (manager_name)

What is returned by the following relational algebra expression

$$\pi_{emp_name}(\sigma_{comp_name=Time \wedge Works.emp_name=livesin.emp_name})$$

(Works \bowtie Livesin)

- (A) Names of all employees who work for TIME
- (B) Names of all employees of TIME who lives in the same city
- (C) Names of people who live in the same city
- (D) None of these

19. Consider the following SQL query:

Select distinct a_1, a_2, \dots, a_n from r_1, r_2, \dots, r_m where P

This query is equivalent to one of the following relational algebra expression:

(A) $\pi_{a_1, a_2, \dots, a_n} \sigma_P(r_1 \times r_2 \times \dots \times r_m)$

(B) $\pi_{a_1, a_2, \dots, a_n} \sigma_P(r_1 r_2 \times \dots \times r_m)$

(C) $\pi_{a_1, a_2, \dots, a_n} \sigma_P(r_1 \cup r_2 \cup \dots \cup r_m)$

(D) $\pi_{a_1, a_2, \dots, a_n} \sigma_P(r_1 r_2 \times \dots \times r_m)$

20. Let $R_1(A, B, C)$ and $R_2(D, E)$ be two relation schemas with primary keys A and D and C be a foreign key in R_1 referring to R_2 . Suppose there is no violation of the above referential integrity constraint in the instances r_1 and r_2 , which of the following relational algebra expression would necessarily produce an empty relation?

(A) $\pi_D(r_2) - \pi_C(r_1)$

(B) $\pi_C(r_1) - \pi_E(r_2)$

(C) $\pi_D(r_1 \bowtie_{C=D} r_2) - \pi_B(r_1)$

(D) $\pi_C(r_1 \bowtie_{C=E} r_2)$

21. Let r be an instance for the schema $R = (A, B, C, D)$. Let $r_1 = \pi_{A, B, C}(r)$ $r_2 = \pi_{A, D}(r)$ and $S = r_1 \bowtie r_2$. Also given that the decomposition of r into r_1 and r_2 is lossy, which of the following is true?

(A) $S \subset r$

(B) $r \cup S = r$

(C) $r \subset S$

(D) $r \bowtie S = S$

22. Which of the following is/are logical database structures?

(A) Network

(B) Tree

(C) Chain

(D) All of the above

23. A relational database management system manages data in more than one file at a time by using which of the following combinations?

(A) Tables and tuples

(B) Relations and tuples

(C) Tables and Relations

(D) Attributes and tuples

24. Let Emp = (Name, ID, ADDRESS, PHONE, SPOUSE, LIVINGAT) be a relation scheme with following FDs, which one of the following is a key

ADDRESS \rightarrow Phone

SPOUSE \rightarrow NAME

SPOUSE, ADDRESS \rightarrow PHONE

NAME \rightarrow ID

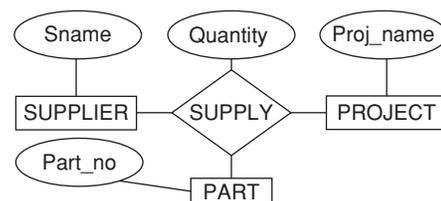
(A) ADDRESS, PHONE

(B) SPOUSE, ADDRESS

(C) NAME, SPOUSE

(D) NAME, ADDRESS

25. Consider the following E-R diagram



Select the most appropriate statement from the following for the above ER diagram:

- (A) Represents a ternary relationship
 - (B) Represents a binary relationship
 - (C) Represents a ternary relationship with instances of the form (s, j, p)
 - (D) Represents 1 – to – many relationships
26. If two relations R_1 and R_2 are such that they are of the same degree and domain of the corresponding fields are also the same, then which one of the following is true about R_1 and R_2 ?
- (A) $R_1 \subset R_2$
 - (B) $R_1 \cup R_2 = R_2 \cup R_1$
 - (C) R_1 and R_2 are union compatible
 - (D) None of these

Common data questions for 27 and 28: Let Employee and Guests be two relations with attributes (id, mobil_no, name, address) and (id, mob-no, comps_working, shifts) Relations respectively {id, mob_no} is the key for both.

27. Which of the following queries are equivalent?
- (i) $\pi_{id}(\text{Employee} \bowtie \text{Guests})$
 - (ii) $\pi_{id}(\text{Employee}) \bowtie \pi_{id}(\text{Guests})$
 - (iii) $\pi_{id}\{(\text{Employee-Guest}) \cap \text{Guest-Employee}\}$
 - (iv) $\pi_{id}\{\pi_{id, mob}(\text{Employee}) \cap \pi_{id, mob}(\text{Guest})\}$
- (A) (ii) and (iii)
 - (B) (ii), (iii) and (iv)
 - (C) (i), (ii) and (iv)
 - (D) (ii) and (iv) only

28. What does the following relational algebra expression represent?

$\pi_{id}(\pi_{id, mobil_no}(\text{Employee-Guests}))$

- (A) Id of all employees working with the company
- (B) Id of all permanent employees
- (C) Id of part time employees
- (D) None of these

Common data for questions 29 and 30:

29. Let R_1 and R_2 be two relations with attributes a_1 and a_2 . P_1 and P_2 be two predicates.

Select the expression from the following which is wrong:

- (A) $\sigma_{P_1}(\sigma_{P_1}(R_1)) \rightarrow \sigma_{P_2}(\sigma_{P_2}(R_1))$
- (B) $\sigma_{P_1}(\pi_{a_1}(R_1)) \rightarrow \pi_{a_1}(\sigma_{P_1}(R_1))$
- (C) $\sigma_{P_1}(R_1 \cup R_2) \rightarrow \sigma_{P_1}(R_1) \cup \sigma_{P_2}(R_2)$
- (D) $\pi_{a_2}(\sigma_{a_1}(R_1)) \rightarrow \sigma_{P_1}(\pi_{a_2}(R_1))$

30. Select from the following corresponding TRC for the wrong expression in the above question:

- (A) $\{t/\exists u, R_1(t[P_1]) = R_2(u[P_1])\}$
- (B) $\{t/\forall u, R_1(t[P_1]) = R_1(u[P_1])\}$
- (C) $\{t/\exists u, R_1(t[P_1]) \neq R_2(u[P_1])\}$
- (D) $\{t/\forall t \in R_1\}$

ANSWER KEYS

- | | | | | | | | | | |
|-------|-------|-------|-------|-------|-------|-------|-------|-------|-------|
| 1. D | 2. D | 3. C | 4. C | 5. B | 6. B | 7. C | 8. B | 9. B | 10. A |
| 11. B | 12. C | 13. B | 14. A | 15. C | 16. A | 17. D | 18. C | 19. A | 20. A |
| 21. C | 22. D | 23. C | 24. B | 25. C | 26. C | 27. C | 28. B | 29. A | 30. B |

Theory of Computation

Chapter 1: Finite Automata and Regular Languages	5.3
Chapter 2: Context Free Languages and Push Down Automata	5.24
Chapter 3: Recursively Enumerable Sets and Turing Machines, Decidability	5.37

U

N

I

T

5

This page is intentionally left blank

Chapter 1

Finite Automata and Regular Languages

LEARNING OBJECTIVES

- ☞ Fundamentals
- ☞ Languages
- ☞ Operations
- ☞ Finite state machine
- ☞ NFA with ϵ -moves
- ☞ Conversion of NFA to DFA
- ☞ Minimization of DFA
- ☞ Equivalence between NFA and DFA
- ☞ Mealy and Moore machines
- ☞ Equivalence of Moore and Mealy machine
- ☞ Regular languages
- ☞ Constructing FA for given RE
- ☞ Pumping lemma for regular sets
- ☞ Closure properties of regular sets
- ☞ Regular grammar

FUNDAMENTALS

Alphabet: An alphabet is a finite non-empty set of symbols.

Example: Portion of a calculator: $\{0, 1, 2, 3 \dots 9, \div, =, -, +, \times, (,)\}$

Note: 1. At least one symbol is necessary.

2. ' Σ ' denote Alphabet.

String: A string over an alphabet ' A ' is a finite ordered sequence of symbols from ' A '. The length of string is number of symbols in string, with repetitions counted.

Example: If $\Sigma = \{0 - 9, \div, =, -, +, \times, (,)\}$ then Strings valid: $12 + 34, 90 \times 10, (1 + 2) \times (1 \div 3)$

Strings Invalid: $\sin(45), \log(10)$ etc. These strings are not valid because $\sin(), \log()$ are not defined over the alphabet set.

Note: Repetitions are allowed.

Length of $|12 + 34| = 5(1, 2, +, 3, 4)$

• The Empty string denoted by ' ϵ ', is the (unique) string of length zero.

Note: Empty string, $\epsilon \neq$ empty set, \emptyset .

• If S and T are sets of strings, then $ST = \{xy | x \in S \text{ and } y \in T\}$

Given an alphabet A ,

$$A^0 = \{\epsilon\}$$

$$A^{n+1} = A.A^n$$

...

$$A^* = \bigcup_{n=0}^{\infty} A^n$$

Languages

- A language ' L ' over Σ is any finite or infinite set of strings over Σ .
- The elements in L are strings – finite sequences of symbols.
- A language which does not contain any elements is called 'empty language'.

Note: Empty language, $\{\} \neq \{\epsilon\}$, empty string because $\{\} = \emptyset \neq \epsilon$ i.e., Empty language resembles empty set i.e., \emptyset .

• A language L over an alphabet A is subset of A^* i.e., $L \subset A^*$.

Example 1: Language (L) for strings that consists of only 0's or only 1's and have an odd length over alphabet $\{0, 1\}$ is

(A) $\{0, 1, 00, 11, 000, 111 \dots\}$

(B) $\{00, 11, 01, 10 \dots\}$

(C) $\{000, 101, 110, 111 \dots\}$

(D) $\{0, 1, 000, 111, 11111, 00000 \dots\}$

Solution: (D)

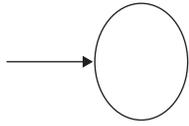
Only 0's \rightarrow should have only 0's. It should not be combination of 0's and 1's.

Only 1's \rightarrow should have only 1's. It should not be combination of 0's and 1's.

Odd length \rightarrow only odd number of 0's or odd number of 1's i.e., length of string should be odd.

5.4 | Unit 5 • Theory of Computation

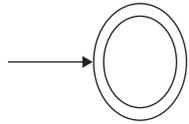
An Empty Languages An empty language is a language which does not accept any strings including. The Finite automata for empty language can be represented as



(i.e., One state, non-accepting and no transitions).

A language which only accepts (ϵ)

E: The language which only accepts ' ϵ ' can be represented as



This machine accepts E – only.

Σ^* : The set of all strings over an alphabet Σ will be denoted by Σ^* .

Σ^+ : This will denote the set $\Sigma^+ = \{\epsilon\}$.

Ex: If $\Sigma = \{0, 1\}$ then

$\Sigma^* = \{\epsilon, 0, 1, 00, 01, 10, 11, 000, 001, \dots\}$

$\Sigma^+ = \{0, 1, 00, 01, 10, 11, 000, 001, \dots\}$

Operations

Operations on strings

1. Concatenation: Combines two strings by putting one after other.

Example 2: Two strings are defined as $x = \text{java}$, $y = \text{script}$. The concatenation ($x.y$) of two strings results in _____.

- (A) scriptjava (B) javascript
(C) jascriptva (D) scrijavapt

Solution: (B)

$x.y = \text{java.script} = \text{javascript}$

Note: Concatenation of empty string with any other string gives string itself.

i.e., $x.\epsilon = \epsilon.x = x$

2. Substring: If ' w ' is a string, then ' v ' is a substring of ' w ' if there exists string x and y such that $w = xvy$.

' x ' is called 'prefix' and y is called suffix of w .

Example 3: String, $w = \text{'gymnastics'}$ is defined with prefix, $x = \text{'gym'}$ and suffix, $y = \text{'cs'}$. The substring of the given string is _____

- (A) nasti (B) mnas
(C) gymnastics (D) ics

Solution: (A)

Because, $w = xvy$

$\Rightarrow \text{gymnastics} = \text{gymvcs}$

$\therefore v = \text{nasti}$

3. Kleen star operation: Let ' w ' be a string, w^* is set of strings obtained by applying any number of concatenations of w with itself, including empty string.

Example: $a^* = \{\epsilon, a, aa, aaa, \dots\}$

4. Reversal: If ' w ' is a string, then w^R is reversal of string spelled backwards.

Rules:

- $x = (x^R)^R$
- $(xz)^R = z^R \cdot x^R$

Example 4: A string, x is defined as, $x = \text{butter}$. Then $(x^R)^R$ is _____

- (A) butter (B) rettub
(C) butret (D) retbut

Solution: (A)

$x \rightarrow \text{butter}$

$x^R \rightarrow \text{rettub}$

$(x^R)^R \rightarrow \text{butter}$.

Operations on languages

1. Union: Given some alphabet Σ , for any two languages, L_1, L_2 over Σ , the union $L_1 \cup L_2$ of L_1 and L_2 is the language, $L_1 \cup L_2 = \{w \in \Sigma^* | w \in L_1 \text{ or } w \in L_2\}$

2. Intersection: Given some alphabet Σ , for any two languages L_1, L_2 over Σ , the intersection $L_1 \cap L_2$ of L_1 and L_2 is language, $L_1 \cap L_2 = \{w \in \Sigma^* | w \in L_1 \text{ and } w \in L_2\}$

3. Difference: Given some alphabet Σ , for any two languages L_1, L_2 over Σ , the difference $L_1 - L_2$ of L_1 and L_2 is language, $L_1 - L_2 = \{w \in \Sigma^* | w \in L_1 \text{ and } w \notin L_2\}$

Note: Difference is also called 'Relative Complement.'

A special case of difference is obtained when $L_1 = \Sigma^*$, in which case. Complement \bar{L} of language, L is defined as, $\bar{L} = \{w \in \Sigma^* | w \notin L\}$

4. Concatenation: Given an alphabet Σ , for any two languages L_1, L_2 over Σ , the concatenation $L_1 L_2$ of L_1 and L_2 is language

$L_1 L_2 = \{w \in \Sigma^* | \exists u \in L_1, \exists v \in L_2, w = uv\}$

Properties:

$L\emptyset = \emptyset = \emptyset L$

$L\{\epsilon\} = L = \{\epsilon\}L$

$(L_1 \cup \{\epsilon\})L_2 = L_1 L_2 \cup L_2$

$L_1(L_2 \cup \{\epsilon\}) = L_1 L_2 \cup L_1$

$L^n L = LL^n = L^{n+1}$

Note: $L_1 L_2 \neq L_2 L_1$

Example 5: Let $L_1 = \{00, 11\}$, $L_2 = \{01, 10\}$. Then $L_1 o L_2 =$ _____

- (A) $\{00, 11, 01, 10\}$
(B) $\{0001, 0010, 1101, 1110\}$
(C) $\{0001, 0010, 11, 01, 10\}$
(D) $\{00, 1101, 1110, 11, 10\}$

Solution: (B)

$L_1 o L_2 = \{00, 11\} o \{01, 10\} = \{00.01, 00.10, 11.01, 11.10\} = \{0001, 0010, 1101, 1110\}$

5. Kleen * closure (L^*): Given an alphabet Σ , for any language L over Σ , the * closure L^* of L is language,

$L^* = \bigcup_{n \geq 0} L^n$

6. Kleen + closure (L^+): The kleen +closure, L^+ of L is the language, $L^+ = U_{n \geq 1} L^n$

$$L^* = L^0 \cup L^1 \cup L^2 \cup \dots \cup L^n \cup \dots$$

$$L^+ = L^1 \cup L^2 \cup L^3 \dots \cup L^n \cup \dots$$

Properties:

$$\emptyset^* = \{\epsilon\}$$

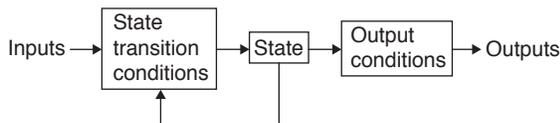
$$L^+ = L^*L$$

$$(L^*)^* = L^*$$

$$L^*L^* = L^*$$

Finite State Machine (FSM)

- FSM is simplest computational model of limited memory computers.
- FSM is designed to solve decision problems i.e., to decide whether given input satisfies certain conditions.
- The next state and output of a FSM is a function of input and of current state.

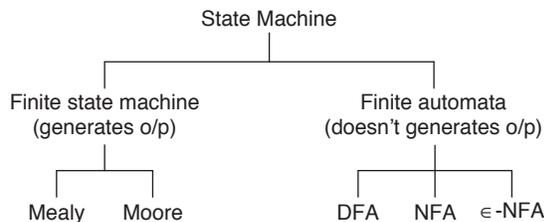


Types of FSM:

1. Melay machine.
2. Moore machine

Finite Automata (FA):

- FA is a state machine that comprehensively captures all possible states and transitions that a machine can take while responding to a stream (sequence) of input symbols.
- FA is recognizer of 'regular languages'.



Types of FA

1. Deterministic Finite Automata (DFA):

- DFA machine can exist in only one state at any given time.
- DFA is defined by 5-tuple: $\{Q, \Sigma, q_0, F, \delta\}$, where

- $Q \rightarrow$ Finite number of states (elements)
- $\Sigma \rightarrow$ Finite set of symbols (alphabets)
- $q_0 \rightarrow$ Start/Initial state
- $F \rightarrow$ Set of final states.
- $\delta \rightarrow$ Transition function, which is a mapping between $\delta: Q \times \Sigma \rightarrow Q$.

How to use DFA:

Input: A word w in Σ^*

Question: Is w acceptable by DFA?

Steps:

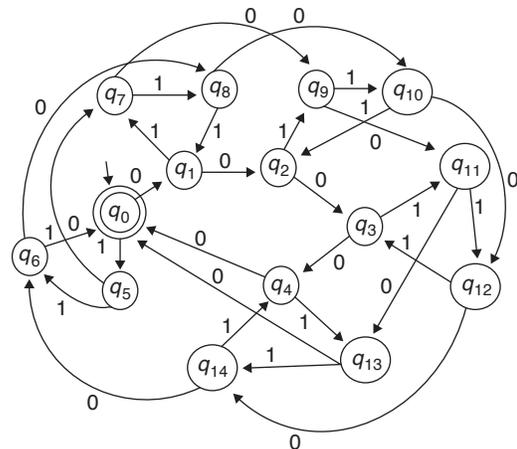
- Start at 'initial state', q_0 .
- For every input symbol in sequence w , do.
- Compute the next state from current state, given the current input symbol in w and transition function.
- If after all symbols in 'w' are consumed, the current state is one of the final states (f) then accept 'w';
- Otherwise, reject w .

Transition diagram: State machines are represented by directed graphs called transition (state) diagrams.

- The vertices denoted by single circle represent the state and arcs labeled with input symbol correspond to transition.
- The final states are represented with double circles.

Transition Table: Transition function can be represented by tables.

Example 6: The following finite state machine accepts all those binary strings in which the numbers of 0's and 1's are respectively.



- (A) Divisible by 3 and 2
- (B) Odd and even
- (C) Divisible by 5 and 3
- (D) Divisible by 2 and 3

Solution: (C)

Number of 0's is divisible by 5.

Number of 1's is divisible by 3.

Table Transition Table

Current State	0	1
$\rightarrow q_0$	q_1	q_5
q_1	q_2	q_7
q_2	q_3	q_9
q_3	q_4	q_{11}
q_4	q_0	q_{13}
q_5	q_7	q_6
q_6	q_8	q_0
q_7	q_9	q_8
q_8	q_{10}	q_1

q_9	q_{11}	q_{10}
q_{10}	q_{12}	q_2
q_{11}	q_{13}	q_{12}
q_{12}	q_{14}	q_3
q_{13}	q_0	q_{14}
q_{14}	q_6	q_4

Note: Minimum number of states for k-divisibility is k-states.

In above example, $q_0 - q_{14} \rightarrow 15$ - states.

$\therefore 5 \times 3 = 15$

The given binary strings have number of 0's divisible by 5 and number of 1's divisible by 3.

2. Non-deterministic finite Automata (NFA):

- The machine can exist in multiple states at the same time.
- Each transition function maps to a set of states.
- NFA is defined by 5-tuple: $\{Q, \Sigma, q_0, F, \delta\}$, where

$Q \rightarrow$ Finite number of states (elements)

$\Sigma \rightarrow$ Finite set of symbols. (Alphabets)

$q_0 \rightarrow$ Start/Initial state

$F \rightarrow$ Set of final states.

$\delta \rightarrow$ Transition function which is a mapping between

$\delta = Q \times \Sigma \rightarrow 2^Q$

How to use NFA:

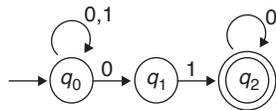
Input: a word w in Σ^*

Question: Is w accepted by NFA?

Steps:

- Start at 'start state' q_0 .
- For every input symbol in the sequence, w does.
- Determine all possible next states from current state, given the current input symbol in w and transition function.
- If after all symbols in w are consumed, at least one of the current states is a final state then accept w .
- Otherwise, reject w .

Example 7: What is the language, L generated by the below NFA, given strings defined over alphabet, $\Sigma = \{0, 1\}$.



- (A) Strings that end with '0'
- (B) Strings that start with '0' and end with '0'
- (C) Strings that contain '01' as substring
- (D) Strings that contain '01' as substring and end with '0'

Solution: (D)

State	0	1
$\rightarrow q_0$	$\{q_0, q_1\}$	$\{q_0\}$
q_1	\emptyset	q_2
q_2	$\{q_2\}$	\emptyset

String: 0100100

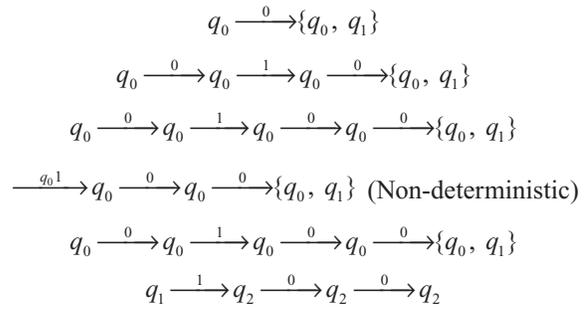


Table 2 Difference between NFA and DFA

DFA	NFA
1. All transitions are deterministic i.e., each transition leads to exactly one state.	1. Transitions could be non-deterministic i.e., a transition could lead to a subset of states.
2. For each state, the transition on all possible symbols should be defined.	2. For each state, not all symbols necessarily have to be defined.
3. Accepts input if last state is in 'F'.	3. Accepts input if one of last states is in 'F'.
4. Practical implementation is feasible.	4. Practical implementation has to be deterministic (so needs conversion to DFA).

Relation between DFA and NFA

- A language 'L' is accepted by a DFA if and only if it is accepted by a NFA.
- Every DFA is special case of a NFA.

Example 8: Let N_f and D_f denote the classes of languages accepted by non-deterministic finite automata and deterministic finite automata respectively. Which one of following is true?

- (A) $D_f \subset N_f$
- (B) $D_f \supset N_f$
- (C) $D_f = N_f$
- (D) $D_f \in N_f$

Solution: (C)

According to 'subset construction', every language accepted by NFA is also accepted by some DFA.

$\therefore D_f = N_f$

NFA WITH ϵ -MOVES

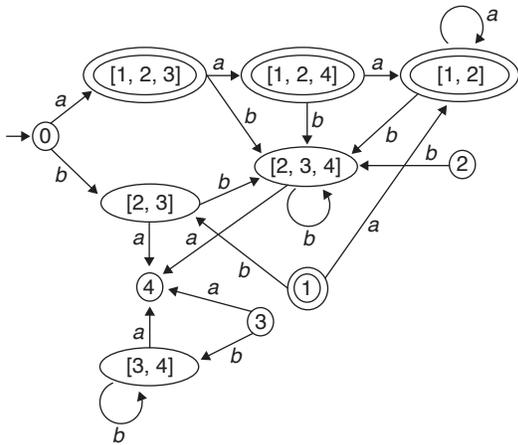
- ϵ -transitions in finite automata allows a state to jump to another state without consuming any input symbol.

Conversion and Equivalence:

ϵ -NFA \rightarrow NFA \rightarrow DFA

NFA without ϵ -moves:

- Two FA, N_ϵ and N are said to be equivalent, if $L(N_\epsilon) = L(N)$ i.e., any language described by some N_ϵ , there is an N that accepts the same language.
- For $N_\epsilon = (Q, Z, \delta, q_0, F)$ and $N = (Q, \Sigma', \delta', q_0, F')$, Find
- $\delta'(q, a) = \epsilon$ -closure ($\delta(\epsilon$ -closure($q, a)$))

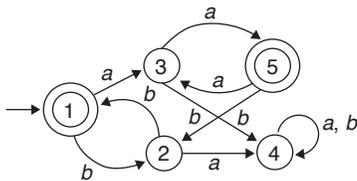


Minimization of DFA

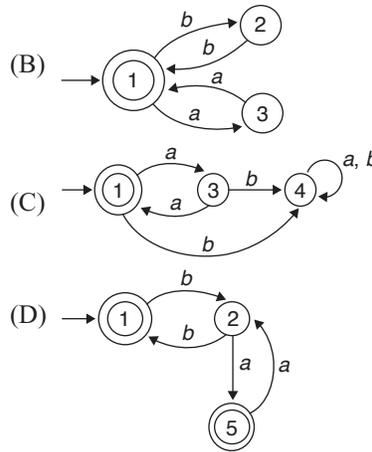
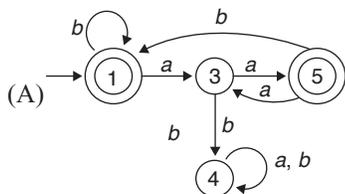
Given a DFA, $M = (Q, \Sigma, \delta, q_0, F)$, we construct a reduced DFA, $M' = (Q', \Sigma', \delta', q'_0, F')$ as follows

1. Remove all inaccessible states. All states that are unreachable from the initial state are removed.
2. Consider all pairs of states (p, q) , If $p \in F$ and $q \notin F$ or vice versa mark the pair (p, q) as distinguishable.
3. Repeat until no previously unmarked pairs are marked. For all pairs (p, q) and all $a \in \Sigma$, compute $\delta(p, a) = p_a$ and $\delta(q, a) = q_a$. If the pair (p_a, q_a) is marked as distinguishable mark (p, q) as distinguishable.
4. Find the sets of all indistinguishable states, say $\{q_i, q_j, \dots, q_k\}$, $\{q_\ell, q_m, \dots, q_n\}$, etc. For each set $\{q_i, q_j, \dots, q_k\}$ of such indistinguishable states, create a state labelled $ij \dots k$ for M .
5. For each transition rule of M of the form $\delta(q_r, a) = q_p$, find the sets to which q_r and q_p belong. If $q_r \in \{q_i, q_j, \dots, q_k\}$ and $q_p \in \{q_\ell, q_m, \dots, q_n\}$, add a rule to δ' : $\delta'(ij \dots k, a) = \ell m \dots n$.

Example 11: A DFA with alphabet $\Sigma = \{a, b\}$ is given below:



Which of the following is valid minimal DFA which accepts same language as given DFA?



Solution: (B)

Initially, $\{1, 5\}$, $\{2, 3, 4\}$

Depending on next states and inputs, the partitions of states can be as: $\{\{1, 5\}, \{2\}, \{3\}, \text{ and } \{4\}\}$

Since, 1 to 5 have same transition, unite $\{1, 5\}$

State 4 is dead state \rightarrow It has transition only to itself. Since, $\{2\}, \{3\}$ are singletons, they exist.

\therefore States in minimized DFA are $\{1, 2, \text{ and } 3\}$

$\{1\} \rightarrow \{1, 5\}$

For transitions, since $1 \xrightarrow{a} 3$, $1 \xrightarrow{b} 2$ in given DFA, in minimized DFA, transitions are added from $1 \xrightarrow{a} 3$, $1 \xrightarrow{b} 2$. Also, since $2 \xrightarrow{b} 1$, $3 \xrightarrow{a} 1$ in given DFA, the minimized DFA, transitions are added from $2 \xrightarrow{b} 1$, $3 \xrightarrow{a} 1$.

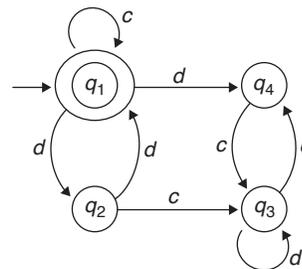
Equivalence Between NFA and DFA

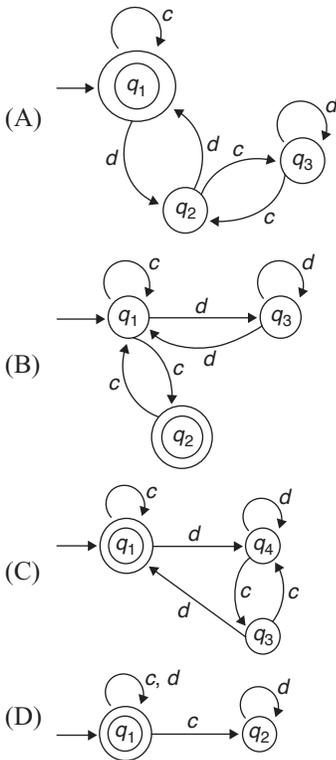
There is a DFA_D for any NFA_N i.e., $L(D) = L(N)$.

Construction:

- In DFA or NFA, whenever an arrow is followed, there is a set of possible states. This set of states is a subset of Q .
- Track the information about subsets of states that can be reached from initial state after following arrows.
- Consider each subset of states of NFA as a state of DFA and every subset of states containing a final state as a final state of DFA.

Example 12: Which of following is equivalent DFA for the NFA given below:





Solution: (A)

Table 6 Transition Table of NFA

δ	c	d
$\rightarrow q_1$	q_1	$\{q_2, q_4\}$
q_2	q_3	q_1
q_3	q_4	q_3
q_4	q_3	\emptyset

Table 7 Transition Table of DFA

δ	c	d
$\rightarrow q_1$	q_1	q_2
q_2	q_3	q_1
q_3	q_2	q_1

Table 8 Common Table

δ	c	d
(q_1, q_1)	(q_1, q_1)	(q_2, q_4, q_2)
(q_2, q_2)	(q_3, q_3)	(q_1, q_1)
(q_3, q_3)	(q_4, q_2)	(q_3, q_3)
(q_4)	q_3	\emptyset

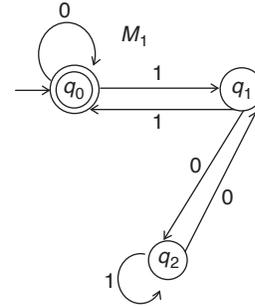
Equivalence of Finite Automatas:

- Two automatas A and B are said to be equivalent if both accept exactly the same set of input strings.
- If two automatas M_1 and M_2 are equivalent then
 - If there is a path from the start state of M_1 to a final state of M_1 labeled $a_1 a_2 \dots a_k$ then there is a path from

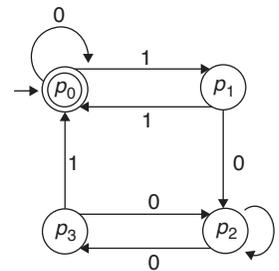
the start state of M_2 to the final state of M_2 labeled $a_1 a_2 \dots a_k$.

- If there is a path from the start state of M_2 to a final state M_2 labeled $b_1 b_2 \dots b_i$ then there is a path from the start state of M_1 to the final state of M_1 labeled $b_1 b_2 \dots b_i$.

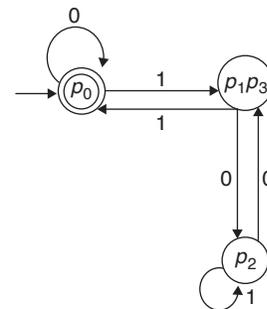
Example:



M_2 :



In M_2 , states p_1 and p_3 are equivalent (as both are reaching either final or non-final states with same input). After minimizing M_2 , we will get



$\therefore M_1$ and M_2 are equivalent.

Union: The union of two languages L and M is the set of strings that are in both L and M .

Ex: $L = \{0, 1\}$, $M = \{111\}$
 $L \cup M = \{0, 1, 111\}$.

Concatenation: The concatenation of Languages L and M is the set of strings that can be formed by taking any string in L and concatenating it with any string in M .

Example: $L = \{0, 1\}$, $M = \{\epsilon, 010\}$
 $LM = \{0, 1, 0010, 1010\}$.

Closure, Star or Kleen star of a language L:

Kleen star is denoted as L^* . It represents the set of strings that can be formed by taking any number of strings from L with repetition and concatenating them. It is a Unary operator. L^0 is the set; we can make selecting zero strings from L .

$L^0 = \{\epsilon\}$

L^1 is the language consisting of selecting one string from L .

L^2 is the language consisting of concatenations selecting two strings from L .

...

L^* is the union of $L^0, L^1, \dots, L^\infty$.

Ex: $L = \{0,10\}$

$L^* = \{0,00,000,10,010, \dots\}$

Intersection:

Let two DFAs M_1 and M_2 accept the languages L_1 and L_2 .

$M_1 = (Q_1, \Sigma, \delta_1, q_1^1, F_1)$

$M_2 = (Q_2, \Sigma, \delta_2, q_2^1, F_2)$

The intersection of M_1 and M_2 can be given as

$M = (Q, \Sigma, \delta, q_0, F)$

Q = Pairs of states, one from M_1 and one from M_2 i.e.,

$Q = \{(q_1, q_2) \mid q_1 \in Q_1 \text{ and } q_2 \in Q_2\}$

$Q = Q_1 \times Q_2$.

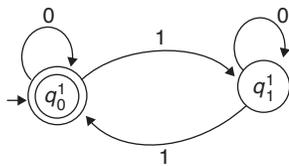
$q_0 = (q_0^1, q_0^2)$

$\delta((q_1^1, q_2^2), x) = (\delta_1(q_1^1, x), \delta_2(q_2^2, x))$

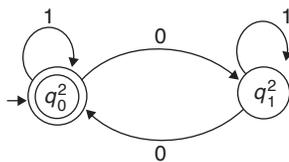
$F = \{(q_1, q_2) \mid q_1 \in F_1 \text{ and } q_2 \in F_2\}$

Example:

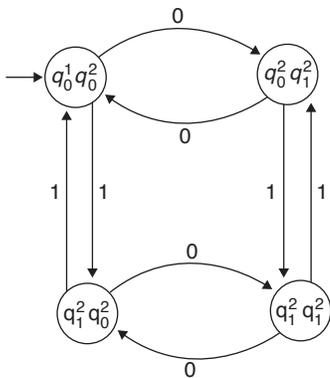
M_1 : Strings with even number of 1's.



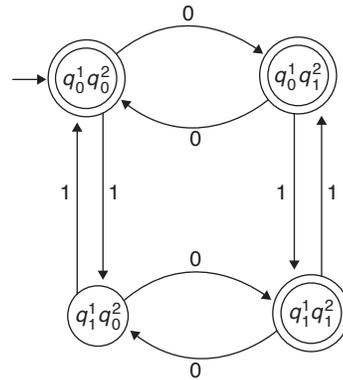
M_2 : Strings with odd number of 0's.



$M_1 \cap M_2$: Strings with even number of 1's and odd number of 0's.



Union of M_1 and M_2 :



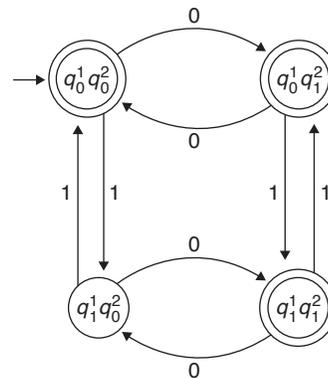
Difference: The difference of L_1 and L_2 can be given as $L_1 - L_2$ with $M = (Q, \Sigma, \delta, q_0, F)$.

$Q = Q_1 \times Q_2$

$q_0 = (q_0^1, q_0^2)$

$\delta((q_i^1, q_j^2), x) = (\delta_1(q_i^1, x), \delta_2(q_j^2, x))$

$F = \{(q_1, q_2) \mid q_1 \in F_1 \text{ and } q_2 \notin F_2\}$



Reversing a DFA:

- M is a DFA which recognizes the language L .
- M^R will accept the language L^R .

To construct M^R :

- Reverse all transitions
- Turn the start state to final state
- Turn the final states to start state.
- Merge states and modify the FA, such that the resultant contain a single start state.

MEALY AND MOORE MACHINES

Moore Machine

A moore machine is a finite state machine, where outputs are determined by current state alone.

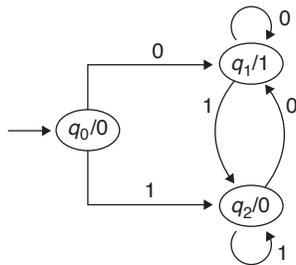
A Moore machine associates an output symbol with each state and each time a state is entered, an output is obtained simultaneously. So, first output always occurs as soon as machine starts.

Moore machine is defined by 6-tuples:

- $(Q, \Sigma, \delta, q_0, \Delta, \lambda)$, where
- $Q \rightarrow$ Finite set of states
- $\Sigma \rightarrow$ Finite set of input symbols
- $\Delta \rightarrow$ It is an output alphabet
- $\delta \rightarrow$ Transition function, $Q \times \Sigma \rightarrow Q$ (state function)
- $\lambda \rightarrow$ Output function, $Q \rightarrow \Delta$ (machine function)
- $q_0 \rightarrow$ Initial state of machine

Note: The output symbol at a given time depends only on present state of moore machine.

Example 13: The language generated by the following moore machine is:



- (A) 2's complement of binary number.
- (B) 1's complement of binary number.
- (C) Has a substring 101.
- (D) Has a substring 110.

Solution: (B)
 Binary number: 1011
 1's complement: 0100

$$q_0 \xrightarrow{1/0} q_2, \quad q_2 \xrightarrow{0/1} q_1, \quad q_1 \xrightarrow{1/0} q_2, \quad q_2 \xrightarrow{1/0} q_2$$

$$1 \rightarrow 0, 0 \rightarrow 1, 1 \rightarrow 0, 1 \rightarrow 0$$

Mealy Machine

- A mealy machine is a FSM, where outputs are determined by current state and input.
- It associates an output symbol with each transition and the output depends on current input.
- Mealy machine is defined on 6-tuples: $(Q, \Sigma, \delta, q_0, \Delta, \lambda)$, where

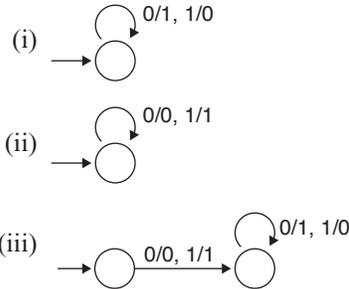
- Q – Finite set of states.
- Σ – Finite set of input symbols.
- $\delta - (Q \times \Sigma \rightarrow Q)$ is transition function.
- $q_0 \rightarrow q_0 \in Q$ is initial state.
- $\Delta \rightarrow$ Finite set of output symbols.
- $\lambda \rightarrow$ Output function, $\lambda(Q \rightarrow \Delta)$

Note: In Moore machine, for input string of length n , the output sequence consists of $(n + 1)$ symbols.

In Mealy machine, for input string of length n , the output sequence also consists of ' n ' symbols.

Example 14: Let $(Me)^2$ mean that given a Mealy machine, an input string is processed and then output string is immediately fed into the machine (as input) and reprocessed.

Only this second resultant output is considered as the final output of $(Me)^2$. If final output string is same as original input string then $(Me)^2$ has an identity property. Consider following machines.

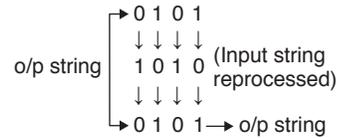


Which of above machines have identity property?

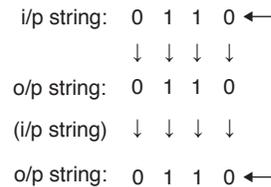
- (A) (i) only
- (B) (i) and (ii) but not (iii)
- (C) (i) and (iii) but not (ii)
- (D) All have identity property

Solution: (D)

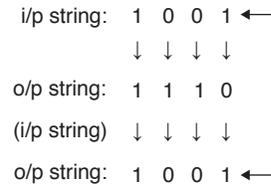
(i) Consider i/p string



(ii)



(iii)



Equivalence of Moore and Mealy machine

(a) Mealy machine equivalent to Moore machine:

If $M_1 = (Q, \Sigma, \Delta, \delta, \lambda, q_0)$ is a Moore machine, then there is a Mealy machine M_2 equivalent to M_1 .

Proof: Let $M_2 = (Q, \Sigma, \Delta, \delta, \lambda^1, q_0)$ and define $\lambda^1(q, a)$ to be $\lambda(\delta(q, a))$ for all states q and input symbol ' a '.

Then M_1 and M_2 enter the same sequence of states on the same input, and with each transition M_2 emits the o/p that M_1 associates with the state entered.

Let us consider Mealy Machine

Present State	Next State			
	Input State	a = 0 Output	Input State	a = 1 Output
→ q ₁	q ₃	0	q ₂	0
q ₂	q ₁	1	q ₄	0
q ₃	q ₂	1	q ₁	1
q ₄	q ₄	1	q ₃	0

To convert the Mealy machine to Moore machine,

- We look into the next state column for any state, say q_i and determine the number of different outputs associated with q_i in next column.
- Split q_i into several different states, the number of such states being equal to the number of different outputs associated with q_i.

Present State	Next State			
	Input State	a = 0 Output	Input State	a = 1 Output
→ q ₁	q ₃	0	q ₂₀	0
q ₂₀	q ₁	1	q ₄₀	0
q ₂₁	q ₁	1	q ₄₀	0
q ₃	q ₂₁	1	q ₁	1
q ₄₀	q ₄₁	1	q ₃	0
q ₄₁	q ₄₁	1	q ₃	0

- The pair of states and outputs in the next state column can be rearranged as:

Present state	Next State			output
	a = 0	a = 1	output	
→ q ₁	q ₃	q ₂₀	1	
q ₂₀	q ₁	q ₄₀	0	
q ₂₁	q ₁	q ₄₀	1	
q ₃	q ₂₁	q ₁	0	
q ₄₀	q ₄₁	q ₃	0	
q ₄₁	q ₄₁	q ₃	1	

Moore machine equivalent to Mealy machine

Let M₁ = (Q, Σ, Δ, δ, λ, q₀) be a Mealy machine. Then there is a machine M₂ equivalent to M₁

Proof: Let M₂ = (QXΔ, Σ, Δ, δ', λ', [q₀, b₀]), where b₀ is an arbitrary selected member of Δ.

That is, the states of M₂ are pairs [q, b] consisting of a state of M₁ and output symbol, Define δ' ([q, b], a) = [δ(q, a), λ(q, a)] and λ' ([q, b]) = b.

The second component of a state [q, b] of M₂ is the output made by M₁ on some transition into state q.

Only the first components of M₂'s states determine the moves made by M₂.

Every induction on 'n' shows that if M₁ enters states q₀, q₁ ... q_n on inputs a₁, a₂ ... a_n and emits output b₁, b₂, ... b_n then M₂ enters states [q₀, b₀], [q₁, b₁] ... [q_n, b_n] and emits outputs b₀, b₁ ... b_n.

Let us consider the Moore machine

Present State	Next State		
	a = 0	a = 1	Output
→ q ₀	q ₃	q ₁	0
q ₁	q ₁	q ₂	1
q ₂	q ₂	q ₃	0
q ₃	q ₃	q ₀	0

- To convert Moore into Mealy machine, we must follow the reverse procedure of converting Mealy machine into Moore machine.
- For every input symbol we form, the pair consisting of the next state and the corresponding output and reconstruct the table for Mealy machine.
- For example, the state q₃ and q₁ in the next state column should be associated with outputs 0 and 1, respectively.

The Transition table for Mealy machine is:

Present state	Next State			
	a = 0 state output	a = 1 state output	a = 1 state output	a = 1 state output
→ q ₀	q ₃ 0	q ₁ 1		
q ₁	q ₁ 1	q ₂ 0		
q ₂	q ₂ 0	q ₃ 0		
q ₃	q ₃ 0	q ₀ 0		

REGULAR LANGUAGES

The set of regular languages over an alphabet Σ is defined recursively as below. Any language belonging to this set is a regular language over Σ.

Definition of set of regular languages

- Basis clause: ∅, {ε}, {a} for any symbol a ∈ Σ, are regular languages.
- Inductive clause: If L_r and L_s are regular languages, then L_r ∪ L_s, L_r · L_s, L_r^{*} are regular languages.
- External clause: Nothing is a regular language, unless it is obtained from above two clauses.

Regular language: Any language represented by regular expression(s) is called a regular language.

Ex: The regular expression a* denotes a language which has {ε, a, aa, aaa, ...}

Regular expression

- Regular expressions are used to denote regular languages.
- The set of regular expressions over an alphabet Σ is defined recursively as below. Any element of that set is a regular expression.

- Basis clause: \emptyset, ϵ, a are regular expression corresponding to languages $\emptyset, \{\epsilon\}, \{a\}$ respectively where a is an element of Σ .
- Inductive clause: If r and s are regular expression corresponding to languages L_r and L_s then $(r + s), (rs)$ and (r^*) are regular expressions corresponding to the languages $L_r \cup L_s, L_r \cdot L_s$ and L_r^* respectively.
- External clause: Nothing is a regular expression, unless it is obtained from above two clauses.

Closure property of regular expressions The iteration or closure of a regular expression R , written as R^* is also a regular expression.

Ex: $\Sigma = \{a\}$ then a^* denotes the closure of Σ .

$$a^* = \{\epsilon, a, aa, aaa, \dots\}$$

Conventions on regular expressions

1. The operation ** has highest precedence over concatenation, which has precedence over union (+).
i.e., $RE(a + (b(c^*))) = a + bc^*$
2. The concatenation of K r's, where r is a regular expression is written as r^k . The language corresponding to r^k is L_r^k . Where L_r is language corresponding to regular expression r i.e., $rr = r^2$
3. r^+ is a regular expression to represent L_r^+

Note: A regular expression is not unique for a language i.e., regular language corresponds to more than one regular expression.

Example 15: Give regular expression for set of strings which either have 'a' followed by some b's or all b's also containing 'ε'.

- (A) $b^* + ab^*$ (B) $a^* + ba^*$
 (C) $(\epsilon) + (\epsilon + a) b^+$ (D) $b^* + ab^* + \epsilon$

Solution: (C)

The regular expression is, $r = ab^+ + b^+ + \epsilon = b^+(a + \epsilon) + \epsilon$.

Identity rules for regular expressions:

1. $\emptyset + R = R$
2. $\emptyset \cdot R = R\emptyset = \emptyset$
3. $\epsilon R = R\epsilon = R$
4. $\emptyset^* = \epsilon$ and $\epsilon^* = \epsilon$
5. $R + R = R$
6. $RR^* = R^*R = R^+$
7. $\epsilon + RR^* = R^*$ and $\epsilon + R^*R = R^*$
8. $(R^*)^* = R^*$
9. $R^*R^* = R^*$
10. $\epsilon + R^* = R^*$
11. $(R + \epsilon)^* = R^*$
12. $R^*(\epsilon + R)^* = (\epsilon + R)^*R^* = R^*$
13. $R^*R + R = R^*R$
14. $(P + Q)R = PQ + QR$ and $R(P + Q) = RP + RQ$
15. $(P + Q)^* = (P^*Q^*)^* = (P^* + Q^*)^*$
16. $(PQ)^*P = P(QP)^*$
17. R is given as, $R = Q + RP$ has unique solution, $R = QP^*$. This is Arden's theorem.
18. $(P + Q)^* = (P^* + Q)^* = (P + Q^*)^*$

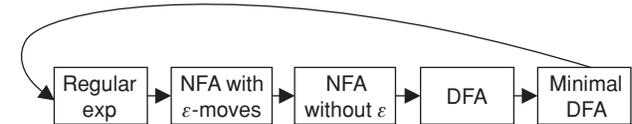
Example 16: If r_1 and r_2 are regular expressions denoting languages L_1 and L_2 respectively then which of following is false?

- (A) $(r_1)|(r_2)$ is regular expression denoting $L_1 \cup L_2$.
 (B) $(r_1)(r_2)$ is regular expression denoting $L_1 \cdot L_2$.
 (C) \emptyset is not a regular expression.
 (D) $\{r_1\}^*$ is regular expression denoting L_1^* .

Solution: (C)

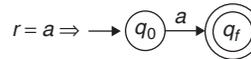
CONSTRUCTING FA FOR GIVEN RE

- Relationship between FA and RE.



Identities:

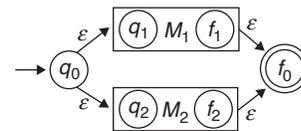
Basis:



Induction:

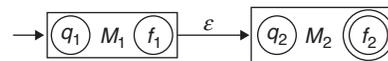
- **Union:** $L(r) = L(r_1) + L(r_2)$ i.e., $L(M) = L(M_1) \cup L(M_2)$

Let $M_1 = (Q_1, \Sigma_1, \delta_1, q_1, \{f_1\})$, $M_2 = (Q_2, \Sigma_2, \delta_2, q_2, \{f_2\})$ with $L(M_1) = L(r_1)$ and $L(M_2) = L(r_2)$, then $M = (Q_1 \cup Q_2 \cup \{q_0, f_0\}, \Sigma_1 \cup \Sigma_2, \delta, q_0, \{f_0\})$



- **Concatenation:**

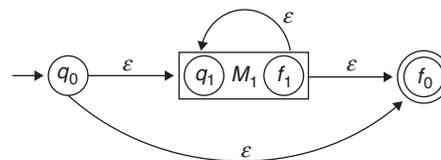
$L(r) = L(r_1) \cdot L(r_2)$ i.e., $L(M) = L(M_1) \cdot L(M_2)$



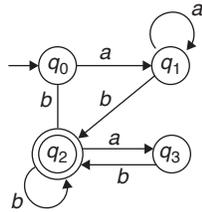
- **Closure:**

$L(r) = L(r)^*$ i.e., $L(M) = L(M_1)^*$

Let $M_1 = (Q_1, \Sigma_1, \delta_1, q_1, \{f_1\})$ then $L(M) = (Q_1 \cup \{q_0, f_0\}, \Sigma_1, \delta, q_0, \{f_0\})$



Example 17: The regular expression generated by the given FA.



- (A) $(a + ba^*) b^*$
- (B) $(aa^*b + bb^*) b^*$
- (C) $(b + ab^*) a^*$
- (D) $(ab + ba^*) a^*$

Solution: (B)

q_2 is final state which is obtained with input symbol only 'b'. So, (C) or (D) is not true.

In (A) $\rightarrow ba^*$ is not defined in given FA. Instead bb^* is defined.

Pumping Lemma for Regular Sets

Theorem Let 'L' be an arbitrary regular language. Then there exists a positive integer, P with following property:

Given an arbitrary member, w of L having length at least P (i.e., $|w| \geq P$), w can be divided into 3-parts, $w = xyz \exists$

- $|y| \geq 1$ (the middle part is non-empty)
- $|xy| \leq P$ (the first two parts have length at most P)
- For each, $i \geq 0$, $xy^iz \in L$ (removing or repeating the middle part produces member of L)

Proof Let L be an arbitrary regular language. Then there is a FA, say M that decides L.

Let P be the number of states of M.

Let w be an arbitrary member of L, having length 'n' with $n \geq P$.

Let q_0, q_1, \dots, q_n be states that M on input w. That is, for each i, after reading the first i symbols of w, M is at q_i .

q_0 is initial state of M. Also, since $w \in L$, q_n is a final state of M.

Let $x = w_1 \dots w_c, y = w_{c+1} \dots w_d, z = w_{d+1} \dots w_n$. Then:

- $|y| \geq 1$
- $|xy| \leq P$
- M transitions from q_0 to q_c on x.
- M transitions from q_c to q_d on y.
- M transitions from q_d to q_n on z.

Thus, for every $i \geq 0$, M transitions from q_0 to q_n on xy^iz and so, xy^iz is a member of L.

Note:

- Pumping lemma is used to verify that given language is not regular.
- Pumping lemma follows pigeon hole principle.

Example 18: The language, L is defined as:

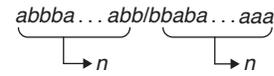
$L = \{w_1w_2 : w_1, w_2 \in \{a, b\}^*, |w_1| = |w_2|\}$. Is the language regular?

- (A) Regular
- (B) Not regular
- (C) Cannot be determined
- (D) None of these

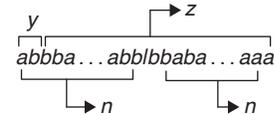
Solution: (A)

Fix pumping length, $K = 2$

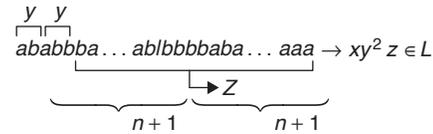
For every proper strings in L, ($2n \geq 2$)



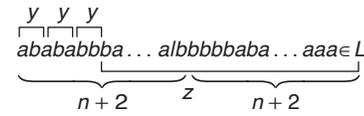
- Split in x, y, z with desired properties.



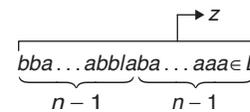
- Let $x = \epsilon, y =$ first two symbols, $z =$ rest.



- xy^3z :



- $xy^0z \rightarrow$



\therefore For every $i \geq 0, xy^iz \in L$. Hence given language is regular.

CLOSURE PROPERTIES OF REGULAR SETS

1. **Union:** If L and M are regular languages, LUM is regular language closed under union.
2. **Concatenation and Kleen closure:** If L and M are regular languages, L.M is regular language and L^* is also regular.
3. **Intersection:** $L \cap M$ is regular, if L and M are regular languages.
4. **Difference:** $L - M$ contains strings in 'L' but not M, where L and M are regular languages.
5. **Complementation:** The complement of language L is $\Sigma^* - L$.

Note: Since Σ^* is surely regular, the complement of a regular language is always regular. Where Σ^* is a universal language.

6. **Homomorphism:** If L is a regular language, h is homomorphism on its alphabet then $h(L) = \{h(w) | w \text{ is in } L\}$ is also a regular language.

Regular grammar

- **Grammar:** Generative description of a language.
- **Automaton:** Analytical description.
- A grammar is a 4-tuple, $G = (V, \Sigma, R, S)$ where V : alphabet (variable) (non-terminals)

$\Sigma \subseteq V$ is set of terminal symbols.

$R \subseteq (V^+ \times V^*)$ is a finite set of production rules.

$S \in V - \Sigma$ is start symbol.

Notation

- Elements of $V - \Sigma$: A, B, \dots
- Elements of Σ : $a, b \dots$
- Rules $(\alpha, \beta) \in R$: $\alpha \rightarrow \beta$ or $\alpha \xrightarrow{G} \beta$
- Start symbol is written as S .
- Empty word: ϵ

Example 19: The regular expression that describe the language generated by grammar, $G = (\{S, A, B\}, \{a, b\}, S, \{S \rightarrow Aab, A \rightarrow Aab|B, B \rightarrow a\})$

- (A) $(ab)^* a$
- (B) $aab(ab)^*$
- (C) $ab^* aa$
- (D) $(a + ba)^*$

Solution: (B)

$S \rightarrow Aab \rightarrow Aab ab \rightarrow A ab abab \rightarrow Bababab \rightarrow aababab \rightarrow aab(ab)^*$

Union of two Regular languages:

If L_1 and L_2 are two languages then

$$L_1 \cup L_2 = \{w/w \in L_1 \text{ or } w \in L_2\}$$

The union of two regular languages is also a regular language.

Let $M_1 = (Q_1, \Sigma, \delta_1, q_1, f_1)$

$M_2 = (Q_2, \Sigma, \delta_2, q_2, f_2)$

$M = M_1 \cup M_2$ can be given as

$M = (Q, \Sigma, \delta, q_0, f)$

Where $Q = \{(r_1, r_2) \mid r_1 \in Q_1 \text{ and } r_2 \in Q_2\}$

i.e., Q is the Cartesian product of sets Q_1 and Q_2 .

Σ is the alphabet, is the same in M_1 and M_2 .

$\Sigma = \Sigma_1 \cup \Sigma_2$.

δ is the transition function given as:

$$\delta(r_1, r_2), a = (\delta_1(r_1, a) \delta_2(r_2, a))$$

q_0 is the pair (q_1, q_2) .

F is the set of pairs in which either member is an accept state of M_1 or M_2 .

$$F = \{(r_1, r_2) \mid r_1 \in F_1 \text{ or } r_2 \in F_2\}$$

TYPES OF GRAMMARS

- Type 0: Unrestricted, recursively enumerable languages.
- Type 1: Context-sensitive grammar.
- Type 2: Context free grammar.
- Type 3: Regular grammar.

Type 0: Recursively enumerable grammar: (Turing Machine) (TM):

Every production rule is of form: $\alpha \rightarrow \beta$, where α and β are in $(V \cup T)^*$, i.e., there can be any strings of terminals and non-terminals (no-restriction).

Type 1: Context-sensitive Grammar: (Linear bounded automaton) (LBA):

Every production rule is of form, $\alpha \rightarrow \beta$ are in $(V \cup T)^*$ and $\alpha \neq \epsilon$ and $|\beta| \geq |\alpha|$ i.e., any strings of terminals and non-terminals and length of string that can appear on RHS of production must be greater than or equal to length of string that can appear on LHS of production.

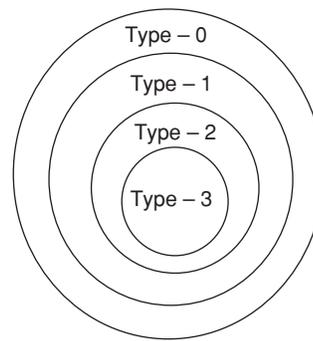
Type 2: Context-free grammar: (Push down automaton) (PDA):

Every production rule is of form, $A \rightarrow \alpha$ where α is in $(V \cup T)^*$ i.e., LHS of rule is single non-terminal and RHS can be any string of terminals and non-terminals.

Type 3: Regular grammar: (Finite automaton) (FA):

Every production is of form, $A \rightarrow aB$ or $A \rightarrow a$ where A and $B \in V$ and $a \in T$. That is, LHS of rule is non-terminal and RHS can be terminal (or) terminal followed by non-terminal.

Relationship between types of grammar:



- Regular sets are properly contained in CFL (Context Free Languages).
- The CFLs not containing empty string ϵ , are properly contained in CSL. (Context sensitive language).
- The CSLs are properly contained in Recursively enumerable languages.
- $RG \subset CFG \subset CSL \subset REG$

Left-linear Grammar:

All productions have form: $A \rightarrow Bx$ or $A \rightarrow x$

Right-linear Grammar:

All productions have the form: $A \rightarrow xB$ or $A \rightarrow x$.

Note:

- The regular grammars characterize the regular sets i.e., a language is regular if and only if it has a left-linear grammar or if and only if it has a right-linear grammar.
- If L has a regular grammar, then L is a regular set.
- If L is a regular set, then L is generated by some left-linear grammar and by some right-linear grammar.

Arden's theorem: Let P and Q be two regular expressions over Σ . If P does not contain ' ϵ ' then the following equation in R , namely $R = Q + RP$ has a unique solution given by $R = QP^*$.

Arden's Theorem to obtain regular expression from given transition diagram: The following steps are used to find the RE recognized by transition system.

The following assumptions are made regarding the transition system.

- (i) The transition graph does not have ϵ -moves
- (ii) It has only one initial state, q_0 .
- (iii) The states in the transition diagram are $q_0, q_1, q_2, \dots, q_n$.
- (iv) Q_i , the regular expression represents the set of strings accepted by a system even though q_i is the final state.
- (v) ${}^a ij$ denotes the regular expression representing the set of labels of edges from q_i to q_j . When there is no such edge ${}^a ij = \phi$.

We will get the following set of equations.

$$Q_1 = Q_1 \alpha_{11} + Q_2 \alpha_{12} + \dots + Q_n \alpha_{n1} + \epsilon$$

$$Q_2 = Q_1 \alpha_{12} + Q_2 \alpha_{22} + \dots + Q_n \alpha_{n2}$$

⋮
⋮
⋮

$$Q_n = Q_1 \alpha_{1n} + Q_2 \alpha_{2n} + \dots + Q_n \alpha_{nn}$$

By Repeatedly applying substitutions and Arden's theorem, we can express Q_i in terms of α_{ij} 's.

For getting the set of strings recognized by the transition system, we have to take the union of all Q_i 's corresponding to final states.

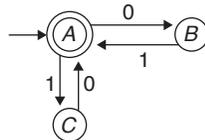
Construction of Regular Grammar from FA

Step I: Associate suitable variables like $A, B, C \dots$ with states of automata.

Step II: Obtain the productions of the grammar as:
If $\delta(A, a) = B$ then add production $A \rightarrow aB$ to list of productions of grammar, if B is a final state, then add either $A \rightarrow a$ or $B \rightarrow \epsilon$, to list of productions of grammar.

Step III: The variable associated with initial state of automata is start symbol of grammar.

Example 20: Regular grammar generating language accepted by below automata is



- (A) $A \rightarrow 0B|1C|\epsilon$
 $B \rightarrow 1A$
 $C \rightarrow 0A$
- (B) $A \rightarrow 1B|0C|\epsilon$
 $B \rightarrow 1A$
 $C \rightarrow 0A$
- (C) $A \rightarrow B|C|\epsilon$
 $B \rightarrow 1$
 $C \rightarrow 0$

- (D) $A \rightarrow 0A|1B|\epsilon$
 $B \rightarrow 1C$
 $C \rightarrow 0A$

Solution: (A)
 $A \rightarrow 0B, A \rightarrow 1C, B \rightarrow 1A, C \rightarrow 0A$

$\therefore A$ is final state, $A \rightarrow \epsilon$

$$\begin{aligned} \therefore A &\rightarrow 0B|1C|\epsilon & A &\rightarrow 0B|1C \\ B &\rightarrow 1A & \text{(or)} & B \rightarrow 1A|1 \\ C &\rightarrow 0A & & C \rightarrow 0A|0 \end{aligned}$$

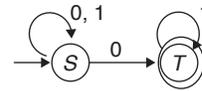
Construction of FA from given regular grammar

Given a regular grammar, G ; a regular expression specifying $L(G)$ can be obtained directly as follows:

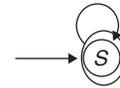
- Replace the ' \rightarrow ' symbol in productions of grammar by '=' symbol, to get set of equations.
- Solve the set of equations obtained above to get the value of variable, S , where S is start symbol of grammar, result is regular expression specifying $L(G)$.

Example 21: The Regular grammar and FA for given regular expression $\phi^* 1^* U (0\phi)^*$ is ____

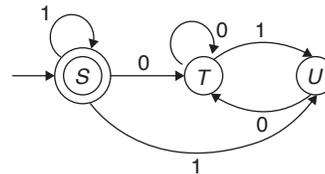
- (A) $S \rightarrow 0S|1S|0$
 $T \rightarrow 1T|\epsilon$



- (B) $S \rightarrow 1S|\epsilon$

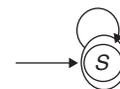


- (C) $S \rightarrow 0T|1S|\epsilon$
 $T \rightarrow 0T|1U|\epsilon$
 $U \rightarrow 0T|1S$



- (D) Cannot be determined

Solution: (B)
 $\emptyset^* 1^* \cup (0\emptyset)^* = \emptyset^* \cdot 1^* \cup \emptyset^* = \epsilon \cdot 1^* \cup \epsilon = 1^*$.



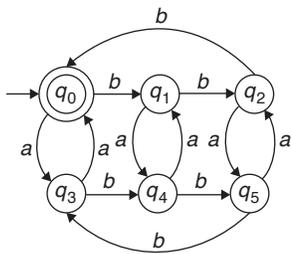
15. What will be number of final states obtained in DFA for language $L = \{w/w \text{ contains at least two } 0\text{'s and at most one } 1\}$ over $\Sigma = \{0, 1\}$.

- (A) 2 (B) 1
(C) 3 (D) 4

Practice Problems 2

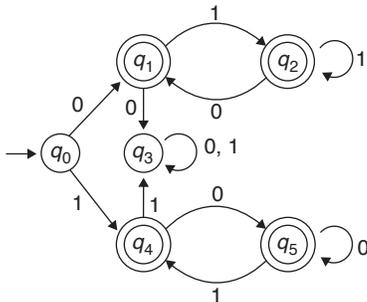
Directions for questions 1 to 15: Select the correct alternative from the given choices.

- Which of the following Regular expression is equal to given regular expression: $(b + aa^*b) + (b + aa^*b)(a + ba^*b)^*(a + ba^*b)$
 (A) $Ab(b + baa^*)$ (B) $a^*b(a + ba^*b)$
 (C) $a^*b(a + ba^*b)^*$ (D) $ab(b + aa^*b)^*$
- The following DFA accepts set of all strings over $\{a, b\}$ that contain

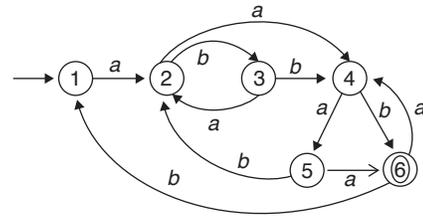


- (A) Number of a's even and number of b's odd.
 (B) Consecutive a's and b's
 (C) Contain *bbb* as substring
 (D) Number of a's even and number of b's divisible by three.

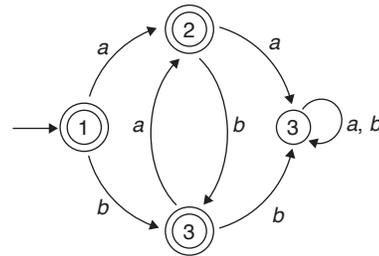
3. The regular language $L(r)$ for the given FSM is:



- (A) It can start with zero followed by any number of 1's but no two consecutive 0's.
 (B) It can start with 1, followed by any number of 0's but no two consecutive 1's.
 (C) It is a combination of 0's and 1's but no two consecutive 0's or 1's.
 (D) Both (A) and (B).
4. The language, L is defined as a set of non-palindromes over $\{a, b\}$. Is L regular?
 (A) Yes (B) No
 (C) Cannot be determined (D) None of above
5. The DFA, for language, L over $\Sigma = \{a, b\}$ is given below. What will be number of states in minimized DFA.

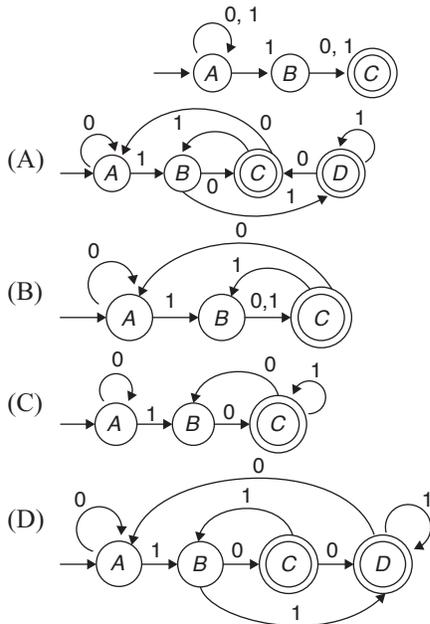


- (A) 4 (B) 6
(C) 2 (D) 3
6. The minimal DFA given below is defined for language, $L = \{w \in \{a, b\}^*\}$ over $\Sigma = \{a, b\}$. The ' L ' is:



- (A) Strings that contain equal number of a's and b's that have adjacent characters same.
 (B) Contains adjacent characters same
 (C) No two adjacent characters are same
 (D) Starts and ends with same character that have adjacent character same.
7. The regular grammar $L(G)$ contains productions, P for language, $L = \{w \in \{a, b\}^* / \text{there is at least one } a\}$ are:
 (A) $S \rightarrow aS | bS | aT$
 $T \rightarrow aT | bT | a | b$
 (B) $S \rightarrow aS | bS | \epsilon$
 (C) $S \rightarrow aBb | bB$
 $B \rightarrow a | b$
 (D) $S \rightarrow bB$
 $B \rightarrow b | \epsilon$
8. The regular expression for a language is defined as $((a^*b)^*(bc^*)^*)$. The total number of final states obtained in both NFA and DFA are respectively:
 (A) 4, 2 (B) 1, 3
 (C) 1, 5 (D) 2, 3
9. The language, L is defined as $\{w/w \text{ has } n \text{ occurrences of } 0\text{'s where } n \text{ mod } 5 \text{ is } 3\}$ over $\Sigma = \{0, 1\}$. The number of final states obtained in the DFA for L is:
 (A) 4 (B) 5
 (C) 1 (D) 2

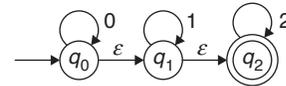
10. Which of the following is an equivalent DFA for the following NFA?



11. A regular grammar over alphabet $\Sigma = \{a, b, c, d\}$ whose language, is set of strings that contain exactly two b's is:

- (A) $S \rightarrow aS|bS|cS|dA$
 $A \rightarrow aA|bA|cA|dA|\epsilon$
- (B) $S \rightarrow aS|cS|dS|bB$
 $B \rightarrow aB|cB|dB|bC$,
 $C \rightarrow aC|cC|dC|\epsilon$
- (C) $S \rightarrow aS|bS|cS|dA$
 $A \rightarrow aA|bB|cC$
 $B \rightarrow b$
 $C \rightarrow c$
- (D) None of above

12. The following NFA contains ϵ -moves with 5, transitions. If this NFA with ϵ -moves is converted to NFA without ϵ -moves, what will be total number of transitions in obtained NFA?



- (A) 5
- (B) 4
- (C) 6
- (D) 3

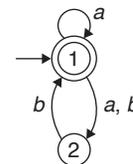
13. The regular expression, $r = (a + b)^*$. One more regular expression which represents same regular expression 'r' is:

- (A) $a^* + b^*$
- (B) $a^* \cdot b^*$
- (C) $a^*(ba^*)^*$
- (D) $(a + b)^*(a + b)$

14. The Regular grammar, $L(G)$ is defined for L with productions as $S \rightarrow Aab$, $A \rightarrow Aab|aB$, $B \rightarrow a$. What is Language generated by $L(G)$?

- (A) Containing alternative a's and b's
- (B) Containing alternative a's and b's, begins with an 'a' and ends with a 'b'.
- (C) 'aa' followed by at least one set of alternating ab's.
- (D) Consecutive aa's followed by 'b'.

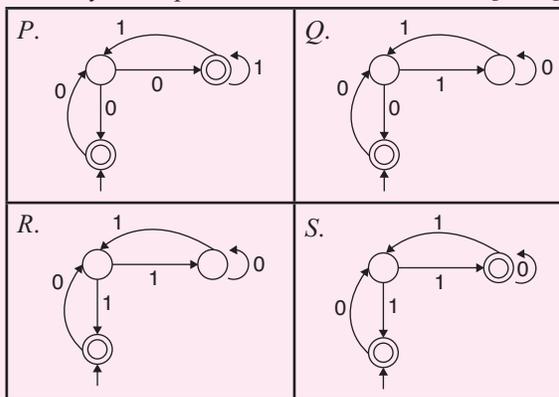
15. The number of final states in DFA after converting the NFA given below is:



- (A) 4
- (B) 2
- (C) 3
- (D) 1

PREVIOUS YEARS' QUESTIONS

1. Match the following NFAs with the regular expressions they correspond to [2008]



- 1. $\epsilon + 0(01^*1 + 00)^*01^*$
- 2. $\epsilon + 0(10^*1 + 00)^*0$
- 3. $\epsilon + 0(10^*1 + 10)^*1$
- 4. $\epsilon + 0(10^*1 + 10)^*10^*$

- (A) P-2, Q-1, R-3, S-4
- (B) P-1, Q-3, R-2, S-4
- (C) P-1, Q-2, R-3, S-4
- (D) P-3, Q-2, R-1, S-4

2. Which of the following are regular sets?

- I. $\{a^n b^{2m} \mid n \geq 0, m \geq 0\}$
- II. $\{a^n b^m \mid n = 2m\}$
- III. $\{a^n b^m \mid n \neq m\}$
- IV. $\{xycy \mid x, y \in \{a, b\}^*\}$

- (A) I and IV only
- (B) I and III only
- (C) I only
- (D) IV only

3. Which one of the following languages over the alphabet $\{0, 1\}$ is described by the regular expression:

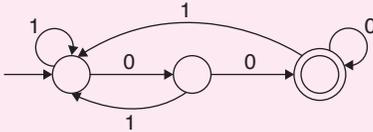
- (A) The set of all strings containing the substring 00.
- (B) The set of all strings containing atmost two 0's.
- (C) The set of all strings containing at least two 0's.

- (D) The set of all strings that begin and end with either 0 or 1.
4. Which one of the following is FALSE? [2009]
- (A) There is a unique minimal DFA for every regular language.
- (B) Every NFA can be converted to an equivalent PDA.
- (C) Complement of every context-free language is recursive.
- (D) Every non-deterministic PDA can be converted to an equivalent deterministic PDA.
5. Match all items in Group 1 with correct options from those given in Group 2. [2009]

Group 1		Group 2	
P.	Regular expression	1.	Syntax analysis
Q.	Pushdown automata	2.	Code generation
R.	Dataflow analysis	3.	Lexical analysis
S.	Register allocation	4.	Code optimization

- (A) P-4, Q-1, R-2, S-3 (B) P-3, Q-1, R-4, S-2
 (C) P-3, Q-4, R-1, S-2 (D) P-2, Q-1, R-4, S-3

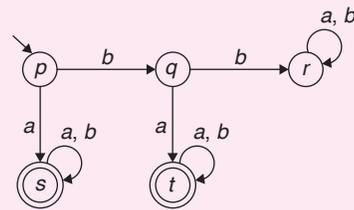
6.



The above DFA accepts the set of all strings over $\{0, 1\}$ that [2009]

- (A) Begin either with 0 or 1
 (B) End with 0
 (C) End with 00
 (D) Contain the substring 00.
7. Let $L = \{w \in (0 + 1)^* \mid w \text{ has even number of 1's}\}$, i.e., L is the set of all bit strings with even number of 1's. Which one of the regular expressions below represents L ? [2010]
- (A) $(0^*10^*1)^*$ (B) $0^*(10^*10^*)^*$
 (C) $0^*(10^*1)^*0^*$ (D) $0^*1(10^*1)^*10^*$
8. Consider the languages $L_1 = \{0^i1^j \mid i \neq j\}$, $L_2 = \{0^i1^j \mid i = j\}$, $L_3 = \{0^i1^j \mid i = 2j + 1\}$, $L_4 = \{0^i1^j \mid i \neq 2j\}$. Which one of the following statements is true? [2010]
- (A) Only L_2 is context free
 (B) Only L_2 and L_3 are context free
 (C) Only L_1 and L_2 are context free
 (D) All are context free
9. Let w be any string of length n in $\{0, 1\}^*$. Let L be the set of all substrings of w . What is the minimum number of states in a non-deterministic finite automaton that accepts L ? [2010]

- (A) $n-1$ (B) n
 (C) $n+1$ (D) 2^{n-1}
10. Let P be a regular language and Q be a context-free language such that $Q \subseteq P$ (For example let P be the language represented by the regular expression p^*q^* and Q be $\{p^nq^n \mid n \in N\}$), Then which of the following is ALWAYS regular? [2011]
- (A) $P \cap Q$ (B) $P - Q$
 (C) $\Sigma^* - P$ (D) $\Sigma^* - Q$
11. A deterministic finite automaton (DFA) D with alphabet $\Sigma = \{a, b\}$ is given below.

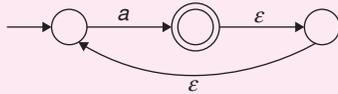


Which of the following finite state machines is a valid minimal DFA which accepts the same language as D ? [2011]

- (A)
- (B)
- (C)
- (D)

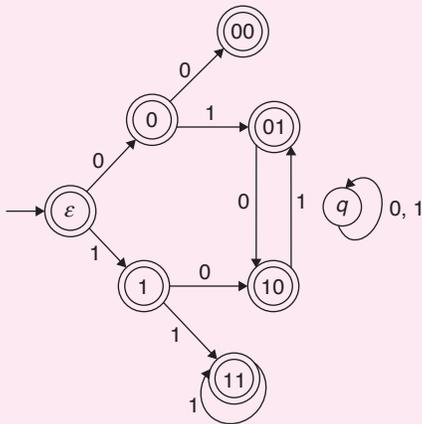
12. Given the language $L = \{ab, aa, baa\}$, which of the following strings are in L^* ? [2012]
- (1) $abaabaaabaa$ (2) $aaaabaaaa$
 (3) $baaaaabaaaaab$ (4) $baaaaabaa$
 (A) 1, 2 and 3 (B) 2, 3 and 4
 (C) 1, 2 and 4 (D) 1, 3 and 4

13. What is the complement of the language accepted by the NFA shown below?



Assume $\Sigma = \{a\}$ and ϵ is the empty string. [2012]

- (A) \emptyset (B) $\{\epsilon\}$
 (C) a^* (D) $\{a, \epsilon\}$
14. Consider the set of strings on $\{0, 1\}$ in which, every substring of 3 symbols has at most two zeros. For example, 001110 and 011001 are in the language, but 100010 are not. All strings of length less than 3 are also in the language. A partially completed DFA that accepts this language is shown below.



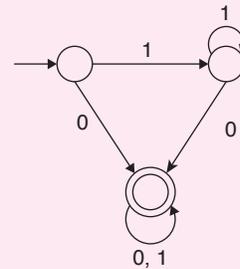
The missing arcs in the DFA are [2012]

- (A)
- | | 00 | 01 | 10 | 11 | q |
|----|----|----|----|----|---|
| 00 | 1 | 0 | | | |
| 01 | | | | 1 | |
| 10 | 0 | | | | |
| 11 | | | 0 | | |
- (B)
- | | 00 | 01 | 10 | 11 | q |
|----|----|----|----|----|---|
| 00 | 0 | | | | 1 |
| 01 | | 1 | | | |
| 10 | | | | 0 | |
| 11 | | 0 | | | |
- (C)
- | | 00 | 01 | 10 | 11 | q |
|----|----|----|----|----|---|
| 00 | | 1 | | | 0 |
| 01 | | 1 | | | |
| 10 | | | 0 | | |
| 11 | | 0 | | | |

- (D)

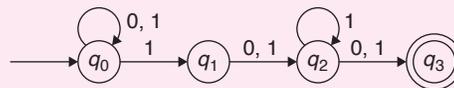
	00	01	10	11	q
00		1			0
01				1	
10	0				
11			0		

15. Consider the languages $L_1 = \Phi$ and $L_2 = \{a\}$. Which one of the following represents $L_1 L_2^* U L_1^*$? [2013]
- (A) $\{\epsilon\}$ (B) Φ
 (C) a^* (D) $\{\epsilon, a\}$
16. Consider the DFA A given below



Which of the following are FALSE?

- Complement of $L(A)$ is context-free.
 - $L(A) = L((11^*0 + 0)(0 + 1)^*0^*1^*)$
 - For the language accepted by A , A is the minimal DFA.
 - A accepts all strings over $\{0, 1\}$ of length at least 2. [2013]
- (A) 1 and 3 only (B) 2 and 4 only
 (C) 2 and 3 only (D) 3 and 4 only
17. Consider the finite automaton in the following figure. [2014]



What is the set of reachable states for the input string 0011?

- (A) $\{q_0, q_1, q_2\}$ (B) $\{q_0, q_1\}$
 (C) $\{q_0, q_1, q_2, q_3\}$ (D) $\{q_3\}$
18. If $L_1 = \{a^n | n \geq 0\}$ and $L_2 = \{b^n | n \geq 0\}$, consider the statements [2014]

- (I) $L_1 \cdot L_2$ is a regular language
 (II) $L_1 \cdot L_2 = \{a^n b^n | n \geq 0\}$
- Which one of the following is CORRECT?

- (A) Only (I) (B) Only (II)
 (C) Both (I) and (II) (D) Neither (I) nor (II)
19. Let $L_1 = \{w \in \{0, 1\}^* | w \text{ has at least as many occurrences of } (110)\text{'s as } (011)\text{'s}\}$. Let $L_2 = \{w \in \{0, 1\}^* | w \text{ has at least}$

as many occurrences of (000)'s as (111)'s}. Which one of the following is TRUE? [2014]

- (A) L_1 is regular but not L_2
- (B) L_2 is regular but not L_1
- (C) Both L_1 and L_2 are regular
- (D) Neither L_1 nor L_2 are regular

20. The length of the shortest string NOT in the language (over $\Sigma = \{a, b\}$) of the following regular expression is _____. [2014]

$$a^*b^*(ba)^*a^*$$

21. Let Σ be finite non-empty alphabet and let 2^{Σ^*} be the power set of Σ^* . Which one of the following is TRUE? [2014]

- (A) Both 2^{Σ^*} and Σ^* are countable
- (B) 2^{Σ^*} is countable and Σ^* is uncountable
- (C) 2^{Σ^*} is uncountable and Σ^* is countable
- (D) Both 2^{Σ^*} and Σ^* are uncountable

1. $\epsilon + 0(01^*1 + 00)^*01^*$
2. $\epsilon + 0(10^*1 + 00)^*0$
3. $\epsilon + 0(10^*1 + 10)^*1$
4. $\epsilon + 0(10^*1 + 10)^*10^*$

- (A) $P - 2, Q - 1, R - 3, S - 4$
- (B) $P - 1, Q - 3, R - 2, S - 4$
- (C) $P - 1, Q - 2, R - 3, S - 4$
- (D) $P - 3, Q - 2, R - 1, S - 4$

22. Consider the DFAs M and N given above. The number of states in a minimal DFA that accepts the language $L(M) \cap L(N)$ is _____. [2015]



23. The number of states in the minimal deterministic finite automaton corresponding to the regular expression $(0 + 1)^*(10)$ is _____. [2015]

24. Which of the following languages is/are regular? [2015]

$L_1: \{wxw^R \mid w_1 x \in \{a, b\}^* \text{ and } |w|, |x| > 0\}$, w^R is the reverse of string w

$L_2: \{a^n b^m \mid m \neq n \text{ and } n, m \geq 0\}$

$L_3: \{a^p b^q c^r \mid p, q, r \geq 0\}$

- (A) L_1 and L_3 only
- (B) L_2 only
- (C) L_2 and L_3 only
- (D) L_3 only

25. Consider the alphabet $\Sigma = \{0, 1\}$, the null/empty string λ and the sets of strings X_0, X_1 and X_2 generated by the corresponding non-terminals of a regular grammar. X_0, X_1 and X_2 are related as follows

$$X_0 = 1X_1$$

$$X_1 = 0X_1 + 1X_2$$

$$X_2 = 0X_1 + \{\lambda\}$$

Which one of the following choices precisely represents the strings in X_0 ? [2015]

- (A) $10(0^* + (10)^*)1$
- (B) $10(0^* + (10)^*)^*1$
- (C) $1(0 + 10)^*1$
- (D) $10(0 + 10)^*1 + 110(0 + 10)^*1$

26. Let L be the language represented by the regular expression $\Sigma^* 0011 \Sigma^*$ where $\Sigma = \{0, 1\}$. What is the minimum number of states in a DFA that recognizes \bar{L} (complement of L)? [2015]

- (A) 4
- (B) 5
- (C) 6
- (D) 8

27. Which of the following languages is generated by the given grammar? [2016]

$$S \rightarrow aS \mid bS \mid \epsilon$$

- (A) $\{a^n b^m \mid n, m \geq 0\}$
- (B) $\{w \in \{a, b\}^* \mid w \text{ has equal number of } a\text{'s and } b\text{'s}\}$
- (C) $\{a^n \mid n \geq 0\} \cup \{b^n \mid n \geq 0\} \cup \{a^n b^n \mid n \geq 0\}$
- (D) $\{a, b\}^*$

28. Which of the following decision problems are undecidable? [2016]

I. Given NFAs N_1 and N_2 , is

$$L(N_1) \cap L(N_2) = \Phi?$$

II. Given a CFG $G = (N, \Sigma, P, S)$ and a string $x \in \Sigma^*$, does $x \in L(G)$?

III. Given CFGs G_1 and G_2 , is

$$L(G_1) = L(G_2)?$$

IV. Given a TM M , is $L(M) = \Phi$?

- (A) I and IV only
- (B) II and III only
- (C) III and IV only
- (D) II and IV only

29. Which one of the following regular expressions represents the language: the set of all binary strings having two consecutive 0's and two consecutive 1's? [2016]

- (A) $(0+1)^* 0011 (0+1)^* + (0+1)^* 1100 (0+1)^*$
- (B) $(0+1)^* (00(0+1)^*11 + 11(0+1)^*00) (0+1)^*$
- (C) $(0+1)^* 00 (0+1)^* + (0+1)^* 11 (0+1)^*$
- (D) $00 (0+1)^* 11 + 11 (0+1)^* 00$

30. The number of states in the minimum sized DFA that accepts the language defined by the regular expression $(0+1)^* (0+1) (0+1)^*$ is _____. [2016]

31. Language L_1 is defined by the grammar: $S_1 \rightarrow aS_1b \mid \epsilon$
Language L_2 is defined by the grammar: $S_2 \rightarrow abS_2 \mid \epsilon$

Consider the following statements:

$P: L_1$ is regular

$Q: L_2$ is regular

Which one of the following is TRUE? [2016]

- (A) Both P and Q are true
- (B) P is true and Q is false
- (C) P is false and Q is true
- (D) Both P and Q are false

32. Consider the following two statements:

- I. If all states of an NFA are accepting states then the language accepted by the NFA is Σ^* .
- II. There exists a regular language A such that for all languages B , $A \cap B$ is regular.

Which one of the following is CORRECT? [2016]

- (A) Only I is true
- (B) Only II is true
- (C) Both I and II are true
- (D) Both I and II are false

33. Consider the language L given by the regular expression $(a + b)^*b(a + b)$ over the alphabet $\{a, b\}$. The smallest number of states needed in a deterministic finite-state automaton (DFA) accepting L is _____.

[2017]

34. The minimum possible number of states of a deterministic finite automaton that accepts the regular language $L = \{w_1aw_2 \mid w_1, w_2 \in \{a, b\}^*, |w_1| = 2, |w_2| \geq 3\}$ is _____.

[2017]

35. Let δ denote the transition function and $\hat{\delta}$ denote the extended transition function of the ϵ -NFA whose transition table is given below:

δ	ϵ	a	b
$\rightarrow q_0$	$\{q_2\}$	$\{q_1\}$	$\{q_0\}$
q_1	$\{q_2\}$	$\{q_2\}$	$\{q_3\}$
q_2	$\{q_0\}$	\emptyset	\emptyset
q_3	\emptyset	\emptyset	$\{q_2\}$

Then $\hat{\delta}(q_2, aba)$ is [2017]

- (A) \emptyset
- (B) $\{q_0, q_1, q_3\}$
- (C) $\{q_0, q_1, q_2\}$
- (D) $\{q_0, q_2, q_3\}$

36. Let N be an NFA with n states. Let k be the number of states of a minimal DFA which is equivalent to N . Which one of the following is necessarily true? [2018]

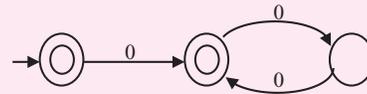
- (A) $k \geq 2^n$
- (B) $k \geq n$
- (C) $k \leq n^2$
- (D) $k \leq 2^n$

37. Given a language L , define L^i as follows:

$$L^0 = \{\epsilon\}$$

$$L^i = L^{i-1} \cdot L \text{ for all } i > 0$$

The order of a language L is defined as the smallest k such that $L^k = L^{k+1}$. Consider the language L_1 (over alphabet 0) accepted by the following automaton.



The order of L_1 is _____. [2018]

ANSWER KEYS

EXERCISES

Practice Problems 1

1. D 2. B 3. C 4. D 5. B 6. B 7. D 8. D 9. C 10. C
 11. A 12. D 13. C 14. A 15. A

Practice Problems 2

1. C 2. D 3. D 4. B 5. B 6. C 7. A 8. C 9. C 10. A
 11. B 12. C 13. C 14. C 15. B

Previous Years' Questions

1. C 2. A 3. C 4. D 5. B 6. C 7. B 8. D 9. C 10. C
 11. A 12. C 13. B 14. D 15. A 16. D 17. A 18. A 19. A 20. C
 21. C 22. 1 23. 3 24. A 25. C 26. B 27. D 28. C 29. B 30. 2
 31. C 32. B 33. 4 34. 8 35. C 36. D 37. 2

Chapter 2

Context Free Languages and Push Down Automata

LEARNING OBJECTIVES

- Context free grammar
- Context free language
- Ambiguity in context free grammars
- Removing ϵ -productions
- Removing unit productions
- Normal forms
- Chomsky's normal form
- Greiback normal form
- Closure properties of CFL's
- Push down automata
- PDAs accepting by final state and empty stack are equivalent
- Converting CFG to PDA
- Deterministic PDA

CONTEXT FREE GRAMMAR

- A context free grammar (CFG) is a finite set of variables (non-terminals) each of which represents a language. The language represented by variables is described recursively in terms of each other. The primitive symbols are called terminals.
- The rules relating variables are called productions. A typical production states that the language associated with a given variable contains strings that are formed by concatenating strings from languages of certain other variables.
- CFG is a collection of three things;
 - An alphabet Z of letters called terminals.
 - A set of symbols called non-terminals, one of which is a start symbol, S .
 - A finite set of productions of the form:
 - One terminal \rightarrow finite set of terminals and/or non-terminals.
- A CFG is defined as: $G = (V, T, P, S)$
 - Where
 - $V \rightarrow$ Finite set of variables (non-terminals)
 - $T \rightarrow$ Finite set of terminals (symbols)
 - $P \rightarrow$ Finite set of productions, each, production is of the form, $A \rightarrow \alpha, A \in V, \alpha \in (V \cup T)^*$
 - $S \rightarrow$ Start symbol

CONTEXT FREE LANGUAGE (CFL)

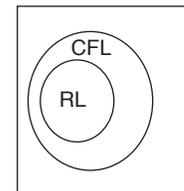
The language generated by CFG is a set of all strings of terminals that can be produced from start symbols, using the productions as

substitutions. A language generated by a CFG is called context free language (CFL).

Note: Every regular grammar is context free, so a regular language (RL) is also context free.

Family of RL's is proper subset of CFL's.

i.e., $RL \subset CFL$



Solved Examples

Example 1: What is the language that is generated by CFG, $G = S \rightarrow AB|A \rightarrow +|-|B \rightarrow CB/C|C \rightarrow 0/1/2/ \dots 9$.

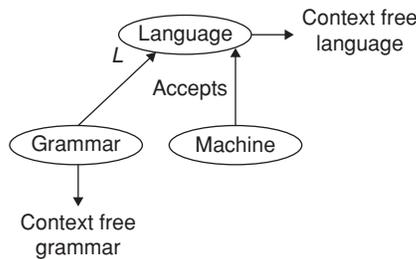
- (A) Set of all rational numbers
- (B) Set of all integers
- (C) Set of all natural numbers
- (D) Set of all complex numbers

Solution: (B)

$S \rightarrow AB|A \rightarrow +|-|B \rightarrow CB/C|C \rightarrow 0/1/2/ \dots 9$

Consider-18 (integer)

$S \rightarrow AB$
 $\rightarrow -B$
 $\rightarrow -CB$
 $\rightarrow -1B$
 $\rightarrow -18$



AMBIGUITY IN CONTEXT FREE GRAMMARS

A CFG, G is called ambiguous if there is $w \in L(G)$ such that w has (at least) two different parse trees with respect to G .

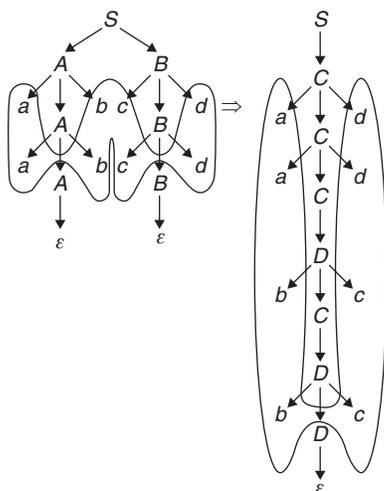
Example 2: The language, $L = \{a^n b^n c^m d^m / n \geq 0, m \geq 0\} \cup \{a^n b^m c^m d^n / n \geq 0, m \geq 0\}$ is designed in CFG, G . The Grammar is

- (A) Ambiguous
- (B) Unambiguous
- (C) Cannot be determined
- (D) None of above

Solution: (A)
CFG G for given language L is:

$S \rightarrow AB|C$
 $A \rightarrow aAb|\epsilon$
 $B \rightarrow cBd|\epsilon$
 $C \rightarrow aCd|D$
 $D \rightarrow bDc|\epsilon$

It's an inherently ambiguous grammar.
Consider string, aabbccdd



Note:

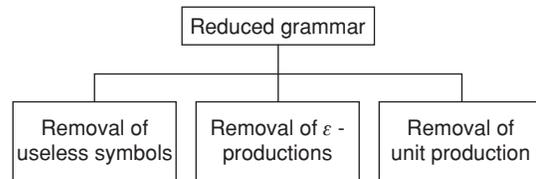
- A context free language with property that all grammars that generate it are ambiguous is inherently ambiguous.
- Inherently ambiguous grammars cannot convert to unambiguous grammars.

MINIMIZATION OF CONTEXT FREE GRAMMAR

- Grammar may consist of some extra symbols (non-terminals). Having extra symbols unnecessarily increases the length of grammar.
- Simplification of grammar means reduction of grammar.

The properties of reduced grammar are:

1. Each variable (non-terminal) and each terminal of G appears in the derivation of some word in L .
2. There should not be any production as $X \rightarrow Y$ where X and Y are non-terminals.
3. If ϵ is not in language L , then there need not be production $X \rightarrow \epsilon$.



Removal of Useless Symbols

- Any symbol is useful when it appears on right hand side, in the production rule and generates some terminal string. If no such derivation exists, then it is supposed to be a useless symbol.
- A symbol P is useful, if there exists some derivation

$$S^* \Rightarrow \alpha P B \text{ and } \alpha P B \overset{*}{\Rightarrow} W$$

Then P is said to be useful symbol.

Example 3: A grammar G' , is generated by removing useless symbols from G defined below. The obtained G' contains productions:

- $S \rightarrow aA|bB$
 $A \rightarrow aA|a$
 $B \rightarrow bB$
 $D \rightarrow ab|Ea$
 $E \rightarrow aC|d$
- (A) $S \rightarrow aA$
 $A \rightarrow aA|a$
 - (B) $S \rightarrow aS|bA|C$
 $A \rightarrow a$
 $C \rightarrow aCd$
 - (C) $S \rightarrow aA|bB$
 $A \rightarrow aA|a$
 $B \rightarrow bB$
 - (D) Cannot remove useless symbols

Solution:

$S \rightarrow aA \rightarrow aaA \rightarrow aaaA \rightarrow aaaa \checkmark$
 $B \rightarrow bB \rightarrow bbB \rightarrow bbbB \rightarrow bbbbB \dots$ (string cannot be generated)
 $\therefore B$ is useless
 D and E cannot be generated from 'S'. So, eliminate. Hence G' contains
 $\therefore S \rightarrow aA$
 $A \rightarrow aA|a$

Removing ϵ -Productions

A production of the form $A \rightarrow \epsilon$ is called an ϵ -production. If A is a non-terminal and $A \rightarrow (^*) \epsilon$, then A is called a 'nullable non-terminal'. So eliminate such productions without changing meaning of grammar.

Example 4: The grammar, G is given below. The CFG generated after eliminating ϵ -production is:

- $S \rightarrow ABC$
- $A \rightarrow BC|a$
- $B \rightarrow bAC|\epsilon$
- $C \rightarrow cAB|\epsilon$
- (A) $S \rightarrow ABC|AB|BC|CA$
 $A \rightarrow BC|B|C$
 $B \rightarrow bAC|bA|bC$
 $C \rightarrow cAB|cA|cB$
- (B) $S \rightarrow ABC$
 $A \rightarrow BC$
 $B \rightarrow bAC$
 $C \rightarrow cAB$
- (C) $S \rightarrow ABC|BC|AC|AB|A|B|C$
 $A \rightarrow BC|B|C|a$
 $B \rightarrow bAC|bA|bC|b$
 $C \rightarrow cAB|cA|cB|c$
- (D) None of these

Solution (C)

$B \rightarrow \epsilon, C \rightarrow \epsilon$
 $\Rightarrow A \rightarrow \epsilon$
 \therefore Remove ϵ -productions and obtained CFG is Choice (C).

Removing Unit Productions

- A production of form $A \rightarrow B$, where A and B are both non-terminals, is called a 'unit production'.
- Presence of unit production in a grammar increases the cost of derivation.

Example 5: The total number of productions obtained by removing unit production from the Grammar,

$A \rightarrow PQ$
 $P \rightarrow 0$
 $Q \rightarrow R|1$
 $R \rightarrow S$
 $S \rightarrow W|1R$
 $W \rightarrow 2|P1$

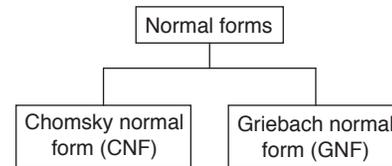
- (A) 9
- (B) 2
- (C) 3
- (D) 5

Solution: (A)

$A \rightarrow PQ$
 $Q \rightarrow R \rightarrow S \rightarrow W \rightarrow 2$
 $\Rightarrow Q, R, S \rightarrow$ Unit production
 $A \rightarrow PQ$
 $P \rightarrow 0$
 $Q \rightarrow R|1$
 \downarrow
 $\Rightarrow Q \rightarrow 2|P1|1R|1$ (substitute the production of R, S, W)
 $\therefore A \rightarrow PQ$
 $P \rightarrow 0$
 $Q \rightarrow 2|P1|1R|1$
 $R \rightarrow 2|P1|1R$
 $\therefore 9$ – Productions

NORMAL FORMS

- It is necessary to have a grammar in some specific form so, grammar normalization is needed.



That is, There should be fixed number of terminals and non-terminals, in CFG.

Chomsky's Normal Form (CNF)

- A context free grammar (CFG), $G = (V, \Sigma, R, S)$ is said to be in CNF, if and only if every rule in R is of one of the following forms
 1. $A \rightarrow a$, for some $A \in V$ and some $a \in \Sigma$
 2. $A \rightarrow BC$, for some $A \in V$ and $B, C \in V \cup \{S\}$
 3. $S \rightarrow \epsilon$
- Every rule either replaces a variable by a single character or by a pair of variables except the start symbol and the only rule that can have the empty word as it's right hand side must have start symbol as it's left hand side.

Note: Every parse tree for a grammar in CNF must be a binary tree and the parse tree for any non-empty word cannot have any leaves labeled with ϵ in it.

Transforming of a grammar to CNF

- In order to construct the grammar G in CNF that is equivalent to a given grammar G , first identify how exactly G can violate the rules for a CNF. Since CNF only restricts the rules in G , see only at R . The 'bad' cases of rules are:
- $A \rightarrow uSv$ where $A \in V$ and $u, v \in (V \cup \Sigma)^*$. The start symbol must not appear on the right-hand side of any rule. This is called 'start symbol rule'.

- **To remove ‘start symbol rule’**, add a new symbol, so make it the start symbol in new grammar G_1 , and add the single rule $S_0 \rightarrow S$ to R to get the rules for G_1 . Since S_0 does not appear in any rules, the new grammar has no start symbol rules.

- $A \rightarrow \epsilon$ where $A \in V \cup \{S\}$. The only symbol that can be replaced by the word is start symbol. This is called ‘ ϵ -rules’.

- **To remove ‘ ϵ -rules’**, identify all variables that can yield the empty string, either directly or indirectly.

These variables are ‘nullable’. Remove all direct rules $A \rightarrow \epsilon$ from the grammar and fix up the grammar by removing all occurrences of nullable variables from the right hand sides of all rules.

$A \rightarrow B$ where $A, B \in V$. The only rules involving variables on the right-hand side must have exactly two of them. This is called ‘unit rules’.

- **To remove ‘Unit rules’**, identify a set of unit pairs.

These are pairs of symbols (A, B) , where $A \xRightarrow{*} B$. Then remove all unit rules by copying right-hand sides. If there is a rule $A \rightarrow B$, (A, B) is a unit pair. Then, if there is a rule $B \rightarrow W$, derive W from A by $A \rightarrow B, B \rightarrow W$. To remove the unit rule and still generate an equivalent grammar, add the right-hand side W to the rules for A directly, $A \rightarrow W$.

$A \rightarrow W$ where $A \in V, W \in (V \cup \Sigma)^*$ and W contains at least one character and at least one variable. The only rules where character appear on right-hand side must have exactly one character as right-hand side. This is called ‘mixed rules’.

- **To remove ‘mixed rules’**, Let $A \rightarrow W \in R_3$ is a mixed rules. Then write W as $W = V_0 C_1 V_1 \dots V_{n-1} C_n V_n$, where $C_i \in \Sigma$ are occurrences of characters, and the $V_i \in V^*$ are strings of only variables. Then add a new symbol, C_i to V_4 for every character C_i and add the rules $C_i \rightarrow c_i$ to R_4 . Finally define $W^1 = V_0 C_1 V_1 \dots V_{n-1} C_n V_n \in V^*$ and add rules $A \rightarrow W^1$ to R_4 . If the rule $A \rightarrow W$ is part of the derivation for some word, replace that single rule by applying rule $A \rightarrow W^1$ first and then replacing all C_i by c_i using their respective rules.

$A \rightarrow w$ Where $A \in V$ and $W \in (V \cup \Sigma)^*$ with $|w| > 2$. Rules must have one symbol (character) or two variables (two variables as right hand side). These are called long rules.

- **To remove ‘long rules’**, Let $A \rightarrow B_1 \dots B_n$ be a long rule, i.e., $n > 2$. B_i is all variables. Break up every single long rule, into several ‘short’ rules, by introducing new ‘helper variables’ and splitting right hand side from left to right: add new symbols A_1, \dots, A_{n-2} to set of variables and add following rules to R_5 : $A \rightarrow B_1 A_1, A_1 \rightarrow B_2 A_2, \dots, A_{n-2} \rightarrow B_{n-1} B_n$.

Example 6: Consider grammar, $G = S \rightarrow ASB, A \rightarrow aAS|a\epsilon, B \rightarrow SbS|A|bb$. The CNF generated contains ____ non-terminals.

- (A) 5 (B) 6
(C) 9 (D) 11

Solution: (C)

Add new start state:

$$S_0 \rightarrow S$$

$$S \rightarrow ASB$$

$$A \rightarrow aAS|a\epsilon$$

$$B \rightarrow SbS|A|bb$$

Eliminate ϵ -rules

$$A \rightarrow \epsilon:$$

$$S_0 \rightarrow S$$

$$S \rightarrow ASB|SB$$

$$A \rightarrow aAS|aAS$$

$$B \rightarrow SbS|A|bb|\epsilon$$

Eliminate $B \rightarrow \epsilon$:

$$S_0 \rightarrow S$$

$$S \rightarrow ASB|SB|S|AS$$

$$A \rightarrow aAS|aAS$$

$$B \rightarrow SbS|A|bb$$

Remove Unit rules:

$$B \rightarrow A:$$

$$S_0 \rightarrow S$$

$$S \rightarrow ASB|SB|S|AS$$

$$A \rightarrow aAS|aAS$$

$$B \rightarrow SbS|bb|aAS|aAS$$

$$S \rightarrow S:$$

$$S_0 \rightarrow S$$

$$S_0 \rightarrow ASB|SB|AS$$

$$A \rightarrow aAS|aAS$$

$$B \rightarrow SbS|bb|aAS|aAS$$

$$S_0 \rightarrow S:$$

$$S_0 \rightarrow ASB|SB|AS$$

$$S \rightarrow ASB|SB|AS$$

$$A \rightarrow aAS|aAS$$

$$B \rightarrow SbS|bb|aAS|aAS$$

Replace rules which have more than two symbols:

$$S_0 \rightarrow ASB: S_0 \rightarrow AU_1 \text{ and } U_1 \rightarrow SB$$

$$\therefore S_0 \rightarrow AU_1|SB|AS$$

$$S \rightarrow AU_1|SB|AS$$

$$A \rightarrow aAS|aAS$$

$$B \rightarrow SbS|bb|aAS|aAS$$

$$U_1 \rightarrow SB$$

$$A \rightarrow aAS \Rightarrow A \rightarrow aU_2 \text{ and } U_2 \rightarrow AS \text{ and } B \rightarrow SbS$$

$$\Rightarrow B \rightarrow SU_3 \text{ and } U_3 \rightarrow bS$$

$$\therefore S_0 \rightarrow AU_1|SB|AS$$

$$S \rightarrow AU_1|SB|AS$$

$$A \rightarrow aU_2|aAS$$

$$B \rightarrow SU_3|bb|aU_2|aAS$$

$$U_1 \rightarrow SB$$

$$U_2 \rightarrow AS$$

$$U_3 \rightarrow bS$$

Eliminate rules which have terminals and variables or two terminals.

$$\text{Let } V_1 \rightarrow a, V_2 \rightarrow b$$

$$\therefore S_0 \rightarrow AU_1|AS|SB$$

$$S \rightarrow AU_1|SB|AS$$

$$A \rightarrow V_1 U_2 | a | V_1 S$$

$B \rightarrow SU_3|V_2V_2|V_1U_2|a|V_1S$
 $U_1 \rightarrow SB$
 $U_2 \rightarrow AS$
 $U_3 \rightarrow V_2S$
 $V_1 \rightarrow a$
 $V_2 \rightarrow b$
 \therefore Nine non-terminals.

Greiback Normal Form (GNF)

- A CFG, $G = (V, T, R, S)$ is said to be in GNF, if every production is of form $A \rightarrow a\alpha$ where $a \in T, \alpha \in V^*$, i.e., α is a string of zero or more variables.
- Left recursion in R can be eliminated by following schema: If $A \rightarrow A\alpha_1 | A\alpha_2 | \dots | A\alpha_r | \beta_1 | \beta_2 | \dots | \beta_s$, then replace the above rules by
 - $A \rightarrow \beta_i | \beta_i Z, 1 \leq i \leq s$
 - $Z \rightarrow \alpha_i | \alpha_i Z, 1 \leq i \leq r$
- If $G = (V, T, R, S)$ is a CFG, then another CFG, $G_1 = (V_1, T, R_1, S)$ can be constructed in GNF $\exists L(G_1) = L(G) - \{\epsilon\}$.

The step wise algorithm is as follows:

- Eliminate null production, unit productions and useless symbols from the grammar G and then construct a $G^1 = (V^1, T, R^1, S)$ in CNF generating the language $L(G^1) = L(G) - \{\epsilon\}$.
- Rename the variables like A_1, A_2, \dots, A_n starting with $S = A_1$.
- Modify the rules in R^1 , so that if $A_i \rightarrow A_j \gamma \in R^1$ then $j > i$.
- Starting with A_1 and proceeding to A_n , can be obtained as:
 - Assume that productions have been modified so that for $1 \leq i \leq k, A_i \rightarrow A_j \gamma \in R^1$ only if $j > i$
 - If $A_k \rightarrow A_j \gamma$ is a production with $j < k$, generate a new set of productions substituting for A_j , the body of each A_j production.
 - Repeating (b) atleast $k - 1$ times, obtains rules of the form $A_k \rightarrow A_p \gamma, p \geq k$.
 - Replace rules $A_k \rightarrow A_k \gamma$ by removing left-recursion.
- Modify the $A_i \rightarrow A_j \gamma$ for $i = n - 1, n - 2, \dots, 1$ in desired form at same time change z production rules.

Example 7: A grammar G is defined with rules $S \rightarrow XA|BB, B \rightarrow b|SB, X \rightarrow b, A \rightarrow a$. The normalized GNF of G contains ____ productions.

- (A) 17 (B) 19
(C) 5 (D) 16

Solution: (B)

- The Grammar, G is already in CNF.
- Re-label with variables
 - S with A_1
 - X with A_2
 - A with A_3
 - B with A_4

Grammar, G now is:

$A_1 \rightarrow A_2 A_3 | A_4 A_4$
 $A_4 \rightarrow b | A_1 A_4$
 $A_2 \rightarrow b$
 $A_3 \rightarrow a$

- Identify all productions which do not conform to any of the types listed below:

$A_i \rightarrow A_j x_k \exists j > i$
 $Z_i \rightarrow A_j x_k \exists j \leq n$
 $A_i \rightarrow a x_k \exists x_k \in V^* \text{ and } a \in T$

- $A_4 \rightarrow A_1 A_4 \dots$ identified
- $A_4 \rightarrow A_1 A_4 | b$

To eliminate A_1 , use substitution rule, $A_1 \rightarrow A_2 A_3 | A_4 A_4$
 $\therefore A_4 \rightarrow A_2 A_3 A_4 | A_4 A_4 | b$
 Substitute $A_2 \rightarrow b$

$\therefore A_4 \rightarrow b A_3 A_4 | A_4 A_4 | b$
 $A_4 \rightarrow A_4 A_4 A_4$ is left recursive. So, remove left recursion i.e., $A_4 \rightarrow b A_3 A_4 | b | b A_3 A_4 Z | b Z$
 $Z \rightarrow A_4 A_4 | A_4 A_4 Z$

- Now, $G = A_1 \rightarrow A_2 A_3 | A_4 A_4$
 $A_4 \rightarrow b A_3 A_4 | b | b A_3 A_4 Z | b Z$
 $Z \rightarrow A_4 A_4 | A_4 A_4 Z$
 $A_2 \rightarrow b$
 $A_3 \rightarrow a$

- A_1, Z are not in GNF. So,

For $A_1 \rightarrow A_2 A_3 | A_4 A_4$:
 Substitute for A_2 and A_4 to convert it to GNF
 $A_1 \rightarrow b A_3 | b A_3 A_4 A_4 | b A_4 | b A_3 A_4 Z A_4 | b Z A_4$
 For $Z \rightarrow A_4 A_4 | A_4 A_4 Z$

substitute for A_4 to convert it to GNF
 $Z \rightarrow b A_3 A_4 A_4 | b A_4 | b A_3 A_4 Z A_4 | b Z A_4 | b A_3 A_4 A_4 Z | b A_4 Z | b A_3 A_4 Z A_4 Z | b Z A_4 Z$

\therefore Final GNF is:
 $A_1 \rightarrow b A_3 | b A_3 A_4 A_4 | b A_4 | b A_3 A_4 Z A_4 | b Z A_4$
 $A_4 \rightarrow b A_3 A_4 | b | b A_3 A_4 Z | b Z$
 $A_2 \rightarrow b$
 $A_3 \rightarrow a$
 $Z \rightarrow b A_3 A_4 A_4 | b A_4 | b A_3 A_4 Z A_4 | b Z A_4 | b A_3 A_4 A_4 Z | b A_4 Z | b A_3 A_4 Z A_4 Z | b Z A_4 Z$
 \therefore 19 productions.

PUMPING LEMMA FOR CONTEXT FREE LANGUAGES

Let 'L' be context free language. There exists some integer, $m \exists \forall w \in L$, with $|w| \geq m, w = uvxyz$ with $|vxy| \leq m$ and $|vy| \geq 1 \exists u v^i x y^j z \in L \forall i = 0, 1, 2, 3, \dots$

Note: Pumping lemma is used to show that a language is Not context free.

Example 8: The language $\{a^n b^m c^n d^{(n+m)}; m, n \geq 0\}$ is

- (A) Regular
(B) Context free but not regular
(C) Neither context free nor regular
(D) Cannot be determined

Solution: (C)

$$L = \{a^n b^m c^n d^{(n+m)} : m, n \geq 0\}$$

Clearly, L is not regular because, number of a's and number of b's must be known to compute number of d's.

' L ' is not context free because, Let $w = a^M b^M c^M d^{2M}$. Clearly neither v nor y can cross regions and include more than one letter, since if that happened; letters obtained will be out of order when pumped.

So, consider cases, where v and y fall within a single region. Consider 4-regions corresponding to a, b, c and d.

(1, 1) → change number of a's and they won't match c's any more.

(1, 2) → If v is not empty, change a's and they won't match with c's. If y is non-empty, number of b's changed won't have right number of d's.

(1, 3), (1, 4) → ruled out $\because |vxy| \leq M$

(2, 2) → Change number of b's and they won't match right number of d's.

(2, 3) → If v is non-empty, change number of b's without changing number of d's. If y is not empty, change c's and they'll no longer match a's.

(2, 4) → ruled out $\because |vxy| \leq M$

(3, 3) → Change number of c's and they won't match a's.

(3, 4) → If v is not empty change c's and they won't match a's. If y is not empty, change d's without changing b's.

(4, 4) → change d's without changing a's or b's.

$\therefore L$ is not context free.

(A) Regular

(B) Context free

(C) Regular but not context free

(D) Cannot be determined

Solution: (B)

$L_1 = \{a^n b^n : n > 0\}$ is context free

$L_2 = \{a^{100} b^{100}\}$ is regular

$\overline{L_2} = \{(a+b)^*\} - \{a^{100} b^{100}\}$ is regular

$\{a^n b^n\}$ context free

$\overline{L_2} = \{(a+b)^*\} - \{a^{100} b^{100}\}$ is regular

$\{a^n b^n\} \cap \overline{L_2} \rightarrow$ context free

$\{a^n b^n\} \cap \overline{L_2} = \{a^n b^n : n \neq 100, n \geq 0\} = L$ is context free:

Table 1 Comparing Regular and Context free Languages:

Regular Language	CFL
Regular expression or regular grammar	Context free grammar
Recognize the language	Parses the language
These are DFSA's	These are NDPDA's
Minimize FSA's	Find deterministic grammar.
Closed under: Concatenation Union Kleen star Complement Intersection	Closed under: Concatenation Union Kleen star

CLOSURE PROPERTIES OF CFL'S

1. CFL's are closed under union: For CFL's L_1, L_2 with CFG's G_1, G_2 and start variables S_1, S_2 . The grammar of Union $L_1 \cup L_2$ has new start symbol S and additional production $S \rightarrow S_1 | S_2$

2. CFL's are closed under concatenation: For CFL's L_1, L_2 with CFG's G_1, G_2 and start variables S_1, S_2 . The grammar of concatenation $L_1 L_2$ has new start variables S and additional production: $S \rightarrow S_1 S_2$

3. CFL's are closed under star operation: For CFL L , with CFG G and start variable S . The grammar of the start operation L^* has new start variable S_1 and additional production:

$$S_1 \rightarrow S S_1 | \epsilon$$

4. CFL's are not closed under intersection: If L_1, L_2 are two context free languages, $L_1 \cap L_2$ not necessarily be context free.

5. CFL's are not closed under complement: If L is context free language, \overline{L} not necessarily be context free.

6. Intersection of CFL's and regular language: (regular closure): If L_1 is a CFL and R_2 is a regular language then $L_1 \cap L_2$ is a CFL.

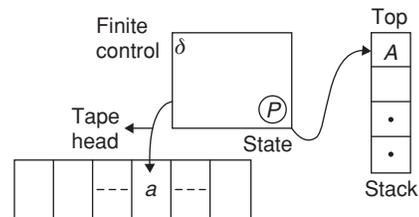
Example 9: The language, $L_1 = \{a^n b^n : n \geq 0\}$ and $L_2 = \{a^{100} b^{100}\}$. The relation $L_1 \cap \overline{L_2}$ is _____

PUSH DOWN AUTOMATA (PDA)

A push down automata is merely a finite automata with a stack added to it.

PDA is used to generate context free language.

The stack allows for unbounded memorization.



Input tape: The tape is divided into finitely many cells. Each cell contains a symbol in an alphabet, Σ .

Stack: The stack head always scans the top symbol of the stack. It performs two basic operations.

- Push: Add a new symbol at the top
- Pop: Read and remove the top symbol

Tape head: The head scans at a cell on the tape and can read a symbol on the cell. In each move, the head can move to the right cell.

Finite control: The finite control has finitely many states which form a set Q . For each move, the state is changed according to the evaluation of transition function.

A PDA is defined as: $P = (Q, \Sigma, \Gamma, \delta, q_0, z_0, F)$

Where Q : set of States

Σ : Input alphabet

Γ : Stack symbol

δ : Transition function

q_0 : Start state

z_0 : Initial stack top symbol

F : Final/accepting states

Transition functions $\delta: Q \times \Gamma \times \Sigma \Rightarrow Q \times \Gamma$

Q : Old state

Γ : Stack top

Σ : Input symbol

Q : New state, Γ : New stack top

PDA's instantaneous description (IDs): A PDA has a configuration at any given instance: (q, w, y)

$q \rightarrow$ current state

$w \rightarrow$ remainder of input (i.e., unconsumed part)

$y \rightarrow$ current stack contents as a string from top to bottom of the stack.

If $\delta(q, a, x) = \{P, A\}$ is a transition, then following are also true:

- $(q, a, x) \vdash (P, \varepsilon, A)$
- $(q, aw, xB) \vdash (p, w, AB)$

Note: 1. \rightarrow : Turnstile notation and represents one move.
2. \vdash^* : represents sequence of moves.

Principles about IDs:

1. If for a PDA, $(q, x, A) \vdash^*(p, y, B)$, then for any string $w \in \Sigma^*$ and $\gamma \in \Gamma^*$, it is also true that:
 $(q, xw, A\gamma) \vdash^*(p, yw, B\gamma)$
2. If for a PDA, $(q, xw, A) \vdash^*(p, yw, B)$, then it is also true that: $(q, x, A) \vdash^*(p, y, B)$

Acceptance by PDA: There are two types of PDAs that one can design:

- Those that accept by final state or
- Those that accept by empty stack

PDAs that accept by final state: For a PDA, P , the language accepted by P , denoted by $L(P)$ by final state, is:

$$\{w \mid (q_0, w, z_0) \vdash^*(q, \varepsilon, A)\} \exists q \in F$$

PDAs that accept by empty stack: For a PDA P , the language accepted by P , denoted by $N(P)$ by empty stack, is:

$$\{w \mid (q_0, w, z_0) \vdash^*(q, \varepsilon, \varepsilon)\}, \text{ for any } q \in Q.$$

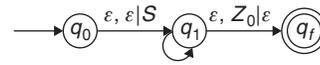
Example 10: Consider the grammar $S \rightarrow aTb \mid b, T \rightarrow Ta \mid \varepsilon$. The PDA constructed contains _____ states.

- (A) 4 (B) 3 (C) 5 (D) 2

Solution: (B)

$S \rightarrow aTb \mid b$

$T \rightarrow Ta \mid \varepsilon$



$\varepsilon, S \mid aTb$

$\varepsilon, T \mid Ta$

$\varepsilon, S \mid b$

$\varepsilon, T \mid \varepsilon$

$a, a \mid \varepsilon$

$b, b \mid \varepsilon$

Let $S \rightarrow q_0, T \rightarrow q_1$

Consider string "aab" $S \rightarrow aTb \rightarrow aTab \rightarrow aab$

$\delta(q_0, aab, z_0) \vdash \delta(q_0, \varepsilon aab, z_0)$

$\vdash \delta(q_1, aab, q_0z_0)$

$\vdash \delta(q_1, aab, aTbz_0)$

$\vdash \delta(q_1, ab, Tbz_0)$

$\vdash \delta(q_1, ab, aTbz_0)$

$\vdash \delta(q_1, b, Tbz_0)$

$\vdash \delta(q_1, b, \varepsilon bz_0)$

$\vdash \delta(q_1, b, bz_0)$

$\vdash \delta(q_1, \varepsilon, z_0)$

$\vdash \delta(q_f, \varepsilon) \rightarrow$ acceptance

PDAs accepting by final state and empty stack are equivalent:

$P_F \rightarrow$ PDA accepting by final state,

$P_F = (Q_F, \Sigma, \Gamma, \delta_F, q_0, z_0, F)$

$P_N \rightarrow$ PDA accepting by empty stack

$P_N = (Q_N, \Sigma, \Gamma, \delta_N, q_0, z_0)$

- For every $P_N, \exists P_F \exists L(P_F) = L(P_N)$
- For every $P_F, \exists P_N \exists L(P_N) = L(P_F)$

CONVERTING CFG TO PDA

The PDA simulates the left most derivation on a given w , and upon consuming it fully it either arrives at acceptance (by empty stack) or non-acceptance.

The steps to convert CFG to PDA are:

1. Push right hand side of the production on to stack, with left most symbol at the stack top.
2. If stack top is the left most variable, then replace it by all its productions (each possible substitution will represent a distinct path taken by non-deterministic PDA (NPDA)).
3. If stack top has a terminal symbol and if it matches with the next symbol in the input string, then pop it. Follow from step-1 again to complete all productions.

Example 11: The CFG, G of a language L is $S \rightarrow AB, A \rightarrow aAb \mid \varepsilon, B \rightarrow cB \mid \varepsilon$. The PDA generated by G contains _____ states.

- (A) 5 (B) 4 (C) 3 (D) 1

Solution: (C)

$$\begin{aligned}
 S &\rightarrow AB \\
 A &\rightarrow aAb|\epsilon \\
 B &\rightarrow cB|\epsilon \\
 \Rightarrow \delta(q_0, w, S) &= (q_1, AB) \\
 \delta(q_1, w, A) &= (q_1, aAb) \\
 \delta(q_1, \epsilon, A) &= \delta(q_1, \epsilon) \\
 \delta(q_1, w, B) &= (q_1, cB) \\
 \delta(q_1, \epsilon, B) &= \delta(q_2, \epsilon) \rightarrow \text{accept} \\
 \therefore \{q_0, q_1, q_2\} &\text{ 3-states.}
 \end{aligned}$$

Converting a PDA into a CFG

Given: $G = (V, T, P, S)$ Initial stack symbol (S) same as start variable in grammar

Output: $P_N = (\{q\}, T, V \cup T, \delta, q, S)$, where δ is

- If q_0 is start state in PDA and q_n is final state of PDA then $[q_0, z, q_n]$ becomes a start state of CFG. Here z represents stack symbol.
- The production rule for the ID of the form $\delta(q_i, a, z_0) = (q_{i+1}, z_1, z_2)$ can be obtained as:

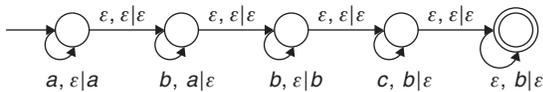
$$\delta(q_i, z_0, q_{i+k}) \rightarrow a(q_{i+1}, z_1, q_m) (q_m, z_2, q_{i+k})$$

Where q_{i+k}, q_m represents the intermediate states, z_0, z_1, z_2 are stack symbols and a is input symbol.

- The production rule for the ID of the form $\delta(q_i, a, z_0) = (q_{i+1}, \epsilon)$ can be converted as

$$(q_i, z_0, q_{i+1}) \rightarrow a$$

Example 12: The PDA, P for language L is generated as:



The CFG for P is:

- (A) $S \rightarrow S_1 S_2$
 $S_1 \rightarrow a S_2 b$
 $S_2 \rightarrow c | \epsilon$
- (B) $S \rightarrow S_1 b c$
 $S_1 \rightarrow a | \epsilon$
- (C) $S \rightarrow S_1 S_2$
 $S_1 \rightarrow a S_1 b | \epsilon$
 $S_2 \rightarrow b S_2 | b S_2 c | \epsilon$
- (D) $S \rightarrow a S_1 c$
 $S_1 \rightarrow b | \epsilon$

Solution: (C)

The language, L generated by given PDA is

$$L = \{a^n b^n b^m c^p : m \geq p \text{ and } n, p \geq 0\}$$

It can be generated by following rules:

$$S \rightarrow S_1 S_2$$

- $S_1 \rightarrow a S_1 b | \epsilon \rightarrow S_1$ generates $a^n b^n$
- $S_2 \rightarrow b S_2 | b S_2 c | \epsilon \rightarrow S_2$ generates $b^m c^p$

DETERMINISTIC PDA (DETERMINISTIC CFL)

$$\left\{ \begin{array}{l} \text{Deterministic} \\ \text{context free} \\ \text{languages} \\ \text{(DPDA)} \end{array} \right\} \subseteq \left\{ \begin{array}{l} \text{Context-free} \\ \text{Languages} \\ \text{PDAs} \end{array} \right\}$$

- Every DPDA is also a PDA.
- A context free language ' L ' accepted by PDA may or may not be accepted by DPDA.

A PDA, $M = (Q, \Sigma, \Gamma, \delta, q_0, F)$ is deterministic if there is no configuration for which M has choice of more than one move. That is, it must satisfy the following conditions:

1. For any $q \in Q, a \in \Sigma \epsilon$ and $s \in \Gamma \epsilon$, the set $\delta(q, a, s)$ has almost one element. (Doesn't allow two or more transitions from same state).
2. For any $q \in Q$, and $s \in \Gamma \epsilon$, if $\delta(q, \epsilon, s) \neq \phi$, then $\delta(q, a, s) = \phi$ for every $a \in \Sigma$ and $\delta(q, a, \epsilon) = \phi$ for all $a \in \Sigma \epsilon$.
3. For any $q \in Q$ and $a \in \Sigma$, if $\delta(q, a, \epsilon) \neq \phi$, then $\delta(q, a, s) = \phi$ for all $s \in \Gamma$ and $\delta(q, \epsilon, t) = \phi$ for all $t \in \Gamma \epsilon$.
4. For any $q \in Q$, if $\delta(q, \epsilon, \epsilon) \neq \phi$, then $\delta(q, a, t) = \phi$ for all $a \in \Sigma \epsilon$ and $t \in \Gamma \epsilon$ (except when $a = \epsilon, t = \epsilon$).

Rule-2 says that if there is a transition from state q that reads character, s from stack but doesn't read other input, other transitions from q , that don't read stack are not allowed and other transitions from q that read s from the stack and read the input are not allowed either.

Rule-3 says that if there is a transition from state q that reads character a , but doesn't read stack, other transitions from q that don't read the input are not allowed and other transitions from q that read ' a ' from input and read the stack are not allowed either.

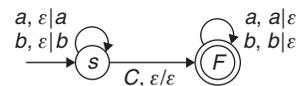
Rule-4 says that if there is a transition from q that doesn't read either input or stack, all other transitions from q are not allowed.

Example 13: A language, L is defined as: $L = \{w c w^R : w \in (a, b)^*\}$. What is Nature of language L ?

- (A) CFL and DCFL
- (B) Only CFL
- (C) Only DCFL
- (D) None of these

Solution: (A)

$$L = \{x = w c w^R \text{ for } w \in (a, b)^*\}$$



Clearly, obtained PDA is also DPDA in sense; there is no choice in transitions.

∴ Hence L is CFL and DCFL.

EXERCISES

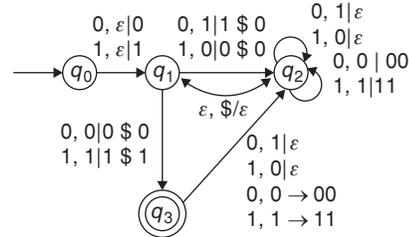
Practice Problems I

Directions for questions 1 to 15: Select the correct alternative from the given choices.

- Consider the grammar, $G = (V, \Sigma, R, S)$ where $V = \{a, b, S, A\}$, $\Sigma = \{a, b\}$, $R = \{S \rightarrow AA, A \rightarrow AAA, A \rightarrow a, A \rightarrow bA, A \rightarrow Ab\}$. How many strings can be generated by $L(G)$ that can be produced by derivations of four or fewer steps?
(A) 5 (B) 10 (C) 14 (D) 8
- Consider the following languages L_1, L_2 and L_3 :
 $L_1 = \{a^n b^m c^{n+m} \mid n, m \geq 0\}$
 $L_2 = \{a^n b^{n+1} c^{n+2} \mid n \geq 0\}$
 $L_3 = \{a^n b^n c^m \mid n, m \geq 0\}$
 Which of following statement is true?
 (A) L_1, L_2, L_3 are context free languages
 (B) L_1, L_2 are context free but not L_3
 (C) L_1, L_3 are context free but not L_2
 (D) L_1, L_2, L_3 are not context free languages.
- The language, $L = \{b_i \# b_{i+1} : b_i \text{ is } i \text{ in binary, } i \geq 1\}$ is:
 (A) Regular
 (B) Context free
 (C) Regular and context free
 (D) Neither context free nor Regular
- The CFG, $G : A \rightarrow BAB|B|E, B \rightarrow 00|E$. The CFG is normalized using CNF. The obtained G' , contains ___ rules.
 (A) 11 (B) 14 (C) 12 (D) 13
- The language $L = \{0^{2^i} : i \geq 1\}$ is:
 (A) Context free
 (B) DCFL
 (C) Both CFL and DCFL
 (D) Not context free language
- The context free grammar, G is defined with production rules $S \rightarrow EcC'|aAE|AU, A \rightarrow aA|E, B \rightarrow bB|E, C' \rightarrow cC'|E, E \rightarrow aEc|F, F \rightarrow bFc|E, U \rightarrow aUc|V, V \rightarrow bVc|bB$. What is the language generated by L ?
 (A) $L = \{a^n b^m c^k : k \neq n + m\}$
 (B) $L = \{a^n b^m c^k : k = n + m\}$
 (C) $L = \{a^n b^m c^k : k > n + m\}$
 (D) $L = \{a^n b^m c^k : k < n + m\}$
- Consider the grammar, $G \equiv S \rightarrow abScB|E, B \rightarrow bB|b$. What language does it generate?
 (A) $L(G) = \{(ab)^n (cb)^m \mid n = m\}$
 (B) $L(G) = \{a^n b^n (cb)^m \mid n \neq m\}$
 (C) $L(G) = \{(ab)^n (cb^m)^n \mid n \geq 0, m > 0\}$
 (D) $L(G) = \{(ab)^n (cb^m)^n \mid n \geq 0, m \geq 0\}$
- The language, $L = \{0^i 1^j 2^k \mid i \neq j \text{ or } j \neq k\}$. The CFG, G generated by L contains ___ rules.
 (A) 23 (B) 20
 (C) 21 (D) 19

- The DPDA constructed to accept language, L with property $L = L_1 \cup L_2$ where $L_1 = \{10^n 1^n \mid n > 0\}$, $L_2 = \{110^n 1^{2n} \mid n > 0\}$ contains ___ states.
 (A) 4 (B) 5
 (C) 6 (D) 7

- The PDA is designed as:



What is the language generated by the above PDA?

- Binary strings that have same number of 0's and 1's.
 - Binary strings that start with 00 and end with 11 and have same number of 0's and 1's.
 - Binary strings that start and end with the same symbol and have same number of 0's and 1's.
 - Binary strings that start with 11 and end with 00 and have same number of 0's and 1's.
- The language, $L = \{ba^{m_1} ba^{m_2} b \dots ba^{m_n} : n \geq 2, m_1, \dots, m_n \geq 0 \text{ and } m_i \neq m_j \text{ for some } i, j\}$. What is nature of ' L '?
 (A) Regular
 (B) Context free but not regular
 (C) Regular but not context free
 (D) Neither context free nor regular
 - Two languages L_1, L_2 are defined as:
 $L_1 = \{a^i b^j c^k : i, j, k \geq 0, i = j\}$
 $L_2 = \{a^i b^j c^k : i, j, k \geq 0, j = k\}$ which of following statements are true?
 (i) $L_1 \cap L_2$ is context free
 (ii) $L_1 \cap L_2 = \{a^n b^n c^n \mid n \geq 0\}$
 (iii) L_1, L_2 are context free
 (iv) Only L_1 is context free
 (A) All are true (B) (i), (ii) are true
 (C) (iii), (iv) are true (D) (ii), (iii) are true
 - The language generated by grammar:
 $S \rightarrow Te|Ue, T \rightarrow cTd|cT|E, U \rightarrow cUd|Ud|dd$. is
 (A) $L = \{c^n d^m e : m \geq n\}$
 (B) $L = \{c^n d^m e : m = n\}$
 (C) $L = \{c^m d^n e^m : m \geq n + 2\}$
 (D) None of these
 - Remove null productions, useless symbols from the following grammar result in:
 $S \rightarrow ABC$
 $A \rightarrow aBC$
 $B \rightarrow C|E$
 $C \rightarrow cd|DCF$
 $D \rightarrow dD|E$

- $E \rightarrow eFE$
 $F \rightarrow eC$
 (A) $S \rightarrow ABC|AC$
 $A \rightarrow aBC|aC$
 $B \rightarrow C$
 $C \rightarrow cd|DCF|CF$
 $D \rightarrow dD|d$
 $F \rightarrow eC$
 (B) $S \rightarrow aBCc$
 $A \rightarrow aBC$
 $B \rightarrow cD|dDEF|dEF$
 $C \rightarrow cD|dDEF|dEF$
 $F \rightarrow eB$
 $D \rightarrow dD|d$
 $E \rightarrow eFE|e$

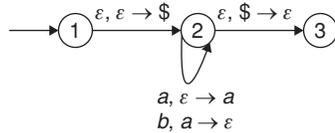
- (C) $S \rightarrow aBCBC|aBC$
 $B \rightarrow cD|dDEF|dEF$
 $F \rightarrow eB$
 $C \rightarrow dD|d$
 $D \rightarrow e$
 $F \rightarrow CD|dDEF$
 (D) None of these

15. Let the language L_1, L_2 are defined as:
 $L_1 = \{a^i b^{2i} c^j \mid i, j \geq 0\}$, $L_2 = \{a^i b^{2i} a^j \mid i \geq 0\}$. Which of following is true?
 (A) L_1, L_2 are context free
 (B) Only L_1 is context free
 (C) Only L_2 is context free
 (D) Neither L_1 nor L_2 is context free

Practice Problems 2

Directions for questions 1 to 15: Select the correct alternative from the given choices.

- Consider the alphabet $\Sigma = \{a, b, c, (,), \cup, *, \phi\}$. Then context free grammar that generates all strings in Σ^* that are regular expressions over $\{a, b\}$ is:
 (A) $S \rightarrow S^*|a|b|SS$
 (B) $S \rightarrow \phi|a|b|S$
 (C) $S \rightarrow \phi|a \cup b|S^*$
 (D) $S \rightarrow \phi|S^*|a|b|(S)|S \cup S|SS$
- The PDA for language, L is designed below. The CFG generated contains ____ productions.



- (A) 5 (B) 4
 (C) 3 (D) 6
- The language, L generated by the following grammar,
 $S \rightarrow SS|AAA|\epsilon, A \rightarrow aA|Aa|b$ is
 (A) $(a^* b^*)^*$ (B) $(a^* b^* b^* a^*)^*$
 (C) $a^* b^* a^*$ (D) $(a^* b a^* b a^* b a^*)^*$
 - The grammar, G is defined with rules $S \rightarrow S_1|S_2, S_1 \rightarrow S_1 b|Ab|\epsilon, A \rightarrow aAb|ab, S_2 \rightarrow S_2 a|Ba|\epsilon, B \rightarrow bBa|ba$. The CNF is applied on G . The obtained grammar, G' contains ____ rules.
 (A) 24 (B) 23
 (C) 21 (D) 20
 - The language, $L = \{b^{n^2} : n \geq 1\}$ is:
 (A) CFL but not DCFL
 (B) DCFL but not CFL
 (C) Only DCFL
 (D) Not CFL

- Consider the grammar, $G = S \rightarrow aSc|B, B \rightarrow bBc|\epsilon$ The language, L generated by G is
 (A) $L = \{a^n b^m c^k : k = n + m\}$
 (B) $L = \{a^n b^m c^k : k \neq n + m\}$
 (C) $L = \{a^n b^m c^k : k > n + m\}$
 (D) $L = \{a^n b^m c^k : k < n + m\}$
- The grammar, G is defined with productions:
 $S \rightarrow 0A|1B, A \rightarrow 0AA|1S|1, B \rightarrow 1BB|0S|0$
 The grammar, G_2 is defined with productions:
 $S \rightarrow AB|aaB, A \rightarrow a|Aa, B \rightarrow b$
 Which grammar is/are ambiguous?
 (A) Only G_1
 (B) Only G_2
 (C) Both G_1 and G_2
 (D) Both G_1 and G_2 are unambiguous
- The language, $L_1 = \{0^n 1^n \mid n > 0\}$ and $L_2 = \{0^n 1^{2n} \mid n > 0\}$. The CFG generated for $L_1 \cup L_2$ is:
 (A) $S \rightarrow 0 A 1|0 A 1 1$
 $A \rightarrow 0|1|\epsilon$
 (B) $S \rightarrow 0 A 1 1$
 $A \rightarrow 0|1|\epsilon$
 (C) $S \rightarrow 0 A 1|0 B 1 1$
 $A \rightarrow 0 A 1|\epsilon$
 $B \rightarrow 0 B 1 1|\epsilon$
 (D) $S \rightarrow 0 A 1 1|0 1 1$
 $A \rightarrow 0|1|\epsilon$
- The NPDA constructed to accept language, L with property, $L = L_1 \cup L_2$, where $L_1 = \{1^n 0^n \mid n > 0\}$, $L_2 = \{0^n 1^{2n} \mid n \geq 0\}$ contains ____ final states.
 (A) 3 (B) 1
 (C) 2 (D) 4
- The DPDA for language, L is designed below. What is the language generated?

- (A) $L = \{a^n b^m : m = n\}$
 (B) $L = \{a^n b^m : m = n + 2\}$
 (C) $L = \{a^n b^m : m \geq n + 2\}$
 (D) $L = \{a^n b^m : m \leq n + 2\}$
11. The CFG, G is defined with rules:
 $S \rightarrow AB|CD, A \rightarrow A00|\epsilon, B \rightarrow B11|1, C \rightarrow C00|0, D \rightarrow D11|\epsilon$. The language generated by G is
 (A) $L = \{0^n 1^n | n \geq 0\}$
 (B) $L = \{0 0^n 1 1^n | n > 0\}$
 (C) $L = \{0^n 1^m | n + m \text{ is odd}\}$
 (D) $L = \{0^n 1^m | n + m \text{ is even}\}$
12. The languages, L_1, L_2, L_3 are defined as:
 $L_1 = \{a^n b^m c^{n+m} | n, m \geq 0\}, L_2 = \{a^n b^n c^m | n, m \geq 0\}, L_3 = \{a^n b^n c^{2n} | n \geq 0\}$. Which of the following statements are true?
 (i) L_1, L_2 are context free
 (ii) L_1, L_3 are context free
 (iii) $L_3 = L_1 \cap L_2$
 (iv) L_1, L_3 are context free but not L_2
 (A) (i), (ii) (B) (i), (iii)
 (C) (ii), (iii) (D) (iii), (iv)

13. The language, L_1 and L_2 are defined as $L_1 = \{a^n b^n : n \geq 0 \text{ and } n \text{ is not a multiple of } 5\}$ and $L_2 = \{0^n \# 0^{2n} \# 0^{3n} | n \geq 0\}$. Which of following is true?
 (A) L_1 and L_2 are context free
 (B) Only L_1 is context free
 (C) Only L_2 is context free
 (D) Neither L_1 nor L_2 is context free
14. The language, L_1 and L_2 are defined as $\overline{L_1} = \{0^n 1^n\}^m | m, n > 0\}$, $L_2 = \{0^n 1^n 0^n 1^n | n \geq 0\}$ which of following is true?
 (A) L_1 and L_2 are context free
 (B) Only $\overline{L_1}$ is context free
 (C) Only L_2 is context free
 (D) Neither L_1 nor L_2 is context free
15. The language L_1, L_2 are defined as $L_1 = \{0^i 1^j 0^i | i, j > 0\}$, $L_2 = \{1^k 0^i 1^j 0^k | i, j, k > 0\}$. Which of following is true?
 (A) L_1 and L_2 are context free
 (B) Only L_1 is context free
 (C) Only L_2 is context free
 (D) Neither L_1 nor L_2 is context free

PREVIOUS YEARS' QUESTIONS

1. Match the following: [2008]

- | | |
|---|---|
| E. Checking that identifiers are declared before their use | P. $L = \{a^n b^m c^n d^m n \geq 1, m \geq 1\}$ |
| F. Number of formal parameters in the declaration of a function agrees with the number of actual parameters in use of that function | Q. $X \rightarrow XbX XcX dXf g$ |
| G. Arithmetic expressions with matched pairs of parentheses | R. $L = \{wcn w \in (a b)^*\}$ |
| H. Palindromes | S. $X \rightarrow bXb cXc \epsilon$ |

- (A) E – P, F – R, G – Q, H – S
 (B) E – R, F – P, G – S, H – Q
 (C) E – R, F – P, G – Q, H – S
 (D) E – P, F – R, G – S, H – Q

2. Consider the languages L_1, L_2 and L_3 as given below.
 $L_1 = \{0^p 1^q | p, q \in N\}$,
 $L_2 = \{0^p 1^q | p, q \in N \text{ and } p = q\}$ and
 $L_3 = \{0^p 1^q 0^r | p, q, r \in N \text{ and } p = q = r\}$. Which of the following statements is NOT TRUE? [2011]
 (A) Push Down Automata (PDA) can be used to recognize L_1 and L_2 .
 (B) L_1 is a regular language.
 (C) All the three languages are context free
 (D) Turing machines can be used to recognize all the languages.

3. Which of the following problems are decidable? [2012]

- (1) Does a given program ever produce an output?
 (2) If L is a context free language, then, is \overline{L} also context free?
 (3) If L is a regular language, then, is \overline{L} also regular?
 (4) If L is recursive language, then, is \overline{L} also recursive?
 (A) 1, 2, 3, 4 (B) 1, 2
 (C) 2, 3, 4 (D) 3, 4

4. Consider the following languages.
 $L_1 = \{0^p 1^q 0^r | p, q, r \geq 0\}$
 $L_2 = \{0^p 1^q 0^r | p, q, r \geq 0, p \neq r\}$
 Which one of the following statements is FALSE? [2013]

- (A) L_2 is context-free
 (B) $L_1 \cap L_2$ is context-free
 (C) Complement of L_2 is recursive
 (D) Complement of L_1 is context-free but not regular

5. Which one of the following is TRUE? [2014]

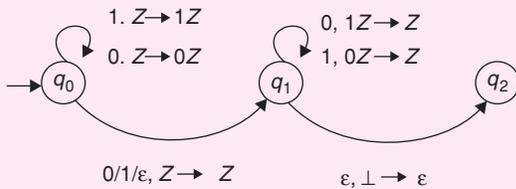
- (A) The language $L = \{a^n b^n | n \geq 0\}$ is regular
 (B) The language $L = \{a^n | n \text{ is prime}\}$ is regular
 (C) The language $L = \{w | w \text{ has } 3k + 1b\text{'s for some } k \in N \text{ with } \Sigma = \{a, b\}\}$ is regular
 (D) The language $L = \{ww^r | w \in \Sigma^*\}$ with $\Sigma = \{0, 1\}$ is regular.

6. Consider the following languages over the alphabet $\Sigma = \{0, 1, c\}$.
 $L_1 = \{0^n 1^n | n \geq 0\}$
 $L_2 = \{wcn^2 | w \in \{0, 1\}^*\}$
 $L_3 = \{ww^r | w \in \{0, 1\}^*\}$

Here w^r is reverse of the string w . Which of these languages are deterministic context-free languages? [2014]

- (A) None of the languages
- (B) Only L_1
- (C) Only L_1 and L_2
- (D) All the three languages

7. Consider the NPDA $\langle Q = \{q_0, q_1, q_2\}, \Sigma = \{0, 1\}, \Gamma = \{0, 1, \perp\}, \delta, q_0, \perp, F = \{q_2\} \rangle$, where (as per usual convention) Q is the set of states, Σ is the input alphabet, Γ is the stack alphabet, δ is the state transition function, q_0 is the initial state, \perp is the initial stack symbol, and F is the set of accepting states. The state transition is as follows:



Which one of the following sequences must follow the string 1011 00 so that the overall string is accepted by the automation? [2015]

- (A) 10110
- (B) 10010
- (C) 01010
- (D) 01001

8. Which of the following languages are context-free? [2015]

- $L_1 = \{a^m b^n a^n b^m \mid m, n \geq 1\}$
- $L_2 = \{a^m b^n a^m b^n \mid m, n \geq 1\}$
- $L_3 = \{a^m b^n \mid m = 2n + 1\}$
- (A) L_1 and L_2 only
- (B) L_1 and L_3 only
- (C) L_2 and L_3 only
- (D) L_3 only

9. Consider the following context-free grammars:

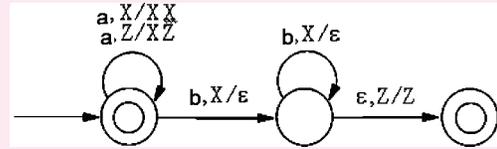
$$G_1: S \rightarrow aS|B, B \rightarrow b|bB$$

$$G_2: S \rightarrow aA|bB, A \rightarrow aA|B| \epsilon, B|bB\epsilon$$

Which one of the following pairs of languages is generated by G_1 and G_2 , respectively? [2016]

- (A) $\{a^m b^n \mid m > 0 \text{ or } n > 0\}$ and $\{a^m b^n \mid m > 0 \text{ and } n > 0\}$
- (B) $\{a^m b^n \mid m > 0 \text{ and } n > 0\}$ and $\{a^m b^n \mid m > 0 \text{ or } n \geq 0\}$
- (C) $\{a^m b^n \mid m \geq 0 \text{ or } n > 0\}$ and $\{a^m b^n \mid m > 0 \text{ and } n > 0\}$
- (D) $\{a^m b^n \mid m \geq 0 \text{ and } n > 0\}$ and $\{a^m b^n \mid m > 0 \text{ or } n > 0\}$

10. Consider the transition diagram of a PDA given below with input alphabet $\Sigma = \{a, b\}$ and stack alphabet = $\{X, Z\}$. Z is the initial stack symbol. Let L denote the language accepted by the PDA.



Which one of the following is TRUE? [2016]

- (A) $L = \{a^n b^n \mid n \geq 0\}$ and is not accepted by any finite automata.
- (B) $L = \{a^n \mid n \geq 0\} \cup \{a^n b^n \mid n \geq 0\}$ and is not accepted by any deterministic PDA.
- (C) L is not accepted by any Turing machine that halts on every input.
- (D) $L = \{a^n \mid n \geq 0\} \cup \{a^n b^n \mid n \geq 0\}$ and is deterministic context-free.

11. Consider the following languages:

$$L_1 = \{a^m b^m c^{m+n} \mid m, n \geq 1\}$$

$$L_2 = \{a^n b^n c^{2n} \mid n \geq 1\}$$

Which one of the following is TRUE? [2016]

- (A) Both L_1 and L_2 are context-free.
- (B) L_1 is context-free while L_2 is not context-free.
- (C) L_2 is context-free while L_1 is not context-free.
- (D) Neither L_1 nor L_2 is context-free.

12. Consider the following context-free grammar over the alphabet $\Sigma = \{a, b, c\}$ with S as the start symbol:

$$S \rightarrow abScT \mid abcT$$

$$T \rightarrow bT \mid b$$

Which one of the following represents the language generated by the above grammar? [2017]

- (A) $\{(ab)^n (cb)^n \mid n \geq 1\}$
- (B) $\{(ab)^n cb^{m_1} cb^{m_2} \dots cb^{m_n} \mid n, m_1, m_2, \dots, m_n \geq 1\}$
- (C) $\{(ab)^n (cb^m)^n \mid m, n \geq 1\}$
- (D) $\{(ab)^n (cb^n)^m \mid m, n \geq 1\}$

13. If G is a grammar with productions

$$S \rightarrow SaS \mid aSb \mid bSa \mid SS \mid \epsilon$$

Where S is the start variable, then which one of the following strings is not generated by G ? [2017]

- (A) $abab$
- (B) $aaab$
- (C) $abbaa$
- (D) $babba$

14. Consider the context-free grammars over the alphabet $\{a, b, c\}$ given below. S and T are non-terminals.

$$G_1: S \rightarrow aSb|T, T \rightarrow cT| \epsilon$$

$$G_2: S \rightarrow bSa|T, T \rightarrow cT| \epsilon$$

The language $L(G_1) \cap L(G_2)$ is [2017]

- (A) Finite
- (B) Not finite but regular
- (C) Context-Free but not regular
- (D) Recursive but not context-free.

15. Consider the following languages over the alphabet $\Sigma = \{a, b, c\}$.

Let $L_1 = \{a^n b^n c^m \mid m, n \geq 0\}$ and $L_2 = \{a^m b^n c^n \mid m, n \geq 0\}$.

Which of the following are context-free languages? [2017]

- I. $L_1 \cup L_2$
- II. $L_1 \cap L_2$
- (A) I only (B) II only
- (C) I and II (D) Neither I nor II

16. Let L_1, L_2 be any two context-free languages and R be any regular language. Then which of the following is/are CORRECT? [2017]

- I. $L_1 \cup L_2$ is context-free.
- II. \bar{L}_1 is context-free.
- III. $L_1 - R$ is context-free.
- IV. $L_1 \cap L_2$ is context-free.
- (A) I, II and IV only (B) I and III only
- (C) II and IV only (D) I only

17. Identify the language generated by the following grammar, where S is the start variable. [2017]

$$\begin{aligned} S &\rightarrow XY \\ X &\rightarrow aX|a \\ Y &\rightarrow aYb|\epsilon \end{aligned}$$

- (A) $\{a^m b^n \mid m \geq n, n > 0\}$
- (B) $\{a^m b^n \mid m \geq n, n \geq 0\}$
- (C) $\{a^m b^n \mid m > n, n \geq 0\}$
- (D) $\{a^m b^n \mid m > n, n > 0\}$

18. Consider the following languages.

$$\begin{aligned} L_1 &= \{a^p \mid p \text{ is a prime number}\} \\ L_2 &= \{a^n b^m c^{2m} \mid n \geq 0, m \geq 0\} \\ L_3 &= \{a^n b^n c^{2n} \mid n \geq 0\} \\ L_4 &= \{a^n b^n \mid n \geq 1\} \end{aligned}$$

Which of the following are CORRECT? [2017]

- I. L_1 is context-free but not regular.
- II. L_2 is not context-free.
- III. L_3 is not context-free but recursive.
- IV. L_4 is deterministic context-free.
- (A) I, II and IV only (B) II and III only
- (C) I and IV only (D) III and IV only

19. Consider the following languages:

- I. $\{a^m b^n c^p d^q \mid m + p = n + q, \text{ where } m, n, p, q \geq 0\}$
- II. $\{a^m b^n c^p d^q \mid m = n \text{ and } p = q, \text{ where } m, n, p, q \geq 0\}$
- III. $\{a^m b^n c^p d^q \mid m = n = p \text{ and } p \neq q, \text{ where } m, n, p, q \geq 0\}$
- IV. $\{a^m b^n c^p d^q \mid mn = p + q, \text{ where } m, n, p, q \geq 0\}$

Which of the languages above are context-free?

[2018]

- (A) I and IV only (B) I and II only
- (C) II and III only (D) II and IV only

ANSWER KEYS

EXERCISES

Practice Problems 1

- 1. D
- 2. C
- 3. D
- 4. B
- 5. D
- 6. A
- 7. C
- 8. B
- 9. D
- 10. C
- 11. B
- 12. D
- 13. D
- 14. A
- 15. B

Practice Problems 2

- 1. D
- 2. C
- 3. D
- 4. A
- 5. D
- 6. A
- 7. C
- 8. C
- 9. C
- 10. C
- 11. C
- 12. B
- 13. B
- 14. B
- 15. C

Previous Years' Questions

- 1. C
- 2. C
- 3. D
- 4. D
- 5. C
- 6. C
- 7. B
- 8. B
- 9. D
- 10. D
- 11. B
- 12. B
- 13. D
- 14. B
- 15. A
- 16. B
- 17. C
- 18. D
- 19. B

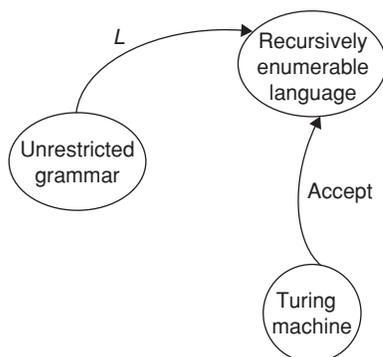
Chapter 3

Recursively Enumerable Sets and Turing Machines, Decidability

LEARNING OBJECTIVES

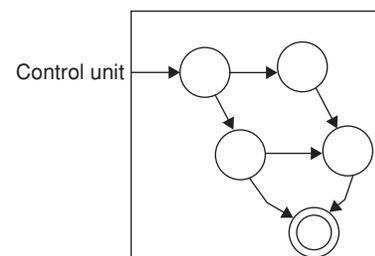
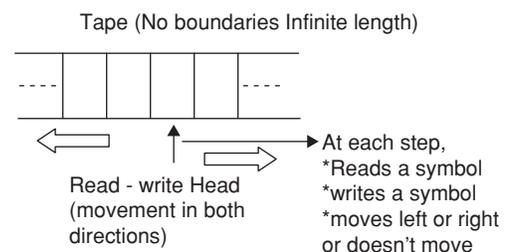
- ☞ Turing machines
- ☞ Model of turing machine
- ☞ Types of turing machines
- ☞ Offline turing machine
- ☞ Universal turing machine
- ☞ Recursively enumerable languages
- ☞ Recursive language
- ☞ Undecidability
- ☞ Church's hypothesis
- ☞ Halting problem
- ☞ Post's correspondence problem
- ☞ 7 P problems
- ☞ NP problems
- ☞ NP – complete problem
- ☞ NP – hard problem
- ☞ Closure properties of formal languages

TURING MACHINES



A Turing machine is a kind of state machine, which is much more powerful in terms of languages it can recognize. At any time, the machine is in any one of the finite number of states. Instructions for a turing machine include the specification of conditions, under which the machine will make transitions from one state to other.

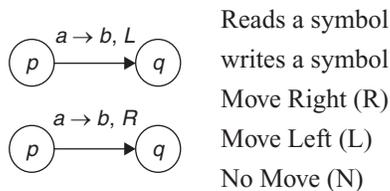
Model of Turing Machine



- A TM (turing machine) consists of Tape, Head, control unit.
- **Tape:** A tape is divided into a sequence of numbered cells, Each cell contains a symbol and cells that have not been written before are assumed to be filled with a blank symbol (B). The set of symbols of tape is denoted by Γ . The tape is assumed to be arbitrarily extensible to the left as well as to the right.
- **Head:** In a single step, a tape head reads the contents of a cell on the tape (reads a symbol), replaces it with some other characters (writes a symbol) and repositions itself to the next cell to the right or to the left of the one it has just read or does not move (moves left or right or does not move).
- **Control unit:** The reading from the tape or writing into the tape is determined by the control unit. It contains a finite set of states, Q . The states are:
 1. Initial state, q_0
 2. Halt state, h : This is state in which TM stops all further operations. There can be one or more halt states in a TM.
 3. Other states.

Note: A TM on entering the halt state stops making moves and whatever string is there on the tape, will be taken as the output, irrespective of whether the position of head is at the end or in the middle of the string on the tape.

Transition Diagram of TM



Specification of TM

5-Tuple specification:

TM = (state1, Read symbol, write symbol, L/R/N, state 2).

7 - Tuple specification of TM:

A TM, M is represented as a 7-tuple:

$M = (Q, \Sigma, \Gamma, \delta, q_0, B, h)$ where

$Q \rightarrow$ Finite set of states

$\Sigma \rightarrow$ Finite set of non-blank symbols

$\Gamma \rightarrow$ Set of tape characters

$q_0 \rightarrow q_0 \in Q$, initial state

$B \rightarrow$ Blank character

$h \rightarrow h \in Q$, final state

$\delta \rightarrow$ Transition function, $Q \times \Gamma \rightarrow Q \times \Gamma \times \{L, R, N\}$

String classes in TM

Every TM, over the alphabet Σ , divides set of input string w into three classes:

1. **Accept (TM):** It is the set of all strings $w \in \Sigma^* \ni$ if the tape initially contains w and the TM is then run, then TM ends in a halt state.

2. **LOOP (TM):** It is the set of all strings, $w \in \Sigma^* \ni$ if the tape initially contains w and the TM is then run, then the TM loops forever (infinite loop).
3. **Reject (TM):** It is the set of all strings $w \in \Sigma^* \ni$, any of the following 3-cases arise.

Case I: There may be a state and a symbol under the tape head, for which δ does not have a value.

Case II: If the head is reading the left most cell (i) containing the symbol x , the state of TM is say q , then $\delta(q, x)$ suggests a move to the left of the current cell. However as there is no cell to the left, no move is possible.

Case III: If TM enters an infinite loop or if a TM rejects a given string w , because of above two cases, TM crashes (terminates unsuccessfully).

LANGUAGES ACCEPTED BY A TM

- The language accepted by TM is the set of accepted strings $w \in \Sigma^*$.

- Formally, let $M = (Q, \Sigma, \Gamma, \delta, q_0, B, h)$ be a TM. The language accepted by M denoted by $L(M)$ is defined as, $L(M) = \{w/w \in \Sigma^* \text{ and if } w = a_1 \dots a_n \text{ then, } (q_0, \epsilon, a_1, a_2, \dots, a_n) (h, b, \dots, b_{i-1}, b_j, b_n) \text{ for some } b_1, b_2 \dots b_n \in N^* \ni\}$

$$L(M) = \{W: q_0 w \vdash^* x_1 h x_2\}$$

- There are three types of turing machine related languages:

1. **Turing Acceptable language:** A language, L over some alphabet is said to be turing acceptable language if there exists a TM, $M \ni L = L(M)$
2. **Turing Decidable Language:** A language L over Σ i.e., $L \subseteq \Sigma^*$ is said to be turing decidable, if both languages, L and its complement $\Sigma^* - L$ are turing acceptable.
3. **Recursively Enumerable Language:** A language L is recursively enumerable, if it is accepted by a TM.

Example 1: Let M be a turing machine has $M = (Q, \Gamma, \Sigma, \delta, S, B, F)$ with $Q = \{q_0, q_1, q_2, q_3, q_4\}$, $\Gamma = \{a, b, X, Y, \#\}$, $\Sigma = \{a, b\}$, $S = q_0$, $B = \#$, δ given by:

	a	b	X	Y	#
q_0	(q_1, X, R)	-	-	(q_3, Y, R)	
q_1	(q_1, a, R)	(q_2, Y, L)	-	(q_1, Y, R)	
q_2	(q_2, a, L)	-	(q_0, X, R)	(q_2, Y, L)	
q_3	-	-	-	(q_3, Y, R)	$(q_4, \#, R)$
q_4	-	-	-	-	-

Which of following is true about M ?

- (A) M halts on L having 'baa' as substring
- (B) M halts on L having 'bab' as substring
- (C) M halts on $L = \{a^n b^n / n \geq 1\}$
- (D) M halts on L not having 'bbaa' as substring.

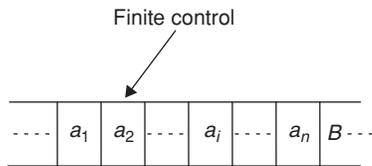
Solution: (C)
M accepts $a^n b^n$.

Example: $aaabbb$

$(q_0, \epsilon, aaabbb)$	\rightarrow	$(q_1, XXXYY, b)$
$(q_1, X, aaabbb)$	\rightarrow	$(q_2, XXXY, YY)$
$(q_1, Xa, abbb)$	\rightarrow	${}^1(q_2, XXX, YYY)$
(q_1, Xaa, bbb)	\rightarrow	$(q_2, XY, XYYY)$
$(q_1, Xa, aYbb)$	\rightarrow	(q_0, XXX, YYY)
$(q_2, X, aaYbb)$	\rightarrow	$(q_3, XXXY, YY)$
$(q_2, \epsilon, XaaYbb)$	\rightarrow	$(q_3, XXXYY, Y)$
$(q_0, X, aaYbb)$	\rightarrow	$(q_3, XXXYYY, \epsilon)$
$(q_1, XX, aYbb)$	\rightarrow	$(q_4, XXXYYY\#, \epsilon)$

TYPES OF TURING MACHINES

Two-way infinite turing machine

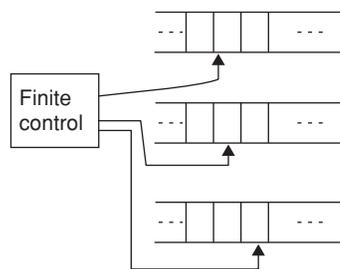


- A TM with a two-way infinite tape is denoted by $M = (Q, \Sigma, \Gamma, \delta, q_0, B, F)$, as in original model.
- The tape is infinite to the left as well as to the right.

If $\delta(q, x) = (p, Y, L)$ then $q \ x \ \alpha \ \vdash_m \ pBY$. The tape, is infinite towards left.

If $\delta(q, x) = (p, B, R)$ then $q \ x \ \alpha \ \vdash_m \ p\alpha$ the is infinite towards right.

Multiple turing machines



- A multiple TM consists of a finite control with k tape heads and k -tapes, each tape is infinite in both directions, on a single move, depending on the state of the finite control and the symbol scanned by each of tape heads, the machine can,
 - change state
 - print new symbol on each of the cells scanned by its tape head
 - move each of its tape heads, independently, one cell to the left or right or keep it stationary.
- Initially, the input appears on the first tape and other tapes are blank.

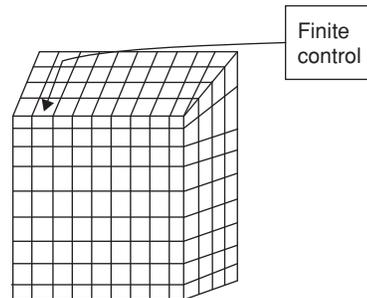
Non-deterministic turing machines

- A non-deterministic turing machine is a device with a finite control and a single one way infinite tape.
- For a given state and a tape symbol scanned by the tape head, the machine has a finite number of choices for next move.

Note: Non-deterministic TM is not permitted to make a move in which the next state is selected from one choice, and the symbol printed and direction of head motion are selected from other choices.

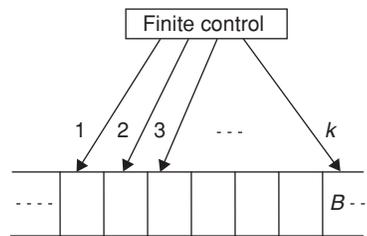
- The non-deterministic TM accepts its input if any sequence of choices of moves leads to an accepting state.

Multi-dimensional TM's



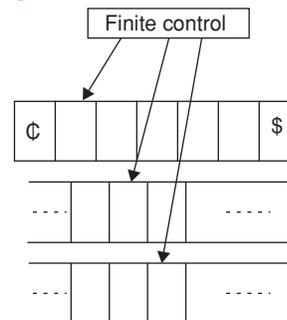
- The tape consists of a k -dimensional array of cells infinite in all $2k$ directions, for some fixed k .
- Depending on the state and the symbol scanned, the device changes its state, prints a new symbol and moves its tape head in one of the $2k$ directions, either positively or negatively, along one of the k -axes.

Multihead TM



- A K -head TM has some fixed ' K ' number of heads. The heads are numbered from 1 through k , and a move of the TM depends on the state and on the symbol scanned by each head.

Offline turing machine



- An offline TM is a multi tape TM, whose input tape is read only. The input is surrounded by end markers, ϵ on left and $\$$ on right. The TM is not allowed to move the input tape head off the region between ϵ and $\$$.

Multi stack machine

- A deterministic two stack machine is a deterministic TM with a read only input and two storage tapes.

Note:

- All these types of TM's does not add any language accepting power and all these are equivalent to the basic model.
- Any language accepted by a 2-PDA can be accepted by some TM and any language accepted by a TM can be accepted by some 2-PDA. Accepting power of a TM = accepting power of a computer.
- Any language accepted by a PDA with n stacks ($n \geq 2$), can also be accepted by some TM.

Example 2: Consider the following statement about L :

1. L is accepted by multi-tape turing machine M_1 .
2. L is also accepted by single tape turing machine M_2 .

Which of following statement is correct?

- (A) Acceptance by M_2 is slower by $O(n^2)$
- (B) Acceptance of M_2 is slower by $O(n)$
- (C) Acceptance of M_2 is faster by $O(n)$
- (D) Acceptance of M_2 is faster by $O(n^2)$

Solution: (A)

While simulating multi-tape TM on a single tape TM the head has to move at least $2k$ cells per move, where k is the number of tracks on single tape TM. Thus for k moves,

$$\sum_{i=1}^k 2i = 2k^2.$$

Which means quadratic slow down?

Thus, acceptance of multi-tape is faster by $O(n^2)$.

Universal turing machine

A Universal turing machine is a turing machine that can simulate an arbitrary turing machine on arbitrary input.

- The machine consists of an input output relation to the machine computes.
- The input is given in binary form on the machine tape and the output consists of the contents of the tape when the machine halts.
- The contents of the tape will change based on the Finite State Machine (FSM) inside the TM.
- The problem with TM is that a different machine will be constructed for every new computation to be performed.
- A UTM can simulate any other machine.

Combining turing machines

If TM_1 and TM_2 are turing machines, then we can combine these machines and create a Turing machine which will first behave like TM_1 and TM_2 .

To combine two turing machines follow below steps:

1. Change all states in TM_2 , so that they do not conflict with the state names in TM_1 .
 2. Change all halts in TM_1 's transition table to the new name of the start state of TM_2 .
 3. Append TM_2 's transition table to the foot of TM_1 's transition table.
- If TM_1 and TM_2 are combined in this way, we will write it as $TM_1 \rightarrow TM_2$.

So this new machine starts off in the initial state of TM_1 , operates as per TM_1 until TM_1 would halt then it launches TM_2 and operates a TM_2 , until TM_2 would halt.

RECURSIVELY ENUMERABLE LANGUAGES

- A language L over the alphabet Σ is called 'recursively enumerable' if there is a TM, M that accept every word in L and either rejects or loops for every word in language L' , the complement of L .
Accept $(M) = L$
Reject $(M) + \text{Loop } (M) = L'$.
- When TM, M is still running on some input of recursively enumerable languages, it is not decided that M will eventually accept, if let it run for long time or M will run forever (in loop).

Recursive language

- A language is said to be recursive, if there exists a TM which will halt and accept when presented with any input string $w \in \Sigma^*$, only if the string is in the language otherwise will halt and reject the string.
- Thus, for turing decidable language L , there is a TM which halts for a large number of inputs w belonging to L .
- A TM that always halts is known as a decider or a total turing machine and is said to decide the recursive language. The recursive language is also called as recursive set of decidable.
- A language accepted by a TM is said to be recursively enumerable language. The subclass of recursively enumerable sets are said to be recursive sets or recursive language.

Note:

- All recursive languages are also recursively enumerable.
- There may be languages which are recursively enumerable but not recursive.
- Set of all possible words over the alphabet of the recursive language is a recursive set.

- Set of all possible words, over the alphabet of the recursive enumerable language, is a recursively enumerable set.

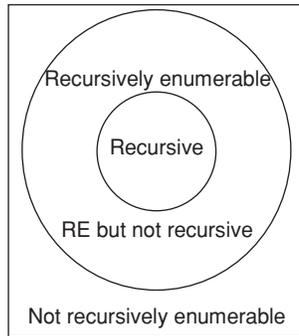


Figure 1 Relationship between the recursive, RE and non-RE languages.

PROPERTIES OF RECURSIVE AND RECURSIVELY ENUMERABLE LANGUAGES

- If a language L is recursive, then there is a TM T that accepts it and always halts.
- If L and L_i are both recursively enumerable, then L and L_i are recursive.
- Union of two recursive languages is recursive.
- Recursively enumerable languages are closed under union.
- If L, L_1 and L_2 are recursive languages, then so are $L_1 \cup L_2, L \cap L_2, L_1 L_2, L^*, L_1 \cap L_2$ and $L_1 - L_2$.
- If L, L_1 and L_2 are recursively enumerable languages, then so are $L_1 \cup L_2, L^*, L_1 \cap L_2, L_1 L_2$.
- If Σ is an alphabet, $L \subseteq \Sigma^*$, is a recursively enumerable language and $\Sigma^* - L$ is recursively enumerable, then L is recursive.

Example 3: If $\Sigma = \{0,1\}$, the canonical order is $\{\epsilon, 0, 1, 00, 01, 10, 11, 000, \dots\}$ where w is the i^{th} word and M_j is TM whose code is the integer j , written in binary. The language generated is $L(M_j)$. The diagonalized language, L_d is a.

- (A) Recursively enumerable language but not recursive
- (B) Recursive language
- (C) Non-recursively enumerable language
- (D) Both (a) and (c)

Solution: (C)
Non-recursively enumerable language.

Non-recursively enumerable language

Non-Recursively Enumerable Language: A language which is not accepted by any Turing machine is non-recursively enumerable.

Example: Power set of an infinite set.

- These languages cannot be defined by any effective procedure.

For any non-empty Σ , there exist languages that are not Recursively Enumerable.

Infinite table for all i and j is:

	$j \rightarrow$			
	1	2	3	4 ...
1	0	1	1	0 ...
2	1	1	0	0 ...
i 3	0	0	1	0 ...
↓ 4	0	1	0	1 ...

↘ Diagonal

To guarantee that no TM accepts L_d :

w_i is in L_d if and only if the (i, i) entry is 0, that is, if M_i does not accept w_i .

Suppose that some TM M_j accepted L_d . Then it contradicts if w_j is in L_d , (j, j) entry is 0, implying that w_j is not in $L(M_j)$ and contradicting $L_d = L(M_j)$.

If w_i is not in L_d , then the (j, j) entry is 1, implying that w_i is in $L(M_j)$, which again contradicts $L_d = L(M_j)$, as w_j is either in or not in L_d , assumption, $L_d = L(M_j)$ is false.

Thus no TM in the list accepts L_d . Hence L_d is non-recursively enumerable language.

Decidable: A problem with two answers (Yes/No) is decidable if the corresponding language is recursive.

Example:

1. $A_{DFA} = \{(M, w) \mid M \text{ accepts the input string } w\}$.

- A Language L is Turing decidable, if there exists a TM M such that on input x , M accepts if $x \in L$ and M rejects otherwise. L is called undecidable if it is not decidable.
- Decidable Languages correspond to algorithmically solvable Decision problems.
- Undecidable language corresponds to algorithmically unsolvable decision problems.

Closure properties of decidable languages

- Decidable Languages are closed under complement, union, intersection, concatenation and star (closure) operations.

Note 1: A language is decidable if both the language and its complement are recognizable.

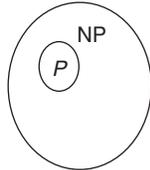
Note 2: Turing Decidable languages are Recursive languages.

UNDECIDABILITY

There are problems that can be computed. There are also problems that cannot be computed. These problems which cannot be computed are called 'computationally undecidable problems'.

NP PROBLEMS

- NP stands for non-deterministic polynomial time.
- A language, L is in class NP, if there is a non-deterministic TM, M is of time complexity $P(n)$ for some polynomial P and M accepts L .
- Class NP consists of problems for which solutions are verified quickly. P consist of problems which can be solved quickly.



- NP languages are closed under union, Intersection, concatenation, Kleen star.
- NP problems are classified into two types:
 1. NP-complete
 2. NP-hard problems.

Example: Vertex (Graph) coloring problem, Travelling salesman problem, the vertex cover problem, the Hamiltonian circuit problem.

NP-COMPLETE PROBLEM

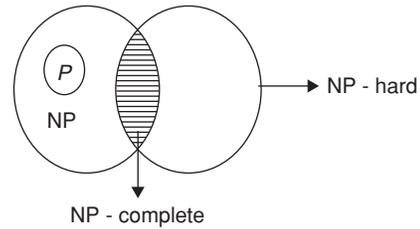
- A class of problems are known as NP-complete problems whose status is unknown. No polynomial time has yet been discovered for NP-complete problems nor has any one been able to prove that no polynomial time exists for any of them. These are hardest of NP-problems. The P and NP-complete problems are disjoint.

Example: (Cook's Theorem) SAT is NP-complete, Bin packing problem, Knapsack Problem.

- A language L is said to be NP-complete if $L \in NP$ and if every $L' \in NP$ is polynomial-time reducible to L .

A language L_1 is said to be polynomial time reducible to some language L_2 if there exists a DTM by which any w_1 in the alphabet of L_1 can be transformed in polynomial time to $a w_2$ in the alphabet of L_2 in such a way that $w_1 \in L_1$ if $w_2 \in L_2$. It follows that if some L_1 is NP-complete and polynomial time reducible to L_2 , then L_2 is also NP-complete.

NP-HARD PROBLEM



- A problem that is NP-hard has a property that all problems that are in NP can be reduced in polynomial time to it.
- A language, L in NP-hard complete if and only if,

Condition 1: For every language, L' in NP, there is a polynomial time reduction of L' to L .

Condition 2: L is not necessarily in NP.

Table 1 NP-Hard versus NP-complete problems:

NP-Hard	NP-Complete
(1) A decision problem P_i is NP-hard if every problem in NP is polynomial time reducible to P_i	(1) A Decision problem P_i is NP-complete if it is NP-hard and is also in class NP itself.
(2) In terms of symbols ' P_i ' is NP-hard if for every $P_j \rightarrow NP$	(2) In terms of symbols, ' P_i ' is NP-complete, if P_i is NP-hard and $P_j \rightarrow NP$
(3) P_i is 'as hard as' all the problem in NP	(3) P_i is one of the hardest problems in NP
(4) If any problem in NP is proved intractable, then P_i must also be intractable	(4) If any one ever shows that as NP-complete problem is also intractable, then every NP-complete problem is also intractable.

Example 5: Which of following is FALSE?

- (A) $\{ \langle x, y \rangle \mid x \text{ and } y \text{ are integers, } \gcd(x, y) = 1 \}$ is a NP class problem.
- (B) CLIQUE is a NP class problem.
- (C) Eulerian PATH is a P class problem
- (D) Dijkstra's algorithm is a problem in P .

Solution: (A)

Choice (A) is a P class problem.

Consider the following table:

D – Decidable, U – Undecidable, ? – Open question, T – Trivially Decidable Question	Regular Sets	DCFL's	CFL's	CSL's	Recursive Sets	Recursively Enumerable Sets
(1) Membership problem?	D	D	D	D	D	D
(2) Emptiness problem?	D	D	D	U	U	U
(3) Completeness problem is $L = \Sigma^*$?	D	D	D	U	U	U
(4) Equality problem?	D	?	U	U	U	U
(5) Subset problem is $L_1 \subseteq L_2$?	D	U	U	U	U	U
(6) Is L Regular?	T	D	U	U	U	U
(7) Is the intersection of two languages, a language, of the same type?	T	U	U	T	T	T
(8) Is the complement of a language, also a language of the same type?	T	T	U	?	T	U
(9) Is L is finite or infinite?	D	D	D	U	U	U

Table 2 Closure properties of formal languages

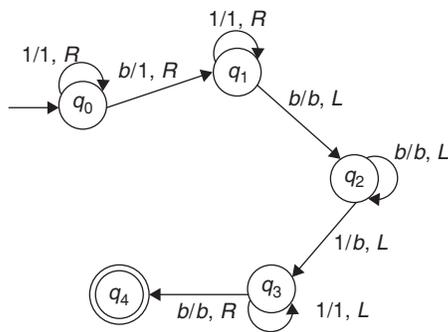
	Regular sets	DCFL'S	CFL'S	CSL'S	Recursive sets	Recursively enumerable sets
(1) Union	Y	N	Y	Y	Y	Y
(2) Concatenation	Y	N	Y	Y	Y	Y
(3) Kleen star	Y	N	Y	Y	Y	Y
(4) Intersection	Y	N	N	Y	Y	Y
(5) Complementation	Y	Y	N	Y	Y	N
(6) Homomorphism	Y	N	Y	N	N	Y
(7) Inverse Homomorphism	Y	Y	Y	Y	Y	Y
(8) Reversal	Y	N	Y	Y	Y	Y
(9) Substitution	Y	N	Y	Y	N	Y
(10) Intersection with regular ets	Y	Y	Y	Y	Y	Y

EXERCISES

Practice Problems I

Directions for questions 1 to 15: Select the correct alternative from the given choices.

1. The TM M over $\Sigma = \{1\}$ is given below



What does M generate?

- (A) The output is total recursive multiplication function.
- (B) The output is addition of two integers.
- (C) The output is subtraction of two integers.
- (D) The output should be w_1w_2 if input = (w_1w_2) a pair of words.

2. Consider language,

$A = \{ \langle M \rangle : M \text{ is a DFA which doesn't accept any string containing odd number 1's} \}$

Which of following is true about A ?

- (A) A is Trivially decidable
- (B) A is undecidable
- (C) A is decidable
- (D) None of these

3. Consider $EQ_{CFG} = \{ \langle G_1G_2 \rangle : G_1, G_2 \text{ are CFGs and } L(G_1) = L(G_2) \}$. Which of following is true about EQ_{CFG} ?

- (A) Recognizable
- (B) Co-Recognizable
- (C) Un-recognizable
- (D) None of the above.

4. A language is given as $INFINITE_{DFA} = \{ \langle A \rangle : A \text{ is a DFA and } L(A) \text{ is an infinite language} \}$. Which of following is true?

- (A) Un-decidable
- (B) Decidable
- (C) Trivially decidable
- (D) None of above.

5. A TM designed over an alphabet $\{0, 1, \#\}$, where 0 indicates blank, which takes a non-null string of 1's and #'s and transfer's the right-most symbol to the left-most end contains-states. (Ex: 000#1#1#1000 ... becomes 0001#1#1#000)

- (A) 4
- (B) 3
- (C) 6
- (D) 5.

6. Which of following statements are true?

- (i) Let K, L be decidable languages. The concatenation of languages, K, L is also decidable language.
 - (ii) Let L be Turing recognizable language. Then the complement, L^c is also Turing recognizable language.
- (A) (i) and (ii)
 - (B) Only (ii)
 - (C) Both are false
 - (D) Only (i)

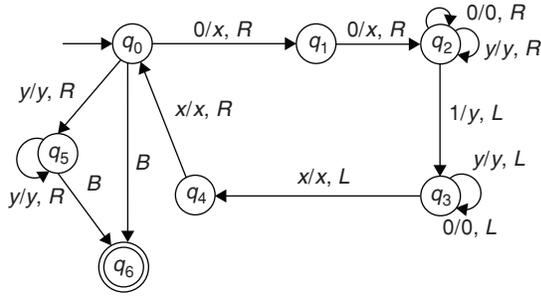
7. Let T_i denote i th TM. Given, X determines whether $X \in S$, Where the set S is defined inductively as follows: If $u \in S$, then $u^2 + 1, 3u + 2$ and $u!$ are all members of S . Which of following is true about the given decision problem?

- (A) Decidable
- (B) Un-decidable
- (C) Trivially decidable
- (D) No solution.

8. Fermat's last theorem asserts that there are no integer solution (x, y, z, n) to equation $x^n + y^n = z^n$ satisfying $x, y > 0$ and $n > 2$. Which of the following is true regarding the halting problem?

- (A) Decidable
- (B) Un-decidable
- (C) Trivially decidable
- (D) May or may not have solution.

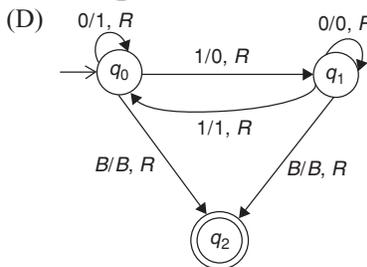
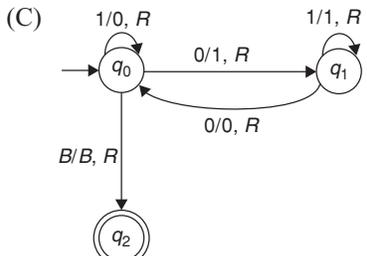
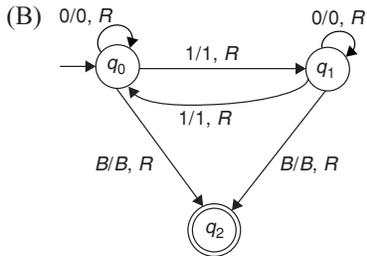
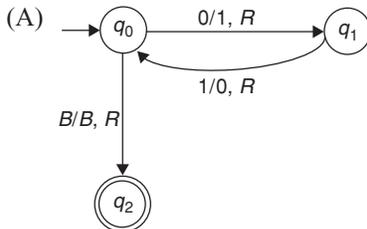
9. The TM, T is designed as



Which of following is true?

- (A) T halts on $0^n 1^n, n \geq 0$
- (B) T halts on $(01)^n (0^n 1^n), n \geq 0$
- (C) T halts on $0^{n^2} 1^{n^2}, n \geq 0$
- (D) T halts on $0^{2^n} 1^n, n \geq 0$

10. Design TM, which reads an input and starts inverting 0's to 1's till the first 1. The first 1 also inverted. After it has inverted first 1, it read the next symbols and keeps them as they are till the next 1. After encountering 1, it starts repeating the cycle by inverting the symbol till next 1. It halts when it encounters a blank symbol?



11. Consider three problems, P_1, P_2 and P_3 . It is known that P_1 has polynomial time solution, P_2 is NP-complete and P_3 is in NP. Which one of the following is true?

- (A) P_3 has polynomial time solution if P_1 is polynomial time reducible to P_3 .
- (B) P_3 is NP-complete if P_3 is polynomial time reducible to P_2 .
- (C) P_3 is NP complete if P_2 is reducible to P_3
- (D) P_3 has polynomial time complexity and P_3 is reducible to P_2 .

12. Let FHAM be the problem of finding a Hamiltonian cycle in a graph G and DHAM be the problem of determining if a Hamiltonian cycle exists in a graph. Which one of the following is true?

- (A) Both FHAM and DHAM are NP-hard.
- (B) FHAM is NP-hard, but DHAM is not.
- (C) DHAM is NP-hard but FHAM is not.
- (D) Neither DHAM nor FHAM is NP-hard.

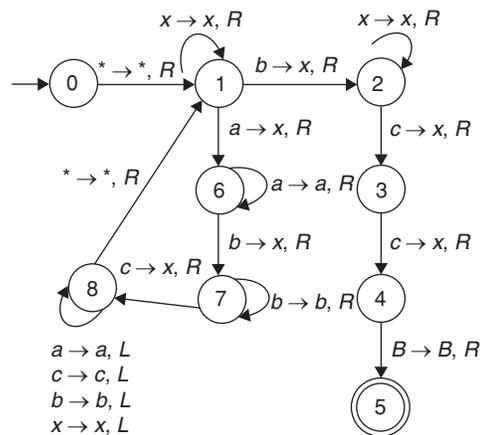
13. The solution for the system of post correspondence problem, $A = \{ba, abb, bab\}, B = \{bab, bb, abb\}$ is

- (A) 1312212
- (B) 15234434
- (C) 1311322
- (D) No solution.

14. A language, prefix_free REX = $\{R/R \text{ is a regular expression where } L(R) \text{ is prefix_free}\}$. Which of following is true about prefix_free REX?

- (A) Decidable
- (B) Un-decidable
- (C) Trivially decidable.
- (D) Can't be determined.

15. The TM, M is designed as:



Which of following is true about M ?

- (A) M is designed for $a^n b^n c^n, n \geq 0$
- (B) M is designed for $a^{n^2} b^{n^3} c^{n^4}, n \geq 0$
- (C) M is designed for $a^n b^{n+1} c^{n+2}, n \geq 0$
- (D) M is designed for $a^n b^n c^n, n > 0$

Practice Problems 2

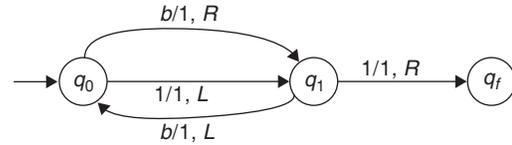
Directions for questions 1 to 15: Select the correct alternative from the given choices.

- Consider the language, $A\epsilon_{-CFG} = \{ \langle G \rangle : G \text{ is a CFG that generates } \epsilon \}$. Which of the following is true?
 - Undecidable
 - Decidable
 - Trivially decidable.
 - None of the above.
- The TM is designed with input and output as binary form. (# represents blank). The turing machine TM (M) is

	0	1	#
q_0	$(q_1, 0, R)$	$(q, 1, R)$	ϕ
q_1	$(q_1, 0, R)$	$(q, 1, R)$	$(q_2, \#, L)$
q_2	$(q_3, \#, L)$	$(q_3, \#, L)$	ϕ
q_3	$(q_3, 0, L)$	$(q_3, 1, L)$	$(q_4, \#, L)$
q_4	ϕ	ϕ	ϕ

Which of following is true?

- M accepts $2n$
 - M accepts n^2
 - M replaces left most symbol with #
 - M replaces right most symbol with #
- The TM is designed with 3-characters 0, 1, # to compute function $f(n) = 2n$. Input and output are to be in binary form and string represented by 'n' is enclosed between two #'s on left and right of it. b is blank symbol. TM contains _____ states.
 - 4
 - 3
 - 2
 - 1
 - The language $\{1^n | n \text{ is a prime number}\}$ is
 - Undecidable
 - Decidable
 - Trivially decidable
 - None of the above
 - Which of following statement(s) are true?
 - Let L be Turing decidable language. Then the complement \bar{L} is also Turing decidable language.
 - Let K and L be two Turing recognizable languages. The intersection, $K \cap L$ is also Turing recognizable language.
 - Both (i) and (ii)
 - Only (i)
 - Only (ii)
 - Neither (i) nor (ii) are true.
 - For the following two-way infinite TM, the equivalent one-way TM contains _____ states.



- 7
- 6
- 5
- 4

- L contains at least two strings. Which of following is true?
 - L has recursively enumerable sets and recursive.
 - L is recursive.
 - L has recursively enumerable sets but not recursive.
 - L does not contain recursively enumerable sets and also is not recursive.

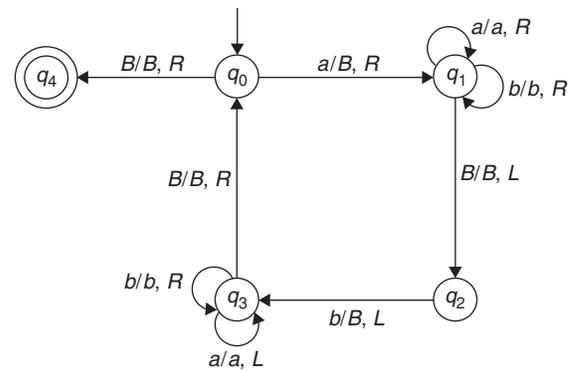
8. Consider the following TM:

Input State	0	1	B
q_0	$(q_0, 1, R)$	$(q_0, 0, R)$	(q_1, B, R)
q_1	-	-	-

What does TM generates?

- It display's the negative of given binary number.
- It computes one's complement of a binary number.
- It computes two's complement of a binary number
- It generates double the 0's as 1's.

9. Consider the following TM, M:



Which of following is true?

- M halts on $a^{n+1} b^n, n \geq 0$.
- M halts on $a^{n^2}, b^{n^3}, n \geq 0$.
- M halts on $(ab)^n, n \geq 0$.
- M halts on $a^n b^n, n \geq 0$.

10. A TM, M is designed generates language

- $L = \{a^n b^m : n \geq 1 \text{ and } n \neq m\}$. The number of states used are _____
- 5
 - 6
 - 7
 - 4

11. Consider three decision problems p_1, p_2 and p_3 . It is known that p_1 is decidable, p_2 is undecidable. Which one of following is true?
 (A) p_3 is decidable if p_1 is reducible to p_3
 (B) p_3 is undecidable if p_3 is reducible to p_2
 (C) p_3 is undecidable if p_2 is reducible to p_3
 (D) p_3 is decidable if p_3 is reducible to p_2 's complement.
12. Which one of following is not decidable?
 (A) Given a TM, M , a string S , and an integer K , M accepts S within K -steps.
 (B) Equivalence of two given Turing machines.
 (C) Language accepted by a given DFSA is non-empty.
 (D) Language accepted by a CFG is non-empty.
13. What is the solution for the correspondence system with two lists $x = \{b, bab^3, ba\}$ and $y = \{b^3, ba, a\}$
 (A) 1312213 (B) 2113
 (C) 3112 (D) No solution.
14. Given a Turing machine M , a state ' q ' and a string ' w '. To determine whether M ever reaches state q when started with input w from its initial state is?
 (A) Decidable
 (B) Un-decidable
 (C) Trivially decidable.
 (D) Can not be determined.
15. Given a Turing machine, M to determine whether M ever moves its head to the left when started with input W is:
 (A) Decidable
 (B) Un-decidable
 (C) Trivially decidable.
 (D) Can not be determined.

PREVIOUS YEARS' QUESTIONS

1. For $s \in (0 + 1)^*$, let $d(s)$ denote the decimal value of s (e.g., $d(101) = 5$). [2006]
 Let $L = \{s \in (0 + 1)^* | d(s) \bmod 5 = 2 \text{ and } d(s) \bmod 7 \neq 4\}$
 Which one of the following statements is true?
 (A) L is recursively enumerable, but not recursive
 (B) L is recursive, but not context-free
 (C) L is context-free, but not regular
 (D) L is regular
2. Which of the following is true for the language $\{a^p | p \text{ is a prime}\}$? [2008]
 (A) It is not accepted by a Turing Machine
 (B) It is regular but not context-free
 (C) It is context-free but not regular
 (D) It is neither regular nor context-free, but accepted by a Turing machine
3. If L and \bar{L} are recursively enumerable then L is [2008]
 (A) regular
 (B) context-free
 (C) context-sensitive
 (D) recursive
4. Let $L = L_1 \cap L_2$, where L_1 and L_2 are languages as defined below:
 $L_1 = \{a^m b^m c a^n b^n | m, n \geq 0\}$
 $L_2 = \{a^i b^j c^k | i, j, k \geq 0\}$
 Then L is [2009]
 (A) Not recursive
 (B) Regular
 (C) Context free but not regular
 (D) Recursively enumerable but not context free.
5. Let L_1 be a recursive language. Let L_2 and L_3 be languages that are recursively enumerable but not recursive. Which of the following statements is not necessarily true? [2010]
 (A) $L_2 - L_1$ is recursively enumerable
 (B) $L_1 - L_3$ is recursively enumerable
 (C) $L_2 \cap L_1$ is recursively enumerable
 (D) $L_2 \cup L_1$ is recursively enumerable
6. Which of the following statements is/are FALSE? [2013]
 1. For every non-deterministic Turing machine, there exists an equivalent deterministic Turing machine.
 2. Turing recognizable languages are closed under union and complementation.
 3. Turing decidable languages are closed under intersection and complementation.
 4. Turing recognizable languages are closed under union and intersection.
 (A) 1 and 4 only (B) 1 and 3 only
 (C) 2 only (D) 3 only
7. Let L be a language and \bar{L} be its complement. Which one of the following is NOT a viable possibility? [2014]
 (A) Neither L nor \bar{L} is recursively enumerable (r. e.)
 (B) One of L and \bar{L} is r.e. but not recursive, the other is not r. e.
 (C) Both L and \bar{L} are r.e. but not recursive
 (D) Both L and \bar{L} are recursive
8. Let $A \leq_m B$ denotes that language A is mapping reducible (also known as many-to-one reducible) to language B . Which one of the following is FALSE? [2014]
 (A) If $A \leq_m B$ and B is recursive then A is recursive.
 (B) If $A \leq_m B$ and A is undecidable then B is undecidable.
 (C) If $A \leq_m B$ and B is recursively enumerable then A is recursively enumerable.

- (D) If $A \leq_m B$ and B is not recursively enumerable then A is not recursively enumerable.
9. Let $\langle M \rangle$ be the encoding of a Turing machine as a string over $\Sigma = \{0, 1\}$. Let $L = \{\langle M \rangle \mid M \text{ is a Turing machine that accepts a string of length } 2014\}$. Then, L is
- (A) Decidable and recursively enumerable
 (B) Undesirable but recursively enumerable
 (C) Undesirable and not recursively enumerable
 (D) Decidable but not recursively enumerable
10. For any two languages L_1 and L_2 such that L_1 is context-free and L_2 is recursively enumerable but not recursive, which of the following is/are necessarily true? [2015]
- I. \bar{L}_1 (complement of L_1) is recursive
 II. \bar{L}_2 (complement of L_2) is recursive
 III. \bar{L}_1 is context-free
 IV. $\bar{L}_1 \cup L_2$ is recursively enumerable
- (A) I only (B) III only
 (C) III and IV only (D) I and IV only
11. Consider the following statements.
- I. The complement of every Turing decidable language is Turing decidable.
 II. There exists some language which is in NP but is not Turing decidable.
 III. If L is a language in NP, L is Turing decidable.
- Which of the above statements is/are true? [2015]
- (A) Only II (B) Only III
 (C) Only I and II (D) Only I and III
12. Let X be a recursive language and Y be a recursively enumerable but not recursive language. Let W and Z be two languages such that \bar{y} reduces to W , and Z reduces to \bar{x} (reduction means the standard many-one reduction). Which one of the following statements is **TRUE**? [2016]
- (A) W can be recursively enumerable and Z is recursive.
 (B) W can be recursive and Z is recursively enumerable.
 (C) W is not recursively enumerable and Z is recursive.
 (D) W is not recursively enumerable and Z is not recursive.
13. Consider the following types of languages: L_1 : Regular, L_2 : Context - free, L_3 : Recursive, L_4 : Recursively enumerable. Which of the following is / are TRUE? [2016]
- I. $\bar{L}_3 \cup L_4$ is recursively enumerable
 II. $\bar{L}_2 \cup L_3$ is recursive
 III. $L_1^* \cap L_2$ is context - free
 IV. $L_1 \cup \bar{L}_2$ is context - free
- (A) I only (B) I and III only
 (C) I and IV only (D) I, II and III only
14. Consider the following languages. [2016]
- $L_1 = \{\langle M \rangle \mid M \text{ takes at least } 2016 \text{ steps on some input}\}$,
 $L_2 = \{\langle M \rangle \mid M \text{ takes at least } 2016 \text{ steps on all inputs}\}$
 and
 $L_3 = \{\langle M \rangle \mid M \text{ accepts } \epsilon\}$
- where for each Turing machine M , $\langle M \rangle$ denotes a specific encoding of M . Which one of the following is TRUE?
- (A) L_1 is recursive and L_2, L_3 are not recursive
 (B) L_2 is recursive and L_1, L_3 are not recursive
 (C) L_1, L_2 are recursive L_3 is not recursive
 (D) L_1, L_2, L_3 are recursive
15. Let A and B be finite alphabets and let $\#$ be a symbol outside both A and B . Let f be a total function from A^* to B^* . We say f is *computable* if there exists a Turing machine M which given an input x in A^* , always halts with $f(x)$ on its tape. Let L_f denote the language $\{x \# f(x) \mid x \in A^*\}$. Which of the following statements is true: [2017]
- (A) f is computable if and only if L_f is recursive.
 (B) f is computable if and only if L_f is recursively enumerable.
 (C) If f is computable then L_f is recursive, but not conversely.
 (D) If f is computable then L_f is recursively enumerable, but not conversely.
16. Let $L(R)$ be the language represented by regular expression R . Let $L(G)$ be the language generated by a context free grammar G . Let $L(M)$ be the language accepted by a Turing machine M . Which of the following decision problems are undecidable? [2017]
- I. Given a regular expression R and a string w , is $w \in L(R)$?
 II. Given a context-free grammar G , is $L(G) = \emptyset$?
 III. Given a context-free grammar G , is $L(G) = \Sigma^*$ for some alphabet Σ ?
 IV. Given a Turing machine M and a string w , is $w \in L(M)$?
- (A) I and IV only (B) II and III only
 (C) II, III and IV only (D) III and IV only
17. The set of all recursively enumerable languages is: [2018]
- (A) Closed under complementation.
 (B) Closed under intersection.
 (C) A subset of the set of all recursive languages.
 (D) An uncountable set.

18. Consider the following problems. $L(G)$ denotes the language generated by a grammar G . $L(M)$ denotes the language accepted by a machine M .

- (I) For an unrestricted grammar G and a string w , whether $w \in L(G)$
- (II) Given a Turing machine M , whether $L(M)$ is regular
- (III) Given two grammars G_1 and G_2 , whether $L(G_1) = L(G_2)$

(IV) Given an NFA N , whether there is a deterministic PDA P such that N and P accept the same language.

Which one of the following statements is correct?

[2018]

- (A) Only I and II are undecidable
- (B) Only III is undecidable
- (C) Only II and IV are undecidable
- (D) Only I, II and III are undecidable

ANSWER KEYS

EXERCISES

Practice Problems 1

- | | | | | | | | | | |
|-------|-------|-------|-------|-------|------|------|------|------|-------|
| 1. D | 2. C | 3. B | 4. B | 5. D | 6. D | 7. A | 8. D | 9. D | 10. D |
| 11. C | 12. A | 13. D | 14. A | 15. C | | | | | |

Practice Problems 2

- | | | | | | | | | | |
|-------|-------|-------|-------|-------|------|------|------|------|-------|
| 1. B | 2. D | 3. B | 4. B | 5. A | 6. B | 7. C | 8. B | 9. D | 10. B |
| 11. C | 12. B | 13. B | 14. B | 15. A | | | | | |

Previous Years' Questions

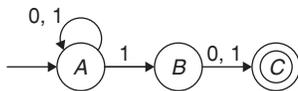
- | | | | | | | | | | |
|-------|-------|-------|-------|-------|-------|-------|-------|------|-------|
| 1. D | 2. D | 3. D | 4. C | 5. B | 6. C | 7. C | 8. D | 9. B | 10. D |
| 11. D | 12. C | 13. D | 14. C | 15. A | 16. D | 17. B | 18. D | | |

THEORY OF COMPUTATION

Time: 60 min.

Directions for questions 1 to 30: Select the correct alternative from the given choices.

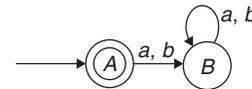
- Phrase structure languages, context-sensitive languages, context-free languages and regular languages are commonly referred to as languages of type-0, 1, 2, and 3 respectively. Then, Chomsky's Hierarchy states that
 - type-0 \supseteq type-1 \supseteq type-2 \supseteq type-3
 - type-0 \supset type-1 \supset type-2 \supset type-3
 - type-0 \subset type-1 \subset type-2 \subset type-3
 - type-0 \subseteq type-1 \subseteq type-2 \subseteq type-3
- Let L be a language recognizable by a finite automaton. The language $\text{Reverse}(L) = \{x \text{ such that } x \text{ is the reverse of } y \text{ where } y \in L\}$ is a
 - Regular language
 - Context-sensitive language
 - Context-free language
 - Phrase-structure language
- Which of the following statement is true?
 - It is possible to construct an NFA with more number of states than its equivalent minimum DFA.
 - There can be a DFA with more than one start state.
 - Both (A) and (B)
 - None of these
- Which of the following is an equivalent DFA for the NFA shown below:



-
-
-
-

- Which one of the following regular expressions over $\{0, 1\}$ denotes the set of strings not containing 100 as a substring?
 - $0^*(1^*0)^*$
 - $0^*1^*01^*$
 - $(0^*(10+1)^*)^*$
 - 0^*1010^*

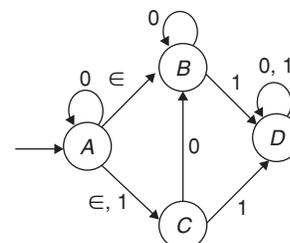
- The following transition diagram of a finite automaton accepts



- All word over sigma (a, b) such that symbol a and b alternate.
 - Only empty string.
 - Only the λ_1 meaning this automaton accepts no string of length greater than zero.
 - All words over sigma (a, b) except λ .
- Sentence that can be generated from the following production grammar is
 $S \rightarrow aS/bA$
 $A \rightarrow d/ccA$
 - $aabccccc$
 - $ababccccc$
 - $bccddd$
 - $aaccdb$
 - Pumping lemma is generally used for proving
 - A given grammar is regular
 - A given grammar is non-regular
 - Whether two given regular expression are equivalent or not
 - Both (A) and (C)
 - Finite state machine can recognize
 - Only context-free grammar
 - Only regular grammar
 - Any unambiguous grammar
 - Any grammar

- Which of the following is false?
 - Regular sets are closed under reversal.
 - Regular sets are closed under substitution.
 - Regular sets are closed under intersection.
 - None of these

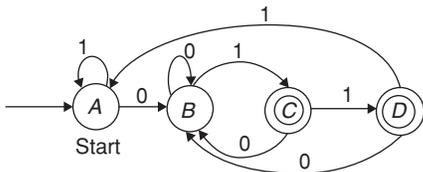
- For the DFA shown below $\hat{\delta}(A, 01)$ will be



- (A) $\{B, D\}$ (B) $\{C, D\}$
 (C) $\{A, B, C\}$ (D) $\{A, B, C, D\}$

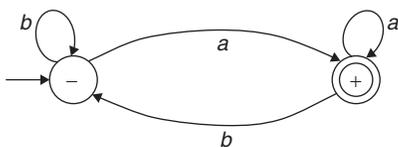
12. 'NFA can be simulated by a DFA'. The statement is
 (A) True (B) False
 (C) Depends on NFA (D) Depends on DFA
13. Given an arbitrary non-deterministic finite automaton (NFA) with N states, the maximum number of states in an equivalent minimized DFA is at least
 (A) N^2 (B) 2^N
 (C) $2N$ (D) $N!$
14. Let $M = (K, \Sigma, \delta, S, F)$ be a finite state automaton, Where
 $K = \{A, B\}$
 $\Sigma = \{a, b\}$
 $S = A$
 $F = \{B\}$,
 $\delta(A, a) = A$
 $\delta(A, b) = B$
 $\delta(B, a) = B$ and
 $\delta(B, b) = A$.
- A grammar to generate the language accepted by M can be specified as $G = (V, \Sigma, R, S)$, where
 $V = K \cup \Sigma$, and $S = A$. Which one of the following set of rules will make $L(G) = L(M)$?
 (A) $\{A \rightarrow aB, A \rightarrow bA, B \rightarrow bA, B \rightarrow aA, B \rightarrow \epsilon\}$
 (B) $\{A \rightarrow aA, A \rightarrow bB, B \rightarrow aB, B \rightarrow bA, B \rightarrow \epsilon\}$
 (C) $\{A \rightarrow bB, A \rightarrow aB, B \rightarrow aA, B \rightarrow bA, B \rightarrow \epsilon\}$
 (D) $\{A \rightarrow aA, A \rightarrow bA, B \rightarrow aB, B \rightarrow bA, A \rightarrow \epsilon\}$

15. A deterministic finite automaton M shown below has a start state A and accepting state D . Which of the following regular expression denotes the set of all words accepted by M ?



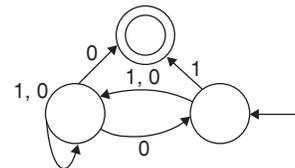
- (A) 001 (B) 10*1*10
 (C) 1*0*001 (D) (0/1)*011

16. Which of the following regular expression is/are true?
 (A) $(x^*)^* = x^*$ (B) $(x + y)^* = x^* + y^*$
 (C) $x^*y^* = x^* + y^*$ (D) All of these
17. Consider the FA shown in the figure given below, where '-' is the start state and '+' is the ending state. The language accepted by the FA is



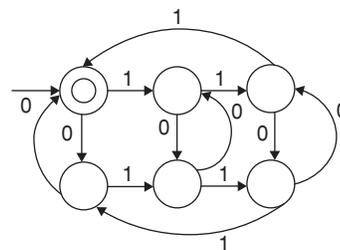
- (A) $(a + b)^*b$ (B) $(a + b)^*a$
 (C) a^*b (D) a^*b^*

18. Which of the following statement is false?
 (A) The family of regular language is closed under the complementary operation.
 (B) If L is a regular language, $L_1 = \{UV : U \in L, |V| = 2\}$ is also regular.
 (C) If L is a regular language, $L_1 = \{UV : U \in L, V \in L^R\}$ is also regular.
 (D) None of these
19. Which of the following is false?
 (A) $L = \{0^i 1^m 2^m : i \geq 1, m \geq 1\}$ over $\Sigma = \{0, 1, 2\}$ is regular.
 (B) $L = \{a^n b^l a^k, k \geq n + l\}$ is not regular.
 (C) $L = \{UWW^2V : U, V, W \in \{a, b\}^+\}$ is regular.
 (D) $L = \{a_n b_k : n > k\} \cup \{a_n b_k : n \neq k - 1\}$ is not regular.
20. Consider a DFA over $\Sigma = \{a, b\}$ accepting all strings which have number of a 's divisible by 6 and number of b 's divisible by 8. What is the minimum number of states that the DFA will have?
 (A) 16 (B) 15
 (C) 48 (D) 8
21. For the NFA M given below. Let the language, accepted by M be L . Let L_1 be the language accepted by the NFA M_1 , obtained by changing the accepting state of M to a non-accepting state and by changing the non-accepting states of M to accepting states. Which of the following statement is true?



- (A) $L_1 = A$ (B) $L_1 \subseteq L$
 (C) $L_1 = \{0, 1\}^*$ (D) $L_1 = (0, 1)^* - L$

22. The following finite state machine accepts all those binary strings in which the number of 1's and 0's are respectively.



- (A) Divisible by 3 and 2
 (B) Odd and even
 (C) Even and odd
 (D) Divisible by 2 and 5

Compiler Design

Chapter 1: Lexical Analysis and Parsing	6.3
Chapter 2: Syntax Directed Translation	6.27
Chapter 3: Intermediate Code Generation	6.36
Chapter 4: Code Optimization	6.56

U

N

I

T

6

This page is intentionally left blank

Chapter 1

Lexical Analysis and Parsing

LEARNING OBJECTIVES

- ☞ Language processing system
- ☞ Lexical analysis
- ☞ Syntax analysis
- ☞ Context free grammars and ambiguity
- ☞ Types of parsing
- ☞ Top down parsing
- ☞ Bottom up parsing
- ☞ Conflicts
- ☞ Operator precedence grammar
- ☞ LR parser
- ☞ Canonical LR parser(CLR)

LANGUAGE PROCESSING SYSTEM

Language Processors

Interpreter

It is a computer program that executes instructions written in a programming language. It either executes the source code directly or translates source code into some efficient intermediate representation and immediately executes this.



Example: Early versions of Lisp programming language, BASIC.

Translator

A software system that converts the source code from one form of the language to another form of language is called as translator. There are 2 types of translators namely (1) Compiler (2) Assembler.

Compiler converts source code of high level language into low level language.

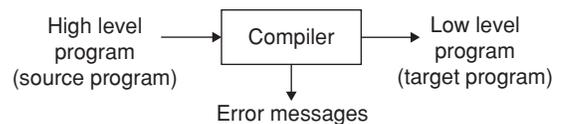
Assembler converts assembly language code into binary code.

Compilers

A compiler is a software that translates code written in high-level language (i.e., source language) into target language.

Example: source languages like C, Java,... etc. Compilers are user friendly.

The target language is like machine language, which is efficient for hardware.



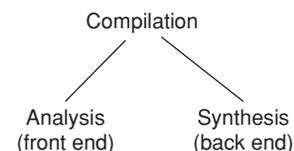
Passes

The number of iterations to scan the source code, till to get the executable code is called as a pass.

Compiler is two pass. Single pass requires more memory and multipass require less memory.

Analysis–synthesis model of compilation

There are two parts of compilation:

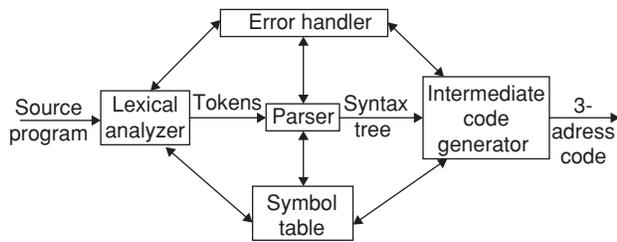


Analysis It breaks up the source program into pieces and creates an intermediate representation of the source program. This is more language specific.

Synthesis It constructs the desired target program from the intermediate representation. The target program will be more machine specific, dealing with registers and memory locations.

Front end vs back end of a compiler

The front end includes all analysis phases and intermediate code generator with part of code optimization.



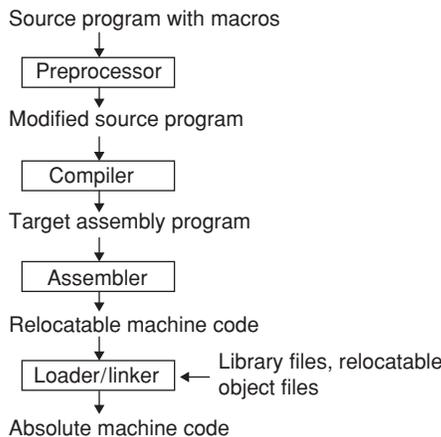
The back end includes code optimization and code generation phases. The back end synthesizes the target program from intermediate code.

Context of a compiler

In addition to a compiler, several other programs may be required to create an executable target program, like pre-processor to expand macros.

The target program created by a compiler may require further processing before it can be run.

The language processing system will be like this:



Phases

Compilation process is partitioned into some subprocesses called phases.

In order to translate a high level code to a machine code, we need to go phase by phase, with each phase doing a particular task and parsing out its output for the next phase.

Lexical analysis or scanning

It is the first phase of a compiler. The lexical analyzer reads the stream of characters making up the source program and groups the characters into meaningful sequences called lexemes.

Example: Consider the statement: `if (a < b)`

In this sentence the tokens are `if, (, a, <, b,)`.

Number of tokens = 6

Identifiers: `a, b`

Keywords: `if`

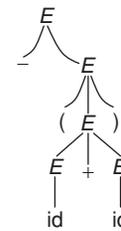
Operators: `<, (,)`

Syntax analyzer or Parser

- Tokens are grouped hierarchically into nested collections with collective meaning.
- A context free grammar (CFG) specifies the rules or productions for identifying constructs that are valid in a programming language. The output is a parse/syntax/ derivation tree.

Example: Parse tree for `-(id + id)` using the following grammar:

- $E \rightarrow E + E$
 - $E \rightarrow E * E$
 - $E \rightarrow -E$
 - $E \rightarrow (E)$
 - $E \rightarrow id$
- (G_1)



Semantic analysis

- It checks the source program for semantic errors.
- Type checking is done in this phase, where the compiler checks that each operator has matching operands for semantic consistency with the language definition.
- Gathers the type information for the next phases.

Example 1: The bicycle rides the boy.

This statement has no meaning, but it is syntactically correct.

Example 2:

```
int a;
bool b;
char c;
c = a + b;
```

We cannot add integer with a Boolean variable and assign it to a character variable.

Intermediate code generation

The intermediate representation should have two important properties:

- It should be easy to produce.
- Easy to translate into the target program

'Three address code' is one of the common forms of Intermediate code.

Three address code consists of a sequence of instructions, each of which has at most three operands.

Example:

```
id1 = id2 + id3 * 10;
t1 = inttoreal(10)
```

$$t_2 := id_3 \times t_1$$

$$t_3 := id_2 + t_2$$

$$id_1 = t_3$$

Code optimization

The output of this phase will result in faster running machine code.

Example: For the above intermediate code the optimized code will be

$$t_1 := id_3 \times 10.0$$

$$id_1 := id_2 + t_1$$

In this we eliminated t_2 and t_3 registers.

Code generation

- In this phase, the target code is generated.
- Generally the target code can be either a relocatable machine code or an assembly code.
- Intermediate instructions are each translated into a sequence of machine instructions.
- Assignment of registers will also be done.

Example:

MOVF	id_3, R_2
MULF	$\neq 60.0, R_2$
MOVF	id_2, R_1
ADDF	R_2, R_1
MOVF	R_1, id_1

Symbol table management

A symbol table is a data structure containing a record for each variable name, with fields for the attributes of the name.

What is the use of a symbol table?

1. To record the identifiers used in the source program.
2. Its type and scope
3. If it is a procedure name then the number of arguments, types of arguments, the method of parsing (by reference) and the type returned.

Error detection and reporting

- (i) Lexical phase can detect errors where the characters remaining in the input 'do not form any token'.
- (ii) Errors of the type, 'violation of syntax' of the language are detected by syntax analysis.
- (iii) Semantic phase tries to detect constructs that have the right syntactic structure but no meaning.

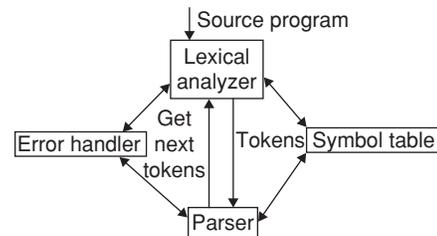
Example: adding two array names etc.

LEXICAL ANALYSIS

Lexical Analysis is the first phase in compiler design. The main task of the lexical analyzer is to read the input characters of the source program, group them into lexemes, and produce as output a sequence of tokens for each lexeme

in the source program. The stream of tokens is sent to the parser for syntax analysis.

There will be interaction with the symbol table as well.



Lexeme: Sequence of characters in the source program that matches the pattern for a token. It is the smallest logical unit of a program.

Example: 10, x, y, <, >, =

Tokens: These are the classes of similar lexemes.

Example: Operators: <, >, =
 Identifiers: x, y
 Constants: 10
 Keywords: if, else, int

Operations performed by lexical analyzer

1. Identification of lexemes and spelling check
2. Stripping out comments and white space (blank, new line, tab etc).
3. Correlating error messages generated by the compiler with the source program.
4. If the source program uses a macro-preprocessor, the expansion of macros may also be performed by lexical analyzer.

Example 1: Take the following example from Fortran
 DO 5 I = 1.25
 Number of tokens = 5
 The 1st lexeme is the keyword DO
 Tokens are DO, 5, I, =, 1.25.

Example 2: An example from C program
 for (int i = 1; i <= 10; i++)
 Here tokens are for, (, int, i, =, 1,;, i, <=, 10,;, i, ++,)
 Number of tokens = 13

LEX compiler

Lexical analyzer divides the source code into tokens. To implement lexical analyzer we have two techniques namely hand code and the other one is LEX tool.

LEX is an automated tool which specifies lexical analyzer, from the rules given by the regular expression.

These rules are also called as pattern recognizing rules.

SYNTAX ANALYSIS

This is the 2nd phase of the compiler, checks the syntax and constructs the syntax/parse tree.

Input of parser is token and output is a parse/ syntax tree.

Constructing parse tree

Construction of derivation tree for a given input string by using the production of grammar is called parse tree.

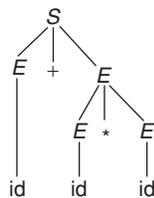
Consider the grammar

$$S \rightarrow E + E/E * E$$

$$E \rightarrow id$$

The parse tree for the string

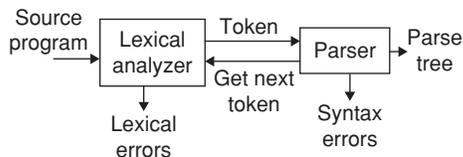
$w = id + id * id$ is



$w = id + id * id$

Role of the parser

1. Construct a parse tree.
2. Error reporting and correcting (or) recovery. A parser can be modeled by using CFG (Context Free Grammar) recognized by using pushdown automata/table driven parser.
3. CFG will only check the correctness of sentence with respect to syntax not the meaning.



How to construct a parse tree?

Parse tree's can be constructed in two ways.

- (i) Top-down parser: It builds parse trees from the top (root) to the bottom (leaves).
- (ii) Bottom-up parser: It starts from the leaves and works up to the root.

In both cases, the input to the parser is scanned from left to right, one symbol at a time.

Parser generator

Parser generator is a tool which creates a parser.

Example: compiler – compiler, YACC

The input of these parser generator is grammar we use and the output will be the parser code.

The parser generator is used for construction of the compilers front end.

Scope of declarations

Declaration scope refers to the certain program text portion, in which rules are defined by the language.

Within the defined scope, entity can access legally to declared entities.

The scope of declaration contains immediate scope always. Immediate scope is a region of declarative portion with enclosure of declaration immediately.

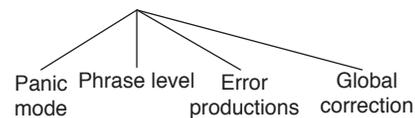
Scope starts at the beginning of declaration and scope continues till the end of declaration. Whereas in the over loadable declaration, the immediate scope will begin, when the callable entity profile was determined.

The visible part refers text portion of declaration, which is visible from outside.

Syntax Error Handling

1. Reports the presence of errors clearly and accurately.
2. Recovers from each error quickly.
3. It should not slow down the processing of correct programs.

Error Recovery Strategies



Panic mode On discovering an error, the parser discards input symbols one at a time until one of the synchronizing tokens is found.

Phrase level A parser may perform local correction on the remaining input. It may replace the prefix of the remaining input.

Error productions Parser can generate appropriate error messages to indicate the erroneous construct that has been recognized in the input.

Global corrections There are algorithms for choosing a minimal sequence of changes to obtain a globally least cost correction.

CONTEXT FREE GRAMMARS AND AMBIGUITY

A grammar is a set of rules or productions which generates a collection of finite/infinite strings.

It is a 4-tuple defined as $G = (V, T, P, S)$

Where

V = set of variables

T = set of terminals

P = set of production rules

S = start symbol

Example: $S \rightarrow (S)e$
 $S \rightarrow (S)$ (1)
 $S \rightarrow e$ (2)

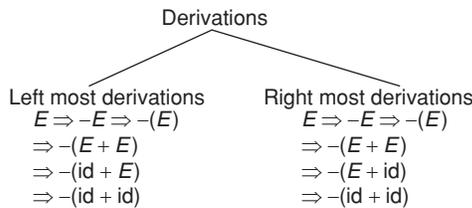
Here S is start symbol and the only variable.
 (,), e is terminals.
 (1) and (2) are production rules.

Sentential forms

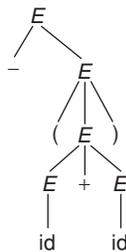
$s \xRightarrow{*} \alpha$, Where α may contain non-terminals, then we say that α is a sentential form of G.

Sentence: A sentence is a sentential form with no non-terminals.

Example: $-(id + id)$ is a sentence of the grammar (G_1).



Right most derivations are also known as canonical derivations.



Ambiguity

A grammar that produces more than one parse tree for some sentence is said to be ambiguous.

Or

A grammar that produces more than one left most or more than one right most derivations is ambiguous.

For example consider the following grammar:

$String \rightarrow String + String / String - String / 0 / 1 / 2 / \dots / 9$

$9 - 5 + 2$ has two parse trees as shown below

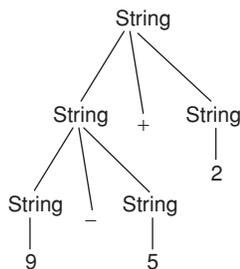


Figure 1 Leftmost derivation

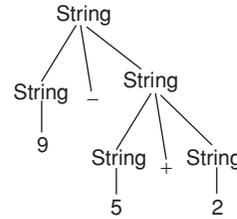


Figure 2 Rightmost derivation

- Ambiguity is problematic because the meaning of the program can be incorrect.
- Ambiguity can be handled in several ways
 - Enforce associativity and precedence
 - Rewrite the grammar by eliminating left recursion and left factoring.

Removal of ambiguity

The grammar is said to be ambiguous if there exists more than one derivation tree for the given input string.

The ambiguity of grammar is undecidable; ambiguity of a grammar can be eliminated by rewriting the grammar.

Example:

$E \rightarrow E + E / id$ } \rightarrow ambiguous grammar
 $E \rightarrow E + T / T$ } rewritten grammar
 $T \rightarrow id$ } (unambiguous grammar)

Left recursion

Left recursion can take the parser into infinite loop so we need to remove left recursion.

Elimination of left recursion

$A \rightarrow A\alpha/\beta$ is a left recursive.

It can be replaced by a non-recursive grammar:

$$A \rightarrow \beta A'$$

$$A' \rightarrow \alpha A' / \epsilon$$

In general

$$A \rightarrow A\alpha_1 / A\alpha_2 / \dots / A\alpha_m / \beta_1 / \beta_2 / \dots / \beta_n$$

We can replace A productions by

$$A \rightarrow \beta_1 A' / \beta_2 A' / \dots / \beta_n A'$$

$$A' \rightarrow \alpha_1 A' / \alpha_2 A' / \dots / \alpha_m A'$$

Example 3: Eliminate left recursion from

$$E \rightarrow E + T / T$$

$$T \rightarrow T * F / F$$

$$F \rightarrow (E) / id$$

Solution $E \rightarrow E + T / T$ it is in the form

$$A \rightarrow A\alpha/\beta$$

So, we can write it as $E \rightarrow TE'$

$$E' \rightarrow +TE' / \epsilon$$

Similarly other productions are written as

$$T \rightarrow FT'$$

$$T^1 \rightarrow * FT' / \epsilon$$

$$F \rightarrow (E) / id$$

Example 4 Eliminate left recursion from the grammar

$$S \rightarrow (L)/a$$

$$L \rightarrow L, S/b$$

Solution: $S \rightarrow (L)/a$
 $L \rightarrow bL'$
 $L' \rightarrow SL'/\epsilon$

Left factoring

A grammar with common prefixes is called non-deterministic grammar. To make it deterministic we need to remove common prefixes. This process is called as Left Factoring.

The grammar: $A \rightarrow \alpha\beta_1/\alpha\beta_2$ can be transformed into
 $A \rightarrow \alpha A'$
 $A' \rightarrow \beta_1/\beta_2$

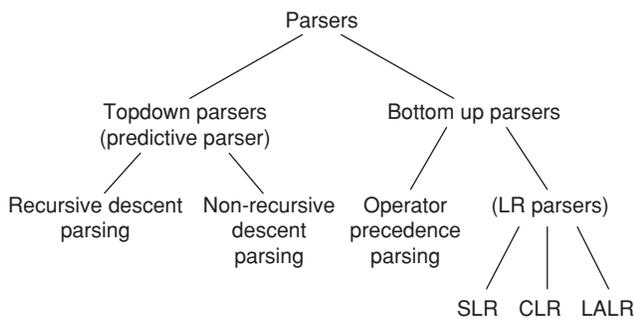
Example 5: What is the resultant grammar after left factoring the following grammar?

$$S \rightarrow iEtS/iEtSeS/a$$

$$E \rightarrow b$$

Solution: $S \rightarrow iEtSS'/a$
 $S' \rightarrow eS/\epsilon$
 $E \rightarrow b$

TYPES OF PARSING



TOPDOWN PARSING

A parse tree is constructed for the input starting from the root and creating the nodes of the parse tree in preorder. It simulates the left most derivation.

Backtracking Parsing

If we make a sequence of erroneous expansions and subsequently discover a mismatch we undo the effects and roll back the input pointer.

This method is also known as brute force parsing.

Example: $S \rightarrow cAd$
 $A \rightarrow ab/a$

Let the string $w = cad$ is to generate:



The string generated from the above parse tree is $cabd$. but, $w = cad$, the third symbol is not matched.

So, report error and go back to A .

Now consider the other alternative for production A .



String generated 'cad' and $w = cad$. Now, it is successful.

In this we have used back tracking. It is costly and time consuming approach. Thus an outdated one.

Predictive Parsers

By eliminating left recursion and by left factoring the grammar, we can have parse tree without backtracking. To construct a predictive parser, we must know,

1. Current input symbol
2. Non-terminal which is to be expanded

A procedure is associated with each non-terminal of the grammar.

Recursive descent parsing

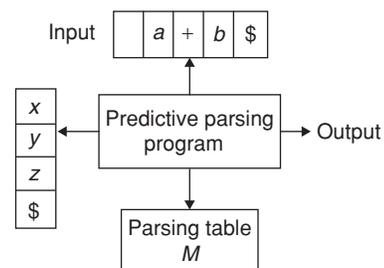
In recursive descent parsing, we execute a set of recursive procedures to process the input.

The sequence of procedures called implicitly, defines a parse tree for the input.

Non-recursive predictive parsing

(table driven parsing)

- It maintains a stack explicitly, rather than implicitly via recursive calls.
- A table driven predictive parser has
 - An input buffer
 - A stack
 - A parsing table
 - Output stream



Constructing a parsing table

To construct a parsing table, we have to learn about two functions:

1. FIRST ()
2. FOLLOW ()

FIRST(X) To compute FIRST(X) for all grammar symbols X, apply the following rules until no more terminals or ε can be added to any FIRST set.

1. If X is a terminal, then FIRST(X) is {X}.
2. If $X \rightarrow \epsilon$ is a production, then add ε to FIRST(X).
3. If X is non-terminal and $X \rightarrow Y_1 Y_2 \dots Y_k$ is a production, then place 'a' in FIRST(X) if for some i, a is an FIRST (Y_i) and ε is in all of FIRST(Y_1), ..., FIRST(Y_{i-1}); that is, $Y_1, \dots, Y_{i-1} \xRightarrow{*} \epsilon$. If ε is in FIRST (Y_j) for all $j = 1, 2, \dots, k$, then add ε to FIRST(X). For example, everything in FIRST (Y_1) is surely in FIRST(X). If Y_1 does not derive ε, then add nothing more to FIRST(X), but if $Y_1 \xRightarrow{*} \epsilon$, then add FIRST (Y_2) and so on.

FOLLOW (A): To compute FOLLOW (A) for all non-terminals A, apply the following rules until nothing can be added to any FOLLOW set.

1. Place \$ in FOLLOW(S), where S is the start symbol and \$ is input right end marker.
2. If there is a production $A \rightarrow \alpha B \beta$, then everything in FIRST (β) except ε is placed in FOLLOW (B).
3. If there is a production $A \rightarrow \alpha B$ or a production $A \rightarrow \alpha B \beta$, where FIRST (β) contains ε, then everything in FOLLOW (A) is in FOLLOW (B).

Example: Consider the grammar

$E \rightarrow TE'$
 $E' \rightarrow +TE'/\epsilon$
 $T \rightarrow FT'$
 $T' \rightarrow *FT'/\epsilon$
 $F \rightarrow (E)/id$. Then
 FIRST (E) = FIRST (T) = FIRST (F) = { (, id }
 FIRST (E') = { +, ε }
 FIRST (T') = { *, ε }
 FOLLOW (E) = FOLLOW (E') = { }, \$ }
 FOLLOW (T) = FOLLOW (T') = { +, }, \$ }
 FOLLOW (F) = { *, +, }, \$ }

Steps for the construction of predictive parsing table

1. For each production $A \rightarrow \alpha$ of the grammar, do steps 2 and 3.
2. For each terminal a in FIRST (α), add $A \rightarrow \alpha$ to $M[A, a]$
3. If ε is in FIRST (α), add $A \rightarrow \alpha$ to $M[A, b]$ for each terminal b in FOLLOW (A). If ε is in FIRST (α) and \$ is in FOLLOW (A), add $A \rightarrow \alpha$ to $M[A, \$]$
4. Make each undefined entry of M be error.

By applying these rules to the above grammar, we will get the following parsing table.

Non-terminal	Input Symbol					
	id	+	*	()	\$
E	$E \rightarrow TE'$			$E \rightarrow TE'$		
E'		$E' \rightarrow +TE'$			$E' \rightarrow \epsilon$	$E' \rightarrow \epsilon$
T	$T \rightarrow FT'$			$T \rightarrow FT'$		
T'		$T' \rightarrow \epsilon$	$T' \rightarrow *FT'$		$T' \rightarrow \epsilon$	$T' \rightarrow \epsilon$
F	$F \rightarrow id$			$F \rightarrow (E)$		

The parser is controlled by a program. The program consider x, the symbol on top of the stack and 'a' the current input symbol.

1. If $x = a = \$$, the parser halts and announces successful completion of parsing.
2. If $x = a \neq \$$, the parser pops x off the stack and advances the input pointer to the next input symbol.
3. If x is a non-terminal, the program consults entry $M[x, a]$ of the parsing table M. This entry will be either an x-production of the grammar or an error entry. If $M[x, a] = \{x \rightarrow UVW\}$, the parser replaces x on top of the stack by WVU with U on the top.

If $M[x, a] = \text{error}$, the parser calls an error recovery routine.

For example, consider the moves made by predictive parser on input id + id * id, which are shown below:

Matched	Stack	Input	Action
	E\$	id+id*id\$	
	TE'\$	id+id*id\$	Output E → TE'
	FT'E'\$	id+id*id\$	Output T → FT'
	idT'E'\$	id+id*id\$	Output F → id
id	T'E'\$	+id*id\$	Match id
id	E'\$	+id*id\$	Output T' → ε
id	+TE'\$	+id*id\$	Output E' → +TE'
id+	TE'\$	id*id\$	Match+
id+	FT'E'\$	id*id\$	Output T → FT'
id+	idT'E'\$	id*id\$	Output F → id
id+id	T'E'\$	*id\$	Match id
id+id	*FT'E'\$	*id\$	Output T' → *FT'
id+id*	FT'E'\$	id\$	Match*
id+id*	idT'E'\$	id\$	Output F → id
id+id*id	T'E'\$	\$	Match id
id+id*id	E'\$	\$	Output T' → ε
id+id*id	\$	\$	Output E' → ε

BOTTOM UP PARSING

- This parsing constructs the parse tree for an input string beginning at the leaves and working up towards the root.
- General style of bottom-up parsing is shift-reduce parsing.

Shift-Reduce Parsing

Reduce a string to the start symbol of the grammar. It simulates the reverse of right most derivation.

In every step a particular substring is matched (in left right fashion) to the right side of some production and then it is substituted by the non-terminal in the left hand side of the production.

For example consider the grammar

$$\begin{aligned} S &\rightarrow aABe \\ A &\rightarrow Abc/b \\ B &\rightarrow d \end{aligned}$$

In bottomup parsing the string ‘abcde’ is verified as

$$\left. \begin{array}{l} abcde \\ aAbcde \\ aAde \\ aABe \\ S \end{array} \right\} \rightarrow \text{reverse order}$$

Stack implementation of shift-reduce parser

The shift reduce parser consists of input buffer, Stack and parse table.

Input buffer consists of strings, with each cell containing only one input symbol.

Stack contains the grammar symbols, the grammar symbols are inserted using shift operation and they are reduced using reduce operation after obtaining handle from the collection of buffer symbols.

Parse table consists of 2 parts goto and action, which are constructed using terminal, non-terminals and compiler items.

Let us illustrate the above stack implementation.

→ Let the grammar be

$$\begin{aligned} S &\rightarrow AA \\ A &\rightarrow aA \\ A &\rightarrow b \end{aligned}$$

Let the input string ‘ ω ’ be abab\$

$\omega = abab\$$

Stack	Input String	Action
\$	abab\$	Shift
\$a	bab\$	Shift
\$ab	ab\$	Reduce ($A \rightarrow b$)
\$aA	ab\$	Reduce ($A \rightarrow aA$)
\$A	ab\$	Shift
\$Aa	b\$	Shift
\$Aab	\$	Reduce ($A \rightarrow b$)
\$AaA	\$	Reduce ($A \rightarrow aA$)
\$AA	\$	Reduce ($S \rightarrow AA$)
\$S	\$	Accept

Rightmost derivation

$$S \Rightarrow aABe \Rightarrow aAde \Rightarrow aAbcde \Rightarrow abcde$$

For bottom up parsing, we are using right most derivation in reverse.

Handle of a string Substring that matches the RHS of some production and whose reduction to the non-terminal on the LHS is a step along the reverse of some rightmost derivation.

$$S \xRightarrow{*} \alpha Ar \Rightarrow \alpha \beta r$$

Right sentential forms of a unambiguous grammar have one unique handle.

Example: For grammar, $S \rightarrow aABe$

$$\begin{aligned} A &\rightarrow Abc/b \\ B &\rightarrow d \end{aligned}$$

$$S \Rightarrow \underline{aABe} \Rightarrow \underline{aAde} \Rightarrow \underline{aAbcde} \Rightarrow \underline{abcde}$$

Note: Handles are underlined.

Handle pruning The process of discovering a handle and reducing it to the appropriate left hand side is called handle pruning. Handle pruning forms the basis for a bottomup parsing.

To construct the rightmost derivation:

$$S = r_0 \Rightarrow r_1 \Rightarrow r_2 \dots \Rightarrow r_n = w$$

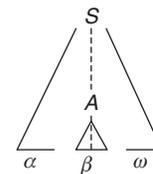
Apply the following simple algorithm:

For $i \leftarrow n$ to 1

Find the handle $A_i \rightarrow B_i$ in r_i

Replace B_i with A_i to generate r_{i-1}

Consider the cut of a parse tree of a certain right sentential form:



Here $A \rightarrow \beta$ is a handle for $\alpha\beta\omega$.

Shift reduce parsing with a stack There are 2 problems with this technique:

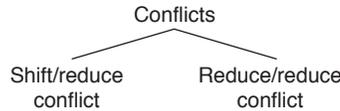
- To locate the handle
- Decide which production to use

General construction using a stack

- ‘Shift’ input symbols onto the stack until a handle is found on top of it.
- ‘Reduce’ the handle to the corresponding non-terminal.
- ‘Accept’ when the input is consumed and only the start symbol is on the stack.
- Errors – call an error reporting/recovery routine.

Viable prefixes The set of prefixes of a right sentential form that can appear on the stack of a shift reduce parser are called viable prefixes.

Conflicts



Shift/reduce conflict

Example: $stmt \rightarrow if\ expr\ then\ stmt\ | if\ expr\ then\ stmt\ else\ stmt\ | any\ other\ statement$

If $exp\ then\ stmt$ is on the stack, in this case we can't tell whether it is a handle. i.e., 'shift/reduce' conflict.

Reduce/reduce conflict

Example: $S \rightarrow aA/bB$

$A \rightarrow c$

$B \rightarrow c$

$W = ac$ it gives reduce/reduce conflict.

Operator Precedence Grammar

In operator grammar, no production rule can have:

- ϵ at the right side.
- two adjacent non-terminals at the right side.

Example 1: $E \rightarrow E + E / E - E / id$ is operator grammar.

Example 2: $E \rightarrow AB$
 $A \rightarrow a$
 $B \rightarrow b$ } not operator grammar

Example 3: $E \rightarrow E0E/id$
 not operator grammar

Precedence relation If

$a < b$ then b has higher precedence than a

$a = b$ then b has same precedence as a

$a > b$ then b has lower precedence than a

Common ways for determining the precedence relation between pair of terminals:

1. Traditional notations of associativity and precedence.

Example: \times has higher precedence than $+$ $\times . > + (or) + < . \times$

2. First construct an unambiguous grammar for the language which reflects correct associativity and precedence in its parse tree.

Operator precedence relations from associativity and precedence

Let us use $\$$ to mark end of each string. Define $\$ < . b$ and $b > \$$ for all terminals b . Consider the grammar is:

$$E \rightarrow E + E / E \times E / id$$

Let the operator precedence table for this grammar is:

	id	+	\times	$\$$
id		\succ	\succ	\succ
+	\prec		\prec	\succ
\times	\prec	\succ		\succ
$\$$	\prec	\prec	\prec	accept

1. Scan the string from left until \succ is encountered
2. Then scan backwards (to left) over any $=$ until \prec is encountered.
3. The handle contains everything to the left of the first \succ and to the right of the \prec is encountered.

After inserting precedence relation is

$$\$ id + id * id \$ \text{ is}$$

$$\$ < id > + < id > * < id > \$$$

Precedence functions Instead of storing the entire table of precedence relations table, we can encode it by precedence functions f and g , which map terminal symbols to integers:

1. $f(a) < f(b)$ whenever $a < b$
2. $f(a) > f(b)$ whenever $a \succ b$
3. $f(a) > f(b)$ whenever $a \succ b$

Finding precedence functions for a table

1. Create symbols $f(a)$ and $g(a)$ for each ' a ' that is a terminal or $\$$.
2. Partition the created symbols into as many groups as possible in such away that $a = b$ then $f(a)$ and $g(b)$ are in the same group
3. Create a directed graph
 If $a < b$ then place an edge from $g(b)$ to $f(a)$
 If $a > b$ then place an edge from $f(a)$ to $g(b)$
4. If the graph constructed has a cycle then no precedence function exists.
 If there are no cycles, let $f(a)$ be the length of the longest path being at the group of $f(a)$.
 Let $g(a)$ be the length of the longest path from the group of $g(a)$.

Disadvantages of operator precedence parsing

- It can not handle unary minus.
- Difficult to decide which language is recognized by grammar.

Advantages

1. Simple
2. Powerful enough for expressions in programming language.

Error cases

1. No relation holds between the terminal on the top of stack and the next input symbol.
2. A handle is found, but there is no production with this handle as the right side.

Error recovery

1. Each empty entry is filled with a pointer to an error routine.
2. Based on the handle tries to recover from the situation.

To recover, we must modify (insert/change)

1. Stack or
2. Input or
3. Both

We must be careful that we don't get into an infinite loop.

LR Parsers

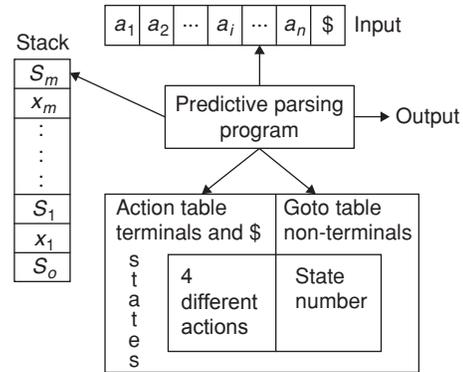
- In LR (K), L stands for Left to Right Scanning, R stands for Right most derivation, K stands for number of look ahead symbols.
- LR parsers are table-driven, much like the non-recursive LL parsers. A grammar which is used in construction of LR parser is LR grammar. For a grammar to be LR it is sufficient that a left-to-right shift-reduce parser be able to recognize handles of right-sentential forms when they appear on the top of the stack.
- The Time complexity for such parsers is $O(n^3)$
- LR parsers are faster than LL (1) parser.
- LR parsing is attractive because
 - The most general non-backtracking shift reduce parser.
 - The class of grammars that can be passed using LR methods is a proper superset of predictive parsers. $LL(1) \text{ grammars} \subset LR(1) \text{ grammars}$.
 - LR parser can detect a syntactic error in the left to right scan of the input.
- LR parsers can be implemented in 3 ways:
 1. Simple LR (SLR): The easiest to implement but the least powerful of the three.
 2. Canonical LR (CLR): most powerful and most expensive.
 3. Look ahead LR (LALR): Intermediate between the remaining two. It works on most programming language grammars.

Disadvantages of LR parser

1. Detecting a handle is an overhead, parse generator is used.
2. The main problem is finding the handle on the stack and it was replaced with the non-terminal with the left hand side of the production.

The LR parsing algorithm

- It consists of an input, an output, a stack, a driver program and a parsing table that has two parts (action and goto).
- The driver/parser program is same for all these LR parsers, only the parsing table changes from parser to another.



Stack: To store the string of the form,

$S_o x_1 S_1 \dots x_m S_m$ where

S_m : state

x_m : grammar symbol

Each state symbol summarizes the information contained in the stack below it.

Parsing table: Parsing table consists of two parts:

1. Action part
2. Goto part

ACTION Part:

Let, $S_m \rightarrow$ top of the stack

$a_i \rightarrow$ current symbol

Then action $[S_m, a_i]$ which can have one of four values:

1. Shift S , where S is a state
2. Reduce by a grammar production $A \rightarrow \beta$
3. Accept
4. Error

GOTO Part:

If $\text{goto}(S, A) = X$ where $S \rightarrow$ state, $A \rightarrow$ non-terminal, then GOTO maps state S and non-terminal A to state X .

Configuration

$$(S_o x_1 S_1 x_2 S_2 \dots x_m S_m, a_i a_{i+1} \dots a_n \$)$$

The next move of the parser is based on action $[S_m, a_i]$

The configurations are as follows.

1. If action $[S_m, a_i] = \text{shift } S$

$$(S_o x_1 S_1 x_2 S_2 \dots x_m S_m, a_i a_{i+1} \dots a_n \$)$$

2. If action $[S_m, a_i] = \text{reduce } A \rightarrow \beta$ then

$$(S_o x_1 S_1 x_2 S_2 \dots x_{m-r} S_{m-r}, AS, a_i a_{i+1} \dots a_n \$)$$

Where $S = \text{goto}[S_{m-r}, A]$

3. If action $[S_m, a_i] = \text{accept}$, parsing is complete.
4. If action $[S_m, a_i] = \text{error}$, it calls an error recovery routine.

Example: Parsing table for the following grammar is shown below:

1. $E \rightarrow E + T$
2. $E \rightarrow T$

- 3. $T \rightarrow T * F$
- 4. $T \rightarrow F$
- 5. $F \rightarrow (E)$
- 6. $F \rightarrow id$

State	Action					Goto		
	id	+	x	()	\$	E	T	F
0	S_5			S_4		1	2	3
1		S_6			<i>acc</i>			
2		r_2	S_7		r_2 r_2			
3		r_4	r_4		r_4 r_4			
4	S_5			S_4		8	2	3
5		r_6	r_6		r_6 r_6			
6	S_5			S_4			9	3
7	S_5			S_4				10
8		S_6		S_1				
9		r_1	S_7		r_1 r_1			
10		r_3	r_3		r_3 r_3			
11		r_5	r_5		r_5 r_5			

Moves of LR parser on input string id*id+id is shown below:

Stack	Input	Action
0	id * id + id\$	Shift 5
0id 5	* id + id\$	reduce 6 means reduce with 6th production $F \rightarrow id$ and goto [0, F] = 3
0F 3	* id + id\$	reduce 4 i.e $T \rightarrow F$ goto [0, T] = 2
0T 2	* id + id\$	Shift 7
0T2 * 7	id + id\$	Shift 5
0T2 * 7 id 5	+ id\$	reduce 6 i.e $F \rightarrow id$ goto [7, F] = 10
0T2 * 7 F 10	+ id\$	reduce 3 i.e $T \rightarrow T * F$
0T 2	+ id\$	goto [0, T] = 2
0E 1	+ id\$	reduce 2 i.e $E \rightarrow T$ & goto [0, E] = 1
0E1 + 6	id\$	Shift 6
0E1 + 6 id 5	\$	Shift 5
0E1 + 6F 3	\$	reduce 6 & goto [6, F] = 3
0E1 + 6T 9	\$	reduce 4 & goto [6, T] = 9
0E1	\$	reduce 1 & goto [0, E] = 1
0E1	\$	accept

Constructing SLR parsing table

LR (0) item: LR (0) item of a grammar G is a production of G with a dot at some position of the right side of production.

Example: $A \rightarrow BCD$

Possible LR (0) items are

- $A \rightarrow .BCD$
- $A \rightarrow B.CD$
- $A \rightarrow BC.D$
- $A \rightarrow BCD.$

$A \rightarrow B.CD$ means we have seen an input string derivable from B and hope to see a string derivable from CD.

The LR (0) items are constructed as a DFA from grammar to recognize viable prefixes.

The items can be viewed as the states of NFA.

The LR (0) item (or) canonical LR (0) collection, provides the basis for constructing SLR parser.

To construct LR (0) items, define

- (a) An augmented grammar
- (b) closure and goto

Augmented grammar (G') If G is a grammar with start symbol S, G' the augmented grammar for G, with new start symbol S' and production $S' \rightarrow S$.

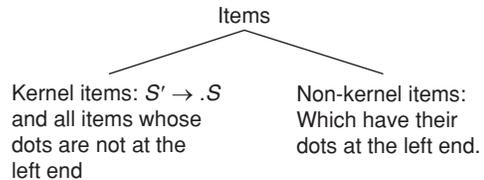
Purpose of G' is to indicate when to stop parsing and announce acceptance of the input.

Closure operation Closure (I) includes

1. Initially, every item in I is added to closure (I)
2. If $A \rightarrow \alpha.B\beta$ is in closure (I) and $\beta \rightarrow \gamma$ is a production then add $B \rightarrow .\gamma$ to I.

Goto operation

Goto (I, x) is defined to be the closure of the set of all items $[A \rightarrow \alpha X \beta]$ such that $[A \rightarrow \alpha X \beta]$ is in I.



Construction of sets of Items

Procedure items (G')

Begin

$C := \text{closure}(\{[S' \rightarrow .S]\})$;

repeat

For each set of items I in C and each grammar symbol x

Such that goto (I, x) is not empty and not in C do add goto (I, x) to C;

Until no more sets of items can be added to C, end;

Example: LR (0) items for the grammar

- $E' \rightarrow E$
- $E \rightarrow E + T/T$
- $T \rightarrow T * F/F$
- $F \rightarrow (E)/id$

is given below:

- $I_0: E' \rightarrow .E$
- $E \rightarrow .E + T$
- $E \rightarrow .T$
- $T \rightarrow .T * F$

$T \rightarrow .F$
 $F \rightarrow .(E)$
 $F \rightarrow .id$

I_1 : got (I_0, E)
 $E' \rightarrow E.$
 $E \rightarrow E. + T$

I_2 : goto (I_0, T)
 $E \rightarrow T.$
 $T \rightarrow T. * F$

I_3 : goto (I_0, F)
 $T \rightarrow F.$

I_4 : goto ($I_0, ($)
 $F \rightarrow .(E)$
 $E \rightarrow .E + T$
 $E \rightarrow .T$
 $E \rightarrow .T * F$
 $T \rightarrow .F$

$F \rightarrow .(E)$
 $F \rightarrow .id$

I_5 : goto (I_0, id)
 $F \rightarrow id.$

I_6 : got ($I_1, +$)
 $E \rightarrow E+ .T$
 $T \rightarrow .T * F$
 $T \rightarrow .F$
 $F \rightarrow .(E)$
 $F \rightarrow .id$

I_7 : goto ($I_2, *$)
 $T \rightarrow T* .F$
 $F \rightarrow .(E)$
 $F \rightarrow .id$

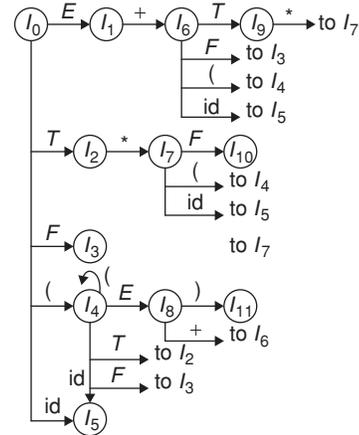
I_8 : goto (I_4, E)
 $F \rightarrow (E.)$

I_9 : goto (I_6, T)
 $E \rightarrow E+ T.$
 $T \rightarrow T. * F$

I_{10} : goto (I_7, F)
 $T \rightarrow T* F.$

I_{11} : goto ($I_8,)$)
 $F \rightarrow (E).$

For viable prefixes construct the DFA as follows:



SLR parsing table construction

1. Construct the canonical collection of sets of LR (0) items for G' .
2. Create the parsing action table as follows:
 - (a) If a is a terminal and $[A \rightarrow \alpha.a\beta]$ is in I_i , goto $(I_i, a) = I_j$ then action (i, a) to shift j . Here 'a' must be a terminal.
 - (b) If $[A \rightarrow \alpha.]$ is in I_i , then set action $[i, a]$ to 'reduce $A \rightarrow \alpha'$ for all a in FOLLOW (A);
 - (c) If $[S' \rightarrow S.]$ is in I_i then set action $[i, \$]$ to 'accept'.
3. Create the parsing goto table for all non-terminals A, if $goto(I_i, A) = I_j$ then $goto[i, A] = j$.
4. All entries not defined by steps 2 and 3 are made errors.
5. Initial state of the parser contains $S' \rightarrow S$.

The parsing table constructed using the above algorithm is known as SLR (1) table for G .

Note: Every SLR (1) grammar is unambiguous, but every unambiguous grammar is not a SLR grammar.

Example 6: Construct SLR parsing table for the following grammar:

1. $S \rightarrow L = R$
2. $S \rightarrow R$
3. $L \rightarrow * R$
4. $L \rightarrow id$
5. $R \rightarrow L$

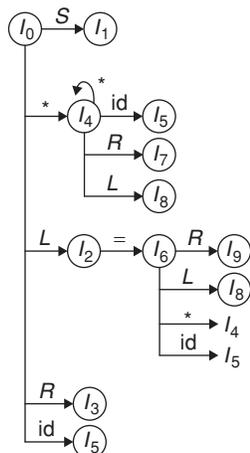
Solution: For the construction of SLR parsing table, add $S' \rightarrow S$ production.

$S' \rightarrow S$
 $S \rightarrow L = R$
 $S \rightarrow R$
 $L \rightarrow * R$
 $L \rightarrow id$
 $R \rightarrow L$

LR (0) items will be

- $I_0: S' \rightarrow .S$
- $S \rightarrow .L = R$
- $S \rightarrow .R$
- $L \rightarrow .*R$
- $L \rightarrow .id$
- $R \rightarrow .L$
- $I_1: \text{goto}(I_0, S)$
- $S' \rightarrow S.$
- $I_2: \text{goto}(I_0, L)$
- $S \rightarrow L. = R$
- $R \rightarrow L.$
- $I_3: \text{got}(I_0, R)$
- $S \rightarrow R.$
- $I_4: \text{goto}(I_0, *)$
- $L \rightarrow *.R$
- $R \rightarrow .L$
- $L \rightarrow .*R$
- $L \rightarrow .id$
- $I_5: \text{goto}(I_0, id)$
- $L \rightarrow id.$
- $I_6: \text{goto}(I_2, =)$
- $S \rightarrow L = .R$
- $R \rightarrow .L$
- $L \rightarrow .*R$
- $L \rightarrow .id$
- $I_7: \text{goto}(I_4, R)$
- $L \rightarrow *R.$
- $I_8: \text{goto}(I_4, L)$
- $R \rightarrow L.$
- $I_9: \text{goto}(I_6, R)$
- $S \rightarrow L = R.$

The DFA of LR(0) items will be



States	Action			Goto			
	=	*	id	\$	S	L	R
0		S_4	S_5		1	2	3
1				acc			
2	S_6, r_5			r_5			
3							
4		S_4	S_5			8	7
5							
6		S_4	S_5			8	9
7							
8							
9							

FOLLOW (S) = { \$ }

FOLLOW (L) = { = }

FOLLOW (R) = { \$, = }

For action [2, =] = S_6 and r_5

∴ Here we are getting shift – reduce conflict, so it is not SLR (1).

Canonical LR Parsing (CLR)

- To avoid some of invalid reductions, the states need to carry more information.
- Extra information input into a state by including a terminal symbol as a second component of an item.
- The general form of an item

$$[A \rightarrow \alpha.\beta, a]$$

Where $A \rightarrow \alpha\beta$ is a production.

a is terminal/right end marker (\$). We will call it as LR (1) item.

LR (1) item

It is a combination of LR (0) items along with look ahead of the item. Here 1 refers to look ahead of the item.

Construction of the sets of LR (1) items Function closure (I):

Begin

Repeat

For each item $[A \rightarrow \alpha.B\beta, a]$ in I,

Each production $B \rightarrow \cdot\gamma$ in G' ,

And each terminal b in FIRST (βa)

Such that $[B \rightarrow \cdot\gamma, b]$ is not in I do

Add $[B \rightarrow \cdot\gamma, b]$ to I;

End;

Until no more items can be added to I;

Example 7: Construct CLR parsing table for the following grammar:

$$S' \rightarrow S$$

$$S \rightarrow CC$$

$$C \rightarrow cC/d$$

6.16 | Unit 6 • Compiler Design

Solution: The initial set of items is

$$I_0: S' \rightarrow .S, \$$$

$$S \rightarrow .CC, \$$$

$$A \rightarrow \alpha.B\beta, a$$

Here $A = S, \alpha = \epsilon, B = C, \beta = C$ and $a = \$$

First (βa) is first $(C\$) = \text{first}(C) = \{c, d\}$

So, add items $[C \rightarrow .cC, c]$

$$[C \rightarrow .cC, d]$$

\therefore Our first set $I_0: S' \rightarrow .S, \$$

$$S \rightarrow .CC, \$$$

$$C \rightarrow .coca, c/d$$

$$C \rightarrow .d, c/d.$$

$I_1: \text{goto}(I_0, X)$ if $X = S$

$$S' \rightarrow S., \$$$

$I_2: \text{goto}(I_0, C)$

$$S \rightarrow C.C, \$$$

$$C \rightarrow .cC, \$$$

$$C \rightarrow .d, \$$$

$I_3: \text{goto}(I_0, c)$

$$C \rightarrow c.C, c/d$$

$$C \rightarrow .cC, c/d$$

$$C \rightarrow .d, c/d$$

$I_4: \text{goto}(I_0, d)$

$$C \rightarrow d., c/d$$

$I_5: \text{goto}(I_2, C)$

$$S \rightarrow CC., \$$$

$I_6: \text{goto}(I_2, c)$

$$C \rightarrow c.C; \$$$

$$C \rightarrow .cC, \$$$

$$C \rightarrow .d, \$$$

$I_7: \text{goto}(I_2, d)$

$$C \rightarrow d. \$$$

$I_8: \text{goto}(I_3, C)$

$$C \rightarrow cC., c/d$$

$I_9: \text{goto}(I_6, C)$

$$C \rightarrow cC., \$$$

CLR table is:

States	Action			Goto	
	c	1	\$	S	C
I_0	S_3	S_4		1	2
I_1			acc		
I_2	S_6	S_7			5
I_3	S_3	S_4			8
I_4	R_3	r_3			
I_5			r_1		
I_6	S_6	S_7			9
I_7			r_3		
I_8	R_2	r_2			
I_9			r_2		

Consider the string derivation 'dcd':

$$S \Rightarrow CC \Rightarrow CcC \Rightarrow Ccd \Rightarrow dcd$$

Stack	Input	Action
0	dcd\$	shift 4
0d4	Cd\$	reduce 3 i.e. $C \rightarrow d$
0C2	Cd\$	shift 6
0C2C6	D\$	shift 7
0C2C6d7	\$	reduce $C \rightarrow d$
0C2C6C9	\$	reduce $C \rightarrow cC$
0C2C5	\$	reduce $S \rightarrow CC$
0S1	\$	

Example 8: Construct CLR parsing table for the grammar:

$$S \rightarrow L = R$$

$$S \rightarrow R$$

$$L \rightarrow *R$$

$$L \rightarrow id$$

$$R \rightarrow L$$

Solution: The canonical set of items is

$$I_0: S' \rightarrow .S, \$$$

$$S \rightarrow .L = R, \$$$

$$S \rightarrow .R, \$$$

$$L \rightarrow .*R, = \quad [\text{first}(=R\$) = \{=\}]$$

$$L \rightarrow .id, =$$

$$R \rightarrow .L, \$$$

$I_1: \text{got}(I_0, S)$

$$S' \rightarrow S., \$$$

$I_2: \text{goto}(I_0, L)$

$$S \rightarrow L. = R, \$$$

$$R \rightarrow L., \$$$

$I_3: \text{goto}(I_0, R)$

$$S \rightarrow R., \$$$

$I_4: \text{got}(I_0, *)$

$$L \rightarrow *.R, =$$

$$R \rightarrow .L, =$$

$$L \rightarrow .*R, =$$

$$L \rightarrow .id, =$$

$I_5: \text{goto}(I_0, id)$

$$L \rightarrow id., =$$

$I_6: \text{goto}(I_7, L)$

$$R \rightarrow L., \$$$

$I_7: \text{goto}(I_2, =)$

$$S \rightarrow L = .R,$$

$$R \rightarrow .L, \$$$

$$L \rightarrow .*R, \$$$

$$L \rightarrow .id, \$$$

$I_8: \text{goto}(I_4, R)$

$$L \rightarrow *.R., =$$

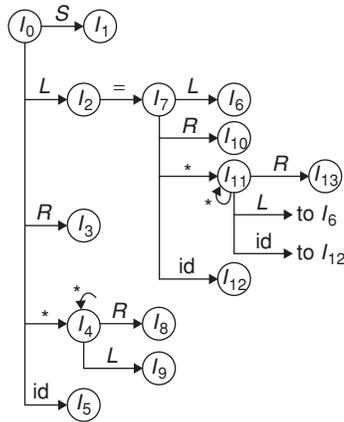
I_9 : goto (I_4 , L)
 $R \rightarrow L., =$

I_{10} : got (I_7 , R)
 $S \rightarrow L = R., \$$

I_{11} : goto (I_7 , $*$)
 $L \rightarrow *.R, \$$
 $R \rightarrow .L, \$$
 $L \rightarrow .*R, \$$
 $L \rightarrow .id, \$$

I_{12} : goto (I_7 , id)
 $L \rightarrow id. , \$$

I_{13} : goto (I_{11} , R)
 $L \rightarrow *R., \$$



We have to construct CLR parsing table based on the above diagram.

In this, we are going to have 13 states
 The shift-reduce conflict in the SLR parser is reduced here.

States	id	*	=	\$	S	L	R
0	S_5	S_4			1	2	3
1				acc			
2			S_7	r_5			
3				r_2			
4	S_5	S_4				9	8
5				r_4			
6				r_5			
7	s_{12}	s_{11}				6	10
8				r_3			
9				r_5			
10				r_1			
11	S_{12}	S_{11}					13
12				r_4			
13				r_3			

Stack	Input
0	$id = id\$$
0id5	$= id\$$
0L2	$= id\$$
0L2 = 7	$id\$$
0L2 = 7! d12	$\$$
0L2 = 7L6	$\$$
0L2 = 7R10	$\$$
0S1 (accept)	$\$$

Every SLR (1) grammar is LR (1) grammar.
 CLR (1) will have 'more number of states' than SLR Parser.

LALR Parsing Table

- The tables obtained by it are considerably smaller than the canonical LR table.
- LALR stands for Lookahead LR.
- The number of states in SLR and LALR parsing tables for a grammar G are equal.
- But LALR parsers recognize more grammars than SLR.
- YACC creates a LALR parser for the given grammar.
- YACC stands for 'Yet another Compiler'.
- An easy, but space-consuming LALR table construction is explained below:
 - Construct $C = \{I_0, I_1, \dots, I_n\}$, the collection of sets of LR (1) items.
 - Find all sets having the common core; replace these sets by their union
 - Let $C' = \{J_0, J_1, \dots, J_m\}$ be the resulting sets of LR (1) items. If there is a parsing action conflict then the grammar is not a LALR (1).
 - Let k be the union of all sets of items having the same core. Then $goto(J, X) = k$
- If there are no parsing action conflicts then the grammar is said to LALR (1) grammar.
- The collection of items constructed is called LALR (1) collection.

Example 9: Construct LALR parsing table for the following grammar:

$S' \rightarrow S$
 $S \rightarrow CC$
 $C \rightarrow cC/d$

Solution: We already got LR (1) items and CLR parsing table for this grammar.
 After merging I3 and I6 are replaced by I36.

I_{36} : $C \rightarrow c.C, c/d\$$
 $C \rightarrow .cC, c/d\$$
 $C \rightarrow .d, c/d\$$

I_{47} : By merging I_4 and I_7
 $C \rightarrow d. c/d/\$$

I_{89} : I_8 and I_9 are replaced by I_{89}
 $C \rightarrow cC., c/d/\$$

The LALR parsing table for this grammar is given below:

State	Action			goto	
	c	d	\$	S	C
0	S_{36}	S_{47}		1	2
1			acc		
2	S_{36}	S_{47}			5
36	S_{36}	S_{47}			89
47	r_3	R_3	r_3		
5			r_1		
89	r_2	r_2	r_2		

Example: Consider the grammar:

$S' \rightarrow S$
 $S \rightarrow aAd$
 $S \rightarrow bBd$
 $S \rightarrow aBe$
 $S \rightarrow bAe$
 $A \rightarrow c$
 $B \rightarrow c$

Which generates strings acd, bcd, ace and bce

LR (1) items are

I_0 : $S' \rightarrow .S, \$$
 $S \rightarrow .aAd, \$$
 $S \rightarrow .bBd, \$$
 $S \rightarrow .aBe, \$$
 $S \rightarrow .bAe, \$$

I_1 : goto (I_0, S)
 $S' \rightarrow S., \$$

I_2 : goto (I_0, a)
 $S \rightarrow a.Ad, c$
 $S \rightarrow a.Be, c$
 $A \rightarrow .c, d$
 $B \rightarrow .c, e$

I_3 : goto (I_0, b)
 $S \rightarrow b.Bd, c$
 $S \rightarrow b.Ae, c$
 $A \rightarrow .c, e$
 $B \rightarrow .c, e$

I_4 : goto (I_2, A)
 $S \rightarrow aA.d, c$

I_5 : goto (I_2, B)
 $S \rightarrow aB.e, c$

I_6 : goto (I_2, c)
 $A \rightarrow c., d$
 $B \rightarrow c., e$

I_7 : goto (I_3, c)
 $A \rightarrow c., e$
 $B \rightarrow c., d$

I_8 : goto (I_4, d)
 $S \rightarrow aAd., c$

I_9 : goto (I_5, e)
 $S \rightarrow aBe., c$

If we union I_6 and I_7
 $A \rightarrow c., d/e$
 $B \rightarrow c., d/e$

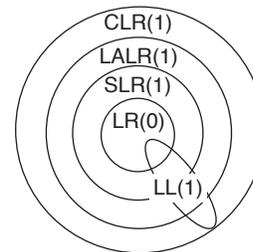
It generates reduce/reduce conflict.

Notes:

1. The merging of states with common cores can never produce a shift/reduce conflict, because shift action depends only on the core, not on the lookahead.
2. SLR and LALR tables for a grammar always have the same number of states (several hundreds) whereas CLR have thousands of states for the same grammar.

Comparison of parsing methods

Method	Item	Goto and Closures	Grammar it Applies to
SLR (1)	LR(0) item	Different from LR(1)	SLR (1) \subset LR(1)
LR (1)	LR(1) item		LR(1) – Largest class of LR grammars
LALR(1)	LR(1) item	Same as LR(1)	LALR(1) \subset LR(1)



Every LR (0) is SLR (1) but vice versa is not true.

Difference between SLR, LALR and CLR parsers

Differences among SLR, LALR and CLR are discussed below in terms of size, efficiency, time and space.

- (A) (i) only (B) (i) and
(C) (ii) and (iii) only (D) (i) and (iv) only

10. The FIRST and FOLLOW sets for the grammar:

$$S \rightarrow SS + /SS^*/a$$

- (A) First (S) = $\{a\}$
Follow (S) = $\{+, *, \$\}$
(B) First (S) = $\{+\}$
Follow (S) = $\{+, *, \$\}$
(C) First (S) = $\{a\}$
Follow (S) = $\{+, *\}$
(D) First (S) = $\{+, *\}$
Follow (S) = $\{+, *, \$\}$
11. A shift reduces parser carries out the actions specified within braces immediately after reducing with the corresponding rule of the grammar:
 $S \rightarrow xxW$ [print '1']
 $S \rightarrow y$ [print '2']
 $W \rightarrow Sz$ [print '3']
 What is the translation of 'x x x x y z z'?
 (A) 1231 (B) 1233
 (C) 2131 (D) 2321

12. After constructing the predictive parsing table for the following grammar:

$$\begin{aligned} Z &\rightarrow d \\ Z &\rightarrow XYZ \\ Y &\rightarrow c/\epsilon \\ X &\rightarrow Y \\ X &\rightarrow a \end{aligned}$$

The entry/entries for $[Z, d]$ is/are

- (A) $Z \rightarrow d$
(B) $Z \rightarrow XYZ$
(C) Both (A) and (B)
(D) $X \rightarrow Y$

13. The following grammar is

$$S \rightarrow AaAb/BbBa$$

$$A \rightarrow \epsilon$$

$$B \rightarrow \epsilon$$

- (A) LL (1) (B) Not LL (1)
(C) Recursive (D) Ambiguous

14. Compute the FIRST (P) for the below grammar:

$$P \rightarrow AQRbe/mn/DE$$

$$A \rightarrow ab/\epsilon$$

$$Q \rightarrow q_1q_2/\epsilon$$

$$R \rightarrow r_1r_2/\epsilon$$

$$D \rightarrow d$$

$$E \rightarrow e$$

- (A) $\{m, a\}$ (B) $\{m, a, q_1, r_1, b, d\}$
(C) $\{d, e\}$ (D) $\{m, n, a, b, d, e, q_1, r_1\}$

15. After constructing the LR(1) parsing table for the augmented grammar

$$S' \rightarrow S$$

$$S \rightarrow BB$$

$$B \rightarrow aB/c$$

What will be the action $[I_3, a]$?

- (A) Accept (B) S_7
(C) r_2 (D) S_5

Practice Problems 2

Directions for questions 1 to 19: Select the correct alternative from the given choices.

1. Consider the grammar

$$S \rightarrow aSb$$

$$S \rightarrow aS$$

$$S \rightarrow \epsilon$$

This grammar is ambiguous by generating which of the following string.

- (A) aa (B) \in
(C) aaa (D) aab
2. To convert the grammar $E \rightarrow E + T$ into LL grammar
 (A) use left factor
 (B) CNF form
 (C) eliminate left recursion
 (D) Both (B) and (C)
3. Given the following expressions of a grammar
 $E \rightarrow E \times F/F + E/F$
 $F \rightarrow F? F/id$
 Which of the following is true?
 (A) \times has higher precedence than $+$
 (B) $?$ has higher precedence than \times

- (C) $+$ and $?$ have same precedence
 (D) $+$ has higher precedence than $*$

4. The action of parsing the source program into the proper syntactic classes is known as

- (A) Lexical analysis
(B) Syntax analysis
(C) Interpretation analysis
(D) Parsing

5. Which of the following is not a bottom up parser?

- (A) LALR (B) Predictive parser
(C) CLR (D) SLR

6. A system program that combines separately compiled modules of a program into a form suitable for execution is

- (A) Assembler.
(B) Linking loader.
(C) Cross compiler.
(D) None of these.

7. Resolution of externally defined symbols is performed by a

- (A) Linker (B) Loader.
(C) Compiler. (D) Interpreter.

8. LR parsers are attractive because
- They can be constructed to recognize CFG corresponding to almost all programming constructs.
 - There is no need of backtracking.
 - Both (A) and (B).
 - None of these
9. YACC builds up
- SLR parsing table
 - Canonical LR parsing table
 - LALR parsing table
 - None of these
10. Language which have many types, but the type of every name and expression must be calculated at compile time are
- Strongly typed languages
 - Weakly typed languages
 - Loosely typed languages
 - None of these
11. Consider the grammar shown below:
- $$S \rightarrow iEtSS'/a/b$$
- $$S' \rightarrow eS/\epsilon$$
- In the predictive parse table M , of this grammar, the entries $M[S', e]$ and $M[S', \$]$ respectively are
- $\{S' \rightarrow eS\}$ and $\{S' \rightarrow \epsilon\}$
 - $\{S' \rightarrow eS\}$ and $\{\}$
 - $\{S' \rightarrow \epsilon\}$ and $\{S' \rightarrow \epsilon\}$
 - $\{S' \rightarrow eS, S' \rightarrow \epsilon\}$ and $\{S' \rightarrow \epsilon\}$
12. Consider the grammar $S \rightarrow CC, C \rightarrow cC/d$. The grammar is
- LL (1)
 - SLR (1) but not LL (1)
 - LALR (1) but not SLR (1)
 - LR (1) but not LALR (1)
13. Consider the grammar
- $$E \rightarrow E + n/E - n/n$$
- For a sentence $n + n - n$, the handles in the right sentential form of the reduction are
- $n, E + n$ and $E + n - n$
 - $n, E + n$ and $E + E - n$
 - $n, n + n$ and $n + n - n$
 - $n, E + n$ and $E - n$
14. A top down parser uses ____ derivation.
- Left most derivation
 - Left most derivation in reverse
 - Right most derivation
 - Right most derivation in reverse
15. Which of the following statement is false?
- An unambiguous grammar has single leftmost derivation.
 - An LL (1) parser is topdown.
 - LALR is more powerful than SLR.
 - An ambiguous grammar can never be LR (K) for any k .
16. Merging states with a common core may produce ____ conflicts in an LALR parser.
- Reduce – reduce
 - Shift – reduce
 - Both (A) and (B)
 - None of these
17. LL (K) grammar
- Has to be CFG
 - Has to be unambiguous
 - Cannot have left recursion
 - All of these
18. The I_0 state of the LR (0) items for the grammar
- $$S \rightarrow AS/b$$
- $$A \rightarrow SA/a.$$
- $S' \rightarrow .S$
 $S \rightarrow .AS$
 $S \rightarrow .b$
 $A \rightarrow .SA$
 $A \rightarrow .a$
 - $S \rightarrow .AS$
 $S \rightarrow .b$
 $A \rightarrow .SA$
 $A \rightarrow .a$
 - $S \rightarrow .AS$
 $S \rightarrow .b$
 - $S \rightarrow A$
 $A \rightarrow .SA$
 $A \rightarrow .a$
19. In the predictive parsing table for the grammar:
- $$S \rightarrow FR$$
- $$R \rightarrow \times S/e$$
- $$F \rightarrow id$$
- What will be the entry for $[S, id]$?
- $S \rightarrow FR$
 - $F \rightarrow id$
 - Both (A) and (B)
 - None of these

PREVIOUS YEARS' QUESTIONS

1. Consider the grammar:
 $S \rightarrow (S) | a$
 Let the number of states in SLR(1), LR(1) and LALR(1) parsers for the grammar be n_1 , n_2 and n_3 respectively. The following relationship holds good: [2005]
 (A) $n_1 < n_2 < n_3$ (B) $n_1 = n_3 < n_2$
 (C) $n_1 = n_2 = n_3$ (D) $n_1 \geq n_2 \geq n_3$
2. Consider the following grammar:
 $S \rightarrow S * E$
 $S \rightarrow E$
 $E \rightarrow F + E$
 $E \rightarrow F$
 $F \rightarrow id$
 Consider the following LR (0) items corresponding to the grammar above.
 (i) $S \rightarrow S * .E$
 (ii) $E \rightarrow F . + E$
 (iii) $E \rightarrow F + .E$
 Given the items above, which two of them will appear in the same set in the canonical sets-of items for the grammar? [2006]
 (A) (i) and (ii) (B) (ii) and (iii)
 (C) (i) and (iii) (D) None of the above
3. Consider the following statements about the context-free grammar
 $G = \{S \rightarrow SS, S \rightarrow ab, S \rightarrow ba, S \rightarrow \epsilon\}$
 (i) G is ambiguous
 (ii) G produces all strings with equal number of a's and b's
 (iii) G can be accepted by a deterministic PDA.
 Which combination below expresses all the true statements about G ? [2006]
 (A) (i) only (B) (i) and (iii) only
 (C) (ii) and (iii) only (D) (i), (ii) and (iii)
4. Consider the following grammar:
 $S \rightarrow FR$
 $R \rightarrow *S | \epsilon$
 $F \rightarrow id$
 In the predictive parser table, M , of the grammar the entries $M[S, id]$ and $M[R, \$]$ respectively. [2006]
 (A) $\{S \rightarrow FR\}$ and $\{R \rightarrow \epsilon\}$
 (B) $\{S \rightarrow FR\}$ and $\{\}$
 (C) $\{S \rightarrow FR\}$ and $\{R \rightarrow *S\}$
 (D) $\{F \rightarrow id\}$ and $\{R \rightarrow \epsilon\}$
5. Which one of the following grammars generates the language $L = \{a^i b^j | i \neq j\}$? [2006]
 (A) $S \rightarrow AC | CB$ (B) $S \rightarrow aS | Sb | a | b$
 $C \rightarrow aCb | a | b$
 $A \rightarrow aA | \epsilon$
 $B \rightarrow Bb | \epsilon$
 (C) $S \rightarrow AC | CB$ (D) $S \rightarrow AC | CB$
 $C \rightarrow aCb | \epsilon$
 $A \rightarrow aA | \epsilon$
 $B \rightarrow Bb | \epsilon$
6. In the correct grammar above, what is the length of the derivation (number of steps starting from S) to generate the string $a^\ell b^m$ with $\ell \neq m$? [2006]
 (A) $\max(1, m) + 2$
 (B) $1 + m + 2$
 (C) $1 + m + 3$
 (D) $\max(1, m) + 3$
7. Which of the following problems is undecidable? [2007]
 (A) Membership problem for CFGs.
 (B) Ambiguity problem for CFGs.
 (C) Finiteness problem for FSAs.
 (D) Equivalence problem for FSAs.
8. Which one of the following is a top-down parser? [2007]
 (A) Recursive descent parser.
 (B) Operator precedence parser.
 (C) An LR (k) parser.
 (D) An LALR (k) parser.
9. Consider the grammar with non-terminals $N = \{S, C, \text{ and } S_1\}$, terminals $T = \{a, b, i, t, e\}$ with S as the start symbol, and the following set of rules: [2007]
 $S \rightarrow iCtSS_1/a$
 $S_1 \rightarrow eS/\epsilon$
 $C \rightarrow b$
 The grammar is NOT LL (1) because:
 (A) It is left recursive
 (B) It is right recursive
 (C) It is ambiguous
 (D) It is not context-free.
10. Consider the following two statements:
 P: Every regular grammar is LL (1)
 Q: Every regular set has a LR (1) grammar
 Which of the following is TRUE? [2007]
 (A) Both P and Q are true
 (B) P is true and Q is false
 (C) P is false and Q is true
 (D) Both P and Q are false
- Common data for questions 11 and 12:** Consider the CFG with $\{S, A, B\}$ as the non-terminal alphabet, $\{a, b\}$ as the terminal alphabet, S as the start symbol and the following set of production rules:
 $S \rightarrow aB$ $S \rightarrow bA$
 $B \rightarrow b$ $A \rightarrow a$
 $B \rightarrow bS$ $A \rightarrow aS$
 $B \rightarrow aBB$ $S \rightarrow bAA$

11. Which of the following strings is generated by the grammar? [2007]
 (A) $aaaabb$ (B) $aabbbb$
 (C) $aabbab$ (D) $abbbba$
12. For the correct answer strings to Q.78, how many derivation trees are there? [2007]
 (A) 1 (B) 2
 (C) 3 (D) 4
13. Which of the following describes a handle (as applicable to LR-parsing) appropriately? [2008]
 (A) It is the position in a sentential form where the next shift or reduce operation will occur.
 (B) It is non-terminal whose production will be used for reduction in the next step.
 (C) It is a production that may be used for reduction in a future step along with a position in the sentential form where the next shift or reduce operation will occur.
 (D) It is the production p that will be used for reduction in the next step along with a position in the sentential form where the right hand side of the production may be found.
14. Which of the following statements are true?
 (i) Every left-recursive grammar can be converted to a right-recursive grammar and vice-versa
 (ii) All ϵ -productions can be removed from any context-free grammar by suitable transformations
 (iii) The language generated by a context-free grammar all of whose productions are of the form $X \rightarrow w$ or $X \rightarrow wY$ (where, w is a string of terminals and Y is a non-terminal), is always regular
 (iv) The derivation trees of strings generated by a context-free grammar in Chomsky Normal Form are always binary trees [2008]
 (A) (i), (ii), (iii) and (iv) (B) (ii), (iii) and (iv) only
 (C) (i), (iii) and (iv) only (D) (i), (ii) and (iv) only
15. An LALR (1) parser for a grammar G can have shift-reduce (S - R) conflicts if and only if [2008]
 (A) The SLR (1) parser for G has S - R conflicts
 (B) The LR (1) parser for G has S - R conflicts
 (C) The LR (0) parser for G has S - R conflicts
 (D) The LALR (1) parser for G has reduce-reduce conflicts
16. $S \rightarrow aSa|bSb|a|b$;
 The language generated by the above grammar over the alphabet $\{a, b\}$ is the set of [2009]
 (A) All palindromes.
 (B) All odd length palindromes.
 (C) Strings that begin and end with the same symbol.
 (D) All even length palindromes.
17. Which data structure in a compiler is used for managing information about variables and their attributes? [2010]

- (A) Abstract syntax tree
 (B) Symbol table
 (C) Semantic stack
 (D) Parse table

18. The grammar $S \rightarrow aSa|bS|c$ is [2010]
 (A) LL (1) but not LR (1)
 (B) LR (1) but not LL (1)
 (C) Both LL (1) and LR (1)
 (D) Neither LL (1) nor LR (1)
19. The lexical analysis for a modern computer language such as Java needs the power of which one of the following machine models in a necessary and sufficient sense? [2011]
 (A) Finite state automata
 (B) Deterministic pushdown automata
 (C) Non-deterministic pushdown automata
 (D) Turing machine

Common data for questions 20 and 21: For the grammar below, a partial LL (1) parsing table is also presented along with the grammar. Entries that need to be filled are indicated as E_1 , E_2 , and E_3 . $\$$ indicates end of input, and, $|$ separates alternate right hand side of productions

$S \rightarrow a A b B|b A a B|e$

$A \rightarrow S$

$B \rightarrow S$

	a	b	$\$$
S	E_1	E_2	$S \rightarrow e$
A	$A \rightarrow S$	$A \rightarrow S$	Error
B	$B \rightarrow S$	$B \rightarrow S$	E_3

20. The FIRST and FOLLOW sets for the non-terminals A and B are [2012]
 (A) $\text{FIRST}(A) = \{a, b, \epsilon\} = \text{FIRST}(B)$
 $\text{FOLLOW}(A) = \{a, b\}$
 $\text{FOLLOW}(B) = \{a, b, \$\}$
 (B) $\text{FIRST}(A) = \{a, b, \$\}$
 $\text{FIRST}(B) = \{a, b, \epsilon\}$
 $\text{FOLLOW}(A) = \{a, b\}$
 $\text{FOLLOW}(B) = \{\$\}$
 (C) $\text{FIRST}(A) = \{a, b, \epsilon\} = \text{FIRST}(B)$
 $\text{FOLLOW}(A) = \{a, b\}$
 $\text{FOLLOW}(B) = \emptyset$
 (D) $\text{FIRST}(A) = \{a, b\} = \text{FIRST}(B)$
 $\text{FOLLOW}(A) = \{a, b\}$
 $\text{FOLLOW}(B) = \{a, b\}$
21. The appropriate entries for E_1 , E_2 , and E_3 are [2012]
 (A) $E_1: S \rightarrow a A b B, A \rightarrow S$
 $E_2: S \rightarrow b A a B, B \rightarrow S$
 $E_3: B \rightarrow S$

- (B) $E_1: S \rightarrow a A b B, S \rightarrow \epsilon$
 $E_2: S \rightarrow b A a B, S \rightarrow \epsilon$
 $E_3: S \rightarrow \epsilon$
- (C) $E_1: S \rightarrow a A b B, S \rightarrow \epsilon$
 $E_2: S \rightarrow b A a B, S \rightarrow \epsilon$
 $E_3: B \rightarrow S$
- (D) $E_1: A \rightarrow S, S \rightarrow \epsilon$
 $E_2: B \rightarrow S, S \rightarrow \epsilon$
 $E_3: B \rightarrow S$
22. What is the maximum number of reduce moves that can be taken by a bottom-up parser for a grammar with no epsilon-and unit-production (i.e., of type $A \rightarrow \epsilon$ and $A \rightarrow a$) to parse a string with n tokens? [2013]
- (A) $n/2$ (B) $n - 1$
 (C) $2n - 1$ (D) 2^n
23. Which of the following is/are undecidable?
- (i) G is a CFG. Is $L(G) = \phi$?
 (ii) G is a CFG, Is $L(G) = \Sigma^*$?
 (iii) M is a Turing machine. Is $L(M)$ regular?
 (iv) A is a DFA and N is an NFA. Is $L(A) = L(N)$? [2013]
- (A) (iii) only
 (B) (iii) and (iv) only
 (C) (i), (ii) and (iii) only
 (D) (ii) and (iii) only
24. Consider the following two sets of LR (1) items of an LR (1) grammar. [2013]
- | | |
|--------------------------|-------------------------|
| $X \rightarrow c.X, c/d$ | $X \rightarrow c.X, \$$ |
| $X \rightarrow .cX, c/d$ | $X \rightarrow .cX, \$$ |
| $X \rightarrow .d, c/d$ | $X \rightarrow .d, \$$ |
- Which of the following statements related to merging of the two sets in the corresponding LALR parser is/are FALSE?
- (i) Cannot be merged since look - ahead are different.
 (ii) Can be merged but will result in S-R conflict.
 (iii) Can be merged but will result in R-R conflict.
 (iv) Cannot be merged since goto on c will lead to two different sets.
- (A) (i) only (B) (ii) only
 (C) (i) and (iv) only (D) (i), (ii), (iii) and (iv)
25. A canonical set of items is given below
- $S \rightarrow L. > R$
 $Q \rightarrow R.$
- On input symbol $<$ the sset has [2014]
- (A) A shift–reduce conflict and a reduce–reduce conflict.
 (B) A shift–reduce conflict but not a reduce–reduce conflict.
 (C) A reduce–reduce conflict but not a shift reduce conflict.
 (D) Neither a shift–reduce nor a reduce–reduce conflict.
26. Consider the grammar defined by the following production rules, with two operators $*$ and $+$
- $S \rightarrow T * P$
 $T \rightarrow U|T * U$
 $P \rightarrow Q + P|Q$
 $Q \rightarrow Id$
 $U \rightarrow Id$
- Which one of the following is TRUE? [2014]
- (A) $+$ is left associative, while $*$ is right associative
 (B) $+$ is right associative, while $*$ is left associative
 (C) Both $+$ and $*$ are right associative.
 (D) Both $+$ and $*$ are left associative
27. Which one of the following problems is undecidable? [2014]
- (A) Deciding if a given context -free grammar is ambiguous.
 (B) Deciding if a given string is generated by a given context-free grammar.
 (C) Deciding if the language generated by a given context-free grammar is empty.
 (D) Deciding if the language generated by a given context free grammar is finite.
28. Which one of the following is TRUE at any valid state in shift-reduce parsing? [2015]
- (A) Viable prefixes appear only at the bottom of the stack and not inside.
 (B) Viable prefixes appear only at the top of the stack and not inside.
 (C) The stack contains only a set of viable prefixes.
 (D) The stack never contains viable prefixes.
29. Among simple LR (SLR), canonical LR, and look-ahead LR (LALR), which of the following pairs identify the method that is very easy to implement and the method that is the most powerful, in that order? [2015]
- (A) SLR, LALR
 (B) Canonical LR, LALR
 (C) SLR, canonical LR
 (D) LALR, canonical LR
30. Consider the following grammar G
- $S \rightarrow F | H$
 $F \rightarrow p | c$
 $H \rightarrow d | c$
- Where S, F and H are non-terminal symbols, p, d and c are terminal symbols. Which of the following statement(s) is/are correct? [2015]
- S_1 . LL(1) can parse all strings that are generated using grammar G
 S_2 . LR(1) can parse all strings that are generated using grammar G

- (A) Only S_1 (B) Only S_2
 (C) Both S_1 and S_2 (D) Neither S_1 nor S_2
31. Match the following: [2016]
- | | |
|--------------------------|---------------------------|
| (P) Lexical analysis | (i) Leftmost derivation |
| (Q) Top down parsing | (ii) Type checking |
| (R) Semantic analysis | (iii) Regular expressions |
| (S) Runtime environments | (iv) Activation records |
- (A) $P \leftrightarrow i, Q \leftrightarrow ii, R \leftrightarrow iv, S \leftrightarrow iii$
 (B) $P \leftrightarrow iii, Q \leftrightarrow i, R \leftrightarrow ii, S \leftrightarrow iv$
 (C) $P \leftrightarrow ii, Q \leftrightarrow iii, R \leftrightarrow i, S \leftrightarrow iv$
 (D) $P \leftrightarrow iv, Q \leftrightarrow i, R \leftrightarrow ii, S \leftrightarrow iii$

32. A student wrote two context - free grammars **G1** and **G2** for generating a single C-like array declaration. The dimension of the array is at least one.

For example, `int a [10] [3];`

The grammars use D as the start symbol, and use six terminal symbols `int; id [] num.` [2016]

Grammar G1	Grammar G2
$D \rightarrow \text{int } L;$	$D \rightarrow \text{int } L;$
$L \rightarrow \text{id } [E$	$L \rightarrow \text{id } E$
$E \rightarrow \text{num }]$	$E \rightarrow E [\text{num}]$
$E \rightarrow \text{num }] [E$	$E \rightarrow [\text{num}]$

Which of the grammars correctly generate the declaration mentioned above?

- (A) Both **G1** and **G2**
 (B) Only **G1**
 (C) Only **G2**
 (D) Neither **G1** nor **G2**
33. Consider the following grammar:

$P \rightarrow xQRS$ $Q \rightarrow yz \mid z$ $R \rightarrow w \mid \epsilon$ $S \rightarrow y$

What is FOLLOW (Q)? [2017]

- (A) $\{R\}$ (B) $\{w\}$
 (C) $\{w, y\}$ (D) $\{w, \$\}$
34. Which of the following statements about parser is/are CORRECT? [2017]
- Canonical LR is more powerful than SLR.
 - SLR is more powerful than LALR.
 - SLR is more powerful than Canonical LR.
- (A) I only (B) II only
 (C) III only (D) I and III only

35. Consider the following expression grammar G :

$E \rightarrow E - T \mid T$
 $T \rightarrow T + F \mid F$
 $F \rightarrow (E) \mid \text{id}$

Which of the following grammars is not left recursive, but is equivalent to G? [2017]

- (A) $E \rightarrow E - T \mid T$ (B) $E \rightarrow TE$
 $T \rightarrow T + F \mid F$ $E' \rightarrow -TE \mid \epsilon$
 $F \rightarrow (E) \mid \text{id}$ $T \rightarrow T + F \mid F$
 $F \rightarrow (E) \mid \text{id}$
- (C) $E \rightarrow TX$ (D) $E \rightarrow TX \mid (TX)$
 $X \rightarrow -TX \mid \epsilon$ $X \rightarrow -TX \mid +TX \mid \epsilon$
 $T \rightarrow FY$ $T \rightarrow \text{id}$
 $Y \rightarrow +FY \mid \epsilon$
 $F \rightarrow (E) \mid \text{id}$

36. Which one of the following statements is FALSE? [2018]

- (A) Context-free grammar can be used to specify both lexical and syntax rules.
 (B) Type checking is done before parsing.
 (C) High-level language programs can be translated to different Intermediate Representations.
 (D) Arguments to a function can be passed using the program stack.

37. A lexical analyzer uses the following patterns to recognize three tokens $T_1, T_2,$ and T_3 over the alphabet $\{a, b, c\}$.

$T_1: a?(b|c)*a$

$T_2: b?(a|c)*b$

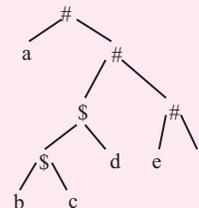
$T_3: c?(b|a)*c$

Note that 'x?' means 0 or 1 occurrence of the symbol x. Note also that the analyzer outputs the token that matches the longest possible prefix.

If the string `baaacabc` is processed by the analyzer, which one of the following is the sequence of tokens it outputs? [2018]

- (A) $T_1T_2T_3$ (B) $T_1T_1T_3$
 (C) $T_2T_1T_3$ (D) T_3T_3

38. Consider the following parse tree for the expression `a#b$cd#e#f`, involving two binary operators \$ and #.



Which one of the following is correct for the given parse tree? [2018]

- (A) \$ has higher precedence and is left associative; # is right associative
 (B) # has higher precedence and is left associative; \$ is right associative
 (C) \$ has higher precedence and is left associative; # is left associative
 (D) # has higher precedence and is right associative; \$ is left associative

ANSWER KEYS**EXERCISES****Practice Problems 1**

1. D 2. A 3. B 4. A 5. D 6. C 7. C 8. D 9. C 10. A
11. C 12. C 13. A 14. B 15. D

Practice Problems 2

1. D 2. C 3. B 4. A 5. B 6. B 7. A 8. C 9. C 10. A
11. D 12. A 13. D 14. A 15. D 16. A 17. C 18. A 19. A

Previous Years' Questions

1. B 2. D 3. B 4. A 5. D 6. A 7. B 8. A 9. C 10. A
11. C 12. B 13. D 14. C 15. B 16. B 17. B 18. C 19. A 20. A
21. C 22. B 23. D 24. D 25. D 26. B 27. A 28. C 29. C 30. D
31. B 32. A 33. C 34. A 35. C 36. B 37. D 38. A

Chapter 2

Syntax Directed Translation

LEARNING OBJECTIVES

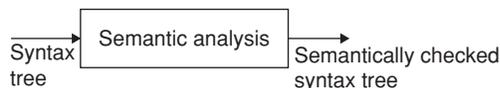
- ☞ Syntax directed translation
- ☞ Syntax directed definition
- ☞ Dependency graph
- ☞ Constructing syntax trees for expressions
- ☞ Types of SDD's
- ☞ S-attributed definition
- ☞ L-attributed definitions
- ☞ Synthesized attributes on the parser
- ☞ Syntax directed translation schemes
- ☞ Bottom up evaluation of inherited attributes

SYNTAX DIRECTED TRANSLATION

To translate a programming language construct, a compiler may need to know the type of construct, the location of the first instruction, and the number of instructions generated... etc. So, we have to use the term 'attributes' associated with constructs.

An attribute may represent type, number of arguments, memory location, compatibility of variables used in a statement which cannot be represented by CFG alone.

So, we need to have one more phase to do this, i.e., 'semantic analysis' phase.



In this phase, for each production CFG, we will give some semantic rule.

Syntax directed translation scheme

A CFG in which a program fragment called output action (semantic action or semantic rule) is associated with each production is known as Syntax Directed Translation Scheme.

These semantic rules are used to

1. Generate intermediate code.
2. Put information into symbol table.
3. Perform type checking.
4. Issues error messages.

Notes:

1. Grammar symbols are associated with attributes.
2. Values of the attributes are evaluated by the semantic rules associated with production rules.

Notations for Associating Semantic Rules

There are two techniques to associate semantic rules:

Syntax directed definition (SDD) It is high level specification for translation. They hide the implementation details, i.e., the order in which translation takes place.

Attributes + CFG + Semantic rules = Syntax directed definition (SDD).

Translation schemes These schemes indicate the order in which semantic rules are to be evaluated. This is an input and output mapping.

SYNTAX DIRECTED DEFINITIONS

A SDD is a generalization of a CFG in which each grammar symbol is associated with a set of attributes.

There are two types of set of attributes for a grammar symbol.

1. Synthesized attributes
2. Inherited attributes

Each production rule is associated with a set of semantic rules.

Semantic rules setup dependencies between attributes which can be represented by a dependency graph.

The dependency graph determines the evaluation order of these semantic rules.

Evaluation of a semantic rule defines the value of an attribute. But a semantic rule may also have some side effects such as printing a value.

Attribute grammar: An attribute grammar is a syntax directed definition in which the functions in semantic rules ‘cannot have side effects’.

Annotated parse tree: A parse tree showing the values of attributes at each node is called an annotated parse tree.

The process of computing the attribute values at the nodes is called annotating (or decorating) of the parse tree.

In a SDD, each production $A \rightarrow \alpha$ is associated with a set of semantic rules of the form:

$$b = f(c_1, c_2, \dots, c_n) \text{ where}$$

f : A function

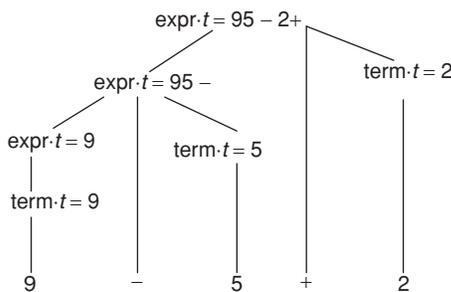
b can be one of the following:

b is a ‘synthesized attribute’ of A and c_1, c_2, \dots, c_n are attributes of the grammar symbols in $A \rightarrow \alpha$.

The value of a ‘synthesized attribute’ at a node is computed from the value of attributes at the children of that node in the parse tree.

Example:

Production	Semantic Rule
$\text{expr} \rightarrow \text{expr1} + \text{term}$	$\text{expr.t} = \text{expr1.t} \text{term.t} '+'$
$\text{expr} \rightarrow \text{expr1} - \text{term}$	$\text{expr.t} = \text{expr1.t} \text{term.t} '-'$
$\text{expr} \rightarrow \text{term}$	$\text{expr.t} = \text{term.t}$
$\text{term} \rightarrow 0$	$\text{term.t} = '0'$
$\text{term} \rightarrow 1$	$\text{term.t} = '1'$
\vdots	\vdots
$\text{term} \rightarrow 9$	$\text{term.t} = '9'$



b is an ‘inherited attribute’ of one of the grammar symbols on the right side of the production.

An ‘inherited attribute’ is one whose value at a node is defined in terms of attributes at the parent and/or siblings of that node. It is used for finding the context in which it appears.

Example: An inherited attribute distributes type information to the various identifiers in a declaration.

For the grammar

- $D \rightarrow TL$
- $T \rightarrow \text{int}$
- $T \rightarrow \text{real}$
- $L \rightarrow L_1, \text{id}$
- $L \rightarrow \text{id}$

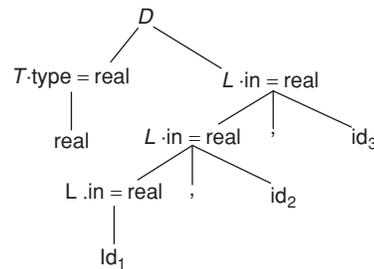
That is, The keyword int or real followed by a list of identifiers.

In this T has synthesized attribute type: $T.type$. L has an inherited attribute in $L.in$

Rules associated with L call for procedure add type to the type of each identifier to its entry in the symbol table.

Production	Semantic Rule
$D \rightarrow TL$	$L.in = T.type$
$T \rightarrow \text{int}$	$T.type = \text{integer}$
$T \rightarrow \text{real}$	$T.type = \text{real}$
$L \rightarrow L_1, \text{id}$	$\text{addtype } L_1.in = L.in(\text{id.entry}, L.in)$
$L \rightarrow \text{id}$	$\text{addtype } (\text{id.entry}, L.in)$

The annotated parse tree for the sentence $\text{real id}_1, \text{id}_2, \text{id}_3$ is shown below:



SYNTHESIZED ATTRIBUTE

The value of a synthesized attribute at a node is computed from the value of attributes at the children of that node in a parse tree. Consider the following grammar:

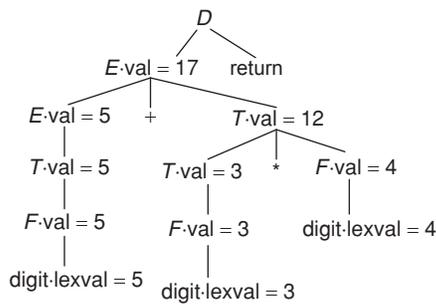
- $L \rightarrow E_n$
- $E \rightarrow E_1 + T$
- $E \rightarrow T$
- $T \rightarrow T_1 * F$
- $T \rightarrow F$
- $F \rightarrow (E)$
- $F \rightarrow \text{digit}$.

Let us consider synthesized attribute value with each of the non-terminals E , T and F .

Token digit has a synthesized attribute lexical supplied by lexical analyzer.

Production	Semantic Rule
$L \rightarrow E_n$	print (E.val)
$E \rightarrow E_1 + T$	$E.val := E_1.val + T.val$
$E \rightarrow T$	$E.val := T_1.val$
$T \rightarrow T_1 * F$	$T.val := T_1.val * F.val$
$T \rightarrow F$	$T.val := F.val$
$F \rightarrow (E)$	$F.val := E.val$
$F \rightarrow \text{digit}$	$F.val := \text{digit.lexval}$

The Annotated parse tree for the expression $5 + 3 * 4$ is shown below:



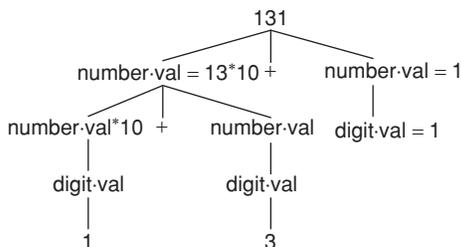
Example 1: Consider an example, which shows semantic rules for Infix to postfix translation:

Production	Semantic Rules
$\text{expr} \rightarrow \text{expr1} + \text{term}$	$\text{expr.t} := \text{expr1.t} \text{term.t} '+'$
$\text{expr} \rightarrow \text{expr1} - \text{term}$	$\text{expr.t} := \text{expr1.t} \text{term.t} '-'$
$\text{expr} \rightarrow \text{term}$	$\text{expr.t} := \text{term.t}$
$\text{term} \rightarrow 0$	$\text{term.t} := '0'$
:	:
$\text{term} \rightarrow 9$	$\text{term.t} := '9'$

Example 2: Write a SDD for the following grammar to determine number.val.

$$\begin{aligned} \text{number} &\rightarrow \text{number digit} \\ \text{digit} &\rightarrow 0|1|\dots|9 \end{aligned} \quad \left\{ \begin{array}{l} \text{digit.val} := '0' \\ \text{digit.val} := '1' \\ \vdots \\ \text{digit.val} := '9' \end{array} \right.$$

number.val := number.val * 10 + digit.val
Annotated tree for 131 is

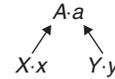


DEPENDENCY GRAPH

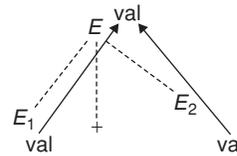
The interdependencies among the attributes at the nodes in a parse tree can be depicted by a directed graph called dependency graph.

- Synthesized attributes have edges pointing upwards.
- Inherited attributes have edges pointing downwards and/or sidewise.

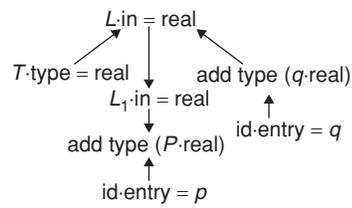
Example 1: $A.a := f(X.x, Y.y)$ is a semantic rule for $A \rightarrow XY$. For each semantic rule that consists of a procedure call:



Example 2:



Example 3: real p, q;



Evaluation order

A topological sort of directed acyclic graph is an ordering m_1, m_2, \dots, m_k of nodes of the graph S . t edges go from nodes earlier in the ordering to later nodes.

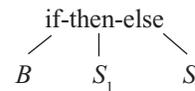
$m_i \rightarrow m_j$ means m_i appears before m_j in the ordering.

If $b := f(c_1, c_2, \dots, c_k)$, the dependent attributes c_1, c_2, \dots, c_k are available at node before f is evaluated.

Abstract syntax tree

It is a condensed form of parse tree useful for representing language constructs.

Example



CONSTRUCTING SYNTAX TREES FOR EXPRESSIONS

Each node in a syntax tree can be implemented as a record with several fields.

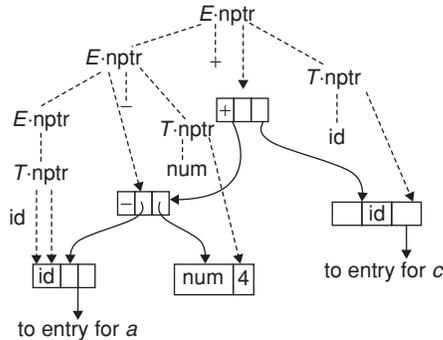
In the node for an operator, one field identifies the operator and the remaining fields contain pointers to the nodes for the operands.

1. mknode (op, left, right)
2. mkleaf (id, entry). Entry is a pointer to symbol table.
3. mkleaf (num, val)

Example:

Production	Semantic Rules
$E \rightarrow E_1 + T$	$E.nptr := \text{mknode} ('+', E_1.nptr, T.nptr)$
$E \rightarrow E_1 - T$	$E.nptr := \text{mknode} ('-', E_1.nptr, T.nptr)$
$E \rightarrow T$	$E.nptr := T.nptr$
$T \rightarrow (E)$	$T.nptr := E.nptr$
$T \rightarrow \text{id}$	$T.nptr := \text{mkleaf}(\text{id}, \text{id.entry})$
$T \rightarrow \text{num}$	$T.nptr := \text{mkleaf}(\text{num}, \text{num.val})$

Construction of a syntax tree for $a - 4 + c$



TYPES OF SDD'S

Syntax Directed definitions (SDD) are used to specify syntax directed translations. There are two types of SDD.

1. S-Attributed Definitions
2. L-Attributed Definitions.

S-attributed definitions

- Only synthesized attributes used in syntax direct definition.
- S-attributed grammars interact well with LR (K) parsers since the evaluation of attributes is bottom-up. They do not permit dependency graphs with cycles.

L-attributed definitions

- Both inherited and synthesized attribute are used.
- L-attributed grammar support the evaluation of attributes associated with a production body, dependency-graph edges can go from left to right only.
- Each S-attributed grammar is also a L-attributed grammar.
- L-attributed grammars can be incorporated conveniently in top down parsing.
- These grammars interact well with LL (K) parsers (both table driven and recursive descent).

Synthesized Attributes on the Parser Stack

A translator for an S-attributed definition often be implemented with LR parser generator. Here the stack is implemented by a pair of array state and val.

- Each state entry is pointed to a LR (1) parsing table.
- Each val[i] holds the value of the attributes associated with the node. For $A \rightarrow xyz$, the stack will be:

State	Val
Top \rightarrow	
Z	Z.z
Y	Y.y
X	X.x

Example: Consider the following grammar:

$S \rightarrow E \$$	{print(E.val)}
$E \rightarrow E + E$	{E.val := E.val + E.val}
$E \rightarrow E * E$	{E.val := E.val * E.val}
$E \rightarrow (E)$	{E.val := E.val}
$E \rightarrow I$	{I.val := I.val * 10 + digit}
$I \rightarrow I \text{ digit}$	
$I \rightarrow \text{digit}$	{I.val := digit}

Implementation

$S \rightarrow E \$$	print (val [top])
$E \rightarrow E + E$	val[ntop] := val[top] + val[top-2]
$E \rightarrow E * E$	val[ntop] := val[top] * val[top-2]
$E \rightarrow (E)$	val[ntop] := val[top-1]
$E \rightarrow I$	val[ntop] := val[top]
$I \rightarrow I \text{ digit}$	val[ntop] := 10*val[top] + digit
$I \rightarrow \text{digit}$	val[ntop] := digit

L-attributed Definitions

A syntax directed definition is L-attributed if each inherited attribute of $X_j, 1 \leq j \leq n$, on the right side of $A \rightarrow X_1 X_2 \dots X_n$, depends only on

1. The attributes of symbols X_1, X_2, \dots, X_{j-1} to the left of X_j in the production.
2. The inherited attributes of A.

Every S-attributed definition is L-attributed, because the above two rules apply only to the inherited attributes.

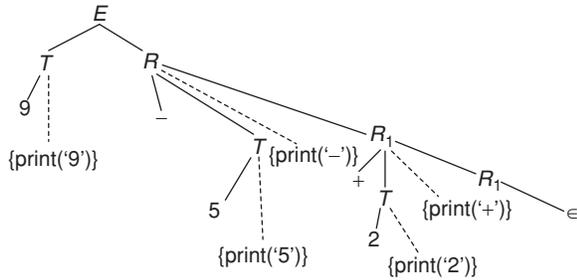
SYNTAX DIRECTED TRANSLATION SCHEMES

A translation scheme is a CFG in which attributes are associated with grammar symbols and semantic actions are enclosed between braces { } are inserted within the right sides of productions.

Example: $E \rightarrow TR$

$R \rightarrow \text{op } T$	{print (op.lexeme)} $R_1 \in$
$T \rightarrow \text{num}$	{print (num.val)}

Using this, the parse tree for $9 - 5 + 2$ is



If we have both inherited and synthesized attributes then we have to follow the following rules:

1. An inherited attribute for a symbol on the right side of a production must be computed in an action before that symbol.
2. An action must not refer to a synthesized attribute of a symbol on the right side of the action.
3. A synthesized attribute for the non-terminal on the left can only be computed after all attributes it references, have been computed.

Note: In the implementation of L-attributed definitions during predictive parsing, instead of syntax directed translations, we will work with translation schemes.

Eliminating left recursion from translation scheme

Consider following grammar, which has left recursion

$$E \rightarrow E + T \{ \text{print} ('+') ; \}$$

When transforming the grammar, treat the actions as if they were terminal symbols. After eliminating recursion from the above grammar.

$$\begin{aligned} E &\rightarrow TR \\ R &\rightarrow +T \{ \text{print} ('+') ; \} R \\ R &\rightarrow \epsilon \end{aligned}$$

BOTTOM-UP EVALUATION OF INHERITED ATTRIBUTES

- Using a bottom up translation scheme, we can implement any L-attributed definition based on LL (1) grammar.
- We can also implement some of L-attributed definitions based on LR (1) using bottom up translations scheme.
 - The semantic actions are evaluated during the reductions.
 - During the bottom up evaluation of S-attributed definitions, we have a parallel stack to hold synthesized attributes.

Where are we going to hold inherited attributes?

We will convert our grammar to an equivalent grammar to guarantee the following:

- All embedding semantic actions in our translation scheme will be moved to the end of the production rules.
- All inherited attributes will be copied into the synthesized attributes (may be new non-terminals).

Thus we will evaluate all semantic actions during reductions, and we find a place to store an inherited attribute. The steps are

1. Remove an embedding semantic action S_i , put new non-terminal M_i instead of that semantic action.
2. Put S_i into the end of a new production rule $M_i \rightarrow \epsilon$.
3. Semantic action S_i will be evaluated when this new production rule is reduced.
4. Evaluation order of semantic rules is not changed. i.e., if

$$A \rightarrow \{S_1\} X_1 \{S_2\} X_2 \dots \{S_n\} X_n$$

After removing embedding semantic actions:

$$A \rightarrow M_1 X_1 M_2 X_2 \dots M_n X_n$$

$$M_1 \rightarrow \epsilon \{S_1\}$$

$$M_2 \rightarrow \epsilon \{S_2\}$$

⋮

$$M_n \rightarrow \epsilon \{S_n\}$$

For example,

$$E \rightarrow TR$$

$$R \rightarrow +T \{ \text{print} ('+') \} R_1$$

$$R \rightarrow \epsilon$$

$$T \rightarrow \text{id} \{ \text{print} (\text{id.name}) \}$$

⇓ remove embedding semantic actions

$$E \rightarrow TR$$

$$R \rightarrow +TMR_1$$

$$R \rightarrow \epsilon$$

$$T \rightarrow \text{id} \{ \text{print} (\text{id.name}) \}$$

$$M \rightarrow \epsilon \{ \text{print} ('+') \}$$

Translation with inherited attributes

Let us assume that every non-terminal A has an inherited attribute $A.i$ and every symbol X has a synthesized attribute $X.s$ in our grammar.

For every production rule $A \rightarrow X_1, X_2, \dots, X_n$, introduce new marker non-terminals

$$M_1, M_2, \dots, M_n \text{ and replace this production rule with } A \rightarrow M_1 X_1 M_2 X_2 \dots M_n X_n$$

The synthesized attribute of X_i will not be changed.

The inherited attribute of X_i will be copied into the synthesized attribute of M_i by the new semantic action added at the end of the new production rule

$$M_i \rightarrow \epsilon$$

Now, the inherited attribute of X_i can be found in the synthesized attribute of M_i .

$$A \rightarrow \{B.i = f_1(\dots)\} B \{c.i = f_2(\dots)\} c \{A.s = f_3(\dots)\}$$

⇓

$$A \rightarrow \{M_1.i = f_1(\dots)\} M_1 \{B.i = M_1.s\} B \{M_2.i = f_2(\dots)\} M_2$$

$$\{c.i = M_2.s\} c \{A.s = f_3(\dots)\}$$

$$M_1 \rightarrow \epsilon \{M_1.s = M_1.i\}$$

$$M_2 \rightarrow \epsilon \{M_2.s = M_2.i\}$$

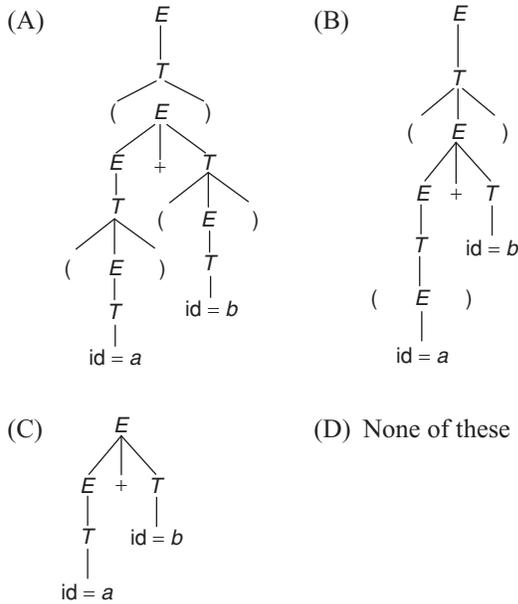
EXERCISES

Practice Problems I

Directions for questions 1 to 13: Select the correct alternative from the given choices.

1. The annotated tree for input ((a) + (b)), for the rules given below is

Production	Semantic Rule
$E \rightarrow E + T$	$\$ \$ = \text{mknode} ('+', \$1, \$3)$
$E \rightarrow E - T$	$\$ \$ = \text{mknode} ('-', \$1, \$3)$
$E \rightarrow T$	$\$ \$ = \$1;$
$T \rightarrow (E)$	$\$ \$ = \$2;$
$T \rightarrow \text{id}$	$\$ \$ = \text{mkleaf} (\text{id}, \$1)$
$T \rightarrow \text{num}$	$\$ \$ = \text{mkleaf} (\text{num}, \$1)$



2. Let synthesized attribute val give the value of the binary number generated by S in the following grammar.

$S \rightarrow L L$
 $S \rightarrow L$
 $L \rightarrow L B$
 $L \rightarrow B$
 $B \rightarrow 0$
 $B \rightarrow 1$

Input 101.101, $S.\text{val} = 5.625$
 use synthesized attributes to determine $S.\text{val}$
 Which of the following are true?

- (A) $S \rightarrow L_1 L_2 \{ S.\text{val} = L_1.\text{val} + L_2.\text{val} / (2^{**}L_2.\text{bits}) \}$
 $\quad | L \{ S.\text{val} = L.\text{val}; S.\text{bits} = L.\text{bits} \}$
 (B) $L \rightarrow L_1 B \{ L.\text{val} = L_1.\text{val} * 2 + B.\text{val}; \}$
 $\quad L.\text{bits} = L_1.\text{bits} + 1 \}$
 $\quad | B \{ L.\text{val} = B.\text{val}; L.\text{bits} = 1 \}$
 (C) $B \rightarrow 0 \{ B.\text{val} = 0 \}$
 $\quad | 1 \{ B.\text{val} = 1 \}$
 (D) All of these

3. Which of the following productions with translation rules converts binary number representation into decimal.

Production	Semantic Rule
$B \rightarrow 0$	$B.\text{trans} = 0$
$B \rightarrow 1$	$B.\text{trans} = 1$
$B \rightarrow B_0$	$B_1.\text{trans} = B_2.\text{trans} * 2$
$B \rightarrow B_1$	$B_1.\text{trans} = B_2.\text{trans} * 2 + 1$

Production	Semantic Rule
$B \rightarrow 0$	$B.\text{trans} = 0$
$B \rightarrow B_0$	$B_1.\text{trans} = B_2.\text{trans} * 4$

Production	Semantic Rule
$B \rightarrow 1$	$B.\text{trans} = 1$
$B \rightarrow B_1$	$B_1.\text{trans} = B_2.\text{trans} * 2$

- (D) None of these

4. The grammar given below is

Production	Semantic Rule
$A \rightarrow LM$	$L.i := l(A.i)$ $M.i := m(L.s)$ $A.s := f(M.s)$
$A \rightarrow QR$	$R.i := r(A.i)$ $Q.i := q(R.s)$ $A.s := f(Q.s)$

- (A) A L-attributed grammar
 (B) Non-L-attributed grammar
 (C) Data insufficient
 (D) None of these

5. Consider the following syntax directed translation:

$S \rightarrow aS \{ m := m + 3; \text{print}(m); \}$
 $\quad | bS \{ m := m * 2; \text{print}(m); \}$
 $\quad | \epsilon \{ m := 0; \}$

A shift reduce parser evaluate semantic action of a production whenever the production is reduced.

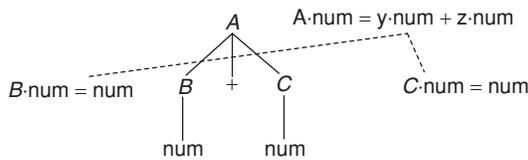
If the string is = a a b a b b then which of the following is printed?

- (A) 0 0 3 6 9 12 (B) 0 0 0 3 6 9 12
 (C) 0 0 0 3 6 9 12 15 (D) 0 0 3 9 6 12

6. Which attribute can be evaluated by shift reduce parser that execute semantic actions only at reduce moves but never at shift moves?

- (A) Synthesized attribute (B) Inherited attribute
 (C) Both (a) and (b) (D) None of these

7. Consider the following annotated parse tree:



Which of the following is true for the given annotated tree?

- (A) There is a specific order for evaluation of attribute on the parse tree.
- (B) Any evaluation order that computes an attribute 'A' after all other attributes which 'A' depends on, is acceptable.
- (C) Both (A) and (B)
- (D) None of these.

Common data for questions 8 and 9: Consider the following grammar and syntax directed translation.

$$\begin{aligned}
 E \rightarrow E + T & \quad E_1.val = E_2.val + T.val \\
 E \rightarrow T & \quad E.val = T.val \\
 T \rightarrow T * P & \quad T_1.val = T_2.val * P.val * \\
 & \quad P.num \\
 T \rightarrow P & \quad T.val = P.val * P.num \\
 P \rightarrow (E) & \quad P.val = E.val \\
 P \rightarrow 0 & \quad P.num = 1 \\
 & \quad P.val = 2 \\
 P \rightarrow 1 & \quad P.num = 2 \\
 & \quad P.val = 1
 \end{aligned}$$

- 8. What is $E.val$ for string $1*0$?
 (A) 8 (B) 6
 (C) 4 (D) 12
- 9. What is the $E.val$ for string $0 * 0 + 1$?
 (A) 8 (B) 6
 (C) 4 (D) 12

10. Consider the following syntax directed definition:

Production	Semantic Rule
$S \rightarrow b$	$S.x = 0$ $S.y = 0$
$S \rightarrow S_1 l$	$S.x = S_1.x + l.dx$ $S.y = S_1.y + l.dy$
$l \rightarrow \text{east}$	$l.dx = 1$ $l.dy = 0$
$l \rightarrow \text{north}$	$l.dx = 0$ $l.dy = 1$
$l \rightarrow \text{west}$	$l.dx = -1$ $l.dy = 0$
$l \rightarrow \text{south}$	$l.dx = 0$ $l.dy = -1$

If Input = begin east south west north, after evaluating this sequence what will be the value of $S.x$ and $S.y$?

- (A) (1, 0) (B) (2, 0)
- (C) (-1, -1) (D) (0, 0)

- 11. What will be the values $s.x, s.y$ if input is 'begin west south west'?
 (A) (-2, -1)
 (B) (2, 1)
 (C) (2, 2)
 (D) (3, 1)

12. Consider the following grammar:

$$\begin{aligned}
 S \rightarrow E & \quad S.val = E.val \\
 & \quad E.num = 1 \\
 E \rightarrow E * T & \quad E_1.val = 2 * E_2.val + 2 * T.val \\
 & \quad E_2.num = E_1.num + 1 \\
 & \quad T.num = E_1.num + 1 \\
 E \rightarrow T & \quad E.val = T.val \\
 & \quad T.num = E.num + 1 \\
 T \rightarrow T + P & \quad T_1.val = T_2.val + P.val \\
 & \quad T_2.num = T_1.num + 1 \\
 & \quad P.num = T_1.num + 1 \\
 T \rightarrow P & \quad T.val = P.val \\
 & \quad P.num = T.num + 1 \\
 P \rightarrow (E) & \quad P.val = E.val \\
 P \rightarrow i & \quad \left\{ \begin{array}{l} E.num = P.num \\ P.val = I | P.num \end{array} \right\}
 \end{aligned}$$

Which attributes are inherited and which are synthesized in the above grammar?

- (A) Num attribute is inherited attribute. Val attribute is synthesized attribute.
- (B) Num is synthesized attribute. Val is inherited attribute.
- (C) Num and val are inherited attributes.
- (D) Num and value are synthesized attributes.

13. Consider the grammar with the following translation rules and E as the start symbol.

$$\begin{aligned}
 E \rightarrow E_1 @ T & \quad \{ E.value = E_1.value * T.value \} \\
 & \quad | T \quad \{ E.value = T.value \} \\
 T \rightarrow T_1 \text{ and } F & \quad \{ T.value = T_1.value + F.value \} \\
 & \quad | F \quad \{ T.value = F.value \} \\
 F \rightarrow \text{num} & \quad \{ F.value = \text{num.value} \}
 \end{aligned}$$

Compute $E.value$ for the root of the parse tree for the expression: $2 @ 3$ and $5 @ 6$ and 4

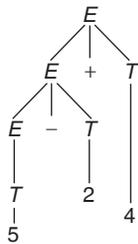
- (A) 200 (B) 180
- (C) 160 (D) 40

Practice Problems 2

Directions for questions 1 to 10: select the correct alternative from the given choices.

1. Consider the following Tree:

Production	Meaning
$E \rightarrow E_1 + T$	$E.t = E_1.t * T.t$
$E \rightarrow E_1 - T$	$E.t = E_1.t + T.t$
$E \rightarrow T$	$E.t = T.t$
$t \rightarrow 0$	$T.t = '0'$
$t \rightarrow 5$	$T.t = '5'$
$t \rightarrow 2$	$T.t = '2'$
$t \rightarrow 4$	$T.t = '4'$



After evaluation of the tree the value at the root will be:

- (A) 28 (B) 32
(C) 14 (D) 7
2. The value of an inherited attribute is computed from the values of attributes at the _____
- (A) Sibling nodes (B) Parent of the node
(C) Children node (D) Both (A) and (B)
3. Consider an action translating expression:

$\text{expr} \rightarrow \text{expr} + \text{term} \quad \{\text{print}(' + ')\}$
 $\text{expr} \rightarrow \text{expr} - \text{term} \quad \{\text{print}(' - ')\}$
 $\text{expr} \rightarrow \text{term}$
 $\text{term} \rightarrow 1 \quad \{\text{print}(' 1 ')\}$
 $\text{term} \rightarrow 2 \quad \{\text{print}(' 2 ')\}$
 $\text{term} \rightarrow 3 \quad \{\text{print}(' 3 ')\}$

Which of the following is true regarding the above translation expression?

- (A) Action translating expression represents infix notation.
 (B) Action translating expression represents prefix notation.
- (C) Action translating expression represents postfix notation.
 (D) None of these
4. In the given problem, what will be the result after evaluating $9 - 5 + 2$?
- (A) $+ - 9 5 2$ (B) $9 - 5 + 2$
 (C) $9 5 - 2 +$ (D) None of these
5. In a syntax directed translation, if the value of an attribute node is a function of the values of attributes of children, then it is called:
- (A) Synthesized attribute (B) Inherited attribute
 (C) Canonical attributes (D) None of these
6. Inherited attribute is a natural choice in:
- (A) Keeping track of variable declaration
 (B) Checking for the correct use of L-values and R-values.
 (C) Both (A) and (B)
 (D) None of these
7. Syntax directed translation scheme is desirable because
- (A) It is based on the syntax
 (B) Its description is independent of any implementation.
 (C) It is easy to modify
 (D) All of these
8. A context free grammar in which program fragments, called semantic actions are embedded within right side of the production is called,
- (A) Syntax directed translation
 (B) Translation schema
 (C) Annotated parse tree
 (D) None of these
9. A syntax directed definition specifies translation of construct in terms of:
- (A) Memory associated with its syntactic component
 (B) Execution time associated with its syntactic component
 (C) Attributes associated with its syntactic component
 (D) None of these
10. If an error is detected within a statement, the type assigned to the Statement is:
- (A) Error type (B) Type expression
 (C) Type error (D) Type constructor

PREVIOUS YEARS' QUESTIONS

Common data for questions 1 (A) and 1 (B): Consider the following expression grammar. The semantic rules for expression evaluation are stated next to each grammar production: [2005]

$E \rightarrow \text{number} \quad E.\text{val} = \text{number.val}$
 $| E ' + ' E \quad E^{(1)}.\text{val} = E^{(2)}.\text{val} + E^{(3)}.\text{val}$
 $| E \rightarrow E \quad E^{(1)}.\text{val} = E^{(2)}.\text{val} \times E^{(3)}.\text{val}$

1. (A) The above grammar and the semantic rules are fed to a yacc tool (which is an LALR (1) parser generator) for parsing and evaluating arithmetic expressions. Which one of the following is true about the action of yacc for the given grammar?
- (A) It detects recursion and eliminates recursion
 (B) It detects reduce-reduce conflict, and resolves

- (C) It detects shift-reduce conflict, and resolves the conflict in favor of a shift over a reduce action.
- (D) It detects shift-reduce conflict, and resolves the conflict in favor of a reduce over a shift action.
- (B) Assume the conflicts in Part (A) of this question are resolved and an LALR (1) parser is generated for parsing arithmetic expressions as per the given grammar. Consider an expression $3 \times 2 + 1$. What precedence and associativity properties does the generated parser realize?
- (A) Equal precedence and left associativity; expression is evaluated to 7
- (B) Equal precedence and right associativity; expression is evaluated to 9
- (C) Precedence of ‘ \times ’ is higher than that of ‘ $+$ ’, and both operators are left associative; expression is evaluated to 7
- (D) Precedence of ‘ $+$ ’ is higher than that of ‘ \times ’, and both operators are left associative; expression is evaluated to 9
2. In the context of abstract-syntax-tree (AST) and control-flow-graph (CFG), which one of the following is TRUE? [2015]
- (A) In both AST and CFG, let node N_2 be the successor of node N_1 . In the input program, the code corresponding to N_2 is present after the code corresponding to N_1 .
- (B) For any input program, neither AST nor CFG will contain a cycle.
- (C) The maximum number of successors of a node in an AST and a CFG depends on the input program.
- (D) Each node in AST and CFG corresponds to at most one statement in the input program.
3. Consider the following Syntax Directed Translation Scheme (SDTS), with non-terminals $\{S, A\}$ and terminals $\{a, b\}$. [2016]
- $$S \rightarrow aA \quad \{ \text{print 1} \}$$
- $$S \rightarrow a \quad \{ \text{print 2} \}$$
- $$A \rightarrow Sb \quad \{ \text{print 3} \}$$
- Using the above SDTS, the output printed by a bottom-up parser, for the input **aab** is:
- (A) 1 3 2 (B) 2 2 3
- (C) 2 3 1 (D) syntax error
4. Which one of the following grammars is free from *left recursion*? [2016]
- (A) $S \rightarrow AB$
 $A \rightarrow Aa/b$
 $B \rightarrow c$
- (B) $S \rightarrow Ab/Bb/c$
 $A \rightarrow Bd/\epsilon$
 $B \rightarrow e$
- (C) $S \rightarrow Aa/B$
 $A \rightarrow Bb/Sc/\epsilon$
 $B \rightarrow d$
- (D) $S \rightarrow Aa/Bb/c$
 $A \rightarrow Bd/\epsilon$
 $B \rightarrow Ae/\epsilon$

ANSWER KEYS

EXERCISES

Practice Problems 1

1. A 2. D 3. A 4. B 5. A 6. A 7. B 8. C 9. B 10. D
11. A 12. A 13. C

Practice Problems 2

1. A 2. D 3. C 4. C 5. A 6. C 7. D 8. B 9. C 10. C

Previous Years' Questions

1. (a) C (b) B 2. C 3. C 4. A

Chapter 3

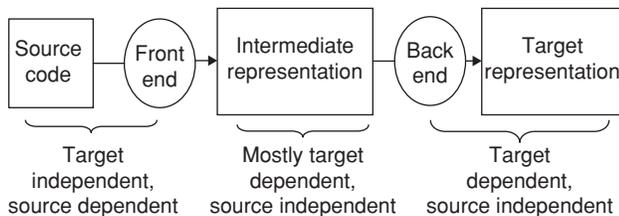
Intermediate Code Generation

LEARNING OBJECTIVES

- Introduction
- Directed Acyclic Graphs (DAG)
- Three address code
- Symbol table operations
- Assignment statements
- Boolean expression
- Flow control of statements
- Procedure calls
- Code generation
- Next use information
- Run-time storage management
- DAG representations of basic blocks
- Peephole optimization

INTRODUCTION

In the analysis–synthesis model, the front end translates a source program into an intermediate representation (IR). From IR the back end generates target code.



There are different types of intermediate representations:

- High level IR, i.e., AST (Abstract Syntax Tree)
- Medium level IR, i.e., Three address code
- Low level IR, i.e., DAG (Directed Acyclic Graph)
- Postfix Notation (Reverse Polish Notation, RPN).

In the previous sections already we have discussed about AST and RPN.

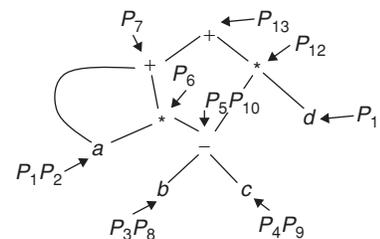
Benefits of Intermediate code generation: The benefits of ICG are

- We can obtain an optimized code.
- Compilers can be created for the different machines by attaching different backend to existing front end of each machine.
- Compilers can be created for the different source languages.

Directed acyclic graphs for expression: (DAG)

- A DAG for an expression identifies the common sub expressions in the given expression.
- A node N in a DAG has more than one parent if N represents a common sub expression.
- DAG gives the compiler, important clues regarding the generation of efficient code to evaluate the expressions.

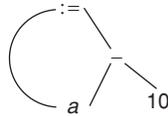
Example 1: DAG for $a + a*(b - c) + (b - c)*d$



- $P_1 = \text{makeleaf}(\text{id}, a)$
- $P_2 = \text{makeleaf}(\text{id}, a) = P_1$
- $P_3 = \text{makeleaf}(\text{id}, b)$
- $P_4 = \text{makeleaf}(\text{id}, c)$
- $P_5 = \text{makenode}(-, P_3, P_4)$
- $P_6 = \text{makenode}(*, P_1, P_5)$
- $P_7 = \text{makenode}(+, P_1, P_6)$
- $P_8 = \text{makeleaf}(\text{id}, b) = P_3$
- $P_9 = \text{makeleaf}(\text{id}, c) = P_4$
- $P_{10} = \text{makenode}(-, P_8, P_9) = P_5$

$P_{11} = \text{makeleaf}(id, d)$
 $P_{12} = \text{makenode}(*, P_{10}, P_{11})$
 $P_{13} = \text{makenode}(+, P_7, P_{12})$

Example 2: $a := a - 10$



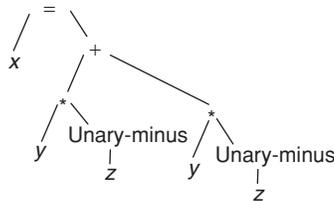
THREE-ADDRESS CODE

In three address codes, each statement usually contains 3 addresses, 2 for operands and 1 for the result.

Example: $-x = y \text{ OP } z$

- x, y, z are names, constants or compiler generated temporaries,
- OP stands for any operator. Any arithmetic operator (or) Logical operator.

Example: Consider the statement $x = y * -z + y * -z$



The corresponding three address code will be like this:

Syntax Tree	DAG
$t_1 = -z$	$t_1 = -z$
$t_2 = y * t_1$	$t_2 = y * t_1$
$t_3 = -z$	$t_5 = t_2 + t_2$
$t_4 = y * t_3$	$X = t_5$
$t_5 = t_4 + t_2$	
$X = t_5$	

The postfix notation for syntax tree is: $xyz \text{ unaryminus } *yz \text{ unaryminus } *+=$.

- Three address code is a 'Linearized representation' of syntax tree.
- Basic data of all variables can be formulated as syntax directed translation. Add attributes whenever necessary.

Example: Consider below SDD with following specifications:

E might have E . place and E .code

E .place: the name that holds the value of E .

E .code: the sequence of intermediate code starts evaluating E .
 Let Newtemp: returns a new temporary variable each time it is called.

New label: returns a new label.

Then the SDD to produce three-address code for expressions is given below:

Production	Semantic Rules
$S \rightarrow id \text{ ASN } E$	$S. \text{code} = E. \text{code} \ \backslash \ \text{gen}(\text{ASN}, id. \text{place}, E. \text{place})$ $E. \text{Place} = \text{newtemp}();$
$E \rightarrow E_1 \text{ PLUS } E_2$	$E. \text{code} = E_1. \text{code} \ \ E_2. \text{code} \ \ \text{gen}(\text{PLUS}, E. \text{place}, E_1. \text{place}, E_2. \text{place});$ $E. \text{place} = \text{newtemp}();$
$E \rightarrow E_1 \text{ MUL } E_2$	$E. \text{code} = E_1. \text{code} \ \ E_2. \text{code} \ \ \text{gen}(\text{MUL}, E. \text{place}, E_1. \text{place}, E_2. \text{place});$ $E. \text{Place} = \text{Newtemp}();$
$E \rightarrow \text{UMINUS } E_1$	$E. \text{code} = E_1. \text{code} \ \ \text{gen}(\text{NEG}, E. \text{Place}, E_1. \text{place});$ $E. \text{code} = E_1. \text{code}$
$E \rightarrow \text{LP } E_1 \text{ RP}$	$E. \text{Place} = E_1. \text{Place}$
$E \rightarrow \text{IDENT}$	$E. \text{place} = id. \text{place}$ $E. \text{code} = \text{empty.list}();$

Types of Three Address Statement

Assignment

- Binary assignment: $x := y \text{ OP } z$ Store the result of $y \text{ OP } z$ to x .
- Unary assignment: $x := \text{op } y$ Store the result of unary operation on y to x .

Copy

- Simple Copy $x := y$ Store y to x
- Indexed Copy $x := y[i]$ Store the contents of $y[i]$ to x
- $x[i] := y$ Store y to $(x + i)$ th address.

Address and pointer manipulation

- $x := \&y$ Store address of y to x
- $x := *y$ Store the contents of y to x
- $*x := y$ Store y to location pointed by x .

Jump

- Unconditional jump:- goto L , jumps to L .
- Conditional:

```
if (x relop y)
goto L1;
else
```

```
{
goto L2;
}
```

Where relop is <, <=, >, >=, = or ≠.

Procedure call

Param x_1 ;

Param x_2 ;

.

.

.

Param x_n ;

Call p, n, x ; Call procedure p with n parameters and store the result in x .

return x Use x as result from procedure.

Declarations

- Global x, n_1, n_2 : Declare a global variable named x at offset n_1 having n_2 bytes of space.
- Proc x, n_1, n_2 : Declare a procedure x with n_1 bytes of parameter space and n_2 bytes of local variable space.
- Local x, m : Declare a local variable named x at offset m from the procedure frame.
- End: Declare the end of the current procedure.

Adaption for object oriented code

- $x = y$ field z : Lookup field named z within y , store address to x
- Class x, n_1, n_2 : declare a class named x with n_1 bytes of class variables and n_2 bytes of class method pointers.
- Field x, n : Declare a field named x at offset n in the class frame.
- New x : Create a new instance of class name x .

Implementation of Three Address Statements

Three address statements can be implemented as records with fields for the operator and the operands. There are 3 types of representations:

1. Quadruples
2. Triples
3. Indirect triples

Quadruples

A quadruple has four fields: op, arg1, arg2 and result.

- Unary operators do not use arg2.
- Param use neither arg2 nor result.
- Jumps put the target label in result.
- The contents of the fields are pointers to the symbol table entries for the names represented by these fields.
- Easier to optimize and move code around.

Example 1: For the expression $x = y * - z + y * - z$, the quadruple representation is

	OP	Arg1	Arg2	Result
(0)	Uminus	z		t_1
(1)	*	y	t_1	t_2
(2)	Uminus	z		t_3
(3)	*	y	t_3	t_4
(4)	+	t_2	t_4	t_5
(5)	=	t_5		x

Example 2: Read (x)

	Op	Arg1	Arg2	Result
(0)	Param	x		
(1)	Call	READ	(x)	

Example 3: WRITE ($A*B, x+5$)

	OP	Arg1	Arg2	Result
(0)	*	A	B	t_1
(1)	+	x	5	t_2
(2)	Param	t_1		
(3)	Param	t_2		
(4)	Call	Write	2	

Triples

Triples have three fields: OP, arg1, arg2.

- Temporaries are not used and instead references to instructions are made.
- Triples are also known as two address code.
- Triples takes less space when compared with Quadruples.
- Optimization by moving code around is difficult.
- The DAG and triple representations of expressions are equivalent.
- For the expression $a = y * - z + y * - z$ the Triple representation is

	Op	Arg1	Arg2
(0)	Uminus	z	
(1)	*	y	(0)
(2)	Uminus	z	
(3)	*	y	(2)
(4)	+	(1)	(3)
(5)	=	a	(4)

Array – references

Example: For $A [I] = B$, the quadruple representation is

	Op	Arg1	Arg2	Result
(0)	[] =	A	I	T_1
(1)	=	B		T_2

The same can be represented by Triple representation also. [] = is called L-value, specifies the address to an element.

	Op	Arg1	Arg2
(0)	[] =	A	I
(1)	=	(0)	B

Example 2: $A: = B [I]$

	Op	Arg1	Arg2
(0)	= []	B	I
(1)	=	A	(0)

$= []$ is called r-value, specifies the value of an element.

Indirect Triples

- In indirect triples, pointers to triples will be there instead of triples.
- Optimization by moving code around is easy.
- Indirect triples takes less space when compared with Quadruples.
- Both indirect triples and Quadruples are almost equally efficient.

Example: Indirect Triple representation of 3-address code

	Statement
(0)	(14)
(1)	(15)
(2)	(16)
(3)	(17)
(4)	(18)
(5)	(19)

	Op	Arg1	Arg2
(14)	Uminus	Z	
(15)	*	y	(14)
(16)	Uminus	Z	
(17)	*	y	(16)
(18)	+	(15)	(17)
(19)	=	x	(18)

SYMBOL TABLE OPERATIONS

Treat symbol tables as objects.

- Mktable (previous);
 - create a new symbol table.
 - Link it to the symbol table previous.
- Enter (table, name, and type, offset)
 - insert a new identifier name with type and offset into table
 - Check for possible duplication.
- Add width (table, width);
 - increase the size of symbol table by width.
- Enterproc (table, name, new table)
 - Enter a procedure name into table.
 - The symbol table of name is new table.
- Lookup (name, table);
 - Check whether name is declared in the symbol table, if it is in the table then return the entry.

Example:

Declaration $\rightarrow M_1 D$

$M_1 \rightarrow \in \{ \text{TOP (Offset):} = 0 ; \}$

$D \rightarrow D \text{ ID}$

$D \rightarrow \text{id: } T \{ \text{enter (top (tblptr), id.name, T.type top (offset)); top (offset):} = \text{top (offset) + T. width ;} \}$

$T \rightarrow \text{integer } \{ T.type : = \text{integer; T. width:} = 4 ; \}$

$T \rightarrow \text{double } \{ T.type: = \text{double; T.width} = 8 ; \}$

$T \rightarrow * T_1 \{ T.type: = \text{pointer (T.type); T.width} = 4 ; \}$

Need to remember the current offset before entering the block, and to restore it after the block is closed.

Example: Block $\rightarrow \text{begin } M_4 \text{ Declarations statements end}$
 $\{ \text{pop (tblptr); pop (offset) ;} \}$

$M_4 \rightarrow \in \{ t: = \text{mktable (top (tblptr); push (t, tblptr); push (top (offset), offset) ;} \}$

Can also use the block number technique to avoid creating a new symbol table.

Field names in records

- A record declaration is treated as entering a block in terms of offset is concerned.
- Need to use a new symbol table.

Example: $T \rightarrow \text{record } M_5 D \text{ end}$

$\{ T.type: = \text{(top (tblptr)); T. width} = \text{top (offset); pop (tblptr); pop (offset) ;} \}$

$M_5 \rightarrow \in \{ t: = \text{mktable (null); push (t, tblptr); push \{(o, offset) ;} \}$

ASSIGNMENT STATEMENTS

Expressions can be of type integer, real, array and record. As part of translation of assignments into three address code, we show how names can be looked up in the symbol table and how elements of array can be accessed.

Code generation for assignment statements gen ([address # 1], [assignment], [address #2], operator, address # 3);

Variable accessing Depending on the type of [address # i], generate different codes.

Types of [address # i]:

- Local temp space
- Parameter
- Local variable
- Non-local variable
- Global variable
- Registers, constants,...

Error handling routine error – msg (error information);

The error messages can be written and stored in other file. Temp space management:

- This is used for generating code for expressions.
- newtemp (): allocates a temp space.
- freetemp (): free t if it is allocated in the temp space

Label management

- This is needed in generating branching statements.
- newlabel (): generate a label in the target code that has never been used.

Names in the symbol table

```
S → id := E {p := lookup (id-name, top (tblptr));
  If p is not null then gen (p, " := ", E.place);
  Else error ("var undefined", id.name);
}
E → E1 + E2 {E.place = newtemp ();
  gen (E.place, " := ", E1.place, "+", E2.place);
  free temp (E1.place);
  freetemp (E2.place);
}
E → -E1 {E.place = newtemp ();
  gen (E.place, " := ", "uminus", E1.place);
  freetemp (E1.place);
}
E → (E1) {E.place = E1.place;
}
E → id {p := lookup (id.name, top (tblptr));
  If p ≠ null then E.place = p.place else error ("var undefined", id.name);
}
```

Type conversions

Assume there are only two data types: integer, float.

For the expression,

$$E \rightarrow E_1 + E_2$$

If E_1 .type = E_2 .type then

generate no conversion code

E .type = E_1 .type;

Else

```
E.type = float;
temp1 = newtemp ();
If E1.type = integer then
  gen (temp1, ':=' int - to - float, E1.place);
  gen (E, ':=' temp1, '+', E2.place);
```

Else

```
gen (temp1, ':=' int - to - float, E2.place);
gen (E, ':=' temp1, '+', E1.place);
Free temp (temp1);
```

Addressing array elements

Let us assume

low: lower bound

w: element data width

Start_addr: starting address

1D Array: A[i]

- $Start_addr + (i - low) * w = i * w + (start_addr - low * w)$
- The value called base, $(start_addr - low * w)$ can be computed at compile time and then stored at the symbol table.

Example: array [-8 ... 100] of integer.

To declare [-8] [-7] ... [100] integer array in Pascal.

2D Array A [i₁, i₂]

Row major order: row by row. $A [i]$ means the i th row.

1st row $A [1, 1]$
 $A [1, 2]$

2nd row $A [2, 1]$
 $A [2, 2]$

$$A [i, j] = A [i] [j]$$

Column major: column by column.

$$A [1, 1] \uparrow A [1, 2]$$

$$A [2, 1] \uparrow A [2, 2]$$

1st Column 2nd column

Address for $A [i_1, i_2]$:

$$Start_addr + ((i_1 - low_1) * n_2 + (i_2 - low_2)) * w$$

Where low_1 and low_2 are the lower bounds of i_1 and i_2 . n_2 is the number of values that i_2 can take. $high_2$ is the upper bound on the value of i_2 . $n_2 = high_2 - low_2 + 1$

We can rewrite address for $A [i_1, i_2]$ as $((i_1 * n_2) + i_2) * w + (start_addr - ((low_1 * n_2) + low_2) * w)$. The value $(start_addr - low_1 * n_2 * w - low_2 * w)$ can be computed at compiler time and then stored in the symbol table.

Multi-Dimensional Array A [i₁, i₂, ... i_k]

Address for $A [i_1, i_2, \dots i_k]$

$$= i_1 * \pi_{i=2}^k n_i + i_2 * \pi_{i=3}^k n_i + \dots + i_k * w + (start_addr - low_1 * w * \pi_{i=2}^k n_i - low_2 * w * \pi_{i=3}^k n_i - \dots - low_k * w)$$

It can be computed incrementally in grammar rules:

$$f(1) = i_1;$$

$$f(j) = f(j-1) * n_j + i_j;$$

$f(k)$ is the value we wanted to compute.

Attributes needed in the translation scheme for addressing array elements:

Elegize: size of each element in the array

Array: a pointer to the symbol table entry containing information about the array declaration.

Ndim: the current dimension index

Base: base address of this array

Place: where a variable is stored.

Limit (array, n) = n_m is the number of elements in the m th coordinate.

Translation scheme for array elements

Consider the grammar

$$S \rightarrow L := E$$

$$E \rightarrow L$$

```

L → id
L → [Elist]
Elist → Elist1, E
Elist → id [E]
E → id
E → E + E
E → (E)

```

- $S \rightarrow L := E$ {if L.offset = null then /* L is a simple id */ gen (L.place, ":", E.place); Else gen (L.place, "[", L.offset, "]", ":", E.place);}
- $E \rightarrow E_1 + E_2$ {E.place = newtemp (); gen (E.place, ":", E1.place, "+", E2.place); }
- $E \rightarrow (E_1)$ {E.place = E₁.place}
- $E \rightarrow L$ {if L.offset = null then /* L is a simple id */ E.place := L.place; Else begin E.place := newtemp (); gen (E.place, ":", L.place, "[", L.offset, "]"); end }
- $L \rightarrow id$ {P! = lookup (id.name, top (tblptr)); If P ≠ null then Begin L.place := P.place; L.offset := null; End Else Error ("Var underfined", id.Name); }
- $L \rightarrow Elist$ {L.offset := newtemp (); gen (L.offset, ":", Elist.elesize, "*", Elist.place); freetemp (Elist.place); L.Place := Elist . base ;}
- $Elist \rightarrow Elist_1, E$ {t: =newtemp (); m: = Elist1.ndim+1; gen (t, ":" Elist1.place, "*", limit (Elist1.array, m)); Gen (t, ":", t"+", E.place); freetemp (E.place); Elist.array := Elist.array; Elist.place := t; Elist.ndim := m ;}
- $Elist \rightarrow id$ [E {Elist.Place := E.place; Elist.ndim := 1; P! = lookup (id.name, top (tblptr)); check for id errors; Elist.elesize := P.size; Elist.base := p.base; Elist.array := p.place ;}
- $E \rightarrow id$ {P := lookup (id,name, top (tblptr); Check for id errors; E.Place := Poplace ;}

BOOLEAN EXPRESSIONS

There are two choices for implementation of Boolean expressions:

1. Numerical representation
2. Flow of control

Numerical representation

Encode true and false values.

Numerically, 1:true 0: false.

Flow of control: Representing the value of a Boolean expression by a position reached in a program.

Short circuit code: Generate the code to evaluate a Boolean expression in such a way that it is not necessary for the code to evaluate the entire expression.

- If a_1 or a_2
 a_1 is true then a_2 is not evaluated.
- If a_1 and a_2
 a_1 is false then a_2 is not evaluated.

Numerical representation

$E \rightarrow id_1 \text{ relop } id_2$

```

{B.place := newtemp ();
gen ("if", id1.place, relop.op, id2.
place, "goto", next stat +3);
gen (B.place, ":", "0");
gen ("goto", nextstat+2);
gen (B.place, ":", "1")' }

```

Example 1: Translate the statement (if $a < b$ or $c < d$ and $e < f$) without short circuit evaluation.

```

100: if a < b goto 103
101: t1 := 0
102: goto 104
103: t1 := 1 /* true */
104: if c < d goto 107
105: t2 := 0 /* false */
106: goto 108
107: t2 := 1
108: if e < f goto 111
109: t3 := 0
110: goto 112
111: t3 := 1
112: t4 := t2 and t3
113: t3 := t1 or t4

```

FLOW OF CONTROL STATEMENTS

$B \rightarrow id_1 \text{ relop } id_2$

```

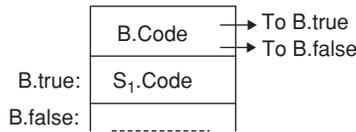
{
B.true := newlabel ();
B.false := newlabel ();
B.code := gen ("if", id1.relop, id2, "goto",

```

```
B.true, "else", "goto", B.false) ||
gen (B.true, ":")
}
S → if B then S1 S.code := B.code || S1.code || gen
(B.false, ':')
|| is the code concatenation operator.
```

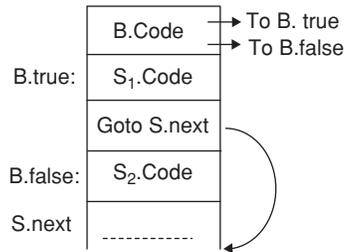
1. If – then implementation:

$S \rightarrow \text{if } B \text{ then } S_1 \{ \text{gen (B.false, " :"); } \}$



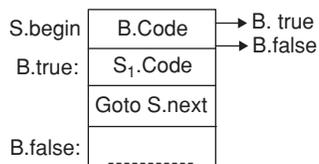
2. If – then – else

```
P → S { S.next := newlabel ();
P.code := S.code || gen (S.next, " :") }
S → if B then S1 else S2 { S1.next := S.next;
S2.next := S.next;
S.code := B.code || S1.code ||
Gen ("goto" S.next) || B.false, " :")
|| S2.code }
Need to use inherited attributes of S to define the
attributes of S1 and S2
```



3. While loop:

```
B → id1 relop id2 B.true := newlabel ();
B.false := newlabel ();
B.code := gen ('if', id.relop,
id2, 'goto', B.true 'else', 'goto', B.false) ||
gen (B.true ':');
S → while B do S1 S.begin := newlabel ();
S.code := gen (S.begin, ':') ||
B.code || S1.code || gen
('goto', S.begin) || gen (B.false, ':');
```



4. Switch/case statement:

The c - like syntax of switch case is
switch epr {
case V [1]: S [1]

```
.
.
.
case V [k]: S[k]
default: S[d]
}
```

Translation sequence

- Evaluate the expression.
- Find which value in the list matches the value of the expression, match default only if there is no match.
- Execute the statement associated with the matched value.

How to find the matched value? The matched value can be found in the following ways:

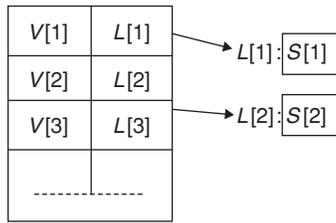
1. Sequential test
2. Lookup table
3. Hash table
4. Back patching

Two different translation schemes for sequential test are shown below:

```
1. Code to evaluate E into t
Goto test
L[i]: code for S [1]
goto next
-----
L[k]: code for S[k]
goto next
L[d]: code for S[d]
Go to next test:
If t = V [1]: goto L [1]
```

```
.
.
.
goto L[d]
Next:
2. Can easily be converted into look up table
If t <> V [i] goto L [1]
Code for S [1]
goto next
-----
L [1]: if t <> V [2] goto L [2]
Code for S [2]
Goto next
L [k - 1]: if t <> V [k] goto L[k]
Code for S[k]
Goto next
.
.
.
L[k]: code for S[d]
Next:
```

Use a table and a loop to find the address to jump



3. Hash table: When there are more than two entries use a hash table to find the correct table entry.

4. Back patching:

- Generate a series of branching statements with the targets of jumps temporarily left unspecified.
- To determine label table: each entry contains a list of places that need to be back patched.
- Can also be used to implement labels and gotos.

PROCEDURE CALLS

- Space must be allocated for the activation record of the called procedure.
- Arguments are evaluated and made available to the called procedure in a known place.
- Save current machine status.
- When a procedure returns:
 - Place returns value in a known place.
 - Restore activation record.

Example: $S \rightarrow \text{call id (Elist)}$

```
{for each item P on the queue Elist.
Queue do gen ('PARAM', q);
gen ('call:', id.place) ;}
Elist  $\rightarrow$  Elist, E {append E.place to the end of
Elist.queue}
Elist  $\rightarrow$  E {initialize Elist.queue to contain only
E.place}
```

Use a queue to hold parameters, then generate codes for params.

Code for E_1 , store in t_1

·
·
·

Code for E_k , store in t_k

PARAM t_1

:
·
·

PARAM t_k

Call P

Terminology:

Procedure declaration:

Parameters, formal parameters

Procedure call:

Arguments, actual parameters.

The values of a variable: $x = y$

r – value: value of the variable, i.e., on the right side of assignment. Ex: y , in above assignment.

l – value: The location/address of the variable, i.e., on the leftside of assignment. Ex: x , in above assignment.

There are different modes of parameter passing

1. call-by-value
2. call-by-reference
3. call-by-value-result (copy-restore)
4. call-by-name

Call by value

Calling procedure copies the r values of the arguments into the called procedure's Activation Record.

Changing a formal parameter has no effect on the actual parameter.

Example: void add (int C)

```
{
C = C+ 10;
printf ('\nc = %d', &C);
}
main ()
{
int a = 5;
printf ('a=%d', &a);
add (a);
printf ('\na = %d', &a);
}
```

In main a will not be affected by calling add (a)

It prints $a = 5$

$a = 5$

Only the value of C in add ($$) will be changed to 15.

Usage:

1. Used by PASCAL and C++ if we use non-var parameters.
2. The only thing used in C.

Advantages:

1. No aliasing.
2. Easier for static optimization analysis.
3. Faster execution because of no need for redirecting.

Call by reference

Calling procedure copies the l -values of the arguments into the called procedure's activation record. i.e., address will be passed to the called procedure.

- Changing formal parameter affects the corresponding actual parameter.
- It will have some side effects.

Example: void add (int *c)

```
{
*c = *c + 10;
printf ('\nc=%d', *c);
}
```

```

}
void main()
{
int a = 5;
printf ("\na = %d", a);
add (&a);
printf ("\na = %d", a);
output: a = 5
          c = 15
          a = 15

```

That is, here the actual parameter is also modified.

Advantages

1. Efficiency in passing large objects.
2. Only need to copy addresses.

Call-by-value-result

Equivalent to call-by-reference except when there is aliasing. That is, the program produces the same result, but not the same code will be generated.

Aliasing: Two expressions that have the same l-values are called aliases. They access the same location from different places.

Aliasing happens through pointer manipulation.

1. Call by reference with global variable as an argument.
2. Call by reference with the same expression as argument twice.

Example: test (x, y, x)

Advantages:

1. If there is no aliasing, we can implement it by using call – by – reference for large objects.
2. No implicit side effect if pointers are not passed.

Call by-name

used in Algol.

- Procedure body is substituted for the call in calling procedure.
- Each occurrence of a parameter in the called procedure is replaced with the corresponding argument.
- Similar to macro expansion.
- A parameter is not evaluated unless its value is needed during computation.

Example:

```

void show (int x)
{
for (int y = 0; y < 10; y++)
x++;
}
main ()
{
int j;
j = -1;
show (j);
}

```

Actually it will be like this
main ()
{

```

int j;
j = - 1;
For (in y= 0; y < 10; y ++)
x ++;
}

```

- Instead of passing values or address as arguments, a function is passed for each argument.
- These functions are called **thunks**.
- Each time a parameter is used, the thunk is called, then the address returned by the thunk is used.

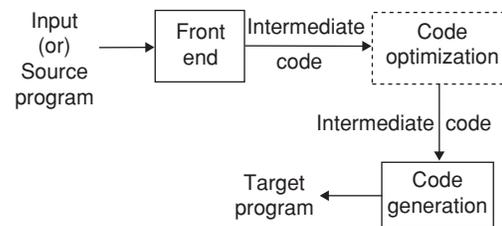
$y = 0$: use return value of thunk for y as the ℓ -value.

Advantages

- More efficient when passing parameters that are never used.
- This saves lot of time because evaluating unused parameter takes a longtime.

CODE GENERATION

Code generation is the final phase of the compiler model.



The requirements imposed on a code generator are

1. Output code must be correct.
2. Output code must be of high quality.
3. Code generator should run efficiently.

Issues in the Design of a Code Generator

The generic issues in the design of code generators are

- Input to the code generator
- Target programs
- Memory Management
- Instruction selection
- Register Allocation
- Choice of Evaluation order

Input to the code generator

Intermediate representation with symbol table will be the input for the code generator.

- High Level Intermediate representation

Example: Abstract Syntax Tree (AST)

- Medium – level intermediate representation

Example: control flow graph of complex operations

- Low – Level Intermediate representation

Example: Quadruples, DAGS

- Code for abstract stack machine, i.e., postfix code.

Target programs

The output of the code generator is the target program. The output may take on a variety of forms:

1. Absolute machine language
2. Relocatable machine language
3. Assembly language

Absolute machine language

- Final memory area for a program is statically known.
- Hard coded addresses.
- Sufficient for very simple systems.

Advantages:

- Fast for small programs
- No separate compilation

Disadvantages: Can not call modules from other languages/compilers.

Relocatable code

It Needs

- Relocation table
- Relocating linker + loader (or) runtime relocation in Memory management Unit (MMU).

Advantage: More flexible.

Assembly language Generates assembly code and use an assembler tool to convert this to binary (object) code. It needs (i) assembler (ii) linker and loader.

Advantage: Easier to handle and closer to machine.

Memory management

Mapping names in the source program to addresses of data objects in runtime memory is done by the front end and the code generator.

- A name in a three address statement refers to a symbol entry for the name.
- Stack, heap, garbage collection is done here.

Instruction selection

Instruction selection depends on the factors like

- Uniformity
- Completeness of the instruction
- Instruction speed
- Machine idioms
- Choose set of instructions equivalent to intermediate representation code.
- Minimize execution time, used registers and code size.

Example: $x = y + z$ in three address statements:

```
MOV y, R0 /* load y into R0 */
ADD z, R0
MOV R0, x /* store R0 into x*/
```

Register allocation

- Instructions with register operands are faster. So, keep frequently used values in registers.
- Some registers are reserved.

Example: *SP, PC* ... etc.

Minimize number of loads and stores.

Evaluation order

- The order of evaluation can affect the efficiency of the target code.
- Some orders require fewer registers to hold intermediate results.

Target Machine

Lets us assume, the target computer is

- Byte addressable with 4 bytes per word
- It has n general purpose registers
 $R_0, R_1, R_2, \dots, R_{n-1}$
- It has 2 address instructions of the form
OP source, destination
[cost: 1 + added]

Example: The op may be MOV, ADD, MUL. Generally cost will be like this

Source	Destination	Cost
Register	Register	1
Register	Memory	2
Memory	Register	2
Memory	Memory	3

Addressing modes:

Mode	Form	Address	Cost
Absolute	M	M	2
Register	R	R	1
Indexed	C(R)	C+contents(R)	2
Indirect register	*R	Contents (R)	1
Indirect indexed	*C(R)	Contents (C+contents (R))	2

Example: $x = y - z$

MOV y, R0 → cost = 2

SUB z, R0 → cost = 2

MOV R₀, x → cost = 2

RUNTIME STORAGE MANAGEMENT

Storage Organization

To run a compiled program, compiler will demand the operating system for the block of memory. This block of memory is called runtime storage.

This run time storage is subdivided into the generated target code, Data objects and Information which keeps track of procedure activations.

The fixed data (generated code) is stored at the statically determined area of the memory. The Target code is placed at the lower end of the memory.

The data objects are stored at the statically determined area as its size is known at the compile time. Compiler stores these data objects at statically determined area because these are compiled into target code. This static data area is placed on the top of the code area.

The runtime storage contains stack and the heap. Stack contains activation records and program counter, data object within this activation record are also stored in this stack with relevant information.

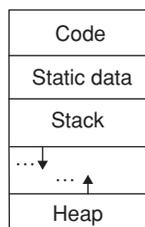
The heap area allocates the memory for the dynamic data (for example some data items are allocated under the program control)

The size of stack and heap will grow or shrink according to the program execution.

Activation Record

Information needed during an execution of a procedure is kept in a block of storage called an activation record.

- Storage for names local to the procedures appears in the activation record.
- Each execution of a procedure is referred as activation of the procedure.
- If the procedure is recursive, several of its activation might be alive at a given time.
- Runtime storage is subdivided into
 1. Generated target code area
 2. Data objects area
 3. Stack
 4. Heap



- Sizes of stack and heap can change during program execution.

For code generation there are two standard storage allocations:

1. **Static allocation:** The position of an activation record in memory is fixed at compile time.
2. **Stack allocation:** A new activation record is pushed on to the stack for each execution of the procedure. The record is popped when the activation ends.

Control stack The control stack is used for managing active procedures, which means when a call occurs, the execution of activation is interrupted and status information of the stack is saved on the stack.

When control is returned from a call, the suspended activation is resumed after storing the values of relevant registers it also includes program counter which sets to point immediately after the call.

The size of stack is not fixed.

Scope of declarations Declaration scope refers to the certain program text portion, in which rules are defined by the language.

Within the defined scope, entity can access legally to declared entities.

The scope of declaration contains immediate scope always. Immediate scope is a region of declarative portion with enclosure of declaration immediately.

Scope starts at the beginning of declaration and scope continues till the end of declaration. Whereas in the over loadable declaration, the immediate scope will begin, when the callable entity profile was determined.

The visible part refers text portion of declaration, which is visible from outside.

Flow Graph

A flow graph is a graph representation of three address statement sequences.

- Useful for code generation algorithms.
- Nodes in the flow graph represents computations.
- Edges represent flow of control.

Basic Blocks

Basic blocks are sequences of consecutive statements in which flow of control enters at the beginning and leaves at the end without a halt or branching.

1. First determine the set of leaders
 - First statement is leader
 - Any target of goto is a leader
 - Any statement that follows a goto is a leader.
2. For each leader its basic block consists of the leader and all statements up to next leader.

Initial node: Block with first statement is leader.

Example: consider the following fragment of code that computes dot product of two vectors x and y of length 10.

```
begin
Prod: = 0;
```

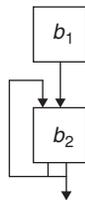
```

i := 1;
repeat
begin
Prod := Prod + x [i] * y [i];
i := i + 1;
end
until i <= 10;
end
    
```

B_1	(1)	Prod := 0
	(2)	i := 1

B_2	(3)	$t_1 := 4 * i$
	(4)	$t_2 := x[t_1]$
	(5)	$t_3 := 4 * i$
	(6)	$t_4 := y[t_3]$
	(7)	$t_5 := t_2 * t_4$
	(8)	$t_6 := \text{Prod} + t_5$
	(9)	Prod := t_6
	(10)	$t_7 := i + 1$
	(11)	$i := t_7$
	(12)	if $i <= 10$ goto (3)

∴ The flow graph for this code will be



Here b_1 is the initial node/block.

- Once the basic blocks have been defined, a number of transformations can be applied to them to improve the quality of code.
 - Global:** Data flow analysis
 - Local:**
 - Structure preserving transformations
 - Algebraic transformations
- Basic blocks compute a set of expressions. These expressions are the values of the names live on exit from the block.
- Two basic blocks are equivalent if they compute the same set of expressions.

Structure preserving transformations:

- Common sub-expression elimination:

$$\begin{array}{l}
 a := b + c \\
 b := a - d \\
 c := b + c \\
 d := a - d
 \end{array}
 \Rightarrow
 \begin{array}{l}
 a := b + c \\
 b := a - d \\
 c := b + c \\
 d := b
 \end{array}$$

- Dead code elimination: Code that computes values for names that will be dead i.e., never subsequently used can be removed.
- Renaming of temporary variables
- Interchange of two independent adjacent statements

Algebraic Transformations

Algebraic identities represent other important class optimizations on basic blocks. For example, we may apply arithmetic identities, such as $x + 0 = 0 + x = x$,

$$\begin{aligned}
 x * 1 &= 1 * x = x \\
 x - 0 &= x \\
 x/1 &= x
 \end{aligned}$$

Next-Use Information

- Next-use info used in code generation and register allocation.
- Remove variables from registers if not used.
- Statement of the form $A = B$ or C defines A and uses B and C .
- Scan each basic block backwards.
- Assume all temporaries are dead or exit and all user variables are live or exit.

Algorithm to compute next use information

Suppose we are scanning

```

i: x := y op z
    
```

in backward scan

- attach to i , information in symbol table about x, y, z .
- set x to not live and no next-use in symbol table
- set y and z to be live and next-use in symbol table.

Consider the following code:

```

1: t1 = a * a
2: t2 = a * b
3: t3 = 2 * t2
4: t4 = t1 + t2
5: t5 = b * b
6: t6 = t4 + t5
7: x = t6
    
```

Statements:

- 7: no temporary is live
- 6: t_6 : use (7) t_4, t_5 not live
- 5: t_5 : use (6)
- 4: t_4 : use (6), t_1, t_3 not live
- 3: t_3 : use (4) t_2 not live

2: t_2 : use (3)

1: t_1 : use (4)

Symbol Table:

t_1 dead use in 4

t_2 dead use in 3
 t_3 dead use in 4
 t_4 dead use in 6
 t_5 dead use in 6
 t_6 dead use in 7

The six temporaries in the basic block can be packed into two locations t_1 and t_2 :

1: $t_1 = a * a$
 2: $t_2 = a * b$
 3: $t_2 = 2 * t_2$
 4: $t_1 = t_1 + t_2$
 5: $t_2 = b * b$
 6: $t_1 = t_1 + t_2$
 7: $x = t_1$

Code Generator

- Consider each statement
- Remember if operand is in a register
- Descriptors are used to keep track of register contents and address for names
- There are 2 types of descriptors
 1. Register Descriptor
 2. Address Descriptor

Register Descriptor

Keep track of what is currently in each register. Initially all registers are empty.

Address Descriptors

- Keep track of location where current value of the name can be found at runtime.
- The location might be a register, stack, memory address or a set of all these.

Issues in design of code generation The issues in the design of code generation are

1. Intermediate representation
2. Target code
3. Address mapping
4. Instruction set.

Intermediate Representation It is represented in post fix, 3-address code (or) quadruples and syntax tree (or) DAG.

Target Code The Target Code could be absolute code, relocatable machine code (or) assembly language code. Absolute code will execute immediately as it is having fixed address relocatable, requires linker and loader to get the code from appropriate location for the assembly code, assemblers are required to convert it into machine level code before execution.

Address mapping In this, mapping is defined between intermediate representations to target code address.

It is based on run time environment like static, stack or heap.

Instruction set It should provide a complete set in such a way that all its operations can be implemented.

Code Generation Algorithm

For each three address statement $x = y \text{ op } z$ do

- Invoke a function getreg to determine location L where x must be stored. Usually L is a register.
- Consult address descriptor of y to determine y' . Prefer a register for y' . If value of y is not already in L generate $\text{MOV } y', L$.
- Generate $\text{OP } z', L$
 Again prefer a register for z . Update address descriptor of x to indicate x is in L . If L is a register update its descriptor to indicate that it contains x and remove x from all other register descriptors.
- If current value of y and/or z have no next use and are dead or exit from block and are in registers then change the register descriptor to indicate that it no longer contain y and /or z .

Function getreg

1. If y is in register and y is not live and has no next use after $x = y \text{ OP } z$ then return register of y for L .
2. Failing (1) return an empty register.
3. Failing (2) if x has a next use in the block or OP requires register then get a register R , store its contents into M and use it.
4. Else select memory location x as L .

Example: $D := (a - b) + (a - c) + (a - c)$

Stmt	Code Generated	reg desc	addr desc
$t = a - b$	$\text{MOV } a, R_0$ $\text{SUB } b, R_0$	R_0 contains t	t in R_0
$u = a - c$	$\text{MOV } a, R_1$ $\text{SUB } c, R_1$	R_0 contains t R_1 contains u	t in R_0 u in R_1
$v = t + u$	$\text{ADD } R_1, R_0$	R_0 contains v R_1 contains u	v in R_0 u in R_1
$d = v + u$	$\text{ADD } R_1, R_0$ $\text{MOV } R_0, d$	R_0 contains d	d in R_0 d in R_0 and memory

Conditional Statements

Machines implement conditional jumps in 2 ways:

1. Based on the value of the designated register (R)
 Branch if values of R meets one of six conditions.
 - (i) Negative
 - (ii) Zero
 - (iii) Positive
 - (iv) Non-negative
 - (v) Non-zero
 - (vi) Non-positive

Example: Three address statement: if $x < y$ goto z
 It can be implemented by subtracting y from x in R , then jump to z if value of R is negative.

- Based on a set of **condition codes** to indicate whether last quantity computed or loaded into a location is negative (or) Zero (or) Positive.
 - compare instruction set codes without actually computing the value.

Example: CMP x, y
 CJL Z .

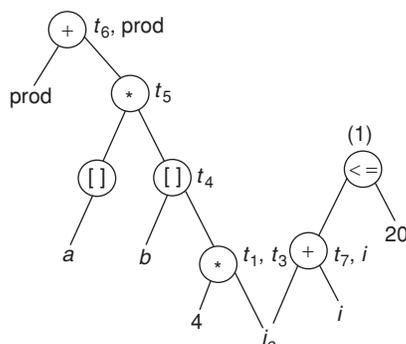
- Maintains a condition code descriptor, which tells the name that last sets the condition codes.

Example: $X := y + z$
 If $x < 0$ goto z
 By
 MOV y, R_o
 ADD z, R_o
 MOV R_o, x
 CJN z .

DAG REPRESENTATION OF BASIC BLOCKS

- DAGS are useful data structures for implementing transformations on basic blocks.
- Tells, how value computed by a statement is used in subsequent statements.
- It is a good way of determining common sub expressions.
- A DAG for a basic block has following labels on the nodes:
 - Leaves are labeled by unique identifiers, either variable names or constants.
 - Interior nodes are labeled by an operator symbol.
 - Nodes are also optionally given as a sequence of identifiers for labels.

- Example:** 1: $t_1 := 4 * i$
 2: $t_2 := a [t_1]$
 3: $t_3 := 4 * i$
 4: $t_4 := b [t_3]$
 5: $t_5 := t_2 * t_4$
 6: $t_6 := \text{prod} + t_5$
 7: $\text{prod} := t_6$
 8: $t_7 := i + 1$
 9: $i := t_7$
 10: if $i \leq 20$ got (1)



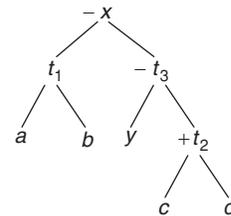
Code Generation from DAG:

$S_1 = 4 * i$	$S_1 = 4 * i$
$S_2 = \text{add}(A) - 4$	$S_2 = \text{add}(A) - 4$
$S_3 = S_2 [S_1]$	$S_3 = S_2 [S_1]$
$S_4 = 4 * i$	
$S_5 = \text{add}(B) - 4$	$S_5 = \text{add}(B) - 4$
$S_6 = S_5[S_4]$	$S_6 = S_5[S_4]$
$S_7 = S_3 * S_6$	$S_7 = S_3 * S_6$
$S_8 = \text{prod} + S_7$	$\text{prod} = \text{prod} + S_7$
$\text{prod} = S_8$	
$S_9 = I + 1$	
$I = S_9$	$I = I + 1$
if $I \leq 20$ got (1)	if $I \leq 20$ got (1)

Rearranging order of the code

Consider the following basic block

$t_1 := a + b$
 $t_2 := c + d$
 $t_3 := e - t_2$
 $x = t_1 - t_3$ and its DAG



Three address code for the DAG:
 (Assuming only two registers are available)

```
MOV a, R_o
ADD b, R_o
MOV c, R_1
MOV R_o, t_1      Register Spilling
MOV e, R_o        Register Reloading
SUB R_1, R_o
MOV t_1, R_1
SUB R_o, R_1
MOV R_1, x
```

Rearranging the code as

$t_2 := c + d$
 $t_3 := e - t_2$
 $t_1 := a + b$
 $x = t_1 - t_3$

The rearrangement gives the code:

```
MOV c, R_o
ADD d, R_o
MOV e, R_1
SUB R_o, R_1
```

```
MOV a, Ro
ADD b, Ro
SUB R1, Ro
MOV R1, x
```

Error detection and Recovery The errors that arise while compiling

1. Lexical errors
2. Syntactic errors
3. Semantic errors
4. Run-time errors

Lexical errors If the variable (or) constants are declared (or) defined, not according to the rules of language, special symbols are included which were not part of the language, etc is the lexical error.

Lexical analyzer is constructed based on pattern recognizing rules to form a token, when a source code is made into tokens and if these tokens are not according to rules then errors are generated.

Consider a c program statement

```
printf ('Hello World');
```

Main printf, (, ' , Hello world, ' ,); are tokens.

Printf is not recognizable pattern, actually it should be printf. It generates an error.

Syntactic error These errors include semi colons, missing braces etc. which are according to language rules.

The parser reports the errors

Semantic errors This type of errors arises, when operation is performed over incompatible type of variables, double declaration, assigning values to undefined variables etc.

Runtime errors The Runtime errors are the one which are detected at runtime. These include pointers assigned with NULL values and accessing a variable which is out of its boundary, unlegible arithmetic operations etc.

After the detection of errors. The following recovery strategies should be implemented.

1. Panic mode recovery
2. Phrase level recovery
3. Error production
4. Global correction.

PEEPHOLE OPTIMIZATION

- Target code often contains redundant instructions and suboptimal constructs.
- Improving the performance of the target program by examining a short sequence of target instructions (peephole) and replacing these instructions by a shorter or faster sequence is peephole optimization.
- The peephole is a small, moving window on the target program. Some well known peephole optimizations are

1. Eliminating redundant instructions
2. Eliminating unreachable code
3. Flow of control optimizations or Eliminating jumps over jumps
4. Algebraic simplifications
5. Strength reduction
6. Use of machine idioms

Elimination of Redundant Loads and stores

```
Example 1: (1) MOV Ro, a
           (2) MOV a, Ro
```

We can delete instruction (2), because the value of a is already in R_o.

```
Example 2: Load x, R0
           Store R0, x
```

If no modifications to R₀/x then store instruction can be deleted

```
Example 3: (1) Load x, R0
           (2) Store R0, x
```

```
Example 4: (1) store R0, x
           (2) Load x, R0
```

Second instruction can be deleted from both examples 3 and 4.

```
Example 5: Store R0, x
           Load x, R0
```

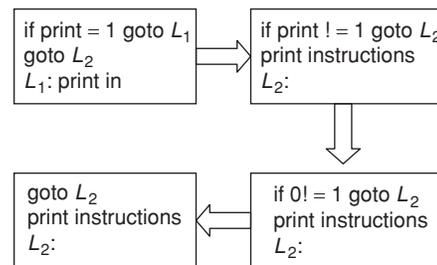
Here load instruction can be deleted.

Eliminating Unreachable code

An unlabeled instruction immediately following and unconditional jump may be removed.

- May be produced due to debugging code introduced during development.
- May be due to updates in programs without considering the whole program segment.

Example: Let print = 0



In all of the above cases print instructions are unreachable. ∴ Print instructions can be eliminated.

```
Example: goto L2
...
L2:
```

Flow of control optimizations The unnecessary jumps can be eliminated.

- Jumps like:
- Jumps to jumps,
 - Jumps to conditional jumps,
 - Conditional jumps to jumps.

3. Consider the program fragment:

```
sum = 0;
For (i = 1; i <= 20; i++)
sum = sum + a[i] + b[i];
```

How many instructions are there in the three-address code for this?

- (A) 15 (B) 16
(C) 17 (D) 18

4. Suppose the instruction set of the processor has only two registers. The code optimization allowed is code motion. What is the minimum number of spills to memory in the compiled code?

```
c = a + b;
d = c * a;
e = c + a;
x = c * c;
If (x > a)
{
y = a * a;
Else
{
d = d * d; e = e * e;
}
```

- (A) 0 (B) 1
(C) 2 (D) 3

5. What is the minimum number of registers needed to compile the above problem's code segment without any spill to memory?

- (A) 3 (B) 4
(C) 5 (D) 6

6. Convert the following expression into postfix notation:

$$a = (-a + 2 * b) / a$$

- (A) $aa - 2b * + a / =$ (B) $a - 2ba * / + =$
(C) $a2b * a / +$ (D) $a2b - * a / +$

7. In the quadruple representation of the following program, how many temporaries are used?

```
int a = 2, b = 8, c = 4, d;
For (j = 0; j <= 10; j++)
```

```
a = a * (j * (b/c));
```

```
d = a * (j * (b/c));
```

- (A) 4 (B) 7
(C) 8 (D) 10

8. Let $A = 2, B = 3, C = 4$ and $D = 5$, what is the final value of the prefix expression: $+ * AB - CD$

- (A) 5 (B) 10
(C) -10 (D) -5

9. Which of the following is a valid expression?

- (A) $BC * D - +$ (B) $* ABC -$
(C) $BBB *** - +$ (D) $- */ bc$

10. What is the final value of the postfix expression $BCD AD - + - +$ where $A = 2, B = 3, C = 4, D = 5$?

- (A) 5 (B) 4
(C) 6 (D) 7

11. Consider the expression $x = (a + b)^* - C/D$. In the quadruple representation of this expression in which instruction '/' operation is used?

- (A) 3rd (B) 4th
(C) 5th (D) 8th

12. In the triple representation of $x = (a + b)^* - c/d$, in which instruction $(a + b)^* - c/d$ result will be assigned to x ?

- (A) 3rd (B) 4th
(C) 5th (D) 8th

13. Consider the three address code for the following program:

```
While (A < C and B > D) do
```

```
If (A == 1) then C = C + 1;
```

```
Else
```

```
While (A <= D) do
```

```
A = A + 3;
```

How many temporaries are used?

- (A) 2 (B) 3
(C) 4 (D) 0

14. Code generation can be done by

- (A) DAG (B) Labeled tree
(C) Both (A) and (B) (D) None of these

15. Live variables analysis is used as a technique for

- (A) Code generation (B) Code optimization
(C) Type checking (D) Run time management

Practice Problems 2

Directions for questions 1 to 19: Select the correct alternative from the given choices

1. Match the correct code optimization technique to the corresponding code:

- | | | |
|------------------|-------------------------|---------------------------|
| (i) $i = i * 1$ | $\Rightarrow j = 2 * i$ | (p) Reduction in strength |
| $j = 2 * i$ | | |
| (ii) $A = B + C$ | $\Rightarrow A = B + C$ | (q) Machine Idioms |
| $D = 10 + B + C$ | $D = 10 + A$ | |

- (iii) For $i = 1$ to 10 \Rightarrow for $i = 1$ to 10 (r) Common sub expression elimination.
 $A[j] = B + C$ $t = B + C$
 $A[j] = t$

- (iv) $x = 2 * y \Rightarrow y << 2$; (s) Code motion

- (A) i - r, iii - s, iv - p, ii - q
(B) i - q, ii - r, iii - s, iv - p
(C) i - s, iii - p, iii - q, iv - r
(D) i - q, ii - p, iii - r, iv - s

2. What will be the optimized code for the following expression represented in DAG?

$$a = q * -r + q * -r$$

- (A) $t_1 = -r$
 $t_2 = q * t_1$
 $t_3 = a * t_1$
 $t_4 = t_2 + t_3$
 $a = t_4$
- (B) $t_1 = -r$
 $t_2 = q * t_1$
 $t_3 = t_2 + t_2$
 $a = t_3$
- (C) $t_1 = -r$
 $t_2 = q$
 $t_3 = t_1 * t_2$
 $t_4 = t_3 + t_3$
 $a = t_4$
- (D) All of these

3. In static allocation, names are bound to storage at _____ time.

- (A) Compile (B) Runtime
 (C) Debugging (D) Both (A) and (B)

4. The actual parameters are evaluate d and their r-values are passed to the called procedure is known as

- (A) call-by-reference
 (B) call-by-name
 (C) call-by-value
 (D) copy-restore

5. If the expression $-(a + b) * (c + d) + (a + b + c)$ is translated into quadruple representation, then how many temporaries are required?

- (A) 5 (B) 6
 (C) 7 (D) 8

6. If the above expression is translated into triples representation, then how many instructions are there?

- (A) 6 (B) 10
 (C) 5 (D) 8

7. In the indirect triple representation for the expression $A = (E/F) * (C - D)$. The first pointer address refers to

- (A) $C - D$
 (B) E/F
 (C) Both (A) and (B)
 (D) $(E/F) * (C - D)$

8. For the given assembly language, what is the cost for it?

- MOV b, a
 ADD c, a
- (A) 3 (B) 4
 (C) 6 (D) 2

9. Consider the expression

$((4 + 2 * 3 + 7) + 8 * 5)$. The polish postfix notation for this expression is

- (A) $423* + 7 + 85*+$ (B) $423* + 7 + 8 + 5*$
 (C) $42 + 37 + *85*+$ (D) $42 + 37 + 85**+$

Common data for questions 10 to 15: Consider the following basic block, in which all variables are integers, and ** denotes exponentiation.

$$a: = b + c$$

$$z: = a ** 2$$

$$x: = 0 * b$$

$$y: = b + c$$

$$w: = y * y$$

$$u: = x + 3$$

$$v: = u + w$$

Assume that the only variables that are live at the exit of this block are v and z . In order, apply the following optimization to this basic block.

10. After applying algebraic simplification, how many instructions will be modified?

- (A) 1 (B) 2
 (C) 4 (D) 5

11. After applying common sub expression elimination to the above code. Which of the following are true?

- (A) $a: = b + c$ (B) $y: = a$
 (C) $z = a + a$ (D) None of these

12. Among the following instructions, which will be modified after applying copy propagation?

- (A) $a: = b + c$ (B) $z: = a * a$
 (C) $y: = a$ (D) $w: = y * y$

13. Which of the following is obtained after constant folding?

- (A) $u: = 3$ (B) $v: = u + w$
 (C) $x: = 0$ (D) Both (A) and (C)

14. In order to apply dead code elimination, what are the statements to be eliminated?

- (A) $x = 0$
 (B) $y = b + c$
 (C) Both (A) and (B)
 (D) None of these

15. How many instructions will be there after optimizing the above result further?

- (A) 1 (B) 2
 (C) 3 (D) 4

16. Consider the following program:

$$L_0: e: = 0$$

$$b: = 1$$

$$d: = 2$$

$$L_1: a: = b + 2$$

$$c: = d + 5$$

$$e: = e + c$$

$$f: a * a$$

$$\text{If } f < c \text{ goto } L_3$$

$$L_2: e: = e + f$$

$$\text{goto } L_4$$

$$L_3: e: = e + 2$$

$$L_4: d: = d + 4$$

$$b: = b - 4$$

$$\text{If } b! = d \text{ goto } 4$$

$$L_5:$$

How many blocks are there in the flow graph for the above code?

- (A) 5
- (B) 6
- (C) 8
- (D) 7

17. A basic block can be analyzed by

- (A) Flow graph
- (B) A graph with cycles
- (C) DAG
- (D) None of these

18. In call by value the actual parameters are evaluated. What type of values is passed to the called procedure?

- (A) l-values
- (B) r-values
- (C) Text of actual parameters
- (D) None of these

19. Which of the following is FALSE regarding a Block?

- (A) The first statement is a leader.
- (B) Any statement that is a target of conditional / unconditional goto is a leader.
- (C) Immediately next statement of goto is a leader.
- (D) The last statement is a leader.

PREVIOUS YEARS' QUESTIONS

1. The least number of temporary variables required to create a three-address code in static single assignment form for the expression $q + r/3 + s - t * 5 + u * v/w$ is _____ [2015]

2. Consider the intermediate code given below.

- (1) $i = 1$
- (2) $j = 1$
- (3) $t_1 = 5 * i$
- (4) $t_2 = t_1 + j$
- (5) $t_3 = 4 * t_2$
- (6) $t_4 = t_3$
- (7) $a[t_4] = -1$
- (8) $j = j + 1$
- (9) if $j \leq 5$ goto (3)
- (10) $i = i + 1$
- (11) if $i < 5$ goto (2)

The number of nodes and edges in the control-flow-graph constructed for the above code, respectively, are [2015]

- (A) 5 and 7
- (B) 6 and 7
- (C) 5 and 5
- (D) 7 and 8

3. Consider the following code segment. [2016]

```
x = u - t;
y = x * v;
x = y + w;
y = t - z;
y = x * y;
```

The minimum number of total variables required to convert the above code segment to *static single assignment form* is _____.

4. What will be the output of the following pseudo-code when parameters are passed by reference and dynamic scoping is assumed? [2016]

$a = 3;$

```
void n(x) { x = x * a; print (x); }
```

```
void m(y) { a = 1; a = y - a; n(a) ; print (a) }
```

```
void main( ) { m(a); }
```

- (A) 6,2
- (B) 6,6
- (C) 4,2
- (D) 4,4

5. Consider the following intermediate program in three address code

$$p = a - b$$

$$q = p * c$$

$$p = u * v$$

$$q = p + q$$

Which one of the following corresponds to a *static single assignment* form of the above code? [2017]

- (A) $p_1 = a - b$
- (B) $p_3 = a - b$
- $q_1 = p_1 * c$
- $q_4 = p_3 * c$
- $p_1 = u * v$
- $p_4 = u * v$
- $q_1 = p_1 + q_1$
- $q_5 = p_4 + q_4$
- (C) $p_1 = a - b$
- (D) $p_1 = a - b$
- $q_1 = p_2 * c$
- $q_1 = p * c$
- $p_3 = u * v$
- $p_2 = u * v$
- $q_2 = p_4 + q_3$
- $q_2 = p + q$

ANSWER KEYS**EXERCISES****Practice Problems 1**

1. D 2. D 3. C 4. C 5. B 6. A 7. B 8. A 9. A 10. A
11. B 12. C 13. A 14. C 15. B

Practice Problems 2

1. B 2. B 3. A 4. B 5. B 6. A 7. B 8. C 9. A 10. A
11. B 12. D 13. A 14. C 15. C 16. A 17. C 18. B 19. D

Previous Years' Questions

1. 8 2. B 3. 10 4. D 5. B

Chapter 4

Code Optimization

LEARNING OBJECTIVES

- ☞ Code optimization basics
- ☞ Principle sources of optimization
- ☞ Loop invariant code motion
- ☞ Strength reduction on induction variables
- ☞ Loops in flow graphs
- ☞ Pre-header
- ☞ Global data flow analysis
- ☞ Definition and usage of variables
- ☞ Use-definition (u-d) chaining
- ☞ Data flow equations

CODE OPTIMIZATION BASICS

The process of improving the intermediate code and the target code in terms of both speed and the amount of memory required for execution is known as code optimization.

Compilers that apply code-improving transformations are called optimizing compilers.

Properties of the transformations of an optimizing compiler are

1. A transformation must preserve the meaning of programs.
2. It must speed up programs by a measurable amount.
3. A transformation must be worth the effort.

Places for improvements

1. Source Code:
 - User can
 - profile a program
 - change an algorithm
 - transform loops
2. Intermediate code can be improved by improving
 - Loops
 - Procedure calls
 - Address calculations
3. Target code can be improved by
 - Using registers
 - Selecting instructions
 - Peephole transformations

Optimizing compiler organization

This applies

- Control flow analysis
- Data flow analysis
- Transformations

Issues in design of code optimization The issues in the design of code optimization are

1. Target machine characteristics
2. Target CPU architecture
3. Functional units

Target machine Optimization is done, according to the target machine characteristics. Altering the machine description parameters, one can optimize single piece of compiler code.

Target CPU architecture The issues to be considered for the optimization with respect to CPU architecture

1. Number of CPU registers
2. RISC Instruction set
3. CISC instruction set
4. Pipelining

Functional units Based on number of functional units, optimization is done. So that instructions can be executed simultaneously.

PRINCIPLE SOURCES OF OPTIMIZATION

Some code improving transformation is Local transformations and some are Global transformations.

Local Transformations can be performed by looking only at a statement in a basic block. Otherwise it is global transformation.

Function Preserving Transformations

These transformations improve the program without changing the function it computes. Some of these transformations are

1. Common sub expression elimination
2. Copy propagation
3. Dead-code elimination
4. Loop optimization
 - Code motion
 - Induction variable elimination
 - Reduction in strength

Common sub expression elimination The process of identifying common sub expressions and eliminating their computation multiple times is known as common sub expression elimination.

Example: Consider the following program segment:

```
int sum_n, sum_n2, sum_n3;
int sum (int n)
{
  Sum_n = ((n)*(n+1))/2;
  sum_n2 = ((n)*(n+1)*(2n+1))/6;
  sum_n3 = (((n)*(n+1))/2)*(((n)*(n+1))/2);
}
```

Three Address code for the above input is

- (0) Proc-begin sum
- (1) $t_0 := n + 1$
- (2) $t_1 := n * t_0$
- (3) $t_2 := t_1 / 2$
- (4) $sum_n = t_2$
- (5) $t_3 := n + 1$
- (6) $t_4 := n * t_3$
- (7) $t_5 := 2 * n$
- (8) $t_6 := t_5 + 1$
- (9) $t_7 := t_4 * t_6$
- (10) $t_8 := t_7 / 6$
- (11) $sum_n2 := t_8$
- (12) $t_9 := n + 1$
- (13) $t_{10} := n * t_9$
- (14) $t_{11} := t_{10} / 2$
- (15) $t_{12} := n + 1$
- (16) $t_{13} := n * t_{12}$
- (17) $t_{14} := t_{13} / 2$
- (18) $t_{15} := t_{11} * t_{14}$
- (19) $sum_n3 := t_{15}$
- (20) label L_0
- (21) Proc end sum

The computations made in quadruples

(1) – (3), (12) – (14), (15) – (17) are essentially same. That is, $((n)*(n+1))/2$ is computed.

It is the common sub expression.

This common sub expression is computed four times in the above example.

It is possible to optimize the code to have common sub expressions computed only once and then reuse the computed values further.

∴ Optimized intermediate code will be

- (0) proc-begin sum
- (1) $t_0 := n + 1$
- (2) $t_1 := n * t_0$
- (3) $sultan := t_1 / 2$
- (4) $t_5 := 2 * n$
- (5) $t_6 := t_5 + 1$
- (6) $t_7 := t_1 * t_6$
- (7) $sum_n2 := t_7 / 6$
- (8) $sum_n3 := sum_n * sum_n$
- (9) proc-end sum

Constant folding The constant expressions in the input source are evaluated and replaced by the equivalent values at the time of compilation.

For example $10*3$, $6 + 101$ are constant expressions and they are replaced by 30, 107 respectively.

Example: Consider the following ‘C’ code:

```
int arr1 [10];
int main ( )
{
  arr1 [0] = 3;
  arr1 [1] = 4;
}
```

Unoptimized three address code equivalent to the above ‘C’ code is

- (0) proc-begin main
- (1) $t_0 := 0 * 4$
- (2) $t_1 := \&arr1$
- (3) $t_1 [t_0] := 3$
- (4) $t_2 := 1 * 4$
- (5) $t_3 := \&arr1$
- (6) $t_3 [t_2] := 4$
- (7) Label L_0
- (8) Proc – end main

In the above code, $0*4$ is a constant expression its value = 0. $1*4$ is a constant expression, its value = 4.

∴ After applying constant folding, optimized code will be

- (0) proc-begin main
- (1) $t_0 := 0$
- (2) $t_1 := \&arr1$
- (3) $t_1 [t_0] := 3$
- (4) $t_2 := 4$

- (5) $t_3 := \&arr1$
- (6) $t_3[t_2] := 4$
- (7) label L_0
- (8) proc – end main

Copy propagation In copy propagation, if there is an expression $x = y$ then use the variable ‘y’ instead of ‘x’. This propagated in the statements following $x = y$.

Example: In the previous example, there are two copy statements.

- (1) $t_0 = 0$
- (2) $t_2 = 4$

After applying copy propagation, the optimized code will be

- (0) proc-begin main
- (1) $t_0 := 0$
- (2) $t_1 := \&arr1$
- (3) $t_1[0] := 3$
- (4) $t_2 := 4$
- (5) $t_3 := \&arr1$
- (6) $t_3[4] := 4$
- (7) Label L_0
- (8) proc-end main

In the three address code shown above, quadruples (1) and (4) are no longer used in any of the following statements.

∴ (1) and (4) can be eliminated.

Three address code after dead store elimination

- (0) proc-begin main
- (1) $t_1 := \&arr1$
- (2) $t_1[0] := 3$
- (3) $t_3 := \&arr1$
- (4) $t_3[4] := 4$
- (5) Label L_0
- (6) proc-end main

In the above example, we are propagating constant values. It is also known as constant propagation.

Variable propagation Propagating another variable instead of the existing one is known as variable propagation.

Example: `int func(int a, int b, int c)`

```
{
    int d, e, f;
    d = a;
    If (a > 10)
    {
        e = d + b;
    }
    Else
    {
        e = d + c;
    }
    f = d*e;
    return (f);
}
```

Three address code (unoptimized):

- (0) proc-begin func
- (1) $d := a$
- (2) if $a > 10$ goto L_0
- (3) goto L_1
- (4) label : L_0
- (5) $e := d + b$
- (6) goto L_2
- (7) label : L_1
- (8) $e := d + c$
- (9) label : L_2
- (10) $f := d * e$
- (11) return f
- (12) goto L_3
- (13) label : L_3
- (14) proc-end func

Three address code after variable (copy) propagation:

- (0) proc-begin func
- (1) $d := a$
- (2) If $a > 10$ goto $.L_0$
- (3) goto L_1
- (4) label: L_0
- (5) $e := a + b$
- (6) goto L_2
- (7) label: L_1
- (8) $e := a + c$
- (9) label: L_2
- (10) $f := a * e$
- (11) return f
- (12) goto L_3
- (13) label: L_3
- (14) proc-end func

After dead store elimination:

In the above code (1) $d := a$ is no more used

∴ Eliminate the dead store $d := a$

- (0) proc-begin func
- (1) If $a > 10$ goto L_0
- (2) goto L_1
- (3) label: L_0
- (4) $e := a + b$
- (5) goto L_2
- (6) label: L_1
- (7) $e := a + c$
- (8) label: L_2
- (9) $f := a * e$
- (10) return f
- (11) goto L_3
- (12) label: L_3
- (13) proc-end func

Dead code elimination Eliminating the code that never gets executed by the program is known as Dead code elimination. It reduces the memory required by the program

Example: Consider the following Unoptimized Intermediate code:

- (0) proc-begin func
- (1) debug: = 0
- (2) If debug == 1 goto L_0
- (3) goto L_1
- (4) label: L_0
- (5) param c
- (6) param b
- (7) param a
- (8) param lcl
- (9) call printf 16
- (10) retrieve to
- (11) label: L_1
- (12) $t_1 := a + b$
- (13) $t_2 := t_1 + c$
- (14) $v_1 := t_2$
- (15) Return v_1
- (16) goto L_2
- (17) label: L_2
- (18) proc-end func

In copy propagation, debug is replaced with 0, wherever debug is used after that assignment.

∴ Statement 2 will be changed as

If 0 == 1 goto L_0

0 == 1, always returns false.

∴ The control cannot flow to label: L_0

This makes the statements (4) through (10) as dead code. (2) Can also be removed as part of dead code elimination. (1) Cannot be eliminated, because ‘debug’ is a global variable. The optimized code after elimination of dead code is shown below.

- (0) proc-begin func
- (1) debug: = 0
- (2) goto L_1
- (3) label: L_1
- (4) $t_1 := a + b$
- (5) $t_2 := t_1 + c$
- (6) $v_1 := t_2$
- (7) return v_1
- (8) goto L_2
- (9) label: L_2
- (10) proc-end func

Algebraic transformations We can use algebraic identities to optimize the code further. For example

Additive Identity: $a + 0 = a$

Multiplicative Identity: $a * 1 = a$

Multiplication with 0: $a * 0 = 0$

Example: Consider the following code fragment:

```
struct mystruct
{
int a [20];
int b;
} xyz;
int func(int i)
{
xyz.a[i] = 34;
}
```

The Unoptimized three address code:

- (0) proc-begin func
- (1) $t_0 := \&xyz$
- (2) $t_1 := 0$
- (3) $t_2 := i * 4$
- (4) $t_1 := t_2 + t_1$
- (5) $t_0[t_1] = 34$
- (6) label: L_0
- (7) proc-end func

Optimized code after copy propagation and dead code elimination is shown below:

The statement $t_1 := 0$ is eliminated.

- (0) proc-being func
- (1) $t_0 := \&xyz$
- (2) $t_2 := i * 4$
- (3) $t_1 := t_2 + 0$
- (4) $t_0[t_1] := 34$
- (5) label: L_0
- (6) proc-end func

After applying additive identity:

- (0) proc-begin func
- (1) $t_0 := \&xyz$
- (2) $t_2 := i * 4$
- (3) $t_1 := t_2$
- (4) $t_0[t_1] := 34$
- (5) label: L_0
- (6) proc-end func

After copy propagation and dead store elimination:

- (0) proc-begin func
- (1) $t_0 := \&xyz$
- (2) $t_2 := i * 4$
- (3) $t_0[t_2] := 34$
- (4) label: L_0
- (5) proc-end func

Strength reduction transformation This transformation replaces expensive operators by equivalent cheaper ones on the target machine.

For example $y := x * 2$ is replaced by $y := x + x$ as addition is less expensive than multiplication.

Similarly

Replace $y := x * 32$ by $y := x \ll 5$

Replace $y := x / 8$ by $y := x \gg 3$

Loop optimization We can optimize loops by

- (1) Loop invariant code motion transformation.
- (2) Strength reduction on induction variable transformation.

Loop invariant code motion

The statements within a loop that compute value, which do not vary throughout the life of the loop are called loop invariant statements.

Consider the following program fragment:

```
int a [100];
int func(int x, int y)
{
  int i;
  int n1, n2;
  i = 0;
  n1 = x*y;
  n2 = x - y;
  while (a[i] > (n1*n2))
  i = i + 1;
  return(i);
}
```

The Three Address code for above program is

- (0) proc-begin func
- (1) $i := 0$
- (2) $n_1 := x * y$
- (3) $n_2 := x - y$
- (4) label : L_0
- (5) $t_2 := i * 4$
- (6) $t_3 := \&arr$
- (7) $t_4 := t_3[t_2]$
- (8) $t_5 := n_1 * n_2$
- (9) if $t_4 > t_5$ goto L_1
- (10) goto L_2
- (11) label : L_1
- (12) $i := i + 1$
- (13) goto L_0
- (14) label : L_2
- (15) return i
- (16) goto L_3
- (17) label : L_3
- (18) proc-end func

In the above code statements (6) and (8) are invariant.

After loop invariant code motion transformation the code will be

- (0) proc-begin func
- (1) $i := 0$
- (2) $n_1 := x * y$
- (3) $n_2 := x - y$
- (4) $t_3 := \&arr$
- (5) $t_5 := n_1 * n_2$
- (6) label : L_0
- (7) $t_2 := i * 4$
- (8) $t_4 := t_3[t_2]$
- (9) if $t_4 > t_5$ goto L_1
- (10) goto L_2
- (11) label : L_1
- (12) $i := i + 1$
- (13) goto L_0
- (14) label : L_2
- (15) return i
- (16) goto L_3
- (17) label : L_3
- (18) proc-end func

Strength reduction on induction variables

Induction variable: A variable that changes by a fixed quantity on each of the iterations of a loop is an induction variable.

Example: Consider the following code fragment:

```
int i;
int a[20];
int func ( )
{
  while(i<20)
  {
    a[i] = 10;
    i = i + 1;
  }
}
```

The three-address code will be

- (0) proc-begin func
- (1) label : L_0
- (2) if $i < 20$ goto L_1
- (3) goto L_2
- (4) label : L_1
- (5) $t_0 := i * 4$
- (6) $t_1 := \&a$
- (7) $t_1[t_0] := 10$
- (8) $i := i + 1$
- (9) goto L_0
- (10) label : L_2
- (11) label : L_3
- (12) proc-end func

After reduction of strength the code will be
Here (5) $t_0 = i*4$ is moved out of the loop and (8) is followed by $t_0 = t_0 + 4$.

```
(0) proc-begin func
(0a)  $t_0 := i*4$ 
(1) label :  $L_0$ 
(2) if  $i < 20$  goto  $L_1$ 
(3) goto  $L_2$ 
(4) label: $L_1$ 
(5)
(6)  $t_1 := &a$ 
(7)  $t_1[t_0] := 10$ 
(8)  $i := i + 1$ 
(8a)  $t_0 := t_0 + 4$ 
(9) goto  $L_0$ 
(10) label :  $L_2$ 
(11) label :  $L_3$ 
(12) proc-end func
```

LOOPS IN FLOW GRAPHS

Loops in the code are detected during the data flow analysis by using the concept called ‘dominators’ in the flow graph.

Dominators

A node d of a flow graph dominates node n , if every path from the initial node to ‘ n ’ goes through ‘ d ’.

It is represented as $d \text{ dom } n$.

Notes:

1. Each and every node dominates itself.
2. Entry of the loop dominates all nodes in the loop.

Example: Consider the following code fragment:

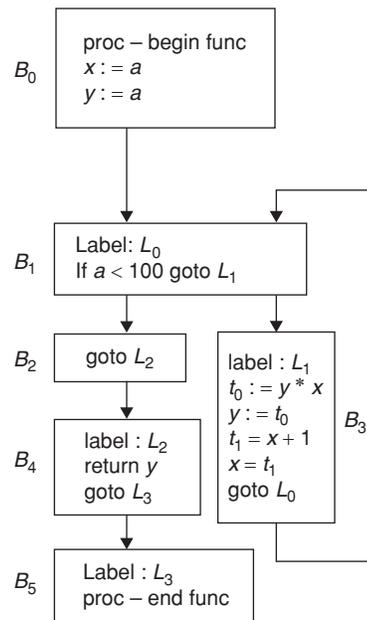
```
int func(int a)
{
int x, y;
x = a;
y = a;
While (a < 100)
{
y = y*x;
x = x+1;
}
return (y);
}
```

The Three Address code after local optimization will be

```
(0) proc-begin func
(1)  $x := a$ 
(2)  $y := a$ 
(3) label:  $L_0$ 
(4) if  $a < 100$  goto  $L_1$ 
(5) goto  $L_2$ 
```

```
(6) label:  $L_1$ 
(7)  $t_0 := y*x$ 
(8)  $y := t_0$ 
(9)  $t_1 := x + 1$ 
(10)  $x := t_1$ 
(11) goto  $L_0$ 
(12) label:  $L_2$ 
(13) return  $y$ 
(14) goto  $L_3$ 
(15) label:  $L_3$ 
(16) proc-end func
```

The Flow Graph for above code will be:



To reach B_2 , it must pass through B_1
∴ B_1 dominates B_2 . Also B_0 dominates B_2 .

dominators $[B_1] = \{B_0, B_1\}$ (or) dominators $[1] = \{0, 1\}$

The dominators for each of the nodes in the flow graph are

```
dominators [0] = {0}
dominators [1] = {0, 1}
dominators [2] = {0, 1, 2}
dominators [3] = {0, 1, 3}
dominators [4] = {0, 1, 2, 4}
dominators [5] = {0, 1, 2, 4, 5}
```

Edge

An edge in a flow graph represents a possible flow of control.

In the flow graph, B_0 to B_1 edge is represented as $0 \rightarrow 1$.

Head and tail: In the edge $a \rightarrow b$, the node b is called head and the node a is called as tail.

Back edges: There are some edges in which dominators [tail] contains the head.

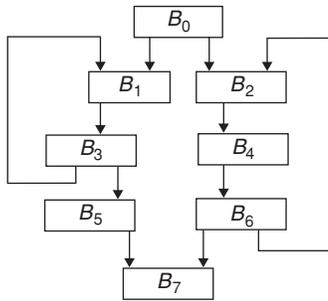
The presence of a back edge indicates the existence of a loop in a flow graph.

In the previous graph, $3 \rightarrow 1$ is a back edge.

Consider the following table:

Edge	Head	Tail	Dominators [head]	Dominators [tail]
$0 \rightarrow 1$	1	0	{0, 1}	{0}
$1 \rightarrow 2$	2	1	{0, 1, 2}	{0, 1}
$1 \rightarrow 3$	3	1	{0, 1, 3}	{0, 1}
$3 \rightarrow 1$	1	3	{0, 1}	{0, 1, 3}
$2 \rightarrow 4$	4	2	{0, 1, 2, 4}	{0, 1, 2}
$4 \rightarrow 5$	5	4	{0, 1, 2, 4, 5}	{0, 1, 2, 4}

Example: Consider below flow graph:



The dominators of each node are

- dominators [0] = {0}
- dominators [1] = {0, 1}
- dominators [2] = {0, 2}
- dominators [3] = {0, 1, 3}
- dominators [4] = {0, 2, 4}
- dominators [5] = {0, 1, 3, 5}
- dominators [6] = {0, 2, 4, 6}
- dominators [7] = {0, 7}

Edge	Head	Tail	Dominators [head]	Dominators [tail]
$0 \rightarrow 1$	1	0	{0, 1}	{0}
$0 \rightarrow 2$	2	0	{0, 2}	{0}
$1 \rightarrow 3$	3	1	{0, 1, 3}	{0, 1}
$3 \rightarrow 1$	1	3	{0, 1}	{0, 1, 3}
$3 \rightarrow 5$	5	3	{0, 1, 3, 5}	{0, 1, 3}
$5 \rightarrow 7$	7	5	{0, 7}	{0, 1, 3, 5}
$2 \rightarrow 4$	4	2	{0, 2, 4}	{0, 2}
$6 \rightarrow 2$	2	6	{0, 2}	{0, 2, 4, 6}
$4 \rightarrow 6$	6	4	{0, 2, 4, 6}	{0, 2, 4}
$6 \rightarrow 7$	7	6	{0, 7}	{0, 2, 4, 6}

Here $\{B_6, B_2, B_4\}$ form a loop (L_1), $\{B_3, B_1\}$ form another loop (L_2)

In a loop, the entry of the loop dominates all nodes in the loop.

Header of the loop The entry of the loop is also called as the header of the loop.

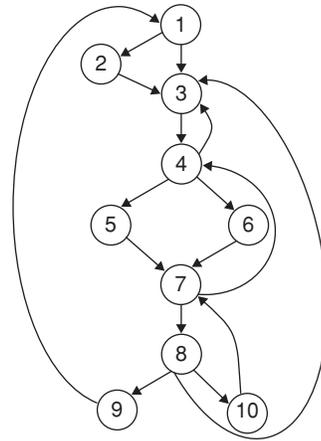
Loop exit block In loop L_1 can be exited from the basic block B_6 . It is called loop exit block. The block B_3 is the loop exit block for the loop L_2 . It is possible to have multiple exit blocks in a loop.

Dominator tree

A tree, which represents dominate information in the form of tree is a dominator tree. In this,

- The initial node is the root.
- Each node d dominates only its descendents in the tree.

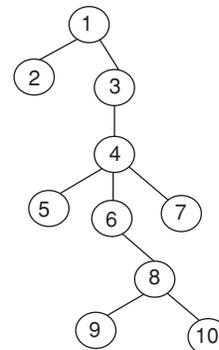
Consider the flow graph



The dominators of each node are

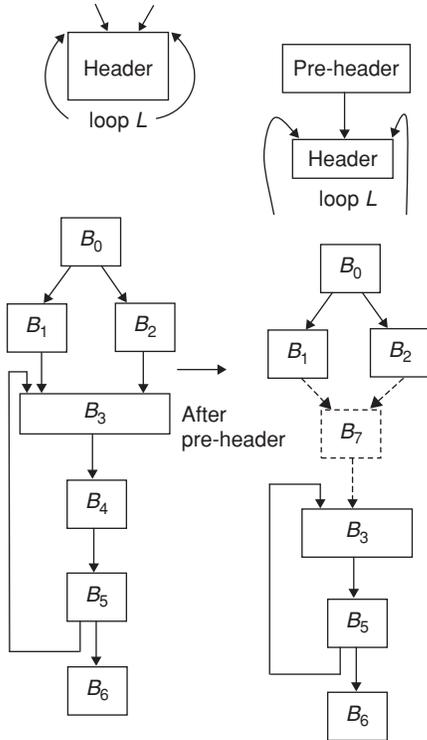
- dominators [1] = {1}
- dominators [2] = {1, 2}
- dominators [3] = {1, 3}
- dominators [4] = {1, 3, 4}
- dominators [5] = {1, 3, 4, 5}
- dominators [6] = {1, 3, 4, 6}
- dominators [7] = {1, 3, 4, 7}
- dominators [8] = {1, 3, 4, 7, 8}
- dominators [9] = {1, 3, 4, 7, 8, 9}
- dominators [10] = {1, 3, 4, 7, 8, 10}

The dominator tree will be:



Pre-header

A pre-header is a basic block introduced during the loop optimization to hold the statements that are moved from within the loop. It is a predecessor to the header block.

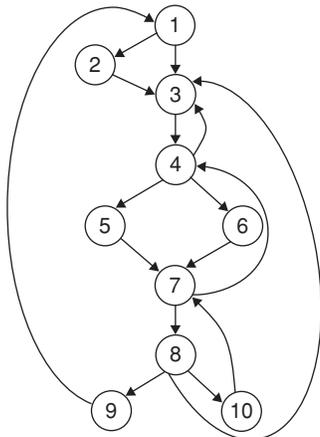


Reducible Flow Graphs

A flow graph G is reducible if and only if we can partition the edges into two disjoint groups:

- (1) Forward edges
- (2) Backward edges with the following properties.
 - (i) The forward edges form an acyclic graph in which every node can be reached from the initial node of G .
 - (ii) The back edges consist only of edges whose heads dominates their tails.

Example: Consider previous flow graph



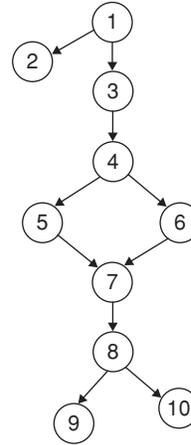
In the above flow graph, there are five back edges

$$4 \rightarrow 3, 7 \rightarrow 4, 8 \rightarrow 3, 9 \rightarrow 1 \text{ and } 10 \rightarrow 7$$

Remove all backedges.

The remaining edges must be the forward edges.

The remaining graph is acyclic.



∴ It is reducible.

GLOBAL DATAFLOW ANALYSIS

Point: A point is a place of reference that can be found at

1. Before the first statement in a basic block.
2. After the last statement in a basic block.
3. In between two adjacent statements within a basic block.

Example 1:

$$\begin{matrix} a^* = 10 \\ b^* = 20 \\ c^* = a * b \end{matrix} B_1$$

Here, In B_1 there are 4 points

Example 2:

$$B_1 \begin{matrix} \bullet P_1 - B_1 \\ \text{proc-begin func} \\ \bullet P_2 - B_1 \\ v_3 = v_1 + v_2 \\ \bullet P_3 - B_1 \\ \text{if } c > 100 \text{ goto } L_0 \\ \bullet P_4 - B_1 \end{matrix}$$

There is 4 point in the basic block B_1 , given by $P_1 - B_1$, $P_2 - B_1$, $P_3 - B_1$ and $P_4 - B_1$.

Path: A path is a sequence of points in which the control can flow.

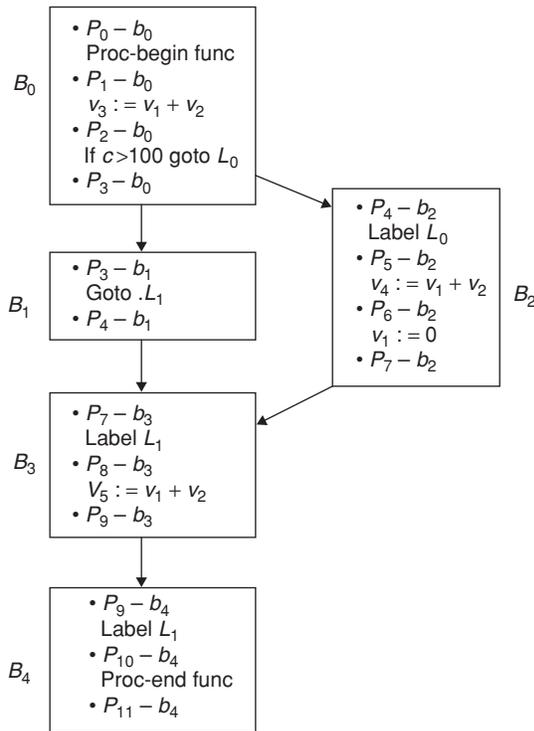
A path from P_1 to P_n is a sequence of points P_1, P_2, \dots, P_n such that for each i between 1 and $n-1$, either

- (a) P_i is the point immediately preceding a statement and P_{i+1} is the point immediately following that statement in the same block.

(OR)

- (b) P_i is the end of some block and P_{i+1} is the beginning of a successor block.

Example:



Path is between the points $P_0 - b_0$ and $P_6 - b_2$:

The sequence of points $P_0 - b_0, P_1 - b_0, P_2 - b_0, P_3 - b_0, P_4 - b_2, P_5 - b_2$ and $P_6 - b_2$.

Path between $P_3 - b_1$ and $P_6 - b_2$: There is no sequence of points.

Path between $P_0 - b_0$ and $P_7 - b_3$: There are two paths.

- (1) Path 1 consists of the sequence of points, $P_0 - b_0, P_1 - b_0, P_2 - b_0, P_3 - b_0, P_4 - b_1$ and $P_7 - b_3$.
- (2) Path 2 consists of the sequence of points $P_0 - b_0, P_1 - b_0, P_2 - b_0, P_3 - b_0, P_4 - b_2, P_5 - b_2, P_6 - b_2, P_7 - b_2$ and $P_7 - b_3$.

Definition and Usage of Variables

Definitions

It is either an assignment to the variable or reading of a value for the variable.

Use

Use of identifier x means any occurrence of x as an operand.

Example: Consider the statement

$$x = y + z;$$

In this statement some value is assigned to x . It defines x and used y and z values.

Global Data-Flow-Analysis

Data Flow Analysis (DFA) is a technique for gathering information about the possible set of values calculated at various points in a program.

- An example of a data-flow analysis is reaching definitions.
- A single way to perform data-flow analysis of program is to setup data flow equations for each node of the control flow graph.

Use definition (U-d) chaining

The use of a value is any point where that variable or constant is used in the right hand side of an assignment or is evaluating an expression.

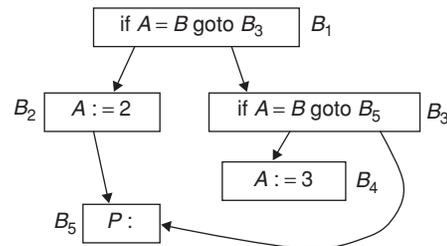
The definition of a value occurs implicitly at the beginning of the whole program for a variable.

A point is defined either prior to or immediately after a statement.

Reaching definitions

A definition of a variable A reaches a point P if there is a path in the flow graph from that definition to P , such that no other definitions of A appear on the path.

Example:



The definition $A := 3$ can reach point p in B_5 .

To determine the definitions that can reach a given program first assign distinct numbers to each definition, since it is associated with a unique quadruple.

- For each simple variable A , make a list of all definitions of A anywhere in the program.
- Compute two sets for each basic block B .

Gen $[B]$ is the set of generated definitions within block B and that reach the end of the block.

1. Kill $[B]$, which is the set of definitions outside of B that define identifiers that also have definitions within B .
2. IN $[B]$, which are all definitions reaching the point just before B 's first statement.

Once this is known, the definitions reaching any use of A within B are found by:

Let u be the statement being examined, which uses A .

1. If there are definitions of A within B before u , the last is the only one reaching u .
2. If there is no definition of A within B prior to u , those reaching u are in IN $[B]$.

Data Flow Equations

1. For all blocks B ,

$$OUT [B] = (IN [B] - KILL [B]) \cup GEN [B]$$

A definition d , reaches the end of B if

- (a) $d \in \text{IN}[B]$ and is not killed by B .
(or)
(b) d is generated in B and is not subsequently redefined here.

$$2. \text{IN}[B] = \text{U OUT}[P]$$

$\forall P$ preceding B

A definition reaches the beginning of B iff it reaches the end of one of its predecessors.

Computing U-d Chains

If a use of variable ' a ' is preceded in its block by a definition of ' a ', this is the only one reaching it.

If no such definition precedes its use, all definitions of ' a ' in $\text{IN}[B]$ are on its chain.

Uses of U-d Chains

- If the only definition of ' a ' reaching this statement involves a constant, we can substitute that constant for ' a '.
- If no definitions of ' a ' reaches this point, a warning can be given.
- If a definition reaches nowhere, it can be eliminated. This is part of dead code elimination.

EXERCISES

Practice Problems I

Directions for questions 1 to 15: Select the correct alternative from the given choices.

- Replacing the expression $2 * 3.14$ by 6.28 is
 - Constant folding
 - Induction variable
 - Strength reduction
 - Code reduction
- The expression $(a*b)*c$ op ... where 'op' is one of '+', '*', and '^' (exponentiation) can be evaluated on CPU with a single register without storing the value of $(a*b)$ if
 - 'op' is '+' or '*'
 - 'op' is '^' or '+'
 - 'op' is '^' or '*'
 - not possible to evaluate without storing
- Machine independent code optimization can be applied to
 - Source code
 - Intermediate representation
 - Runtime output
 - Object code
- In block B if S occurs in B and there is no subsequent assignment to y within B , then the copy statement $S: x = y$ is

(A) Generated	(B) Killed
(C) Blocked	(D) Dead
- If E was previously computed and the value of variable in E have not changed since previous computation, then an occurrence of an expression E is
 - Copy propagation
 - Common sub expression
 - Dead code
 - Constant folding
- In block B , if x or y is assigned there and s is not in B , then $s: x = y$ is

(A) Generated	(B) Killed
(C) Blocked	(D) Dead
- Given the following code


```
A = x + y;
B = x + y;
Then the corresponding optimized code as
_____
_____
C = x + y;
_____
A = C;
_____
B = C;
When will be optimized code pose a problem?
(A) When C is undefined.
(B) When memory is consideration.
(C) C may not remain same after some statements.
(D) Both (A) and (C).
```
- Can the loop invariant $X = A - B$ from the following code be moved out?


```
For i = 1 to 10
{
A = B * C;
X = A - B;
}
(A) No
(B) Yes
(C)  $X = A - B$  is not invariant
(D) Data insufficient
```
- If every path from the initial node goes through a particular node, then that node is said to be a

(A) Header	(B) Dominator
(C) Parent	(D) Descendant

Common data for questions 10 and 11: Consider the following statements of a block:

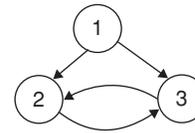
$a: = b + c$

$b: = a - d$

$c: = b + c$

$d: = a - d$

10. The above basic block contains, the value of b in 3rd statement is
 (A) Same as b in 1st statement
 (B) Different from b in 1st statement
 (C) 0
 (D) 1
11. The above basic block contains
 (A) Two common sub expression
 (B) Only one common sub expression
 (C) Dead code
 (D) Temporary variable
12. Find the induction variable from the following code:
 $A = -0.2;$
 $B = A + 5.0;$
 (A) A
 (B) B
 (C) Both A and B are induction variables
 (D) No induction variables
13. The analysis that cannot be implemented by forward operating data flow equations mechanism is
 (A) Interprocedural
 (B) Procedural
 (C) Live variable analysis
 (D) Data
14. Which of the following consist of a definition, of a variable and all the uses, U , reachable from that definition without any other intervening definitions?
 (A) Ud-chaining (B) Du-chaining
 (C) Spanning (D) Searching
15. Consider the graph



The graph is

- (A) Reducible graph
 (B) Non-reducible graph
 (C) Data insufficient
 (D) None of these

Practice Problems 2

Directions for questions 1 to 15: Select the correct alternative from the given choices.

1. In labeling algorithm, let n is a binary node and its children have L_1 and L_2 , if $L_1 = L_2$ then LABEL (n):
 (A) $L_1 - 1$ (B) $L_2 + 1$
 (C) $L_1 + L_1$ (D) $L_1 + 1$
2. The input for the code generator is a:
 (A) Tree at lexical level
 (B) Tree at semantic level
 (C) Sequence of assembly language instructions
 (D) Sequence of machine idioms
3. In labeling algorithm, let n is a binary node and its children have i_1 and i_2 , LABEL (n) if $i_1 \neq i_2$ is
 (A) Max (i_1, i_2)
 (B) $i_2 + 1$
 (C) $i_2 - 1$
 (D) $i_2 - i_1$
4. The following tries to keep frequently used value in a fixed register throughout a loop is:
 (A) Usage counts
 (B) Global register allocation
 (C) Conditional statement
 (D) Pointer assignment
5. Substitute y for x for copy statement $s: x = y$ if the following condition is met

- (A) Statements s may be the only definition of x reaching u
 (B) x is dead
 (C) y is dead
 (D) x and y are aliases

6. Consider the following code

```

for (i=0; i<m; i++)
{
for (j=0; j<m; j++)
If (i%2)
{
a = a + (14*j+5*i);
b = b + (9 + 4*j);
}
}
}
  
```

Which of the following is false?

- (A) There is a scope of common reduction in this code
 (B) There is a scope of strength reduction in this code.
 (C) There is scope of dead code elimination in this code
 (D) Both (A) and (C)
7. S_1 : In dominance tree, the initial node is the root.
 S_2 : Each node d dominates only its ancestors in the tree.
 S_3 : if $d \neq n$ and $d \text{ dom } n$ then $d \text{ dom } m$.
 Which of the statements is/are true?
 (A) S_1, S_2 are true
 (B) S_1, S_2 and S_3 are true

- (C) Only S_3 is true
 (D) Only S_1 is true
8. The specific task storage manager performs:
 (A) Allocation/Deallocation of storage to programs
 (B) Protection of storage area allocated to a program from illegal access by other programs in the system
 (C) The status of each program
 (D) Both (A) and (B)
9. Concept which can be used to identify loops is:
 (A) Dominators
 (B) Reducible graphs
 (C) Depth first ordering
 (D) All of these
10. A point cannot be found:
 (A) Between two adjacent statements
 (B) Before the first statement
 (C) After the last statement
 (D) Between any two statements
11. In the statement, $x = y*10 + z$; which is/are defined?
 (A) x (B) y
 (C) z (D) Both (B) and (C)
12. Consider the following program:

```
void main ( )
{
    int x, y;
    x = 3; y = 7;
    -----
    -----
    if (x<y)
    {
        int x;
```

```
{
int y;
y = 9;
-----
x = 2*y;
}
-----
-----
x = x + y;
printf ("%d", x);
}
-----
printf ("%d", x);
}
```

The output is

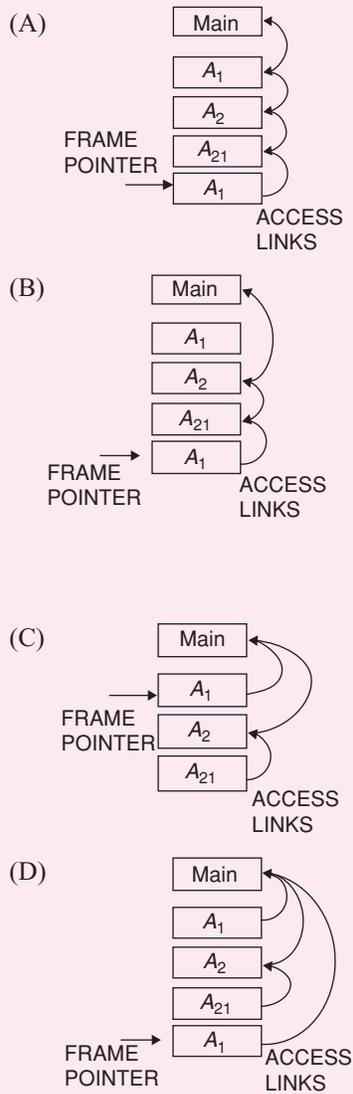
- (A) 3 – 25 (B) 25 – 3
 (C) 3 – 3 (D) 25 – 25
13. The evaluation strategy which delays the evaluation of an expression until its value is needed and which avoids repeated evaluations is:
 (A) Early evaluation (B) Late evaluation
 (C) Lazy evaluation (D) Critical evaluation
14. If two or more expressions denote same memory address, then the expressions are:
 (A) Aliases (B) Definitions
 (C) Superiors (D) Inferiors
15. Operations that can be removed completely are called:
 (A) Strength reduction
 (B) Null sequences
 (C) Constant folding
 (D) None of these

PREVIOUS YEARS' QUESTIONS

1. In a compiler, keywords of a language are recognized during: **[2011]**
 (A) parsing of the program
 (B) the code generation
 (C) the lexical analysis of the program
 (D) dataflow analysis
2. Consider the program given below, in a block structured pseudo-language with lexical scoping and nesting of procedures permitted. **[2012]**
 Program main;
 Var ...
 Procedure A_1 ;
 Var ...
 Call A_2 ;

```
End  $A_1$ 
Procedure  $A_2$ ;
Var ...
Procedure  $A_{21}$ ;
Var ...
Call  $A_1$ ;
End  $A_{21}$ 
Call  $A_{21}$ ;
End  $A_2$ 
Call  $A_1$ ;
End main
```

Consider the calling chain: $\text{Main} \rightarrow A_1 \rightarrow A_2 \rightarrow A_{21} \rightarrow A_1$
 The correct set of activation records along with their access links is given by:



Common data for questions 3 and 4: The following code segment is executed on a processor which allows only register operands in its instructions. Each instruction can have at most two source operands and one destination operand. Assume that all variables are dead after this code segment.

```

c = a + b;
d = c * a;
e = c + a;
x = c * c;
If (x > a) {
    y = a * a;
}
Else {
    d = d * d;
    e = e * e;
}
    
```

3. What is the minimum number of registers needed in the instruction set architecture of the processor to

compile this code segment without any spill to memory? Do not apply any optimization other than optimizing register allocation. [2013]

- (A) 3
- (B) 4
- (C) 5
- (D) 6

4. Suppose the instruction set architecture of the processor has only two registers. The only allowed compiler optimization is code motion, which moves statements from one place to another while preserving correctness. What is the minimum number of spills to memory in the compiled code? [2013]

- (A) 0
- (B) 1
- (C) 2
- (D) 3

5. Which one of the following is NOT performed during compilation? [2014]

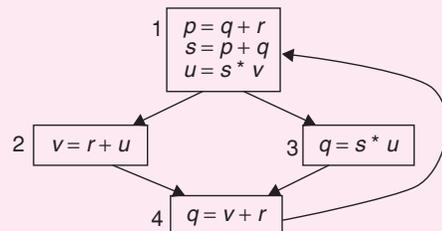
- (A) Dynamic memory allocation
- (B) Type checking
- (C) Symbol table management
- (D) Inline expansion

6. Which of the following statements are CORRECT? [2014]

- (i) Static allocation of all data areas by a compiler makes it impossible to implement recursion.
 - (ii) Automatic garbage collection is essential to implement recursion.
 - (iii) Dynamic allocation of activation records is essential to implement recursion.
 - (iv) Both heap and stack are essential to implement recursion.
- (A) (i) and (ii) only
 - (B) (ii) and (iii) only
 - (C) (iii) and (iv) only
 - (D) (i) and (iii) only

7. A variable x is said to be live at a statement S_i in a program if the following three conditions hold simultaneously: [2015]

1. There exists a statement S_j that uses x
2. There is a path from S_i to S_j in the flow graph corresponding to the program.
3. The path has no intervening assignment to x including at S_i and S_j .



The variables which are live both at the statement in basic block 2 and at the statement in basic block 3 of the above control flow graph are

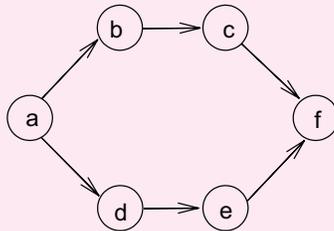
- (A) p, s, u
- (B) r, s, u
- (C) r, u
- (D) q, v

8. Match the following [2015]

P. Lexical analysis	1. Graph coloring
Q. Parsing	2. DFA minimization
R. Register allocation	3. Post-order traversal
S. Expression evaluation	4. Production tree

- (A) P-2, Q-3, R-1, S-4 (B) P-2, Q-1, R-4, S-3
 (C) P-2, Q-4, R-1, S-3 (D) P-2, Q-3, R-4, S-1

9. Consider the following directed graph:



The number of different topological orderings of the vertices of the graph is _____. [2016]

10. Consider the following grammar:

stmt → if expr then expr else expr; stmt | Ø
 expr → term relop term | term
 term → id | number
 id → a | b | c
 number → [0 - 9]

where **relop** is a relational operator (e.g., <, >, ...), Ø refers to the empty statement, and **if**, **then**, **else** are terminals.

Consider a program P following the above grammar containing ten **if** terminals. The number of control flow paths in P is _____. For example, the program

if e_1 **then** e_2 **else** e_3

has 2 control flow paths, $e_1 \rightarrow e_2$ and $e_1 \rightarrow e_3$. [2017]

11. Consider the expression $(a-1) * (((b+c)/3) + d)$. Let X be the minimum number of registers required by an *optimal* code generation (without any register spill) algorithm for a load/store architecture, in which (i) only load and store instruction can have memory operands and (ii) arithmetic instructions can have only register or immediate operands. The value of X is _____. [2017]

12. Match the following according to input (from the left column) to the compiler phase (in the right column) that processes it: [2017]

(P) Syntax tree	(i) Code generator
(Q) Character stream	(ii) Syntax analyzer
(R) Intermediate representation	(iii) Semantic analyzer
(S) Token stream	(iv) Lexical analyzer

- (A) $P \rightarrow$ (ii), $Q \rightarrow$ (iii), $R \rightarrow$ (iv), $S \rightarrow$ (i)
 (B) $P \rightarrow$ (ii), $Q \rightarrow$ (i), $R \rightarrow$ (iii), $S \rightarrow$ (iv)
 (C) $P \rightarrow$ (iii), $Q \rightarrow$ (iv), $R \rightarrow$ (i), $S \rightarrow$ (ii)
 (D) $P \rightarrow$ (i), $Q \rightarrow$ (iv), $R \rightarrow$ (ii), $S \rightarrow$ (iii)

ANSWER KEYS

EXERCISES

Practice Problems 1

1. A 2. C 3. B 4. A 5. B 6. B 7. C 8. B 9. B 10. B
 11. B 12. D 13. C 14. B 15. B

Practice Problems 2

1. D 2. B 3. A 4. B 5. A 6. D 7. D 8. D 9. D 10. D
 11. A 12. B 13. C 14. A 15. B

Previous Years' Questions

1. C 2. D 3. B 4. B 5. A 6. D 7. C 8. C 9. 6 10. 1024
 11. 2 12. C

COMPILER DESIGN

Time: 45 min.

Directions for questions 1 to 30: Select the correct alternative from the given choices.

- The most powerful parsing method is
(A) LALR (B) LR
(C) CLR (D) LL (1)
- In which phase 'type checking' is done?
(A) Lexical analysis
(B) Code optimization
(C) Syntax analysis
(D) Semantic analysis
- A shift reduces parser carries out the actions specified within braces immediately after reducing the corresponding rule of grammar, as below:
 $S \rightarrow aaD$ {Print "1"}.
 $S \rightarrow b$ {Print "2"}
 $D \rightarrow Sc$ {Print "3"}
 What is the translation of 'aaaabcc' using the syntax directed translation scheme described by the above rules?
 (A) 33211 (B) 11233
 (C) 11231 (D) 23131
- $E \rightarrow TE'$
 $E' \rightarrow + TE' / \in$
 $T \rightarrow FT'$
 $T' \rightarrow * FT' / \in$
 $F \rightarrow (E) / id$
 From above grammar, FOLLOW (E) is
 (A) { }, \$ (B) { \$, * }
 (C) { (, id } (D) { +, }, \$ }
- To eliminate backtracking, which one is used?
(A) Left Recursion
(B) Left Factoring
(C) Right Recursion
(D) Right Factoring
- Consider the grammar
 $T \rightarrow (T) / \in$
 Let the number of states in SLR (1), LR (1) and LALR (1) parsers for the grammar be n_1 , n_2 and n_3 respectively. Which relationship holds well?
 (A) $n_1 = n_2 = n_3$
 (B) $n_1 \geq n_3 \geq n_2$
 (C) $n_1 = n_3 < n_2$
 (D) $n_1 < n_2 < n_3$
- If w is a string of terminals and A, B are two non-terminals then which of the following are left-linear grammars?
 (A) $A \rightarrow wB/w$
 (B) $A \rightarrow Bw/w$
 (C) $A \rightarrow wB$
 (D) None of the above
- The grammar $E \rightarrow E * E / E + E / a$, is
 (A) Ambiguous
 (B) Unambiguous
 (C) Will not depend on the given sentence
 (D) None of these
- Shift-reduce parsers are
 (A) Bottom up parsers
 (B) Top down parsers
 (C) Both (A) and (B)
 (D) None of these
- Consider the following grammars:
 I. $E \rightarrow TE'$
 $E' \rightarrow + TE' / \in$
 $T \rightarrow FT'$
 $T' \rightarrow * FT' / \in$
 $F \rightarrow (E) / id$
 II. $S \rightarrow iCtSS' / a$
 $S' \rightarrow eS / \in$
 $C \rightarrow b$
 Which of the following is true?
 (A) II is LL (1) (B) I is LL (1)
 (C) Both (A) and (B) (D) None of these
- Consider the following grammar:
 $S \rightarrow iCtSS' / a$
 $S' \rightarrow eS / \in$
 $C \rightarrow b$
 First (S') is
 (A) { i, a } (B) { \$, e }
 (C) { e, \in } (D) { b }
- From the above grammar Follow(S) is.
 (A) { \$, e } (B) { \$ }
 (C) { e } (D) { \$, e }
- Find the LEADING (S) from the following grammar:
 $S \rightarrow a | ^ | (T)$
 $T \rightarrow T, S / S$
 (A) { $a, ^, ($ } (B) { $, a,$ }
 (C) { $, a, ($ } (D) { $, a, ^,$ }
- From above grammar find the TRAILING (T).
 (A) { $a,$ } (B) { $a, ^,$ }
 (C) { $,,$ } (D) { $, a,$ }
- Which of the following remarks logically follows?
 (A) $FIRST(\in) = \{\in\}$.
 (B) If FOLLOW (A) contains \$, then A may or may not be the start symbol.
 (C) If $A \rightarrow w$, is a production in the given grammar G, then $FIRST_k(A)$ contains $FIRST_k(w)$.
 (D) All of the above

16. Consider the following grammar:

$$S \rightarrow AB$$

$$B \rightarrow ab$$

$$A \rightarrow aa$$

$$A \rightarrow a$$

$$B \rightarrow b.$$

The grammar is

- (A) Ambiguous
- (B) Unambiguous
- (C) Can't predictable
- (D) None of these

17. If a handle has been found but there is no production with this handle as a right side, then we discover

- (A) Logical error
- (B) Runtime error
- (C) Syntactic error
- (D) All of the above

18. The function of syntax phase is

- (A) To build a literal table
- (B) To build an uniform symbol table
- (C) To parse the tokens produced by lexical analyzer
- (D) None of these

19. Which of the following are cousins of compilers?

- (A) Pre-processor and Assembler
- (B) Assembler and LEX
- (C) Pre-processor and YACC
- (D) LEX and YACC.

20. Error is detected in predictive parsing when ____ hold(s).

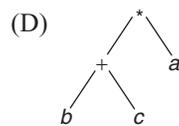
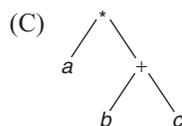
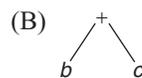
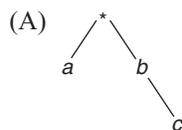
- (i) 'a' on top of stack and next input symbol is 'b'.
- (ii) When 'a' is on top of stack, 'a' is next input symbol and parsing table entry $M[A, a]$ is empty.
- (A) Neither (i) nor (ii)
- (B) Both (i) and (ii)
- (C) only (i)
- (D) only (ii)

21. Which one indicates abstract syntax tree (AST) of "a * b + c" with following grammar:

$$E \rightarrow E * T/T$$

$$T \rightarrow T + F/F$$

$$F \rightarrow id$$



22. The parse tree is constructed and then it is traversed and the semantic rules are evaluated in a particular order by a

- (A) Recursive evaluator
- (B) Bottom up translation
- (C) Top down translation
- (D) Phase tree method

23. The following grammar indicates

$$S \rightarrow a \alpha b | b \alpha c | a b$$

$$S \rightarrow \alpha S | b$$

$$S \rightarrow \alpha b b/a b$$

$$S \rightarrow \alpha b d b/b$$

- (A) LR (0) grammar
- (B) SLR grammar
- (C) Regular grammar
- (D) None of these

24. If the attributes of the child depends on the attributes of the parent node then it is ____ attribute.

- (A) Inherited
- (B) Directed
- (C) Synthesised
- (D) TAC

25. The semantic rule is evaluated and the intermediate code is generated when the production is expanded in ____

- (A) Parse tree method
- (B) Bottom up translation
- (C) Top down translation
- (D) Recursive evaluator model

26. Consider the grammar shown below:

$$S \rightarrow CC$$

$$C \rightarrow cC/a$$

The grammar is

- (A) LL (1)
- (B) SLR (1) But not LL (1)
- (C) LALR (1) but not SLR (1)
- (D) LR (1) but not LALR

27. The class of grammars for which we can construct predictive parsers looking k-symbols ahead in the input is called

- (A) LR (k)
- (B) CLR (k)
- (C) LALR (k)
- (D) LL (k)

28. A compiler is a program that

- (A) Places programs into memory and prepares them for execution.
- (B) Automates the translation of assembly language into machine language.
- (C) Accepts a program written in a high level language and produces an object program.
- (D) Appears to execute a source program as if it were machine language.

Common data for questions 29 and 30:

Consider the grammar

$$E \rightarrow TE'$$

$$E' \rightarrow + TE' \mid \epsilon$$

$$T \rightarrow FT'$$

$$T' \rightarrow * FT' \mid \epsilon$$

$$F \rightarrow (E) \mid id.$$

29. Which one is FOLLOW (F)?

(A) {+, }, \$}

(B) {+, (, *, }

(C) {*, }, \$}

(D) {+, *, }, \$}

30. FIRST (E) will be as same as

(A) FIRST (T)

(B) FIRST (F)

(C) Both (A) and (B)

(D) None of these

ANSWERS KEYS

- | | | | | | | | | | |
|-------|-------|-------|-------|-------|-------|-------|-------|-------|-------|
| 1. A | 2. D | 3. D | 4. A | 5. B | 6. C | 7. B | 8. A | 9. A | 10. B |
| 11. C | 12. A | 13. A | 14. C | 15. D | 16. A | 17. C | 18. C | 19. A | 20. B |
| 21. C | 22. A | 23. D | 24. A | 25. C | 26. A | 27. D | 28. C | 29. D | 30. C |

Operating System

Chapter 1: Processes and Threads	7.3
Chapter 2: Interprocess Communication, Concurrency and Synchronization	7.17
Chapter 3: Deadlock and CPU Scheduling	7.35
Chapter 4: Memory Management and Virtual Memory	7.54
Chapter 5: File Systems, I/O Systems, Protection and Security	7.74

U

N

I

T

7

This page is intentionally left blank

Chapter 1

Processes and Threads

LEARNING OBJECTIVES

- Basics of operating systems
- Services of OS
- Evolution of operating systems
- Processes
- Processes and process control blocks
- Process states
- Process creation
- Suspended processes
- OS control structures
- Process attributes
- Modes of execution
- Threads
- Thread functionality
- Thread synchronization

BASICS OF OPERATING SYSTEM

An operating system (OS) is a program that controls the execution of application programs and acts as an interface between applications and the computer hardware. The three objectives of an OS are as follows:

- Convenience:** An OS makes a computer more convenient to use.
- Efficiency:** An OS allows the computer system resources to be used in an efficient manner.
- Ability to evolve:** An OS should be constructed in such a way as to permit the effective development, testing and introduction of new system functions without interfering with service.

OS as a User-Computer Interface

Consider the below figure (Figure 1) which shows the hardware and software used in providing applications to a user in a layered fashion.

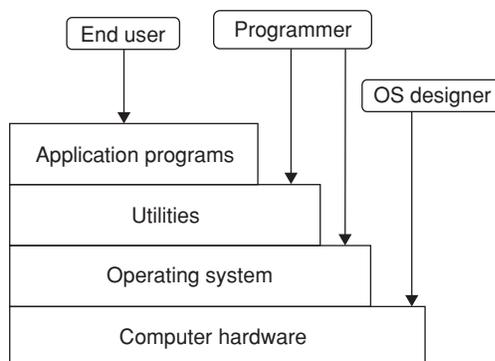


Figure 1 Layers and views of a computer system.

The end user of the application is not concerned with the details of computer hardware. Utilities implement the frequently used functions that assist in program creation, the management of files and control of input/output (I/O) devices. The most important collection of system programs comprises the OS. The OS masks the details of the hardware from the programmer and provides the programmer with a convenient interface for using the system.

Services of an OS

- Program development:** An OS provides a variety of facilities and services as editors and debuggers to assist programmers in creating programs.
- Program execution:** An OS handles the scheduling duties of program execution for the user.
- Access to I/O devices:** An OS provides a uniform interface that hides the details of I/O devices so that the programmers can access the I/O devices using simple reads and writes.
- Controlled access to files:** In a system with multiple users, an OS provides protection mechanism to control access to the files.
- System access:** For shared or public systems, an OS controls access to the system as a whole and to specific system resources.
- Error detection and response:** An OS must provide a response that clears the error condition with the least impact on running applications.
- Accounting:** A good OS will collect usage statistics for various resources and monitor performance parameters (viz., response time).

OS as Resource Manager

A computer is a set of resources for the movement, storage and processing of data and for the control of these functions. An OS is responsible for managing these resources.

1. An OS functions in the same way as ordinary computer software, that is, it is a program or suite of programs executed by the processor.
2. An OS frequently relinquishes control and must depend on the processor to allow it to regain control.

The OS directs the processor in the use of the other system resources and in the timing of its execution of other programs.

Figure 2 shows the resources that are managed by an OS.

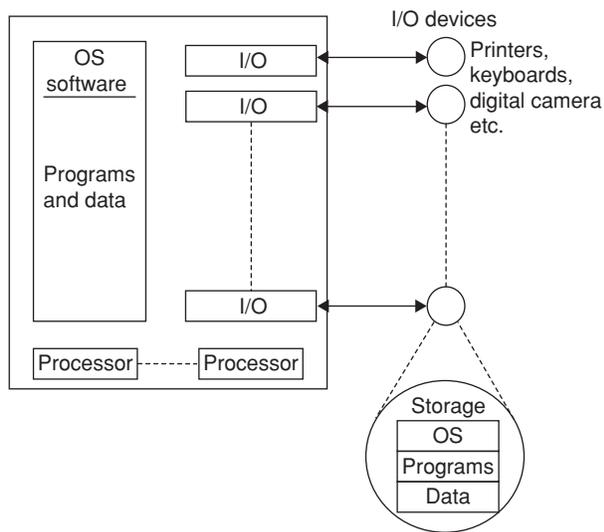


Figure 2 OS as a resource manager.

A portion of the OS lies in the main memory. This includes the kernel or nucleus, which contains the most frequently used functions in the OS and, at a given time, other portions of the OS currently in use. The remainder of main memory contains user programs and data. The allocation of this resource is controlled jointly by the OS and memory management hardware in the processor.

Evolution of OSs

Serial processing

With the earliest computers, the programmer interacted directly with the computer hardware. There was no OS. Programs in machine code were loaded via the input devices. This mode of operation is termed as *serial processing*. Problems with this system are scheduling and setup time.

Simple batch systems

- It requires the grouping up of similar jobs, which consist of programs, data and system commands.
- Users have no control over results of a program.
- Off-line debugging.

Multiprogrammed batch systems

In view of simultaneous execution of multiple programs, it improves system throughput and resource utilization.

Example: Windows XP, 98

- **Multitasking OS:** A running state of a program is called a *process* or a *task*. The concept of managing a multitude of simultaneously active programs, competing with each other for accessing the system resources is called *multitasking*.
- Serial multitasking or context switching is the simplest form of multitasking.

Example: Windows NT, Linux

- **Multuser OS:** It is defined as multiprogramming OS that supports simultaneous interaction with multiple users.

Example: Linux, Unix, a dedicated transaction processing system (viz., railway reservation system).

- **Multiprocessing OS:** The term *multiprocessing* means multiple CPUs performing more than one job at one time.

The term 'multiprogramming' means situation in which a single CPU divides its time between more than one job.

Time sharing systems

In this kind of OS, the processor time is shared among multiple users. The CPU switches rapidly from one user to another user; each user is given an impression that he/she has his own computer while it is actually one computer shared among many users.

If there are n users actively requesting service at one time, each user will only see on the average $1/n$ of the effective computer capacity, not counting OS overhead.

Bootstrap Bootstrap is an initial program which runs, when a computer is powered up (or) restarted. The task is to initialize system aspects (CPU registers to device controllers to memory contents). It is stored within the computer hardware known as *firm ware* (EEPROM).

PROCESSES

Processes and Process Control Blocks

Process A process is an instance of a program in execution. Two essential elements of a process are as follows:

1. Program code
2. Set of data

At any given point in time, while the program is executing, the process can be uniquely characterized by a number of elements, including the following:

1. Identifier
2. State
3. Priority

4. Program counter
5. Memory pointers
6. Context data
7. I/O status information
8. Accounting information

This information is stored in a data structure, typically called a *process control block*, that is created and managed by the OS.

Process control block (PCB) It contains sufficient information so that it is possible to interrupt a running process and later resume execution as if the interruption has not occurred.

Process States

The behaviour of an individual process can be characterized by listing the sequence of instructions that executed for that process. This listing is referred to as a *trace* of the process. Also the behaviour of a processor is shown by listing the traces of the various processes that are interleaved.

Dispatcher A dispatcher is a small program that switches the processor from one process to another.

Two-state process model In the simplest possible process model (Figure 3), at any time, a process is either being executed by a processor or not, that is, a processor may be in one of two states: *running* or *not running*.

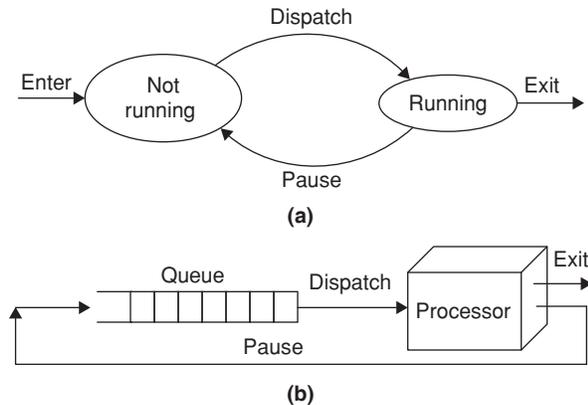


Figure 3 Two-state process model. (a) State transition diagram, (b) Queuing diagram

When the OS creates a new process, it creates the PCB for the process and enters that process into the system in the *not running* state. The process exists, is known to the OS, and is waiting for an opportunity to execute. From time to time, the currently running process will be interrupted and the dispatcher portion of the OS will select some other process to run. The former process moves from the *running* state to the *not running* state and one of the other processes moves to the *running* state.

Processes that are not running must be kept in some sort of queue, waiting their turn to execute. Figure 3(b) shows the structure. There is a single *queue* in which each entry is a pointer to the PCB of a particular process.

Creation and Termination of Processes

Process creation

When a new process is to be added to those currently being managed, the OS builds the data structures that are used to manage the process and allocates address space in main memory to the process. The common events which lead to process creation are as follows:

1. New batch job
2. Interactive logon
3. Created by OS to provide a service
4. Spawned by existing process

When the OS creates a process at the explicit request of another process, the action is referred to as *process spawning*. When one process spawns another, the former is referred to as the *parent process* and the spawned process is referred to as the *child process*.

Process termination

The following are the reasons for process termination:

1. *Normal completion*: Process executes an OS service call to indicate its completion.
2. *Time limit exceeded*: The process has run longer than the specified total time limit.
3. *Memory unavailable*: The process requires more memory than the system can provide.
4. *Bounds violation*: The process tries to access a memory location that it is not allowed to access.
5. *Protection error*: The process attempts to use a resource that is not allowed to access.
6. *Arithmetic error*: The process tries a prohibited computation.
7. *Time overrun*: The process has waited longer than a specified maximum for a certain event to occur.
8. *I/O failure*: Error occurs during input or output.
9. *Invalid instruction*: The process attempts to execute a non-existent instruction.
10. *Privileged instruction*: The process attempts to use an instruction reserved for OS.
11. *Data misuse*: A piece of data is of the wrong type or is not initialized.
12. *Operator or OS intervention*.
13. *Parent termination*.
14. *Parent request*.

Five-state model

The five states in Figure 4 are as follows:

1. *New*: The process is created but not admitted to the pool of executable processes.
2. *Running*: Process in execution, that is, it is using CPU.
3. *Blocked*: Waiting for some event to occur (i.e., I/O) before it can continue execution.

4. *Ready*: Process is ready for execution. Just it is waiting.
5. *Exit*: The process has been aborted by parent process or has finished its execution.

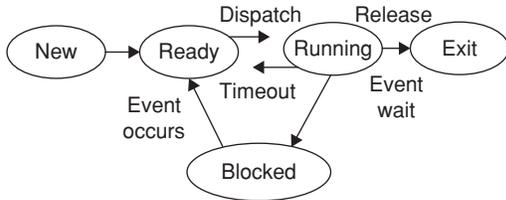


Figure 4 Process states.

Figure 4 indicates the types of events that lead to each state transition for a process. The possible transitions are as follows:

1. NULL → New: A new process is created to execute a program.
2. New → Ready: The OS will move a process from the New state to the Ready state when it is prepared to take on an additional process.
3. Ready → Running: When it is time to select a process to run, the OS chooses one of the processes in the Ready state.
4. Running → Exit: The currently running process is terminated by the OS if the process indicates that it has completed or if it aborts.
5. Running → Ready: The reasons for this transition are
 - Running process has reached the maximum allowable time for uninterrupted execution.
 - As the OS assigns different levels of priority to different processes, there will be pre-emption.
 - A process may voluntarily release control of the processor.
6. Running → Blocked: A process is put in the blocked state if it requests something for which it must wait.
7. Blocked → Ready: This transition occurs when the event for which the process has been waiting occurs.
8. Ready → Exit: A parent may terminate a child process at any time.
9. Blocked → Exit: Parent may terminate any blocked process.

Queuing model for five-state model

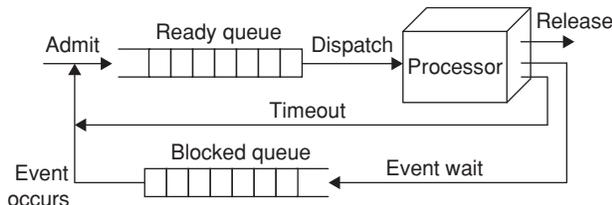


Figure 5 Single blocked queue.

Suspended Processes

Need for swapping In five-state process model using multiple blocked queues, the memory holds multiple processes. Moreover, the processor can move to another process when one process is blocked. But the processor is so much faster than I/O that it will be common for all of the processes in memory to be waiting for I/O. Thus, even with multiprogramming, a processor could be idle most of the time.

Then we can extend the main memory to accommodate more processes, but it is not an efficient solution. Another solution to this problem is swapping. Swapping involves moving part or all of a process from main memory to disk. When none of the processes in main memory is in the ready state, the OS swaps one of the blocked processes out onto disk into a suspend queue. The OS then brings in another process from the suspend queue or it honours a new process request. Then execution continues with the newly arrived process. With the use of swapping, another state is added to the process in the behaviour model.

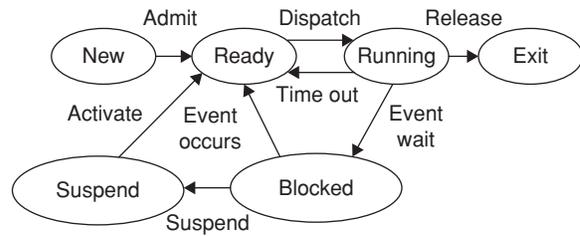


Figure 6 Process state-transition diagram with suspend state.

The four distinguishable states in this process model are as follows:

1. *Ready*: The process is in main memory and is available for execution.
2. *Blocked*: The process is in main memory and awaiting an event.
3. *Blocked/suspend*: The process is in secondary memory and awaiting an event.
4. *Ready/suspend*: The process is in secondary memory but is available for execution as soon as it is loaded into main memory.

Figure 7 shows the process state model with two suspend states:

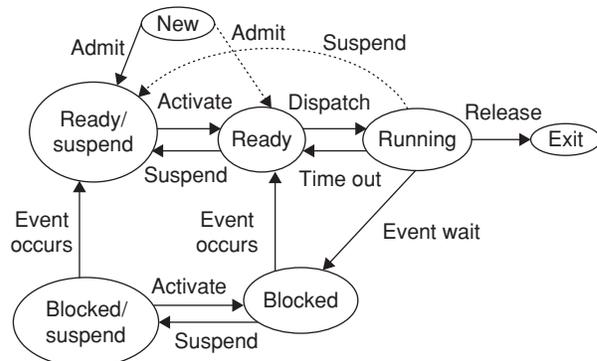


Figure 7 Process state transition diagram with suspend state.

Uses of Suspension

Characteristics of suspended process are as follows:

1. The process is not immediately available for execution.
2. The process may or may not wait on an event.
3. The process was placed in a suspend state by either itself, a parent or the OS.

Reasons for process suspension

- *Swapping*: To Release sufficient main memory.
- *Other OS reason*: OS may suspend a background process.
- *Interactive User Request*: A user may wish to suspend execution of a program.
- *Timing*: A process may be executed periodically and may be suspended.
- *Parent Process request*: A parent process may wish to suspend execution of a descendent.

OS CONTROL STRUCTURES

If the OS is to manage processes and resources, it must have information about the current status of each process and resources. The OS constructs and maintains tables of information about each entity that it is managing. Figure 8 shows the general structure of OS control tables:

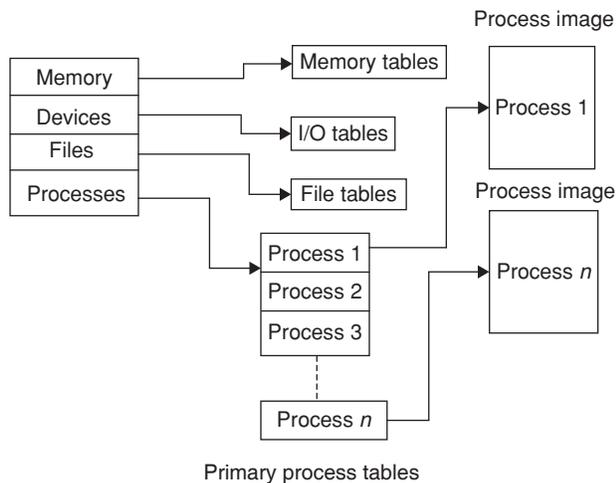


Figure 8 OS control tables.

Different tables maintained by the OS are

1. Memory
2. I/O
3. File
4. Process

Memory tables: These tables are used to keep track of both main and secondary memory.

I/O tables: These are used by the OS to manage the I/O devices and channels of the computer system.

File tables: These tables provide information about the existence of files, their location on secondary memory, their current status and other attributes.

Process tables: An OS must maintain process tables to manage processes.

PROCESS CONTROL STRUCTURES

The OS must know about

1. Process location
2. Process attributes

Process Location

The collection of program, data, stack and attributes is referred as process image.

The location of a process image will depend on the memory management scheme being used. The process image is maintained as a contiguous or continuous block of memory. This block is maintained in secondary memory, usually disk, so that the OS can manage the process, at least a small portion of its image must be maintained in main memory. To execute the process, the entire process image must be loaded into main memory or at least virtual memory. Thus the OS needs to know the location of each process on disk and for each such process that is in the main memory, the location of that process is in main memory.

For this, the OS maintains process tables. There is a primary process table with one entry for each process. Each entry contains, at least, a pointer to a process image.

Process Attributes

The typical information required by the OS for each process is as follows:

1. Process identification
2. Process state information
3. Process control information

Process identification Each process is assigned a unique numeric identifier, which may simply be an index into the primary process table. The identifier for a PCB includes the following:

1. Identifier of the process
2. Identifier of the process that created current process
3. User identifier

Process state information It consists of the contents of processor registers. It includes details of

1. User-visible register
2. Control and status registers
3. Stack pointers

Process control information It consists of the additional information needed by the OS to control and coordinate the various active processes. It includes the following:

1. Scheduling and state information
2. Data structuring
3. Interprocess communication

4. Process privileges
5. Memory management
6. Resource ownership and utilization

PROCESS CONTROL

Modes of Execution

Most processors support at least two modes of execution as follows:

1. More-privileged mode
2. Less-privileged mode

Two modes are required to protect the OS and key OS tables from interference by user programs.

1. **More-privileged mode:** This is also referred as *system mode*, *control mode* or *kernel mode*. Certain instructions can only be executed in Kernel mode (e.g., reading or altering a control register, viz., PSW, primitive I/O instructions, etc.).

The Kernel of the OS is a portion of the OS and encompasses the important system functions.

The functions of an OS kernel are as follows:

- Process management
- Memory management
- I/O Management
- Support functions

2. **Less-privileged mode:** This is also referred as *user mode*, because user programs typically would execute in this mode.

In this mode, the software has complete control of the processor and all its instructions, registers and memory.

Process creation If the OS decides to create a table, it has to proceed as follows:

1. Assign a unique process identifier to the new process.
2. Allocate space for the process.
3. Initialize the PCB
4. Set the appropriate linkages.
5. Create or expand other data structures.

Process switching In process switching, a running process is interrupted and the OS assigns another process to the running state and turns control over to that process. The design issues are as follows:

1. When to switch processes
2. Mode switching
3. Change of process state

When to switch processes A process switch may occur anytime that the OS has gained control from the currently running process. The mechanisms for interrupting the execution of a process are as follows:

1. Interrupt
2. Trap
3. Supervisor call

Interrupt When an interrupt occurs, the control is first transferred to an interrupt handler, which does some basic housekeeping and then branches to an OS routine that is concerned with the particular type of interrupt that has occurred (e.g., clock interrupt, I/O interrupt, memory fault, etc.).

Trap Trap related to an error or exception condition gets generated within the currently running process. If the error is fatal, the currently running process is moved to exit state and a process switch occurs, otherwise the action of the OS will depend on the nature of the error and design of the OS.

Supervisor call The OS may be activated by a supervisor call from the program being executed. The case of system call may place the user process in blocked state.

Mode switching If the processor identifies that any interrupt is pending, then

1. it sets the PC to the starting address of an interrupt handler program.
2. it switches from user mode to Kernel mode so that the interrupt processing code may include privileged instructions.

During this process, the context of the process, that has been interrupted, is saved into that PCB of the interrupted program. The context of a program includes PC, other processor registers and stack information.

The occurrence of an interrupt does not necessarily mean a process switch.

Change of process state The mode switch is a concept distinct from that of the process switch.

A mode switch may occur without changing the state of the process that is currently in Running state. In that case, the context saving and subsequent restoral involve little overhead. However, if the currently running process is to be moved to another state then the OS must make substantial changes in its environment. Thus, the process switch, which involves a state change, requires more effort than a mode switch.

System call In order to access the OS services, an interface is required which is provided by the system call.

All the system call routines are executed in Kernel mode. Whenever the system call is invoked, the process status word is changed from user mode to Kernel mode (0 → 1).

System calls are of six types as follows:

1. File system
2. Process
3. Scheduling
4. Interprocess communications
5. Socket
6. Miscellaneous

Execution of the OS

There are three possibilities to consider about OS execution:

1. Separate Kernel (Figure 9)
2. OS functions execute within user processes (Figure 10)
3. OS functions execute as separate processes (Figure 11)

Separate Kernel

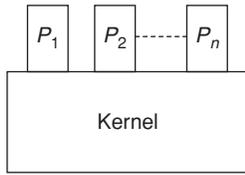


Figure 9 OS as separate Kernel.

Here the Kernel of the OS is executed outside of any process. When currently running process is interrupted or issues a supervisor call, the mode context of this process is saved and control is passed to the Kernel.

Execution within user process

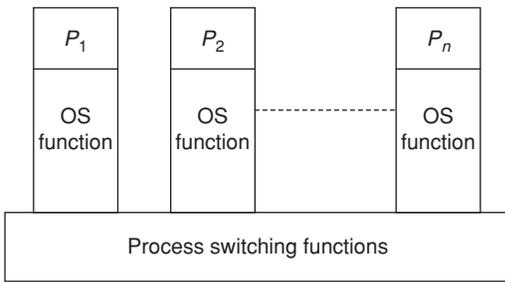


Figure 10 OS functions execute within user processes.

We execute virtually all OS software in the context of user process. To pass control from a user program to the OS, the mode context is saved and a mode switch takes place to an OS routine, that is, a process switch is not performed, just a mode switch within the same process.

Process-based OS

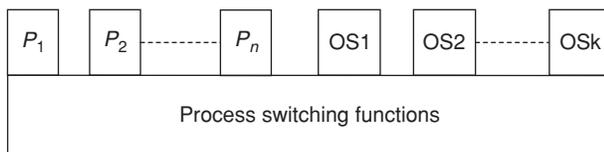


Figure 11 OS functions execute as separate process.

Here the OS is a collection of system processes. This approach encourages the use of modular OS with minimal, clean interface between the modules.

THREADS

A *thread* is a basic unit of CPU utilization. It comprises a thread ID, a program counter, a register set and a stack.

Multithreading

It refers to the ability of an OS to support multiple, concurrent paths of execution within a single process.

The threads which belong to same process can share their

1. Code section
2. Data section
3. Other OS resources

If a process has multiple threads of control, it can perform more than one task at a time. Figure 12 shows single threaded and multithreaded process models:

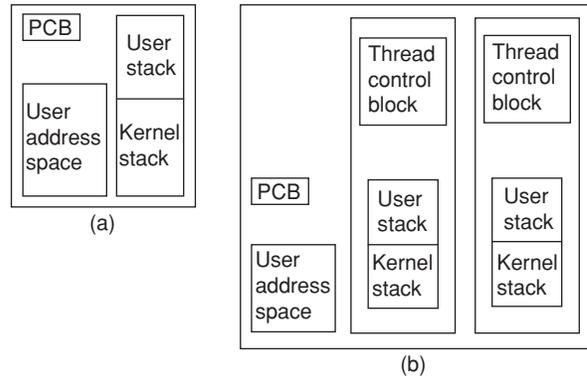


Figure 12 Process models. (a) Single-thread process model (b) Multithreaded process model.

As OS based on its design will be in one of the following manners (Figure 13):

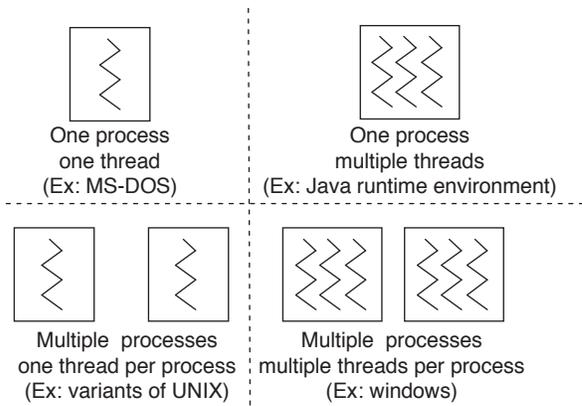


Figure 13 Threads and processes.

The threads of a process consist of the following:

1. Thread execution state.
2. Saved thread context when not running
3. An execution stack
4. Some pre-thread static storage for local variables.
5. Access to the memory and resources of its process, shared with all other threads in that process.

All the threads of a process share the state and resources of that process.

Benefits of multithreaded programming

1. *Responsiveness*: Multithreading an interactive application may allow program to continue running even if

part of it is blocked or is performing a lengthy operation, thereby increasing responsiveness to the user.

2. *Resource sharing*: Threads share the memory and the resources of the process to which they belong by default.
3. *Economy*: It is economical to create a new thread in an existing process than to create a brand-new process. It takes less time to context switch between two threads of same process than to switch between processes. Also the time to terminate a thread is less than process termination.
4. *Scalability*: Multithreading on a multi-CPU machine increases parallelism.

Applications that benefit from thread As the threads take advantages of multiple processors, image processing which can be done in parallel, will execute in threads.

Animation rendering is another thread application, where each frame can be rendered in parallel, as each one is independent of other GUI programming will execute at least two threads when it is processing large number of files.

Applications that cannot benefit from thread The main drawbacks of threads is if kernel is single threaded, system call of one thread will block the whole process, in which CPU will be idle during the blocking period.

The other major drawback is security as it is possible that a thread can overwrite the stack as the other thread, as they were meant to cooperate on a single task. Applications that are developed using PHP does not support multithreading at the server side.

THREAD FUNCTIONALITY

The key states for threads are as follows:

1. Running
2. Ready
3. Blocked

There are four basic thread operations associated with a change in thread state:

1. *Spawn*: When a new process is spawned, a thread for that process is also spawned.
Also, a thread within a process may spawn another thread within the same process.
2. *Block*: When a thread needs to wait for an event, it will block. Then the processor may turn to the execution of another ready thread in the same or different process.
3. *Unblock*: When an event for which a thread is blocked occurs, the thread is moved to Ready queue.
4. *Finish*: When a thread completes, its register context and stacks are deallocated.

Multithreading on a Uni-processor

On a uni-processor, multiprogramming enables the interleaving of multiple threads within multiple processes. For

example, consider the execution of three threads *A*, *B*, *C* in two processes on a single processor which are interleaved (Figure 14).

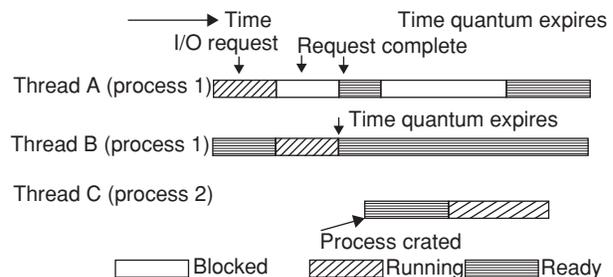


Figure 14 Multithreading on a uni-processor.

Execution passes from one thread to another, either when the currently running thread is blocked or its time slice is exhausted.

Resources used in thread creation and process creations

As process has heavy weight, when it is created, new address space is required, which includes stack, heap and data section, etc. If a process shares the memory, then the IPC is expensive.

The thread is a light-weight process, if it doesn't require any new resources, as it will share the process resources to which it belongs. The major benefit of this is, several threads belong to same activity and can run under same address space.

Thread Synchronization

All of the threads of a process share the same address space and other resources, such as open files. Any alteration of a resource by one thread affects the environment of the other threads in the same process. So, synchronization mechanism is required to coordinate the activities of all the threads within a process.

The techniques used for thread synchronization is the same as process synchronization techniques, which has been discussed later in this book.

TYPES OF THREADS

User-Level Threads

All thread management is done by the application. The Kernel is not aware of the existence of threads. Thread creation and scheduling are done in user space. User-level threads (Figure 15) are fast to create and manage. User-level library provides support for creating, managing and scheduling threads. In single-threaded Kernel, blocking system call from user level thread will block the entire process, even if other threads are ready to run.

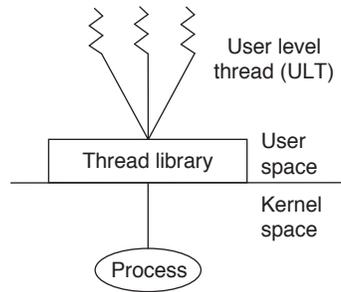


Figure 15 Pure user-level threads.

Advantages of ULTs

1. ULT creation does not require Kernel mode privileges.
2. ULT scheduling can be application specific.
3. ULTs can run on any OS.

Disadvantages of ULTs

1. When a ULT executes a system call, not only is that thread blocked, but also all of the threads within the process are blocked.
2. In a pure ULT strategy, a multithreaded application cannot take advantage of multiprocessing.

Kernel-level Threads

Kernel-level threads (KLTs) are supported directly by the OS. The creation, scheduling, management are done by kernel in kernel space. They are slower to create and manage. In a multiprocessor, Kernel can schedule threads on different processors.

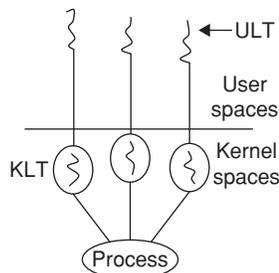


Figure 16 Pure kernel-level threads.

Advantages of KLTs

1. Kernel can simultaneously schedule multiple threads from the same process on multiple processors.
2. If one thread in a process is blocked, the kernel can schedule another thread of the same process.
3. Kernel routines themselves multithreaded.

Disadvantages of KLTs

The transfer of control from one thread to another within the same process requires a mode switch to the Kernel.

Combined Approach

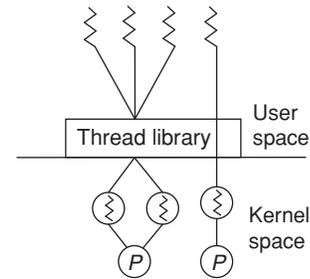


Figure 17 Combined approach.

In combined approach (Figure 17), multiple threads within the same application can run in parallel on multiple processors and a blocking system call need not block the entire process.

Relationship between the threads and processes

1. *One-to-one relationship*: Each thread of execution is a unique process with its own address and resources.
Example: Traditional UNIX.
2. *Many-to-one relationship*: A process defines an address space and dynamic resources ownership. Multiple threads may be created and executed within that process.
Example: Windows NT, Solaris, Linux.
3. *One-to-many relationship*: A thread may migrate from one process environment to another. This allows a thread to be easily moved among distinct systems.
Example: Emerald
4. *Many-to-many relationship*: It combines the attributes of M:1 and 1:M cases.
Example: TRIX

THREADING ISSUES

fork() and exec() System Calls

A `fork()` system call is used to create a separate, duplicate process. In UNIX, each process is identified by its process identifier, which is a unique integer. A new process is created by `fork()` system call.

The new process consists of a copy of the address space of the original process. This mechanism allows the parent process to communicate easily with its child process. Both processes continue execution at the instruction after the `fork()`, with one difference: the return code for the `fork()` is zero for the new process, whereas the process identifier of the child is returned to the parent. The `exec()` system call is used after a `fork()` system call by one of the two processes to replace the processes memory space with a new program.

7.12 | Unit 7 • Operating System

If one thread in a program calls `fork()`, then UNIX chooses two alternatives as follows:

1. Duplicates all the threads
2. Duplicates only the thread that invoked the `fork()` system call.

If a thread invokes the `exec()` system call, the program specified in the parameter to `exec()` will replace the entire process including all threads.

Cancellation

Thread cancellation is the task of terminating a thread before it has completed. A thread that is to be cancelled is often referred to as the *target thread*. Cancellation of a thread may occur in two different scenarios as follows:

1. *Asynchronous cancellation*: One thread immediately terminates the target thread.
2. *Deferred cancellation*: The target thread periodically checks whether it should terminate, allowing it an opportunity to terminate itself in an orderly fashion.

Microkernels

Microkernel (Figure 18) is a small OS core that provides the foundation for modular extensions.

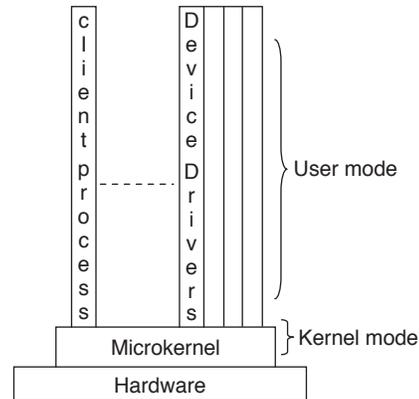


Figure 18 Microkernel architecture.

Advantages of Microkernel Organization

- Uniform interfaces
- Extensibility
- Flexibility
- Portability
- Reliability
- Distributed system support
- Support for object oriented OS

EXERCISES

Practice Problems I

Directions for questions 1 to 15: Select the correct alternative from the given choices.

1. In `fork()` system call, the return value to the parent process and the child process are respectively
(A) PID of child process, 1
(B) PID of child process, 0
(C) PID of child, PID of parent process
(D) PID of parent process, PID of child
2. Which of the following is not an advantage of thread?
(A) Inter process communication
(B) Less memory space occupied by thread
(C) Less time to create and terminate than a process
(D) Context switching is faster
3. A process executes the following segment of code
`for (i = 1; i < 10; i++) fork();`
The number of new processes created is
(A) 1024 (B) 1023
(C) 1025 (D) 1028
4. For each of the following transitions, between process states, which transition is not possible?
(A) Running → Ready
(B) Blocked → Suspend
(C) Ready → Ready/Suspend
(D) Blocked → Running
5. An operating system can be mapped to a five-state process model. A new event has been designated as capable to pre-empt the existing processes in order to trigger a new process to complete. Select the correct statement from below:
(A) A new state need to be added to the existing transition model to accommodate the changes.
(B) The existing model still holds good.
(C) Both the states and transitions of the existing model have to be changed.
(D) Only the transitions need to be modified.
6. The advantage of having multiple threads over multiple processes is
(i) Less time for creation
(ii) Less time for termination
(iii) Less time for switching
(iv) Kernel not involved in communication among threads
(A) (i), (ii), (iii) (B) (i), (ii), (iv)
(C) (ii), (iii), (iv) (D) (i), (ii), (iii), (iv)
7. Select the correct sequence of steps taken by the processor when an interrupt occurs
(i) Switch from user mode to kernel mode.
(ii) Set the program counter to the first instruction of the interrupt handling routine.
(iii) Save the current context.

Practice Problems 2

Directions for questions 1 to 14: Select the correct alternative from the given choices.

1. Many-to-many multithreading model is used in which of the following operating system?
 - (A) Windows NT/2000 with Thread Fibre
 - (B) Windows 95
 - (C) Windows 98
 - (D) Solaris Green Threads
2. Which of the following does not interrupt a running process?
 - (A) Device
 - (B) Timer
 - (C) Scheduler
 - (D) Power failure
3. Which of the following need not be saved on a context switch between processes?
 - (A) General purpose registers
 - (B) Translation look aside buffer
 - (C) Program counter
 - (D) All of the above
4. Which of the following actions is/are typically not performed by the OS when switching context from process A to process B?
 - (A) Saving current register values and restoring the register values for process B
 - (B) Changing address translation tables
 - (C) Swapping out the memory image of process A to the disk
 - (D) Both (B) and (C)
5. For each thread in a multithreaded process, there is a separate
 - (A) Process control block
 - (B) User address space
 - (C) User and kernel stack
 - (D) Kernel space only
6. When a supervisor call is received
 - (A) Mode switch happens
 - (B) Process switch happens
 - (C) Both (A) and (B)
 - (D) Neither (A) nor (B)
7. What is the purpose of jacketing?
 - (A) Convert non-blocking system call to blocking system call
 - (B) Convert blocking system call to non-blocking system call
 - (C) Convert blocking system call into a new thread
 - (D) Convert non-blocking system call into a new thread
8. Which of the following statements is/are always true?
 - (i) Time taken for mode switch is always greater than process switch.
 - (ii) Time taken for mode switch is always less than process switch.
 - (iii) Time taken for mode switch is always equal to process switch.
 - (A) (i) and (iii)
 - (B) (ii) and (iii)
 - (C) Only (i)
 - (D) Only (ii)
9. Which of the following is the property of time sharing systems?
 - (i) Multiple user access
 - (ii) Multiprogramming
 - (A) (i) only
 - (B) (ii) only
 - (C) Both (i) and (ii)
 - (D) Neither (i) nor (ii)
10. Which of the following is/are not a valid reason for process creation?
 - (i) Created by OS
 - (ii) Interactive logon
 - (iii) Privileged instruction
 - (A) (i), (ii)
 - (B) (ii), (iii)
 - (C) (i), (iii)
 - (D) (iii) only
11. Which of the following is/are reason(s) for blocking a running process?
 - (i) A call from the running program to a procedure that is a part of OS code.
 - (ii) A running process may initiate an I/O operation.
 - (iii) A user may block a running process.
 - (A) (i), (ii) only
 - (B) (ii), (iii) only
 - (C) (i), (iii) only
 - (D) (iii) only
12. If the OS is pre-empting a running process because a higher priority process on blocked/suspend queue has just become unblocked, then the running process moved to _____ queue.
 - (A) Suspend
 - (B) Ready/suspend
 - (C) Blocked
 - (D) Blocked/suspend
13. Which of the following is used to call an OS function?
 - (A) Interrupt
 - (B) Trap
 - (C) Supervisor call
 - (D) All of these
14. Which of the following is a general component of a thread?
 - (i) Thread ID
 - (ii) Register set
 - (iii) User stack
 - (iv) Kernel stack
 - (A) (i), (iii), (iv)
 - (B) (i), (ii), (iv)
 - (C) (i), (ii), (iii)
 - (D) (i), (ii), (iii), (iv)

PREVIOUS YEARS' QUESTIONS

1. Consider the following statements with respect to User-level threads and Kernel-supported threads. **[2004]**
 - (1) Context switch is faster with Kernel-supported threads
 - (2) For user-level threads, a system call can block the entire process
 - (3) Kernel supported threads can be scheduled independently
 - (4) User level threads are transparent to the Kernel

Which of the above statements are true?

(A) 2, 3 and 4 only (B) 2 and 3 only
(C) 1 and 3 only (D) 1 and 2 only
2. Which one of the following is true for a CPU having a single interrupt request line and a single interrupt grant line? **[2005]**
 - (A) Neither vectored interrupt nor multiple interrupting devices are possible
 - (B) Vectored interrupts are not possible but multiple interrupting devices are possible
 - (C) Vectored interrupts and multiple interrupting devices are both possible
 - (D) Vectored interrupt is possible but multiple interrupting devices are not possible
3. Normally user programs are prevented from handling I/O directly by I/O instructions in them. For CPUs having explicit I/O instructions, such I/O protection is ensured by having the I/O instructions privileged. In a CPU with memory mapped I/O, there is no explicit I/O instruction. Which one of the following is true for a CPU with memory mapped I/O? **[2005]**
 - (A) I/O protection is ensured by operating system routine(s)
 - (B) I/O protection is ensured by a hardware trap
 - (C) I/O protection is ensured during system configuration
 - (D) I/O protection is not possible
4. Consider the following code fragment: `if (fork ()==0)` **[2005]**

```
{ a = a + 5; printf ("%d, %d\n", a, &a) ; }
else { a=a-5; printf ("%d, %d\n", a, &a) ; }
```

Let u , v be the values printed by the parent process, and x , y be the values printed by the child process. Which one of the following is *true*?

 - (A) $u = x + 10$ and $v = y$
 - (B) $u = x + 10$ and $v \neq y$
 - (C) $u + 10 = x$ and $v = y$
 - (D) $u + 10 = x$ and $v \neq y$
5. Consider the following statements about user-level threads and Kernel-level threads. Which one of the following statements is *false*? **[2007]**
 - (A) Context switch time is longer for Kernel-level threads than for user level threads.
 - (B) User level threads do not need any hardware support.
 - (C) Related Kernel-level threads can be scheduled on different processors in a multi-processor system.
 - (D) Blocking one Kernel-level thread blocks all related threads.
6. Which of the following statements about synchronous and asynchronous I/O is NOT *true*? **[2008]**
 - (A) An ISR is invoked on completion of I/O in synchronous I/O but not in asynchronous I/O
 - (B) In both synchronous and asynchronous I/O, an ISR (Interrupt Service Routine) is invoked after completion of the I/O
 - (C) A process making a synchronous I/O call waits until I/O is complete, but a process making an asynchronous I/O call does not wait for completion of the I/O
 - (D) In the case of synchronous I/O, the process waiting for the completion of I/O is woken up by the ISR that is invoked after the completion of I/O
7. A process executes the following code


```
for (i=0; i<n; i++) fork( );
```

The total number of child processes created is **[2008]**

 - (A) n (B) $2^n - 1$
 - (C) 2^n (D) $2^{n+1} - 1$
8. A CPU generally handles an interrupt by executing an interrupt service routine **[2009]**
 - (A) As soon as an interrupt is raised.
 - (B) By checking the interrupt register at the end of fetch cycle.
 - (C) By checking the interrupt register after finishing the execution of the current instruction.
 - (D) By checking the interrupt register at fixed time intervals.
9. A thread is usually defined as 'light weight process' because an operating system (OS) maintains smaller data structures for a thread than for a process. In relation to this, which of the following is true? **[2011]**
 - (A) On per-thread basis, the OS maintains only CPU register state
 - (B) The OS does not maintain a separate stack for each thread
 - (C) On per thread basis, the OS does not maintain virtual memory state.
 - (D) On per thread basis, the OS maintains only scheduling and accounting information.
10. Let the time taken to switch between user and kernel modes of execution be t_1 while the time taken to switch between two processes be t_2 . Which of the following is *true*? **[2011]**

- (A) $t_1 > t_2$
 (B) $t_1 = t_2$
 (C) $t_1 < t_2$
 (D) Nothing can be said about the relation between t_1 and t_2

11. A process executes the code

```
fork();
fork();
fork();
```

The total number of **child** processes created is [2012]

- (A) 3 (B) 4
 (C) 7 (D) 8

12. Which one of the following is **FALSE**? [2014]

- (A) User level threads are not scheduled by the Kernel.
 (B) When a user level thread is blocked, all other threads of its process are blocked.

- (C) Context switching between user level threads is faster than context switching between Kernel-level threads.
 (D) Kernel-level threads cannot share the code segment.

13. Threads of a process share [2017]

- (A) global variables but not heap.
 (B) heap but not global variables.
 (C) neither global variables nor heap.
 (D) both heap and global variables.

14. Which of the following is/are shared by all the threads in a process? [2017]

- I. Program counter
 II. Stack
 III. Address space
 IV. Registers

- (A) I and II only (B) III only
 (C) IV only (D) III and IV only

ANSWER KEYS

EXERCISES

Practice Problems 1

1. B 2. A 3. B 4. D 5. D 6. D 7. C 8. B 9. C 10. D
 11. A 12. C 13. C 14. C 15. A

Practice Problems 2

1. A 2. C 3. B 4. C 5. C 6. A 7. B 8. D 9. C 10. D
 11. A 12. B 13. C 14. D

Previous Years' Questions

1. A 2. C 3. A 4. D 5. D 6. A 7. B 8. C 9. C 10. C
 11. C 12. D 13. D 14. B

Chapter 2

Interprocess Communication, Concurrency and Synchronization

LEARNING OBJECTIVES

- Principles of concurrency
- Process interaction
- Mutual exclusion
- Semaphores
- Binary semaphore
- Mutual exclusion using semaphores
- Progress using semaphores
- Classical problems of synchronization
- Dining philosophers problem
- Monitors
- Message passing
- Indirect addressing
- Mutual exclusion using message passing

BASIC CONCEPTS

Multiprogramming: It deals with the management of multiple processes within a uniprocessor system.

Multiprocessing: It deals with the management of multiple processes within a multiprocessor.

The fundamental operating system (OS) design is *concurrency*. Concurrency encompasses a host of design issues, including communication among processes, sharing of and competing for resources, synchronization of the activities of multiple processes and allocation of processor time to processes.

PRINCIPLES OF CONCURRENCY

There are two examples for concurrent processing as follows:

- In a single-processor multiprogramming system, processes are interleaved in time to yield the appearance of simultaneous execution.
- In a multiprocessor system, it is possible not only to interleave the execution of multiple processes but also to overlap them.

There are two problems with these techniques:

- Problem with sharing of global resources
- Problem with allocation of resources optimally.
- Problem with locating a programming error as results is not deterministic and reproducible.

Example:

```
void process( )
{
    in = getchar( );
    out = in;
    putchar(out);
}
```

The procedure 'process' reads a character and prints it. Let us suppose that we have a uniprocessor system, with single user. Let the user running multiple applications and all applications use the procedure for reading and printing, that is, all the applications share common procedure for efficient and close interaction among them. But this sharing leads to problems. For example,

- Let the process P_1 invokes 'process' and is interrupted immediately after 'getchar' returns its value and stores it in 'in'. Here the most recently entered character 'C' is stored in variable 'in'.
- Now, suppose the process P_2 is activated and it invokes 'process', which runs to conclusion, inputting and then displaying a single character, D , on the screen.
- The process P_1 is resumed. By this time, the value 'C' has been overwritten in 'in' and therefore lost. Instead 'in' contains 'D', which is transferred to 'out' and displayed.

Here the problem is with sharing a global variable. To avoid these types of problems, we impose some rules like, only one process

at a time may enter 'process' and that once in 'process' the procedure must run to completion before it is available for another process.

This problem is also applicable to multiprocessor systems.

Race condition: A race condition occurs when multiple processes or threads read and write data items so that the final result depends on the order of execution of instructions in the multiple processes.

Example: Let P_1 and P_2 be two processes that share global variables a and b , with initial values $a = 0$, $b = 1$. At some point in its execution, P_1 executes $a = a + b$ and P_2 executes $b = a + b$.

If P_1 executes before P_2 , then $a = 1$, $b = 2$.

If P_2 executes before P_1 , then $b = 1$, $a = 1$.

OS Concerns for Concurrency

1. The OS must be able to keep track of various processes using PCBs.
2. The OS must allocate and deallocate various resources for each active process.
3. The OS must protect the data and physical resources of each process.
4. The functioning of a process and the output it produces must be independent of the speed at which its execution is carried out relative to the speed of other concurrent processes.

PROCESS INTERACTION (IPC)

Process classification There are two types of processes as follows:

1. Independent/isolated
2. Cooperating

Independent Process: It cannot affect or be affected by the execution of another process.

Cooperating Process: It can affect or be affected by the execution of another process.

We can classify the ways in which processes interact on the basis of the degree to which they are aware of each other's existence. There are three types of process interaction as follows:

1. Process unaware of each other
2. Process indirectly aware of each other
3. Processes directly aware of each other

Competition Among Processes for Resources

1. This situation arises when processes unaware of each other.
2. There is no exchange of information between the competing processes.

3. But the execution of one process may affect the behaviour of competing processes.
4. With competing processes, there will be three control problems as follows:

Need for mutual exclusion

Example: Suppose two or more processes require access to a single non-sharable resource, such as a printer. Then that resource is referred as *critical resource* and the portion of the program that uses it is called *critical section* of the program. In the case of printer, only one process will have the control of printer while it prints an entire file.

Possibility of deadlock

Example: Two processes waiting for each other indefinitely for the release of resources.

Possibility of starvation

Example: One process is denied access to a particular resource which is required for the execution of that process.

Control of competition inevitably involves the OS, because it is the OS that allocates resources.

Cooperation Among Processes by Sharing

This situation arises when the processes are indirectly aware of each other. Processes may use and update the shared data without reference to other processes but know that other processes may have access to same data. So the control mechanism must ensure the integrity of the shared data. Problems with this type of sharing are

1. Mutual exclusion
2. Deadlock
3. Starvation
4. Data coherence

Data coherence

Suppose two items of data p and q are maintained in the relationship $p = q$, that is, any program that updates p and q values must maintain the relationship.

$$\begin{aligned} \text{Let } P_1 : p &= p + 1; \\ & q = q + 1; \\ P_2 : p &= p * 2; \\ & q = q * 2; \end{aligned}$$

Let initially the state is consistent, that is, $p = 2$, $q = 2$

Then the concurrent execution of P_1 and P_2 with mutual exclusion on p , q will be $p = p + 1$;

$$\begin{aligned} q &= q * 2; \\ q &= q + 1 \\ p &= p * 2; \end{aligned}$$

The final values of p and q will be $p = 6$, $q = 5$.

So the consistency is not maintained.

Cooperation Among Processes by Communication

This situation arises when processes are aware of each other directly. All the processes communicate with each other to synchronize or coordinate the various activities. Problems with this communication are as follows:

1. Deadlock
2. Starvation

BASIC DEFINITIONS

Atomic operations: A sequence of one or more statements that appears to be indivisible, that is, no process will interrupt the operation.

Critical section A section of code within a process that requires access to shared resources and that must not be executed while another process is in a corresponding section of code.

Deadlock A situation in which two or more processes are unable to proceed because each is waiting for one of the others to do something.

Mutual exclusion The requirement that when one process is in a critical section that accesses shared resources, no other process may be in a critical section that accesses any of those shared resources.

Race condition A situation in which multiple threads or processes read and write a shared data item and the final result depends on the relative timing of their execution.

Starvation A situation in which a runnable process is overlooked indefinitely by the scheduler; although it is able to proceed, it is never chosen.

CRITICAL SECTION

1. A section of code or set of operations, in which process may be changing shared variables, updating a common file or a table etc.
2. For the process that execute concurrently, it should ensure that execution of critical section should be made atomic. Atomic means that either an operation in the critical section should happen in its entirety or not at all.
3. Critical section of a process should not be executed concurrently with the critical section of another process.
4. To avoid Race Condition, we must have the following: 'if one process is in critical section, other competing process must be excluded to enter their critical sections, that is, a process must enter the critical section in a mutually exclusive way'.
This is called *problem of mutual exclusion*.

Region of code that updates or uses shared data to provide a consistent view of objects need to make sure an update is not in progress when reading the data.

5. Need to provide mutual exclusion for a critical section.

Requirements for Critical Section Problem

Mutual exclusion: No two contending processes should be simultaneously executing inside their critical section.

Bounded waiting: No process should have to wait forever to enter its critical section.

Progress: If no process is executing in its critical section and there exists some processes that wish to enter their critical sections, then only those processes that are not executing in the critical section can participate in the decision of which will enter its critical section next and this selection cannot be postponed indefinitely.

No Assumption: No assumption should be made about relative speed and properties of contenting processes.

MUTUAL EXCLUSION

1. Only one process at a time can be updating shared objects.
2. Successful use of concurrency among processes requires the ability to define critical sections and enforce mutual exclusion.
3. Mutual exclusion in the use of a shared resource is provided by making its access mutually exclusive among the processes that share the resources.

Any facility that provides mutual exclusion should meet the following requirements:

1. No assumption regarding the relative speed of the process.
2. A process is in its critical section for a finite time only.
3. Only one process allowed in the critical section.
4. Process requesting access to critical section should not wait indefinitely.
5. A process waiting to enter critical section cannot be blocking a process in critical section or any other process.

Mutual exclusion can be satisfied in one of the following ways:

1. **Software approach:** This approach leaves the mutual exclusion responsibility with the process that wish to execute concurrently. This approach is prone to high processing overhead and bugs.
2. **Hardware support:** Use special-purpose machine instructions. This approach involves less overhead.
3. **Provide some level of support within the OS or a programming language:** Such techniques are as follows:
 - Semaphores
 - Monitors
 - Message passing

Hardware support for mutual exclusion: Use one of the following techniques:

1. Interrupt disabling
2. Special machine instructions
 - Compare and swap instructions
 - Exchange instructions

Interrupt Disabling

To provide mutual exclusion, it is sufficient to prevent a process from being interrupted. A process can enforce mutual exclusion in the following way:

```
while(true)
{
/* disable interrupts */;
/* critical section */;
/* enable interrupts */;
/* remainder */;
}
```

Here, the critical section is not interrupted, so mutual exclusion is guaranteed.

Problems with Interrupt Disabling

The efficiency of execution could be noticeably degraded, because the processor is limited in its ability to interleave processes. Interrupt disabling does not work on in a multi-processor architecture.

Special Machine Instructions

Two of the most commonly used special machine instructions are as follows:

1. Compare and swap
2. Exchange instruction

Compare and swap instructions It is defined as below:

```
int compare_and_swap(int *word, int
testval, int newval)
{
int oldval;
oldval = *word;
if (oldval == testval)
*word = newval;
return oldval;
}
```

- ‘Access to a memory location excludes any other access to that same location.’ On the basis of this principle, special machine instructions provide mutual exclusion.
- The above compare and swap instructions will check a memory location (*word) against a test value. If current value is test value, it is replaced with new value. Always the old value is returned.

- The below code provides mutual exclusion using compare and swap instructions:

```
/* Mutual Exclusion */
const int n = /* number of processes */
int S;
void P(int i)
{
while (true)
{
while (compare_and_swap(S, 0, 1)==1)
/* do nothing */;
/* critical section */;
S = 0;
/* remainder */;
}
}
void main( )
{
S = 0;
begin (P(1), P(2) .... P(n));
}
```

Here, a shared variable ‘S’ is initialized to ‘0’. The only process that may enter its critical section is one that finds ‘S’ equal to 0.

All other processes that enter their critical sections go into a busy waiting mode.

Busy waiting or *spin waiting* is a technique in which a process can do nothing but continue to execute an instruction or set of instructions that tests the appropriate variable to gain entrance.

When a process leaves, its critical section sets ‘S’ to 0. Then the one of the waiting process will get access to enter its critical section.

Exchange instructions The exchange instructions can be defined as follows:

```
void exchange (int reg, int mem)
{
int temp;
temp = mem;
mem = reg;
reg = temp;
}
```

Mutual Exclusion Using Exchange Instructions

```
/* Mutual Exclusion */;
int const n = /* number of processes */;
int S;
void P(int i)
{
int ki = 1;
```

```

while (true)
{
do
exchange(ki, S);
while(ki! = 0);
/* critical section */;
S = 0;
/* remainder */;
}
}
void main( )
{
S = 0;
begin (P(1), P(2), ... P(n));
}

```

A shared variable 'S' is initialized to '0'. Each process uses a local variable *ki* that is initialized to 1. The only process that may enter its critical section is one that finds 'S' equal to 0. It excludes all other processes from critical section by setting 'S' to 1. When a process leaves its critical section, it resets 'S' to 0 allowing another process to its critical section.

Advantages of using machine instruction approach

1. Applicable to any number of processes on either a single processor or multiple processors, sharing the main memory.
2. Simple and easy to verify.
3. Used to support multiple critical sections.

Disadvantages

1. Busy waiting
2. Starvation is possible
3. Deadlock is possible

OTHER MECHANISMS FOR MUTUAL EXCLUSION

Let us discuss OS and programming language mechanisms that are used to provide concurrency.

Semaphores

Semaphore is an integer value used for signalling among processes.

There are two types of semaphores as follows:

1. Binary semaphore
2. Counting (or) general semaphore

Counting or general semaphore

Three operations may be performed on a semaphore all of which are atomic:

- Initialize
- Decrement
- Increment

The working of a counting semaphore with its operations is defined as below:

1. The semaphore may be initialized to a non-negative integer value.
2. The semwait operation decrements the semaphore value. If the value becomes negative, then the process executing the semwait is blocked, otherwise the process continues execution.
3. The semsignal operation increments the semaphore value. If the resulting value is less than or equal to zero, then one of the processes blocked by a semwait operation, if any, is unblocked.

Example: Let the semaphore value $S = 3$.

If the semaphore value is positive, then that value gives the number of processes that can issue a wait and immediately continue to execute.

Let five processes P_1, P_2, P_3, P_4, P_5 are going to execute a critical section code based on the semaphore value $S = 3$.

$$\begin{aligned}
 S &= \boxed{3} \leftarrow \text{Initially} \\
 &\downarrow P_1 \text{ sem wait} \\
 S &= \boxed{2} \geq 0, P_1 \text{ executes} \\
 &\downarrow P_2 \text{ sem wait} \\
 S &= \boxed{1} \geq 0, P_2 \text{ executes} \\
 &\downarrow P_3 \text{ sem wait} \\
 S &= \boxed{0} \geq 0, P_3 \text{ executes} \\
 &\downarrow P_4 \text{ sem wait} \\
 S &= \boxed{-1} \not\geq 0, P_4 \text{ Blocked} \\
 &\downarrow P_5 \text{ sem wait} \\
 S &= \boxed{-2} \not\geq 0, P_5 \text{ Blocked}
 \end{aligned}$$

Initially, the 'S' value 3 means that at a time three processes can issue a 'wait' signal and continue execution.

Whenever S becomes 0, the next process which executes 'wait' operation will be blocked.

Here P_4 is blocked as it operates on the semaphore when $S = 0$.

If the semaphore value becomes negative, it specifies the number of processes waiting to be unblocked.

$S = -2$ means two processes are waiting to be unblocked.

Definition of semwait and semsignal operations

```

struct semaphore
{
int semvalue;
Queue type Queue;
};
void semwait(semaphore S)

```

```

{
S.semvalue--;
if (S.semvalue < 0)
{
/* place the process in S. Queue */;
/* Block this process */;
}
}
void semsignal(semaphore S)
{
S.semvalue++;
if (S.semvalue <= 0)
{
/* remove a process from S. Queue */;
/* Place the process in ready queue */;
}
}

```

Advantages

1. Because the waiting processes will be permitted to enter their Critical Section in a FCFS order, so the requirement of bounded waiting is met.
2. CPU cycles are saved here as waiting process does not perform any busy waiting.

Disadvantages

1. Complex to implement, since it involves implementation of FCFS.
2. Context switching is more, so more overheads are involved.

Binary semaphore It is a semaphore that takes on only the values 0 and 1.

The operations performed on a binary semaphore are as follows:

1. A binary semaphore may be initialized to 0 or 1.
2. The semwaitB operation checks the semaphore value. If the value is 0, then the process executing the semwaitB is blocked. If the value is 1, then the value is changed to 0 and the process continues execution.
3. The semsignalB operation checks to see if any processes are blocked on this semaphore. If so, then a process blocked by a semwaitB operation is unblocked. If no processes are blocked, then the value of the semaphore is set to 1.

Definition of semwaitB and semsignalB

```

struct binary-semaphore
{
enum {zero, one} value;
Queue-type Queue;
};
void semwaitB(binary-semaphore S)
{

```

```

if (S.value == one)
S.value = zero;
else
{
/* Place this process in S.Queue */;
/* Block this process */;
}
}
void semsignalB(binary-semaphore S)
{
if(S.Queue is empty)
S.value = one;
else
{
/* remove a process from S.Queue */;
/* Place process in ready list */;
}
}

```

Advantages

1. The implementation of binary semaphore is extremely simple.

Disadvantages

1. It does not meet the requirement of Bounded waiting.
2. A process, waiting to enter its Critical Section, will perform Busy waiting, thus wasting CPU cycles.

Notes:

1. Binary semaphores have the same expressive power as general semaphores.
2. *MUTEX*: It is similar to binary semaphore. The key difference between the two is that the process that locks the mutex must be the one to unlock it.
3. Both counting semaphores and binary semaphores use a queue to hold processes waiting on the semaphore. The order in which the processes removed from a Queue is FIFO, that is, the process that has been blocked the longest is released from Queue first.
4. A semaphore whose definition includes the order of removal of Blocked processes is referred as a strong semaphore otherwise it is a weak semaphore.
5. Strong semaphores guarantee freedom from starvation.

Mutual Exclusion Using Semaphores

```

const int n = /* number of processes */;
Semaphore S = 1;
void P(int i)
{
while (true)
{
Semwait(S);
/* critical section */;
Semsignal(S);

```

```

/* remainder */;
}
}
void main( )
{
begin (P(1), P(2), ... P(n));
}

```

If no process is executing in Critical Section, then the semaphore value is 1. The first process that is executing 'wait' operation will decrement value to 0 and enter its critical section. The process which execute 'wait' operation while a cooperating process is executing in its critical section, will find the semaphore value to 0 and keep looping in the 'while-loop' of 'wait' operation. Spinning of a waiting process in the while-loop, the binary semaphores are also known as spin locks. When the process executing in the CS makes an exit (from CS), it will execute the 'signal' operation and increment the semaphore value to 1.

At a time, only one of the cooperating processes can enter critical section with the condition that the wait operation is executed automatically. Mutual Exclusion is satisfied.

Example: Consider the following figure (Figure 1) which shows the possible sequence of three processes using mutual exclusion with a semaphore, S. Processes P, Q, R accesses a shared resource protected by the semaphore S.

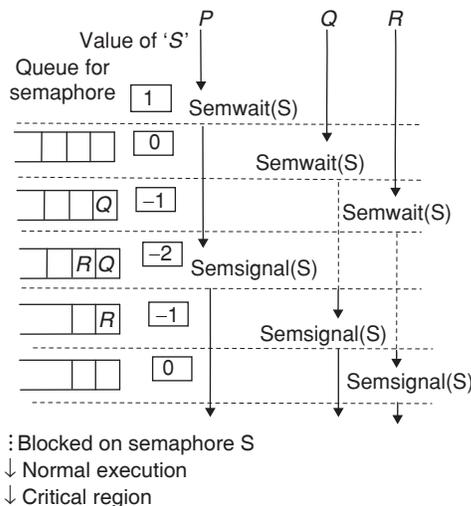


Figure 1 Mutual exclusion using semaphore.

Progress Using Semaphores

When no process is executing in the CS, the semaphore value will be 1. Then, one of the waiting process looping in the while loop of wait-operation will find the semaphore value of 1, exit from the while-loop, decrement the semaphore value to 0 and enter CS. Thus, if no process is executing in critical section and some are waiting to enter, then one of the waiting processes will enter its critical section immediately.

Bounded waiting using semaphores

Actually, one of the waiting processes will get entry into its CS when an operating process executing in its critical section exits. As this selection process is arbitrary, so a process waiting to enter its CS is likely to face starvation. So, the requirement of bounded waiting is not met.

CLASSICAL PROBLEMS OF SYNCHRONIZATION

We will discuss three problems of synchronization:

1. Bounded buffer problem
2. Readers/writers problem
3. Dining philosophers problem

Producer-Consumer Problem

1. Producer inserts item in the buffer
2. Updates Insertion pointer
3. Consumer consumes items in the buffer
4. Updates removal pointer
5. Both update information about how full, how empty the buffer.
6. Prevents buffer overflow, prevents buffer underflow, proper synchronization.

Producer	Consumer
repeat	repeat
produce item v;	while (in <= out);
b[in] = v;	w = b[out];
in = in + 1;	out = out + 1;
forever;	consume w;
	forever;

Table 1 Producer consumer problem solution using semaphores

Producer	Consumer
repeat	repeat
produce item v;	while (in <= out);
SemwaitB(S);	SemsignalB(S);
b[in] = v;	w = b[out];
in = in + 1;	out = out + 1;
SemsignalB(S);	SemsignalB(S);
forever;	consume w;
	forever;

If producer is slow or late, then consumer will busy at the while statement.

Table 2 Improved solution

Producer	Consumer
repeat	repeat
produce item v;	Semwait(n);
Semwait(S);	Semwait(S);
b[in] = v;	w = b[out];
in = in + 1;	out = out + 1;
Semsignal(S);	Semsignal(S);
Semsignal(n);	consume w;
forever;	forever;

The initial value of n and S are $n = 0, S = 1$. (n is the number of items in the buffer).

Table 3 *Producer consumer bounded buffer problem*

Producer	Consumer
repeat	repeat
produce item v ;	while($in == out$)
while($(in + 1) \% n == out$)	no operation;
no operation;	$w = b[out]$;
$b[in] = v$;	$out = (out + 1) \% n$;
$in = (in + 1) \% n$;	consume w ;
forever;	forever;

The buffer size is enforced using another counting semaphore.

Table 4 *Producer consumer bounded buffer problem solution*

Producer	Consumers
Repeat	repeat
Produce item v ;	Semwait(e);
Semwait(e);	Semwait(S);
Semwait(S);	$w = b[out]$;
$b[in] = v$;	$out = (out + 1) \% n$;
$in = (in + 1) \% n$	Semsignal(S);
Semsignal(S);	Semsignal(e);
Semsignal(e);	consume w ;
forever;	forever;

The initial value of buffer size, e is the size of the bounded buffer.

Observations on semaphores:

1. Semaphores are easy to use.
2. wait() and signal() are to be implemented as atomic operations.

Problems:

1. signal() and wait() may be exchanged by the programmer, this may result in deadlock or violation of mutual exclusion.

Readers/Writers Problem

1. A reader reads data.
2. A writer writes data.
3. Data is shared among a number of processes.
4. Multiple readers may read the data simultaneously, that is, concurrently.
5. Only one writer can write the data any time, that is, no reader should be present.
6. A reader and writer cannot access data simultaneously.
7. Locking table: Whether any two can be in the critical section simultaneously is shown in the table.

	Reader	Writer
Reader	OK	NO
Writer	NO	NO

Solution: Readers have priority; if a reader is in CS, any number of readers could enter irrespective of any writer waiting to enter critical section

Writer	Reader
while (true)	while (true)
{	{
Semwait (S);	Semwait (x);
writeunit ();	Num = Num + 1;
Semsignal (S);	if (Num == 1)
}	Semwait (S);
	Semsignal (x);
	Readunit ();
	Semwait (x);
	Num = Num - 1;
	if (Num == 0)
	Semsignal (S);
	Semsignal (x);
	}

Semaphore ‘ S ’ is used to enforce mutual exclusion.

Semaphore ‘ x ’ is used to assure that ‘Num’ is updated properly.

Solution: If a writer wants critical section as soon as the critical section is available, writer enters it.

Dining Philosophers Problem

N philosophers are sitting around a dining table. There are N plates placed on the table such that each plate is in front of a philosopher and N forks placed between the plates. There is a bowl of Noodles placed at the centre of the table. Whenever a philosopher feels hungry, he tries to pick two forks which are shared with his nearest neighbour. If any of his neighbours happens to be eating at the time, the philosopher has to wait. Whenever a hungry philosopher gets two forks, he pours noodles into his plate. After he finishes, he places the chopsticks back onto the table and starts thinking. Now forks are available for neighbours.

Solution:

```
# define N 5 /* Number of philosophers*/
voidphilosopher(int i) /* philosopher number,
from 0 to 4*/
{
while (true)
{
think(); /*philosopher is thinking*/
take_fork(i) /*take left fork*/
take_fork ((i+ 1)% N); /* take right fork; %
is modulo operator*/
eat( );
put _fork ( ); /* put left back on the table*/
put _fork ((i+1) % N); /* put right fork
back on the table */
```

Notes:

1. This solution leads to deadlock.
2. Everyone picks the left fork and indefinitely wait for right fork causing starvation.

MONITORS

Monitor is a programming language construct that encapsulates variables, access procedures and initialization code within an abstract data type. The monitor's variable may only be accessed via its access procedures and only one process may be actively accessing the monitor at any one time. The access procedures are critical sections. A monitor may have a queue of processes that are waiting to access it.

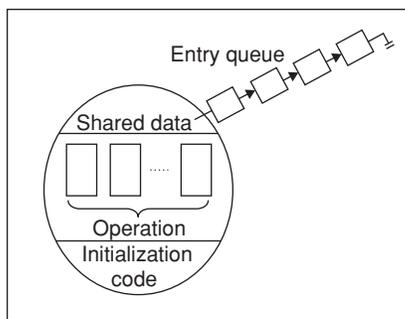
1. If the data in a monitor represent some resource, then the monitor provides a mutual exclusion facility for accessing the resource.
2. A monitor supports synchronization by the use of condition variables that are contained within the monitor and accessible only within the monitor.
3. Operations on conditional variables:
 - `cwait(c)`: Suspend execution of calling process on condition *c*. The monitor is now available for use by another process.
 - `csignal(c)`: Resume execution of some process blocked after a wait on the same condition. If there are several such processes, choose one of them; if there is no such process, do nothing.

Monitor syntax is as follows:

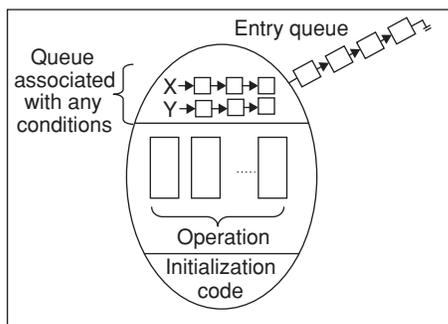
Monitor monitor - name

```
{
shared variable declarations;
Procedure body P1 (...) { }
Procedure body P2 (...) { }
Procedure body Pn (...) { }
initialization code { }
}
```

Schematic view of a monitor:



Schematic View of a monitor with condition variables:



MESSAGE PASSING

When processes interact with one another, two fundamental requirements must be satisfied:

1. Synchronization
2. Communication

One approach to provide both of these is *message passing*.

The primitive functions in message passing are send (destination, message) receive (source, message)

Design Characteristics of Message Systems for IPC and Synchronization

Synchronization There must be some synchronization existing between two processes to communicate with each other.

- **Send:** When a 'send' primitive is executed in a process, then the sender may
 - blocked or
 - non-blocked
 - until the message is received.
- **Receive:** When a process issues a 'Receive' primitive there are two possibilities:
 1. If a message has previously been sent, the message is received and execution continues.
 2. If there is no waiting message then either
 - (i) The process is blocked until a message arrives or
 - (ii) The process continues to execute, abandoning the attempt to receive.

Thus, both the sender and receiver may be in one of

1. Blocking send, blocking receive: Allows tight synchronization.
2. Non-blocking send, blocking receive:
 - Useful synchronization
 - Possibility of generating repeated messages
3. Non-blocking send, non-blocking receive: No need to wait

Addressing Two types of addressing methods:

1. Direct addressing
2. Indirect addressing

Direct Addressing

The send primitive includes a specific identifier of the destination process. The 'receive' primitive can be handled in one of two ways:

Explicit

The process must know ahead of time from which process a message is expected. Useful for cooperating concurrent processes.

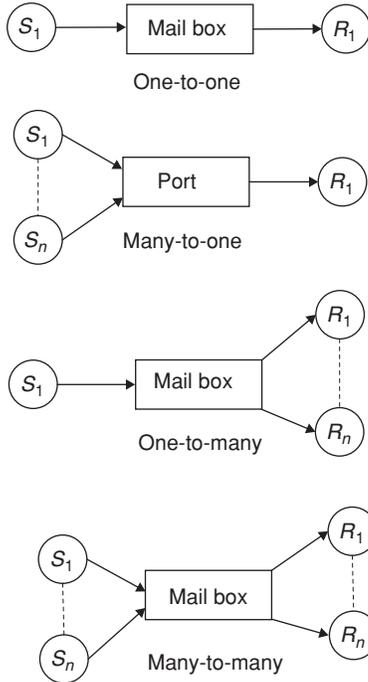
Implicit

The 'source' parameter of 'receive' primitive possesses a value returned when the receive operation has been performed. Example, Printer server.

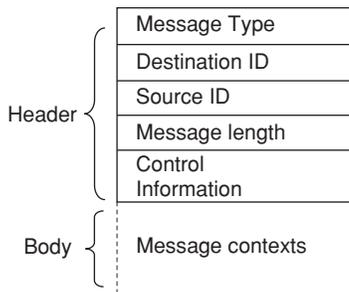
Indirect Addressing

Messages are not sent directly from sender to receiver but rather are sent to a shared data structure consisting of queues that can temporarily hold messages. These queues are referred to as mailboxes. The relationship between the sender and receiver is

- one-to-one
- many-to-one
- one-to-many
- many-to-many



Message format The general message format will be



Queuing discipline The queuing discipline may be

- FIFO
- Priority

Mutual Exclusion using message passing:

```
const int n = /* number of processes */;
void P(int i)
{
    message msg;
    while (true)
    {
        receive (box, msg);
        /* critical section */;
        send (box, msg);
        /* remainder */;
    }
}

void main( )
{
    create mailbox(box);
    send(box, null);
    begin (P(1), P(2) ... P(n));
}
```

1. Here a set of concurrent processes that share a mailbox 'box', which can be used by all processes to send and receive.
2. The mailbox is initialized to contain a single message with null content. A process wishing to enter its critical section first attempts to receive a message.
3. If the mail box is empty, then the process is blocked.
4. Once a process acquired the message, it performs its critical section and then places the message back into the mailbox.
5. Hence, the message functions as a token that is passed from process to process.

EXERCISES

Practice Problems I

Directions for questions 1 to 17: Select the correct alternative from the given choices.

1. The value of a counting semaphore is 7. Then 15 wait operations and 10 signal operations were completed on this semaphore. The resulting value of semaphore is
 (A) 5 (B) 7
 (C) 2 (D) 0
2. At a particular time of computation, the value of counting semaphore is 7. Then 20 wait operations and 'x'

signal operations were completed on this semaphore. If the final value of the semaphore is 5, what is x?

- (A) 18 (B) 13
 (C) 5 (D) 0
3. A process using a semaphore has a start value of 1 for its semaphore. Since the start of execution of the program, 12 signal operations were completed. How many wait operations have been completed so far if the current value of semaphore is 6?
 (A) 1 (B) 5
 (C) 7 (D) 11

4. It is found that a program has multiple critical sections. Choose correct statements from below:
- (i) Multiple semaphores are needed for handling the situation.
 - (ii) A single semaphore that uncompresses all the critical section is sufficient and is also more efficient.
 - (iii) To get better control of the code, monitors need to be implemented.
- (A) (i) and (ii) (B) (ii) and (iii)
 (C) (i) and (iii) (D) (i), (ii), (iii)

5. Consider the below pseudocode:

```
semaphore S = 1;
semaphore E = 1;
if(thread_count++ < 100)
spawnnewthread( );
wait(E);

// critical section - begin
-----
-----
// critical section - end
signal(S);
```

Assume that above pseudocode gets called a hundred times, what is the count of semaphore *E*?

- (A) 0 (B) 1
 (C) -99 (D) -100
6. Consider the below code for a process *i*:

```
-----
-----
flag[i] = true;
if(turn == i and flag [i] == true)
/* critical section begin */
counter++;
/* critical section end */
turn = x;
-----
-----
```

If the value of a counter started, what would be the value of 'counter' count at the end of the program:

- (A) Semaphore count
 (B) Thread count
 (C) Concurrency count
 (D) Deadlock process count
7. Consider the below pseudocode:

```
function waitB(s)
{
if(s.value ==1)
s.value = 0;
else
place the process in the Queue;
```

```
}
function signalB(s)
{
s.value = 1;
}
```

What does the code most likely behave as

- (A) general semaphore (B) weak semaphore
 (C) binary semaphore (D) mutex
8. A shared variable *x*, initialized to 0 is operated on by four concurrent processes *P, Q, R, S* as follows:

<pre>P(x) { wait(); read(x); increment x by 1; store(x); signal(); } R(x) { wait(); read (x); decrement x by 2; store (x); signal(); }</pre>	<pre>Q(x) { wait(); read(x); increment x by 1; store(x); signal(); } S(x) { wait(); read(x); decrement x by 2; store(x); signal(); }</pre>
---	--

A counting semaphore '*N*' is used by the processes whose value is initialized to 2. What is the maximum possible value of '*x*' after all processes complete execution?

- (A) -2 (B) -1
 (C) 1 (D) 2
9. Consider the following code:

```
Program concurrency;
Var x: Integer (: = 0);
    y: Integer (: = 0);
Procedure threadA( );
begin
x = 1; /*S1*/
y = y + x; /*S2*/
end;
Procedure threadB( );
begin
y = 4;        /*S3*/
x = x + 5; /*S4*/
end;
begin /*mainprogram*/
parbegin
threadA( );
threadB( );
parend;
end.
```

Suppose a process has two concurrent threads: one thread executes statements S_1 and S_2 , and the other thread executes statements S_3 and S_4 . What are the maximum possible values of x and y when the code finishes execution? (All the statements S_1, S_2, S_3 and S_4 are atomic).

- (A) $x = 6, y = 4$ (B) $x = 6, y = 5$
 (C) $x = 1, y = 5$ (D) $x = 6, y = 10$

10. Consider the following program:

```
boolean lock[2];
int turn;
void P(int id)
{
while(true)
{
lock[id] = true;
while (turn != id)
{
while (lock [1 - id])
/*do nothing*/
turn = id;
}
/*critical section*/
lock[id] = false;
/*remainder*/
}
}
void main( )
{
lock[0] = false;
lock[1] = false;
turn = 0;
parbegin (P(0), P(1));
}
```

Which of the following statements is correct for two processes executing this code?

- (A) Given program provides mutual exclusion.
 (B) Given program does not provide mutual exclusion.
 (C) Given program provides mutual exclusion and also solves starvation problem.
 (D) Given program provides mutual exclusion but does not prevent from starvation.

11. Consider two process P_0 and P_1 which share the following variables:

```
boolean flag [2]; /*initially false*/
int turn;
```

These two processes, $P_i (i = 0 \text{ or } 1)$, $P_j (j = 1 \text{ or } 0)$ execute the following code:

```
do
{
flag[i] = TRUE;
while(flag [j])
```

```
{
if (turn == j)
{
flag[i] = false;
while (turn == j);
flag [i] = TRUE;
}
}
// critical section
turn = j;
flag[i] = FALSE;
// remainder.
}
```

while(TRUE);

The code satisfies

- (i) Mutual exclusion
 (ii) Progress
 (iii) Bounded waiting
 (A) (i), (ii) only (B) (ii), (iii) only
 (C) (i), (iii) only (D) (i), (ii), (iii)

12. Which of the following statement(s) is false?

- (i) Spinlocks are not appropriate for single-processor systems.
 (ii) Mailboxes may be used for synchronization.
 (iii) Message passing and semaphores do not have equivalent functionality.
 (A) (i) only (B) (iii) only
 (C) (i), (iii) (D) (i), (ii), (iii)

13. Consider the following code:

```
signal (mutex);
.....
Critical section
.....
wait (mutex);
```

Here 'mutex' is a semaphore variable, which is initialized to 1. Then

- (A) Mutual exclusion is provided
 (B) Mutual exclusion violated, if several processes are simultaneously active in their critical section.
 (C) Deadlock will occur
 (D) Starvation is possible

14. Which of the following sequence of 'wait' and 'signal' operations leads to deadlock?

(Here 'mutex' is a semaphore variable initialized to 1.)

- (A) wait (mutex);

 Critical section

 Signal (mutex);

```
(B) wait (mutex);
.....
Critical section
.....
Wait (mutex);
(C) Signal (mutex);
.....
Critical section
.....
Wait (mutex);
(D) signal (mutex);
.....
Critical section
.....
Signal (mutex);
```

15. Which of the following situation arises if a process omits the wait(S) or the signal(S) on a semaphore variable 'S' (Initially $S = 1$).

- (i) Mutual exclusion violated
(ii) Deadlock will occur
(A) (i) only (B) (ii) only
(C) Both (i) and (ii) (D) Neither (i) nor (ii)

16. consider the following shared data and code:

```
data:
int turn;
Boolean flag[2];
Code:
do
{
flag [i] = TRUE;
```

```
turn = j;
while (flag [j] && turn == j);
//critical section
flag [i] = FALSE;
//remainder
}
```

```
while (TRUE);
```

Let two processes P_i ($i = 0$ or 1) and P_j ($j = 1$ or 0) use the shared data and executes the code. Then the code provides

- (A) a solution to critical section problem
(B) mutual exclusion but not progress.
(C) progress but not mutual exclusion
(D) both mutual exclusion, progress but no bounded waiting.

17. Consider the following code that shows the structure of a process in an algorithm to solve the critical section problem for two processes.

```
var flag[2] of Boolean; /* initialized to false */
repeat
flag[i] = true;
while flag[j] do no-op;
//critical section
flag[i] = false;
// remainder
until false
```

Then which of the following statements is true?

- (A) The algorithm satisfies all the requirements of critical section problem.
(B) The algorithm satisfies only mutual exclusion and progress.
(C) The algorithm only satisfies progress requirement.
(D) The algorithm does not satisfy critical section problem requirements.

Practice Problems 2

Directions for questions 1 to 15: Select the correct alternative from the given choices.

- The result of a computation depends upon the speed of the processes involved, is said to be:

(A) Cycle stealing (B) Race condition
 (C) A time lock (D) A deadlock
- A relation between processes such that each has some part which must not be executed, while the critical section of another is being executed is known as

(A) Mutual exclusion (B) Semaphore
 (C) Multi-tasking (D) Mutli-programming
- Producer–consumer problem can be solved using

(A) Semaphores (B) Event counters
 (C) Monitors (D) All of the above
- To avoid the race condition, the number of processes allowed in critical section is

(A) 0 (B) 1
 (C) 2 (D) 3

- Mutual exclusion problem occurs between

(A) Two disjoint processes that unaware of each other
 (B) Processes that share resources
 (C) Processes directly aware of each other.
 (D) Both (A) and (B)
- Semaphores are used to solve the problem of

(A) Race condition (B) Multitasking
 (C) Mutual exclusion (D) Both (A) and (C)
- At a particular time, the value of a counting semaphore is 10. It will become 7 after

(A) 3 signal operations
 (B) 3 wait operations
 (C) 5 signal operations and 2 wait operations
 (D) None of the above
- Critical region is

(A) A part of the OS which is not allowed to be accessed by any process
 (B) A set of instructions that accesses common shared resource, which exclude one another in time

- (C) The portion of main memory, which can be accessed only by one process at a time
(D) Both (A) and (C)
9. Concurrent processes are:
(A) Processes that don't overlap in time
(B) Processes that overlap in time
(C) Processes that are executed
(D) Processes that are executed by a processor at the same time
10. Semaphore operations are atomic because they are implemented within the _____.
(A) Kernel (B) Shell
(C) User process (D) Normal process space
11. The programming language construct that provides equivalent functionality of a semaphore and better control is
(A) Signal (B) Monitor
(C) Mutex (D) Critical section.
12. What is the ideal way of emptying the queue of a strong semaphore?
(A) Random (B) LIFO
(C) FIFO (D) binary
13. What are the disadvantages of machine instruction approach?
(i) While a process is waiting for entering a critical section, process still consumes resources.
(ii) There could be starvation
(iii) There could be deadlocks
(A) (i), (ii) only (B) (ii), (iii) only
(C) (iii), (i) only (D) (i), (ii), (iii)
14. Select from below the advantages of Machine Instruction approach?
(i) Applicable to any number of processes either uni-processor or multi-processor system
(ii) Simple and easy to verify
(iii) Supports multiple critical sections
(A) (i), (ii) only (B) (ii), (iii) only
(C) (iii), (i) only (D) (i), (ii), (iii)
15. Which of the below are requirements for mutual exclusion?
(i) Only one process is allowed into critical section.
(ii) A process remains inside its critical section for finite time only.
(iii) It must be possible for a process accessing critical section to be delayed indefinitely.
(iv) A process halting in critical section must do so without interfering with other processes.
(A) (i), (ii), (iii) (B) (ii), (i), (iv)
(C) (i), (ii), (iv) (D) (i), (ii), (iii), (iv)

PREVIOUS YEARS' QUESTIONS

1. Consider these two functions and two statements S_1 and S_2 about them: [2006]

<pre>int work1 (int *a,int i, int j) { int x=a[i+2]; a[j]=x+1; return a[i+2]-3; }</pre>	<pre>int work2 (int *a,int i, int j) { int t1=i+2; int t2=a[t1]; a[j]=t2+1; return t2 - 3; }</pre>
---	--

- S1:** The transformation from work1 to work2 is valid, that is, for any program state and input arguments, work 2 will compute the same output and have the same effect on program state as work 1
- S2:** All the transformations applied to work 1 to get work 2 will always improve the performance (i.e., reduce CPU time) of work 2 compared to work 1
- (A) S_1 is false and S_2 is false
(B) S_1 is false and S_2 is true
(C) S_1 is true and S_2 is false
(D) S_1 is true and S_2 is true
2. The atomic fetch-and-set x, y instructions unconditionally sets the memory location x to 1 and fetches the old value of the of x in y without allowing any

intervening access to the memory location x . Consider the following implementation of P and V functions on binary semaphore S .

```
void P (binary-semaphore *s) {
    unsigned y;
    unsigned *x = & (s → value);
    do {
        fetch-and-set x, y ;
    } while (y) ;
}

void V (binary-semaphore *s) {
    s → value = 0;
}
```

Which one of the following is true?

- [2006]
- (A) The implementation may not work if context switching is disabled in P
(B) Instead of using fetch-and-set, a pair of normal load/store can be used
(C) The implementation of V is wrong
(D) The code does not implement a binary semaphore
3. The P and V operations on counting semaphores, where s is a counting semaphore, are defined as follows:

```
P(s) : s = s - 1;
      if s < 0 then wait;
V(s) : s = s + 1;
if s <= 0 then wake up a process waiting on s;
```

Assume that P_b and V_b , the wait and signal operations on binary semaphores, are provided. Two binary semaphores x_b and y_b are used to implement the semaphore operations $P(s)$ and $V(s)$ as follows:

```
P(s) :      P_b (x_b) ;
          s = s - 1;
          if (s < 0) {
              V_b (x_b) ;
              P_b (y_b) ;
          }
          else V_b (x_b) ;
```

```
V(s) :      P_b (x_b) ;
          s = s + 1;
          if (s <= 0) V_b (y_b) ;
          V_b (x_b) ;
```

The initial values of x_b and y_b are respectively

[2008]

- (A) 0 and 0
- (B) 0 and 1
- (C) 1 and 0
- (D) 1 and 1

4. Consider a system with four types of resources R_1 (3 units), R_2 (2 units), R_3 (3 units), R_4 (2 units). A non-pre-emptive resource allocation policy is used. At any given instance, a request is not entertained if it cannot be completely satisfied. Three processes P_1, P_2, P_3 request the resources as follows if executed independently.

[2009]

Process P_1 :	Process P_2 :	Process P_3 :
$t=0$: requests 2 units of R_2	$t=0$: requests 2 units of R_3	$t=0$: requests 1 unit of R_4
$t=1$: requests 1 unit of R_3	$t=2$: requests 1 unit of R_4	$t=2$: requests 2 units of R_1
$t=3$: requests 2 units of R_1	$t=4$: requests 1 unit of R_1	$t=5$: releases 2 units of R_1
$t=5$: releases 1 unit of R_2 and 1 unit of R_1 .	$t=6$: releases 1 unit of R_3	$t=7$: requests 1 unit of R_2
$t=7$: releases 1 unit of R_3	$t=8$: Finishes	$t=8$: requests 1 unit of R_3
$t=8$: requests 2 units of R_4		$t=9$: Finishes
$t=10$: Finishes		

Which one of the following statements is true if all three processes run concurrently starting at time $t=0$?

- (A) All processes will finish without any deadlock
 - (B) Only P_1 and P_2 will be in deadlock.
 - (C) Only P_1 and P_3 will be in a deadlock.
 - (D) All three processes will be in deadlock.
5. The enter_CS() and leave_CS() functions to implement critical section of a process are realized using test-and-set instruction as follows: [2009]

```
void enter_CS(X)
{
    while (test-and-set(X)) ;
}
void (leave_CS(X))
{
    X=0;
}
```

In the above solution, X is a memory location associated with the CS and is initialized to 0. Now consider the following statements:

- (i) The above solution to CS problem is deadlock-free.
- (ii) The solution is starvation free.
- (iii) The processes enter CS in FIFO order.
- (iv) More than one process can enter CS at the same time.

Which of the above statements is true?

- (A) (i) only
- (B) (i) and (ii)

- (C) (ii) and (iii)
- (D) (iv) only

6. Consider the methods used by processes P_1 and P_2 for accessing their critical sections whenever needed, as given below. The initial values of shared Boolean variables S_1 and S_2 are randomly assigned.

[2010]

Method used by P_1	Method used by P_2
while ($S_1 = S_2$);	while ($S_1 \neq S_2$);
Critical section	Critical section
$S_1 = S_2$;	$S_2 = \text{not } (S_1)$;

Which one of the following statements describes the properties achieved?

- (A) Mutual exclusion but not progress
 - (B) Progress but not mutual exclusion
 - (C) Neither mutual exclusion nor progress
 - (D) Both mutual exclusion and progress
7. The following program consists of three concurrent processes and three binary semaphores. The semaphores are initialized as $S_0 = 1, S_1 = 0, S_2 = 0$.

Process P_0	Process P_1	Process P_2
while (true) {	wait (S1);	wait (S2);
wait (S0);	Release (S0);	release (S0);
print '0'		
release (S1);		
release (S2);		
}		

How many times will process P_0 print '0'? [2010]

- (A) At least twice (B) Exactly twice
(C) Exactly thrice (D) Exactly once

8. Fetch_And_Add(X, i) is an atomic Read-Modify-Write instruction that reads the value of memory location X , increments it by the value i , and returns the old value of X . It is used in the pseudocode shown below to implement a busy wait lock. L is an unsigned integer shared variable initialized to 0. The value of 0 corresponds to lock being available, while any non-zero value corresponds to the lock being not available.

[2012]

```
AcquireLock(L) {
while (Fetch_And_Add(L, 1))
L = 1;
}
ReleaseLock (L) {
L = 0;
}
```

This implementation

- (A) fails as L can overflow
(B) fails as L can take on a non-zero value when the lock is actually available
(C) works correctly but may starve some processes
(D) works correctly without starvation
9. A shared variable x , initialized to 0, is operated on by four concurrent processes W, X, Y, Z as follows. Each of the processes W and X reads x from memory, increments by one, stores it to memory, and then terminates. Each of the processes Y and Z reads x from memory, decrements by two, stores it to memory, and then terminates. Each process before reading x invokes the P operation (i.e., wait) on a counting semaphore S and invokes the V operation (i.e., signal) on the semaphore S after storing x to memory. Semaphore S is initialized to 2. What is the maximum possible value of x after all processes complete execution? [2013]

- (A) -2 (B) -1
(C) 1 (D) 2

10. A certain computation generates two arrays 'a' and 'b' such that $a[i] = f(i)$ for $0 \leq i < n$ and $b[i] = g(a[i])$ for $0 \leq i < n$. Suppose this computation is decomposed into two concurrent processes X and Y such that X computes the array 'a' and Y computes the array 'b'. The processes employ two binary semaphores R and S , both initialized to zero. The array 'a' is shared by the two processes. The structures of the processes are shown below.

Process X:

```
private i;
for (i=0; i<n; i++) {
    a[i] = f(i);
    ExitX(R, S);
}
```

Process Y:

```
private i;
for (i=0; i<n; i++) {
    EntryY(R, S);
    b[i] = g(a[i]);
}
```

Which one of the following represents the correct implementations of Exit X and Entry Y ?

[2013]

- (A) ExitX(R, S) {
 P(R);
 V(S);
}
EntryY(R, S) {
 P(S);
 V(R);
}
- (B) ExitX(R, S) {
 V(R);
 V(S);
}
EntryY(R, S) {
 P(R);
 P(S);
}
- (C) ExitX (R, S) {
 P(S);
 V(R);
}
EntryY (R, S) {
 V(S);
 P(R);
}
- (D) ExitX (R, S) {
 V(R);
 P(S);
}
EntryY (R, S) {
 V(S);
 P(R);
}

11. Consider the procedure below for the *producer-consumer* problem which uses semaphores; [2014]

```
semaphore n = 0;
semaphore s = 1;
void producer ()
{
while (true)
{
produce ( )
semWait (s);
addToBuffer ();
semSignal (s);
semSignal (n);
}
}
void consumer ()
{
while (true)
{
semWait (s);
semWait (n);
remove FromBuffer ();
semSignal (s);
consume ( ) ;
}
}
```

Which one of the following is true?

- (A) The producer will be able to add an item to the buffer, but the consumer can never consume it.
(B) The consumer will remove no more than one item from the buffer.
(C) Deadlock occurs if the consumer succeeds in acquiring semaphore s when the buffer is empty.
(D) The starting value for the semaphore n must be 1 and not 0 for deadlock free operation.

12. The following two function $P1$ and $P2$ that share a variable B with an initial value of 2 execute concurrently. [2015]

```

P1 ( ) {
    C = B - 1;
    B = 2 * C;
}

P2 ( ) {
    D = 2 * B;
    B = D - 1;
}
    
```

The number of distinct values that B can possibly take after the execution is _____

13. Two processes X and Y need to access a critical section. Consider the following synchronization construct used by both the processes [2015]

Process X	Process Y
<pre> /* other code for process X */ while (true) { varP = true; while(varQ == true) { /* Critical Section */ varP = false; } } /* other code for process X */ </pre>	<pre> /* other code for process Y */ while (true) { varQ = true; while (varP == true) { /* Critical Section */ varQ = false; } } /* other code for process Y */ </pre>

Here, $varP$ and $varQ$ are shared variables and both are initialized to false. Which one of the following statements is true?

- (A) The proposed solution prevents deadlock but fails to guarantee mutual exclusion.
 (B) The proposed solution guarantees mutual exclusion but fails to prevent deadlock.
 (C) The proposed solution guarantees mutual exclusion and prevents deadlock.
 (D) The proposed solution fails to prevent deadlock and fails to guarantee mutual exclusion.
14. Consider the following proposed solution for the critical section problem. There are n process: $P_0 \dots P_{n-1}$. In the code, function $pmax$ returns an integer not smaller than any of its arguments. For all i , $t[i]$ is initialized to zero. [2016]
- ```

do {
 c[i] = 1; t[i] = pmax (t[i],,
 t[n - 1]) + 1; c[i] = 0;
 for every $j \neq i$ in $(0, \dots, n - 1)$ {
 while (c[j]);
 while (t[j] != 0 && t[j] <= t[i]);
 }
 Critical Section;
 t[i] = 0;
 Remainder Section;
} while (true);

```

Which one of the following is TRUE about the above solution?

- (A) At most one process can be in the critical section at any time.  
 (B) The bounded wait condition is satisfied.  
 (C) The progress condition is satisfied.  
 (D) It cannot cause a deadlock.
15. Consider the following two - process synchronization Solution.

| Process 0                     | Process 1                     |
|-------------------------------|-------------------------------|
| Entry: loop while (turn = 1); | Entry: loop while (turn = 0); |
| (Critical section)            | (Critical section)            |
| Exit: turn = 1;               | Exit: turn = 0;               |

The shared variable  $turn$  is initialized to zero. Which one of the following is TRUE? [2016]

- (A) This is a correct two - process synchronization Solution.  
 (B) This Solution violates mutual exclusion requirement.  
 (C) This Solution violates progress requirement.  
 (D) This Solution violates bounded wait requirement.
16. Consider a non-negative counting semaphore  $S$ . The operation  $P(S)$  decrements  $S$ , and  $V(S)$  increments  $S$ . During an execution, 20  $P(S)$  operations and 12  $V(S)$  operations are issued in some order. The largest initial value of  $S$  for which at least one  $P(S)$  operation will remain blocked is \_\_\_\_\_. [2016]
17. A multithreaded program  $P$  executes with  $x$  number of threads and uses  $y$  number of locks for ensuring mutual exclusion while operating on shared memory locations. All locks in the program are *non-reentrant*, i.e., if a thread holds a lock  $l$ , then it cannot re-acquire lock  $l$  without releasing it. If a thread is unable to acquire a lock, it blocks until the lock becomes available. The *minimum* value of  $x$  and the *minimum* value of  $y$  together for which execution of  $P$  can result in a deadlock are: [2017]

- (A)  $x = 1, y = 2$                       (B)  $x = 2, y = 1$   
 (C)  $x = 2, y = 2$                       (D)  $x = 1, y = 1$

18. Consider the following solution to the producer-consumer synchronization problem. The shared buffer size is  $N$ . Three semaphores empty, full and mutex are defined with respective initial values of 0,  $N$  and 1. Semaphore empty denotes the number of available slots in the buffer, for the consumer to read from. Semaphore full denotes the number of available slots in the buffer, for the producer to write to. The placeholder variables, denoted by  $P, Q, R$ , and  $S$ , in the code below can be assigned either empty or full. The valid semaphore operations are: wait () and signal ().

| Producer                                                                                                         | Consumer                                                                                                               |
|------------------------------------------------------------------------------------------------------------------|------------------------------------------------------------------------------------------------------------------------|
| <pre>do {   wait(P);   wait (mutex);   //Add item to   buffer   signal (mutex);   signal (Q); } while (1);</pre> | <pre>do {   wait(R);   wait (mutex);   //Consume item   from buffer   signal (mutex);   signal (S); } while (1);</pre> |

Which one of the following assignments to  $P$ ,  $Q$ ,  $R$  and  $S$  will yield the correct solution? **[2018]**

- (A)  $P$ : full,  $Q$ : full,  $R$ : empty,  $S$ : empty  
 (B)  $P$ : empty,  $Q$ : empty,  $R$ : full,  $S$ : full  
 (C)  $P$ : full,  $Q$ : empty,  $R$ : empty,  $S$ : full  
 (D)  $P$ : empty,  $Q$ : full,  $R$ : full,  $S$ : empty

## ANSWER KEYS

### Practice Problems 1

1. C    2. A    3. C    4. D    5. C    6. C    7. C    8. D    9. D    10. B  
 11. D    12. B    13. B    14. B    15. C    16. A    17. D

### Practice Problems 2

1. B    2. A    3. D    4. B    5. D    6. C    7. B    8. B    9. B    10. A  
 11. B    12. C    13. D    14. D    15. C

### Previous Years' Questions

1. D    2. A    3. C    4. A    5. A    6. A    7. A    8. B    9. D    10. C  
 11. C    12. 3    13. A    14. A    15. C    16. 7    17. D    18. C

# Chapter 3

## Deadlock and CPU Scheduling

### LEARNING OBJECTIVES

- Deadlock
- System model
- Bridge crossing
- Resources
- Resource allocation graph
- Methods of handling deadlocks
- Deadlock prevention
- Resource allocation denial (OR) banker's algorithm
- Deadlock detection
- Dining philosophers problem
- Scheduling algorithms
- Scheduling policies
- Round Robin scheduling
- Shortest remaining time
- Highest response ratio next
- Multilevel feedback queue scheduling

### DEADLOCK

It is a situation where a process or set of processes is blocked, waiting for some resource that is held by other waiting processes.

### System Model

Let the resource types be  $R_1, R_2 \dots R_m$  (like CPU cycles, memory space, input/output (I/O) devices, etc.). Each resource type  $R_i$  has  $W_i$  instances; each process utilizes a resource as follows:

**Request** A process, needing a resource, will request the operating system (OS) for assignment of the needed resource. Then the process waits, till operating system assigns it an instance of the requested resource.

**Assignment** The OS will assign to the requesting process an instance of the requested resource, whenever, it is available. Then, the process comes out of its waiting state.

**Use** The process will use the assigned resource. In case, the resource is non-sharable, the process will have exclusive access to it.

**Release** After the process finished with the use assigned resource, it will return the resource to the system pool. The released resource can now be assigned to another waiting process.

**Example:**

### Bridge crossing



Traffic is allowed only in one direction. Each section of a bridge can be viewed as a resource.

If a deadlock occurs, it can be resolved if one car backs up (pre-empt resources and rollback). Several cars may have to be backed up if a deadlock occurs.

Problem of starvation (infinite wait) is possible.

### Resources

Types of resources:

- Reusable resources
- Consumable resources

**Reusable resources** These resources can be safely used by only one process at a time, and are not depleted by that use.

**Examples:** Processors, I/O channels, main and secondary memory, devices and files, etc.

Consider two processes  $P$  and  $Q$  that compete for exclusive access to a disk file  $D$  and tape drive  $T$ . Let their implementation is as shown below:

Table 1 Process P

| Step  | Action           |
|-------|------------------|
| $P_0$ | Request ( $D$ )  |
| $P_1$ | Lock ( $D$ )     |
| $P_2$ | Request ( $T$ )  |
| $P_3$ | Lock ( $T$ )     |
| $P_4$ | Perform function |
| $P_5$ | Unlock ( $D$ )   |
| $P_6$ | Unlock ( $T$ )   |

Table 2 Process Q

| Step  | Action           |
|-------|------------------|
| $Q_0$ | Request ( $T$ )  |
| $Q_1$ | Lock ( $T$ )     |
| $Q_2$ | Request ( $D$ )  |
| $Q_3$ | Lock ( $D$ )     |
| $Q_4$ | Perform Function |
| $Q_5$ | Unlock ( $T$ )   |
| $Q_6$ | Unlock ( $D$ )   |

$P$  and  $Q$  are executing on a single processor in interleaved fashion. Then deadlock occurs if each process holds one resource and requests the other.

For example, deadlock occurs if the multiprogramming system interleaves the execution of the two processes as follows:

$$P_0, P_1, Q_0, Q_1, P_2, Q_2$$

One strategy to deal with this type of deadlocks is to impose system design constraints concerning the order in which resources can be requested.

**Consumable resources** A consumable resource is one that can be created and destroyed. There is no limit on the number of consumable resources of a particular type.

**Examples:** Interrupts, signals, messages, etc. Consider the following pair of processes, in which each process attempts to receive a message from the other process and then send a message to the other process:

| $P_1$                | $P_2$                |
|----------------------|----------------------|
| .....                | .....                |
| receive ( $P_2$ );   | receive ( $P_1$ );   |
| .....                | .....                |
| send ( $P_2, M_1$ ); | send ( $P_1, M_2$ ); |

Deadlock occurs in above case, if the receive is blocking. There is no single effective strategy that can deal with all types of deadlocks.

### Deadlock Characteristics

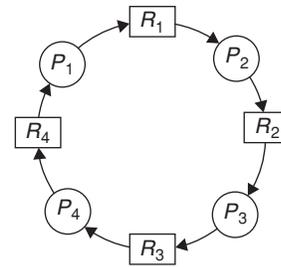
Deadlock is an undesirable state of the system. The following are the four conditions that must hold simultaneously for a deadlock to occur:

**Mutual exclusion** A resource can be used by only one process at a time. If another process requests for that resource then the requesting process must be delayed until the resource has been released.

**Hold-and-wait** Some processes must be holding some resources in a non-sharable mode and at the same time must be waiting to acquire some more resources, which are currently held by other processes in a non-sharable mode.

**No pre-emption** Resources granted to a process can be released back to the system only as a result of the voluntary action of that process, after the process has completed its task.

**Circular wait** Deadlocked processes are involved in a circular chain such that each process holds one or more resources being requested by the next process in the chain.



### Resource Allocation Graph

A deadlock is described in terms of a directed graph called a system Resource Allocation Graph (RAG). It consists of two sets:

1. The set of vertices,  $V$
2. The set of edges,  $E$

The set of vertices is again divided into two categories.

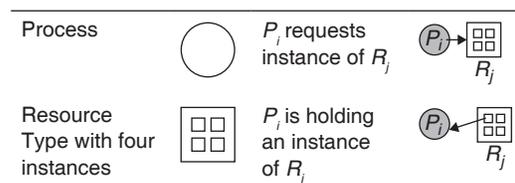
The set of all active processes in the system is  $P = \{P_1, P_2, \dots, P_n\}$  and the set of all different type of resources i.e.,  $R = \{R_1, R_2, \dots, R_m\}$

There are two types of edges in the RAG:

1. A directed edge from the process  $P_i$  to resources type  $R_j$  and is denoted by  $P_i \rightarrow R_j$ . It signifies that the  $i$ th process is requesting one unit of the resource type  $j$ . This edge is request edge.
2. A directed edge from the resource  $R_i$  to process  $P_j$  denoted by  $R_i \rightarrow P_j$ . It signifies that one unit of  $i$ th resource is held by the process  $j$ . This edge is also called as an *allocation edge/assignment edge*.

### Notations Used in RAGs

We denote a process by a circle and each resource by a rectangle.



However, if we have more number of instances of a resource type, then it is denoted by more dots.

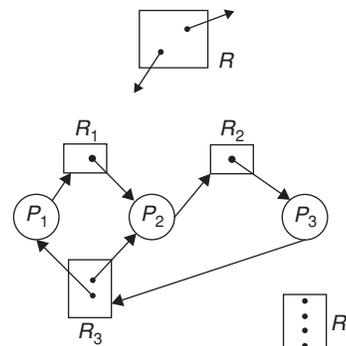


Figure 1 RAG with deadlock.

**Notes:**

1. If a cycle exists in the RAG, there may or may not be a deadlock.
2. Acyclic RAG implies no deadlock.
3. No deadlock implies acyclic RAG. This means that cycles can be there even if there is no deadlock.

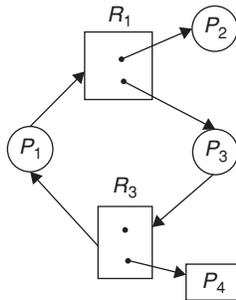


Figure 2 RAG with no deadlock but contains cycle.

**Notes:**

1. If graph contains no cycles, then no deadlock.
2. If graph contains a cycle

If there is only one instance per resource type, then deadlock occurs.

If there are several instances per resource type, deadlock may occur.

## METHODS OF HANDLING DEADLOCKS

There are three approaches to deal with deadlocks. They are

1. Deadlock prevention
2. Deadlock avoidance
3. Deadlock detection

### Deadlock Prevention

1. The strategy of deadlock prevention is to design a system in such a way that the possibility of deadlock is excluded.
2. Two classes of deadlock prevention are
  - Indirect method
  - Direct method

**Indirect method** Prevent the occurrence of one of three necessary conditions of deadlock i.e., mutual exclusion, No pre-emption and hold and wait.

**Direct method** Prevent the occurrence of circular wait.

### Prevention Techniques

**Mutual exclusion** This is supported by the OS.

#### Hold and wait

1. This condition can be prevented by requiring that a process request all of its required resources at one time and blocking the process until all requests can be granted simultaneously.
2. But this prevention does not yield good results because

- Long waiting time required
- Not efficient use of allocated resources
- A process may not know all the required resources in advance.

### Advantages

1. Works well for processes that perform a single burst of activity.
2. No pre-emption necessary.

**No pre-emption** Prevention strategies for ‘no pre-emption’ are

1. If a process that is holding some resources, requests another resource that cannot be immediately allocated to it, then all resources currently being held are released and if necessary request them again together with the additional resources.
2. If a process requests a resource that is currently held by another process, the OS may pre-empt the second process and require it to release its resources. This technique works only when two processes do not have same priority.

### Advantages

Convenient when applied to resources whose state can be saved and restored easily.

### Disadvantage

Pre-empts more often than necessary.

**Circular wait** One way to ensure that this condition never holds is to impose a total ordering of all resource types and to require that each process requests resource in an increasing order of enumeration, i.e., if a process has been allocated resources of type  $R$ , then it may subsequently request only those resources of types following  $R$  in the ordering.

### Advantages

1. Feasible to enforce via compile time checks.
2. No run-time computation required.

### Disadvantages

1. Disallows incremental resources requests.

**Note:** The deadlock prevention strategies are conservative and undercommits resources.

### Deadlock Avoidance

1. This approach allows the three necessary conditions of deadlock but makes judicious choices to assure that deadlock point is never reached.
2. Deadlock avoidance allows more concurrency than prevention.
3. A decision is made dynamically whether the current resource allocation request will, if granted, potentially lead to a deadlock.

4. It requires the knowledge of future process requests.
5. Two techniques to avoid deadlock:
  - Process Initiation Denial
  - Resource Allocation Denial

**Process initiation denial** In this technique, do not start a process if its demands might lead to deadlock.

Consider a system of 'n' processes and 'm' different types of resources. Let us define the following vectors and matrices:

Resources  $R = (R_1, R_2, \dots, R_m)$   
 $R_1$ : amount of type 1 resources  
 $R_2$ : amount of type 2 resources  
 Available =  $V = (V_1, V_2, \dots, V_m)$

'V' specifies total amount of each resource not allocated to any process.

$$\text{Claim } C = \begin{matrix} P_1 \\ P_2 \\ \vdots \\ P_n \end{matrix} \begin{bmatrix} C_{11} & C_{12} & \dots & C_{1m} \\ C_{21} & C_{22} & \dots & C_{2m} \\ \vdots & \vdots & \ddots & \vdots \\ C_{n1} & C_{n2} & \dots & C_{nm} \end{bmatrix}$$

$C_{ij}$  = requirement of process  $i$  for resources  $j$

$$\text{Allocation } A = \begin{matrix} P_1 \\ P_2 \\ \vdots \\ P_n \end{matrix} \begin{bmatrix} A_{11} & A_{12} & \dots & A_{1m} \\ A_{21} & A_{22} & \dots & A_{2m} \\ \vdots & \vdots & \ddots & \vdots \\ A_{n1} & A_{n2} & \dots & A_{nm} \end{bmatrix}$$

$A_{ij}$  = Current allocation to process  $i$  of resource  $j$ .

The following relationships must hold:

1. All resources are either available or allocated, that is,

$$R_j = V_j + \sum_{i=1}^n A_{ij}, \forall j$$

2. No process can claim more than the total amount of resources in the system, that is,

$$C_{ij} \leq R_j, \quad i, j$$

3. No process is allocated more resources of any type than the process originally claimed to need, that is,

$$A_{ij} \leq C_{ij}, \quad i, j$$

With these properties satisfied, we can define a deadlock avoidance policy that refuses to start a new process if its resource requirements might lead to deadlock. Start a new process  $P_{n+1}$  only if

$$R_j \geq C_{(n+1)j} + \sum_{i=1}^n C_{ij}, \forall j$$

that is, a process is only started if the maximum claim of all current processes plus those of the new process can be met.

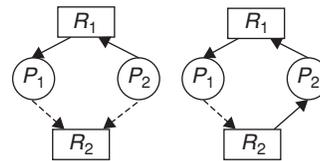
**Resource allocation denial (OR) banker's algorithm**

Consider a system with a fixed number of processes and a fixed number of resources. At any time, a process may have zero or more resources allocated to it.

**State:** The state of a system reflects the current allocation of resources to processes.

**Safe state**

1. When a process requests an available resource, the system must decide if immediate allocation leaves the system in a safe state.
2. A state is safe if the system can allocate resources to each process in some order and still avoid deadlock.
3. More formally, a system is in safe state if there exists a safe sequence of all processes.
4. A sequence of processes  $\langle P_1, P_2, \dots, P_n \rangle$  is a safe sequence for the current allocation state, if for each  $P_i$ , the resource requests that  $P_i$  can still make can be satisfied by the currently available resources plus the resources held by all the  $P_j$ , with  $j < i$ .
5. When  $P_j$  is finished,  $P_i$  can obtain the needed resources, completed its designated task, return its allocated resources and terminates.
6. When  $P_i$  terminates,  $P_{i+1}$  can obtain its needed resources and so on.
7. If a system is in safe state, no deadlock occurs.
8. If a system is in unsafe state deadlock may occur.
9. Avoidance ensures that a system will never enter an unsafe state.



10. The dotted line in the above graph represents a claim edge, i.e., a process may request that resource sometime in the future.
11. A request can only be granted if it does not result in the formation of a cycle in the graph.
12. If  $P_2$  request  $R_2$ , we cannot allocate it, since this would create a cycle.
13. A cycle indicates the system is in *unsafe state*.

**Example 1:** Consider the following state of a system consisting of three processes and two resources:

$$R = (R_1, R_2) = (5, 3)$$

$$V = (V_1, V_2) = (2, 1)$$

$$C = \begin{matrix} & R_1 & R_2 \\ P_1 & 4 & 3 \\ P_2 & 2 & 1 \\ P_3 & 3 & 3 \end{matrix}$$

$$A = \begin{matrix} & R_1 & R_2 \\ P_1 & 2 & 1 \\ P_2 & 1 & 1 \\ P_3 & 0 & 0 \end{matrix}$$

Is this is a safe state?

**Solution:**

To check whether the state is safe or not, identify whether any one of the three process can run to completion with the resources available, that is,  $C_{ij} - A_{ij} \leq V_j, j$

$$C - A = \begin{matrix} & R_1 & R_2 \\ P_1 & (2 & 2) \\ P_2 & (1 & 0) \\ P_3 & (3 & 3) \end{matrix}$$

We can identify that

$[1 \ 0] < [2 \ 1]$ .

$\therefore P_2$  can execute first.

After  $P_2$  execution it will release all its resources then

$V = (2 \ 1) + (1 \ 1) = (3 \ 2)$ .

Now  $P_1$  can execute as  $(2 \ 2) < (3 \ 2)$ .

After that  $P_1$  can release its resources then

$V = (3 \ 2) + (2 \ 1) = (5 \ 3)$

Now  $P_3$  can execute and release the resources after completion.

Hence, the safe sequence is  $\langle P_2, P_1, P_3 \rangle$ .

**Example 2:** Now suppose for the above system the allocation matrix

$$A = \begin{pmatrix} 2 & 1 \\ 1 & 1 \\ 2 & 1 \end{pmatrix} \text{ and } V = (00)$$

Then no process can run to completion as no  $C_{ij} - A_{ij} \leq V_j, j$   
Hence, the system is in unsafe state.

## Detection Algorithm for Several Instances of a Resource Type

### Safety algorithm

To find out whether or not a system is in a safe state.

**Step I:** Let 'work' and 'finish' be the two vectors of length  $m$  and  $n$ .

Initialize: Work = Available and Finish  $[i] = \text{false}$ ;

**Step II:** Find  $i$  such that both:

(a) Finish  $[i] = \text{false}$ ;

(b) Need  $\leq$  Work (Need = claim - Allocation)

If no such  $i$  exists, go to step 4.

**Step III:** Work = Work + Allocation

Finish  $[i] = \text{true}$

Go to step 2

**Step IV:** If Finish  $[i] = \text{true}$  for all  $i$ , then the system is in safe state else it is in unsafe state. This algorithm takes  $O(m \cdot n^2)$  operations to decide whether a state is safe.

### Resource-Request algorithm

Let Request $_i$  be the request vector for process,  $P_i$ . If Request $_i[j] = k$ , the process  $P_i$  wants  $k$  instances of resource

type  $R_j$ . When a request for resources is made by process,  $P_i$ , the following actions are taken.

**Step I:** If Request $_i \leq$  Need $_i$ , go to step 2. Else raise an exception (error) as the process has exceeded its maximum claim.

**Step II:** If Request $_i \leq$  Available, go to step 3. Else  $P_i$  must wait, since the resources are not available.

**Step III:** Have the system pretend to have allocated the requested resources to process,  $P_i$ , by modifying the state as follows.

Available = Available - Request $_i$

Allocation = Allocation + Request $_i$

Need $_i$  = Need $_i$  - Request $_i$

If the resulting resource-allocation is safe, then the transaction is completed and process,  $P_i$  is allocated its resources. However, if the new state is unsafe then  $P_i$  must wait for Request $_i$  and the old resource allocation state is restored.

### Advantages of deadlock avoidance technique

1. Not necessary to pre-empt and rollback processes.
2. Less restrictive than deadlock prevention.

### Disadvantages

1. Future resource requirement must be known in advance.
2. Processes can be blocked for long periods.
3. Exists fixed number of resources for allocation.

### Deadlock Detection

1. This technique does not limit resource access or restrict process action.
2. Requested resources are granted to processes whenever possible.

Deadlock detection is used by employing an algorithm that tracks the circular waiting and killing one or more processes so that the deadlock is removed.

The system state is examined periodically to determine if a set of processes is deadlocked.

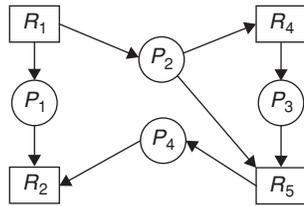
A deadlock is resolved by aborting and restarting a process, relinquishing all the resources that the process held.

**For single instance of each resource type** If in the RAG, every resource has only one instance (or single instance) then we define a deadlock detection algorithm that uses a variant of the RAG and is called a wait-for-graph.

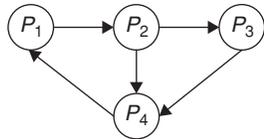
How can we get this graph from RAG? : We can get this by removing the nodes of type resource and collapsing the appropriate edges. Wait-for-graph has a cycle, then there is deadlock in the system.

To detect deadlocks, the system needs to maintain the wait-for-graph and to periodically invoke an algorithm. The complexity of this algorithm is  $O(n^2)$  where  $n$  is the number of vertices in the graph.

Consider the RAG:



We draw the wait-for-graph by removing all nodes that represents resources and collapsing their edges.



Wait-for-graph

The system is in deadlock state.

Cycle  $P_1, P_2, P_4, P_1$   
 Cycle  $P_1, P_2, P_3, P_4, P_1$

### Deadlock Detection Algorithm for Several Instances of Resource Type

Consider the Allocation matrix, A, Request matrix  $Q$  ( $Q_{ij}$  represents the amount of resource of type  $j$  requested by process  $i$ ), Resource vector  $R$  and available vector  $V$ .

The algorithm proceeds by marking processes that are not deadlocked. Initially all processes are unmarked.

1. Mark each process that has a row in the Allocation matrix of all zeros.
2. Initialize a temporary vector  $W$  to equal the Available vector.
3. Find an index  $i$  such that process  $i$  is currently unmarked and the  $i$ th row of  $Q$  is less than or equal to  $W$ , that is,  $Q_{ik} \leq W_k$ , for  $1 \leq k \leq M$ . If no such row is found, terminate the algorithm.
4. If such a row is found, mark process  $i$  and add the corresponding row of the allocation matrix to  $W$ . That is, set  $W_k = W_k + A_{ik}$  for  $1 \leq k \leq M$ . Return to step 3.

A deadlock exists if and only if there are unmarked processes at the end of the algorithm. Each unmarked process is deadlocked.

**Example 3:** Let the

|                       | $R_1$ | $R_2$ | $R_3$ |   |
|-----------------------|-------|-------|-------|---|
| Request matrix Q =    | $P_1$ | 1     | 1     | 1 |
|                       | $P_2$ | 1     | 0     | 0 |
|                       | $P_3$ | 1     | 1     | 1 |
|                       | $P_4$ | 1     | 1     | 1 |
| Allocation matrix A = | $P_1$ | 1     | 1     | 0 |
|                       | $P_2$ | 1     | 0     | 0 |
|                       | $P_3$ | 0     | 1     | 1 |
|                       | $P_4$ | 0     | 1     | 1 |

Resource vector  $R = (3 \ 3 \ 3)$   
 Available Vector  $V = (1 \ 0 \ 1)$   
 Is deadlock existing in this system?

**Solution:**

$W = (1 \ 0 \ 1)$

The request of  $P_2$  is less than  $W$ . So  $W = W + (1 \ 0 \ 0) = (2 \ 0 \ 1)$   
 So mark  $P_2$ . No other unmarked process has a row  $Q$  that is less than or equal to  $W$ .

Terminate the algorithm.

$\therefore P_1, P_3, P_4$  are in deadlock.

### Advantages

1. Never delays process initiation
2. Facilitates online handling

### Disadvantages

1. Inherent pre-emption losses

### Deadlock Recovery

The possible deadlock recovery strategies are as follows:

1. Abort all deadlocked processes.
2. Back up each deadlocked process to some previously defined checkpoint and restart all processes.
3. Successively abort deadlocked processes until deadlock no longer exists.
4. Successively pre-empt resources until deadlock no longer exists.

### DINING PHILOSOPHERS PROBLEM

Consider the following solution for dining philosophers problem using semaphores:

```

Semaphore fork[5] = {1};
int i;
void philosopher(int i)
{
 while (true)
 {
 think();
 wait(fork[i]);
 wait(fork[(i + 1) mod 5]);
 eat();
 signal(fork[(i + 1) mod 5]);
 signal(fork[i]);
 }
}
void main()
{
 Begin (Philosopher(0), Philosopher(1),
 Philosopher(2), Philosopher(3),
 Philosopher(4));
}

```

Here, each philosopher picks up first the fork on the left and then the fork on the right. After the philosopher is finished eating, the two forks are placed on the table. But this

solution leads to deadlock, if all of the philosophers are hungry at the same time, they all sit down, they all pick up the fork on their left and they all reach out for the other fork, which is not there.

A refined solution to dining philosophers problem which is deadlock free is shown below:

```
Semaphore fork[5] = {1};
Semaphore room = {4};

int i;
void philosopher(int i)
{
while(true)
{
think();
wait(room);
wait(fork[i]);
wait(fork[(i + 1) mod 5]);
eat();
signal(fork[(i + 1) mod 5]);
signal(fork[i]);
signal(room);
}
}
void main ()
{
Begin
{
(Philosopher (0), Philosopher (1), Philosopher
(2), Philosopher (3), Philosopher (4));
}
}
```

This solution is free from deadlock and starvation.

## CPU SCHEDULING

1. The objective of multi programmed OS is to maximize CPU utilization by having some process running at all times.
2. The objective of time shared OS is to switch the CPU among processes so frequently that the users can interact with each program while it is executing.
3. When there are more than one process ready to execute with the processor, a selection decision needs to be made to pick a process for execution from among the ready processes. This activity is called *process scheduling*.

**Scheduling queue:** It maintains information of all ready processes for CPU devices. It is maintained as a linked list.

### Types of Scheduling Queue

1. **Job queue:** It consists of all processes in the system.
2. **Ready queue:** It consists of all processes that are residing in the main memory and are ready but waiting to execute on CPU.
3. **Device queue:** It consists of processes waiting for a particular I/O device. Each device has its own queue.

## Process CPU-I/O Burst Cycle

The execution of process consists of CPU burst and I/O burst. The execution of process starts with CPU burst and I/O burst, which are executed alternatively.

The alternating sequence of CPU and I/O burst are shown below:

```
Read a }
Inc a } CPU Burst
Read x }
```

```
I/O waiting } I/O Burst
```

```
Dec x }
Store x } CPU Burst
```

```
I/O waiting } I/O Burst
```

```
.
.
.
```

There should be proper balance between CPU bound process and I/O bound process in a schedule.

## Scheduler

A process migrates between various scheduling queues throughout its lifetime. The process of selecting processes from the queues is carried out by scheduler.

## TYPES OF PROCESSOR SCHEDULING

There are three types of processor scheduling:

1. Long-term scheduling
2. Medium-term scheduling
3. Short-term scheduling

The following figure relates the scheduling functions to the process state transition diagram:

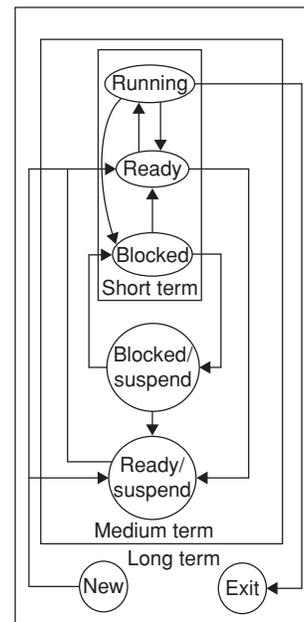


Figure 3 Levels of Scheduling

## Long-term Scheduling

1. This is performed when a new process is created. This is also called *job scheduling*.
2. This is a decision whether to add a new process to the set of process that is currently active.
3. It controls the degree of multiprogramming.
4. The long-term scheduler creates processes from the queue when it can.
5. This involves two decisions:
  - The Scheduler must decide when the OS can take on one or more additional processes.
  - The scheduler must decide which job(s) to accept and turn into processes.

## Medium-term Scheduling

1. It is a part of swapping function.
2. This is a decision whether to add a process to those that are at least partially in main memory and therefore avail for execution.

## Short-term Scheduling

1. It is the decision regarding which ready process to be executed next.
2. This is also known as CPU scheduler.
3. This is invoked whenever an event occurs that may lead to the blocking of the current process or that may provide an opportunity to pre-empt a currently running process in favour of another.
4. Another term involved in short-term scheduling is *dispatcher* which is a module that gives control of the CPU to the process selected by short-term scheduler.

**Examples:** Clock interrupts, I/O interrupts, OS calls, etc.

## SCHEDULING ALGORITHMS

### Scheduling Criteria

The commonly used scheduling criteria can be categorized along two dimensions:

1. User oriented versus system oriented.
2. Performance related versus others.

### User-oriented, performance-related criteria

1. **Turnaround time:** It is the time taken to execute a process. It is calculated as the interval from the time of submission of a process to the time of completion. Turnaround time = waiting time + execution time + time spent in I/O + time spent to get into memory
2. **Response time:** Amount of time it takes from when a request was submitted until the first response is produced. It should be minimum.
3. **Deadlines:** When process completion deadlines can be specified, the scheduling discipline should subordinate other goals to that of maximizing the percentage of deadlines met.

**Waiting time:** It is the amount of time that a process spends in ready queue and doing I/O, and it should be minimum.

### User-oriented, other criteria

**Predictability:** A given job should run in about the same amount of time and at about the same cost regardless of the load on the system.

### System-oriented, performance-related criteria

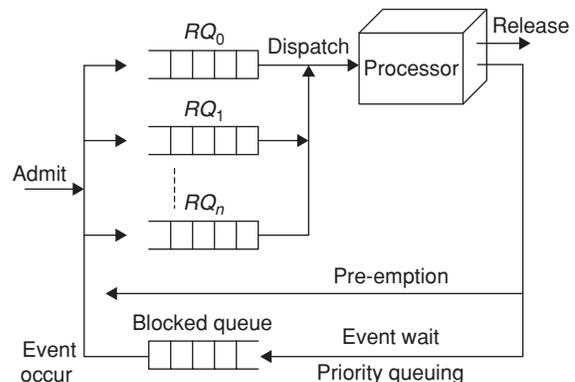
1. **Throughput:** The scheduling policy should attempt to maximize the number of processes completed per unit of time. This is a measure of how much work is being performed.
2. **Processor (CPU) utilization:** This is the percentage of time that the processor is busy. In real system, it should range from 40–90%.

### System-oriented, other criteria

1. **Fairness:** No process should suffer from starvation.
2. **Enforcing priorities:** Should favour higher priority processes and use Aging technique in order to increase the priority of processes is that wait in the system for long time.
3. **Balancing resources:** Should keep the resources of the system busy.

### Use of priorities

1. Each process is assigned a priority, and the scheduler will always choose a process of higher priority over one of lower priority.



Here  $RQ_0, RQ_1, \dots, RQ_n$  are ready queues with priority ( $RQ_i > \text{priority}(RQ_j)$  for  $i < j$ ). The scheduler starts with  $RQ_0$  processes, if it is empty choose a process from  $RQ_1$  and so on.

### Scheduling Policies

We will discuss the following scheduling algorithms:

1. FCFS
2. Round Robin
3. SPN
4. SRN
5. HRRN
6. Feedback

Before discussing these algorithms, let us discuss some basic concepts:

**Selection function**

1. Determines which process among ready process, is selected next for execution.
2. This may depend on priority, resource requirement, execution characteristics of process.
3. Execution characteristics of a process include  
 $w$  = waiting time  
 $e$  = execution time,  
 $S = w + e$

**Decision mode** It is of the following two types:

1. Non-pre-emptive
2. Pre-emptive
  - In non-pre-emptive scheduling, if once the CPU has been allocated to a process, the process can keep the CPU until it releases it, either by terminating or switching to waiting state.
  - In pre-emptive scheduling, CPU can be taken away from a process during execution.

**Comparison of non-pre-emptive and pre-emptive scheduling**

|    | Non-pre-emptive Scheduling                                                                                                   | Pre-emptive Scheduling                                                                                                                                        |
|----|------------------------------------------------------------------------------------------------------------------------------|---------------------------------------------------------------------------------------------------------------------------------------------------------------|
| 1. | In non-pre-emptive scheduling, if once a process has been allocated CPU then the CPU cannot be taken away from that process. | In pre-emptive scheduling, the CPU can be taken away before the completion of the process.                                                                    |
| 2. | No preference is given when a higher priority job comes.                                                                     | It is useful when a higher priority job comes as here the CPU can be snatched from a lower priority process.                                                  |
| 3. | The treatment of all processes is fairer.                                                                                    | The treatment of all processes is not fairer as CPU snatching is done either due to constraints or due to higher priority, process request for its execution. |
| 4. | It is a cheaper scheduling method. First come first served is an example.                                                    | It is a costlier scheduling method. Round Robin is an example.                                                                                                |

**First come, first served scheduling (FCFS)** The process that requests the CPU first is allocated the CPU first. It is non-pre-emptive scheduling and average waiting time is quite long.

**Example 4:** Consider the following processes:

| Process | Cpu Burst Time (Millisecond) |
|---------|------------------------------|
| $P_1$   | 20                           |
| $P_2$   | 5                            |
| $P_3$   | 3                            |

Find the average waiting time.

**Solution:**

Suppose they are in the order  $P_1, P_2, P_3$  at time 0. So, Gantt chart is

|       |       |       |
|-------|-------|-------|
| $P_1$ | $P_2$ | $P_3$ |
| 0     | 20    | 25 28 |

$$\text{Average waiting time} = \frac{(0+20+25)}{3} = 15 \text{ ms}$$

If they arrived in order  $P_3, P_2, P_1$  then, Gantt chart is

|       |       |       |
|-------|-------|-------|
| $P_3$ | $P_2$ | $P_1$ |
| 0     | 3     | 8 28  |

$$\text{Average waiting time} = \frac{0+3+8}{3} = \frac{11}{3} = 3.67 \text{ ms}$$

**Notes:**

1. Throughput is not that much emphasized.
2. Response time may be high especially if there is a large variance in process execution times.
3. Minimum overhead required.
4. It penalizes short processes, also penalizes I/O bound processes.
5. There is no possibility of starvation.

**Advantages**

1. Simple and brutally fair.
2. It is suitable for batch systems.

**Disadvantages**

1. The average waiting time is not minimal.
2. Not suitable for time sharing systems like Unix.
3. **Convoy effect:** Short process behind long process results in lower CPU utilization.

**Round Robin scheduling**

1. It is designed for time sharing system.
2. Similar to FCFS with pre-emption added.
3. Each process gets a small central CPU time (a time slice) usually 10 – 100 ms.
4. After time slice has elapsed and added to the end of the ready queue.
5. The scheduler picks the first process from the ready queue, sets a timer to interrupt after one time quantum and then dispatches the process. One of the following happens.
6. The process may have a CPU burst of less than 1 time quantum. (or)
7. CPU burst time of the currently executing process is longer than one time quantum. In this case, the timer will go off, cause an interrupt, a context switch is then

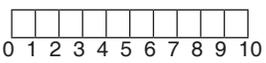
7.44 | Unit 7 • Operating System

executed and the process is put at the tail of the ready queue.

8. Average waiting time is quite long.

**Performance of Round Robin scheduling**

Let us assume that we have only one process of 10 time units.

| Process time = 10                                                                 | Quantum | Switch |
|-----------------------------------------------------------------------------------|---------|--------|
|  | 12      | 0      |
|  | 6       | 1      |
|  | 1       | 9      |

**Example 5:** Consider the following processes, arrival times and CPU processing requirements with Round Robin scheduling algorithm.

| Process | CPU time | Arrival time |
|---------|----------|--------------|
| A       | 8        | 0            |
| B       | 1        | 1            |
| C       | 2        | 3            |
| D       | 1        | 4            |
| E       | 5        | 2            |

What will be the mean turnaround time if time quantum is 4 msec?

**Solution:**

Plotting the Gantt chart

|   |   |   |   |    |    |    |
|---|---|---|---|----|----|----|
| A | B | E | C | D  | A  | E  |
| 0 | 4 | 5 | 9 | 11 | 12 | 16 |
|   |   |   |   |    |    | 17 |

Turnaround time = Finish time – Arrival time

TAT of A = 16 – 0 = 16

B = 5 – 1 = 4

C = 11 – 3 = 8

D = 12 – 4 = 8

E = 17 – 2 = 15

$$\text{Mean Turnaround time} = \frac{16+4+8+8+15}{5} = 10.2 \text{ msec}$$

- If there are  $n$  processes in the ready queue and time quantum  $q$ , then each process gets  $\frac{1}{n}$  of the CPU time in chunks of at most  $q$  time units at once.
- No process waits more than  $(n - 1)q$  time units until the next time quantum.
- The performance of Round Robin depends on time slice. If it is larger it is same as FCFS. If  $q$  is very small overhead is too high as the number of context switches increases.

**Notes:**

- Throughput is low if quantum is too small.
- Provides good response time for short processes.

- Minimum overhead.
- All processes treated fairly.
- No starvation.

**Shortest process next (SPN)**

- This is a non-pre-emptive policy in which the process with the shortest expected processing time is selected next.
- A short process will jump to the head of the queue past longer jobs.

**Example 6:** Consider the following process, such that all have arrived at time = 0

| Process | Burst time |
|---------|------------|
| $P_1$   | 5          |
| $P_2$   | 9          |
| $P_3$   | 6          |
| $P_4$   | 3          |

Find the average waiting time using SPN.

**Solution:** Gantt chart is

|       |       |       |       |
|-------|-------|-------|-------|
| $P_4$ | $P_1$ | $P_3$ | $P_2$ |
| 0     | 3     | 8     | 14    |
|       |       |       | 23    |

$$\text{Average waiting time} = \frac{0+3+8+14}{4} = 6.25 \text{ ms}$$

**Notes:**

- Difficulty with this policy is that we need to know or at least estimate the required processing time of each process.
- High throughput is possible.
- Provides good response time for short processes.
- High overhead.
- It penalizes long processes.
- There is a possibility of starvation.

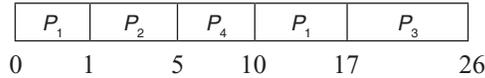
**Shortest remaining time (SRT)**

- It is a pre-emptive version of SPN.
- Here the scheduler always chooses the process that has the shortest expected remaining processing time.
- When a new process joins the ready queue, it may in fact have a shorter remaining time than the currently running process.
- Accordingly, the scheduler may pre-empt the current process when a new process becomes ready.

**Example 7:** Consider the following process. Find the average waiting time for SRT.

| Process | Arrival time | Burst time |
|---------|--------------|------------|
| $P_1$   | 0            | 8          |
| $P_2$   | 1            | 4          |
| $P_3$   | 2            | 9          |
| $P_4$   | 3            | 5          |

**Solution:** Gantt chart



So average waiting time

$$\frac{(10-1) + (1-1) + (17-2) + (5-3)}{4}$$

$$\frac{9+0+15+2}{4} = \frac{26}{4}$$

$$= 6.5 \text{ ms}$$

Here  $P_3$  starts at time 17 but the arrival time was at 2. So waiting time of  $P_3$  will be  $(17-2)$ .

**Notes:**

1. High throughput.
2. Provides good response time.
3. High overhead.
4. Penalizes long processes.
5. Starvation is possible.

### Highest Response Ratio Next (HRRN)

This algorithm works on the principle the executes the job first which has the highest response ratio. We define response ratio as the ratio between turnaround time and response time.

$$\text{Response Ratio} = \frac{(W+S)}{S}, \text{ where}$$

$W$  – Time spend waiting for the processor

$S$  – Service time

This response ratio is also named as normalized turnaround time.

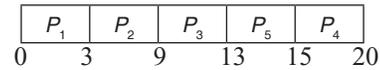
**Example 8:** Consider 5 processes with their arrival and service times:

| Process | Arrival time | Service time |
|---------|--------------|--------------|
| $P_1$   | 0            | 3            |
| $P_2$   | 2            | 6            |
| $P_3$   | 4            | 4            |
| $P_4$   | 6            | 5            |
| $P_5$   | 8            | 2            |

What is the average turnaround time using HRRN technique?

**Solution:**

Gantt chart



Average turnaround time

$$= \frac{3 + (6+1) + (4+5) + (5+9) + (2+5)}{5} = \frac{40}{5} = 8 \text{ ms}$$

**Notes:**

1. It is a non-pre-emptive scheduling algorithm.
2. High throughput.
3. Provides good response time.
4. High overhead.
5. Good balance of any type processes.
6. No starvation.

### Multilevel feedback queue scheduling

1. In multilevel queue scheduling algorithm, processes are permanently assigned to a queue on entry to the system. Processes cannot move between queues.
2. Processes can move between queues. If a process uses too much CPU time, it will be moved to a lower priority queue.
3. I/O bound and interactive processes are put into higher priority queue.
4. A process that waits too long in a lower priority queue may be moved to a higher priority queue. This form of aging prevents starvation.

In general, a multilevel feedback queue scheduler is defined by following parameters:

1. The number of queues.
2. The scheduling algorithm for each queue.
3. The method used to determine when to upgrade a process to a higher priority queue.
4. The method used to determine when to denote a process to a lower priority queue.
5. The method used to determine in which queue a process will enter when process needs service.

**Notes:**

1. Pre-emptive at time quantum.
2. Throughput is not that much emphasized.
3. Response time is not that much emphasized.
4. High overhead.
5. Favours I/O bound processes.
6. Starvation is possible.

**EXERCISES**

**Practice Problem I**

**Directions for questions 1 to 19:** Select the correct alternative from the given choices.

- Consider the following processes; find the average waiting time using non-pre-emptive priority scheduling?

| Process | Arrival time (ms) | Burst time (ms) | Priority |
|---------|-------------------|-----------------|----------|
| $P_0$   | 0                 | 10              | 5        |
| $P_1$   | 1                 | 6               | 4        |
| $P_2$   | 3                 | 2               | 2        |
| $P_3$   | 5                 | 4               | 0        |

- (A) 2.36 ms                      (B) 0.31 ms  
 (C) 7.75 ms                      (D) 13.25 ms
- Consider a set of five processes whose arrival time, CPU times needed are given below.

| Process | CPU time (in m sec) | Arrival time (in msec) |
|---------|---------------------|------------------------|
| $P_1$   | 10                  | 5                      |
| $P_2$   | 5                   | 2                      |
| $P_3$   | 3                   | 0                      |
| $P_4$   | 20                  | 4                      |
| $P_5$   | 2                   | 3                      |

If the CPU scheduling policy is SJF, find the average waiting time (with pre-emption).

- (A) 4.8 ms                      (B) 5.6 ms  
 (C) 2.16 ms                      (D) 2.8 ms
- Consider the following snapshot of a system:

|       | Allocation |   |   |   | Max |   |   |   | Available |   |   |   |
|-------|------------|---|---|---|-----|---|---|---|-----------|---|---|---|
|       | A          | B | C | D | A   | B | C | D | A         | B | C | D |
| $P_0$ | 0          | 0 | 1 | 2 | 0   | 0 | 1 | 2 | 1         | 5 | 2 | 0 |
| $P_1$ | 1          | 0 | 0 | 0 | 1   | 7 | 5 | 0 |           |   |   |   |
| $P_2$ | 1          | 3 | 5 | 4 | 2   | 3 | 5 | 6 |           |   |   |   |
| $P_3$ | 0          | 6 | 3 | 2 | 0   | 6 | 5 | 2 |           |   |   |   |
| $P_4$ | 0          | 0 | 1 | 4 | 0   | 6 | 5 | 6 |           |   |   |   |

- Which of the following is true?
- The system is in a safe state.
  - If a request from process  $P$ , arrives for (0, 4, 2, 0). Then the request can be granted.
- (A) Only (i)                      (B) Only (ii)  
 (C) Both (i) and (ii)                      (D) Neither (i) nor (ii)
- In a Round Robin scheduling context switch time is 4 units, the average process running time before blocking is 6 units then CPU efficiency is
 

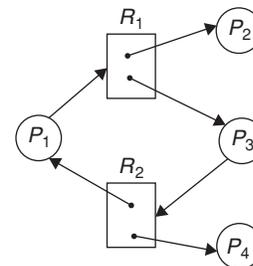
(A) 0.2                      (B) 0.4  
 (C) 0.6                      (D) 0.1
  - A short-term scheduler executes at least once every 20 msec. If it takes 2 msec to decide to execute a process

for 2 msec, what is the percentage of CPU time wasted?  
 (A) 8%                      (B) 9%  
 (C) 10%                      (D) 11%

- Consider a system which has  $n$  resources of the same type. The  $n$  resources are shared among three processes  $A, B, C$ , which have high demands of 3, 5, 6, respectively. For what value of  $n$  will deadlock not occur?  
 (A) 11                      (B) 10  
 (C) 9                      (D) 15
- A comparative study of scheduling algorithm was performed, the average arrival time in the queue is 5 m sec and waiting time of the processes is 10 msec. What is the average queue length of the waiting processes?  
 (A) 50                      (B) 60  
 (C) 70                      (D) 80
- A CPU scheduling algorithm determines an order for the execution of its scheduled processes. Given five processes to be scheduled on one processor, how many possible different schedules are there?  
 (A) 50                      (B) 100  
 (C) 120                      (D) 150
- Consider the following set of jobs (processes) along with their Arrival Time (AT), start time (ST) and Finish Time (FT). Find weighted turnaround time.

| Job no. | AT   | ST   | FT   |
|---------|------|------|------|
| 1       | 10.0 | 10.0 | 10.3 |
| 2       | 10.2 | 10.3 | 10.8 |
| 3       | 10.4 | 10.8 | 10.9 |
| 4       | 10.5 | 10.9 | 11.3 |
| 5       | 10.8 | 11.3 | 11.4 |

- (A) 3.04                      (B) 2.04  
 (C) 4.04                      (D) 0.56
- Is the following resource allocation graph in a deadlock state?



- (A) Yes                      (B) No  
 (C) Not predictable                      (D) Insufficient data
- Starvation of longer jobs happens in one of the following scheduling algorithm?  
 (A) Shortest run remaining time first  
 (B) Round Robin  
 (C) Highest response ratio next  
 (D) First-come first-served

12. Suppose  $n$  processes,  $P_1 \dots P_n$  share  $n$  identical resource units, which can be reserved and released one at a time. The maximum resource requirement of process  $P_i$  is  $S_i$ , where  $S_i > 0$ . Which one of the following is a sufficient condition for ensuring that deadlock does not occur?
- (A)  $+i, S_i < m$                       (B)  $+i, S_i < n$   
 (C)  $\sum_{i=1}^n S_i < (m+n)$             (D)  $\sum_{i=1}^n S_i < (m * n)$
13. A system with following processes and resources exists. Check the system for safe state and find the safe sequence of processes

|       | Allocation |   |   | Max |   |   | Available |   |   |
|-------|------------|---|---|-----|---|---|-----------|---|---|
|       | X          | Y | Z | X   | Y | Z | X         | Y | Z |
| $P_0$ | 0          | 1 | 0 | 7   | 5 | 3 | 3         | 3 | 2 |
| $P_1$ | 2          | 0 | 0 | 3   | 2 | 2 |           |   |   |
| $P_2$ | 3          | 0 | 2 | 9   | 0 | 2 |           |   |   |
| $P_3$ | 2          | 1 | 1 | 2   | 2 | 2 |           |   |   |
| $P_4$ | 0          | 0 | 2 | 4   | 3 | 3 |           |   |   |

- (A)  $\langle P_1, P_3, P_4, P_2, P_0 \rangle$   
 (B)  $\langle P_3, P_4, P_2, P_0, P_1 \rangle$   
 (C)  $\langle P_2, P_4, P_0, P_1, P_3 \rangle$   
 (D) The system is in unsafe state.
14. Does the below statements be executed concurrently?  
 $S_1: a = x + y$   
 $S_2: b = z + 1$
- (A) Yes                                      (B) No  
 (C) Not predictable                      (D) None of the above
15. Let  $A, B, C$  be three jobs. Their arrival time and execution time are shown below. By applying monoprogramming and multiprogramming (use Round Robin with time slice 1 unit) approaches, calculate the amount of reduction in turnaround time?

| Job | Arrival time | Execution time |
|-----|--------------|----------------|
| A   | 1            | 2              |
| B   | 2            | 6              |
| C   | 3            | 1              |

- (A) 3.33                                      (B) 4.33  
 (C) 5.33                                      (D) 2.33
16. Consider a system with three processes  $A, B, C$  with 15 tape drivers. Process A has 4 tape drives but requires 14 tape drives.

Process B has 5 tape drives but requires 9 tape drives.  
 Process C has 3 tape drives but requires 7 tape drives.  
 Among the following processes which will enter the deadlock state?

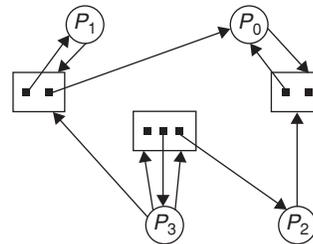
- (A) A, B only                              (B) A, B, C  
 (C) A, C only                              (D) B, C only
17. Assume that the following jobs are to be executed on a uniprocessor system:

| Job id | CPU burst time |
|--------|----------------|
| P      | 4              |
| Q      | 1              |
| R      | 8              |
| S      | 1              |
| T      | 2              |

The jobs are assumed to have arrived at 0, and in the order P, Q, R, S and T. Calculate the departure with time slice (completion time) for job P if scheduling is Round Robin with time slices of 1 unit (slice).

(A) 4                                          (B) 10  
 (C) 11                                        (D) 12

**Common data for questions 18 and 19:** Consider the following Resource Allocation Graph:



18. The system is in a deadlock state. This remark is:  
 (A) True                                      (B) False  
 (C) Impossible to determine  
 (D) Unpredictable
19. Which one is a safe sequence?  
 (A)  $P_0, P_1, P_2, P_3$                       (B)  $P_1, P_0, P_2, P_3$   
 (C)  $P_2, P_0, P_1, P_3$                       (D) Both (A) and (C)

**Practice Problem 2**

**Directions for questions 1 to 15:** Select the correct alternative from the given choices.

1. Let there be five processes ( $P_1$  to  $P_5$ ) and three resource types A, B, C.  
 Resource type A has 10 instances,  
 Resource type B has 5 instances,  
 Resource type C has 7 instances.

Suppose that at time  $T_0$ , the following snapshot of the system has been taken.

|       | Allocation | Max   | Available |
|-------|------------|-------|-----------|
| $P_1$ | 0 1 0      | 7 5 3 | 3 3 2     |
| $P_2$ | 2 0 0      | 3 2 2 |           |
| $P_3$ | 3 0 2      | 9 0 2 |           |
| $P_4$ | 2 1 1      | 2 8 2 |           |
| $P_5$ | 0 0 2      | 4 3 3 |           |

Which of the following statement is true?

- (A) The system is in safe state.  
 (B) The system has process initiation denial problem.  
 (C) No process causes initiation denial  
 (D) Both (A) and (C)
2. Consider three CPU intensive processes, which require 20, 30 and 40 time units and arrive at times 0, 2 and 4, respectively. How many context switches are needed if the operating system implements a shortest remaining time first scheduling algorithm? Do not count the context switches at time zero and at the end.  
 (A) 0 (B) 1  
 (C) 2 (D) 3
3. Consider the set of processes  $P_1$  to  $P_5$  with the following CPU burst times. Find the average turnaround time using shortest remaining time first.

| Process | CPU burst time | Arrival time |
|---------|----------------|--------------|
| $P_1$   | 3              | 0            |
| $P_2$   | 6              | 2            |
| $P_3$   | 4              | 4            |
| $P_4$   | 5              | 6            |
| $P_5$   | 2              | 8            |

- (A) 1.3 ms (B) 3.5 ms  
 (C) 5.8 ms (D) 7.2 ms
4. All processes are arriving at time 0, find the average waiting time.

| Process | Burst time | Priority |
|---------|------------|----------|
| $P_1$   | 10         | 3        |
| $P_2$   | 1          | 1        |
| $P_3$   | 2          | 3        |
| $P_4$   | 1          | 4        |
| $P_5$   | 5          | 2        |

- (A) 8.2 ms (B) 4.1 ms  
 (C) 2.0 ms (D) 1.3 ms
5. Consider a set of three processes  $P_1$ ,  $P_2$  and  $P_3$  with their priorities and arrival times as given below.

| Process | Burst time | Priority   | Arrival time |
|---------|------------|------------|--------------|
| $P_1$   | 10         | 3          | 0            |
| $P_2$   | 5          | 2          | 1            |
| $P_3$   | 2          | 1(highest) | 2            |

- Find the average waiting time.  
 (A) 1 ms (B) 2 ms  
 (C) 3 ms (D) 4 ms
6. The portion of the process scheduler in an OS that dispatches processes is concerned with:  
 (A) assigning ready processes to the CPU  
 (B) activating suspended I/O bound processes  
 (C) temporarily suspending processes when the CPU load is too great.  
 (D) All the above
7. In a time-sharing OS, when the time slot given to a process is completed, the process goes from the running state to the

- (A) blocked state (B) ready state  
 (C) suspended state (D) terminated state
8. On a system with  $n$  CPUs, what is the maximum number of processes that can be in the ready state?  
 (A)  $n$  processes  
 (B) No process can be in ready state  
 (C) There is no limit to the number of processes in the ready state  
 (D) None of the above
9. Consider a set of  $n$  tasks with known runtimes  $r_1, r_2, \dots, r_n$  to be run on a uniprocessor machine. Which of the following processor scheduling algorithms will result in the maximum throughput?  
 (A) Round Robin  
 (B) SJF  
 (C) Highest response ratio next  
 (D) First-come first-served

10. Match the following:

|                   |                           |
|-------------------|---------------------------|
| A Critical region | I Hoare's monitor         |
| B Wait/signal     | II Mutual exclusion       |
| C Working set     | III Principle of locality |
| D Deadlock        | IV Circular wait          |

- (A) A – II, B – I, C – III, D – IV  
 (B) A – I, B – II, C – III, D – IV  
 (C) A – II, B – I, C – IV, D – III  
 (D) A – I, B – II, C – IV, D – III
11. Consider three processes  $A$ ,  $B$ ,  $C$  to be scheduled as per SRT algorithm.  $A$  is known to be scheduled first and when  $A$  has been running for 7 units of time,  $C$  has arrived.  $C$  has run for 1 unit of time when  $B$  has arrived and completed running in 2 units of time, what could be the minimum time of executions for  $A$  and  $C$ ?  
 (A) 11 and 4 (B) 11 and 3  
 (C) 12 and 3 (D) 12 and 4
12. Select the correct statements from below:  
 (i) SRT and SPN can cause starvation for larger processes.  
 (ii) FCFS can potentially block small processes in favour of much larger processes.  
 (iii) Round Robin algorithm gives fair treatment to all the processes.  
 (iv) FCFS is a pre-emptive algorithm.  
 (v) The throughput for Round Robin is high even for small time slices.  
 (A) (i), (ii), (iii) (B) (i), (ii), (iv)  
 (C) (iii), (iv), (v) (D) (i), (ii), (iii), (v)
13. Three processes share four resource units that can be reserved and released only one at a time. Each process needs a maximum of two units. Then  
 (A) there is a possibility of deadlock  
 (B) no deadlock will occur

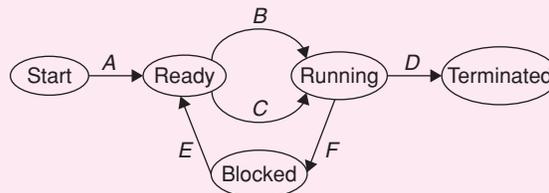
- (C) there will be a circular wait  
(D) nothing can be predicted about dead lock.
14.  $N$  processes share  $M$  resource units that can be reserved and released only one at a time. The maximum need of each process does not exceed  $M$  and the sum of all maximum needs is less than  $M + N$ , then  
(A) there is a possibility of deadlock  
(B) there will be no deadlock
- (C) circular wait exists  
(D) nothing can be predicted about dead lock.
15. Which of the following scheduling algorithms could result in starvation?  
(A) First-come, first-served  
(B) Shortest job first  
(C) Round Robin  
(D) Highest response ratio next

## PREVIOUS YEARS' QUESTIONS

1. Consider three processes (process id 0, 1, 2, respectively) with compute bursts 2, 4 and 8 time units. All processes arrive at time zero. Consider the longest remaining time first (LRTF) scheduling algorithm. In LRTF ties are broken by giving priority to the process with the lowest process id. The average turnaround time is: [2006]  
(A) 13 units (B) 14 units  
(C) 15 units (D) 16 units
2. Consider three processes, all arriving at zero, with total execution time of 10, 20 and 30 units, respectively. Each process spends the first 20% of execution time doing I/O, the next 70% of time doing computation, and the last 10% of time doing I/O again. The operating system uses a shortest remaining compute time first scheduling algorithm and schedules a new process either when the running process gets blocked on I/O or when the running process finishes its compute burst. Assume that all I/O operations can be overlapped as much as possible. For what percentage of time does the CPU remain idle? [2006]  
(A) 0% (B) 10.6%  
(C) 30.0% (D) 89.4%
3. A single processor system has three resource types  $X$ ,  $Y$  and  $Z$ , which are shared by three processes. There are five units of each resource type. Consider the following scenario, where the column *alloc* denotes the number of units of each resource type allocated to each process, and the column *request* denotes the number of units of each resource type requested by a process in order to complete execution. Which of these processes will finish *last*? [2007]

|       | alloc |     |     | request |     |     |
|-------|-------|-----|-----|---------|-----|-----|
|       | $X$   | $Y$ | $Z$ | $X$     | $Y$ | $Z$ |
| $P_0$ | 1     | 2   | 1   | 1       | 0   | 3   |
| $P_1$ | 2     | 0   | 1   | 0       | 1   | 2   |
| $P_2$ | 2     | 2   | 1   | 1       | 2   | 0   |

- (A)  $P_0$   
(B)  $P_1$   
(C)  $P_2$   
(D) None of the above, since the system is in a deadlock.
4. Which of the following is *not* true of deadlock prevention and deadlock avoidance schemes? [2008]  
(A) In deadlock prevention, the request for resources is always granted if the resulting state is safe  
(B) In deadlock avoidance, the request for resources is always granted if the result state is safe  
(C) Deadlock avoidance is less restrictive than deadlock prevention  
(D) Deadlock avoidance requires knowledge of resource requirements *a priori*
5. In the following process state transition diagram for a uniprocessor system, assume that there are always some processes in the ready state: [2009]



Now consider the following statements:

- I. If a process makes a transition  $D$ , it would result in another process making transition  $A$  immediately.  
II. A process  $P_2$  in blocked state can make transition  $E$ , while another process  $P_1$  is in running state.  
III. The OS uses pre-emptive scheduling.  
IV. The OS uses non-pre-emptive scheduling.
- Which of the above statements are true?  
(A) I and II (B) I and III  
(C) II and III (D) II and IV
6. Which of the following statements are true?  
I. Shortest remaining time first scheduling may cause starvation  
II. Pre-emptive scheduling may cause starvation  
III. Round Robin is better than FCFS in terms of response time [2010]  
(A) I only (B) I and III only  
(C) II and III only (D) I, II and III
7. A system has  $n$  resources  $R_0, \dots, R_{n-1}$ , and  $k$  processes  $P_0, \dots, P_{k-1}$ . The implementation of the resource request logic of each process  $P_i$  is as follows:

```

if (i%2 == 0) {
if (i<n) request Ri;
if (i+2<n) request Ri+2;
}
else {
if (i<n) request Rn-i;
if (i+2<n) request Rn-i-2;
}

```

In which one of the following situations is a deadlock possible? [2010]

- (A)  $n = 40, k = 26$  (B)  $n = 21, k = 12$   
 (C)  $n = 20, k = 10$  (D)  $n = 41, k = 19$

8. Consider the following table of arrival time and burst time for three processes  $P_0, P_1$  and  $P_2$ . [2011]

| Process | Arrival time | Burst time |
|---------|--------------|------------|
| $P_0$   | 0 ms         | 9 ms       |
| $P_1$   | 1 ms         | 4 ms       |
| $P_2$   | 2 ms         | 9 ms       |

The pre-emptive shortest job first scheduling algorithm is used. Scheduling is carried out only at arrival or completion of processes. What is the average waiting time for the three processes?

- (A) 5.0 ms (B) 4.33 ms  
 (C) 6.33 ms (D) 7.33 ms

9. Consider the three processes,  $P_1, P_2$  and  $P_3$  as shown in the table. [2012]

| Process | Arrival time | Time units required |
|---------|--------------|---------------------|
| $P_1$   | 0            | 5                   |
| $P_2$   | 1            | 7                   |
| $P_3$   | 3            | 4                   |

The completion order of the three processes under the policies FCFS and  $RR_2$  (Round Robin scheduling with CPU quantum of 2 time units) are

- (A) FCFS:  $P_1, P_2, P_3$  RR2:  $P_1, P_2, P_3$   
 (B) FCFS:  $P_1, P_3, P_2$  RR2:  $P_1, P_3, P_2$   
 (C) FCFS:  $P_1, P_2, P_3$  RR2:  $P_1, P_3, P_2$   
 (D) FCFS:  $P_1, P_3, P_2$  RR2:  $P_1, P_2, P_3$

10. A scheduling algorithm assigns priority proportional to the waiting time of a process. Every process starts with priority zero (the lowest priority). The scheduler re-evaluates the process priorities every  $T$  time units and decides the next process to schedule. Which one of the following is TRUE if the processes have no I/O operations and all arrive at time zero? [2013]

- (A) This algorithm is equivalent to the first-come-first-serve algorithm.  
 (B) This algorithm is equivalent to the Round Robin algorithm.

(C) This algorithm is equivalent to the shortest-job-first algorithm.

(D) This algorithm is equivalent to the shortest-remaining-time-first algorithm.

11. An operating system uses the *Banker's algorithm* for deadlock avoidance when managing the allocation of three resource types  $X, Y$  and  $Z$  to three processes  $P_0, P_1$ , and  $P_2$ . The table given below presents the current system state. Here, the *allocation* matrix shows the current number of resources of each type allocated to each process and the *Max* matrix shows the maximum number of resources of each type required by each process during its execution. [2014]

|       | Allocation |   |   | Max |   |   |
|-------|------------|---|---|-----|---|---|
|       | X          | Y | Z | X   | Y | Z |
| $P_0$ | 0          | 0 | 1 | 8   | 4 | 3 |
| $P_1$ | 3          | 2 | 0 | 6   | 2 | 0 |
| $P_2$ | 2          | 1 | 1 | 3   | 3 | 3 |

There are three units of type  $X$ , two units of type  $Y$  and two units of type  $Z$  still available. The system is currently in a *safe* state. Consider the following independent requests for additional resources in the current state:

REQ1:  $P_0$  requests 0 units of  $X$ , 0 units of  $Y$  and two units of  $Z$

REQ 2:  $P_1$  requests two units of  $X$ , 0 units of  $Y$  and 0 units of  $Z$

Which one of the following is true?

- (A) Only REQ1 can be permitted  
 (B) Only REQ2 can be permitted  
 (C) Both REQ1 and REQ2 can be permitted  
 (D) Neither REQ1 nor REQ2 can be permitted

12. Consider the following set of processes that need to be scheduled on a single CPU. All the times are given in milliseconds.

| Process name | Arrival time | Execution time |
|--------------|--------------|----------------|
| A            | 0            | 6              |
| B            | 3            | 2              |
| C            | 5            | 4              |
| D            | 7            | 6              |
| E            | 10           | 3              |

Using the *shortest remaining time first* scheduling algorithm, the average process turnaround time (in msec) is \_\_\_\_\_. [2014]

13. Three processes  $A, B$  and  $C$  each execute a loop of 100 iterations. In each iteration of the loop, a process performs a single computation that requires  $t_c$  CPU milliseconds and then initiates a single I/O operation that lasts for  $t_{io}$  milliseconds. It is assumed that the computer where the processes execute has sufficient number of I/O devices and the OS of the computer assigns different I/O devices to each process. Also, the scheduling overhead of the OS is negligible. The processes have the following characteristics:

| Process id | $t_c$  | $t_o$  |
|------------|--------|--------|
| A          | 100 ms | 500 ms |
| B          | 350 ms | 500 ms |
| C          | 200 ms | 500 ms |

The processes  $A$ ,  $B$  and  $C$  are started at times 0, 5 and 10 milliseconds, respectively, in a pure time sharing system (Round Robin scheduling) that uses a time slice of 50 milliseconds. The time in milliseconds at which process  $C$  would complete its first I/O operation is \_\_\_\_\_ [2014]

14. A system contains three programs and each requires three tape units for its operation. The minimum number of tape units which the system must have such that deadlocks never arise is \_\_\_\_\_. [2014]
15. An operating system uses *shortest remaining time first* scheduling algorithm for pre-emptive scheduling of processes. Consider the following set of processes with their arrival times and CPU burst times (in milliseconds):

| Process | Arrival time | Burst time |
|---------|--------------|------------|
| $P_1$   | 0            | 12         |
| $P_2$   | 2            | 4          |
| $P_3$   | 3            | 6          |
| $P_4$   | 8            | 5          |

The average waiting time (in milliseconds) of the processes is \_\_\_\_\_. [2014]

16. Consider a uniprocessor system executing three tasks  $T_1$ ,  $T_2$  and  $T_3$ , each of which is composed of an infinite sequence of jobs (or instances) which arrive periodically at intervals of 3, 7 and 20 milliseconds, respectively. The priority of each task is the inverse of its period, and the available tasks are scheduled in order of priority, with the highest priority task schedule first. Each instance of  $T_1$ ,  $T_2$  and  $T_3$  requires an execution time of 1, 2 and 4 milliseconds, respectively. Given that all tasks initially arrive at the beginning of the 1st millisecond and task preemptions are allowed, the first instance of  $T_3$  completes its execution at the end of \_\_\_\_\_ milliseconds. [2015]
17. A system has 6 identical resources and  $N$  processes competing for them. Each process can request at most 2 resources. Which one of the following values of  $N$  could lead to a deadlock? [2015]
- (A) 1 (B) 2  
(C) 3 (D) 4
18. The maximum number of processes that can be in Ready state for a computer system with  $n$  CPUs is [2015]
- (A)  $n$  (B)  $n^2$   
(C)  $2^n$  (D) Independent of  $n$
19. Consider the following policies for preventing deadlock in a system with mutually exclusive resources. [2015]

- (1) Processes should acquire all their resources at the beginning of execution. If any resource is not available, all resources acquired so far are released.
- (2) The resources are numbered uniquely, and processes are allowed to request for resources only in increasing resource numbers.
- (3) The resources are numbered uniquely, and processes are allowed to request for resources only in decreasing resource numbers.
- (4) The resources are numbered uniquely. A process is allowed to request only for a resource with resource number larger than its currently held resources.

Which of the above policies can be used for preventing deadlock?

- (A) Any one of 1 and 3 but not 2 or 4  
(B) Any one of 1, 3 and 4 but not 2  
(C) Any one of 2 and 3 but not 1 or 4  
(D) Any one of 1, 2, 3 and 4
20. For the processes listed in the following table, which of the following scheduling schemes will give the lowest average turnaround time? [2015]

| Process | Arrival Time | Processing Time |
|---------|--------------|-----------------|
| A       | 0            | 3               |
| B       | 1            | 6               |
| C       | 4            | 4               |
| D       | 6            | 2               |

- (A) First Come First Serve  
(B) Non-preemptive Shortest Job First  
(C) Shortest Remaining Time  
(D) Round Robin with Quantum value two
21. Consider an arbitrary set of CPU - bound processes with unequal CPU burst lengths submitted at the same time to a computer system. Which one of the following process scheduling algorithms would minimize the average waiting time in the ready queue? [2016]
- (A) Shortest remaining time first  
(B) Round-robin with time quantum less than the shortest CPU burst  
(C) Uniform random  
(D) Highest priority first with priority proportional to CPU burst length
22. Consider the following processes, with the arrival time and the length of the CPU burst given in milliseconds. The scheduling algorithm used is preemptive shortest remaining - time first.

| Process | Arrival Time | Burst Time |
|---------|--------------|------------|
| $P_1$   | 0            | 10         |
| $P_2$   | 3            | 6          |

|                |   |   |
|----------------|---|---|
| P <sub>3</sub> | 7 | 1 |
| P <sub>4</sub> | 8 | 3 |

The average turn around time of these processes is \_\_\_\_\_ milliseconds. **[2016]**

23. Consider the following CPU processes with arrival times (in milliseconds) and length of CPU bursts (in milliseconds) as given below :

| Process | Arrival time | Burst time |
|---------|--------------|------------|
| P1      | 0            | 7          |
| P2      | 3            | 3          |
| P3      | 5            | 5          |
| P4      | 6            | 2          |

If the pre-emptive shortest remaining time first scheduling algorithm is used to schedule the processes, then the average waiting time across all processes is milliseconds. **[2017]**

24. A system shares 9 tape drives. The current allocation and maximum requirement of tape drives for three processes are shown below:

| Process | Current Allocation | Maximum Requirement |
|---------|--------------------|---------------------|
| P1      | 3                  | 7                   |
| P2      | 1                  | 6                   |
| P3      | 3                  | 5                   |

Which of the following best describes current state of the system? **[2017]**

- (A) Safe. Deadlocked
  - (B) Safe. Not Deadlocked
  - (C) Not Safe. Deadlocked
  - (D) Not Safe, Not Deadlocked
25. Consider the set of processes with arrival time (in milliseconds), CPU burst time (in milliseconds), and priority (0 is the highest priority) shown below. None of the processes have I/O burst time.

| Process        | Arrival Time | Burst Time | Priority |
|----------------|--------------|------------|----------|
| P <sub>1</sub> | 0            | 11         | 2        |
| P <sub>2</sub> | 5            | 28         | 0        |

|                |    |    |   |
|----------------|----|----|---|
| P <sub>3</sub> | 12 | 2  | 3 |
| P <sub>4</sub> | 2  | 10 | 1 |
| P <sub>5</sub> | 9  | 16 | 4 |

The average waiting time (in milliseconds) of all the processes using preemptive priority scheduling algorithm is \_\_\_\_\_. **[2017]**

26. Consider a system with 3 processes that share 4 instances of the same resource type. Each process can request a maximum of  $K$  instances. Resource instances can be requested and released only one at a time. The largest value of  $K$  that will always avoid deadlock is \_\_\_\_\_. **[2018]**

27. In a system, there are three types of resources:  $E$ ,  $F$  and  $G$ . Four processes  $P_0, P_1, P_2$  and  $P_3$  execute concurrently. At the outset, the processes have declared their maximum resource requirements using a matrix named Max as given below. For example,  $Max[P_2, F]$  is the maximum number of instances of  $F$  that  $P_2$  would require. The number of instances of the resources allocated to the various processes at any given state is given by a matrix named Allocation.

Consider a state of the system with the Allocation matrix as shown below, and in which 3 instances of  $E$  and 3 instances of  $F$  are the only resources available.

| Allocation     |   |   |   | Max            |   |   |   |
|----------------|---|---|---|----------------|---|---|---|
|                | E | F | G |                | E | F | G |
| P <sub>0</sub> | 1 | 0 | 1 | P <sub>0</sub> | 4 | 3 | 1 |
| P <sub>1</sub> | 1 | 1 | 2 | P <sub>1</sub> | 2 | 1 | 4 |
| P <sub>2</sub> | 1 | 0 | 3 | P <sub>2</sub> | 1 | 3 | 3 |
| P <sub>3</sub> | 2 | 0 | 0 | P <sub>3</sub> | 5 | 4 | 1 |

From the perspective of deadlock avoidance, which one of the following is true? **[2018]**

- (A) The system is in safe state.
- (B) The system is not in safe state, but would be safe if one more instance of  $E$  were available.
- (C) The system is not in safe state, but would be safe if one more instance of  $F$  were available.
- (D) The system is not in safe state, but would be safe if one more instance of  $G$  were available.

**ANSWER KEYS****EXERCISES****Practice Problem 1**

1. C    2. A    3. C    4. C    5. B    6. D    7. A    8. C    9. A    10. B  
 11. A    12. C    13. A    14. A    15. B    16. B    17. C    18. B    19. D

**Practice Problem 2**

1. B    2. C    3. D    4. A    5. C    6. A    7. B    8. C    9. B    10. A  
 11. D    12. A    13. B    14. B    15. B

**Previous Years' Questions**

1. A    2. B    3. C    4. A    5. C    6. D    7. B    8. A    9. C    10. B  
 11. B    12. 7.2    13. 1000    14. 7    15. 5.5    16. 12    17. none    18. D    19. D    20. C  
 21. A    22. 8.2 to 8.3    23. 3    24. B    25. 29    26. 2    27. A

# Chapter 4

## Memory Management and Virtual Memory

### LEARNING OBJECTIVES

- Basic concepts
- Memory management requirements
- Relocation and memory mapping techniques
- Placement algorithm
- Dynamic partitioning
- Placement algorithm
- Buddy system
- Non-contiguous storage allocation methods
- Paging
- Segmentation
- Page table structure
- Hierarchical page table
- Inverted page table
- Address translation in a segmentation system

### BASIC CONCEPTS

**Uniprogramming system** Main memory is divided into two parts as follows:

1. Operating system (OS) part
2. Program part (which is currently being executed)

**Multiprogramming system** Here the user part of memory must be further subdivided to accommodate multiple processes.

The task of subdivision is carried out dynamically by the OS and is known as memory management.

### Memory Hierarchy

The triangle in Figure 1 gives the hierarchy of memory. The memory hierarchy shows the performance issues.

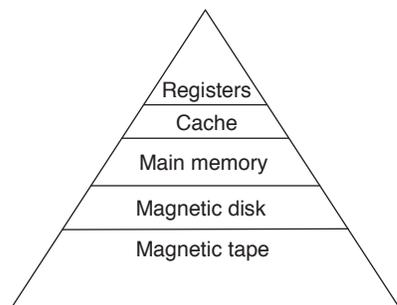


Figure 1 Memory hierarchy.

The memory hierarchy has different types of storage system in computers which are arranged in hierarchy, with respect to speed and cost.

If one moves down the hierarchy, access time increases, the cost per bit decreases, the memory capacity increases and memory access frequency by the processor decreases.

The registers, cache and main memory are volatile, whereas magnetic disc and magnetic tapes are non-volatile storage devices.

### MEMORY MANAGEMENT REQUIREMENTS

Memory management requirements are as follows:

1. Relocation
2. Protection
3. Sharing
4. Logical organization
5. Physical organization

### Relocation

1. The role of relocation, the ability to execute processes independently from their physical location in memory, is central for memory management.
2. In a general purpose multiprogramming environment, a program cannot know in advance what processes will be running in memory when it is executed, nor how much memory the system has available for it, nor where it is located.
3. Hence program relocation is required such that a program must be compiled and linked in such a way that it can later be loaded starting from an unpredictable address in memory, an address that can even change during the execution of the process itself, if any swapping occurs.

4. The basic requirement for program relocation is that all the references to memory it makes during execution must not contain absolute (physical) address of memory cells, but must be generated relatively, that is, as a distance measured in number of contiguous memory words, from some known point.

## Protection

Each process should be protected against unwanted interference by other processes, whether accidental or intentional.

Thus, programs in other processes should not be able to reference memory locations in a process for reading or writing purpose without permission.

## Sharing

1. Any protection mechanism must have the flexibility to allow several processes to access the same portion of main memory.
2. Processes that are cooperating on some task may need to share access to the same data structure.
3. The memory management system must therefore allow controlled access to shared areas of memory without compromising essential protection.

## Logical Organization

Main memory in a computer system is organized as a linear address space consisting of a sequence of bytes or words. But most of the programs are organized into modules.

If the OS and hardware can effectively deal with user programs and data in the form of modules, then there are some advantages.

1. Modules can be written and compiled independently.
2. Different degrees of protection can be given to different modules.
3. It is better to share modules among processes.

## Physical Organization

Computer memory is organized in two levels:

1. Main memory
2. Secondary memory

The flow of information between these two modules is a major concern. If this is assigned to user, then there are some problems:

1. **Overlaying** may be possible. In overlaying concept, the various modules of a program can be assigned to the same region of memory, which causes wastage of programmer time.
2. The programmer does not know at the time of coding how much space will be available or where that space will be. So it must be handled by the system.

**Address binding** Addresses may be represented in different ways during the program execution:

1. Addresses in source program are generally symbolic.
2. A compiler will typically bind these symbolic addresses to relocatable addresses.
3. The linkage editor or loader will in turn bind the relocatable addresses to absolute addresses.

So the binding of instructions and data to memory addresses can be done at any step along the way:

1. Compile time
2. Load time
3. Execution time

**Compile time** If you know at compile time where the process will reside in memory, then absolute code can be generated.

**Load time** If it is not known at compile time where the process will reside in memory, then the compiler must generate relocatable code.

**Execution time** If the process can be moved during its execution from one memory segment to another, then binding must be delayed until run time.

## Logical versus Physical Address Apave

**Logical address** An address generated by the CPU is commonly referred to as a logical address.

**Physical address** An address seen by the memory unit, that is, one loaded into MAR is referred as a *physical address*.

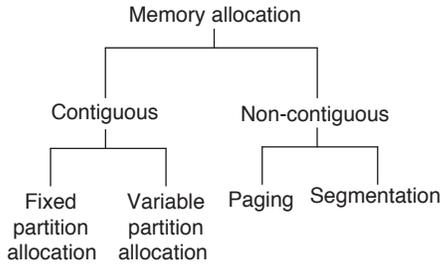
### Notes:

1. Logical and physical addresses differ in execution time address-binding scheme.
2. Logical and physical addresses are same in compile time and load time address-binding schemes.
3. The run-time mapping from logical to physical address is done by a hardware device called the *memory management unit* (MMU).

## MEMORY MAPPING TECHNIQUES

The principle operation of memory management is to bring processes into main memory for execution by the processor. Let's now discuss various memory management techniques as follows:

1. Fixed partitioning
2. Dynamic partitioning
3. Simple paging
4. Simple segmentation
5. Virtual memory paging
6. Virtual memory segmentation.



### Contiguous Storage Allocation

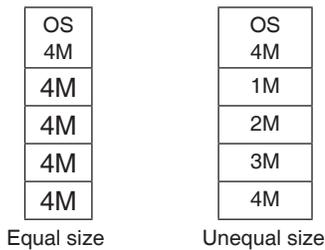
In this allocation, a memory resident program occupies a single contiguous block of memory.

### Fixed/Static Partitioning

The main memory is divided into a number of static partitions at system generation time. Moreover, a process may be loaded into a partition of equal or greater size.

**Partition size** Two alternatives of fixed partition are as follows:

1. Equal-size partitions
2. Unequal-size partitions



**Equal-size partitions:** Any process whose size is less than or equal to the partition size can be loaded into any available partition.

Two problems with this technique are as follows:

1. A program may be too big to fit into a partition. Use overlaying to solve this problem.
2. Main memory utilization is extremely inefficient, as there is a possibility of internal fragmentation.

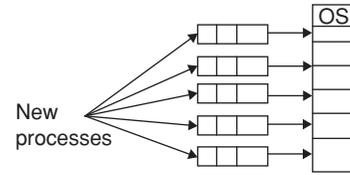
In *internal fragmentation*, there is a space wastage internal to a partition due to the fact that the block of data loaded is smaller than the partition.

**Unequal-sized partition:** Both the problems with equal-size partition can be lessened by using unequal-sized partitions.

**Placement algorithm:** With equal-size partitions, the placement of processes in memory is trivial. As all partitions are of equal size, it doesn't matter which partition is used.

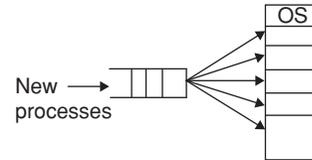
With Unequal-size partitions, there are two possible ways to assign processes to partitions:

1. Assign each process to the smallest partition within which it will fit.



- Figure shows one process queue for partition.
- Minimized internal fragmentation.
- Possibility of unused partitions.

2. Employ a single queue for all processes.



- When it is time to load a process into main memory, the smallest available partition that will hold the process is selected.

### Advantages

1. Simple to implement.
2. Little OS overhead.

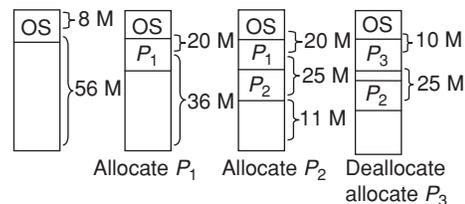
### Disadvantages

1. Inefficient use of memory due to internal fragmentation.
2. Maximum number of active processes is fixed.

### Dynamic Partitioning

With dynamic partitioning, the partitions are of variable length and number. When a process is brought into main memory, it is allocated exactly as much memory as it requires and no more.

**Example:**



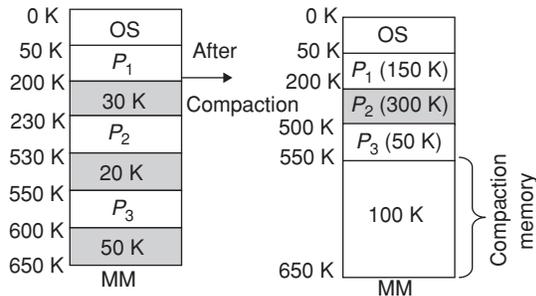
- This method starts out well, but eventually it leads to a situation in which there are a lot of small holes in memory.
- As time goes on, memory becomes more and more fragmented and memory utilization declines. This phenomenon is referred to as *external fragmentation*.

It indicates the memory that is external to all partitions becomes increasingly fragmented.

### Compaction

Compaction is a technique by which the resident program are relocated in such a way that the small chunks of free

memory are made contiguous to each other and clubbed together into a single free partition that may be big enough to accommodate more programs.



It should be noted that compaction involves dynamic relocation of a program.

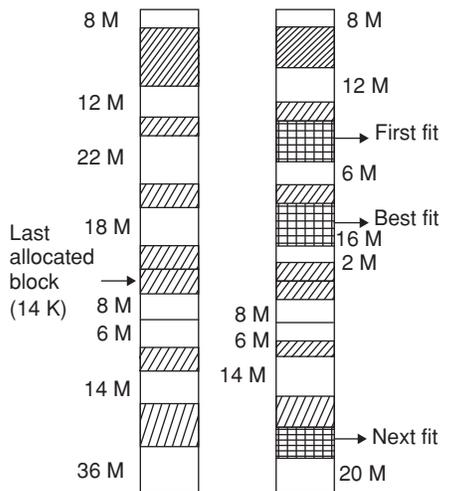
### Placement algorithm

Memory compaction is a time-consuming process, and hence the OS uses some placement algorithms.

The three most common strategies to allocate free partitions to the new processes are as follows:

1. *First fit*: Allocate the first free partition, large enough to accommodate the process. IT executes faster.
2. *Best fit*: Allocate the smallest free partition that meets the requirement of the process. It achieves higher utilization of memory by searching smallest free partition.
3. *Worst fit*: Allocate the largest available partition to the newly entered process in the system.
4. *Next fit*: Start from current location in the list.

**Example:** Consider the following memory configuration after a number of placement and swapping out operations. The last block that was used was a 22 MB block from which a 14 MB partition was created. The figure (b) shows 16 MB allocation request.



(a) Before allocation

(b) After allocation

- Allocated block
- Free block
- Possible new allocation

### Advantages of dynamic partitioning

1. Memory utilization is generally better as partitions are created dynamically.
2. No internal fragmentation as partitions are changed dynamically.
3. The process of merging adjacent holes to form a single larger hole is called coalescing.

### Disadvantages

1. Lots of OS space, time, complex memory management algorithms are required.
2. Compaction time is very high.

**Buddy system:** Both fixed and dynamic partitioning schemes have drawbacks.

In Buddy system, memory blocks are available of size  $2^k$  words,  $L \leq k \leq U$ , where,

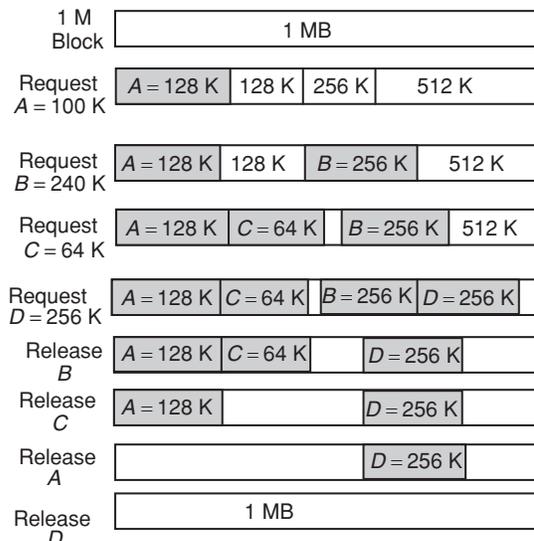
$2^L$  = Smallest-size block that is allocated.

$2^U$  = Largest-size block that is allocated.

Generally,  $2^U$  is the size of the entire memory available for allocation. If a request of size 'S' such that  $2^{U-1} < S \leq 2^U$  is made, then the entire block is allocated. Otherwise the block is split into two equal buddies of size  $2^{U-1}$ . If  $2^{U-2} < S \leq 2^{U-1}$ , then the request is allocated to one of the two buddies. Otherwise, one of the buddies is split in half again. This process continues until the smallest block greater than or equal to 'S' is generated and allocated to the request.

1. At any time, the buddy system maintains a list of holes of each size  $2^i$ .
2. A hole may be removed from the  $(i + 1)$  list by splitting it in half to create two buddies of size  $2^i$  in the 'i' list.
3. Whenever a pair of buddies on the  $i$  list both become unallocated, they are removed from the list and coalesced into a single block on the  $(i + 1)$  list.

### Example:



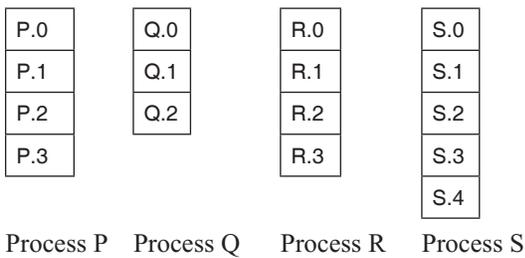
## Non-contiguous Storage Allocation Methods

### Paging

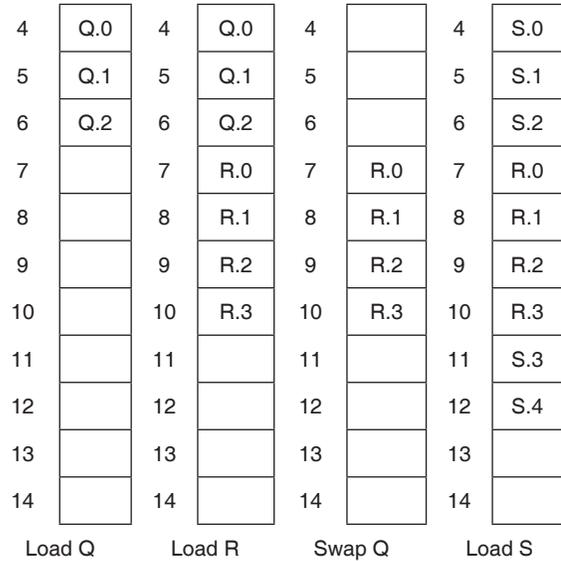
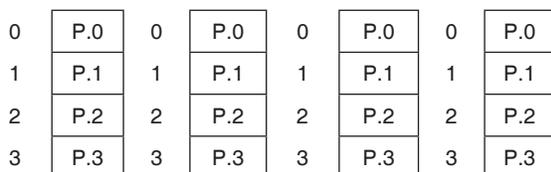
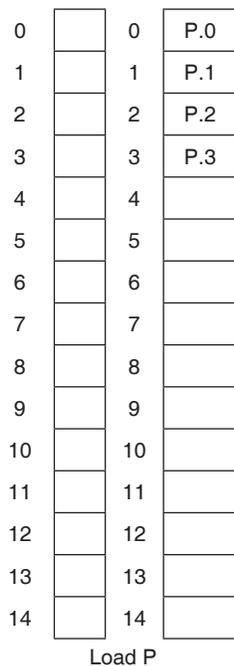
In simple paging, the main memory is divided into a number of equal-size frames. Each process is divided into a number of equal-size frames. The chunks of processes are referred as *pages*. A process is loaded by loading all of its pages into available, not necessarily contiguous frames.

**Example:** At a point in time, some of the frames in memory are in use and some are free. A list of free frames is maintained by the OS.

Consider four processes with their pages as displayed below:



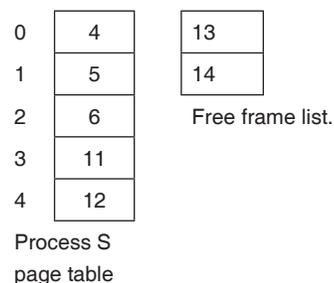
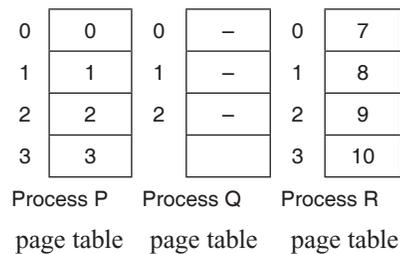
Let the main memory consist of 15 frames:  
Main memory



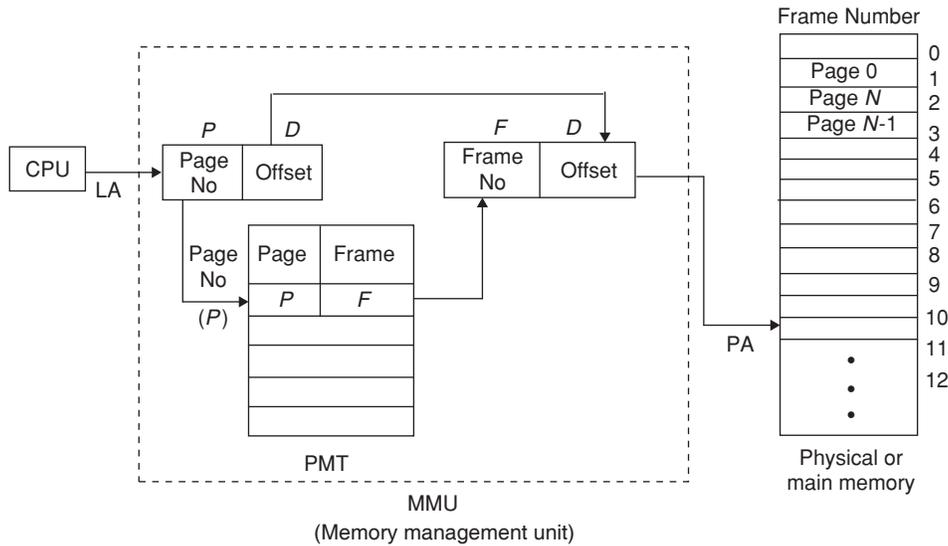
- The OS maintains a page table for each process.
- The page table shows the frame location for each page of the process.
- Within a program, each logical address consists of a page number and an offset within the page.
- Here a logical address is the location of a word relative to the beginning of the program; the processor translates that into a physical address.
- For this, the processor must know the following details:

- Logical address:** Consists page number and offset.
- Page table:** Used to produce physical address (Frame number, offset).

In the previous example, the page tables of each process will be:

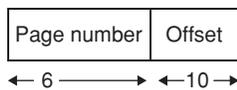


**Address mapping in paging**



**Example:** Let 16-bit address is used by the processor and the page size is 1 KB.

Then number of pages in main memory =  $\frac{2^{16}}{2^{10}} = 2^6$ .  
 $\therefore$  Page number = 6-bits  
 Offset = 10-bits



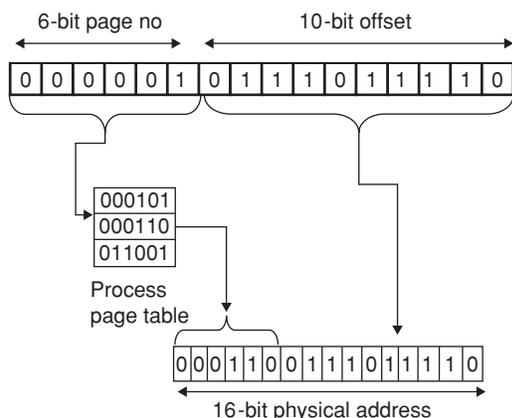
For the relative address 1502 =

0000010111011110

The page number = 000001 = 1

and offset = 0111011110 = 478, that is, the physical location will be an offset (478)<sub>10</sub> on page 1.

Let the page 1 is present in frame 6. Then the physical address will be 0001100111011110.



**Figure 2** Paging.

**Steps for address translation**

1. Extract the page number from the logical address.
2. Use the page number as an index into the process page table to find the frame number.
3. The physical address will be constructed by appending the frame number to the offset.

**Advantage**

There is no external fragmentation.

**Disadvantage**

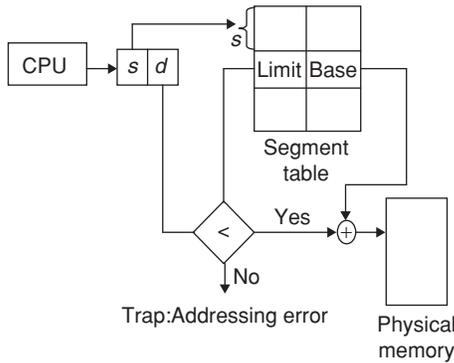
There is a small amount of internal fragmentation.

**Segmentation**

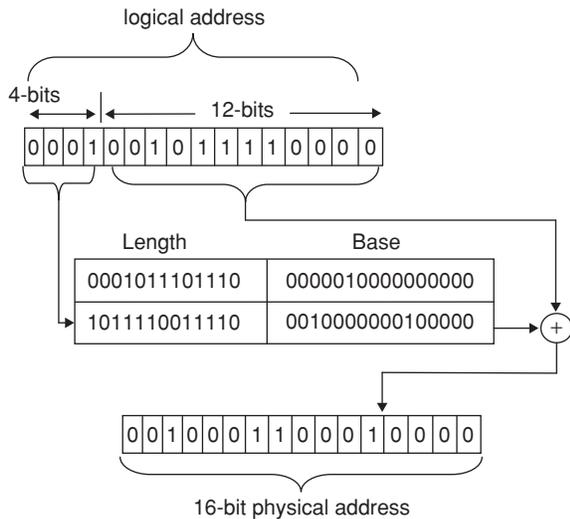
1. Each process is divided into a number of unequal-size segments.
2. A process is loaded by loading all of its segments into dynamic partitions that need not be contiguous.
3. The logical address using segmentation consists of two parts: *segment number* and an *offset*.
4. The principle inconvenience of segmentation is that the programmer must be aware of the maximum segment size limitation.
5. It makes use of a segment table for each process and a list of free blocks of main memory.
6. Each segment table entry would have to give the starting address in main memory of the corresponding segment. It also contains the length of the segment.
7. Steps for address translation:
  - Extract the segment number from the logical address.
  - Use the segment number as an index into the process segment table to find the starting physical address of the segment.

- Compare the offset to the length of the segment. If the offset is greater than or equal to the length, the address is invalid.
- The desired physical address is the sum of the starting physical address of the segment plus the offset.

**Hardware support for segmentation**



**Example:** Consider the logical address 0001001011110000.  
 Let the segment number consists of 4-bits. Then segment number = 0001 = 1  
 Offset = 001011110000 = 752.



**Advantages**

1. No internal fragmentation
2. Improved memory utilization.
3. Reduced overhead compared to dynamic partitioning.

**Disadvantage**

1. External fragmentation.

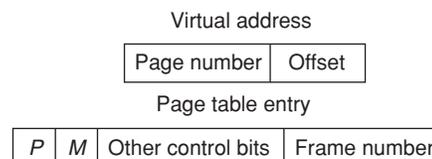
**VIRTUAL MEMORY**

1. In simple paging/segmentation, it is not necessary that all of the pages or all of the segments of a process be in main memory during execution.

2. Suppose that it is time to bring new process into memory. The OS begins by bringing in only one or a few pages to include the initial program page and initial data page to which those instructions refer.
3. The portion of a process that is actually in main memory at any time is defined to be the resident set of the process.
4. If the processor encounters a logical address that is not in main memory, it generates an interrupt indicating a memory access fault.
5. Then the OS brings the required page to the main memory.
6. With virtual memory,
  - More processes may be maintained in main memory.
  - A process may be larger than all of main memory, then also it will be executed.
7. Virtual memory is a storage allocation scheme in which secondary memory can be addressed as though it were part of main memory.
8. **Thrashing:** When the OS brings one page in, it must throw another out. If it throws out a page just before it is used, then it will just have to go get that piece again almost immediately. Too much of this leads to a condition known as *thrashing*.  
 If the system spends most of its time in swapping rather than executing instructions then that situation refers to thrashing.
9. Principle of locality suggests that a virtual memory scheme may work.
10. Virtual memory will be practical and effective if
  - There is a hardware support for paging/segmentation.
  - The OS includes software for managing the movement of pages/segments.

**Paging with Virtual Memory**

1. The main difference between paging and virtual memory paging is that in virtual memory paging concept, not all pages of a process need to be in main memory frames for the process to run. Pages may be read in as needed.
2. A page table is also needed for a virtual memory scheme based on paging. Also it is typical to associate a unique page table with each process.
3. The virtual address and page table entries for virtual memory paging are shown below:



*P*: Present bit. This bit specifies whether that particular page is present in main memory or not.

*M*: Modified bit. This bit indicates whether the contents of the corresponding page have been altered since the page was last loaded into main memory.

### Page table structure

1. To read a word from memory, translate the virtual or logical address consisting of page number and offset into physical address consisting of frame number and offset, using a page table.
2. Page table must be stored in main memory to access it.

### Hardware implementation for virtual memory paging

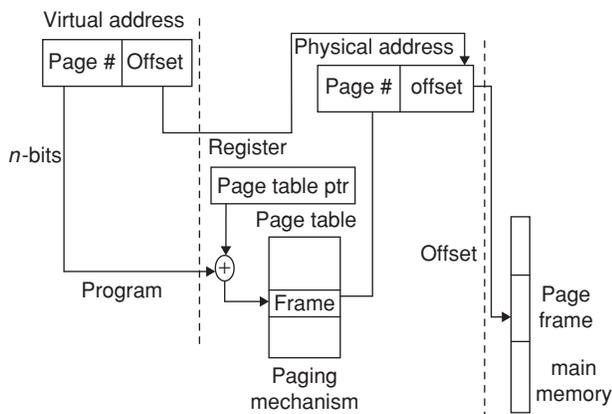
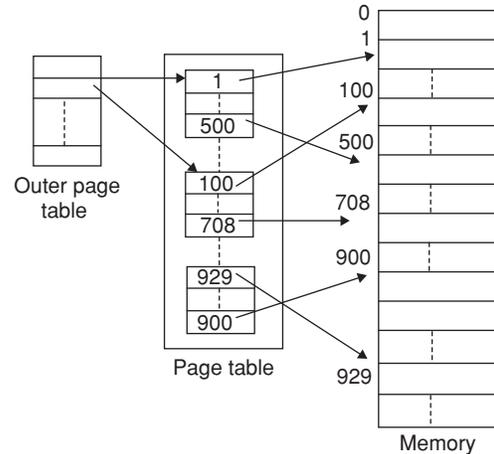


Figure 3 Address translation in paging system.

1. The amount of memory used by the page tables could be high.
2. To overcome this problem, most virtual memory schemes store page tables in virtual memory rather than real memory, that is, page tables are subject to paging just as other pages are.
3. When a process is running, at least a part of its page table must be in main memory, including the page table entry of the currently executing page.
4. Some processors make use of a two-level scheme to organize large page tables.

### Hierarchical page table

1. If page table size is large, then use hierarchical page table.
2. The logical address space is broken up into multiple page tables.

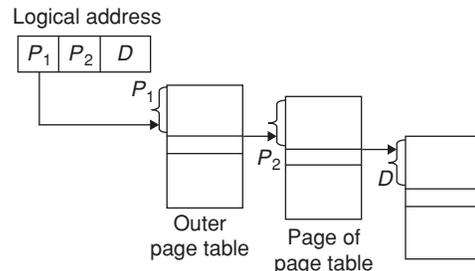


A logical address space (on 32-bit machine with 1 K page size) is divided into a page number consisting of 22-bits, page offset consisting of 10-bits. The page number is paged, the page number is divided into 12-bit page number and 10-bit page offset.

| Page number |       | Page offset |
|-------------|-------|-------------|
| $P_1$       | $P_2$ | $D$         |
| 12          | 10    | 10          |

Here,  $P_1$  is an index into the outer page table,  $P_2$  is the displacement within the page of the outer page table.

### Address translation (diagrammatic)



**Drawback:** Page table size is proportional to that of the virtual address space.

### Inverted page table

1. Here, the page number portion of a virtual address is mapped into a hash value, using simple hash function.
2. The hash value is a pointer to the inverted page table, which contains the page table entries.
3. There is one entry in the inverted page table for each real memory page frame rather than one per virtual page.
4. Thus, a fixed proportion of real memory is required for the tables.
5. One virtual address may map into the same hash table entry, so a chaining technique is used for managing the overflow.
6. The page table's structure is called *inverted*, because it indexes page table entries by frame number rather than by virtual page number.

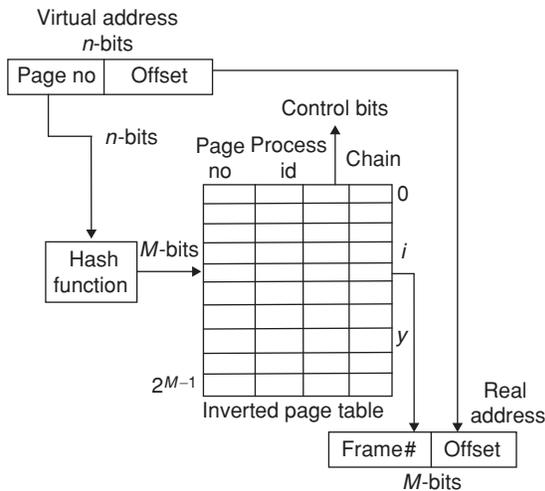
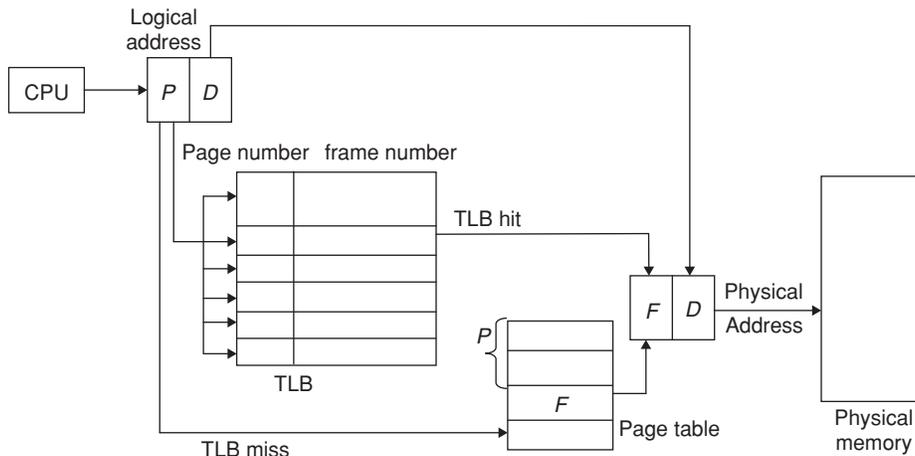


Figure 4 Inverted page table structure.

### Translation look-a-side buffer (TLB)

1. The straight-forward virtual memory scheme would have the effect of doubling the memory access time.
2. To overcome this problem, most virtual memory schemes make use of a special high speed cache for page table entries, usually called TLB.



### Organization of TLB

1. Each entry in TLB must include the page number and the complete page table entry.
2. The TLB may organized its entries either in
  - Direct mapping
  - Associative Mapping
3. The hardware must also consider the ways in which entries are organized and which entry to replace.

**Note:** The virtual memory mechanism must interact with the main memory cache system also.

**Page size** The factors to be considered for page size are as follows:

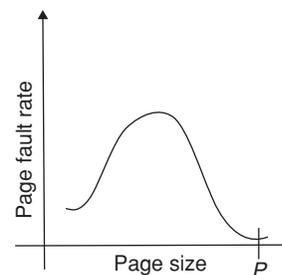
1. If page size is smaller: Then internal fragmentation is less. But it results in larger page tables. For large programs, the page fault rate increases.

3. TLB contains the page table entries that have been most recently used.

### Paging hardware with TLB

1. Given a virtual address, the processor will first examine the TLB. If the desired page table entry is present, that is, TLB hit, then the frame number is retrieved and real address is formed.
2. If there is a TLB miss, the processor uses the page number to index the process page table and examine the corresponding page table entry.
3. If present bit is set, then the page is in main memory and the processor can retrieve the frame number from the page table entry to form the real address.
4. The processor also updates the TLB to include this new page table entry.
5. If the page is not in main memory, then page fault is issued.
6. Then the OS will load the needed page and updates the page table.

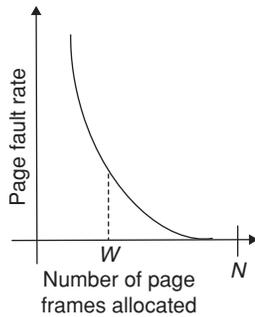
2. Rotational secondary devices favour a larger page size for more efficient block transfer.



where  $P$  = size of entire process.

The above figure shows the relationship between page size and page fault rate.

The page fault rate is also determined by the number of frames allocated to a process. This relation is shown below.



where  $W$  = working set size  
 $N$  = Total number of pages in process

**Note:** The design issue of page size is related to the size of physical main memory and program size.

### Advantages of virtual memory paging

1. No external fragmentation
2. Higher degree of multiprogramming
3. Large virtual address space

### Disadvantage

Overhead of complex memory management.

### Segmentation with Virtual Memory

1. Memory consists of multiple segments.
2. Segments are of unequal and dynamic in size.
3. It simplifies the handling of growing data structures.
4. It allows programs to be altered and recompiled independently.
5. It lends itself to sharing among processes.
6. It lends itself to protection.
7. A unique segment table is associated with each process.
8. The virtual address and segment table entries are as shown below:

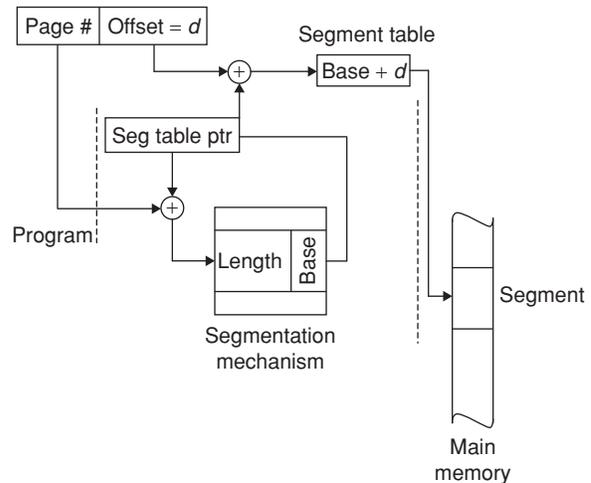
| Virtual address |        |
|-----------------|--------|
| Segment number  | Offset |

| Segment table entry |     |                    |        |              |
|---------------------|-----|--------------------|--------|--------------|
| $P$                 | $M$ | Other control bits | Length | Segment base |

9. Only some of the segments of a process may be in main memory. To identify which segment is present in the main memory, use present bit  $P$ .

10. To know whether the segment is modified or not, use  $M$ -bit.

### Address translation in a segmentation system (using virtual memory)



### Advantages

1. No internal fragmentation
2. Higher degree of multiprogramming
3. Large virtual address space.
4. Protection and sharing support.

### Disadvantage

1. Overhead of complex memory management.

### Combined paging and segmentation

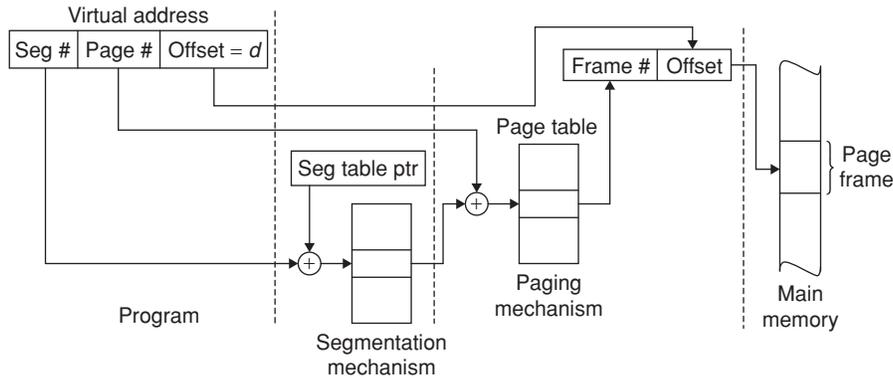
1. Here, the users address space is broken up into a number of segments by the programmer.
2. Each segment is, in turn, broken up into a number of fixed size pages, which are equal in length to a main memory frame.
3. If a segment has length less than that of a page, the segment occupies just one page.
4. The virtual address, segment table and page table entries are as shown below:

| Virtual address |             |        |
|-----------------|-------------|--------|
| Segment number  | Page number | Offset |

| Segment table entry |        |              |
|---------------------|--------|--------------|
| Control bits        | Length | Segment base |

| Page table entry |     |                    |              |
|------------------|-----|--------------------|--------------|
| $P$              | $M$ | Other control bits | Frame number |

### Structure of combined segmentation/paging system



#### Notes:

1. Segmentation lends itself to the implementation of protection and sharing policies.
2. To achieve sharing, it is possible for a segment to be referenced in the segment tables of more than one process.

## OS SOFTWARE FOR MEMORY MANAGEMENT

We consider the following software policies for virtual memory:

1. Fetch policy:
  - Demand
  - Prepaging
2. Placement policy
3. Replacement policy:
  - Optimal
  - LRU
  - FIFO
  - Clock
4. Resident set management
  - Resident set size
    - I. Fixed
    - II. Variable
  - Replacement scope
    - I. Local
    - II. Global
5. Cleaning policy
  - Demand
  - Pre-cleaning
6. Load control
  - Degree of multi-programming

**Fetch policy** Determines when a page should be brought into main memory.

1. **Demand paging:** Here a page is brought into main memory only when a reference is made to a location on that page.

2. **Prepaging:** It is a technique that reduces the large number of page faults at process start up.

- Prepaging is used to get before all or some of the pages a process will need, before they are referenced.
- If prepagged pages are unused, I/O and memory would be wasted.
- Assume 's' pages are prepagged and  $\alpha$  of the pages are used.
- Cost of  $(s * \alpha)$  to save page faults greater or less than the cost of prepagging  $s \times (1 - \alpha)$  unnecessary page. If  $\alpha$  near zero  $\Rightarrow$  prepagging is lost.

### Placement Policy

1. Determines where in real memory a process piece is to reside.
2. In pure segmentation system, the policies like best fit, first fit, etc., are used.
3. For a system that uses either pure paging or paging combined with segmentation, placement is usually irrelevant.

**Replacement policy** This deals with the selection of a page in main memory to be replaced when a new page must be brought in. In a replacement policy, we have to consider the following:

1. How many page frames are to be allocated to each active process.
2. Whether the set of pages to be considered for replacement should be limited to those of the process that caused the page fault or encompass all the page frames in main memory.
3. Among the set of pages considered, which particular page should be selected for replacement.

### Page Fault

1. Whenever a processor needs to execute a particular page and that page is not available in main memory, this situation is said to be *page fault*.

- When the page fault occurs, the page replacement will be done.
- 'Page Replacement' means select a victim page in the main memory, replace that page with the required page from the backing store (disk).
- Some of the replacement algorithms are as follows:
  - FIFO
  - Optimal
  - LRU
  - Clock

### FIFO (First-in-First-Out Algorithm)

- Replace a page that is the oldest page of all the pages of the main memory.
- Focuses on the length of time a page has been in memory rather than how much the page is being used.

#### Example:

Consider the reference string: 0 1 2 3 0 1 2 3 0 1 2 3 1

|  |   |   |   |   |   |   |
|--|---|---|---|---|---|---|
|  | 0 | 1 | 2 | 3 | 0 | 1 |
|  | F | F | F | F | F | F |
|  | 0 | 0 | 0 | 3 | 3 | 3 |
|  |   |   | 1 | 1 | 0 | 0 |
|  |   |   | 2 | 2 | 2 | 1 |

|  |   |   |   |   |   |   |    |
|--|---|---|---|---|---|---|----|
|  | 2 | 3 | 0 | 1 | 2 | 3 | 1  |
|  | F | F | F | F | F | F | H  |
|  | 2 | 2 | 2 | 1 | 1 | 1 | 1* |
|  | 0 | 3 | 3 | 3 | 2 | 2 | 2  |
|  | 1 | 1 | 0 | 0 | 0 | 3 | 3  |

Here the symbol 'F' indicates page fault.

The number of page faults = 12

'H' indicates the page is already in the memory. The remaining pages are not present in memory that is why page fault occurs. In general, the more frames there are, the less page fault.

$$\text{Page fault rate} = \frac{\text{Number of page faults}}{\text{Number of bits in reference string}}$$

$$= \frac{12}{13} = 0.923 = 92.3\%$$

### Belady's Anomaly

**Example:** Consider the reference string: 1, 2, 3, 4, 1, 2, 5, 1, 2, 3, 4, 5

Number of frames = 4

|  |   |   |   |   |   |   |
|--|---|---|---|---|---|---|
|  | F | F | F | F | H | H |
|  | 1 | 1 | 1 | 1 | 1 | 1 |
|  |   | 2 | 2 | 2 | 2 | 2 |
|  |   |   | 3 | 3 | 3 | 3 |
|  |   |   | 4 | 4 | 4 | 4 |
|  | 1 | 2 | 3 | 4 | 1 | 2 |

|   |   |   |   |   |   |
|---|---|---|---|---|---|
| F | F | F | F | F | F |
| 5 | 5 | 5 | 5 | 4 | 4 |
| 2 | 1 | 1 | 1 | 1 | 5 |
| 3 | 3 | 2 | 2 | 2 | 2 |
| 4 | 4 | 4 | 3 | 3 | 3 |
| 5 | 1 | 2 | 3 | 4 | 5 |

The number of page faults = 10

Consider the same reference string with three frames.

|   |   |   |   |   |   |
|---|---|---|---|---|---|
| 1 | 2 | 3 | 4 | 1 | 2 |
| F | F | F | F | F | F |
| 1 | 1 | 1 | 4 | 4 | 4 |
|   | 2 | 2 | 2 | 1 | 1 |
|   |   | 3 | 3 | 3 | 2 |

|   |   |   |   |   |   |
|---|---|---|---|---|---|
| 5 | 1 | 2 | 3 | 4 | 5 |
| F | H | H | F | F | H |
| 5 | 5 | 5 | 5 | 5 | 5 |
| 1 | 1 | 1 | 3 | 3 | 3 |
| 2 | 2 | 2 | 2 | 4 | 4 |

Here number of page faults = 9

Here, as the number of frames increases the page fault also increases. This is known as *Belady's anomaly*.

### Optimal page replacement algorithm

- Replace the page that will not be used for the longest period of time.

**Example:** Consider the reference string: 1, 2, 3, 4, 1, 2, 5, 1, 2, 3, 4, 5 using four frames.

|  |   |   |   |   |   |   |
|--|---|---|---|---|---|---|
|  | F | F | F | F | H | H |
|  | 1 | 1 | 1 | 1 | 1 | 1 |
|  |   | 2 | 2 | 2 | 2 | 2 |
|  |   |   | 3 | 3 | 3 | 3 |
|  |   |   | 4 | 4 | 4 | 4 |
|  | 1 | 2 | 3 | 4 | 1 | 2 |

|  |   |   |   |   |   |   |
|--|---|---|---|---|---|---|
|  | F | H | H | H | F | H |
|  | 1 | 1 | 1 | 1 | 4 | 4 |
|  | 2 | 2 | 2 | 2 | 2 | 2 |
|  | 3 | 3 | 3 | 3 | 3 | 3 |
|  | 5 | 5 | 5 | 5 | 5 | 5 |
|  | 5 | 1 | 2 | 3 | 4 | 5 |

Number of page faults = 6

### Disadvantage

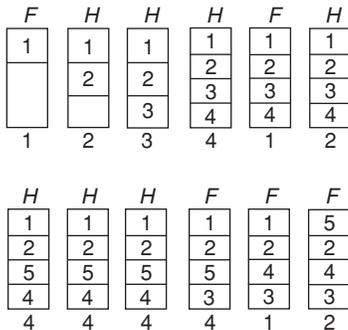
It requires future knowledge of reference string, so used for comparison studies.

### LRU (least recently used) algorithm

- Replace a page that has not been used for the longest period of time.
- It looks backward in time rather than forward.
- It associates with each page the time of that page last used.
- Two methods of implementation:
  - Counters
  - Stack

**Example:** Consider reference String: 1, 2, 3, 4, 1, 2, 5, 1, 2, 3, 4, 5

Number of frames = 4



Number of page faults = 8 (less than FIFO)

### Approach to Implement LRU Replacement

LRU is a good page replacement policy, but the problem with these is how to determine the frame used for the last time.

This is implemented by using two approaches:

1. Using counters
2. Using stacks

Using counters, LRU is implemented as follows: Every page entry has a counter, whenever a page is referenced, the clock value is copied into the counter. If the page has to be replaced, then it refers to the look up of the counter, which ever is having oldest time that is changed.

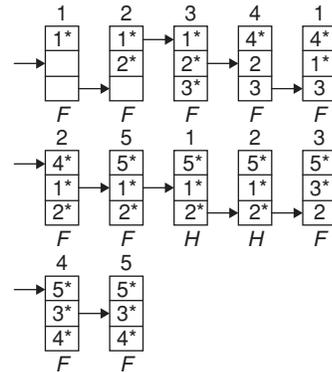
Using stack, LRU is implemented as follows: Form a doubly linked list of page numbers and keep it in stack whenever a page is referenced it is moved to the top of the stack that is top of the stack contains recently referenced page. Bottom of the stack will have least recently used one.

### Clock replacement algorithm

1. The simplest form of clock policy requires the association of an additional bit with each frame, referred to as the use bit.
2. When a page is first loaded into frame in memory, the use bit for that frame is set to 1.
3. Whenever the page is subsequently referenced, its use bit is set to 1.
4. The set of frames that are candidates for replacement is considered to be a circular buffer, with which a pointer is associated.
5. When a page is replaced, the pointer is set to indicate the next frame in the buffer after the one just updated.
6. When it comes time to replace a page, the OS scans the buffer to find a frame with a use bit set to 0.
7. Each time it encounters a frame with a use bit of 1, it resets that bit to 0 and continues on.
8. If any of the frames in the buffer have a use bit of 0 at the beginning of this process, the first such frame encountered is chosen for replacement.

9. If all of the frames have a use bit of 1, then the pointer will make one complete cycle through the buffer, setting all the use bits to 0 and stop at its original position, replacing the page in that frame.

**Example:** Consider the reference string 1, 2, 3, 4, 1, 2, 5, 1, 2, 3, 4, 5, with three main memory frames:



∴ Number of misses = 9

**Note:** The clock algorithm was approximately closer in performance to LRU.

### Effective memory access time

The percentage of times a page number is found in the associative registers is called the *hit ratio*. If we fail to find the page number in the associative registers, then we must first access memory for the page table and frame number, and then access the required byte in memory. To find the effective access time, we should weigh each case by its probability.

EMAT: Is given as  $= p * s + (1 - p) * m$ .

Where

- $p$  = Page fault rate
- $s$  = Page fault service time
- $m$  = Main memory access time
- $(1 - p)$  = page hit ratio.

**Frame locking** Some of the frames in main memory may be locked. When a frame is locked, the page currently stored in that frame may not be replaced.

**Page buffering** To improve performance, a replaced page is not lost but rather is assigned to one of two lists:

1. The free page list if the page has not been modified or
2. The modified page list if it has modified.

### Resident set management

**Resident set size** The OS must decide how many pages to bring in, that is, how much main memory to allocate to a particular process. Two policies are there:

1. Fixed allocation
2. Variable allocation

**Fixed allocation** This policy gives to a process a fixed number of frames in main memory within which to execute.

**Variable allocation** This policy allows the number of page frames allocated to a process to be varied over the life time of the process.

**Replacement scope** This is of two types as follows:

**Local page replacement:** When a process requests for a new page to be brought in and there are no free frames in the memory, we choose a frame allocated to only that process for replacement.

**Global replacement:** It allows a process to select a replacement frame from the set of all frames, even if that frame is currently allocated to some other processes. So, one process can take a frame from another.

#### Fixed allocation, local replacement

1. Number of frames allocated to a process is fixed.
2. Page to be replaced is chosen from among the frames allocated to that process.

#### Variable allocation, global scope

1. Page to be replaced is chosen from all available frames in main memory.
2. Size of resident set of processes varies.

#### Variable allocation, local scope

1. The number of frames allocated to a process may be changed from time to time.
2. Page to be replaced is chosen from among the frames allocated to that process.

**Working set:** This strategy is used to determine the resident set size and the timing of changes.

The working set with parameter  $\Delta$  for a process at virtual time  $t$ , which we designated as  $W(t, \Delta)$ , is the set of pages of that process that have been referenced in the last  $\Delta$  virtual time units.

**Virtual time:** Consider a sequence of memory references,  $r(1), r(2), \dots$  in which  $r(i)$  is the page that contains the  $i$ th virtual address generated by a given process.

Time is measured in memory references; thus,  $t = 1, 2, 3, \dots$  measures the processes internal virtual time.

The variable ' $\Delta$ ' is a window of virtual time over which the process is observed.

The working set size will be a non-decreasing function of the window size.

$$W(t, \Delta + 1) \supseteq W(t, \Delta).$$

For the sequence of page references 24, 15, 18, 23, 17, 15, 24, 18, 17, 17, 15. And window size = 2 then working set will be

{24, {24, 15}}, {15, 18}, {18, 23}, {23, 24}, {24, 17}, {17, 18}, {18, 24}, {18, 24}, {18, 17}, {17}, {17, 15}}

#### Page fault frequency (PFF) algorithm

1. The algorithm requires a use bit to be associated with each page in memory.
2. The bit is set to 1, when that page is accessed.
3. When a page fault occurs, the OS notes the virtual time since the last page fault for the process.
4. A threshold  $F$  is defined. If the amount of time since the last page fault is less than  $F$ , then a page added to resident set of the process.
5. Otherwise, discard all pages with a use bit of 0 and shrink the resident set according.

**Note:** PFF does not perform well during the transient periods when there is a shift to a new locality.

#### Variable interval sampled working set (VSWS)

1. The VSWS policy evaluates the working set of a process at sampling instances based on elapsed virtual time.
2. VSWS considers three parameters:  
 $M$ : The minimum duration of sampling interval  
 $L$ : The maximum duration of sampling interval  
 $Q$ : The number of page faults that are allowed to occur between sampling instances.
3. The policy works as following:
  - If virtual time since the last sampling instance reaches  $L$ , then suspend the process and scan the use bits.
  - If, prior to an elapsed virtual time of  $L$ ,  $Q$  page faults occur,
    - I. If the virtual time since the last sampling instance is less than  $M$ , then wait until the elapsed virtual time reaches  $M$  to suspend the process and scan the use bits.
    - II. If the virtual time since the last sampling instance is greater than or equal to  $M$ , suspend the process and scan the use bits.

**Cleaning policy** It determines when a modified page should be written out to secondary memory. Two approaches are as follows:

1. Demand cleaning
2. Pre-cleaning

**Demand cleaning** A page is written out to secondary memory only when it has been selected for replacement.

**Pre-cleaning** This policy writes modified pages before their page frames are needed so that pages can be written out in batches.

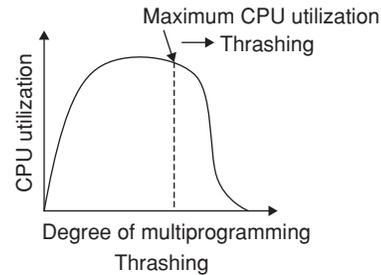
**Load control** It is concerned with determining the number of processes that will be resident in main memory, which has been referred to as the multiprogramming level.

If too few processes are resident at any one time, it leads to swapping. If too many processes are present, thrashing will occur.

**Cause of thrashing** Consider the following scenario:

The OS monitors CPU utilization. If the utilization is too low, we increase the degree of multiprogramming by introducing a new process to the system. A global page replacement algorithm is used, which replaces pages with no regard to the process to which they belong. Now, say that a process enters a new phase in its execution and needs more frames. It starts faulting and taking frames away from other processes. These processes need those pages, however and so they also fault, taking frames from other processes. These faulting processes must use the paging device to swap pages in and out. As they queue up for paging device, the ready queue empties. As processes wait, for the paging device, CPU utilization decreases the CPU scheduler sees the decreasing CPU utilization; so it increases the degree of multiprogramming. The new process tries to get started by taking frames from running processes, causing

more page faults, and a longer queue for paging device. As a result, CPU utilization drops even, further. The CPU scheduler tries to increase the degree of multiprogramming even more. Thrashing occurs, and the system throughput plunges. The page fault rate (FT) increases tremendously. Effective memory access time increases. No work is getting done because the processes are spending all their time in paging.



## EXERCISES

### Practice Problems I

**Directions for questions 1 to 20:** Select the correct alternative from the given choices.

1. Consider a 3-level memory hierarchy shown in the following table, with access times to memory:

| Hierarchy level | Cache hit ratio | Page transfer time |
|-----------------|-----------------|--------------------|
| $M_1$           | 0.55            | 0.003 ms           |
| $M_2$           | 0.92            | 0.3 ms             |
| $M_3$           | –               | 1.0 ms             |

When a miss occurs, data is fetched from the next level. Calculate the average time required for a process to read one word from the memory system.

- (A) 0.379 (B) 0.162  
(C) 0.2798 (D) 0.172
2. Consider a memory system with FIFO page replacement algorithm policy. For an arbitrary page access pattern, increasing the number of page frames in main memory will
- (A) Always decrease the number of page faults  
(B) Always increase the number of page faults  
(C) Sometimes increase the number of page faults  
(D) Never effect the number of page faults
3. Consider the below page address stream generated by executing a program:
- 4 5 4 3 7 4
- Assuming that LRU is used for page replacement and at most three frames are available in the memory for the process, find the number of page faults that can occur (initially all frames empty).
- (A) 0 (B) 1  
(C) 3 (D) 4

4. Consider a logical address space of four pages of 2048 words each mapped into a physical memory of 32 frames. How many bits in logical address?
- (A) 12-bits (B) 14-bits  
(C) 13-bits (D) 11-bits
5. The time taken to service a page fault is on average 10 ms and the memory access time is 20  $\mu$ s. If the hit ratio is 70%, calculate the average access time.
- (A) 3018  $\mu$ s (B) 4014  $\mu$ s  
(C) 3014  $\mu$ s (D) 4024  $\mu$ s
6. Consider the following page trace:
- 4, 3, 2, 1, 4, 3, 5, 4, 3, 2, 1, 5
- Number of frames for the Job  $M = 4$ . Then the page fault ratio using FIFO technique will be
- (A) 63% (B) 75%  
(C) 83% (D) 94%
7. The available main memory for loading pages is 64 MB with a frame size of 8 MB. If pages of size 6 MB, and 4 MB are loaded into memory, what is the percentage of the internal fragmentation resulted?
- (A) 42.5 (B) 37.5  
(C) 57.5 (D) 62.45
8. A demand paging system takes 50 time units to handle a page fault and 200 time units to replace a dirty page. Access time of memory is 2 time units. Probability of page fault and dirty page is  $P$ . Average access time is 4 time units. Then what is the value of  $P$ ?
- (A) 0.037 (B) 0.027  
(C) 0.012 (D) 0.042
9. Assume that a total memory 20 KB is available with no partition. If Buddy system technique is used and there are total of four partitions to serve the request, the closest range of the requested size is

- (A) 12.5 and 25                      (B) 14.5 and 30  
(C) 15.5 and 25                      (D) 14.5 and 40
10. A system uses FIFO page replacement algorithm. It has three page frames with no pages loaded. First 50 pages are accessed in some order and the same pages are accessed in the reverse order. What is the number of page faults?  
(A) 98                                      (B) 96  
(C) 97                                      (D) 95
11. If 32-bit addressing is used for pages whose maximum size is 512 KB, what is the maximum number of pages that can be addressed?  
(A) 4096                                  (B) 2048  
(C) 8192                                  (D) 16384
12. Calculate the overhead due to page table if given the average process in bytes is 16-bytes, the page size is 32-bytes and the page entry is 2-bytes.  
(A) 15                                      (B) 17  
(C) 18                                      (D) 19
13. In a system with 32-bit virtual address and 1 KB page size, use of one-level page tables for virtual to physical address translation is not practical because of  
(A) The large amount of internal fragmentation  
(B) The large amount of external fragmentation  
(C) The large memory overhead in maintaining page tables  
(D) The large computation overhead in the translation process
14. If an instruction takes time 10 m sec if there is no page fault and time 20 m sec if there is a page fault, what is the effective instruction time if page fault occurs once every 5 instructions?  
(A) 12.5 msec                          (B) 12 msec  
(C) 14 msec                              (D) 15 msec
15. Let an instruction take 10 ms and page fault takes an additional 5 ms. If the average page fault occurs after 20 instructions, the effective instruction time will be  
(A) 10 ms                                  (B) 10.25 ms  
(C) 0.25 ms                              (D) 10.75 ms
16. A 0.8 MB-sized memory is managed using variable partitions, while the rest of the memory is occupied by a 0.26 MB partition, 0.27 MB partition and 0.25 MB partitions in the order. Best-fit strategy is being adopted where would be a 0.18 MB allocation request is fulfilled?  
(A) 0.26  
(B) 0.27  
(C) 0.25  
(D) Request will be denied
- Common data for questions 17 and 18:** Suppose that the OS uses variable length partitions for memory management. At some particular time, the running process occupies a partition between physical addresses 20,000 and 40,000.
17. The values of base and limit register are respectively  
(A) 20, 000, 40,000  
(B) 20,000, 20,000  
(C) 0, 10,000  
(D) 0, 40,000
18. What physical address corresponds to a virtual address of 13,000?  
(A) 13,000  
(B) 43,000  
(C) 33,000  
(D) Out of range
19. Consider a page table where translation look ahead buffer is used. TLB hit ratio is 0.95 and generally takes 1 nanosecond to retrieve the frame number. If a miss is recorded by the TLB then an additional overhead of 10 nanoseconds should be taken into consideration, further cache and main memory reference takes 100 ns on average, what is the average memory fetch time using the TLB? Assume main memory accesses are always a success.  
(A) 101 ns                                  (B) 105 ns  
(C) 200 ns                                  (D) 205 ns
20. Consider a 1 MB process which is divided into five segments. Each segment is further divided into pages whose size is 4 KB. What is the maximum segments possible? Assume that the system is byte addressable.  
(A) 64                                          (B) 128  
(C) 256                                      (D) 512

## Practice Problems 2

**Directions for questions 1 to 20:** Select the correct alternative from the given choices.

1. Consider a logical address space of 32 pages of 2048 words mapped into memory of 64 frames. Then the number of bits required for logical address are  
(A) 16-bits                                  (B) 17-bits  
(C) 18-bits                                  (D) 20-bits
2. In which of the page table techniques the logical address space is broken into multiple page table?  
(A) Inverted Page Table  
(B) Hierarchical Page Table  
(C) Hashed Page Table  
(D) None of the above
3. Consider a system with 70% hit ratio, 60 nanoseconds time to search the associative registers, 800 nanoseconds to access memory. What is the effective memory access time?  
(A) 1200 nsec                              (B) 1100 nsec  
(C) 1300 nsec                              (D) 2200 nsec

4. On a system using fixed partitions, all of size  $2^8$ , the number of bits used by the limit register is  
 (A) 128 (B) 256  
 (C) 8 (D) 1024
5. Working set  $(t, k)$  at an instant of time,  $t$ , is the set of  
 (A)  $k$  future references that the operating system will make  
 (B) future references that the operating system will make in the next ' $k$ ' time units  
 (C)  $k$  references with high frequency  
 (D) pages that have been referenced in the last  $k$  time units.
6. Cache and interleaved memories are ways of speeding up memory access between CPUs and slower RAM. Which of the following memory models are best suited (i.e., improves performance the most) for which programs?  
 (i) Cached memory is best suited for small loops.  
 (ii) Interleaved memory is best suited for small loops.  
 (iii) Interleaved memory is best suited for large loops.  
 (iv) Cache memory is best suited for large sequential code.  
 (A) (i) and (ii) are true  
 (B) (i) and (iii) are true  
 (C) (iv) and (ii) are true  
 (D) (iv) and (iii) are true
7. A paging system with a page table in memory every reference to memory takes 100 ns. The TLB hit ratio is 85% and the time needed for searching TLB is almost negligible. What is the effective memory access time?  
 (A) 115 ns (B) 135 ns  
 (C) 145 ns (D) 125 ns
8. Consider the page sequence 4, 2, 1, 5, 3, 2, 1, 5, 0, 2, 5. If FIFO page replacement algorithm is used and frame size is 3, then the percentage of page fault is  
 (A) 99% (B) 90.9%  
 (C) 80.8% (D) 89.9%
9. If the page size is 32 KB, primary page table contains 4096 entries and the secondary page table contains 256 entries, then what is the size of logical address in bits?  
 (A) 15-bits (B) 20-bits  
 (C) 32-bits (D) 35-bits
10. If page size is 2 KB and logical address is 20-bit, then the number of entries in the page table is  
 (A) 2048 B (B) 256 B  
 (C) 512 B (D) 1 MB
11. Consider a paging system with the page table stored in memory. If a memory reference takes 200 ns, how long does a paged memory reference take?  
 (A) 100 ns (B) 200 ns  
 (C) 300 ns (D) 400 ns

12. Consider a logical address space of eight pages of 1024 words each mapped onto a physical memory of 32 frames. How many bits are there in the logical address and in the physical address?  
 (A) 10, 18 (B) 13, 18  
 (C) 13, 15 (D) 10, 5
13. For a paged system, TLB hit ratio is 0.9. Let the RAM access time ' $t$ ' be 20 ns and the TLB access time ' $T$ ' be 100 ns. Then effective memory access (with TLB) will be  
 (A) 120 ns (B) 200 ns  
 (C) 130 ns (D) 150 ns
14. Assume that a user program is 100 K words and secondary storage device is a fixed hard disk with an average latency of 8 ms and a transfer rate of 2,50,000 words/second. Then find the swap time of a transfer of 100 K words to or from memory.  
 (A) 816 ms (B) 408 ms  
 (C) 204 ms (D) 8 ms
15. Consider the following segment table:

| Segment | Limit | Base |
|---------|-------|------|
| 0       | 1000  | 1400 |
| 1       | 400   | 6300 |
| 2       | 400   | 4300 |
| 3       | 1100  | 3200 |
| 4       | 1000  | 4700 |

The physical address for a logical address which is in segment 2 with offset 253 is

- (A) 4553 (B) 6353  
 (C) 6253 (D) 4453
16. Consider a process of size 2 MB. If the page size is 0.5 KB, what is the size of the page table (assuming that each page is mapped by a 32-bit size page table entry)?  
 (A) 8 KB (B) 16 KB  
 (C) 24 KB (D) 32 MB
17. A CPU generates 32-bit virtual address. The page size is 2 KB. The translation look-aside buffer (TLB) which can hold 256 page table entries and is two-way set associative mapping. The number of bits in the TLB tag is  
 (A) 10-bits (B) 12-bits  
 (C) 14-bits (D) 15-bits
18. Assume that a total memory  $M$  is available with no partitions made yet. If Buddy system strategy is being used and a total of  $n$  partitions have been made to serve the request. The closest range of the requested size is

- (A)  $\frac{M}{2^n + 1}$  and  $\frac{M}{2^n}$  (B)  $\frac{M}{n}$  and  $\frac{M}{n-1}$   
 (C)  $\frac{M}{2^n}$  and  $\frac{M}{2^{n-1}}$  (D)  $\frac{M}{n+1}$  and  $\frac{M}{n}$

19. The memory has four free blocks of sizes 2K, 6K, 20K, 4K. The request blocks are allocated according to best fit allocation method. The allocation requests are stored in queue as shown:

| Request no.   | $P_1$ | $P_2$ | $P_3$ | $P_4$ | $P_5$ | $P_6$ | $P_7$ |
|---------------|-------|-------|-------|-------|-------|-------|-------|
| Request sizes | 4 K   | 10 K  | 2 K   | 3 K   | 5 K   | 4 K   | 2 K   |
| Usage time    | 1     | 4     | 2     | 6     | 3     | 1     | 8     |

The time at which request for  $P_7$  will be completed is  
 (A) 10 unit time  
 (B) 14 unit time  
 (C) 20 unit time  
 (D) 15 unit time

20. A memory page containing a heavily used variable that was initialized very early and is in constant use is removed when \_\_\_\_\_ page replacement is used.  
 (A) LRU  
 (B) FIFO  
 (C) LFU  
 (D) Optimal

### PREVIOUS YEARS' QUESTIONS

1. A processor uses 36-bit physical addresses and 32-bit virtual addresses, with a page frame size of 4 K bytes. Each page table entry is of size 4 bytes. A three-level page table is used for virtual to physical address translation, where the virtual address is used as follows:

[2008]

- Bits 30–31 are used to index into the first level page table
- Bits 21–29 are used to index into the second level page table
- Bits 12–20 are used to index into the third level page table, and
- Bits 0–11 are used as offset within the page

The number of bits required for addressing the next level page table (or page frame) in the page table entry of the first, second and third level page tables are, respectively,

- (A) 20, 20 and 20                      (B) 24, 24 and 24  
 (C) 24, 24 and 20                      (D) 25, 25 and 24
2. How many 32 K × 1 RAM chips are needed to provide a memory capacity of 256 K bytes? [2009]  
 (A) 8                                          (B) 32  
 (C) 64                                        (D) 128
3. In which one of the following page replacement policies, Belady's anomaly may occur? [2009]  
 (A) FIFO                                      (B) Optimal  
 (C) LRU                                        (D) MRU
4. The essential content(s) in each entry of a page table is/are [2009]  
 (A) Virtual page number  
 (B) Page frame number  
 (C) Both virtual page number and page frame number  
 (D) Access right information
5. A multilevel page table is preferred in comparison to a single-level page table for translating virtual address to physical address because [2009]

- (A) It reduces the memory access time to read or write a memory location.  
 (B) It helps to reduce the size of page table needed to implement the virtual address space of a process.  
 (C) It is required by the translation look-aside buffer.  
 (D) It helps to reduce the number of page faults in page replacement algorithms.

6. A system uses FIFO policy for page replacement. It has four-page frames with no pages loaded to begin with. The system first accesses 100 distinct pages in some order and then accesses the same 100 pages but now in the reverse order. How many page faults will occur? [2010]

(A) 196                                      (B) 192  
 (C) 197                                      (D) 195

7. Let the page fault service time be 10ms in a computer with average memory access time being 20 ns. If one page fault is generated for every  $10^6$  memory accesses, what is the effective access time for the memory? [2011]

(A) 21 ns                                      (B) 30 ns  
 (C) 23 ns                                      (D) 35 ns

8. Consider the virtual page reference string

1, 2, 3, 2, 4, 1, 3, 2, 4, 1

on a demand paged virtual memory system running on a computer system that has main memory size of three-page frames which are initially empty. Let LRU, FIFO and OPTIMAL denote the number of page faults under the corresponding page replacement policy. Then [2012]

- (A) OPTIMAL < LRU < FIFO  
 (B) OPTIMAL < FIFO < LRU  
 (C) OPTIMAL = LRU  
 (D) OPTIMAL = FIFO

9. A RAM chip has a capacity of 1024 words of 8 bits each ( $1 \text{ K} \times 8$ ). The number of  $2 \times 4$  decoders with

enable line needed to construct a  $16\text{ K} \times 16\text{ RAM}$  from  $1\text{ K} \times 8\text{ RAM}$  is [2013]

- (A) 4 (B) 5  
(C) 6 (D) 7

**Common data for Questions 10 and 11:** A computer uses 46-bit virtual address, 32-bit physical address, and a three-level paged page table organization. The page table base register stores the base address of the first-level table ( $T_1$ ), which occupies exactly one page. Each entry of  $T_1$  stores the base address of a page of the second-level table ( $T_2$ ). Each entry of  $T_2$  stores the base address of a page of the third-level table ( $T_3$ ). Each entry of  $T_3$  stores a page table entry (PTE). The PTE is 32-bits in size. The processor used in the computer has a 1 MB 16-way set associative virtually indexed physically tagged cache. The cache block size is 64 bytes.

10. What is the size of a page in KB in this computer? [2013]

- (A) 2 (B) 4  
(C) 8 (D) 16

11. What is the minimum number of page colours needed to guarantee that no two synonyms map to different sets in the processor cache of this computer? [2013]

- (A) 2 (B) 4  
(C) 8 (D) 16

12. Assume that there are three page frames which are initially empty. If the page reference string is 1, 2, 3, 4, 2, 1, 5, 3, 2, 4, 6, the number of page faults using the *optimal replacement policy* is \_\_\_\_\_. [2014]

13. A computer has 20 physical page frames which contain pages numbered 101 through 120. Now a program accesses the pages numbered 1, 2, ... 100 in that order, and repeats the access sequence THRICE. Which one of the following page replacement policies experiences the same number of page faults as the optimal page replacement policy for this program? [2014]

- (A) Least-recently used (B) First-in-first-out  
(C) Last-in-first-out (D) Most-recently-used

14. A system uses three page frames for storing process pages in main memory. It uses the Least Recently Used (LRU) page replacement policy. Assume that all the page frames are initially empty. What is the total number of page faults that will occur while processing the page reference string given below? [2014]

4, 7, 6, 1, 7, 6, 1, 2, 7, 2

15. Consider a paging hardware with a TLB. Assume that the entire page table and all the pages are in the physical memory. It takes 10 milliseconds to search the TLB and 80 milliseconds to access the physical

memory. If the TLB hit ratio is 0.6, the effective memory access time (in milliseconds) is \_\_\_\_\_. [2014]

16. Consider a system with byte-addressable memory, 32-bit logical addresses, 4 kilobyte page size and page table entries of 4 bytes each. The size of the page table in the system in megabytes is \_\_\_\_\_. [2015]

17. Consider a main memory with five page frames and the following sequence of page references: 3, 8, 2, 3, 9, 1, 6, 3, 8, 9, 3, 6, 2, 1, 3. Which one of the following is true with respect to page replacement policies First In First Out (FIFO) and Least Recently Used (LRU)? [2015]

- (A) Both incur the same number of page faults  
(B) FIFO incurs 2 more page faults than LRU  
(C) LRU incurs 2 more pages faults than FIFO  
(D) FIFO incurs 1 more page faults than LRU

18. Consider six memory partitions of sizes 200 KB, 400 KB, 600 KB, 500 KB, 300 KB and 250 KB, where KB refers to kilobyte. These partitions need to be allotted to four processes of sizes 357 KB, 210 KB, 468 KB and 491 KB in that order. If the best fit algorithm is used, which partitions are NOT allotted to any process? [2015]

- (A) 200 KB and 300 KB  
(B) 200 KB and 250 KB  
(C) 250 KB and 300 KB  
(D) 300 KB and 400 KB

19. A computer system implements 8 kilobyte pages and a 32-bit physical address space. Each page table entry contains a valid bit, a dirty bit, three permission bits, and the translation. If the maximum size of the page table of a process is 24 megabytes, the length of the virtual address supported by the system is \_\_\_\_\_ bits. [2015]

20. Consider the following two C code segments.  $Y$  and  $X$  are one and two dimensional arrays of size  $n$  and  $n \times n$  respectively, where  $2 \leq n \leq 10$ . Assume that in both code segments, elements of  $Y$  are initialized to 0 and each element  $X[i][j]$  of array  $X$  is initialized to  $i + j$ . Further assume that when stored in main memory all elements of  $X$  are in same main memory page frame.

Code segment 1: [2015]

```
//initialize elements of Y to 0
//initialize elements X[i][j] of X to i+j
for (i = 0; i < n; i++)
 Y[i] += X[0][i];
```

Code Segment 2:

```
//initialize elements of Y to 0
```

```
//initialize elements X[i] [j] of X
to i + j
for (i = 0; i < n; i++)
 Y[i] += X[i] [0];
```

Which of the following statements is/are correct?

- $S_1$ : Final contents of array  $Y$  will be same in both code segments
- $S_2$ : Elements of array  $X$  accessed inside the for loop shown in code segment 1 are contiguous in main memory
- $S_3$ : Elements of array  $X$  accessed inside the for loop shown in code segment 2 are contiguous in main memory.
- (A) Only  $S_2$  is correct  
 (B) Only  $S_3$  is correct  
 (C) Only  $S_1$  and  $S_2$  are correct  
 (D) Only  $S_1$  and  $S_3$  are correct
21. Consider a computer system with 40-bit virtual addressing and page size of sixteen kilobytes. If the computer system has a one-level page table per process and each page table entry requires 48 bits, then the size of the per-process table is \_\_\_\_\_ megabytes. [2016]
22. Consider a computer system with ten physical page frames. The system is provided with an access sequence  $(a_1, a_2, \dots, a_{20}, a_1, a_2, \dots, a_{20})$ , where each  $a_i$  is a distinct virtual page number. The difference in the number of page faults between the last-in-first-out page replacement policy and the optimal page replacement policy is \_\_\_\_\_. [2016]
23. In which one of the following page replacement algorithms it is possible for the page fault rate to increase even when the number of allocated frames increases? [2016]
- (A) **LRU** (Least Recently Used)  
 (B) **OPT** (Optimal Page Replacement)  
 (C) **MRU** (Most Recently Used)  
 (D) **FIFO** (First In First Out)
24. Recall that Belady's anomaly is that the page-fault rate may *increase* as the number of allocated frames increases. Now, consider the following statements:  
 $S_1$ : *Random page replacement* algorithm (where a page chosen at random is replaced) suffers from Belady's anomaly  
 $S_2$ : *LRU page replacement* algorithm suffers from Belady's anomaly
- Which of the following is CORRECT? [2017]
- (A)  $S_1$  is true,  $S_2$  is true  
 (B)  $S_1$  is true,  $S_2$  is false  
 (C)  $S_1$  is false,  $S_2$  is true  
 (D)  $S_1$  is false,  $S_2$  is false
25. Consider a process executing on an operating system that uses demand paging. The average time for a memory access in the system is  $M$  units if the corresponding memory page is available in memory and  $D$  units if the memory access causes a page fault. It has been experimentally measured that the average time taken for a memory access in the process is  $X$  units.
- Which one of the following is the correct expression for the page fault rate experienced by the process? [2018]
- (A)  $(D - M)/(X - M)$       (B)  $(X - M)/(D - M)$   
 (C)  $(D - X)/(D - M)$       (D)  $(X - M)/(D - X)$

## ANSWER KEYS

### EXERCISES

#### Practice Problems 1

1. B    2. B    3. D    4. C    5. C    6. C    7. B    8. A    9. A    10. C  
 11. C    12. B    13. C    14. B    15. B    16. C    17. B    18. C    19. A    20. C

#### Practice Problems 2

1. B    2. B    3. B    4. C    5. D    6. B    7. A    8. B    9. D    10. C  
 11. D    12. B    13. C    14. A    15. A    16. B    17. C    18. C    19. A    20. B

#### Previous Years' Questions

1. D    2. C    3. A    4. B    5. B    6. A    7. B    8. B    9. B    10. C  
 11. C    12. 7    13. D    14. 6    15. 122    16. 4    17. A    18. A    19. 36    20. C  
 21. 384    22. 1    23. D    24. B    25. B

# Chapter 5

## File Systems, I/O Systems, Protection and Security

### LEARNING OBJECTIVES

- 🔗 File systems
- 🔗 File management systems
- 🔗 File system architecture
- 🔗 Device drivers
- 🔗 Basic input/output supervisor
- 🔗 Logical input/output
- 🔗 Access methods

### FILE SYSTEMS

The file system consists of two distinct parts:

1. Collection of files
2. Directory structure

**File** A file is a named collection of related information that is recorded on secondary storage. The files must have

1. Long-term existence
2. Sharable between processes
3. Structure

**File attributes** A typical file attributes are

1. Name
2. Identifier
3. Type
4. Location
5. Size
6. Protection
7. Time, date and user identification

**File operations** The operations that are applied on files are

1. Creation
2. Deletion
3. Closing
4. Reading
5. Writing

**File types** A common technique for implementing file types is to include the type as part of the file name. The name is split into two parts:

1. A name
2. An extension (usually separated by a period character)

The type of a file may be

1. Executable (exe, com, bin)
2. Object (obj, o)
3. Source code (c, cc, java)
4. Batch (bat, sh)
5. Text (txt, doc)
6. Word processor (wp, text, doc)
7. Library (lib)
8. Print or view (ps, pdf, jpg)
9. Archive (zip, tar)
10. Multimedia (mpeg, mov, rm)

**File structure** The four common terms of file systems are

1. Field
  - Basic element of data
  - Contains a single value
  - Has a particular length and data type
2. Record
  - It is a collection of related fields
  - It is treated as a unit
3. File
  - It is a collection of similar records
  - Treated as a single entity
  - Has file names
  - Access to file may be restricted or unrestricted
4. Database
  - Collection of related data.
  - Relationship exists among elements.

## File Management Systems

It is a set of system software that provides services to users and applications in the use of files. Objectives of file management systems are as follows:

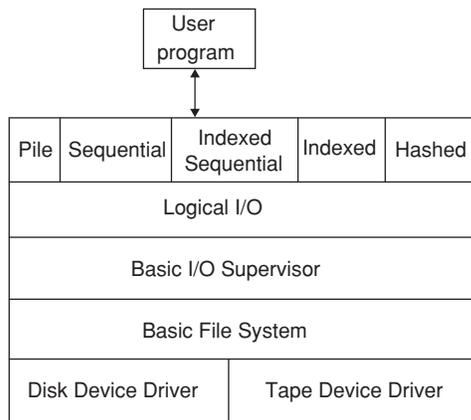
1. To meet the data management needs and requirements of the user including storage of data.
2. To guarantee, that the data in the file are valid.
3. To optimize performance (i.e., throughput, response time).
4. To provide Input/Output (I/O) support for a variety of storage device types.
5. To minimize or eliminate the potential for lost or destroyed data.
6. To provide a standardized set of I/O interface routines.
7. To provide I/O support for multiple users.

The minimal set of requirements from user's point of view for an interactive, general purpose, file system are as follows:

1. Ability to create, delete, read, write and modify files.
2. Controlled access to other users files.
3. Control the type of accesses to files.
4. Restructure the user's files.
5. Able to move data between files.
6. Ability to back up and recover user's files.
7. Able to access a file by name.

## FILE SYSTEM ARCHITECTURE

The file system architecture is shown below:



## Device Drivers

1. It is at the lower part.
2. Communicates directly with peripheral devices.
3. Initiative to start I/O operations on a device.
4. Processes the completion of an I/O request.

## Basic File System

1. Mainly concerned with I/O.
2. Exchanges blocks of data.
3. Deals with placement of blocks.
4. Deals with buffering blocks in main memory.

## Basic I/O Supervisor

1. Responsible for file initiation and termination.
2. Maintain control structures.
3. Selects the device while file I/O is to be performed.
4. Deals with scheduling access to optimize performance.
5. Part of the operating system (OS).

## Logical I/O

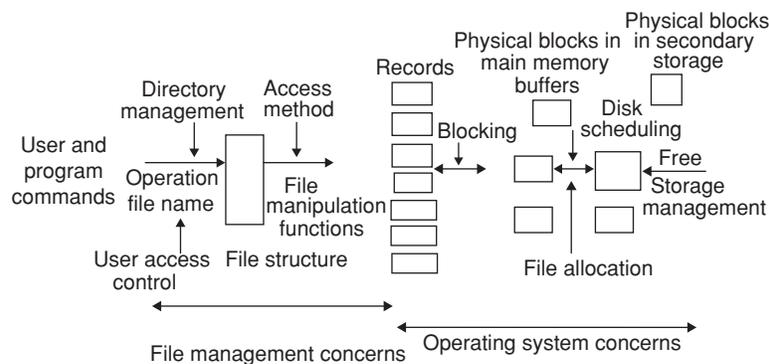
1. Enables users and application to access records.
2. Provides general purpose record I/O capability.
3. Maintains basic data about file.

## Access Methods

1. They reflect different file structures.
2. Provide different ways to access and process data.
3. Access methods are as follows:
  - Sequential: read next, write next
  - Direct: read block n, write block n

## File Management Functions

The functions of a file system is shown below:



### File Organization and Access

File organization refers to the logical structuring of the records as determined by the way in which they are accessed. We choose a particular file organization based on

1. Short access time
2. Ease of update
3. Economy of storage
4. Simple maintenance
5. Reliability

We will discuss five types of file organizations:

1. The pile
2. The sequential file
3. The indexed sequential file
4. The indexed file
5. The direct or hashed file

**The pile** The pile organization is shown in Figure 1:

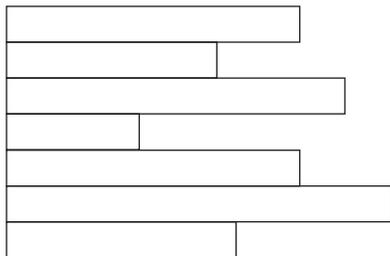


Figure 1 Pile file.

1. Data are collected in order they arrive.
2. The aim is to accumulate mass of data and save.
3. Records may have different fields.
4. There is no structure.
5. Record are accessed by exhaustive search.
6. There can be variable length records.
7. This type of files are encountered when data are collected and stored prior to processing or when data are not easy to organize.

**The sequential file** The sequential file is shown in Figure 2.

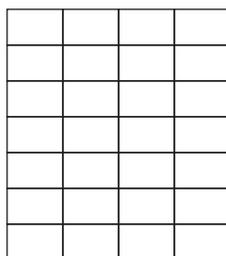


Figure 2 Sequential file.

1. Fixed format used for records.
2. Records are of same length.
3. All fields are of the same order and length.
4. One field is the key field. It uniquely identifies the record. Records are stored in a key sequence.

5. New records are placed in a log file or transaction file.
6. Batch update is performed to merge the log file with the master file.
7. Used in batch applications.
8. Not suitable for interactive applications.

### Indexed sequential file

1. Index provides a look up capability to quickly reach the vicinity of the desired record.
2. It contains key field and a pointer to the main file.
3. Index is searched to find highest key value that is equal to or precedes the desired key value.
4. Search continues in the main file at the location indicated by the pointer.
5. New records are added to an overflow file.
6. Record in main file that precedes it is updated to contain a pointer to the new record.
7. The overflow is merged with main file during a batch update.
8. Multiple indexes for the same key can be set up to increase efficiency.

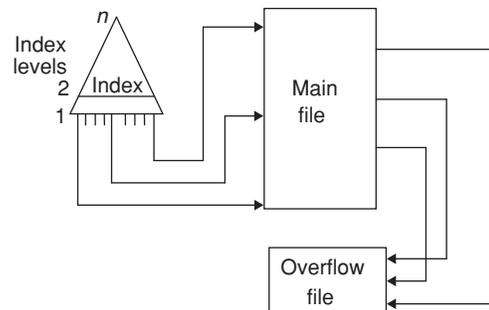
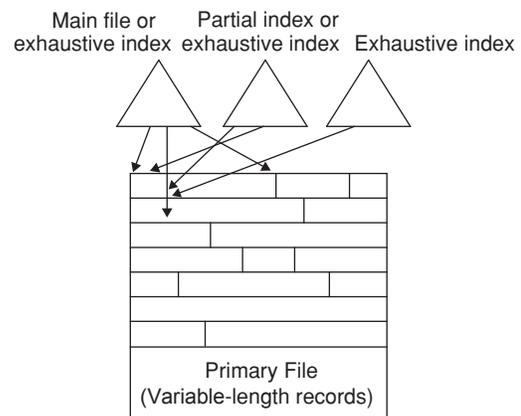


Figure 3 Indexed file.

9. Key field required for each record.
10. It uses multiple indexes for different key fields.
11. It may contain exhaustive index that contains one entry for every record in the main file or partial index.

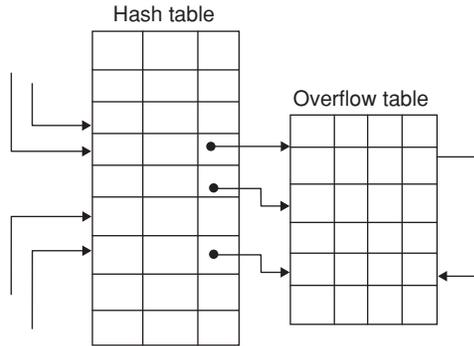


12. Used in the applications where timeliness of information is critical and where data are rarely processed exhaustively.

**Direct or Hashed file**

1. This file has the capability to access any block of a known address.
2. Key field required in each record.
3. No concept of sequential ordering.

Figure 4 shows hashed file organization:



**Figure 4** Hashed file organization.

**FILE DIRECTORIES**

The collection of files is a file directory.

**Contents**

1. Contain information about files, such as attributes, location, and ownership.
2. Itself a file owned by the OS.
3. It provides mapping between file names and the file themselves.

**Directory Structure**

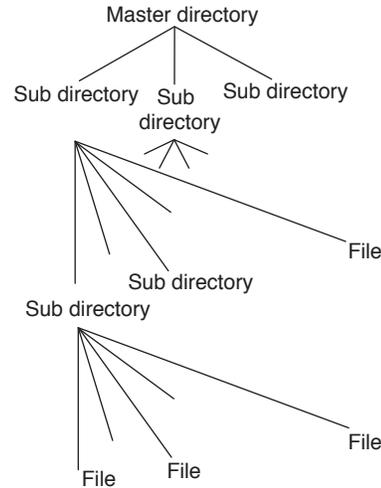
1. It consists of list of entries, one for each file.
2. Sequential file with the name of the file serving as the key.
3. Provides no help in organizing the files.
4. Not scalable (same name cannot be used for two different files).
5. As directory grows in size, searching is too time consuming.

**Two-level Directory Structure**

1. Contains one directory for each user and a master directory.
2. The master directory contains entry for each user
3. User directory is a simple list of files for that user.
4. File naming conflict is solved.

**Hierarchical or Tree-structured Directory**

1. Contains master directory with user directory underneath it.
2. Each user directory may have subdirectories and files as entries.



3. Files can be located by following a path from the root, or master directory down various branches which is called *pathname* for the file.
4. Current directory is called the *working directory*.
5. Files are referenced relative to the working directory.

**Naming** The use of a tree-structured directory minimizes the difficulty in assigning unique names. Any file in the system can be located by following a path from the root or master directory down various braches until the file is reached.

**Path name** The series of directory names, culminating in the file name itself, constitutes a path name for the file.

**Example:** Time/Gate/Exam/OS

The slash is used to delimit names in the sequence. The name of master directory is implicit, because all paths starts at that directory. Files can also be referenced from the working directory.

**File sharing** It has got two issues, such as:

1. Access rights
2. Management of simultaneous access

**Access rights** Users or groups of users are granted certain access rights to a file. A wide range of access rights has been used. The access rights may be

1. None
2. Knowledge
3. Execution
4. Reading
5. Appending
6. Updating
7. Changing protection
8. Deletion

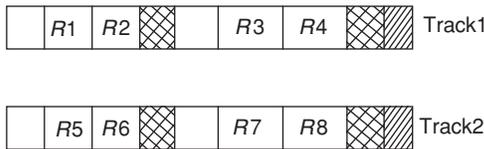
These access rights can be specified to specific users or user groups or all users.

**Record Blocking** For I/O to be performed, records must be organized as blocks. Given the size of a block, there are three methods of blocking that can be used:

1. Fixed blocking
2. Variable length spanned blocking
3. Variable length unspanned blocking.

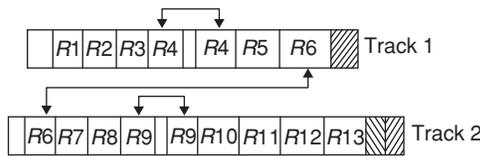
**Fixed blocking**

1. Fixed-length records (Figure 5) are used and an integral number of records are stored in a block.
2. Possibility of internal fragmentation.
3. Used for sequential files.



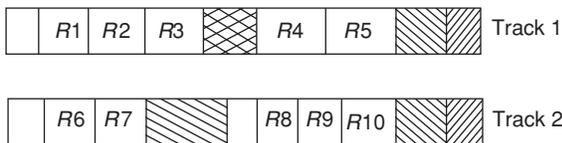
**Figure 5** Fixed blocking.

**Variable length spanned blocking** Variable length records are used and are packed into blocks with no unused space. Some records can span two blocks. These do not limit the size of records.



**Figure 6** Variable blocking: spanned.

**Variable length unspanned blocking** Variable length records (Figure 7) are used, but spanning is not employed. There is wasted space in most blocks and limits record size.



**Figure 7** Variable blocking: unspanned.

- : Waste due to record fit to block size
- : Gaps due to hardware design
- : Waste due to block size constraint from fixed record size.
- : Waste due to block fit to track size

**SECONDARY STORAGE MANAGEMENT**

1. Space (or blocks) must be allocated to files on disk.
2. Need to keep track of the space available (free blocks) for allocation to files.
3. Preallocation of blocks to files can be used to allocate space for files. For this, it needs to know the maximum size of the file at the time of creation.

**File Allocation**

**Preallocation Versus Dynamic Allocation**

1. A preallocation policy requires that the maximum size of a file be declared at the time of the file creation request. For many applications, it is difficult to estimate the file size. It is better to use dynamic allocation, which allocates space to a file in portions as needed.

**Portion size** The portion size which is allocated to a file may be

1. Variable, large contiguous portions.
2. Blocks

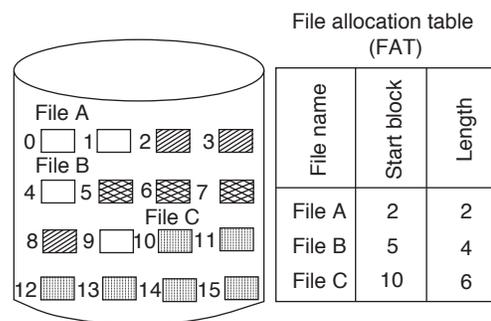
Some strategies for dealing with fragmentation of free space are as follows:

1. *First-fit*: Choose the first unused contiguous group of blocks of sufficient size.
2. *Best-fit*: Choose the smallest unused group that is of sufficient size.
3. *Nearest-fit*: Choose the unused group of sufficient size that is closest to the previous allocation.

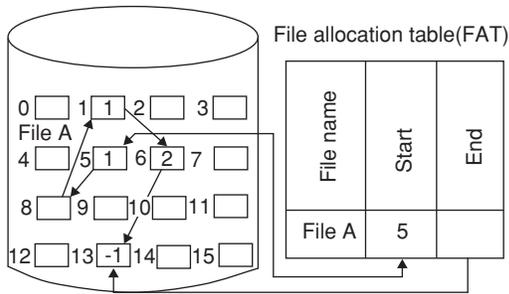
**File allocation methods** It has three methods as follows:

1. Contiguous allocation
2. Chained allocation
3. Indexed allocation

**Contiguous allocation** Here, a single set of blocks is allocated to a file at the time of creation. Only a single entry in the file allocation table is created consisting of starting block and length of the file. It exhibits external fragmentation and performs compaction.

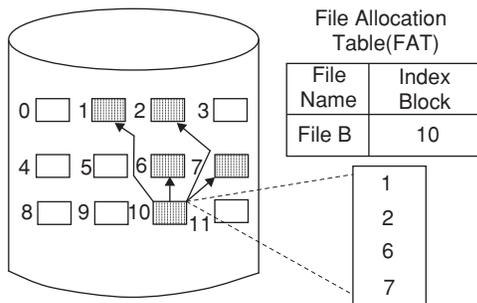


**Linked or chained allocation** The allocation is done on basis of individual block. Each block contains a pointer to the next block in the chain. Only single entry is created in the file allocation table consisting of starting block and length of file. There occurs no external fragmentation and it is best for sequential files. There is no accommodation of principle of locality. If block size is  $n$ , then only  $n - 1$  units of data are stored and 1 unit stores the link information.

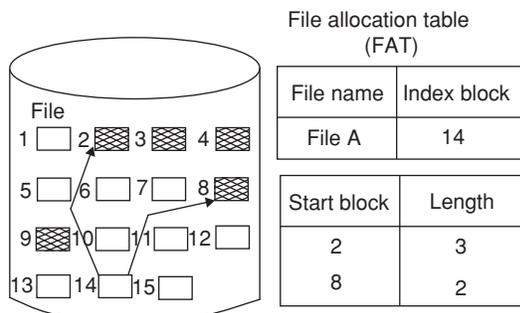


**Indexed file allocation** The file allocation table contains a separate one level index for each file. The index has one entry for each portion allocated to the file. The file allocation table contains block number for the index. If a file requires  $n$  blocks, then  $n + 1$  blocks are used, where the first block contains index information (pointers to data blocks).

**Index Allocation with Block Pointers**



**Indexed Allocation with Variable length Portion**



**Example 1:** A direct access of file has fixed size 50 byte records. Assuming the first record is record 1, the first byte of record 10 will be at what logical location?

**Solution:**

Total records =  $50 \times 10 = 500$

First record is record 1. This record is already read.

Logical location of first byte =  $500 - 50 = 450$

The correct logical location =  $450 + 1 = 451$ .

**Example 2:** A sequential access file has fixed-size 32-byte records. Assuming that the first record is record 0, the first byte of record 20 will be at what location?

**Solution:** Since the first record is record 0, the first byte of record 20 will be at logical location =  $32 \times 20 = 640$

**FREE SPACE MANAGEMENT**

In addition to file allocation table, disk allocation table is also required to know what blocks on the disk are available. Some of the free space management techniques are as follows:

1. Bit tables
2. Chained free portions
3. Indexing
4. Free block list

**Bit tables** This method uses a vector containing, one bit for each block on the disk. Each entry of a '0' corresponds to a free block and each '1' corresponds to a block in use.

**Advantage**

1. Easy to find one or a contiguous group of free blocks.
2. Smaller in size.

The amount of memory required for a block bitmap will be

$$\frac{\text{Disk size in (bytes)}}{8 \times \text{file system block size}}$$

**Chained free portion** The free portions may be chained together by using a pointer and length value in each free portion. This method has negligible space overhead. This method is suitable for all file allocation methods. The disk will become quite fragmented, after some use. It is slower for individual block file creation and also for deletion.

**Indexing** It treats the free space as a file and uses an index table (same as in file allocation). The index should be on the basis of variable size portions rather than blocks.

**Free block list** Here, each block is assigned a number sequentially and the list of the numbers of all free blocks is maintained in a reserved portion of the disk.

**Volumes**

It is a collection of addressable sectors in a secondary memory that an OS or application can use for data storage. The sectors in a volume need not be consecutive on a physical storage device. (a single disk equals one volume).

**UNIX FILE MANAGEMENT**

**I-nodes (Index Node)**

UNIX files are administered by the OS by means of i-node. An i-node (index node) is a control structure that contains the key information needed by the OS for a particular file.

The attributes of the file as well as its permissions and other control information are stored in the i-node. The

exact i-node structure varies from UNIX implementation to another. The FreeBSD i-node structure is shown in Figure 8.

**File Allocation**

1. It is done on a block basis.
2. Allocation is dynamic.

3. The blocks of a file on disk are not necessarily contiguous.
4. An indexed method is used to keep track of each file, with i-node includes a number of direct pointers and three indirect pointers.

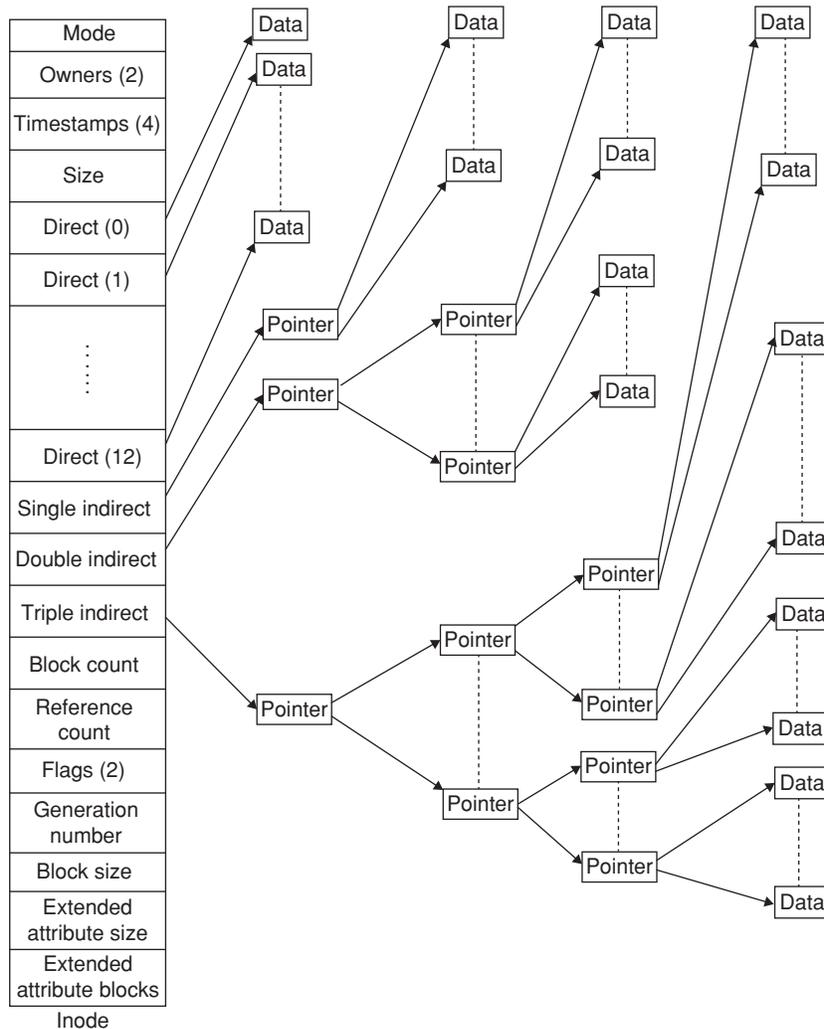


Figure 8 Structure of free BSD i-node and file.

5. The free BSD i-node includes 120 bytes of address information that is organized as fifteen 64-bit addresses.
6. The first 12 addresses point to the first 12 data blocks of the file.
7. If the file requires more than 12 data blocks, one or more levels of indirection are used as follows:
  - The thirteenth address in the i-node points to a block on disk that contains the next portion of the index. This is referred to as the *single indirect block*.
  - If the file contains more blocks, the fourteenth address in the i-node points to a double indirect block. Each block consists of single indirect blocks, each of which contains pointers to file blocks.

- If the file contains still more blocks, the fifteenth address in the i-node points to a triple indirect block that is a third level of indexing. This block points to additional double indirect blocks.

The capacity of FreeBSD file with 4 kB block size is shown below:

| Level           | Number of Blocks   | Number of Bytes |
|-----------------|--------------------|-----------------|
| Direct          | 12                 | 48 K            |
| Single indirect | 512                | 2 M             |
| Double indirect | 512 × 512 = 256 K  | 1 G             |
| Triple indirect | 512 × 256K = 128 M | 512 G           |

The total number of data blocks in a file depends on the capacity of the fixed-size blocks in the system. In FreeBSD, the minimum block size is 4 kB, and each block can hold a total of 512 block addresses. Thus, the maximum size of a file with this block size is over 500 GB.

### Windows NT File System

The windows NT file system provides a combination of reliability, compatibility and performance, which are not available in the FAT file system.

1. It will quickly perform standard file operations, such as write, read and search.
2. It also performs file-system recovery on very large hard disks.
3. NTFS file system formatting on a volume results in the creation of several system files and the master file table (MST), which contains information about all the files and folders on the NTFS (Figure 9).

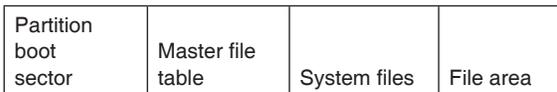


Figure 9 NTFS after formatting.

4. It includes security features required for high-end personal computers and file servers.
5. It supports data access control and ownership privileges.
6. Folders shared on a Windows NT are assigned specific permissions.
7. It allows users to assign permissions to individual files.

### I/O SYSTEMS

#### Design Objectives of OS

Two objectives of OS for I/O systems are as follows:

1. Efficiency
2. Generality

#### Logical Structure of the I/O Function

The three most important logical structures are as follows:

1. Local peripheral devices
2. Communication ports
3. File system

#### Local Peripheral Device



#### Logical I/O

Manages general I/O functions on behalf of user processes, allowing them to deal with the device in terms of a device identifier and simple commands, such as open, close, read, write.

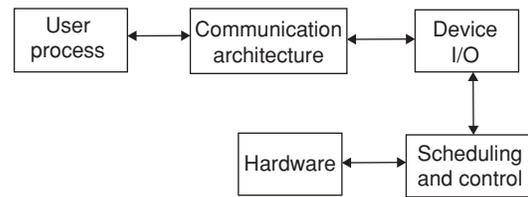
### Device I/O

The requested operations and data are converted into appropriate sequences of I/O instructions, channel commands and controller orders.

### Scheduling and Control

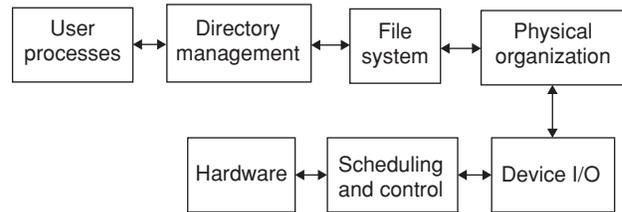
The actual queuing and scheduling of I/O operations occurs at this layer as well as the control of the operations. Interrupts are handled. I/O status is collected and reported.

### Communication Port



Here, the communication architecture may itself consist of a number of layers.

### File System



**Directory management** Here, the symbolic file names are converted to identifiers that either reference the file directly or indirectly through a file descriptor or index table.

**File system** This layer deals with the logical structure of files and with the operations that can be specified by users, such as open, close, read, write.

**Physical organization** Allocation of secondary storage space and main storage buffers is generally treated here.

**I/O buffering** In buffering, we perform input transfers in advance of requests being made and perform output transfers sometime after the request is made. We will discuss

1. Single buffering
2. Double buffering
3. Circular buffering

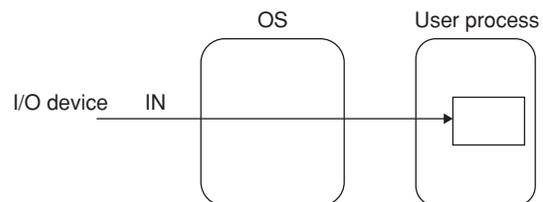


Figure 10 No buffering.

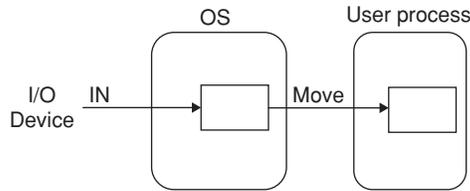


Figure 11 Single buffering.

- When a user process issues an I/O request, the OS assigns a buffer in the system portion of main memory to the operation.
- Provides a speed up compared to the lack of system buffering.
- For block oriented devices, input transfers are made to the system buffer. When the transfer is complete, the process moves the block into user space and immediately requests another block. This is same for block-oriented output.
- Suppose that  $T$  is the time required to input one block and that  $C$  is the computation time that intervenes between input requests. Without buffering, the execution time per block =  $T + C$ . With a single buffer, the execution time per block =  $\max [C, T] + M$ , where  $M$  = Time required to move the data from the system buffer to user memory.
- For stream-oriented I/O, the single buffering scheme can be used in a line-get-a-time fashion or a byte-at-a-time fashion.

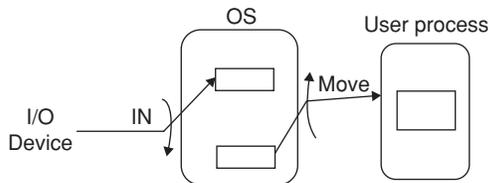


Figure 12 Double buffer

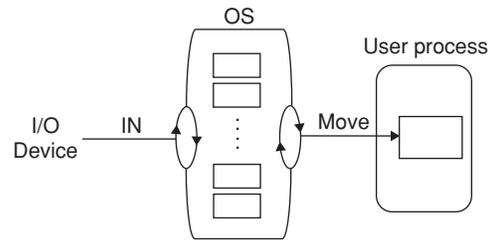
- An improvement over single buffer.
- A process now transfers data to one buffer while the OS empties the other. This is known as *double buffering* or *buffer swapping*.
- For block-oriented transfer, the estimated execution time is  $\max [C, T]$ .
- The stream-oriented transfer may be line-at-a-time or byte-at-a-time.

### Automatic and Explicit Buffering

An indefinite length queue is provided in automatic buffering. Sender never blocks to copy a message. No specifications are present for providing automatic buffering. One method is to reserve large amount of memory. In this method, most of the memory is wasted.

The size of the queue is provided in explicit buffering. Sender will block if the requested space is more than available. In this scheme, memory wastage is less likely to happen.

### Circular Buffer



1. Here, more than two buffers are used. Each individual buffer is one unit of circular buffer.

### Kernel I/O Subsystem

Kernel I/O provides I/O-related services. It is built on the hardware and device driver infrastructure.

One of the responsibilities of Kernel I/O subsystem is to protect itself from the erroneous process and malice users.

Services provided by the I/O subsystem are as follows:

1. Scheduling
2. Buffering
3. Caching
4. Spooling
5. Device reservation
6. Error handling

### DISK SCHEDULING

#### Hard Disk Performance Parameters (Terminologies)

**Seek time** It is defined as the time required to move the disk arm to the required track. It consists of

1. Initial start-up time
2. The time taken to traverse the tracks that has to be crossed once the access arm is up to speed.

$$\text{Seek time, } T_s = m \times n + S$$

where  $T_s$  = Estimated seek time

$n$  = Number of track traversed

$m$  = Constant that depends on the disk drive

$S$  = Start-up time

**Rotational delay** Time required to reach the desired sector by read/write head. Rotational speed ranges from 5400 to 10,000 rpm.

1. Floppy disks typically rotate between 300 and 600 rpm.
2. For 1,00,000 rpm, the average rotational delay will be 3 ms.

**Transfer time** The transfer time to or from the disk depends on the speed of the disk.

$$T = \frac{b}{rN}$$

where,  $T$  = Transfer time

$b$  = Number of bytes to be transferred

$N$  = Number of bytes on a track

$r$  = Rotation speed in revolutions/second

Total average access time

$$T_a = T_s + \frac{1}{2r} + \frac{b}{rN}$$

$T_s \rightarrow$  average seek time

### Hard Disk Scheduling Algorithms

Disk bandwidth and fast access time are to be considered. Bandwidth is the total number of bytes transferred divided by the total time between the first request for service and completion of the last transfer.

**FCFS (First Come First Served)** The disk in controller processes the I/O requests the order in which they arrive, thereby moving backwards and forwards across the surface of the disk to get the next requested location each time (Figure 10).

**Example 3:** A disk queue has the following requests to read tracks:

87, 170, 40, 150, 36, 72, 66, 15

Consider the disk head is initially at cylinder 60. Total head movement =  $(87 - 60) + (170 - 87) + (170 - 40) + (150 - 40) + (150 - 36) + (72 - 36) + (72 - 66) + (66 - 15) = 27 + 83 + 130 + 110 + 114 + 36 + 6 + 51 = 557$  cylinders

Average head movement =  $\frac{557}{8} = 69.6$  cylinders

#### Advantage

Improved response time as a request gets response in fair amount of time.

#### Disadvantages

1. Involves a lot of random head movements and disk rotations.
2. Throughput is not efficient.
3. Used in small systems only where I/O efficiency is not very important.

**Shortest Seek Time First (SSTF)** When a disk operation finishes, choose the request that is closer to the current head position or choose the request that has minimum seek time from the current head position.

**Example 4:** Consider the following requests: 87, 170, 40, 150, 36, 72, 66, 15. Find the average head movement for SSTF.

The initial head position is say 60. Now, closest to the head position is the request at cylinder 66. Then, the closest to 66 is 72, closest request to 72 is 87, and so on. Total head movements =  $(66 - 60) + (72 - 66) + (87 - 72) + (87 - 40) + (40 - 36) + (36 - 15) + (150 - 15) + (170 - 150) = 6 + 6 + 15 + 47 + 4 + 21 + 135 + 20 = 254$  cylinders

Average head movements =  $\frac{254}{8} = 31.75$  cylinders

#### Advantages

1. It minimizes latency
2. Better throughput than FIFO method

#### Disadvantages

1. Starvation occurs if some process has to wait for long time until its requests are satisfied.
2. SSTF services requests for those tracks which are highly localized.

**SCAN/elevator algorithm** The disk head constantly moves from the most inner cylinder to the outer cylinder and then it changes its direction back towards the centre. As the head moves, if there is a request for the current disk position then it is satisfied.

1. It is known as *elevator algorithm* because it services all the request of going up and then reaching at the top, it goes downward.
2. It needs two information:
  - Direction of head movement
  - Last position of the disk head

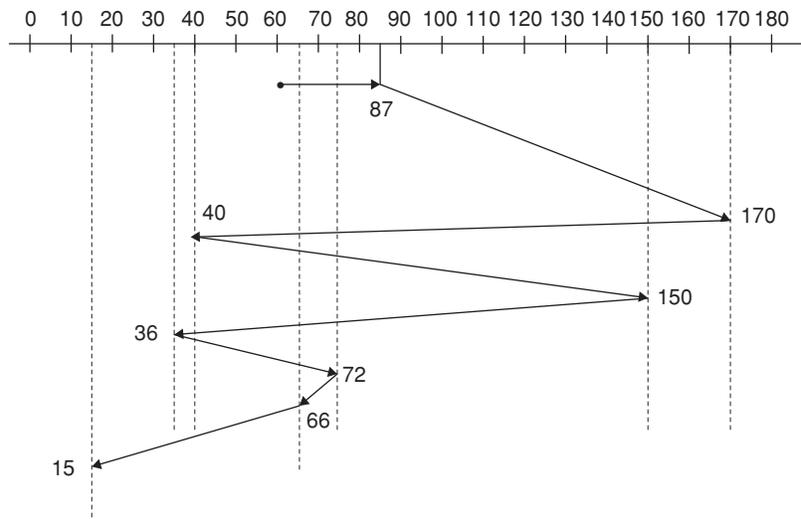


Figure 13 FCFS.

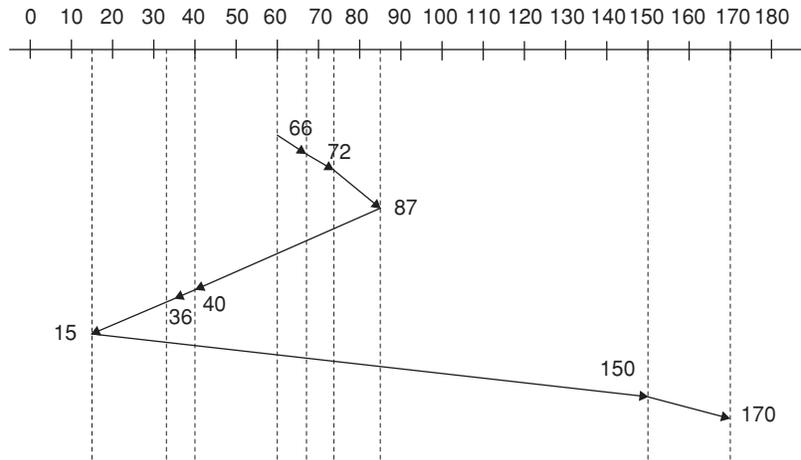


Figure 14 SSTF

**Example 5:** For the following track requests: 87, 170, 40, 150, 36, 72, 66, 15. (Initially head is at track 60 to the arm is moving outwards.

$$\begin{aligned} \text{Total head movement} &= (66 - 60) + (72 - 66) + (87 - 72) + (150 - 87) + (170 - 150) + (180 - 170) + (180 - 40) + (40 - 36) + (36 - 15) \\ &= 6 + 6 + 15 + 63 + 20 + 10 + 140 + 4 + 21 \\ &= 285 \text{ cylinders} \end{aligned}$$

$$\text{Average head movement} = \frac{285}{8} = 35.6 \text{ cylinders}$$

**Advantages**

1. Throughput better than FIFO.
2. Basic for most scheduling algorithms.
3. Eliminates the discrimination.
4. No starvation.

**Disadvantages**

1. Because of the continuous scanning of disk from end to end, the outer tracks are visited less often than the mid-range tracks.

2. Disk arm keeps scanning between two extremes; this may result in wear and tear of the disk assembly.
3. Certain requests arriving ahead of the arm position would get immediate service but some other requests that arrive behind the arm position will have to wait for the arm to return back.

**C-SCAN algorithm (one-way elevator algorithm)** It treats the cylinder as a circular list. The head sweeps from the innermost cylinder to the outermost cylinder, satisfying the waiting requests in order of their locations. When it reaches the outermost cylinder, it sweeps back to the innermost cylinder without satisfying any requests and then starts again.

**Example 6:** Consider the cylinders requests: 87, 170, 40, 150, 36, 72, 66, 15 Starting cylinder = 60th (arm moving outwards)

$$\begin{aligned} \text{Total head movement} &= (66 - 60) + (72 - 66) + (87 - 72) + (150 - 87) + (170 - 150) + (180 - 170) + (180 - 0) + (15 - 0) + (36 - 15) + (40 - 36) \\ &= 6 + 6 + 15 + 63 + 20 + 10 + 180 + 15 + 21 + 4 = 340 \end{aligned}$$

$$\text{Average head movement} = \frac{340}{8} = 42.5$$

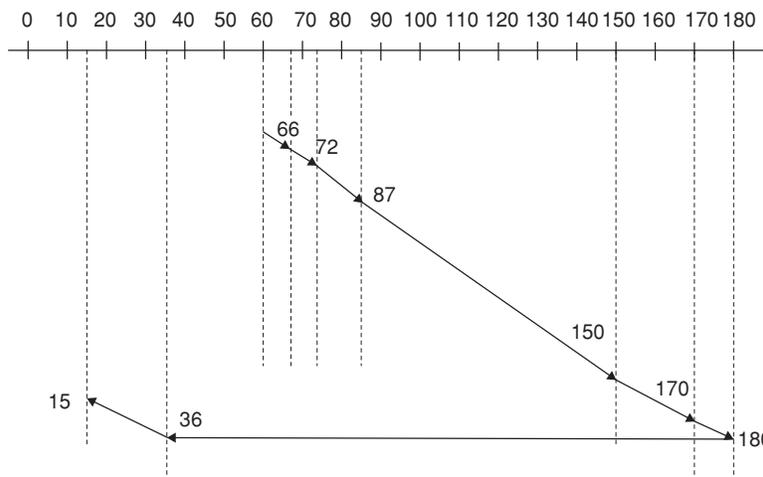


Figure 15 SSTF

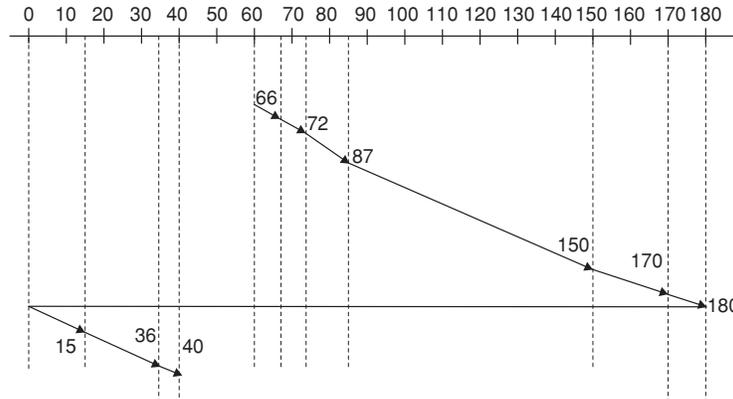


Figure 16 C-SCAN

**Advantage**

Lower service variability.

**Disadvantages**

1. An average head movement is more compared to SCAN algorithm.
2. Increase in the total seek time

**LOOK/SEEK algorithm** Look is similar to SCAN, but stops moving inwards (or) outwards when there are no more requests in that direction.

**Example 7:** Consider the following requests: 87, 170, 40, 150, 36, 72, 66, 15 initially head is at track 60 (moves outwards).

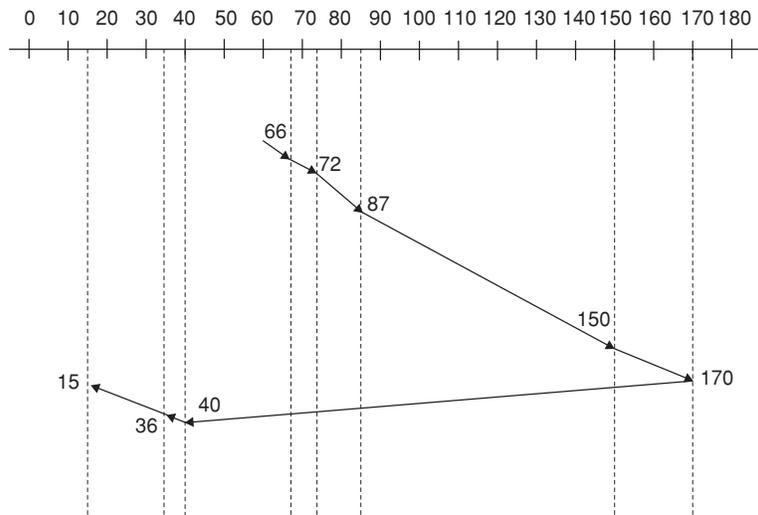


Figure 17 LOOK

$$\text{Average head movement} = \frac{275}{8} = 34.37$$

**C-LOOK/C-SEEK algorithm** The head moves inwards servicing requests until there are no more requests in that direction. Then it jumps to the outer most outstanding request.

**Magnetic Disk and Factors that Determine Access Speed**

The factors that determine the access speed are rotational speed, which is measured in revolutions per minute are seek time, the time taken to read or write the particular sector of the disk.

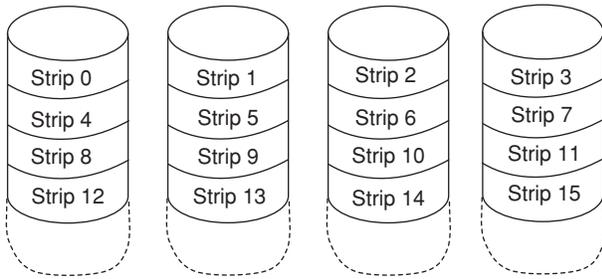
The other factors are sequential read, sequential write, random read and random write. These vary enormously, but for spinning disks one can expect 25 Mbps to 150 Mbps for

sequential read and write it is about 3 Mbps to 50 Mbps for the random read and write.

**RAID (REDUNDANT ARRAY OF INDEPENDENT DISKS)**

1. It is multiple disk database design.
2. It consists of seven levels, zero through six.
3. Characteristics of RAID Levels:
  - RAID is a set of physical disk drives viewed by the OS as a single logical drive.
  - Data are distributed across the physical drives of an array in a scheme known as striping.
  - Redundant disk capacity is used to store parity information, which guarantees data recoverability in case of a disk failure.

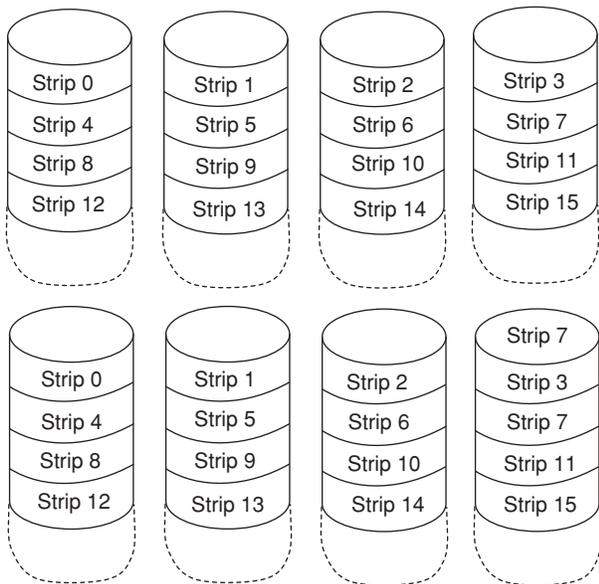
**RAID Level 0:**



**Figure 18** Non-redundant (RAID0)

1. It does not include redundancy.
2.  $N$  disks are required.
3. Data available in RAID level 0 is lower than single disk.
4. It has high data transfer capacity.
5. It has high I/O request rate.

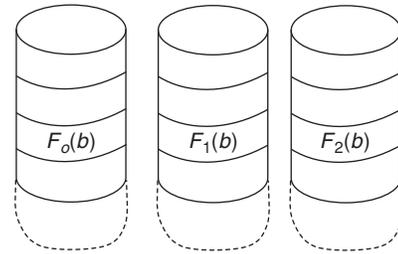
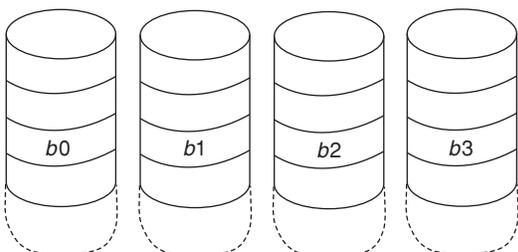
**RAID Level 1:**



**Figure 19** RAID 1 (Mirrored)

1. Redundancy is achieved by the simple expedient of duplicating all the data.
2.  $2N$  disks required.
3. Data availability is higher than RAID 2, 3, 4, or 5. But lower than RAID 6.
4. Recovery from failure is simple.
5. RAID1 costs more.

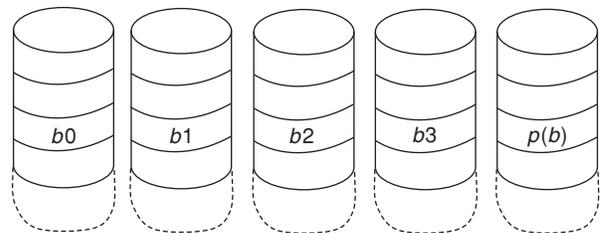
**RAID Level 2:**



**Figure 20** RAID 2 (Redundancy through Hamming code).

6. Here redundancy is achieved through hamming code.
7.  $N + m$  disks required.

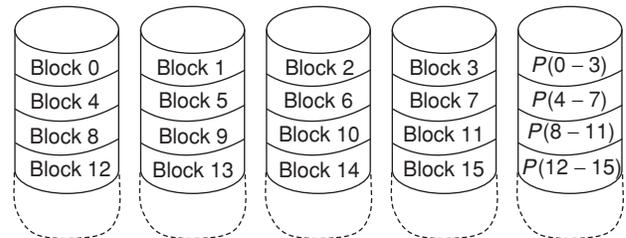
**RAID Level 3:**



**Figure 21** RAID 3 (bit-interleaved parity)

1. It has bit interleaved parity.
2. Provides parallel access.
3.  $N + 1$  disks required.

**RAID Level 4:**



**Figure 22** RAID 4 (Block-level parity)

1.  $N + 1$  disks required.
2. Provides independent access.
3. Consider an array of five drives in which  $X_0$  through  $X_3$  contains data and  $X_4$  is the parity disk.

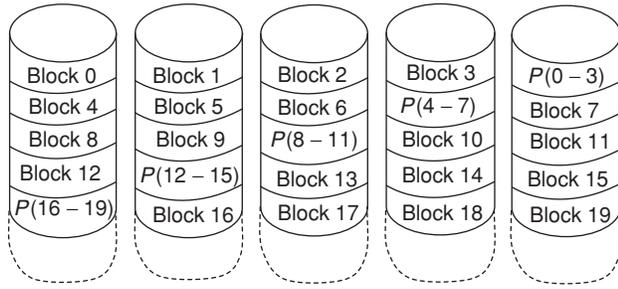
Suppose that a write is performed that only involves a strip on disk  $X_1$ . Initially for each bit  $i$ , we have the following relationship:

$$X_4(i) = X_3(i) \oplus X_2(i) \oplus X_1(i) \oplus X_0(i)$$

After the update, with potentially altered bits indicated by a prime symbol:

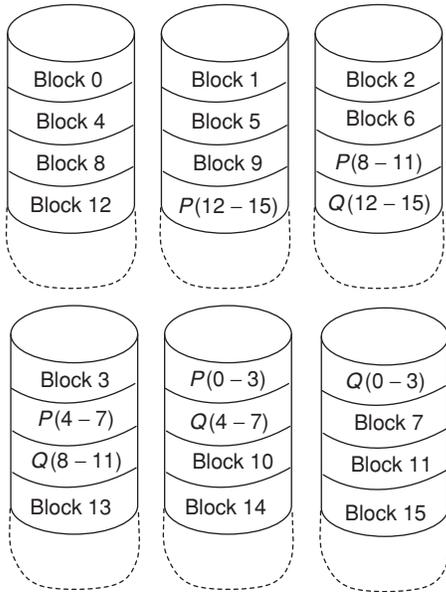
$$\begin{aligned} X_4'(i) &= X_3(i) \oplus X_2(i) \oplus X_1'(i) \oplus X_0(i) \\ &= X_3(i) \oplus X_2(i) \oplus X_1'(i) \oplus X_0(i) \oplus X_1(i) \oplus X_1(i) \\ &= X_4(i) \oplus X_1(i) \oplus X_1'(i) \end{aligned}$$

**RAID Level 5:**



1. It has block level distributed parity.
2. Provides independent access.
3. It has  $N + 1$  disks.

**RAID Level 6:**

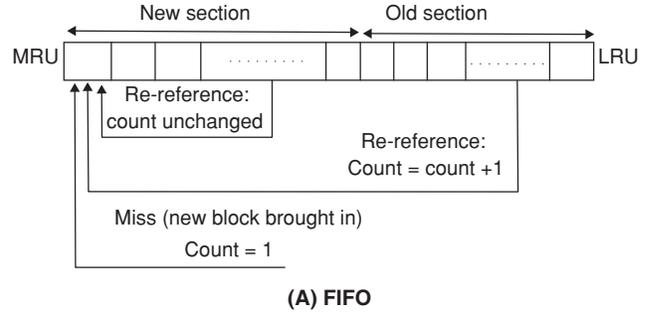


1. It has dual redundancy.
2. It provides independent access.
3. It has  $N + 2$  disks.

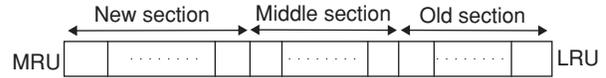
**Disk Cache:**

1. It is a buffer in main memory for disk sectors.
2. The cache contains a copy of some of the sectors on the disk.
3. Replacement policy:
  - LRU
  - LFU

But these two replacement policies lead to poor performance. So a new technique frequency-based replacement is proposed. Two alternatives of frequency-based technique are shown below:



(A) FIFO



(B) Use of three sections.

In FIFO, the blocks are logically organized in a stack; where the top part of the stack is the new section. When there is a cache hit, the referenced block is moved to top of the stack. If the block was already in the new section, its reference count is not incremented; otherwise it is incremented by 1.

In another technique, we divide the stack into three sections: New, middle and old. Here only blocks in the old section are eligible for replacement.

**PROTECTION AND SECURITY**

**System protection** Protection refers to a mechanism for controlling the access of programs, processes or users to the resources defined by a computer system.

**Principles of Protection**

Programs, users and even systems are given just enough privileges to perform their tasks. This is the principle of *least privilege*.

**Domain of protection** A computer system is a collection of processes and objects. The objects may be software or hardware objects. The operations that are possible may depend on the object.

A process should be allowed to access only those resources for which it has authorization. At any time, a process should be able to access only those resources that it currently requires to complete its task. This requirement is referred as *need to know principle* and is useful in limiting the amount of damage, which a faulty process can cause in the system.

**Domain structure** Each domain defines a set of objects and the types of operations that may be invoked on each object.

The ability to execute an operation on an object is an *access right*. A domain is a collection of access rights, each of which is an ordered pair <obj-name, rights-set>

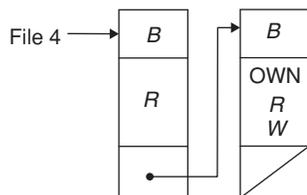
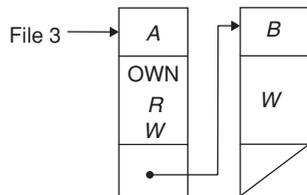
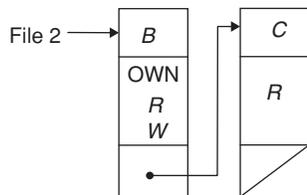
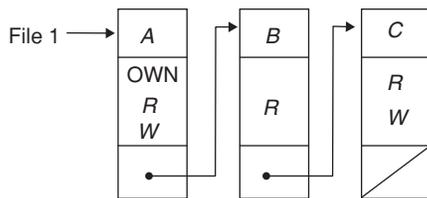
### File System Security

1. A general model of access control for file management is an access matrix.
2. The basic elements of the model are as follows:
  - *Subject*: An entity capable of accessing objects.
  - *Object*: Anything to which access is controlled.
  - *Access right*: The way in which an object is accessed by a subject.

Access matrix format is shown below:

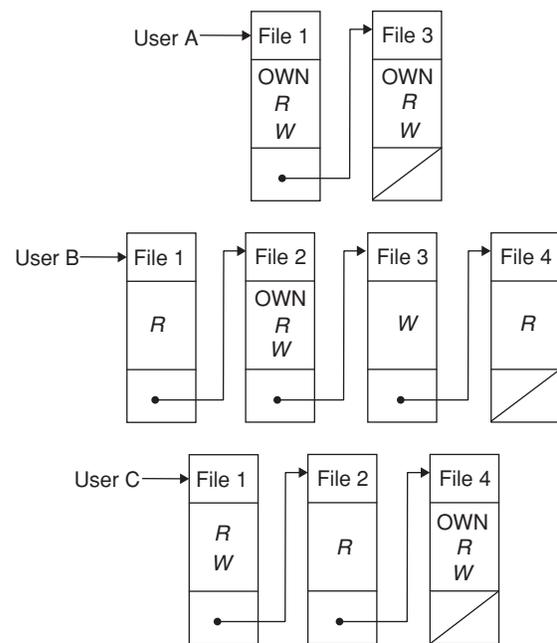
|        | File 1        | File 2        | File 3        | File 4        | Account 1         | Account 2         |
|--------|---------------|---------------|---------------|---------------|-------------------|-------------------|
| User A | OWN<br>R<br>W |               | OWN<br>R<br>W |               | Inquiry<br>credit |                   |
| User B | R             | OWN<br>R<br>W | W             | R             | Inquiry<br>debit  | Inquiry<br>credit |
| User C | R<br>W        | R             |               | OWN<br>R<br>W |                   | Inquiry<br>Debit  |

The access matrix may be decomposed by columns, yielding access control lists. The access control list for above matrix is



**Note:** This allows users that are not explicitly listed as having special rights to have a default set of rights.

Decomposition by row yields. Capability tickets.



**Note:** These tickets would have to be held in a region of memory inaccessible to users.

**System security** This is the protection afforded to an automated information system in order to attain objectives of preserving the integrity, availability and confidentiality of information system resources.

### Objectives of Computer Security

**Confidentiality** Preserving authorized restrictions on information access and disclosure, including means for protecting personal privacy and proprietary information.

**Integrity** Guarding against improper information modification or destruction, including ensuring information non-repudiation and authenticity.

**Availability** Ensuring timely and reliable access to and use of information.

**Authenticity** The property of being genuine and being able to be verified and trusted.

**Accountability** The security goal that generates the requirement for actions of an entity to be traced uniquely to that entity.

### Threats, Attacks and Assets

**Threats and attacks** *Unauthorized disclosure* is the threat to confidentiality. The attacks of following this threat are as follows:

1. *Exposure*: Sensitive data are directly released to an unauthorized entity.
2. *Interception*: An unauthorized entity directly accesses sensitive data traveling between authorized sources and destinations.

3. *Inference*: A threat action where an unauthorized entity indirectly accesses sensitive data by reasoning from characteristics or by-products of communications.
4. *Intrusion*: An unauthorized entity gains access to sensitive data by circumventing a system's security protections.

**Deception** Threat to either system integrity or data integrity. Types of attacks that can result are as follows:

1. *Masquerade*: An unauthorized entity gains access to a system or performs a malicious act by posing as an authorized entity.
2. *Falsification*: False data deceive an authorized entity.
3. *Repudiation*: An entity deceives another by falsely denying responsibility for an act.

**Disruption** A circumstance or event that interrupts or prevents the correct operation of system services and functions. Attacks for this threat are as follows:

1. *Incapacitation*: Prevents or interrupts system operation by disabling a system component.
2. *Corruption*: Undesirably alters system operation by adversely modifying system functions or data.
3. *Obstruction*: A threat action that interrupts delivery of system service by hindering system operation.

**Usurpation** A circumstance or event that results in control of system services or functions by an unauthorized entity. Attacks with this threat are as follows:

1. *Misappropriation*: An entity assumes unauthorized logical or physical control of a system resource.
2. *Misuse*: Causes a system component to perform a function or service that is detrimental to system security.

**Threats and assets** The assets of a computer are as follows:

1. Hardware
2. Software
3. Data
4. Communication lines

**Hardware** A major threat to computer system hardware is the threat to availability (e.g., theft of CD-ROMS).

#### Software

1. A key threat to software is an attack on availability (e.g., deletion of software).
2. A threat to integrity.
3. A threat to confidentiality.

**Data** Threats to data are an attack on

1. Availability
2. Confidentiality
3. Integrity

**Communication lines and networks** Two types of attacks:

1. Passive attacks
2. Active attacks

#### Passive attacks

1. These are in the nature of monitoring of the transmissions.
2. Attackers obtain information that is being transmitted.
3. Two types of passive attacks:
  - Release of message contents.
  - Traffic analysis
4. These are very difficult to detect because they do not involve any alteration of the data.

#### Active attacks

1. These attacks involves some modification of the data stream or the creation of a false stream and can be subdivided into four categories:
  - Replay
  - Masquerade
  - Modification of messages
  - Denial of service
2. It is difficult to prevent active attacks absolutely.

**Intruders** Three types of intruders:

**Masquerader** An individual who is not authorized to use the computer and who penetrates a systems access controls to exploit a legitimate user's account. He is likely to be an outsider.

**Misfeasor** A legitimate user who accesses data, programs or resources for which such access is not authorized or who is authorized for such access but misuses his or her privileges (generally an insider).

**Clandestine user** An individual who seizes supervisory control of the system and uses this control to evade auditing and access controls or to suppress audit collection (either outsider or insider).

**Hackers** Those who hack into computers do so for the thrill of it or for status. Attackers often took for targets of opportunity and then share the information with others.

**Criminals** Organized group of hackers have become a widespread and common threat to internet based systems.

**Malicious software overview** The most sophisticated types of threats to computer systems are presented by programs that exploit vulnerabilities in computing systems. These threats are referred as malicious software (or) malware.

1. It is designed to cause damage to or use up the resources of a target computer.
2. There are two types of malicious software:
  - Those that need a host program.  
**Example:** Viruses, logic bombs.
  - Those that are independent.  
**Example:** Worms, bot programs.
3. We can also differentiate between two types of software threats:
  - That do not replicate. These programs are activated by a trigger.  
**Example:** Logic bombs, backdoors.

- Those that replicate.

**Example:** Viruses, worms

**Backdoor (trapdoor)** It is a secret entry point into a program that allows someone who is aware of the backdoor to gain access without going through usual security access procedures.

**Logic bomb** It is a program inserted into software by an intruder. It lies dormant until a predefined condition is met; the program then triggers an unauthorized act.

**Trojan Horse** It is an apparently useful program, containing hidden code that, when invoked, performs some unwanted or harmful function.

**Mobile code** It is a software that can be shipped unchanged to a heterogeneous collection of platforms and execute with identical semantics.

**Viruses** A computer virus is a piece of software that can infect other programs by modifying them.

### Nature of Viruses

A virus can do anything that other programs do. It attaches itself to another program and executes secretly when the host program is run. Three parts of computer virus are as follows:

1. Infection mechanism
2. Trigger
3. Payload

Phases of computer virus are

1. Dormant phase
2. Propagation phase
3. Triggering phase
4. Execution phase

### Types of Virus

1. Encrypted virus
2. Stealth virus
3. Polymorphic virus
4. Metamorphic virus

**Worms** A worm is a program that can replicate itself and send copies from computer to computer across network connections.

**Bots** A bot is a program that secretly takes over another internet attached computer and then uses that computer to launch attacks that are difficult to trace to the bot's creator.

## EXERCISES

### Practice Problems I

**Directions for questions 1 to 15:** Select the correct alternative from the given choices.

1. Given a system using unspanned blocking and 100 byte blocks. A file contains records 30, 40, 55, 80, 30, 40. What percentage of space will be wasted in the blocks allocated for the file?  
(A) 31.25% (B) 41.25%  
(C) 51.25% (D) 62.15%
2. Disk requests come into the disk driver for cylinders 15, 25, 10, 2, 35, 9, 42 in that order. The disk head is currently positioned over cylinder 15. A seek takes 6 msec per cylinder moved. What is the total seek time using First Come First Served Algorithm?  
(A) 750 msec (B) 650 msec  
(C) 550 msec (D) 450 msec
3. A Java application needs to load 50 libraries. To load each library, one disk access is required. Seek time to access the location is 10 ms. Rotational speed is 6000 rpm. The total time needed to load all libraries is  
(A) 0.65 sec (B) 0.75 sec  
(C) 0.85 sec (D) 1 sec
4. A program has just read the 13th record in a sequential access file. If it wants to read the 10th record next, how many records must the program read to input the tenth record?

- (A) 5 (B) 0  
(C) 10 (D) 13

5. A disk is formatted into 40 sectors and 20 tracks. The disk rotates at 200 ms in one revolution. The time taken by the head to move from the centre to the rim is 10 ms. There are three different files stored on the disk:  
File *P* : Sector 2, track 4  
File *Q* : Sector 5, track 1  
File *R* : Sector 6, track 2  
Calculate the average latency time required for the three files.  
(A) 22.55 ms (B) 32.22 ms  
(C) 21.66 ms (D) 30.22 ms
6. Match the following  
(a) RAID0 (1) Parallel access  
(b) RAID1 (2) Striping  
(c) RAID2 (3) Use hamming code  
(d) RAID3 (4) Mirrored  
(A) a – 2, b – 4, c – 3, d – 1  
(B) a – 1, b – 2, c – 3, d – 4  
(C) a – 3, b – 2, c – 4, d – 1  
(D) a – 4, b – 1, c – 2, d – 3
7. The correct matching for the following pairs is  
(A) Disk scheduling (1) Round Robin  
(B) Batch processing (2) SCAN  
(C) Time sharing (3) LIFO  
(D) Interrupt processing (4) FIFO

- (A) A – 3, B – 4, C – 2, D – 1  
 (B) A – 4, B – 3, C – 2, D – 1  
 (C) A – 2, B – 4, C – 1, D – 3  
 (D) A – 2, B – 1, C – 4, D – 3
8. A program  $P$  reads and processes 2000 consecutive records from a sequential file stored on device  $R$  without using any file system facilities. Given the following:  
 Size of each record = 3500 bytes  
 Access time of  $R = 20$  ms  
 Data transfer rate of  $R = 500 \times 10^3$  bytes/sec  
 CPU time to process each record = 5 m sec  
 What is the elapsed time of  $P$  if  $R$  contains unblocked records and  $P$  does not use buffering?  
 (A) 64 sec (B) 46 sec  
 (C) 34 sec (D) 17 sec
9. A disk has 19456 cylinders, 16 heads and 63 sectors per track. The disk spins at 5400 rpm. Seek time between adjacent tracks is 2 ms. Assuming the read/write head is already positioned at track 0, how long does it take to read the entire disk?  
 (A) 35 min (B) 68 min  
 (C) 58 min (D) 53 min
10. On a disk with 1000 cylinders, numbered 0 to 999, compute the number of tracks the disk arm must move to satisfy all the requests in the disk queue using SCAN algorithm. Assume the last request serviced was at track 345 and the head is moving towards 0. The queue in FIFO order contains requests for the following tracks:  
 123, 874, 692, 475, 105, 376  
 (A) 219 (B) 635  
 (C) 845 (D) 1219

**Common data for questions 11 to 14:** In the operation of a certain disk drive mechanism, a disk is formatted into 20 sectors and 10 tracks. The disk can be rotated either

clockwise or anti-clockwise. The times required to perform certain operations are as follows:

- I. Rotate the disk through one revolution = 200 ms  
 II. Move the disk head from the centre to the rim = 20 ms  
 III. Read and transmit one block of data = 0.3 ms
- Three files are stored on the disk:  
 File  $A$ : 2 blocks at track 6  
 File  $B$ : 5 blocks at track 2  
 File  $C$ : 1 block at track 5
11. The disk head is initially at sector 0, track 0. If all three files  $A$ ,  $B$  and  $C$  are to be read in the minimum amount of time, they should be read in the following order:  
 (A)  $A, B, C$  (B)  $A, C, B$   
 (C)  $B, C, A$  (D)  $C, A, B$
12. The disk head is initially at sector 0, track 0. The files are read in the order  $C, B, A$ . The total time to read the files is  
 (A) 143.9 ms (B) 100.4 ms  
 (C) 114.0 ms (D) 102.6 ms
13. The most nearly average latency time for the sequence  $CBA$  is:  
 (A) 1 ms (B) 7 ms  
 (C) 27 ms (D) 50 ms
14. The most nearly average seek time for  $CBA$  is  
 (A) 1 ms (B) 8 ms  
 (C) 30 ms (D) 50 ms
15. A CD has 150 tracks rotating at 3500 rpm. Average seek time for consecutive tracks is 0.1 ms, the disk is subjected to read data from the track numbers 89, 75, 112, 5. What is the total seek time if the requests are served unidirectionally (C-Scan) and the first request determines initial direction? Assume that the current position of the head is at track 100.  
 (A) 20.2 ms (B) 16.9 ms  
 (C) 21.3 ms (D) 14 ms

## Practice Problems 2

**Directions for questions 1 to 15:** Select the correct alternative from the given choices.

1. The strategy that allocates the smallest possible chunk of disk space that is sufficient to the file is  
 (A) Nearest fit (B) Best fit  
 (C) Worst fit (D) First fit
2. If a process of 200 kB is transferred from backing store to memory and average disk latency is 10 ms, then what would be the total swap time, if transfer ratio is 2 Mbps?  
 (A) 10 ms  
 (B) 20 ms  
 (C) 30 ms  
 (D) 40 ms
3. Let us assume that the user process is 10 MB in size and backing store is a standard hard disk with a transfer rate of 40 MB per second. Let the average latency is 8 millisecond. Find actual time transfer of the 10 MB process to or from main memory?  
 (A) 8 ms (B) 250 ms  
 (C) 258 ms (D) 516 ms
4. Disk scheduling involves deciding  
 (A) which disk should be accessed next  
 (B) the order in which disk access requests must be serviced.  
 (C) the physical location when files should be accessed in the disk  
 (D) disk access time and an unused space.

5. The root directory of a disk should be placed
  - (A) at a fixed address in the memory
  - (B) anywhere on the disk
  - (C) at a fixed location on the system disk
  - (D) at a location on floppy.
6. Direct access methods are not effectively supported by
  - (A) Contiguous allocation
  - (B) Linked allocation
  - (C) Indexed allocation
  - (D) Sequential allocation
7. In which of the following directory systems, it is possible to have multiple paths for a file, starting from the root directory?
  - (A) Single-level directory
  - (B) Two-level directory
  - (C) Tree-structured directory
  - (D) A cyclic graph directory
8. The most common system's security method is:
  - (A) Passwords
  - (B) Key card systems
  - (C) Surveillance system
  - (D) Lock system
9. Trojan Horse programs
  - (A) are legitimate programs that allow unauthorized access.
  - (B) are hacker programs that do not show up on the system
  - (C) really do not work
  - (D) are immediately discovered
10. Which of the following is a program that spreads throughout the network?
  - (A) Trojan Horse
  - (B) Virus
  - (C) TSR
  - (D) Worm
11. A program has just read the 15th record in a sequential access file. If it wants to read the 10th record next, how many records must the program read to input the tenth record?
  - (A) 0
  - (B) 5
  - (C) 4
  - (D) 10
12. Formatting of a floppy disk refers to
  - (A) Arranging the data on the disk in contiguous fashion
  - (B) Writing the directory
  - (C) Erasing the system area
  - (D) Writing identification information on all tracks
13. Sector interleaving in disks is done by
  - (A) the disk manufacturer
  - (B) the disk controller card
  - (C) the operating system
  - (D) the user
14. Disk I/O is done in terms of
  - (A) Tracks
  - (B) Blocks
  - (C) Bits
  - (D) Bytes
15. How many six-letter passwords can be constructed using lowercase letters and digits?
  - (A)  $26^6$
  - (B)  $10^6$
  - (C)  $36^6$
  - (D)  $35^6$

### PREVIOUS YEARS' QUESTIONS

1. Consider a disk drive with the following specifications: [2005]  
 16 surfaces, 512 tracks/surfaces, 512 sectors/track, 1 KB/sector, rotation speed 3000 rpm. The disk is operated in cycle stealing mode whereby whenever one 4 byte word is ready it is sent to memory; similarly, for writing, the disk interface reads a 4 byte word from the memory in each DMA cycle. Memory cycle time is 40 nsec. The maximum percentage of time that the CPU gets blocked during DMA operation is:
  - (A) 10
  - (B) 25
  - (C) 40
  - (D) 50
2. Consider a disk pack with 16 surfaces, 128 tracks per surface and 256 sectors per track. 512 bytes of data are stored in a bit serial manner in a sector. The capacity of the disk pack and the number of bits required to specify a particular sector in the disk are, respectively: [2007]
  - (A) 256 Mbyte, 19 bits
  - (B) 256 Mbyte, 28 bits
  - (C) 512 Mbyte, 20 bits
  - (D) 64 Gbyte, 28 bits
3. For a magnetic disk with concentric circular tracks, the seek latency is not linearly proportional to the seek distance due to [2008]
  - (A) non-uniform distribution of requests
  - (B) arm starting and stopping inertia
  - (C) higher capacity of tracks on the periphery of the platter
  - (D) use of unfair arm scheduling policies
4. Consider a file of 16384 records. Each record is 32 bytes long and its key field is of size 6 bytes. The file is ordered on a non-key field, and the file organization is unspanned. The file is stored in a file system with block size 1024 bytes, and the size of a block pointer is 10 bytes. If the secondary index is built on the key field of the file, and a multi-level index scheme is used to store the secondary index, the number of first-level and second-level blocks in the multi-level index are, respectively [2008]

- (A) 8 and 0                      (B) 128 and 6  
(C) 256 and 4                    (D) 512 and 5

5. Consider a disk system with 100 cylinders. The requests to access the cylinders occur in following sequence: [2009]

4, 34, 10, 7, 19, 73, 2, 15, 6, 20

Assuming that the head is currently at cylinder 50, what is the time taken to satisfy all requests if it takes 1ms to move from one cylinder to adjacent one and shortest seek time first policy is used?

- (A) 95 ms                          (B) 119 ms  
(C) 233 ms                        (D) 276 ms

**Common data for questions 6 and 7:** A hard disk has 63 sectors per track, 10 platters each with 2 recording surfaces and 1000 cylinders. The address of a sector is given as a triple  $\langle c, h, s \rangle$ , where  $c$  is the cylinder number,  $h$  is the surface number and  $s$  is the sector number. Thus, the 0th sector is addressed as  $\langle 0, 0, 0 \rangle$ , the 1st sector as  $\langle 0, 0, 1 \rangle$ , and so on

6. The address  $\langle 400, 16, 29 \rangle$  corresponds to the sector number: [2009]

- (A) 505035                        (B) 505036  
(C) 505037                        (D) 505038

7. The address of the 1039th sector is [2009]

- (A)  $\langle 0, 15, 31 \rangle$                     (B)  $\langle 0, 16, 30 \rangle$   
(C)  $\langle 0, 16, 31 \rangle$                     (D)  $\langle 0, 17, 31 \rangle$

8. A file system with 300 GByte disk uses a file descriptor with 8 direct block addresses, 1 indirect block address and 1 doubly indirect block address. The size of each disk block is 128 bytes and the size of each disk block address is 8 bytes. The maximum possible file size in this file system is [2012]

- (A) 3 Kbytes  
(B) 35 Kbytes  
(C) 280 Kbytes  
(D) dependent on the size of the disk

9. Consider a hard disk with 16 recording surfaces (0–15) having 16384 cylinders (0–16383) and each cylinder contains 64 sectors (0–63). Data storage capacity in each sector is 512 bytes. Data are organized cylinder-wise and the addressing format is  $\langle \text{cylinder no.}, \text{surface no.}, \text{sector no.} \rangle$ . A file of size 42797 KB is stored in the disk and the starting disk location of the file is  $\langle 1200, 9, 40 \rangle$ . What is the cylinder number of the last sector of the file, if it is stored in a contiguous manner? [2013]

- (A) 1281                            (B) 1282  
(C) 1283                            (D) 1284

10. A FAT (File allocation table)-based file system is being used, and the total overhead of each entry in

the FAT is 4 bytes in size. Given a  $100 \times 10^6$  bytes disk on which the file system is stored and data block size is  $10^3$  bytes, the maximum size of a file that can be stored on this disk in units of  $10^6$  bytes is \_\_\_\_\_?

[2014]

11. Consider a disk pack with a seek time of 4 milliseconds and rotational speed of 10000 rotations per minute (RPM). It has 600 sectors per track and each sector can store 512 bytes of data. Consider a file stored in the disk. The file contains 2000 sectors. Assume that every sector access necessitates a seek, and the average rotational latency for accessing each sector is half of the time for one complete rotation. The total time (in milliseconds) needed to read the entire file is \_\_\_\_\_.

[2015]

12. Suppose the following disk request sequence (track numbers) for a disk with 100 tracks is given: 45, 20, 90, 10, 50, 60, 80, 25, 70. Assume that the initial position of the R/W head is on track 50. The additional distance that will be traversed by the R/W head when the Shortest Seek Time First (SSTF) algorithm is used compared to the SCAN (Elevator) algorithm (assuming that SCAN algorithm moves towards 100 when it starts execution) is \_\_\_\_\_ tracks. [2015]

13. Consider a typical disk that rotates at 15000 rotations per minute (RPM) and has a transfer rate of  $50 \times 10^6$  bytes/sec. If the average seek time of the disk is twice the average rotational delay and the controller's transfer time is 10 times the disk transfer time, the average time (in milliseconds) to read or write a 512-byte sector of the disk is \_\_\_\_\_ [2015]

14. Consider a disk queue with requests for I/O to blocks on cylinders 47, 38, 121, 191, 87, 11, 92, 10. The C-LOOK scheduling algorithm is used. The head is initially at cylinder number 63, moving towards larger cylinder numbers on its servicing pass. The cylinders are numbered from 0 to 199. The total head movement (in number of cylinders) incurred while servicing these requests is \_\_\_\_\_. [2016]

15. In a file allocation system, which of the following allocation scheme(s) can be used if no external fragmentation is allowed? [2017]

- I. Contiguous  
II. Linked  
III. Indexed

- (A) I and III only                    (B) II only  
(C) III only                            (D) II and III only

16. Consider a storage disk with 4 platters (numbered as 0, 1, 2 and 3), 200 cylinders (numbered as 0, 1, ..., 199), and 256 sectors per track (numbered as 0, 1, ..., 255). The following 6 disk requests of the form [sector number, cylinder number, platter number] are

received by the disk controller at the same time:

[120, 72, 2], [180, 134, 1], [60, 20, 0], [212, 86, 3],  
[56, 116, 2], [118, 16, 1]

Currently the head is positioned at sector number 100 of cylinder 80, and is moving towards higher cylinder numbers. The average power dissipation in moving the head over 100 cylinders is 20 milliwatts and for

reversing the direction of the head movement once is 15 milliwatts. Power dissipation associated with rotational latency and switching of head between different platters is negligible.

The total power consumption in milliwatts to satisfy all of the above disk requests using the Shortest Seek Time First disk scheduling algorithm is \_\_\_\_\_.

[2018]

## ANSWER KEYS

### EXERCISES

#### Practice Problems 1

1. A    2. A    3. B    4. C    5. C    6. A    7. C    8. A    9. C    10. D  
11. C    12. B    13. C    14. B    15. A

#### Practice Problems 2

1. B    2. C    3. D    4. B    5. C    6. D    7. C    8. A    9. A    10. D  
11. D    12. D    13. C    14. B    15. C

#### Previous Years' Questions

1. B    2. A    3. C    4. C    5. B    6. C    7. C    8. B    9. D    10. 99.6  
11. 14020    12. 10    13. 6.1 to 6.2    14. 346    15. D    16. 85

## TEST

## OPERATING SYSTEM

Time: 60 min.

**Directions for questions 1 to 30:** Select the correct alternative from the given choices.

1. A program is \_\_\_\_\_ entity, while a process is \_\_\_\_\_ entity.
  - (A) Active, passive
  - (B) Active, sometimes active
  - (C) Passive, active
  - (D) Both (B) and (C)
2. All the information associated with a specific process is contained in:
  - (A) Process control block
  - (B) Program control block
  - (C) TLB
  - (D) Heap
3. Kernel-level threads and user-level threads are supported, respectively, by \_\_\_\_\_.
  - (A) Operating system and operating system
  - (B) Operating system and user
  - (C) User and user
  - (D) None of these
4. Which of the following is false about user-level threads?
  - (A) User-level threads are visible to the programmer and are unknown to the Kernel.
  - (B) These are faster to create.
  - (C) Kernel never interferes.
  - (D) There is no effect of a system call () on process.
5. Which of the following interprocess communication models are implemented using system calls?
  - (A) Shared memory
  - (B) Message passing
  - (C) Both (A) and (B)
  - (D) Neither (A) nor (B)
6. Peterson's solution
  - (i) is restricted to two processes
  - (ii) share two data items turn and flag [i]
  - (iii) mutual exclusion is achieved
  - (iv) is a hardware solution

Which of the above are true?

  - (A) (i), (ii), (iii)
  - (B) (ii), (iii), (iv)
  - (C) (iv), (i), (ii)
  - (D) (i), (ii), (iii), (iv)
7. Which of the following requires a mode switch from one thread to another?
  - (A) One process multiple thread
  - (B) User-level thread
  - (C) Kernel-level threads
  - (D) Both (B) and (C)
8. When a process is created, its state is
  - (A) New
  - (B) Ready
  - (C) Block
  - (D) Suspend
9. The data section of a process in memory contains
  - (A) Local variables, function parameters
  - (B) Return addresses
  - (C) Global variables
  - (D) None of the above
10. Which one of the following is true about process states?
  - (i) A process which is running must have terminated as next state.
  - (ii) From running state process can go to either waiting, ready or terminated state.
  - (iii) Only one process can run at any instant.
  - (iv) Ready process can go to waiting state.
  - (A) (i), (ii), (iii)
  - (B) (ii) and (iii) only
  - (C) (i) and (iv) only
  - (D) (ii), (iii), (iv)
11. Message passing model of inter process communication can be
  - (A) Blocking only
  - (B) Blocking and non-blocking
  - (C) Synchronous and asynchronous
  - (D) Both (B) and (C)
12. The definition of wait( ) is as follows:
 

```
wait (S) {
while (S <=0);
S - - ;
}
```

The semicolon after while statement, signifies

  - (A) Infinite looping
  - (B) Blank statement
  - (C) Depends on interpretation of compiler
  - (D) No operation
13. To avoid race condition, the number of processes using the critical sections is/are:
  - (A) 1
  - (B) 2
  - (C) 3
  - (D) More than 3
14. The 'Critical Section' is the region in which
  - (A) Any number of processes can enter without any permission
  - (B) Only one process enters at a time and others wait for it.
  - (C) Section is very critical
  - (D) None of these

15. What does process control block contain?

- (A) Process Identification
- (B) Process state information
- (C) Process control information
- (D) All of the above

16. Match the following

|                                      |                                                                    |
|--------------------------------------|--------------------------------------------------------------------|
| (i) Multiprogramming                 | (x) Managing multiple processes executing on multiple computers    |
| (ii) Multiprocessing                 | (y) Management of multiple processes within a uniprocessor system. |
| (iii) Distributed process Management | (z) Management of multiple processes within a multiprocessor.      |

- (A) (i) -y (ii) -z (iii) -x
- (B) (i) -z (ii) -x (iii) -y
- (C) (i) -y (ii) -x (iii) -z
- (D) Ambiguous

17. For 'n' number of fork() system call, how many parent and child processes will be created?

- (A) 1,  $2^n - 1$ , respectively
- (B) 1,  $2^n$ , respectively
- (C)  $2^n - 1$ , 1, respectively
- (D) n,  $2n$ , respectively

18. If the value of binary semaphore is initialized with 1 and three wait() operations are performed, how many processes are there in the block list?

- (A) 1
- (B) 0
- (C) 3
- (D) 2

19. A counting semaphore is initialized with the value 3. A list of 'P' and 'V' operations are performed on the semaphore as: 1P, 2V, 2P, 3V, 5P, 7V, 2P, 3V. The final value of semaphore is?

- (A) 5
- (B) 8
- (C) 7
- (D) 6

20. The final value of semaphore after 10 'P' operations and 23 'V' operations is 1. What will be the initial value of this counting semaphore?

- (A) -14
- (B) -13
- (C) -12
- (D) -11

21. For a machine-instruction approach to enforce mutual exclusion, following are its properties:

- (i) starvation and deadlock free
- (ii) it is applicable to any number of processes.
- (iii) it can be used to support multiple critical sections, each defined by its own variable.
- (iv) it is simple, easy to verify and employed with busy waiting

Which of the above is false?

- (A) (iv) only
- (B) (ii), (iii) only
- (C) (i), (ii) only
- (D) (i) only

22. Consider the following code:

```
if (fork() == 0)
{
a = a + 5;
printf("%d, %d\n", a, &a);
}
else
{
a = a - 5;
printf("%d, %d\n", a, &a);
}
```

Let  $p, q$  be the values printed by the parent process, and  $s, t$  be the values printed by the child process. Which one of the following is true?

- (A)  $p = s + 10$  and  $q = t$
- (B)  $p = s + 10$  and  $q \uparrow t$
- (C)  $p + 10 = s$  and  $q = t$
- (D)  $p + 10 = s$  and  $q \uparrow t$

23. Consider the following statements with respect to user-level threads:

- (i) Context switch is faster with kernel-supported threads.
- (ii) For user-level threads, a system call can block the entire process.
- (iii) Kernel-supported threads can be scheduled independently.
- (iv) User-level threads are transparent to the Kernel.

Which of the above statements are true?

- (A) (i), (iii) and (iv) only
- (B) (ii) and (iii) only
- (C) (i) and (iii) only
- (D) (i) and (ii) only

24. Suppose there are 'n' CPUs and 'm' processes such that  $m > n$ . What will be the minimum and maximum number of ready, running and blocked process, respectively?

- (A) 0, 0, 0 and  $m, n, m$
- (B) 1,  $m, 1$  and  $n, n, n$
- (C)  $m, 1, 0$  and  $m, m, n$
- (D) 0, 0, 0 and  $n, m, m$

25. Consider the following signal semaphore code signal (semaphore \*s)

```
{
s.value++;
if(_____ (I) _____)
{
remove a process P from S.list;
_____ (II) _____
}
}
```

Choose the suitable options for (I) and (II), respectively

- (A)  $S.value = 0$  and wakeup(P);
- (B)  $S.value \leq 0$  and wakeup(P);
- (C)  $S.value < 0$  and block( );
- (D)  $S.value \leq 0$  and block( );

26. Consider the methods used by processes  $P_1$  and  $P_2$  for accessing their critical sections whenever needed. The initial values of shared Boolean variables  $S_1$  and  $S_2$  are randomly assigned.

Method used by  $P_1$

While ( $S_1 == S_2$ );

Critical section

$S_1 = S_2$ ;

Method used by  $P_2$

While ( $S_1 != S_2$ );

Critical section

$S_2 = !(S_1)$ ;

Which of the following statements describes the properties achieved?

- (A) Mutual exclusion but not progress
- (B) Progress only
- (C) Bounded waiting, progress
- (D) Mutual exclusion, progress, bounded waiting

27. Consider the following statements regarding spin locks:

- (i) No context switch is required when a process wait on a lock
- (ii) Spin locks are useful when locks are expected to be held for short times.
- (iii) They are often employed on multiprocessor systems
- (iv) Process 'spins' while waiting for a lock

Choose the correct option:

- (A) (i), (ii), (iii), (iv) are true
- (B) Only (i) and (ii) are true
- (C) (iii) is false
- (D) (ii) is true and (iv) is false

**Common data for questions 28, 29 and 30:** From the Readers-Writers problem, the data structure for reader process is:

```
semaphore mutex, wrt;
int readcount;
```

```
while(1)
{
wait(mutex);
readcount++;
if(readcount == 1)
wait(wrt);
signal(mutex);
- - - - -
- - - - -
wait(mutex);
readcount - -;
if(readcount ==0)
signal(wrt); signal(mutex);
}
```

mutex and wrt are initialized to 1 and readcount is initialized to 0.

28. Mutual exclusion for readers is attained by

- (A) Wrt
- (B) Mutex
- (C) Readcount
- (D) Both (A) and (B)

29. Which of the following semaphore or semaphores is used by the first or last reader that enters or exits the critical section?

- (A) Wrt
- (B) Mutex
- (C) Readcount
- (D) Both (A) and (B)

30. The readcount variable keeps track of how many processes are \_\_\_\_\_.

- (A) Currently reading the object
- (B) Currently writing the object
- (C) Waiting in the queue
- (D) Reading the shared data

### ANSWER KEYS

- |       |       |       |       |       |       |       |       |       |       |
|-------|-------|-------|-------|-------|-------|-------|-------|-------|-------|
| 1. C  | 2. A  | 3. B  | 4. D  | 5. B  | 6. A  | 7. C  | 8. A  | 9. C  | 10. B |
| 11. D | 12. D | 13. A | 14. B | 15. D | 16. A | 17. A | 18. D | 19. B | 20. C |
| 21. D | 22. D | 23. B | 24. A | 25. B | 26. A | 27. A | 28. B | 29. A | 30. A |

*This page is intentionally left blank*

# Networks, Information Systems, Software Engineering and Web Technology

## Part A Network

|                                      |      |
|--------------------------------------|------|
| <b>Chapter 1:</b> OSI Layers         | 8.3  |
| <b>Chapter 2:</b> Routing Algorithms | 8.24 |
| <b>Chapter 3:</b> TCP/UDP            | 8.36 |
| <b>Chapter 4:</b> IP(V4)             | 8.52 |
| <b>Chapter 5:</b> Network Security   | 8.66 |

U

N

I

T

8

*This page is intentionally left blank*

# Chapter 1

## OSI Layers

### LEARNING OBJECTIVES

- Computer network
- LAN
- LAN topologies
- CSMA/CD
- WAN
- The OSI reference model
- LAN technologies
- Physical layer
- Data link layer
- Types of error
- MAC sub layer
- FDM/TDM

### COMPUTER NETWORK

Computer network is the collection of two or more computers that are interconnected with each other to perform data communication using the data communication protocol through communications media (wired or wireless). So these computers can share information, data, programs, and use of hardware together. Data communications that can be done include text data, images, video and sound.

Or

A computer network, often simply referred, as a network is a collection of computers and devices interconnected by communication channels that facilitate communication and allow sharing of resources and information among interconnected devices. There are different networks:

1. LAN
2. MAN
3. WAN

### LAN

A Local Area Network (LAN) is a network that is confined to a relatively small area. It is generally limited to a geographic area such as a lab, school, or building. LAN Computers rarely spans more than a mile apart.

In a typical LAN configuration, one computer is designated as the file server. It stores all the software that controls the network, as well as the software that can be shared by the computers attached to the network. Computers connected to the file server are called workstations. The workstations can be less powerful than the file server, and they may have additional software on their hard drives. On many LANs, cables are used to connect the network interface cards in each

computer. However, one LAN can be connected to other LANs over any distance via telephone lines and radio waves.

Most LANs connect workstations and personal computers. Each node (individual computer) in a LAN has its own CPU which executes programs and it is also able to access data and devices anywhere on the LAN. This means that many users can share expensive devices, such as laser printers, as well as data. Users can also use LAN to communicate with each other, by sending e-mail or engaging in chat sessions.

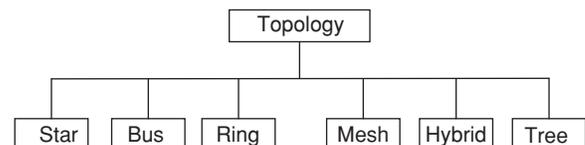
### Three characteristic features of LAN

1. The size of a LAN network.
2. The topology of the local area network.
3. The technology used for transmission.

In simple LAN configuration, a single cable runs through the entire set up and the peripherals and computers are attached to the cable. Traditional LAN speeds are 10 Mbps to 100 Mbps. Modern LAN cables are capable of much higher data transfer per second.

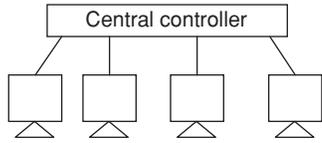
In case two or more systems need to use the LAN at the same time, then an arbitration mechanism is deployed to resolve the conflict. A first come first serve policy or a prioritized approach may be chosen.

### LAN topologies



**Star Topology** Each device has a dedicated point-to-point link to a central controller called a hub. Most used LAN topology.

If one device wants to send data to another, it sends the data to the controller, which then relays the data to the other connected device as shown in the figure.



Each device needs only one link and one I/O port to connect it to any number of others.

**Advantages**

1. Robust, if one link fails, only that link is affected. All other links remain active.
2. As long as hub is working, it can monitor link problems and by pass defective links.

**Disadvantages**

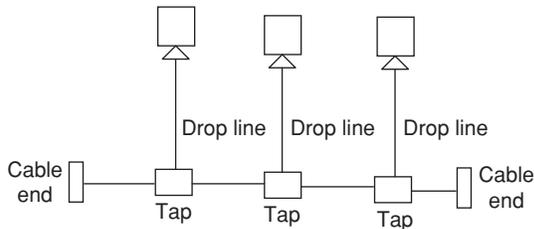
1. If hub goes down, the whole system dead.
2. More cabling is required in a star than ring or bus.

**Bus Topology** A bus topology is multipoint. One long cable acts as a backbone to link all the devices in a network.

Carrier Sense Multiple Access/Collision Detection (CSMA/CD) is the accessing technique used.

The traffic can go in either direction, i.e., it is bidirectional.

Nodes are connected to the bus cable by drop lines and taps as shown in the figure.



**Advantages**

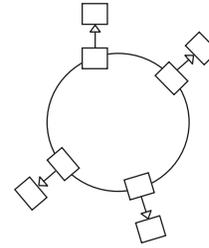
1. Ease of installation.
2. Require less cabling than mesh or star topologies.

**Disadvantages**

1. Difficult to add new devices.
2. A fault in the bus cable stops all transmission. The damaged area reflects signals back in the direction of origin, creating noise in both directions.

**Ring Topology** Each device has a dedicated point-to-point connection with only the two devices on either side of it.

Each device in the ring incorporates a repeater; when a device receives a signal intended for another device, its repeater regenerates the bits and passes them along.



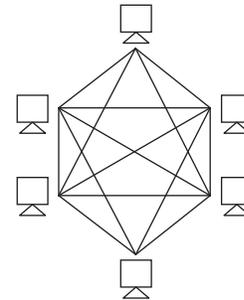
**Advantages**

1. Easy to install and reconfigure.
2. Fault isolation is simplified as it issues alarm which alerts the network operator to the problem and its location.

**Disadvantages**

1. A break in the ring can disable the entire network.
2. It is not relevant for higher-speed LANs.

**Mesh topology** Every station is interconnected to every other station as shown in the figure.



$n(n - 1)/2$  (duplex mode) links are required for communication in both directions. Each device on the network must have  $(n - 1)$  I/O ports to be connected to the other  $(n - 1)$  stations.

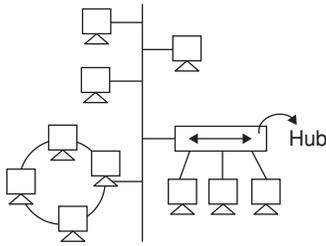
**Advantages**

1. The use of dedicated links guarantees that each connection can carry its own data load, thus eliminating traffic problems.
- 2 This topology is robust. If one link becomes unusable, it does not incapacitate the entire system.
3. There is advantage of security, only the intended recipient sees the message on the dedicated line.
4. Fault identification and fault isolation is easy because of point-to-point links.

**Disadvantages**

1. As the hardware(cables) required for connection is more, it is expensive.
2. Installation and reconnection are difficult.
3. The sheer bulk of the wiring can be greater than the available space.

**Hybrid Topology** More than one topology in a network.



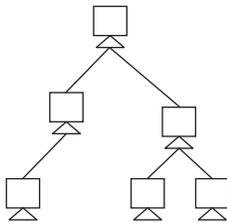
#### Advantages

1. Fault detection is easier.
2. We can add new stations without affecting the original architecture.

#### Disadvantages

1. As different topologies are combined so complexity of design increases. Very less practical implementation.
2. The hub which is used to connect different topologies is very costly. Moreover the cost of whole infrastructure is very high.

**Tree Topology** This topology uses the combination of star and bus topology.



#### Advantages

1. Expansion is easier; one can add new stations easily.
2. Errors can be easily detected.
3. Robust, if one link fails the remaining system is in communication.

#### Disadvantages

1. With the increase in the number of nodes, complexity and maintenance become difficult.

**Examples:** The most common type of local area network is an Ethernet LAN. The smallest home LAN can have exactly two computers; a large LAN can accommodate thousands of computers. Many LANs are divided into logical groups called subnets. An Internet Protocol (IP) 'Class A' LAN can in theory accommodate more than 16 million devices organized into subnets.

## MAN

A metropolitan area network is a computer network that usually spans a city or a large campus. A MAN usually

interconnects a number of local area networks (LANs) using a high-capacity backbone technology, such as fiber-optical links, and provides up-link services to wide area networks (or WAN) and the Internet.

## WAN

Wide Area Networks (WANs) connect larger geographic areas, such as Florida, the United States, or the world. Dedicated transoceanic cabling or satellite uplinks may be used to connect this type of network.

A WAN is complicated; it uses multiplexers to connect local and metropolitan networks to global communications networks like the Internet. To users, however, a WAN will not appear to be much different than a LAN. As the term implies, a WAN spans a large physical distance. The Internet is the largest WAN, spanning the Earth.

A WAN is a geographically-dispersed collection of LANs. A network device called a router connects LAN to a WAN. In IP networking, the router maintains both a LAN address and a WAN address.

A WAN differs from a LAN in several important ways. Most WANs (like the Internet) are not owned by any one organization but rather exist under collective or distributed ownership and management. WANs tend to use technology like ATM, Frame Relay and X.25 for connectivity over the longer distances.

Residences typically employ one LAN and connect to the Internet WAN via an Internet Service Provider (ISP) using a broadband modem. The ISP provides a WAN IP address to the modem, and all of the computers on the home network use LAN (so-called private) IP addresses. All computers on the home LAN can communicate directly with each other but must go through a central gateway, typically a broadband router, to reach the ISP.

## THE OSI REFERENCE MODEL

The concept of how a modern day network operates can be understood by dissecting it into seven layers. This seven layer model is known as the OSI Reference Model and defines how the vast majority of the digital networks on earth function. OSI is the acronym for Open Systems Interconnection. The important concept to realize about the OSI Reference Model is that it does not define a network standard, but rather provides guidelines for the creation of network standards.

### Physical Layer

The first layer of a network is the Physical Layer. The Physical Layer is literally what its name implies: the physical infrastructure of a network.

This includes the cabling or other transmission medium and the network interface hardware placed inside computers

and other devices which enable them to connect to the transmission medium.

The purpose of the Physical Layer is to take binary information from higher layers, translate it into a transmission signal or frequency, transmit the information across the transmission medium, receive this information at the destination and finally translate it back into binary before passing it up to the higher layers.

Transmission signals or frequencies vary between network standards and can be as simple as pulses of electricity over copper wiring or as complex as flickers of light on optical lines or amplified radio frequency transmissions.

The information that enters and exits the Physical Layer must be bits; either 0s or 1s in binary. The higher layers are responsible for providing the Physical Layer with binary information. Since almost all information inside a computer is already digital, this is not difficult to achieve.

The Physical Layer does not examine the binary information nor does it validate it or make changes to it. The Physical Layer is simply intended to transport the binary information between higher layers located at points A and B.

## Data Link Layer

The second layer in the OSI Model is the Data Link Layer, the only layer in the OSI model that specifically addresses both hardware and software.

The Data Link Layer receives information on its software side from higher layers, places this information inside 'frames', and finally gives this frame to the Physical Layer, Layer 1, for transmission as pure binary.

A frame essentially takes the information passed down from a higher layer and surrounds it with Physical Address information. This information is important for the Data Link Layer on the receiving end of the transmission.

When the frame, in binary form, arrives at the destination node, it is passed from the transmission medium to the Data Link Layer (Layer 2) by the Physical Layer (Layer 1).

The Data Link Layer on the receiving node checks the frame surrounding the information received to see if it's Physical Address matches that of its own. If the Physical Address does not match, the frame and its encapsulated data is discarded. If the Physical Address is a match, then the information is removed from the frame and passed up to the next highest layer in the OSI Model.

The Physical Addressing system allows multiple nodes to be on the same network medium, but retain the ability to address only a specific node with a transmission.

The Physical Address used in the Data Link Layer's Physical Addressing system is known as a MAC address and is embedded physically into the node's Network Interface Card during manufacturing.

Every NIC's MAC address is unique in order to prevent addressing conflicts. It is this relationship that causes the

Data Link Layer to be known as the only layer that addresses both hardware and software.

In this layer the information on the network makes the move from the physical infrastructure of the network into the software realm. The remainders of the OSI reference model's layers are entirely software.

## Network Layer

OSI Layer 3 is known as the Network Layer. The purpose of the Network Layer is to direct network traffic to a destination node whose Physical Address is not known. This is achieved through a system known as Logical Addressing.

Logical Addresses are software addresses assigned to a node at Layer 3 of the OSI Model. Since these addresses are able to be defined by software rather than being random and permanent like Physical Addresses, Logical Addresses are able to be hierarchical. This allows extremely large networks to be possible.

A smart device working at Layer 3 that handles network signals from each node directly rather than nodes just blindly repeating packets at Layer 1 until they happen to reach their destination. Such a device is known as a network router.

A network router sits in the center of a network with all nodes having a direct link to it rather than being linked to each other. This strategic position allows the router to intercept and direct all traffic on the network.

A routed network can be illustrated by a star formation, as shown in Diagram 1. On a routed network, Layer 3 packets are no longer broadcasted to all nodes, but rather received by the router and passed on only to the appropriate node. This is a valuable concept because it allows for the collision free-transport of packets across a network.

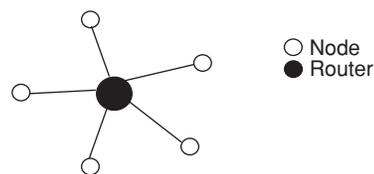


Figure 1

As being linked directly to all nodes in a local network, a router can be linked directly to other routers. This allows groups of nodes separated by distance to communicate with each other in a practical way.

It would not be practical to have nodes separated by a great distance all connected to a single router. The amount of cabling required would be immense and depending on the number of nodes involved, the router may not possess the required number of physical connections.

Routers can be chained in a line, or as shown in Diagram 2, can be connected by a central router. This concept is virtually infinitely scalable and is very efficient.

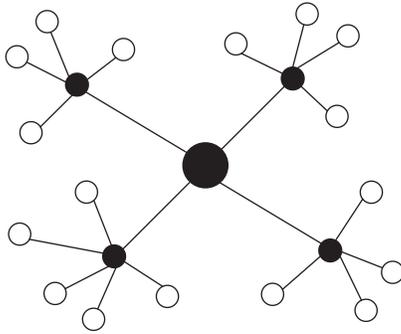


Figure 2

When a node starts a transmission, the OSI Layer 3 protocol takes the information passed down from higher layers and encapsulates it with the logical address of the destination node in a unit called a packet.

This packet, then passes through the remaining lower layer protocols, is transmitted over the network medium from the node to the router. This router reads the logical address that the packet contains and compares it to a list of physical addresses of nodes that are directly connected to it.

If the packet's destination address matches an entry in this list, the packet is transmitted directly on the line that leads straight to the destination node.

If the router does not know of a direct connection to the destination node, the packet is transmitted on a line leading directly to another router. This router then treats the packet much like the first router did upon receipt.

The packet's logical address is checked for matches against the list of logical addresses belonging to nodes directly connected to the router.

If the packet reaches a router with connections only to other routers, as shown in Diagram 2, the router uses the logical address's orderly numbering scheme to try and determine the closest router to the destination node and then transmits the packet to that router.

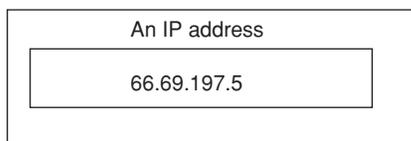


Figure 3

In IP, logical addresses look like four sets. Diagram 3 shows an example of an IP address. IP addresses are orderly on four levels, from left to right. The first section of the IP address refers to a top level router, or a router that is at the highest level of this particular branch of the network. In Diagram 3, the first number is 66. Therefore all IP addresses between 66.0.0.1 and 66.255.255.255 are managed by this router. Only one router is required in a routed network, but more may exist. A router may have a maximum of 255 nodes, which may be either ordinary nodes or other routers. This

effectively means that each branch of a network, a group of nodes that have the first set of numbers in their IP address in common, could theoretically have over sixteen million end nodes.

## Transport Layer

OSI Layer 4 is known as the Transport Layer, all information transferred is assumed to be at the correct destination node and is being passed up to Layer 4.

The Transport Layer is responsible for the reliability of the link between two end users and for dividing the data that is being transmitted by assigning port numbers to its Layer 4 packages, known as segments.

Ports can be thought of as virtual destination mailboxes or outlets. When information reaches a Layer 4 protocol, the segment is examined to determine the destination port of the data it contains. Once the port is determined, just as all of the past layers have done, the wrapper is discarded and the payload data passed up to the next layer's protocol.

Higher layer protocols that provide services such as email, web browsing, text chat, file transfer and more, each operate on their own unique Layer 4 port, allowing all of these protocols to be operated at once without interference.

On the reliability front, Transport Layer protocols are capable of running a checksum on the payload data, which they carry. This allows the protocol to determine the integrity of incoming payload data. If this data has been corrupted, the Layer 4 protocol will request the segment to be retransmitted.

## Session Layer

OSI Layer 5, known as the Session Layer, still serves a purpose in the OSI Reference Model. The Session Layer draws the outline for protocols that manage the combination and synchronization of data from two separate higher layers.

Layer 5 protocols are responsible for ensuring that the data is Synchronized and consistent before transmitted. A good example is the streaming of live multimedia audio and video, where perfect synchronization between video and audio is desired.

## Presentation and Application Layers

The sixth and seventh layers in the OSI Reference Model are the Presentation Layer and the Application Layer. The primary purpose of these layers is to facilitate the movement of formatted information between applications interacting with end users on nodes.

Commonly used top layer protocols are HTTP (for the secure transfer of web page related files), File Transfer Protocol (FTP), Simple Mail Transfer Protocol (SMTP, used for sending email messages), and SSH (Secure Shell), used to secure remote shell access for a computer operating system.

### OSI reference model concept

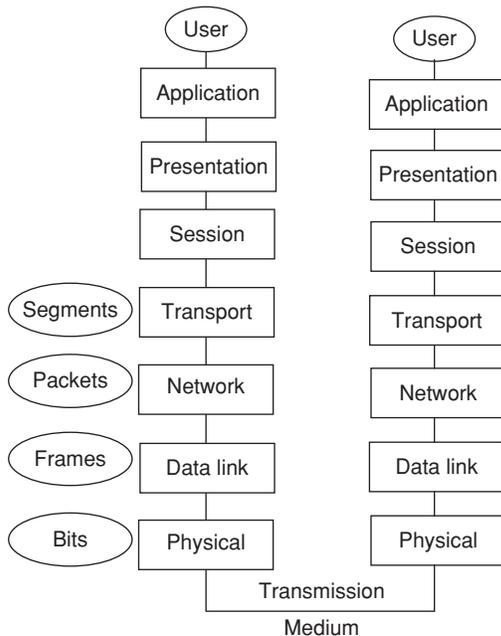


Figure 4

The OSI Reference Model exists not to make hard rules or to shape the industry, but to provide a logical, well-researched, and tested model after which the world's best communication protocol stacks are modeled. The TCP/IP stack is very well-known for being the driving force behind most of the internet, and represents the third (IP) and fourth (TCP) layers of the OSI Model. Every layer in the OSI Model is a reference for a protocol which must facilitate communication between both higher and lower layers. The 'U-shaped' example shown in Diagram 4 provides a visual concept of how two users may be linked on a given network in reference to the OSI Model. Data starts and ends with the user. From the Application Layer of the first user, it must travel down through layers 7 to 1, across the transmission medium, then back up to layers 1 to 7 to be presented at the Application Layer to the user on the end of the transmission. Diagram 4, shows an example of a path between two nodes. Protocols defined by this reference are dependent on the next lowest layer protocol. So, for example, one could not run an Application Layer protocol on a node without the presence of Layer 1 through 6, protocols also being utilized on the node.

## LAN TECHNOLOGIES

### IEEE standard for networking

IEEE standard project 802 is designed for the enter – connectivity between LAN's

IEEE 802 maps to physical and data link layer

**Example:** Ethernet, Token ring etc, the IEEE standards for the different groups are

- |                                                                                                                                                                                                                                                                                                                                                                                                                                     |   |                                        |
|-------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------|---|----------------------------------------|
| <ul style="list-style-type: none"> <li>802.1 – Higher layer LAN Protocol</li> <li>802.3 – Ethernet</li> <li>802.11 – Wireless LAN</li> <li>802.15 – WPAN</li> <li>802.16 – Broad band wireless Access</li> <li>802.17 – Resilient packet Ring</li> <li>802.18 – Radio Regulatory TAG</li> <li>802.19 – Co existence TAG</li> <li>802.20 – Mobile Broad band wireless access</li> <li>802.21 – Media independent Hand off</li> </ul> | } | Active working group                   |
| <ul style="list-style-type: none"> <li>802.2 – Logical link control working group</li> <li>802.4 – Token Bus</li> <li>802.5 – Token Ring</li> <li>802.7 – Broad band area Network</li> <li>802.8 – Fiber optic TAG</li> <li>802.9 – Integrated service LAN</li> <li>802.10 – Security working group</li> <li>802.12 – Demand priority working group</li> <li>802.14 – Cable modern working group</li> </ul>                         | } | In active or dis-banded working groups |

## Ethernet

We have

- 10 Mbps – Ethernet
- 100 Mbps – Fast ethernet
- 1 Gbps – Gigabit Ethernet
- 10 GE – 10 Gigabit Ethernet

Best suited for LAN because it is capable of handling high speed bandwidth.

- Ethernet medium:
  - Thick wire – 10B5
  - Thin wire – 10B2
  - Twisted pair – 10BT, 100BT, 1000BT
  - Fibre – 10BF, 100BF, 1000BF
  - CAT 4 – 10 Mbps
  - CAT 5 – 10/100 Mbps
  - CAT 6 – 10/100/1000 Mbps
- Fundamental is CSMA/CD, Standard is 802.3.
- It defines two categories:
  1. Base band.
  2. Broad band
- Baseband uses digital manchester encoding techniques.
- IP communication in ethernet is of 3 types:
  - (i) Unicast
  - (ii) multi cast
  - (iii) broadcast
- When user sends data he puts destination and source address.
  - (i) In unicast, only intended users responds, however all can get the signal (individual MAC).
  - (ii) In multicast, group of users will get the data (group MAC).
  - (iii) In broadcast, all users on Ethernet can see the data (all MAC).
- Every computer accepts 3 types of packets, to his own, to the group it belongs, to all.

## CSMA/CD

### CSMA (Carrier Sense Multiple Access)

CSMA protocols performance is better than ALOHA—Monitor the channel before and/or during data transmission.

**1-Persistent** Check whether the channel is free before transmitting the data. If busy, wait until it becomes free and then immediately start Re-transmitting.

**Non-Persistent** When the channel is busy, wait for a random period of time before trying again

If the waiting time is too long, the channel utilization decreases.

**P-Persistent** Used in slotted systems, If the channel is idle during the current slot, transmit with probability  $P$ , and defer until next slot with probability  $(1 - P)$

Two or more computers can get connected on same physical medium. All computers can communicate whenever they feel like. Any computer want to communicate, it senses the medium, if medium is free and not used by anyone it captures the medium and puts its data on to the channel.

All computers listens to the sent data but only intended computer/system will respond. At this instance, sending computer is owner of the medium; no other system can be owner or can send the data. When two or more computers try to send data at same time by sensing the medium, collision occurs, which will be sensed by all the computers, then they keep integral wait unit of time for next transmission of data. Once the sending machine gets the corrupted collision message it retransmits using integral time.

### MAC sublayer

The medium Access control (MAC) sub layer is the bottom half of the Data link layer. The upper half is commonly called logical link control (LLC) sublayer.

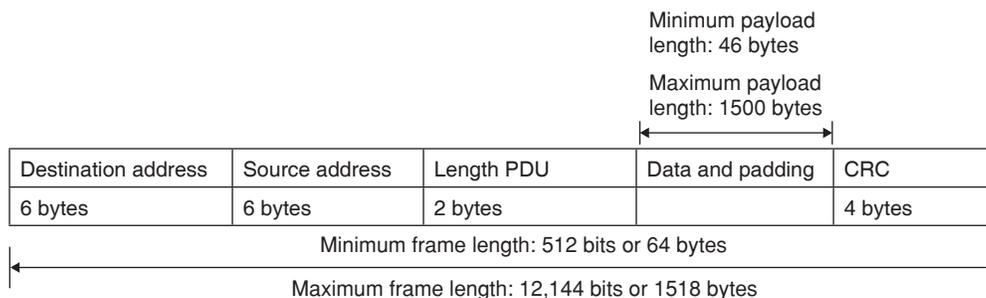
#### Frame format

| Preamble | SFD    | Destination address | Source address | Length or type | Data and padding | CRC     |
|----------|--------|---------------------|----------------|----------------|------------------|---------|
| 7 bytes  | 1 byte | 6 bytes             | 6 bytes        | 2 bytes        | 46 to 1500 bytes | 4 bytes |

**Preamble** The first field of the 802.3 frame contains 7 bytes (56 bits) of alternating 0's and 1's that alerts the receiving

### Frame length

Ethernet has imposed restrictions on both the minimum and maximum lengths of a frame, as shown in figure below:.



system to the coming frame and enables it to synchronize its input timing.

**Start frame delimiter (SFD)** The second field (1 byte: 10101011) signals the beginning of the frame. The SFD warns the station or stations that, this is the last chance for synchronization. The last 2-bits is 11 and alerts the receiver that the next field is the destination address.

**Destination address (DA)** The DA field is 6 bytes and contains the physical address of the destination station to receive the packet.

**Source address (SA)** The SA field is also 6 bytes and contains the physical address of the sender of the packet.

**Length field** The original ethernet used this field as the type field to define the upper-layer protocol using the MAC frame. The IEEE standard uses it as the length field to define the number of bytes in the data field.

**Data** This field contains data encapsulated from the upper layer protocols. It is of minimum 46 bytes and a maximum of 1500 bytes.

Ethernet follows **binary exponential back off** algorithm to give waiting time for stations, which are involved in collisions. After collisions, waiting time for the stations will be  $K * 51.2 \mu \text{ sec}$ , where  $K$  is randomly picked up from 0 to  $2^n - 1$ , 'n' is the collision number. But after 10 collisions, the randomization internal is frozen at a maximum of 1023 slots.

If each station transmits during a contention slot with probability  $p$ , the probability  $A$  that some station acquires the channel in that slot is

$$A = Kp(1 - p)^{k-1}$$

$A$  is maximized when  $p = 1/k$ , with  $A \rightarrow 1/e$  as  $K \rightarrow \infty$

The probability that the contention internal has exactly  $j$  slots in it is  $A(1 - A)^{j-1}$ , hence mean number of slots per contention is given by

$$\sum_{j=0}^{\infty} jA(1 - A)^{j-1} = \frac{1}{A}$$

**CRC** The last field contains error detection information.

Minimum length of frame is 512 bytes or 64 bytes. If we count 18 bytes of header and trailer, then minimum length of data from the upper layer is  $64 - 18 = 46$  bytes.

If the upper layer packet is less than 46, padding is added to make up the difference and used to find out collision. Maximum length of the frame is 1518 bytes. If we subtract 18 bytes of header and trailer, the maximum length of the payload is 1500 bytes.

- The maximum length restriction has two reasons.
- First, memory was very expensive when ethernet was designed, a maximum length restriction helped to reduce the size of the buffer.
- Second, the maximum length restriction prevents one station from monopolizing the shared medium, blocking other stations that have data to send.

Since each slot has a duration  $2T$ , the mean contention interval,  $w$ , is  $2T/A$ . Assuming optional  $p$ , the mean number of contention slots is never more than  $e$ , so  $w$  is at most  $2Te = 5.4 T$ .

### Frame formats

|        |        |        |                     |                |                  |         |        |        |
|--------|--------|--------|---------------------|----------------|------------------|---------|--------|--------|
| SD     | AC     | FC     | Destination address | Source address | Data             | CRC     | ED     | FS     |
| 1 byte | 1 byte | 1 byte | 2-6 bytes           | 2-6 bytes      | Up to 4500 bytes | 4 bytes | 1 byte | 1 byte |

|             |    |    |
|-------------|----|----|
| SD          | AC | ED |
| Token frame |    |    |

|             |    |
|-------------|----|
| SD          | ED |
| Abort frame |    |

802.5 Token ring uses differential Manchester digital signal encoding. It supports data rates upto 16 mbps. Tokens ring protocol specifies three types of frames: Data, token, and abort.

The token and abort frames are both truncated data frames.

#### Data frame

**Start delimiter (SD)** It is one byte long and is used to alert the receiving station for the arrival of a frame as well as to synchronize it's timing.

**Access control (AC)** It is one byte long and includes sub-fields. It has the format PPPTMRRR. First 3-bits are priority field. T denotes whether this is a data frame, token or an abort frame. Token bit is followed by monitor bit. The last 3 bits are the reservation field that can be set by stations wishing to resume access to ring.

**Frame control** This field is one byte long and contains two fields. The first is a one bit field used to indicate the type of information (whether it is a control information or data). The second uses the remaining seven bits of the byte and contains information used by the token ring logic.

**Destination address (DA)** The six byte DA field contains the physical address of the frame's next destination.

$$\text{Channel efficiency} = \frac{p}{p + 2T/A}$$

$$= \frac{1}{1 + 2B \frac{Le}{cF}}$$

- Where  $F$  = Frame length
- $B$  = Network bandwidth
- $L$  = Cable length
- $C$  = Speed of signal propagation
- $E$  = Contention slots per frame

### 802.5 TOKEN Ring

Here ring topology is used and devices are physically arranged to form a ring. A token is passed among stations. If a station wants to send data, it must wait and capture the token. Only the token holders are permitted to transmit frame. Token ring allows each station to send one frame per turn.

**Source address (SA)** The six byte SA field contains the physical address of the sending station.

**Data contains LLC data unit** Data contains 0 or more bytes, maximum size of the data depends upon taken holding time.

**CRC** The CRC field is 4 byte long and contains a CRC – 32 bit error detection sequence.

**End delimiter (ED)** ED is a second flag field of one byte and indicates the end of the sender's data and control information.

**Frame status** It is one byte long

|     |  |  |     |  |  |
|-----|--|--|-----|--|--|
| A/C |  |  | A/C |  |  |
|-----|--|--|-----|--|--|

$A$ : Addressed recognized bit

$C$ : Copies bit

It can be set by the receiver to indicate that the frame has been read/copied etc.

When a frame arrives at the station with the destination address, the station turns  $A$  bits to 1. If station copies the frame to the station it also turns on the  $C$  bit. A station might fail to copy a frame due to lack of frame buffer or other reason.

When the sending station receives the frame, it examines the  $A$  and  $C$  bits.

Three combinations are possible:

1.  $A = 0, C = 0$ : destination not ready /present.
2.  $A = 1, C = 0$  : destination present byte frame not accepted.
3.  $A = 1, C = 1$  : destination present and frame copied.

### Token frame

It includes only 3 fields: SD, AC and ED

1. The SD indicates, the frame is coming
2. The AC indicates that the frame is a token and includes priority and reservation fields.  $T = 0$  for token in AC.
3. The ED indicates the end of the frame.

### Abort frame

An abort frame contains no information at all just starting and ending delimiters. It can be generated by the sender to stop its own transmission. Each station has a priority code, as a frame passes by, a station waiting to transmit may reserve the next open token by entering its priority code in the Access control field (AC) of the token or data frame. A station with a higher priority may remove a lower priority reservation and replace it with its own. Among stations of equal priority, the process is first come, first served. Through this mechanism, the station holding the reservation gets the opportunity to transmit as soon as the token is free, whether or not it comes next physically on the ring.

**Monitor station** Several problems may appear to disrupt the operation of a token ring network. If the token is destroyed by noise there will be no token on the ring and no station can send data. To solve such a problem, one station on the ring is designated as a monitor. The monitor sets a time, each time the token passes. If the token does not appear in the allotted period of time, it is assumed to be lost and the monitor generates a new token and introduces it to the ring. The monitor detects the orphan frames, by setting the monitor bit in the access control byte.

As the frame passes, the monitor checks the status field. If the monitor bit is set, something is wrong since the frame has passed the monitor twice, so monitor discards it. The monitor then destroys the frame and puts a token on the ring. If monitor fails, the protocol ensures that another station is quickly selected as monitor. Every station has the capability of becoming the monitor. While the monitor is functioning properly, it alone is responsible for seeing that the ring operates correctly.

When station notices that either of its neighbors appears to be dead it transmits BEACON frame giving the address of the dead station.

## PHYSICAL LAYER

Physical layer is concerned with transmitting raw bits over a communication channel. The design issues have to do with making sure that when one side sends a 1-bit, it is

received by the other side as 1-bit and not as 0 bit. In physical layer we deal with the communication medium used for transmission.

### Types of Medium

Medium can be classified into two categories:

1. **Guided Media:** Guided media means that signals are guided by the presence of physical media i.e., signals are under control and remains in the physical wire. For example, copper wire.
2. **Unguided Media:** Unguided media means that there is no physical path for the signal to propagate. Unguided media has essentially electromagnetic waves. There is no control on flow of signal. For example, radio waves.

### Transmission Media

In Guided transmission media, generally two kinds of materials are used.

1. Copper
  - Coaxial cable
  - Twisted pair
2. Optical Fiber

#### Coaxial cable

Coaxial cable consists of an inner conductor and an outer conductor which are separated by an insulator. The inner conductor is usually copper. The outer conductor is covered by a plastic jacket. It is named coaxial because the two conductors are coaxial. Typical diameter of coaxial cable lies between 0.4 inches to 1 inch.

#### Twisted pair

A twisted pair consists of two insulated copper wires, typically 1mm thick. The wires are twisted together in a helical form, the purpose of twisting is to reduce cross talk interference between several pairs. Twisted pair is much cheaper than coaxial cable but it is susceptible to noise and electromagnetic interference and attenuation is large.

#### Optical fiber

In optical Fiber light is used to send data. In general terms presence of light is taken as bit-1 and its absence as bit 0. Optical fiber consists of either glass or plastic core which is surrounded by cladding of the same material but of different refractive index. This cladding is surrounded by a plastic jacket which prevents optical fiber from electromagnetic interference and harshly environments. It uses the principle of total internal reflection to transfer data over optical fibers. Optical fiber is much better in bandwidth as compared to copper wire, since there is hardly any attenuation or electromagnetic interference in optical wires. Hence there is

less requirement to improve quality of signal, in long distance transmission. Disadvantage of optical fiber is that end points are fairly expensive.

## Communication Links

In a network nodes are connected through links. The communication through links can be classified as

**Simplex** Communication can take place only in one direction.

**Example:** TV broadcasting.

**Half duplex** Communication can take place in one direction at a time. Suppose node A and B are connected then half duplex communication means that at a time data can flow from A to B or from B to A but not simultaneously.

**Example:** Two persons talking to each other such that when one speaks the other listens and vice versa, walkie-talkies, citizens band radios.

**Full duplex** Communication can take place simultaneously in both directions.

**Example:** telephone network.

Links can be further classified as:

**Point-to-Point** In this communication only two nodes are connected to each other. When a node sends a packet then it can be received only by the node on the other side and none else.

**Multi-Point** It is a kind of sharing communication in which signal can be received by all nodes. This is also called broadcast.

## Digital Data to Digital Signals

A digital signal is sequence of discrete, discontinuous voltage pulses. Each pulse is a signal element. Encoding scheme is an important factor in knowing that how successfully the receiver interprets the incoming signal.

### Encoding techniques

Following are several ways to map data bits to signal elements:

**Non-return-to-zero (NRZ):** NRZ codes share the property that voltage level is constant during a bit interval. High level voltage = bit 1 and low level voltage = bit 0. A problem arises when there is a long sequence of 0's and 1's and the voltage level is maintained at the same value for a long time.

This creates a problem on the receiving end because now, the clock synchronization is lost due to lack of any transitions and hence, it is difficult to determine the exact number of 0's and 1's in this sequence.

The two variations are as follows:

1. **NRZ-Level:** In NRZ-L encoding, the polarity of the signal changes only when the incoming signal

changes from a '1' to a '0' or from a '0' to a '1'. NRZ-L method, looks just like the RZ method, except for the first input one data bit. This is because NRZ does not consider the first data bit to be a polarity change, where NRZ-L does.

2. **NRZ-Inverted:** Transition at the beginning of bit interval = bit 1 and no transition at the beginning of bit interval = bit 0 or vice versa. This technique is known as differential encoding.

## Digital Data Communication Techniques

For two devices linked by a transmission medium to exchange data, a high degree of co-operation is required. Typically data is transmitted one bit at a time. The timing (rate, duration, spacing) of these bits be same for transmitter and receiver. There are two options for transmission of bits.

**Parallel** All bits of a byte are transferred simultaneously on separate parallel wires. Synchronization between multiple bits is required which becomes difficult over large distance. Parallel communication gives large bandwidth but expensive, possible only for devices which are close to each other.

**Serial** Bits transferred serially one after other. Serial communication gives less bandwidth but cheaper, suitable for transmission over long distances.

### Manchester encoding

Manchester encoding is used in Ethernet (IEEE 802.3) it is a line code in which bit encoding has at least one transition and consumes the same time.

It ensures frequent line voltage transitions, which are directly proportional to clock rate

It is not dependent on data, so it will not carry any information.

## Transmission Techniques

### Asynchronous

Small blocks of bits (generally bytes) are sent at a time without any time relation between consecutive bytes. When no transmission occurs a default state is maintained corresponding to bit 1, due to arbitrary delay between consecutive bytes, the time occurrences of the clock pulses at the receiving end need to be synchronized for each byte. This is achieved by providing two extra bits, start and stop.

**Start Bit** It is prefixed to each byte and equals 0. Thus it ensures a transition from 1 to 0 at onset of transmission of byte. The leading edge of start bit is used as a reference for generating clock pulses at required sampling instants. Thus each onset of a byte results in resynchronization of receiver clock.

**Stop Bit** To ensure that transition from 1 to 0 is always present at beginning of a byte it is necessary that default state be 1, but there may be two bytes one immediately following the other and if last bit of first byte is 0, transition from 1 to 0 will not occur. Therefore a stop bit is suffixed to each byte equaling 1. It's duration is usually 1, 1.5, 2 bits. Asynchronous transmission is simple and cheap but requires an overhead of 3 bits i.e., for 7 bit code 2(start, stop bits) + 1 parity bit implying 30% overhead. However this percentage can be reduced by sending larger blocks of data but then timing errors between receiver and sender cannot be tolerated beyond  $[50/\text{number. of bits in block}]%$ . It will not only result in incorrect sampling but also misaligned bit count. i.e., a data bit can be mistaken for stop bit if receiver's clock is faster.

### Synchronous

Larger blocks of bits are successfully transmitted. Blocks of data are either treated as sequence of bits or bytes. To prevent timing, drift clocks at two ends need to be synchronized. This can be done in two ways.

1. Provide a separate clock line between receiver and transmitter. (or)
2. Clocking information is embedded in data signal i.e., Biphase coding for digital signals.

Still another level of synchronization is required so that receiver determines beginning or end of block of data. Hence each block begins with a start code and end with a stop code. These are in general same, known as flag that is unique sequence of fixed number of bits. In addition some control characters encompass data within these flags. Data and control information is called a frame. Since any arbitrary bit pattern can be transmitted, there is no assurance that bit pattern for flag will not appear inside the frame, thus destroying frame stuffing.

**Channel Allocation** A large class of networks is built on broadcast channels, a number of stations will share the same channel, if one station sends, all other stations have to hear it.

Problem occurs when, 2 stations want to start data transmission at the same time, in this situation 2 frames collide.

To avoid frame collision, allocate the channel to one of the stations.

There are 3-strategies for channel allocation:

1. Let a station try to use the channel, and when the collision occurs, that is taken care of later.
2. Each station in turn is allowed to use the channel. This is applied in token-based systems. Only the station that has the token can use the channel.
3. Reserve the channel in prior, It is used in slotted systems. The problem is how to make a reservation.

## DATA LINK LAYER

Data link layer provides interface to the network layer, determines the number of bits of the physical layer to be grouped into frames, detects transmission error and regulates the flow of frames.

Functions of data link layer:

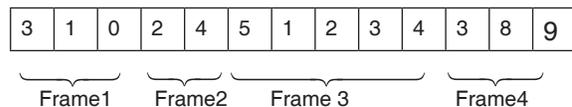
1. Framing
2. Physical addressing
3. Flow control
4. Error control
5. Access control

Various methods of Framing are

1. Time gaps
2. Character count
3. Starting and ending characters, with character stuffing
4. Starting and ending flags, with bit stuffing
5. Physical layer coding violations

**Time gaps** Framing is done by inserting time gaps between frames, very similar to the way of spacing between words in ordinary text. It is risky to count on timing to mark the start and end of each frame.

**Character count** It uses a field in the header to specify the number of characters in the frame. Thus at the destination by seeing the character count it knows how many characters follows and where the end of the frame exists.



**Problem** If count of any frame changes, destination will get out of synchronization and is unable to locate start of next frame.

### Starting and ending characters, with character stuffing

Each frame starts with the ASCII character sequence DLESTX and ends with the sequence DLEETX. If destination loses track of the frame boundaries, all it has to do is to look for DLESTX or DLEETX character

### Starting and ending flags, with bit stuffing

**Bit Stuffing:** Suppose our flag bits are 01111110. So the transmitter will always insert an extra 0 bit after each occurrence of five 1s (except for flags). After detecting a starting flag the receiver monitors the bit stream. If pattern of five 1's appear, the sixth bit is examined and if it is 0 it is deleted; else if it is 1 and next bit is 0 the combination is accepted as a flag. Similarly byte stuffing is used for byte oriented transmission. Here we use an escape sequence to prefix a byte similar to flag and two escape sequences if byte itself is an escape sequence.

Has arbitrary number of bits and allows character codes with an arbitrary number of bits per character. Every frame begins and ends with a special bit pattern, 01111110, called a flag byte.

As soon as the sender's data link layer encounters five consecutive one's in the data, it stuffs a 0 bit into the outgoing bit stream.

Receiver de-shifts the 0 bit of the five consecutive incoming 1 bits, followed by a 0 bit.

If the user data is 01111110, transmitted as 011111010 but stored at receiver as 01111110.

**Physical layer coding violations** Applied to the networks in which the encoding on the physical medium contains some redundancy.

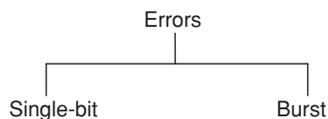
1 → high – low pair

0 → low – high pair

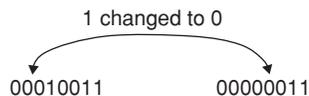
Here high-high, low-low not used for data.

Every data bit has a transition in the middle, thus easy for the receiver to locate the bit boundaries.

### TYPES OF ERRORS



**Single bit error** The term single bit error means that only one bit in the data unit has changed, it can either be from 1 to 0 or from 0 to 1.



### Single bit error correction

A single bit error occurs when a bit changes in value from 0 to 1 (or) from 1 to 0 while storing (or) while performing read (or) write operation. If that error bit is identified, that can be corrected by complementing.

### Hamming codes

In hamming codes,  $K$  parity bits are added to an  $n$ -bit data word, that forms a new word of  $(n + k)$  bits. The bit positions are numbered in sequence from 1 to  $n + k$ . These positions numbered with powers of 2 are reserved for the parity bits; the remaining bits are the data bits.

**Example:** Consider the given 8-bit data word 11000100, we include four parity bits with this word and arrange the bits as follows.

### Bit position

|       |       |   |       |   |   |   |       |   |    |    |    |
|-------|-------|---|-------|---|---|---|-------|---|----|----|----|
| 1     | 2     | 3 | 4     | 5 | 6 | 7 | 8     | 9 | 10 | 11 | 12 |
| $P_1$ | $P_2$ | 1 | $P_4$ | 1 | 0 | 0 | $P_8$ | 0 | 1  | 0  | 0  |

The parity bits are in positions, 1, 2, 4, 8. Each parity bit is calculated as

$$P_1 = \text{XOR of bits (3, 5, 7, 9, 11)} = 1 \oplus 1 \oplus 0 \oplus 0 \oplus 0 = 0$$

$$P_2 = \text{XOR of bits (3, 6, 7, 10, 11)} = 0$$

$$P_4 = \text{XOR of bits (5, 6, 7, 12)} = 1$$

$$P_8 = \text{XOR of bits (9, 10, 11, 12)} = 1$$

⇒ If there is odd number of 1s, XOR gives 0

⇒ If there is even number of 1s, XOR gives 1

The values  $P_1 = 0, P_2 = 0, P_4 = 1, P_8 = 1$  are substituted in 12-bit composed word

### Bit position

|   |   |   |   |   |   |   |   |   |    |    |    |
|---|---|---|---|---|---|---|---|---|----|----|----|
| 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 | 11 | 12 |
| 0 | 0 | 1 | 1 | 1 | 0 | 0 | 1 | 0 | 1  | 0  | 0  |

### Check for errors:

$$C_1 = \text{XOR of bits (1, 3, 5, 7, 9, 11)}$$

$$C_2 = \text{XOR of bits (2, 3, 6, 7, 10, 11)}$$

$$C_4 = \text{XOR of bits (4, 5, 6, 7, 12)}$$

$$C_8 = \text{XOR of bits (8, 9, 10, 11, 12)}$$

Since the bits were written with even parity, the result  $C = C_8 C_4 C_2 C_1 = 0000$

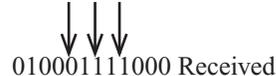
∴ Indicates that no error has occurred.

- The code can be used with words of any length.

**Burst Error** The term burst means that two or more bits in the data unit have changed, either changed, from 1 to 0 or changed from 0 to 1.

Sent:

010011010000-sent bits corrupted by burst error



### Parity bit

Parity bit is an error detecting code. This bit is added to data words depending on number of 1's in the data word; It could be even parity and odd parity.

$n$ -bit data word is transformed to  $(n + 1)$  bit code word with the addition of a bit. Even parity makes even number of 1's in a code word, similarly odd parity makes odd number of 1's in a code word.

Let us illustrate with example

Data word – 1 0 1 1 Parity bit

Code word – 1 0 1 1 1 parity bit (even parity)

Code word: 1 0 1 1 0 (odd parity)

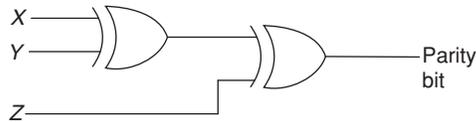
At the receiver side, when the code is received, the receiver checks the same as it is done by the generator. But here it adds all the bits which results in syndrome. If the syndrome is 0 then the number of 1's in code word is even, else number of 1's is odd.

Decision logic analyzer will decide, whether the code word is correct or not, based on syndrome value.

### Parity bit generator

The parity bit generator for a 3-bit data word is given below.

The message is in the form of  $XYZ$



Parity bit generator

When the message is passed through the above circuit, the parity will be generated accordingly.

### Error correction code

When data is transmitted from the source to destination, there is a chance of error introduction into the data. Error detection will detect the errors in data, while error correction will rebuild the original data,

Error correction code can be implemented in 2 ways

1. Forward error correction (FEC)
2. Automatic repeat request (ARQ)

In FEC, when sender is sending data, sender adds redundant data (encoded information) to the original data. At the receiver side this redundant data is used to recover the original data, when original data is tampered [error data].

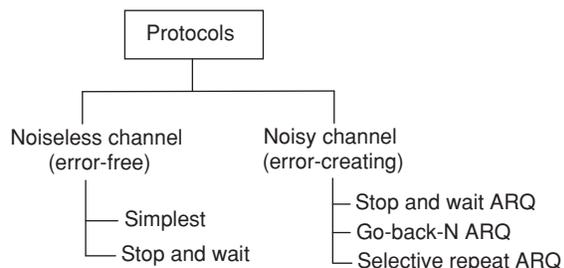
In ARQ, the receiver requests for the retransmission of the data packets, which are corrupted. Receiver will check the data using some error detection code.

### Flow Control

It regulates the flow of frames so that slow receivers are not affected by the fast sender or vice versa.

It tell the sender how much data it should transmit before it waits for an acknowledgement from the receiver. Flow control refers to a set of procedures used to restrict the amount of data that the sender can send before waiting for acknowledgement.

Error control in the data link layer is often implemented simply. Any time an error is detected in an exchange, specified frames are retransmitted. This process is called automatic repeat request (ARQ).



All the protocols we discuss as unidirectional in the sense that data frames travels from sender to receiver. Although special frames called acknowledgment (ACK) and negative acknowledgment (NAK) can flow in the opposite direction for flow and error control purposes, data flow is in only one direction. In real life network, the data link protocols are implemented as bidirectional, data flow in both directions. In these protocols flow and error control information such as ACKs and NAKs are included in the data frames in a technique called piggybacking.

**Stop and wait** Sender sends one frame, stops until it receives confirmation from the receiver. Error correction in stop and wait ARQ is done by keeping a copy of the sent frame and retransmitting the frame when the timer expires.

Only 2 sequence numbers 0 and 1 are used.

Window size is 1.

No ACK for lost or damaged frames.

$$\text{Throughput} = \frac{\text{One packet}}{RTT}$$

$$\text{Utilization} = \frac{L}{L + BR}$$

$L$  = packet length

$B$  = Bandwidth

$R$  = RTT

If  $L < BR$ , Efficiency  $> 50$

$L > BR$ , Efficiency = 50

$$\mu = \frac{1}{1 + 2a}, \quad a = \frac{\text{propagation time}}{\text{Transmission time}}$$

Link utilization is low in stop and wait.

**GBN protocol** We can send several frames before receiving acknowledgements; we keep a copy of these frames with the acknowledgment.

- Sequence numbers ranges from  $2^m - 1$ .
- $m$  – number of bits for sequence numbers.
- The sender window slides one or more slots when a valid acknowledgment arrives.
- It uses cumulative acknowledgement or piggy backing wherever possible to acknowledge the frames.
- It discards duplicate and out of order packets.
- Receiving window size is 1.
- If the sender receives a NAK, it resends all frames in the sender window.

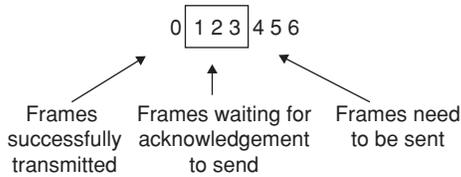
If a single packet is lost, damaged or acknowledgement is lost, it will resend all the packets.

$$\text{Link efficiency} = \frac{1 - p}{1 - p + p^w}$$

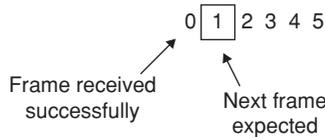
Where,  $p$  is the packet loss probability

$w$  is the sender's window size.

Sender:



Receiver:



- If  $N$  is maximum sequence number, then sender window size =  $N$ , Receiver window size = 1.
- If  $N$  is the number of sequence number, sender window size =  $N - 1$ , Receiver window size = 1.

**Selective repeat** More efficient for noisy links but processing at the receiver is more complex. Receiver window size is same as of sender window size. Sender window maximum size is  $2^{m-1}$ , receiver window maximum size is  $2^{m-1}$ . Sender and receiver window must be at most one half of  $2^m$ .

Receives out of order packets because receiver's window size is greater than 1.

It uses cumulative or independent or piggyback ACK whenever possible. If sender receives a NAK, it resends just the frame specified by the NAK.

If  $N$  is maximum sequence number,

$$\text{Sender window size} = \frac{N+1}{2},$$

$$\text{Receiver window size} = \frac{N+1}{2}$$

If  $N$  is the number of sequence numbers, sender window size =  $\frac{N}{2}$ , Receiver window size =  $\frac{N}{2}$ .

## MEDIUM ACCESS CONTROL SUBLAYER Multiplexing

When two communicating nodes are connected through a media, it generally happens that bandwidth of media is several times greater than that of the communicating nodes. Transferring of a single signal at a time is both slow and expensive. The whole capacity of the link is not utilized in this case. This link can be further exploited by sending several signals combined into one. This combining of signals into one is called multiplexing.

### Frequency Division Multiplexing (FDM)

This is possible in the case where transmission media has a bandwidth higher than the required bandwidth of signals to be transmitted. A number of signals can be transmitted at the same time. Each source is allotted a frequency range in which it can transfer its signals, and a suitable frequency gap is given between two adjacent signals to avoid overlapping. This type of multiplexing is commonly seen in the cable TV networks.

### Time Division Multiplexing (TDM)

This is possible when data transmission rate of the media is much higher than that of the data rate of the source. Multiple signals can be transmitted if each signal is allowed to be transmitted for a definite amount of time. These time slots are so small that all transmissions appear to be in parallel.

**Synchronous TDM** Time slots are pre. Assigned and are fixed. Each source is given its time slot at every turn due to it. This turn may be once per cycle or several turns per cycle, if it has a high data transfer rate, or may be once in a number of cycles if it is slow. This slot is given even if the source is not ready with data. So this slot is transmitted empty.

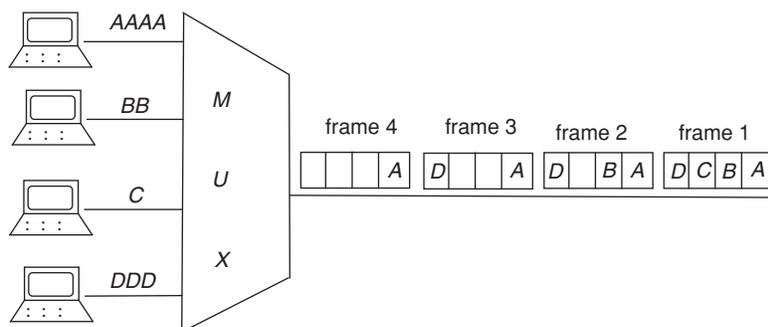


Figure 5 Synchronous TDM: Multiplexing process

**Asynchronous TDM** In this method, slots are not fixed. They are allotted dynamically depending on

speed of sources and whether they are ready for transmission.

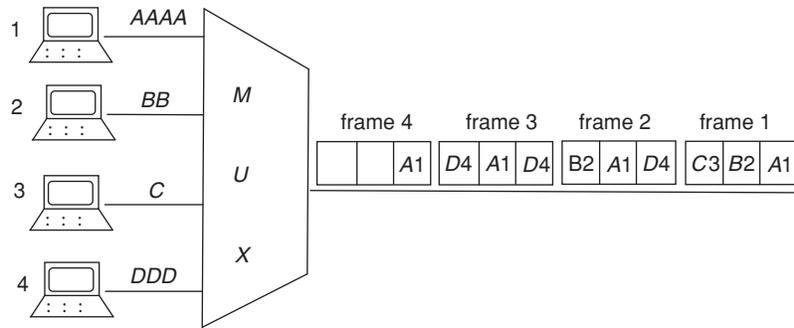


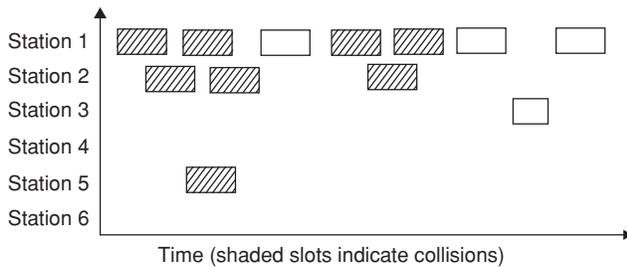
Figure 6 Asynchronous TDM

### Aloha Protocols

The Aloha Protocol was designed to provide data transmission between computers on several islands using radio transmission.

#### Pure aloha

Pure Aloha is an unslotted, fully decentralized protocol. It is extremely simple and trivial to implement. The ground rule is ‘when you want to talk, just talk!’ So, a node which wants to transmit, will go ahead and sends the packet on its broadcast channel, with no consideration of who so ever to any body else is transmitting (or) (not).



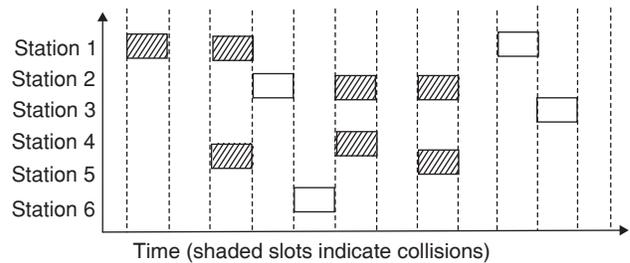
One serious drawback here is that, you don’t know whether what you are sending, has been received properly or not. To resolve this in pure Aloha, when one node finishes speaking it expects an acknowledgement in a finite amount of time otherwise it simply retransmits the data. This scheme works well in small networks where the load is not high. But in large, load intensive networks where many nodes may want to transmit at the same time, this scheme fails miserably. This led to the development of slotted Aloha.

#### Slotted Aloha

This is quite similar to pure Aloha, differing only in the way transmissions take place. Instead of transmitting right at the demand time, the sender waits for some time. This delay is specified as follows—the timeline is divided into equal slots and then it is required that transmission should take place only at slot boundaries. To be more precise, the slotted Aloha makes the following assumptions.

- All frames consist of exactly  $L$  bits.
- Time is divided into slots of size  $L/K$  seconds. (i.e., a slot equals the time to transmit one frame)

- Nodes start to transmit frames only at the beginning of slots.
- The nodes are synchronized so that each node knows when the slots begin.
- If two or more frames collide in a slot, then all the nodes detect the collision event before slot ends.



In this, way the number of collisions that can possibly take place is reduced by a huge margin. And hence, the performance became much better compared to pure Aloha. Collisions may only take place with nodes that are ready to speak at the same time.

### Virtual private network

Virtual Private Networking (VPN) Internet protocol security (IP sec) is one of the most complete, secure, standards-based protocol developed for transporting data.

A VPN is a shared network, where private data can be accessed only by the intended recipient.

The term VPN is used to describe a secure connection over the Internet.

VPN is also used to describe private networks such as Frame Relay and Asynchronous Transfer Mode (ATM).

The purpose of data security is that the data flowing across the network is protected by encryption technologies.

IP sec-based VPNs use encryption to provide data security, that increases the networks resistance to data tampering.

IP sec-based VPNs can be created over any type of IP Network, including Internet, ATM, Frame Relay, among all only Internet is inexpensive.

#### Uses of VPN

**Intranets** Intranets connect an organization’s locations. These locations could be head quarters offices, branch offices, Employees home which is located in some Remote area.

This connectivity is used for e-mails, sharing files etc.

The cost of connecting remote home users is very expensive compared to Internet access technologies because of this organizations have moved their networks to the Internet.

**Remote access** It enables telecommuters and mobile workers to access e-mail and business applications.

A dial-up connection to an organizations modem pool is one method to access remote workers. It is expensive, because of long distance telephone and service costs.

**IP sec**

IP sec is an Internet Engineering Task Force (IETF) standard suite of protocols that provide data authentication, integrity, and confidentiality between 2 communication points across IP-Network.

It provides data security at the IP-packet level.

IP sec protects against possible security exposures by protecting data while in transit.

**Features**

IP sec was designed to provide the following security features when transferring packets across networks.

1. Authentication: Verifies that the packet received is actually from the correct sender or not.
2. Integrity: Ensures that the contents of packet did not change while transmitting data.
3. Confidentiality: Conceals the message content through encryption.

**Components of IP sec**

**ESP:** (Encapsulating security payload), It provides confidentiality, authentication and integrity.

**AH:** (Authentication Header) provides Authentication and Integrity.

**IKE:** (Internet key Exchange) provides key management and security Association (SA) management.

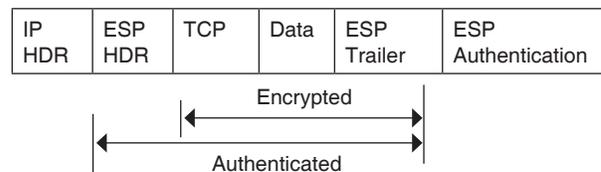
**ESP:**

- Most importantly, it provides message content protection.
- IP sec provides an open frame work for implementing standard algorithms such as SHA and MD5.

- The algorithms IP sec uses produces a unique identifier for each packet, which is a data equivalent to a finger print.
- This Finger Print allows the device to determine whether a packet has been tampered with.
- Packets that are not authenticated are discarded and not delivered to the intended receiver
- ESP also provides all encryption services in IP sec.
- Encryption/decryption allows only the sender and the authorized receiver to read the data.
- The authentication performed by ESP is called ESP authentication.
- ESP provides authentication and integrity for the payload and not for the IP-header



Figure 7 Original packet



The ESP Header is inserted into the packet between the IP-header and any subsequent packet contents.

- ESP encrypts the data, the payload is changed.
- ESP does not encrypt the ESP header, nor does it encrypt the ESP authentication.

**AH:**

Provides optional anti-replay protection, which protects against unauthorized retransmission of packets.

The authentication header is inserted into the packet between the IP-header and any sub sequent packet contents.

AH does not protect the data's confidentiality.

For added protection in certain cases, AH and ESP can be used together.



Original packet

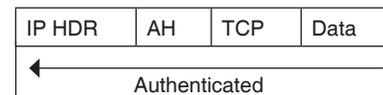


Figure 8 Packet with IP sec Authentication Header.

**EXERCISES**

**Practice Problems I**

**Directions for questions 1 to 15** Select the correct alternative from the given choices.

1. Assume that, in a stop and wait ARQ system, the bandwidth of the line is 1 mbps, and 1 bit takes 20 ms to make a round trip. What is the bandwidth delay product utilization percentage of the link if we send 1000 bits?
 

|         |         |
|---------|---------|
| (A) 1%  | (B) 5%  |
| (C) 10% | (D) 50% |

2. A channel has a bit rate of 20 kbps and propagation delay of 100 msec. For what size does stop and wait gives an efficiency of 50%?
 

|               |               |
|---------------|---------------|
| (A) 2000 bits | (B) 3000 bits |
| (C) 4000 bits | (D) 6000 bits |
3. CSMA/CD LAN of 1 gbps is to be designed over 1 km cable without repeater. The minimum frame size that Data link layer should consider, if cable support signal speed of 20,000 km/sec

- (A) 10 k bits                      (B) 20 k bits  
(C) 30 k bits                      (D) 40 k bits
4. A 20 mbps satellite link has a propagation delay of 400  $\mu$ s. The transmitter employs the 'go-back-n ARQ' scheme with  $n$  set to 10. Assuming each frame is 100 bytes long. What is the maximum data rate possible?  
(A) 1 mbps                      (B) 2 mbps  
(C) 5 mbps                      (D) 10 mbps
5. A satellite channel has capacity of  $B$  bits/sec, the frame size is of  $L$  bits, and round trip propagation time of  $R$  sec, uses stop and wait protocol, what is the channel utilization?  
(A)  $\frac{L}{L - BR}$                       (B)  $\frac{L}{L + BR}$   
(C)  $\frac{L}{B + R}$                       (D)  $\frac{L}{B - R}$
6. Find efficiency of the ring where data rate of link is 4 mbps, number of stations are 20 separated by 100 meters and bit delay in each station is 2.5 bits. (velocity of propagation =  $2 \times 10^8$  m/s)  
(A) 60 bits                      (B) 75 bits  
(C) 90 bits                      (D) 120 bits
7. If you are designing sliding window protocol of 1 mbps which has one way delay of 1.25 seconds. Assuming each frame carries 1 kB of data, what is the minimum number of bits you need for the sequence number?  
(A) 8                      (B) 9  
(C) 10                      (D) 12
8. What are the sequence numbers of sender and receiver windows in Go-back-n and selective repeat if  $m$ -bits are used?  
(A)  $2^m - 1, 1, 2^{m-1}, 2^{m-1}$                       (B)  $2^m, 1, 2^{m-1}, 2^{m-1}$   
(C)  $2^m, 2, 2^m, 2^m$                       (D)  $2^m - 1, 1, 2^m, 2^m$
9. A 100 km long cable runs at 1.536 mbps. The propagation speed in the cable is  $2/3$  of speed of light. Number of bits fit in the cable would be?  
(A) 428 bits                      (B) 526 bits  
(C) 672 bits                      (D) 768 bits
10. If the bandwidth of the link is 256 mbps, Assume that sequence number field consists 32 bits. Find the wrap around time for sequence numbers?  
(A) 128 sec                      (B) 256 sec  
(C) 512 sec                      (D) 1024 sec
11. After a series of collisions a station has selected slot 984. In how many successive collisions, the station was a part of communication?  
(A) 4                      (B) 6  
(C) 8                      (D) 10
12. There are 10 stations in a LAN always having constant load and ready to transmit. During any particular contention slot each station transmits with a probability of 0.1. If the average frame takes 122 ms to transmit, what is the channel efficiency, if round trip time is 51.2 micro secs?  
(A) 0.23                      (B) 0.35  
(C) 0.48                      (D) 0.56
13. Which of the below are issues concerning data link layer?  
(i) Ensures that the transmission facility is free of undetected transmission errors  
(ii) Regulates the transmission rates so as to match the receiver's capabilities  
(iii) Ensures the design of the line such that when a '1' bit is sent it is always received as '1' bit at receivers end.  
(A) (i), (ii)                      (B) (ii), (iii)  
(C) (iii), (i)                      (D) (i), (ii), (iii)
14. An Ethernet LAN has the capability of 100 Mbps. If Manchester encoding is used, what is the rate of signal change?  
(A) 20 million times/sec  
(B) 200 million times/sec  
(C) 50 million times/sec  
(D) 500 million times/sec
15. For 10 Mbps LAN it is found that 64 bytes is the minimum frame size to aid in collision detection. What should be the minimum frame size for a 100 Mbps LAN?  
(A) 6.4 bytes                      (B) 64 bytes  
(C) 640 bytes                      (D) 6400 bytes

## Practice Problems 2

**Directions for questions 1 to 15** Select the correct alternative from the given choices.

1. What is the probability of success for any arbitrary station among ' $N$ ' stations to transmit in CSMA/CD?  
(A)  $Np_s(1 - p_s)^N$                       (B)  $(N - 1)p_s(1 - p_s)$   
(C)  $Np_s(1 - p_s)^{N-1}$                       (D)  $Np_s(1 - p_s)^N$
2. If 4-bits are used to represent sequence numbers for flow control. What are sender and receiver window sizes in Go-back-n and selective repeat?  
(A) 16, 1, 8, 8                      (B) 15, 1, 8, 8  
(C) 15, 2, 8, 8                      (D) 15, 1, 16, 8
3. If the available maximum sequence number is 13, compute sender and receiver window sizes in go-back-n and selective repeat?  
(A) 4, 1, 4, 4                      (B) 4, 1, 7, 7  
(C) 13, 1, 7, 7                      (D) 13, 1, 4, 4
4. In a gigabit ethernet LAN, the receiver couldn't empty the input buffer on some line for 1 millisecond. What is the maximum accumulation of frames possible neglecting propagation delays?

- (A) 1024 frames (B) 2097 frames  
(C) 4096 frames (D) 5120 frames
5. A Token ring LAN is using differential Manchester encoding. If the LAN speed is 10 Mbps. What is the baud rate?  
(A) 10 M baud (B) 20 M baud  
(C) 5 M baud (D) 100 M baud
6. Consider a 100-meter 10 mbps token ring containing 10 stations, each transmitting with equal priority. Each station can transmit 4 bytes before giving up the token. Token holding time per station is 10 ns. Also propagation speed is 200 m/s. Assume that the Ring monitor has created a new token, how long does it take for the token to come back to the Ring monitor if no station uses the token?  
(A) 2.55  $\mu$ sec (B) 3.64  $\mu$ sec  
(C) 4.65  $\mu$ sec (D) 2.93  $\mu$ sec
7. In the above question, if only 6 nodes including Ring monitor are active what is total propagation delay in  $\mu$ sec?  
(A) 3.60 (B) 3.61  
(C) 3.62 (D) 3.63
8. In the above case if bit regeneration time is 1 ns/bit. What is the regeneration overhead caused if a 4 kB token is taken by 1st node and if it uses to transmit 4B data to ring monitor.  
(A) 412 ns (B) 544 ns  
(C) 640 ns (D) 800 ns
9. Which of the below operation is applied to full-duplex mode operation of gigabit Ethernet?  
(i) Traffic is allowed in both directions at any time.  
(ii) CSMA/CD protocol is used.  
(iii) Maximum length of cable segment used to connect stations is limited by CSMA/CD protocol.  
(A) (i) and (ii) (B) (ii), (iii)  
(C) (iii), (i) (D) (i), (ii), (iii)
10. Which of the below are not applied to Token Ring networks?  
(i) Collisions  
(ii) Limits on length of the cable segment  
(iii) Time slots for transmission  
(iv) Usage of repeaters  
(A) (i), (ii) (B) (ii), (iii)  
(C) (iii), (iv) (D) (i), (iv)
11. Select the correct statements from below (pertaining to Ethernet):  
(i) Frame collisions don't occur at a repeater  
(ii) Frame collisions can occur at the hub itself  
(iii) Switch frames are never lost due to collisions  
(iv) Entire bridge is a point of collisions  
(A) (i), (ii), (iii) (B) (ii), (iii), (iv)  
(C) (i), (iii), (iv) (D) (i), (ii), (iv)
12. Match the different layers with possible security methods in those layers.  
(p) Data link layer (i) user authentication  
(q) Network layer (ii) use firewalls  
(r) Transport layer (iii) encryption of connections  
(s) Application layer (iv) point to point encryption of data stream  
(A) p – i, q – ii, r – iii, s – iv  
(B) p – ii, q – iv, r – iii, s – i  
(C) p – iv, q – ii, r – iii, s – i  
(D) p – iii, q – iv, r – ii, s – i
13. The hamming distance between 001111 and 010011 is  
(A) 1 (B) 2  
(C) 3 (D) 4
14. Which of the following represents the polynomial  $x^5 + x^4 + x^0$  using the CRC?  
(A) 110000 (B) 110001  
(C) 110010 (D) 110101
15. For a sliding window of size  $n-1$  ( $n$  sequence number) there can be maximum of how many frames sent but unacknowledged?  
(A) 0 (B)  $n - 1$   
(C)  $n$  (D)  $n + 1$

### PREVIOUS YEARS' QUESTIONS

1. The message 11001001 is to be transmitted using the CRC polynomial  $x^3 + 1$  to protect it from errors. The message that should be transmitted is: [2007]  
(A) 11001001000 (B) 11001001011  
(C) 11001010 (D) 110010010011
2. The distance between two stations  $M$  and  $N$  is  $L$  kilometers. All frames are  $K$  bits long. The propagation delay per kilometer is  $t$  seconds. Let  $R$  bits/second be the channel capacity. Assuming that processing delay is negligible, the minimum number of bits for the sequence number field in a frame for maximum

utilization, when the sliding window protocol is used, is: [2007]

- (A)  $\left\lceil \log_2 \frac{2LtR + 2k}{k} \right\rceil$  (B)  $\left\lceil \log_2 \frac{2LtR}{k} \right\rceil$   
(C)  $\left\lceil \log_2 \frac{2LtR + k}{k} \right\rceil$  (D)  $\left\lceil \log_2 \frac{2LtR + k}{2k} \right\rceil$

3. Match the following:

- (P) SMTP (1) Application layer  
(Q) BGP (2) Transport layer

- (R) TCP (3) Data link layer  
 (S) PPP (4) Network layer  
 (5) Physical layer

[2007]

- (A) P-2 Q-1 R-3 S-5  
 (B) P-1 Q-4 R-2 S-3  
 (C) P-1 Q-4 R-2 S-5  
 (D) P-2 Q-4 R-1 S-3

4. A layer-4 firewall (a device that can look at all protocol headers up to the transport layer) CANNOT

[2011]

- (A) Block entire HTTP traffic during 9.00PM and 5.00PM  
 (B) Block all ICMP traffic  
 (C) Stop incoming traffic from a specific IP address but allow outgoing traffic to the same IP address  
 (D) Block TCP traffic from a specific user on a multi-user system during 9.00PM and 5.00AM.

5. If two fair coins are flipped and at least one of the outcomes is known to be a head, what is the probability that both outcomes are heads? [2011]

- (A) 1/3 (B) 1/4  
 (C) 1/2 (D) 2/3

6. Which of the following transport layer protocols is used to support electronic mail? [2012]

- (A) SMTP (B) IP  
 (C) TCP (D) UDP

7. Consider a source computer ( $S$ ) transmitting a file of size  $10^6$  bits to a destination computer ( $D$ ) over a network of two routers ( $R_1$  and  $R_2$ ) and three links ( $L_1$ ,  $L_2$  and  $L_3$ ).  $L_1$  connects  $S$  to  $R_1$ ;  $L_2$  connects  $R_1$  to  $R_2$ ; and  $L_3$  connects  $R_2$  to  $D$ . Let each link be of length 100 km. Assume signals travel over each link at a speed of  $10^8$  meters per second. Assume that the link bandwidth on each link is 1 Mbps. Let the file be broken down into 1000 packets each of size 1000 bits. Find the total sum of transmission and propagation delays in transmitting the file from  $S$  to  $D$ ? [2012]

- (A) 1005 ms (B) 1010 ms  
 (C) 3000 ms (D) 3003 ms

8. Determine the maximum length of the cable (in km) for transmitting data at a rate of 500 Mbps in an Ethernet LAN with frames of size 10,000 bits. Assume the signal speed in the cable to be 2,00,000 km/s. [2013]

- (A) 1 (B) 2  
 (C) 2.5 (D) 5

9. Consider a token ring network with a length of 2 km having 10 stations including a monitoring station. The propagation speed of the signal is  $2 \times 10^8$  m/s and the token transmission time is ignored. If each station is allowed to hold the token for 2  $\mu$ sec, the minimum time for which the monitoring station should wait (in  $\mu$ sec) before assuming that the token is lost is \_\_\_\_\_

[2014]

10. Consider the store and forward packet switched network given below. Assume that the bandwidth of each link is  $10^6$  bytes/sec. A user on host  $A$  sends a file of size  $10^3$  bytes to host  $B$  through routers  $R_1$  and  $R_2$  in three different ways. In the first case a single packet containing the complete file is transmitted from  $A$  to  $B$ . In the second case, the file is split into 10 equal parts, and these packets are transmitted from  $A$  to  $B$ . In the third case, the file is split into 20 equal parts, and these packets are sent from  $A$  to  $B$ . Each packet contains 100 bytes of header information along with the user data. Consider only transmission time and ignore processing, queuing and propagation delays. Also assume that there are no errors during transmissions. Let  $T_1$ ,  $T_2$  and  $T_3$  be the times taken to transmit the file in the first, second and third case respectively. Which one of the following is CORRECT? [2014]



- (A)  $T_1 < T_2 < T_3$   
 (B)  $T_1 > T_2 > T_3$   
 (C)  $T_2 = T_3, T_3 < T_1$   
 (D)  $T_1 = T_3, T_3 > T_2$

11. In the following pairs of OSI protocol layer/ sub-layer and its functionality the INCORRECT pair is [2014]

- (A) Network layer and Routing  
 (B) Data Link Layer and Bit synchronization  
 (C) Transport layer and End-to-end process communication  
 (D) Medium Access Control sub-layer and Channel sharing

12. A bit-stuffing based framing protocol uses an 8-bit delimiter pattern of 01111110. If the output bit-string after stuffing is 01111100101, then the input bit-string is [2014]

- (A) 0111110100 (B) 0111110101  
 (C) 0111111101 (D) 0111111111

13. Suppose that the stop-and-wait protocol is used on a link with a bit rate of 64 kilobits per second and 20 milliseconds propagation delay. Assume that the transmission time for the acknowledgement and the processing time at nodes are negligible. Then the minimum frame size in bytes to achieve a link utilization of at least 50% is \_\_\_\_\_ [2015]

14. A link has a transmission speed of  $10^6$  bits/sec. It uses data packets of size 1000 bytes each. Assume that the acknowledgement has negligible transmission delay, and that its propagation delay is the same as the data propagation delay. Also assume that the processing delays at nodes are negligible. The efficiency of the stop-and-wait protocol in this setup is exactly 25%. The value of the one-way propagation delay (in milliseconds) is \_\_\_\_\_ [2015]

15. Consider a CSMA/CD network that transmits data at a rate of 100 Mbps ( $10^8$ -bits per second) over a 1 km (kilometer) cable with no repeaters. If the minimum frame size required for this network is 1250 bytes, what is the signal speed (km/sec) in the cable? [2015]
- (A) 8000 (B) 10000  
(C) 16000 (D) 20000
16. Consider a LAN with four nodes  $S_1$ ,  $S_2$ ,  $S_3$  and  $S_4$ . Time is divided into fixed-size slots, and a node can begin its transmission only at the beginning of a slot. A collision is said to have occurred if more than one node transmit in the same slot. The probabilities of generation of a frame in a time slot by  $S_1$ ,  $S_2$ ,  $S_3$  and  $S_4$  are 0.1, 0.2, 0.3 and 0.4, respectively. The probability of sending a frame in the first slot without any collision by any of these four stations is \_\_\_\_\_. [2015]
17. Consider a network connecting two systems located 8000 kilometers apart. The bandwidth of the network is  $500 \times 10^6$ -bits per second. The propagation speed of the media is  $4 \times 10^6$  meters per second. It is needed to design a Go-Back-N sliding window protocol for this network. The average packet size is  $10^7$ -bits. The network is to be used to its full capacity. Assume that processing delays at nodes are negligible. Then, the minimum size in bits of the sequence number field has to be \_\_\_\_\_. [2015]
18. Two hosts are connected via a packet switch with  $10^7$  bits per second links. Each link has a propagation delay of 20 microseconds. The switch begins forwarding a packet 35 microseconds after it receives the same. If 10000-bits of data are to be transmitted between the two hosts using a packet size of 5000-bits, the time elapsed between the transmission of the first bit of data and the reception of the last bit of the data in microseconds is \_\_\_\_\_. [2015]
19. A sender uses the Stop-and-Wait ARQ protocol for reliable transmission of frames. Frames are of size 1000 bytes and the transmission rate at the sender is 80 Kbps (1 Kbps = 1000 bits/second). Size of an acknowledgement is 100 bytes and the transmission rate at the receiver is 8 Kbps. The one-way propagation delay is 100 milliseconds. Assuming no frame is lost, the sender throughput is \_\_\_\_\_ bytes/second. [2016]
20. In an Ethernet local area network, which one of the following statements is **TRUE**? [2016]
- (A) A station stops to sense the channel once it starts transmitting a frame.  
(B) The purpose of the jamming signal is to pad the frames that are smaller than the minimum frame size.  
(C) A station continues to transmit the packet even after the collision is detected  
(D) The exponential back off mechanism reduces the probability of collision on retransmissions.
21. Consider a  $128 \times 10^3$  bits/second satellite communication link with one way propagation delay of 150 milliseconds. Selective retransmission (repeat) protocol is used on this link to send data with a frame size of 1 kilobyte. Neglect the transmission time of acknowledgement. The minimum number of bits required for the sequence number field to achieve 100% utilization is \_\_\_\_\_. [2016]
22. A computer network uses polynomials over  $GF(2)$  for error checking with 8 bits as information bits and uses  $x^3 + x + 1$  as the generator polynomial to generate the check bits. In this network, the message 01011011 is transmitted as [2017]
- (A) 01011011010 (B) 01011011011  
(C) 01011011101 (D) 01011011100
23. The values of parameters for the Stop-and Wait ARQ protocol are as given below:
- Bit rate of the transmission channel = 1Mbps.  
Propagation delay from sender to receiver = 0.75 ms.  
Time to process a frame = 0.25 ms.  
Number of bytes in the information frame = 1980.  
Number of bytes in the acknowledge frame = 20.  
Number of overhead bytes in the information frame = 20.
- Assume that there are no transmission errors. Then, the transmission efficiency (expressed in percentage) of the stop-and-wait ARQ protocol for the above parameters is \_\_\_\_\_ (correct to 2 decimal places) [2017]
24. Consider a binary code that consists of only four valid codewords as given below:
- 00000,01011,10101,11110
- Let the minimum Hamming distance of the code be  $p$  and the maximum number of erroneous bits that can be corrected by the code be  $q$ . Then the values of  $p$  and  $q$  are [2017]
- (A)  $p=3$  and  $q=1$   
(B)  $p=3$  and  $q=2$   
(C)  $p=4$  and  $q=1$   
(D)  $p=4$  and  $q=2$
25. Consider two hosts  $X$  and  $Y$  connected by a single direct link of rate  $10^6$  bits/sec. The distance between the two hosts is 10.000 km and the propagation speed along the link is  $2 \times 10^8$  m/sec. Host  $X$  sends a file of 50,000 bytes as one large message to host  $Y$  continuously. Let the transmission and propagation delays be  $p$  milliseconds and  $q$  milliseconds, respectively. Then the values of  $p$  and  $q$  are [2017]
- (A)  $p=50$  and  $q=100$   
(B)  $p=50$  and  $q=400$

- (C)  $p=100$  and  $q=50$   
 (D)  $p=400$  and  $q=50$

26. Consider a simple communication system where multiple nodes are connected by a shared broadcast medium (like Ethernet or wireless). The nodes in the system use the following carrier-sense based medium access protocol. A node that receives a packet to transmit will carrier-sense the medium for 5 units of time. If the node does not detect any other transmission in this duration, it starts transmitting its packet in the next time unit. If the node detects another transmission, it waits until this other transmission finishes, and then begins to carrier-sense for 5 time units again. Once they start to transmit, nodes do not perform any

collision detection and continue transmission even if a collision occurs. All transmissions last for 20 units of time. Assume that the transmission signal travels at the speed of 10 meters per unit time in the medium.

Assume that the system has two nodes  $P$  and  $Q$ , located at a distance  $d$  meters from each other.  $P$  starts transmitting a packet at time  $t = 0$  after successfully completing its carrier-sense phase. Node  $Q$  has a packet to transmit at time  $t = 0$  and begins to carrier-sense medium.

The maximum distance  $d$  (in meters, rounded to the closest integer) that allows  $Q$  to successfully avoid a collision between its proposed transmission and  $P$ 's ongoing transmission is \_\_\_\_\_. [2018]

## ANSWER KEYS

### EXERCISES

#### Practice Problems I

1. B    2. C    3. A    4. D    5. B    6. C    7. B    8. A    9. D    10. A  
 11. D    12. C    13. A    14. B    15. C

#### Practice Problems II

1. C    2. B    3. C    4. B    5. B    6. B    7. A    8. B    9. D    10. B  
 11. A    12. D    13. C    14. B    15. B

#### Previous Years' Questions

1. B    2. C    3. B    4. A    5. A    6. C    7. A    8. B    9. 28 to 30  
 10. D    11. B    12. B    13. 160    14. 12    15. D    16. 0.40 to 0.46    17. 8  
 18. 1575    19. 2500    20. D    21. 4    22. C    23. 87.3    24. A    25. D    26. 50

## Routing Algorithms

### LEARNING OBJECTIVES

- Routing algorithm basics
- Flooding
- Multipath routing
- Distance vector routing
- Link state routing
- Hierarchical routing
- Rip
- Ospf
- Congestion control techniques
- Traffic shaping

### ROUTING ALGORITHMS BASICS

The main function of network layer is routing packets from the source machine to the destination machine. The routing algorithms are part of the network layer software, responsible for deciding which output line an incoming packet should be transmitted on.

Routing algorithms can be grouped into two major classes: Non-adaptive and Adaptive.

1. Non-adaptive algorithms do not base their routing decisions on measurements or estimates of the current traffic and topology. Instead, the choice of the route to use is downloaded to the routers when the network is booted. This procedure is called static routing.
2. Adaptive algorithms, in contrast, change their routing decisions to reflect changes in the topology, and the traffic as well.

### Store and Forward Packet Switching

In this Technique, the data packet will be stored at the node and it is forwarded to its next appropriate intermediate node. The next intermediate node will first store the packet in the buffer, based on the router decision, it selects an interface, and forwards to receiver.

The technique is most suitable for the networks with unsteady connectivity.

The length of the packet we take shows effect on the file transfer, if the data packet is small, in the store the forward, delay will be less at each node, but causes extra overhead with headers. So, the packet size selection should be done appropriately.

### FLOODING

Static algorithms, in which every incoming packet is sent out on every outgoing line except the one on which it is arrived. Header contains the hop count of each packet. Hop counter is decremented at each hop, with the packet being discarded when the counter reaches zero.

Another way for damming the flood is to keep track of which packets have been flooded, to avoid sending them out a second time. A variation of flooding that is slightly more practical is selective flooding. In this algorithm the routers do not send every incoming packet out on every line, only on those lines that are going approximately in the right direction.

### Multipath Routing

Multipath routing is routing the packets from the source, on multiple paths to the destination. It is nothing but spreading the traffic.

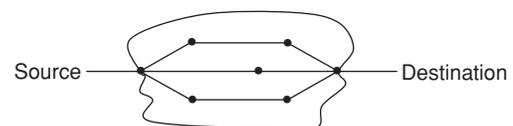


Figure 1 Multipath routing model

Single path routing causes QoS, throughput and delay problems, and multipath routing, improves network performance with sharing of available resources of network.

The components of multipath routing are

1. Multipath calculation algorithm
2. Multipath forwarding algorithm
3. End-Host protocol

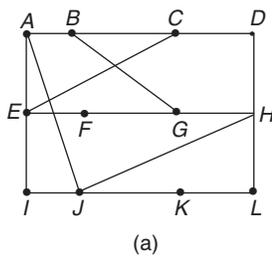
The algorithms specified above are based on Dijkstra's shortest path algorithm they generate paths according to path characteristics and ensure path quality and path independence.

The end-host protocol uses the multipath (determined) effectively performance will be improved if end-users use the multiple paths effectively.

### DISTANCE VECTOR ROUTING

A dynamic routing algorithm, operates by having each router maintain a table (i.e., a vector) giving the best known distance to each destination and which line to use to get there. These tables are updated by exchanging information with the neighbors.

The Metric used might be number of hops, time delay in milliseconds, and total number of packets queued along the path or something similar.



|   | New estimated |    |    |    | Delay from |   |
|---|---------------|----|----|----|------------|---|
|   | A             | I  | H  | K  | J          |   |
| A | 0             | 24 | 20 | 21 | 8          | A |
| B | 12            | 36 | 31 | 28 | 20         | A |
| C | 25            | 18 | 19 | 36 | 28         | I |
| D | 40            | 27 | 8  | 24 | 20         | H |
| E | 14            | 7  | 30 | 22 | 17         | I |
| F | 23            | 20 | 19 | 40 | 30         | I |
| G | 18            | 31 | 6  | 31 | 18         | H |
| H | 17            | 20 | 0  | 19 | 12         | H |
| I | 21            | 0  | 14 | 22 | 10         | I |
| J | 9             | 11 | 7  | 10 | 0          | - |
| K | 24            | 22 | 22 | 0  | 6          | K |
| L | 29            | 33 | 9  | 9  | 15         | k |

| JA delay is | JI delay is | JH delay is | JK delay is | New routing table for J<br>Vectors received from J's four neighbours. |
|-------------|-------------|-------------|-------------|-----------------------------------------------------------------------|
| 8           | 10          | 12          | 6           |                                                                       |

Figure 2 (a) Subnet, (b) Delay vectors of J

Figure 2(a) 'shows a subnet. The first 4 columns of figure 2(b) shows the delay vectors received from the neighbors of router J. A claims to have a 12 m sec delay to B, a 25 m sec delay to C, a 40 m sec delay to D, etc.

Suppose that J has estimated its delay to its neighbour, A, J, H and K as 8, 10, 12 and 6 m sec, respectively.

Now J computes its new route to router G. It knows that it can get to A in 8 m sec, and A claims to be able to get to G in 18 m sec, so J knows it can count on a delay of 26 m sec to G if it forwards packets bound for G to A, similarly, it computes the delay to G via I, H and K as 41 (31 + 10), 18 (6 + 12) and 37 (31 + 6) m sec, respectively.

The best of these values is 18, so it makes an entry in its routing table that the delay to G is 18 m sec and that the route to use is via H.

### Count to Infinity Problem

It reacts rapidly to good news, but leisurely to bad news. Actual network may be down but routers will exchange routes with one another.

Following measures are taken to avoid count-to-infinity problem:

- Hop limit:** Limit number of hops normally 0 hops directly connected, hop 16 is (0–15), 16 hops unreachable.
- Split horizon:** Never send information back in direction where it came from.
- Route poisoning and poison reverse, hold on timer trigger.
- As soon as network goes down, make metric of root infinity to resolve the immediate instability created because of routing updates from neighbor.
- When router sends update with infinite metric to neighbor, neighbor will make it down.
- Now routers will initiate hold on time to learn alternate paths and send update in direction where it came (Poison reverse) from.
- Routers will incorporate final roots in routing table.

### LINK STATE ROUTING

The idea behind link state routing is simple and can be stated as five parts. Each router must do the following:

- Discover its neighbors and learn their network addresses.
- Measure the delay or cost to each of its neighbors.
- Construct a packet telling all it has just learned.
- Send this packet to all other routers.
- Compute the shortest path to every other router.

**Learning about the neighbors** When a router is booted, its first task is to learn who its neighbors are. It accomplishes this goal by sending a special HELLO packet on each point to point line. The router on the other end is expected to send back a reply telling who it is. These names must be globally unique.

**Measuring the cost** The link state routing algorithm requires each router to know, or at least have a reasonable estimate of, the delay to each of its neighbors.

The most direct way to determine this delay is to send over the line a special ECHO packet that the other side is required to send back immediately.

By measuring the round trip time and dividing it by two, the sending router can get a reasonable estimate of the delay. If two paths with same bandwidth exists and one path is heavily loaded then the path which is not heavily loaded is chosen. But this may oscillate in the choice of best path. So to avoid oscillation in the choice of best path, distribute the load over multiple lines with same known fraction going over each line.

**Building link state packets** Once the information needed for the exchange has been collected, the next step is, for each router to build a packet containing all the data. The packet starts with identity of the sender, followed by a sequence number and age, and a list of neighbors. For each neighbor, delay to that neighbor is given.

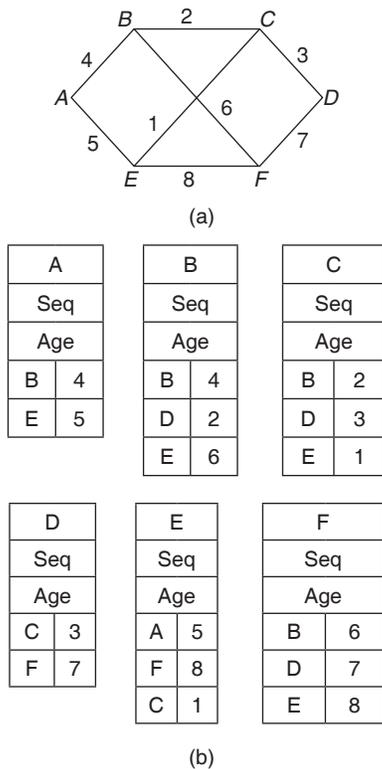


Figure 3 (a) Subnet5, (b) Link state packets for this subnet.

**Distributing the link state packets** As the packets are distributed and installed, the routers getting the routing packet first will change their routes.

Consequently, the different routers may be using different versions of the topology, which can lead to inconsistencies, loops, unreachable machines, and other problems. The fundamental idea is to use flooding to distribute the link state packets. To keep the flood in check, each packet contains a sequence number that is incremented for each new packet sent. Routers keep track of all the (source router, sequence) pairs they see.

When a new link state packet comes in, it is checked against the list of packets already seen. If it is new, it is forwarded on all lines except the one it arrived on. If it is a duplicate, it is discarded.

If a packet with a sequence number lower than the highest one seen so far ever arrives, it is rejected as being obsolete since the router has more recent data.

If a router ever crashes it will lose track of its sequence number. If its starts again at 0, the next packet will be rejected as a duplicate. Also due to bit error, packets may be rejected as obsolete. Solution to these problems is to include the age of each packet after the sequence number and decrement it once per second.

When the age hits zero, the information from that router is discarded.

**Computing new routes** Once a router has accumulated a full set of link state packets, it can construct the entire subnet graph because every link is represented. Every link is, in fact, represented twice, once for each direction. The two values can be averaged or used separately. Now dijkstra's algorithm can be run locally to construct the shortest path to all possible destinations.

### Hierarchical Routing

Hierarchical Routing is mainly designed for large topologies. With increase in the topology there is proportionate increase in the routing tables, which consume more memory for maintaining tables and requires more bandwidth for the status reports.

In this routing, network topology is divided into hierarchies, these will reduce size of routing table. The node at each hierarchy will know about the nodes present in that level. It forwards the packet to its border router (at its level) if destination is not at its level. Hierarchical routing increases efficiency in routing, less traffic, reduction of table size in an order of about (log n).

### RIP

1. It calculates best route based on hop count.
2. RIP cannot handle more than 15 hops, anything above 15 hops away is considered unsearchable by RIP. This fact is used by RIP to prevent routing loops.
3. RIP is a classful routing protocol.
4. Interval between route update advertisements: 30 sec. Time out/hold on times: 180 sec
5. RIP implements the split horizon, route isonning and hold down mechanisms to prevent looping.
6. It is a dynamic distance vector routing protocol.

### OSPF

The open shortest path first is an adaptive routing protocol for IP networking. It uses a link state routing algorithm. OSPF keeps track of the state of all the various network connections between itself and a network it is trying to send

data to. OSPF selects the best route by finding the lowest cost paths to a destination. All router interfaces are given a cost. Its domain is an autonomous system.

**Backbone routers** Backbone routers have one or more interfaces in Area 0 (the backbone area).

**Area border router (ABR)** Routers that belong to multiple areas, and connect these areas to the backbone area are called ABR. It has interfaces in multiple areas.

**Autonomous system boundary router (ASBR)** If the router connects the OSPF autonomous system to another autonomous system, it is called an autonomous system boundary router (ASBF).

OSPF elects two or more routers to manage the link state advertisements.

**Designated router (DR)** Every OSPF will have a DR, a backup DR. The DR is the route to which all other routers within the area, send their link state advertisements.

## OSPF areas

OSPF areas are used to impose a hierarchical structure to the flow of data over the network. A network using OSPF will always have at least one area and if there is more than one area, one of the two areas must be the backbone area. Areas are used to group routers into manageable groups that exchange routing information locally, but summarize the routing information, when advertising the routes externally, ABR's are used to connect the areas.

## CONGESTION CONTROL TECHNIQUES

Objective of congestion control technique is to limit queue lengths at the nodes, so as to avoid throughput collapse.

1. Send a control packet from a congested node to some or all source nodes to stop or slow the rate of transmission from source and thus limit the total number of packets in the network.
2. Allow packet switching nodes to add congestion information to packets as they pass by. The packets carrying such information can go in both the directions i.e., opposite of the congestion and in the same direction of the congestion.

Packets in the opposite direction of congestion quickly reach the source node which can reduce the flow of packets into the network.

Packets going in the same direction as the congestion, reach the destination. The destination asks the source to adjust the load by returning the signal back to the source in the packets.

3. Provides link delay information to other nodes. This information can be used to influence the rate at which new packets are produced. As these delays are influenced by the routing decision, they may vary too rapidly to use effectively for congestion control.

## Congestion Control

Congestion control maintains the number of packets within the network below the level at which performance falls dramatically.

Every node has a queue of packets for each outgoing channel. If, rate at which packets arrive and queue up, exceeds the rate of packet transmission, then size of queue grows without bound and thus delay experienced by a packet goes to infinity.

When the packets arrive they are stored in the input buffer, of the corresponding link. The node examines each incoming packet to make a routing decision and then moves the packet to the appropriate output buffer. Packet queued up for output in output buffer is transmitted as soon as possible. When saturation point is reached, one can do any of the following:

1. Discard incoming packet for which there is no available buffer space.
2. Node should exercise some sort of flow control over its neighbors so that the traffic flow remains manageable.
3. Traffic shaping is about regulating the average rate of data transmission.

## Leaky Bucket Algorithm

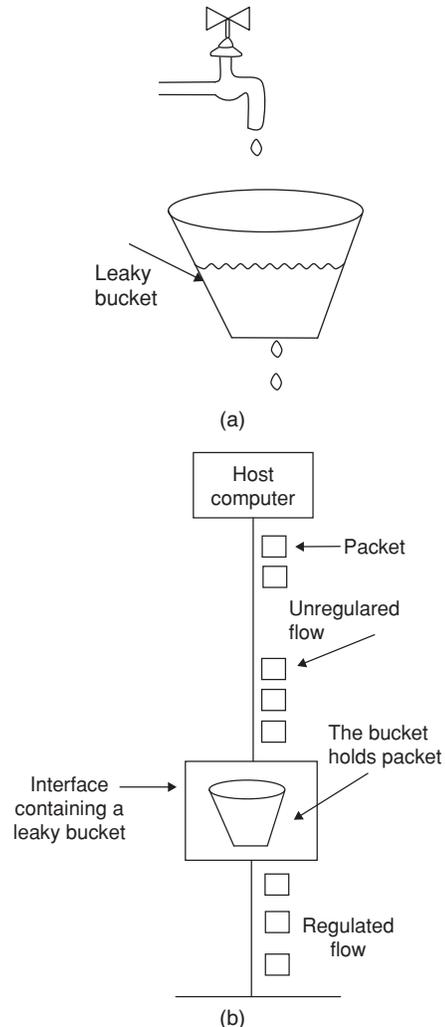
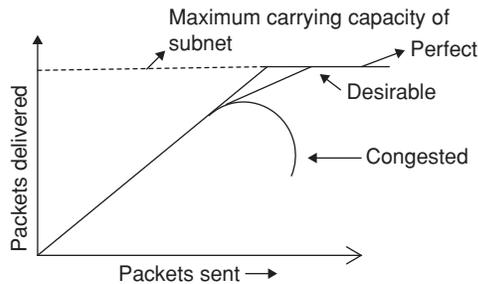


Figure 4 (a) A leaky bucket with water, (b) A leaky bucket with network

A leaky bucket is a bucket with a small hole. No matter at what rate water enters the bucket, the outflow is at constant rate,  $S$ , when there is any water in the bucket and zero when bucket is empty. Once the bucket is full, any additional water entering it, spills over the sides and it is lost. Each host is connected to the network by an interface containing a leaky bucket (i.e., a finite internal queue) congestion control algorithms.



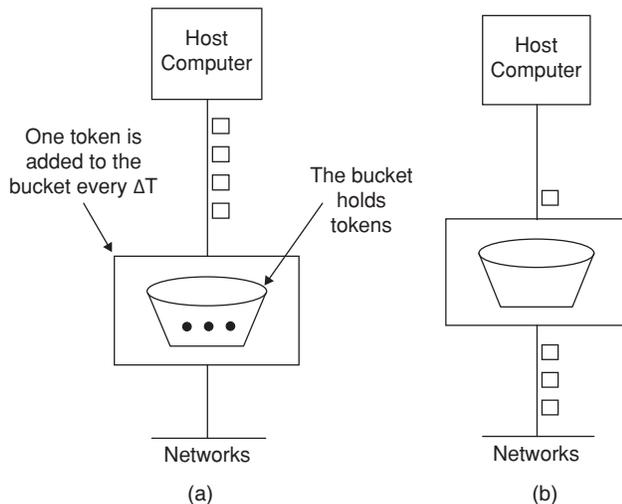
When too many packets are present in the subnet, performance degrades. This situation is called congestion.

### Causes of congestion

1. If all of a sudden, stream of packets are arriving on three or four input lines and all need same output line, a queue will build up.
2. Slow processor.
3. Low bandwidth line.

### Token bucket

Tokens are added at a constant rate. For a packet to be transmitted, it must capture and destroy one token.



**Figure 5** (a) shows that the bucket holds 3 tokens with 4 packets waiting to be transmitted, (b) shows that 3 packets have gotten through but the other one is stuck waiting for tokens to be generated.

Unlike leaky bucket, token bucket allows saving up to maximum size of bucket ‘ $n$ ’.

The bursts of upto ‘ $n$ ’ packets can be sent at once, giving faster response to sudden bursts.

- Leaky bucket discards packets when the bucket is full, whereas token bucket throws away tokens when the bucket is full but never discards packets.
- Let Token bucket capacity be  $c$ (bits), token arrival rate  $\rho$ (bps), maximum output rate  $M$ (bps), and burst length  $S(s)$ .
- During the burst length of  $S(s)$ , tokens generated are  $\rho S$ (bits), output burst contains a maximum of  $C + \rho S$ (bits)
- Output in a maximum burst of length  $S(s)$  is  $MS$ .

$$C + \rho S = MS \quad (\text{or}) \quad S = \frac{C}{M - \rho}$$

- Token bucket still allows large bursts, even though the maximum burst length ‘ $s$ ’ can be regulated by selection of  $\rho$  and  $M$ .
- To reduce the peak rate, put a leaky bucket of a larger rate after the token bucket (To avoid discarding packets)

### Traffic Shaping

1. One of the main causes of congestion is, that traffic is often burst.
2. If hosts could be made to transmit at uniform rate, congestion would be less.

This arrangement can be built into the network interface or simulated by the host OS. The host is allowed to put one packet per tick on the network.

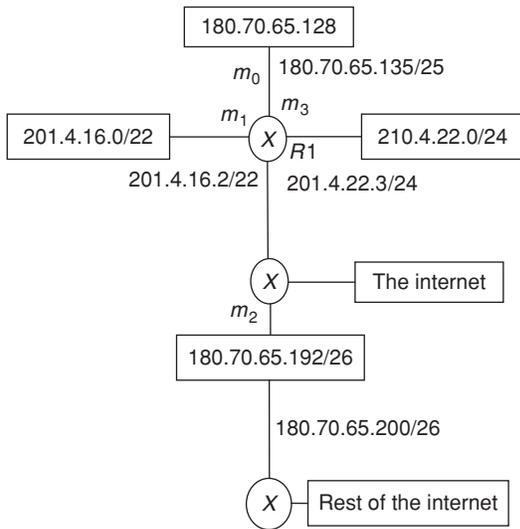
1. When the packets are all of the same size at every clock tick, one packet is transmitted.
2. When variable size packets are used.
  - (i) At every tick, a counter is initialized to  $n$ . If the first packet on the queue has fewer bytes than the current value of the counter, it is transmitted and counter is decremented by that number of bytes.
  - (ii) Additional packets may also be sent, as long as the counter is high enough.
  - (iii) When the counter drops below the length of the next packet on the queue, transmission stops until the next tick, at that time the residual byte count is overwritten and lost.

**EXERCISES**

**Practice Problems I**

**Directions for questions 1 to 15:** Select the correct alternative from the given choices.

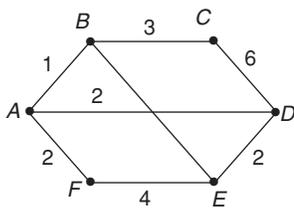
1. Consider below figure:



The network address, 180.70.65.130 goes through which of the following interface?

- (A)  $m_0$
- (B)  $m_1$
- (C)  $m_2$
- (D)  $m_3$

2. Consider below graphical representation of a subnet with each node denoting a router. If all the routers are booted at the same time, what is the number of link state packets that are generated having the cost/delay information?

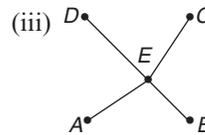
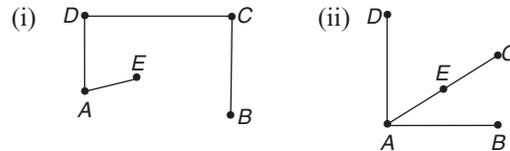
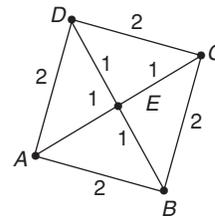


- (A) 3
- (B) 4
- (C) 5
- (D) 6

3. In a TCP connection it is found that burst size of 1024, 2048, 4096 have been transmitted while that of 8192 has resulted in a time out. The receiver has earlier set a window size of 4096. As per slow start algorithm which of the below statement is true?

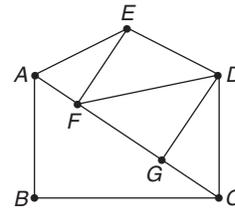
- (i) Congestion window is set to 4096.
  - (ii) Maximum allowed burst size is 8192
- (A) (i) only
  - (B) (ii) only
  - (C) Both (i) and (ii)
  - (D) Neither (i) nor (ii)

4. From the below graph select the sink tree(s):



- (A) (i), (ii)
- (B) (ii), (iii)
- (C) (i), (iii)
- (D) (i), (ii), (iii)

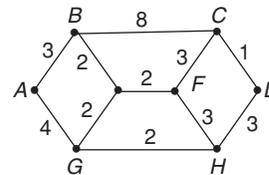
5. Consider the below graph.



It is known that  $D$  is the optimal route from  $A$  to  $C$  and the optimal route from  $A$  to  $C$  has 3 hops. Which of the below statements is certainly true?

- (i)  $B$  is not in the optimal route from  $A$  to  $C$
  - (ii)  $G$  is not in the optimal route from  $B$  to  $C$
  - (iii) Either  $E$  or  $F$  is in the optimal route from  $A$  to  $C$
  - (iv)  $ED, FD$  are both optimal routes
- (A) (i), (ii), (iii)
  - (B) (ii), (iii), (iv)
  - (C) (i), (iii), (iv)
  - (D) (i), (iv), (ii)

6. The shortest path using Dijkstra's algorithm after 3 iterations is



- (A) A G
- (B) A B E
- (C) A B C
- (D) A G H

7. There are totally 20 links among the routers of a subnet. How many rows are needed in all when link state packets

combined together, which are used to notify each other about cost/delay in transmitting data to immediate neighbours. Assume 1 row is needed for each neighbour?

- (A) 10 (B) 20  
(C) 40 (D) 80

8. Below are the link state packets generated by routers in a subnet. What is the shortest distance between *A* and *D*?

| A   |   |
|-----|---|
| Seq |   |
| Age |   |
| B   | 4 |
| E   | 5 |

| B   |   |
|-----|---|
| Seq |   |
| Age |   |
| A   | 4 |
| C   | 2 |
| F   | 6 |

| C   |   |
|-----|---|
| Seq |   |
| Age |   |
| B   | 2 |
| D   | 3 |
| E   | 1 |

| D   |   |
|-----|---|
| Seq |   |
| Age |   |
| C   | 3 |
| F   | 7 |

| E   |   |
|-----|---|
| Seq |   |
| Age |   |
| A   | 5 |
| C   | 1 |
| F   | 8 |

| F   |   |
|-----|---|
| Seq |   |
| Age |   |
| B   | 6 |
| D   | 7 |
| E   | 8 |

- (A) 6 (B) 9  
(C) 10 (D) 11

9. What are the advantages of reverse path forwarding over other broadcasting algorithms like spanning trees, multidestination routing, broadcasting, and flooding?

- (i) Route does not need to know information regarding spanning tree structures  
(ii) Uses destination tables for further forwarding  
(iii) Does not need a halt mechanism to stop packets from further getting routed

- (A) (i), (ii) (B) (ii), (iii)  
(C) (i), (iii) (D) (i), (ii), (iii)

10. Which of the following specifies the correct sequence of steps to route packets to mobile hosts?

- (i) Sender is given foreign agent's address  
(ii) Packet is sent to mobile host's home address  
(iii) Packet is tunneled to foreign agent  
(iv) Subsequent packets are tunneled to the foreign agent

- (A) (i), (ii), (iii), (iv)  
(B) (ii), (iii), (iv), (i)

- (C) (ii), (iii), (i), (iv)  
(D) (iii), (iv), (i), (ii)

11. What are the different parts of congestion control by closed loop methods?

- (i) Design the system in advance to make sure congestion doesn't occur in first place  
(ii) Monitor the system to detect when and where congestion occurs  
(iii) Pass congestion information to places where action can be taken  
(iv) Adjust system operation to correct the problem

- (A) (i), (ii), (iii)  
(B) (ii), (iii), (iv)  
(C) (iii), (iv), (i)  
(D) (i), (ii), (iv)

12. In Selective flooding

- (A) Packets are sent in all outgoing lines.  
(B) Packets are sent in only on those lines that are approximately in the right direction.  
(C) Both (A) and (B)  
(D) None of these

13. There are 5 routers and 6 networks in an inter-networking, using link state routing, how many routing tables are there?

- (A) 1 (B) 5  
(C) 6 (D) 11

14. Congestion control for multicasting flows from multiple sources to multiple destinations, the solution that can handle this is

- (A) RSVP (Resource reSerVation Protocol)  
(B) Load shedding  
(C) Both (A) and (B)  
(D) None of these.

15. Which of the below are part of backward learning algorithm?

- (i) As the bridge starts operating, a hash table to map source addresses to corresponding LANs is constructed.  
(ii) It dynamically updates the hash tables when machines are connected and re connected to the LAN.

- (iii) It encrypts the frames for security reasons.

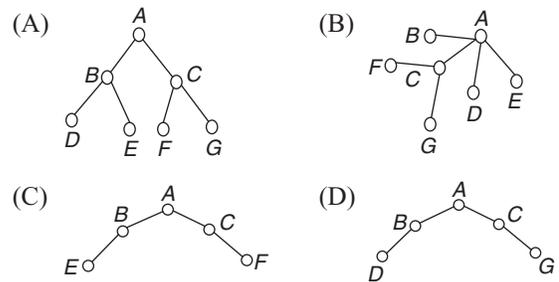
- (A) (i), (ii)  
(B) (ii), (iii)  
(C) (i), (iii)  
(D) (i), (ii), (iii)

## Practice Problems 2

**Directions for questions 1 to 15:** Select the correct alternative from the given choices.

- What does a routing algorithm perform?
  - Decides if incoming packet should be further corrected for transmission errors
  - Adds checksum bits to packets
  - Encrypts the packets
  - Decides the output line on which the incoming packet should be transmitted
- What happens in session routing?
  - User's session variables are managed by the network layer
  - Route remains same throughout the user session
  - Packets change their route for optimization sake during user session
  - Provides special routes for important packets
- What is the type of algorithm that changes their routing decision based on changes in topology and traffic?
  - Adaptive routing
  - Static routing
  - Non-adaptive routing
  - Network routing
- Which of the below routing method always ensures the shortest path even though routers crash during course of routing?
  - Dijkstra Routing
  - Flooding
  - Distance Vector Routing
  - Link State Routing
- What is the root cause for count-to-infinity problem?
  - The routing tables are static and are not updated.
  - The routing tables run out of space to accommodate more entries in table.
  - When router  $X$  tells router  $Y$  that there is a path, it doesn't say if  $Y$  itself is in the path.
  - When router  $X$  tells router  $Y$  that there is a path (to target route  $Z$ ) it doesn't inform  $Z$  about the path.
- In a strict sure security path ABCD, where A, B, C, D are routers, the maximum bandwidth is found to be 500 kbps, 700 kbps, 900 kbps, 300 kbps respectively. What is the effective bandwidth if no buffering is possible?
  - 600 kbps
  - 900 kbps
  - 300 kbps
  - 2400 kbps
- What is the characteristic of Distance Vector Routing?
  - Time taken to reach other routers in the network is maintained in the routing tables.
  - Algorithm is susceptible to count-to-infinity problem.
  - The preferred outgoing line to be used for a particular destination is also stored in tables.
  - (i), (ii)
  - (ii), (iii)
  - (i), (iii)
  - (i), (ii), (iii)

- A subnet using link state algorithm has router, using link state packets with sequence of 16-bit fixed size. If a link state packet is sent every second, how long would it take before wrap around occurs. Assume starting sequence number is 0.
  - 24.5 hours
  - 18.20 hours
  - 17.5 hours
  - 16.4 hours
- Which of the following are features of link state routing?
  - In the first step discover all the routers in the subnet and find their network addresses.
  - Measure cost/delay to the neighbours.
  - Transmit the information as obtained in (ii) across the subnet.
  - Thus by pass the necessity for shortest path algorithm.
  - (i), (ii)
  - (ii), (iii)
  - (iii), (iv)
  - (i), (iv)
- In multidestination routing,
  - Each router makes new copies of the incoming packets.
  - It retains the same destination list in all copies.
  - It places them on appropriate outgoing lines.
  - (i), (ii)
  - (ii), (iii)
  - (iii), (i)
  - (i), (ii), (iii)
- In a subnet which follows reverse path forwarding, routers  $B$  and  $C$  have received packets from  $A$  which have been further forwarded to  $D$  and  $E$  by  $B$  and to  $F$  and  $G$  by  $C$ . Of this  $D$ ,  $G$  has always discarded the valid packets. Construct the preferred routing lines in the subnet.



- Which of the following layers accept services from network layer and provides services to session layer?
  - Data link layer
  - Presentation layer
  - Transport layer
  - Physical layer.
- Which of the below are different metrics for congestion?
  - Packets discarded for lack of buffer space
  - Packets that are retransmitted
  - Average packet delay
  - Average queue length
  - (i), (ii), (iii)
  - (ii), (iii), (iv)
  - (iii), (iv), (i)
  - (i), (ii), (iii), (iv)

14. What are the ways to decrease congestion?

- (i) Put spare routers to use
  - (ii) Increase bandwidth by routing on alternate lines
  - (iii) Increase the size of tables in the routers
  - (iv) Decrease the load
- (A) (i), (ii), (iii)                      (B) (ii), (iii), (iv)  
 (C) (iii), (iv), (i)                      (D) (iv), (i), (ii)

15. The algorithm which tells the routers to maintain certain data structures in their memories for congestion control is

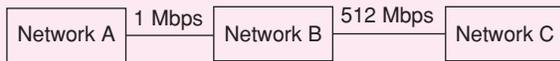
- (A) Resource Reservation Protocol.
- (B) Fair queuing algorithm.
- (C) Token bucket algorithm.
- (D) None of these

**PREVIOUS YEARS' QUESTIONS**

**Common data for questions 1 and 2:** Consider three IP networks, *A*, *B* and *C*. Host  $H_A$  in network *A* sends messages each containing 180 bytes of application data to a host  $H_C$  in network *C*. The TCP layer prefixes a 20 byte header to the message. This passes through an intermediate network *B*. The maximum packet size, including 20 byte IP header, in each network is:

- A: 1000 bytes
- B: 100 bytes
- C: 1000 bytes

The networks *A* and *B* are connected through a 1 Mbps link, while *B* and *C* are connected by a 512 Kbps link (bps = bits per second)



1. Assuming that the packets are correctly delivered, How many bytes, including headers, are delivered to the IP layer at the destination for one application message, in the best case? Consider only data packets. [2004]

- (A) 200                                      (B) 220
- (C) 240                                      (D) 260

2. What is the rate at which the application data is transferred to host  $H_C$ ? Ignore errors, acknowledgements, and other overheads. [2004]

- (A) 325.5 kbps                              (B) 354.5 kbps
- (C) 409.6 kbps                              (D) 512.0 kbps

3. In a packet switching network, packets routed from source to destination along a single path having two intermediate nodes. If the message size is 24 bytes and each packet contains a header of 3 bytes, then the optimum packet size is: [2005]

- (A) 4                                              (B) 6
- (C) 7                                              (D) 9

4. Suppose the round trip propagation delay for a 10 Mbps Ethernet having 48-bit jamming signal is 46.4  $\mu$ s. The minimum frame size is: [2005]

- (A) 94                                              (B) 416
- (C) 464                                              (D) 512

5. Station *A* uses 32 byte packets to transmit messages to station *B* using a sliding window protocol. The round trip delay between *A* and *B* is 80 milliseconds and the bottleneck bandwidth on the path between *A* and *B*

is 128 kbps. What is the optimal window size that *A* should use? [2006]

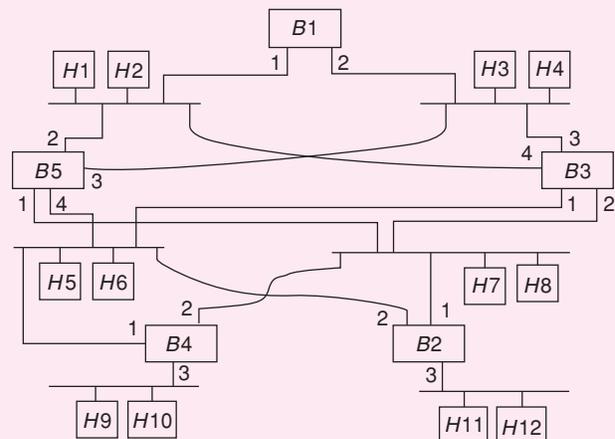
- (A) 20                                              (B) 40
- (C) 160                                              (D) 320

6. Station *A* needs to send a message consisting of 9 packets to station *B* using a sliding window (window size 3) and go-back-n error control strategy. All packets are ready and immediately available for transmission. If every 5th packet that *A* transmits gets lost (but no acks from *B* ever get lost), then what is the number of packets that *A* will transmit for sending the message to *B*? [2006]

- (A) 12                                              (B) 14
- (C) 16                                              (D) 18

**Common data for questions 7 and 8:** Consider the diagram shown, where a number of LANs are connected by (transparent) bridges. In order to avoid packets looping through circuits in the graph, the bridges organize themselves in a spanning tree. First, the root bridge is identified as the bridge with the least serial number. Next, the root sends out (one or more) data units to enable the setting up of shortest paths from the root bridge to each bridge.

Each bridge identifies a port (the root port) through which it will forward frames to the root bridge. Port conflicts are always resolved in favor of the port with the lower index value. When there is possibility of multiple bridges forwarding to the same LAN (But not through the root port), ties are broken as follows: bridges closest to the root get preference and between such bridges, the one with the lowest serial number is preferred.



7. For the given connection of LANs by bridges, which one of the following choices represents the depth first traversal of the spanning tree of bridges? [2006]

(A) B1, B5, B3, B4, B2 (B) B1, B3, B5, B2, B4  
(C) B1, B5, B2, B3, B4 (D) B1, B3, B4, B5, B2

8. Consider the spanning tree for the previous question. let Host H1 send out a broadcast ping packet. Which of the following options represents the correct forwarding table on B3? [2006]

(A)

| Hosts            | Port |
|------------------|------|
| H1, H2, H3, H4   | 3    |
| H5, H6, H9, H10  | 1    |
| H7, H8, H11, H12 | 2    |

(B)

| Hosts                 | Port |
|-----------------------|------|
| H1, H2                | 4    |
| H3, H4                | 3    |
| H5, H6                | 1    |
| H7, H8, H10, H11, H12 | 2    |

(C)

| Hosts            | Port |
|------------------|------|
| H3, H4           | 3    |
| H5, H6, H9, H10  | 1    |
| H1, H2           | 4    |
| H7, H8, H11, H12 | 2    |

(D)

| Hosts           | Port |
|-----------------|------|
| H2, H2, H3, H4  | 3    |
| H5, H7, H9, H10 | 1    |
| H7, H8, 11, H12 | 4    |

9. In a token ring network the transmission speed is  $10^7$  bps and the propagation speed is 200 metres/ $\mu$ s. The 1-bit delay in this network is equivalent to: [2007]

(A) 500 metres of cable.  
(B) 200 metres of cable.  
(C) 20 metres of cable.  
(D) 50 metres of cable.

10. In the slow start phase of the TCP congestion control algorithm, the size of the congestion window [2008]

(A) Does not increase  
(B) Increases linearly  
(C) Increases quadratically  
(D) Increases exponentially

11. A computer on a 10 Mbps network is regulated by a token bucket. The token bucket is filled at a rate of 2 Mbps. It is initially filled to capacity with 16 Megabits. What is the maximum duration for which the computer can transmit at the full 10 Mbps? [2008]

(A) 1.6 seconds (B) 2 seconds  
(C) 5 seconds (D) 8 seconds

12. Let  $G(x)$  be the generator polynomial used for CRC checking. What is the condition that should be satisfied by  $G(x)$  to detect odd number of bits in error? [2009]

(A)  $G(x)$  contains more than two terms  
(B)  $G(x)$  does not divide  $1 + x^k$ , for any  $k$  not exceeding the frame length  
(C)  $1 + x$  is a factor of  $G(x)$   
(D)  $G(x)$  has an odd number of terms.

**Common data for questions 13 and 14:** Frames of 1000 bits are sent over a  $10^6$  bps duplex link between two hosts. The propagation time is 25 ms. Frames are to be transmitted into this link to maximally pack them in transit (within the link).

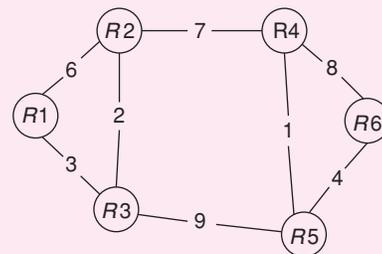
13. What is the minimum number of bits ( $l$ ) that will be required to represent the sequence numbers distinctly? Assume that no time gap needs to be given between transmission of two frames. [2009]

(A)  $l = 2$  (B)  $l = 3$   
(C)  $l = 4$  (D)  $l = 5$

14. Suppose that the sliding window protocol is used with the sender window size of  $2^l$ , where  $l$  is the number of bits identified in the earlier part and acknowledgements are always piggy backed. After sending  $2^l$  frames, what is the minimum time the sender will have to wait before starting transmission of the next frame? (Identify the closest choice ignoring the frame processing time.) [2009]

(A) 16 ms (B) 18 ms  
(C) 20 ms (D) 22 ms

**Common data for questions 15 and 16:** Consider a network with 6 routers  $R1$  to  $R6$  connected with links having weights as shown in the following diagram



15. All the routers use the distance vector based routing algorithm to update their routing tables. Each router starts with its routing table initialized to contain an entry for each neighbour with the weight of the respective connecting link. After all the routing tables stabilize, how many links in the network will never be used for carrying any data? [2010]

(A) 4 (B) 3 (C) 2 (D) 1

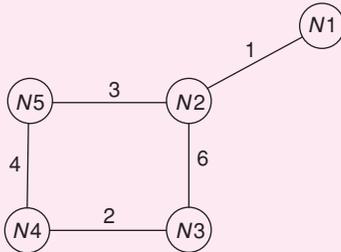
16. Suppose the weights of all unused links in the previous question are changed to 2 and the distance

vector algorithm is used again until all routing tables stabilize. How many links will now remain unused?

[2010]

- (A) 0    (B) 1    (C) 2    (D) 3

**Common data for questions 17 and 18:** Consider a network with five nodes,  $N1$  to  $N5$  as shown below.



The network uses a distance vector routing protocol. Once the routes have stabilized, the distance vectors at different nodes are as following.

- $N1$  : (0, 1, 7, 8, 4)                       $N4$  : (8, 7, 2, 0, 4)  
 $N2$  : (1, 0, 6, 7, 3)                       $N5$  : (4, 3, 6, 4, 0)  
 $N3$  : (7, 6, 0, 2, 6)

Each distance vector is the distance of the best known path at that instance to nodes,  $N1$  to  $N5$ , where the distance to itself is 0. Also, all links are symmetric and the cost is identical in both directions. In each round, all nodes exchange their distance vectors with their respective neighbors. Then all nodes update their distance vectors. In between two rounds, any change in cost of a link will cause the two incident nodes to change only that entry in their distance vectors.

17. The cost of link  $N2$ - $N3$  reduces to 2 (in both directions). After the next round of updates, what will be the new distance vector at node,  $N3$ ? [2011]  
 (A) (3, 2, 0, 2, 5)                      (B) (3, 2, 0, 2, 6)  
 (C) (7, 2, 0, 2, 5)                      (D) (7, 2, 0, 2, 6)
18. After the update in the previous question, the link  $N1$ - $N2$  goes down.  $N2$  will reflect this change immediately in its distance vector as cost,  $\infty$ . After the NEXT ROUND of update, what will be the cost to  $N1$  in the distance vector of  $N3$ ? [2011]  
 (A) 3    (B) 9    (C) 10    (D)  $\infty$
19. Consider an instance of TCP's Additive Increase Multiplicative Decrease (AIMD) algorithm where the window size at the start of the slow start phase is 2 MSS and the threshold at the start of the first transmission is 8 MSS. Assume that a timeout occurs during the fifth transmission. Find the congestion window size at the end of the tenth transmission. [2012]  
 (A) 8 MSS                      (B) 14 MSS  
 (C) 7 MSS                      (D) 12 MSS

20. Assume that source  $S$  and destination  $D$  are connected through two intermediate routers labeled  $R$ . Determine how many times each packet has to visit the network layer and the data link layer during a transmission from  $S$  to  $D$ . [2013]



- (A) Network layer – 4 times and Data link layer – 4 times  
 (B) Network layer – 4 times and Data link layer – 3 times  
 (C) Network layer – 4 times and Data link layer – 6 times  
 (D) Network layer – 2 times and Data link layer – 6 times
21. Consider a selective repeat sliding window protocol that uses a frame size of 1kB to send data on a 1.5 Mbps link with a one-way latency of 50 msec. To achieve a link utilization of 60%, the minimum number of bits required to represent the sequence number field is \_\_\_\_\_. [2014]
22. Consider the following three statements about link state and distance vector routing protocols, for a large network with 500 network nodes and 4000 links.  
 [S1] The computational overhead in link state protocols is higher than in distance vector protocols.  
 [S2] A distance vector protocol (with split horizon) avoids persistent routing loops, but not a link state protocol.  
 [S3] After a topology change, a link state protocol will converge faster than a distance vector protocol.  
 Which one of the following is correct about  $S1$ ,  $S2$  and  $S3$ ? [2014]  
 (A)  $S1$ ,  $S2$  and  $S3$  are all true  
 (B)  $S1$ ,  $S2$  and  $S3$  are all false  
 (C)  $S1$  and  $S2$  are true, but  $S3$  is false  
 (D)  $S1$  and  $S3$  are true, but  $S2$  is false.
23. Let the size of congestion window of a TCP connection be 32 kB when a timeout occurs. The round trip time of the connection is 100 msec and the maximum segment size used is 2 kB. The time taken (in msec) by the TCP connection to get back to 32 kB congestion window is \_\_\_\_\_. [2014]
24. Which one of the following is TRUE about the interior gateway routing protocols-Routing information protocol (RIP) and Open Shortest Path First (OSPF)? [2014]  
 (A) RIP uses distance vector routing and OSPF uses link state routing  
 (B) OSPF uses distance vector routing and RIP uses link state routing  
 (C) Both RIP and OSPF use link state routing  
 (D) Both RIP and OSPF use distance vector routing

25. Consider the store and forward packet switched network given below. Assume that the bandwidth of each link is  $10^6$  bytes/sec. A user on host  $A$  sends a file of size  $10^3$  bytes to host  $B$  through routers  $R_1$  and  $R_2$  in three different ways. In the first case a single packet containing the complete file is transmitted from  $A$  to  $B$ . In the second case, the file is spilt into 10 equal parts, and these packets are transmitted from  $A$  to  $B$ . In the third case, the file is spilt into 20 equal parts, and these packets are sent from  $A$  to  $B$ . Each packet contains 100 bytes of header information along with the user data. Consider only transmission time and ignore processing, queuing and propagation delays. Also assume that there are no errors during transmissions. Let  $T_1$ ,  $T_2$  and  $T_3$  be the times taken to transmit the file in the first, second and third case respectively. Which one of the following is CORRECT? [2014]



- (A)  $T_1 < T_2 < T_3$                       (B)  $T_1 > T_2 > T_3$   
 (C)  $T_2 = T_3, T_3 < T_1$                   (D)  $T_1 = T_3, T_3 > T_2$
26. An IP machine  $Q$  has a path to another IP machine  $H$  via three IP routers  $R_1, R_2$ , and  $R_3$ .  $Q - R_1 - R_2 - R_3 - H$ .  $H$  acts as an HTTP server, and  $Q$  connects to  $H$  via HTTP and downloads a file. Session layer encryption is used with DES as the shared key encryption protocol. Consider the following four pieces of information.
- [I1] The URL of the file downloaded by  $Q$   
 [I2] The TCP port numbers at  $Q$  and  $H$   
 [I3] The IP addresses of  $Q$  and  $H$   
 [I4] The link layer addresses of  $Q$  and  $H$
- Which of  $I1, I2, I3$  and  $I4$  can an intruder learn through sniffing at  $R_2$  alone? [2014]
- (A) Only  $I1$  and  $I2$                       (B) Only  $I1$   
 (C) Only  $I2$  and  $I3$                       (D) Only  $I3$  and  $I4$

27. An IP router with a Maximum Transmission Unit (MTU) of 1500 bytes has received an IP packet of size 4404 bytes with an IP header of length 20 bytes. The values of the relevant fields in the header of the third IP fragment generated by the router for this packet are [2014]
- (A) MF bit : 0, Datagram Length: 1444; Offset: 370  
 (B) MF bit: 1, Datagram Length : 1424; Offset: 185  
 (C) MF Bit: 1, Datagram Length: 1500; Offset: 370  
 (D) MF bit: 0, Datagram Length: 1424; Offset: 2960
28. Identify the correct order in which a server process must invoke the function calls accept, bind, listen, and recv according to UNIX socket API. [2015]
- (A) listen, accept, bind, recv  
 (B) bind, listen, accept, recv  
 (C) bind, accept, listen, recv  
 (D) accept, listen, bind, recv
29. For a host machine that uses the token bucket algorithm for congestion control, the token bucket has a capacity of 1 megabyte and the maximum output rate is 20 megabytes per second. Token arrive at a rate to sustain output at a rate of 10 megabytes per second. The token bucket is currently full and the machine needs to send 12 megabytes of data. The minimum time required to transmit the data is \_\_\_\_\_ seconds. [2016]
30. Consider the following statements about the routing protocols. Routing Information Protocol (RIP) and Open Shortest Path First (OSPF) in an IPv4 network.
- I: RIP uses distance vector routing  
 II: RIP packets are sent using UDP  
 III: OSPF packets are sent using TCP  
 IV: OSPF operation is based on link-state routing
- Which of the statements above are CORRECT? [2017]
- (A) I and IV only                      (B) I, II and III only  
 (C) I, II and IV only                      (D) II, III and IV only

## ANSWER KEYS

### EXERCISES

#### Practice Problems 1

1. A      2. D      3. A      4. B      5. A      6. A      7. C      8. B      9. C      10. C  
 11. B      12. B      13. B      14. A      15. A

#### Practice Problems 2

1. D      2. B      3. A      4. B      5. C      6. C      7. D      8. B      9. B      10. C  
 11. C      12. C      13. D      14. D      15. A

#### Previous Years' Questions

1. D      2. B      3. D      4. D      5. B      6. C      7. C      8. A      9. C      10. D  
 11. B      12. C      13. D      14. B      15. C      16. B      17. A      18. C      19.      20. C  
 21. 5      22. D      23. 1100 to 1300      24. A      25. D      26. C      27. A      28. B      29. 1.1  
 30. C

## TCP/UDP

### LEARNING OBJECTIVES

- Transport layer
- User Datagram Protocol (UDP)
- TCP/IP
- TCP/IP vs OSI reference model
- TCP state transition diagram
- TCP flow control
- Application layer
- ICMP, SMTP, POP3, IMAP 4, HTTP, FTP
- DNS
- Network devices

### TRANSPORT LAYER

Real communication takes place between two applications programs i.e., processes. For this, process-to-process delivery is needed. A mechanism is required in order to deliver data from one of these processes running on the source host to the corresponding process running on the destination host.

The transport layer is responsible for process-to-process delivery.

### Addressing in Transport Layer

#### Port addresses

- A transport layer address is a port number.
- The destination port number is needed for delivery and the source port number is needed for reply.
- The port numbers are 16-bit integers ranging from 0 to 65535.

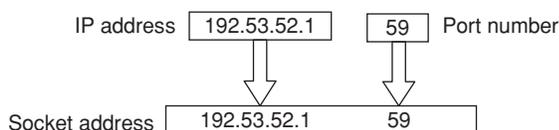
The IANA (Internet Assigned Number Authority) has divided the port numbers as:

- Well-known ports (0 to 1023)
- Registered ports (1024 to 49,151)
- Dynamic or private or ephemeral ports (49,152 to 65,535)

#### Socket address

Process to process delivery needs two identifiers, IP address and port address at each end to make a connection.

The combination of an IP address and a port number is socket address.



### Protocols at transport layer

- UDP
- TCP
- SCTP

### USER DATAGRAM PROTOCOL (UDP)

- UDP is connectionless protocol.
  - There is no mechanism for connection establishment or connection termination.
  - The packets may be delayed or lost or may arrive out of sequence, i.e., there is no acknowledgement.
  - Each user datagram sent by UDP is an independent program. Even if the user datagram's are coming from the same source program and going to the same destination process, there is no relationship between the different datagrams.

Thus, user datagrams can travel on a different path.

- Multicasting capability is embedded in UDP.
- It is a simple, unreliable transport protocol.
  - There is no flow control, no window mechanism.
  - There is no error control as well except for the checksum. The sender does not know if a message has been lost or duplicated. When the receiver detects an error through the checksum, the datagram is discarded silently.
- It is used in real-time applications.
  - The header length is fixed, of 8 bytes. Real time applications require a constant flow of data. Moreover, the unreliability (fast and less complex service) of UDP aids in real-time applications like voice over IP, online games etc.
- It encapsulates and decapsulates messages in an IP datagram.

### User Datagram

UDP packets have other name called user datagrams. They have a fixed size header of 8 bytes. The datagram is divided into 4 fields.

|                                 |                                      |
|---------------------------------|--------------------------------------|
| Source port number<br>(16-bits) | Destination port<br>number (16-bits) |
| Total length<br>(16-bits)       | Checksum<br>(16-bits)                |

Figure 1 User datagram header format

1. **Source Port Number** It is a 16-bit number used by the process running on the source host.
2. **Destination Port Number** It is also a 16 bit number used by the process running on the destination host.
3. **Total length** It is a 16-bit field, it defines the total length of the user datagram header and data. It can define a total length of 0 to 65535 bytes. A UDP packet is encapsulated in an IP packet.

|                                                                    |
|--------------------------------------------------------------------|
| $\text{UDP length} = \text{IP length} - \text{IP header's length}$ |
|--------------------------------------------------------------------|

4. **Checksum:** It is optional field, if not available the field is filled with 1's. It is used to detect errors in user datagram (header plus data).

### Protocols That Take UDP Services

Following are a few protocols that take the services of UDP:

1. Domain Name Service (port – 53): UDP is used to send small data. If the data is less than 512 bytes, then DNS uses UDP else it goes for TCP.
2. Trivial File Transfer Protocol (port – 69): TFTP is used to transfer simple and small files, it uses UDP service.
3. Routing Information protocol: It uses UDP service on port number 520 to update routers.
4. Simple Network Management Protocol (SNMP): The SNMP agent receives requests on UDP port 161 for management process.
5. Bootstrap protocol (BOOTP): For client (port 68) and for server (port – 67).

### UDP Checksum Calculation

- The checksum includes a pseudo header, the UDP header and the data coming from the application layer.

|                               |                     |                         |
|-------------------------------|---------------------|-------------------------|
| 32-bit source IP address      |                     |                         |
| 32-bit destination IP address |                     |                         |
| All 0's                       | 8-bit protocol (17) | 16-bit UDP total length |

Figure 2 Pseudo header of UDP for checksum calculation

- The value of protocol field is 17. If this value is changed during transmission, the checksum calculation at the receiver will detect it and UDP drops the packet.
- If the checksum is not calculated, the field is filled with 0's. This means checksum calculation is optional.
- The calculated checksum can never be all 1's as this implies that the sum is all 0's. But this is impossible because for this the value of fields have to be 0's.

### TCP/IP

TCP/IP is a network model which is used for the internet architecture, its main objectives are

- Connecting the multiple networks.
- Maintaining the intact connection between two machines, which are functioning.

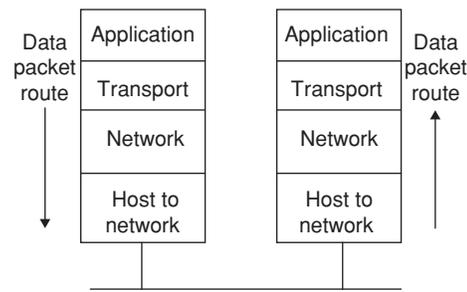


Figure 3 TCP/IP network protocol

### TCP/IP vs OSI Reference Model

| OSI                                                                       | TCP/IP                              |
|---------------------------------------------------------------------------|-------------------------------------|
| (1) There are 7 layers                                                    | (1) There are 5 layers              |
| (2) There is no definition for multicasting                               | (2) Multicasting is clearly defined |
| (3) Less flexibility                                                      | (3) Lot of flexibility              |
| (4) Practically it is not suggestible as it is based on theoretical rules | (4) It is based on practical rules  |

- TCP stands for Transmission Control Protocol.
- It is connection-oriented protocol.
  - It creates a virtual connection between two TCPs to send data then data is transferred and at the end the connection is released.
  - There is acknowledgement mechanism for safe and sound arrival of data.
- It is a reliable transport protocol.
  - Uses flow and error control.
  - Slower and more complex service.
  - Duplicate segments are detected, lost segments are resent, the bytes are delivered to the end process in order.
- It is a stream-oriented protocol.

- Allows the sending process to deliver data as a stream of bytes and allows the receiving process to obtain data as a stream of bytes.
- TCP offers full-duplex service.
  - Data can flow in both directions at the same time.
  - Each TCP has a sending and receiving buffer.
- It cannot be used in real time applications as the header length varies from 20-to-60 bytes, moreover it needs reliability.

### TCP Header Format

- A packet in TCP is called a segment. The segment consists of a 20-to-60 bytes header.
- If there are no options, the header is of 20 bytes.

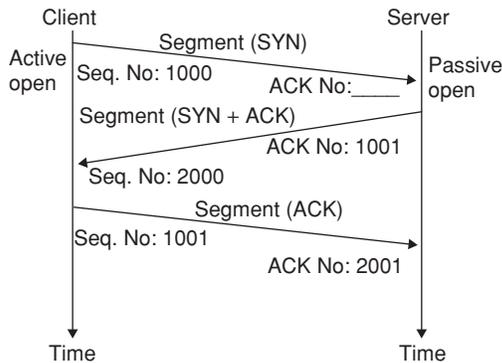
|                                  |                      |     |     |     |                                    |     |     |                          |
|----------------------------------|----------------------|-----|-----|-----|------------------------------------|-----|-----|--------------------------|
| Source port address (16-bits)    |                      |     |     |     | Destination port address (16-bits) |     |     |                          |
| Sequence number (32-bits)        |                      |     |     |     |                                    |     |     |                          |
| Acknowledgement number (32-bits) |                      |     |     |     |                                    |     |     |                          |
| HLEN<br>(4-bits)                 | Reserved<br>(6-bits) | URG | ACK | PSH | RST                                | SYN | FIN | Window size<br>(16-bits) |
| Checksum (16-bits)               |                      |     |     |     | Urgent Pointer (16-bits)           |     |     |                          |
| Options and Padding              |                      |     |     |     |                                    |     |     |                          |

Figure 4 TCP header format

- If there are options, the header goes upto 60 bytes.
  - **Source Port addresses** A 16-bit field that defines the port number of the application program in the host that is sending the segment.
  - **Destination Port address** A 16-bit field that defines the port number of the application program in the host is receiving the segment.
  - **Sequence number** A 32-bit field whose value defines the number of the first data byte contained in that segment. During connection establishment, a random number is generated to create an initial sequence number (ISN) which is usually different in each direction.
  - **Acknowledgement Number** A 32-bit field whose value defines the number of the next byte, a party expects to receive. If the receiver of the segment has successfully received byte number  $x$  from the other party, it defines  $x + 1$  as the acknowledgement number. The acknowledgement number is cumulative.
  - **HLEN(Header Length)** This field is of 4-bit. The header length can be between 20 and 60 bytes. The value of this field can be between  $5(5 \times 4 = 20)$  and  $15(15 \times 4 = 60)$ .
  - **Reserved** This is a 6-bit field which is reserved for future use.
  - **Control** This field contains 6 control flags. These are as follows.
    - URG: Urgent pointer. This flag is set when the value of urgent pointer field is valid.
    - ACK: Acknowledgement pointer. This flag is set when the value of acknowledgement field is valid. It is not set at the start of connection during 3-way handshake.
    - RST: Reset pointer. Used to reset the connection, reject an invalid segment or refuse an attempt to open a connection.
    - PSH: Push pointer. When a data is pushed the flag is set.
    - SYN: Synchronization pointer, used to synchronize sequence numbers during connection. If it is set to 1, then it is ISN. If set to 0, then it is the accumulated sequence number of the first data byte of the segment for the current session.
    - FIN: Finish Pointer. It is used to terminate a connection. It indicates that the sender is not interested in sending any more data.
  - **Window size** The field size is of 16-bits and thus the maximum size of the window is 65,535 bytes. This field is determined by the receiver and thus referred to as the receiving window. The window size is variable.
  - **Checksum** The inclusion of this 16-bits field is mandatory in TCP. The calculation of the checksum for TCP follows the same procedure as in UDP, only the value of protocol field in TCP is 6.
  - **Urgent pointer** This 16-bit field, is valid only if the urgent flag is set. This field is used when the segment contains urgent data.
  - **Options and padding** When the header length is greater than 5, option field is used to make the segment into the multiples of 32. Padding is used to ensure the ending of TCP header, it is composed to 32 zeros.
- ### TCP Connection
- TCP is connection-oriented and the connection is virtual not physical.
  - TCP operates at a higher level. TCP uses the services of IP to deliver individual segments to the receiver, but it controls the connection itself. Lost or corrupted segments are retransmitted.
  - In TCP, connection-oriented transmission requires three phases:
    1. Connection establishment
    2. Data transfer
    3. Connection termination

### Connection establishment

- The connection establishment in TCP is called three-way handshaking.
- The process starts with the server. The server program tells its TCP that it is ready to accept a connection. This is a request for a passive open.
- The client program issues a request for an active open. A client that wishes to connect to an open server tells its TCP that it needs to be connected to a particular server. Hence the TCP can start the three-way handshaking process as shown in the figure.



1. The first segment which is a SYN segment is identified by the randomly generated number and is assigned to a 1 byte dummy data indicating the sequence number.
2. Again from the server side a randomly generated number is assigned for the dummy data indicating the first byte.
3. A SYN segment cannot carry data, but it consumes one sequence number.  
A (SYN + ACK) segment cannot carry data, but consumes one sequence number.  
An ACK segment, if carrying no data, consumes no sequence number.
4. Initial Sequence Number (ISN) 1000 is sent from the client to server. Server receives the segment 1000 and is expecting segment 1001 as the next one.

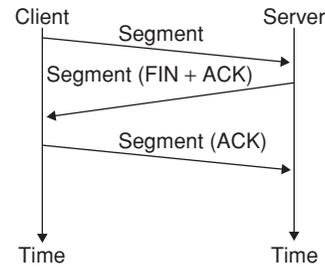
### Data transfer

- After the connection is established, bidirectional data transfer can take place. Both the client and server can send data and acknowledgements.
- The data segments sent by the client have the PSH (push) flag set so that the server TCP knows to deliver data to the server process as soon as they are received.
- Sometimes the sending application program wants a piece of data to be read out of order by the receiving application program that means an application program needs to send urgent bytes then in this case the URG bit is set and the segment is sent. The sending TCP creates a segment and inserts the urgent data at the beginning of the segment.

### Connection termination

There are two options for connection termination.

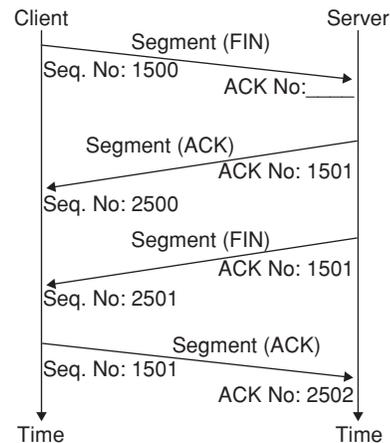
#### Three-way handshaking



- The client process sends the first segment, a FIN segment in which the FIN flag is set. The FIN segment consumes one sequence number if it does not carry data.
- The server TCP sends the second segment, a FIN + ACK segment, to confirm the receipt of the FIN segment from the client and at the same time announce the closing of the connection in the other direction. The FIN + ACK segment consumes one sequence number if it does not carry data.
- The client sends the last ACK segment to the server. This segment contains acknowledgement number which is 1 plus, the sequence number received in the FIN segment from the server.

#### Four-way handshaking

- **Half-close:** In TCP, one end can stop sending data while still receiving data. This is half close.
- The client half-closes the connection by sending a FIN segment.
- The server accepts the half-close by sending the ACK segment. The data transfer from the client to the server stops.
- When the server has sent all the processed data, it sends a FIN segment, which is acknowledged by an ACK from the client.



## TCP State Transition Diagram

The functionality of TCP connection setup, communication phase and termination phase can be easily depicted by the state transition diagram where the TCP will be only at one state at a time with respect to server or client.

A change in the state is only observed after receiving a request for change like ACK (acknowledgement).

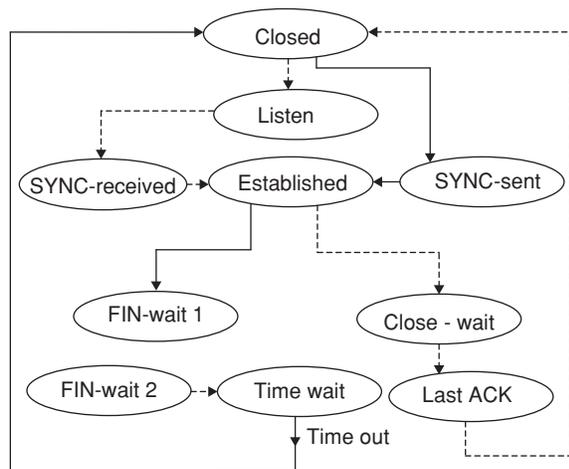


Figure 5 State transition diagram

Here, Solid line ‘—’ is for client states, Break line ‘---’ is for server states

State ‘Closed’ is common for both client and server. Initially the client and the server are in the closed state where no TCP connection is set. When an application request for a TCP connection then the client changes its state from closed to SYNC-sent state.

### Client states

1. **SYNC-sent** After the client sends a SYNC-sent and receives an ACK for the sent SYNC segment, it changes its state to ESTABLISHED STATE.
2. **Established** In this state the client and the server exchange user data. After the requested application is completed, it sends a FIN segment and changes its state to FIN-wait 1.
3. **FIN-wait 1** FIN-wait 1 changes to FIN-wait 2 after receiving an ACK for sent FIN segment.
4. **FIN-wait 2** The client will remain in this state until it receives a FIN segment from the server. When the last ACK is sent by the client, the client changes its state to Time-wait.
5. **Time-wait** A timer is set at this state for any delayed segment from the server which are removed or discarded at the client and after the timeout is reached, the client changes its state from present state to the closed state again.

### Server states

1. **Listen** This is a passive state where the server always listens for the SYNC request segment on different TCP ports.
2. **SYNC-received** After receiving the SYNC request from the client, the server acknowledges its state to the Established state.
3. **Closed-wait** The server changes its state from Established to close-wait after receiving the finish segment from the client. In this state the server sends an ACK and finish segments. Afterwards it changes the state to last-ACK.
4. **Last-ACK** In this state the server expects the last ACK segment from the client, as and when it receives the ACK segment it changes its state to again closed state.

## TCP Congestion Control

- Deals with end-to-end delivery.
- Congestion handling in TCP is based on three phases:
  - Slow start
  - Congestion avoidance
  - Congestion detection

### Slow start (exponential increase)

1. By default the receiver window size is initially set to 1.
2. In the first instance the transmitter receives an ACK for the window size indicating the receiver window size as 2 segments.
3. After 2 segments are sent it is acknowledged with 4 segments.
4. After 4 segments are sent it is acknowledged with 8 segments.
5. This is exponential growth and this growth continues until the window size reaches the threshold value.
6. If there are delayed ACKs, the increase in the size of the window is less than power of 2.

### Congestion avoidance (additive increase)

1. To avoid the congestion before it happens, the exponential growth of slow start algorithm must be slowed down.
2. When the threshold is reached, then the additive phase begins. Here each time the whole window of segments is acknowledged, the size of congestion window is increased by 1.

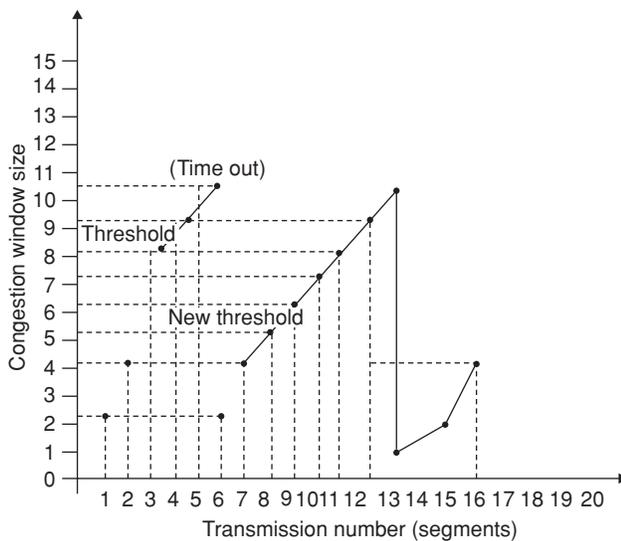
### Congestion detection (multiplicative decrease)

1. If congestion occurs, the congestion window size must be decreased. The only way the sender can guess the congestion has occurred is by the need to retransmit a segment.
2. Retransmission can occur in two cases:
  - (i) When a timer times out.
  - (ii) When 3 ACKs are received.

3. In both the cases the size of threshold is dropped to one-half of the current window size and the window size is decreased to initial window size “1”. This is multiplicative decrease.

**Example:** Let us take an example to explain the TCP congestion control.

Consider an instance of TCP additive increase, multiplicative decrease algorithm where the window size at the start of slow-start phase is 2 MSS (Maximum Segment Size) and threshold value is 8 MSS. The timeout occurs at the fifth transmission. Then what is the congestion window size at the end of the tenth transmission?



Window size is 2 MSS initially.

8 MSS is threshold value, after this there is only increase of 1-1 window size till timeout value which is 10.

The new threshold value becomes half of the value of current congestion window i.e., 5.

Timeout remains the same i.e., 10.

At 10th transmission the window size is 7.

After time-out, at 13th transmission window size = 1 and at 14th transmission window size = 2.

## TCP Flow Control

- For flow control sliding window protocol is used.
- The window size is set by the receiver and is controlled by the receiver. The window size is not fixed (variable).
- The sliding window protocol in TCP looks like the Go-Back-N protocol because it does not use NAKs; it looks like Selective Repeat because the receiver holds the out-of-order segments until the missing ones arrive.
- A sliding window is used to make transmission more efficient as well as to control the flow of data so that the destination does not become overwhelmed with data. TCP sliding windows are byte-oriented.

## TCP Error Control

- TCP provides reliability using error control.
- Error control includes mechanism for detecting corrupted segments, lost segments, out-of-order segments and duplicated segments.
- Error detection and correction in TCP is achieved through the use of three tools:
  - Checksum
  - Acknowledgment
  - Time-out

### Checksum

Each segment includes a checksum field which is used to check for a corrupted segment. A 16-bit checksum is mandatory in every segment.

### Acknowledgement

- There is no negative ACK in TCP.
- There is no ACK for the received ACK.
- Only the correctly received segments are acknowledged, if any segment is found to be corrupted through checksum such segments are not acknowledged.

### Time-out

Different timers are deployed for error control.

1. **Time-awated timer:** This timer is used to handle TCP termination process specially to handle duplicate finish segments. Its value is set to twice the life time of a segment.
2. **Keep-Alive Timer:** This timer is used to handle long idle TCP connections. By default its value is 2 hours, beyond which a probe (1 byte dummy data) is used for 10 consecutive times with a separation of 75 milliseconds. If there is no response beyond this, then the connection is terminated.
3. **Persistence Timer:** This timer is used to handle Zero(0) window size scenario. The sender sends 1 probe every 60 seconds until it receives a non-zero window size from where the communication resumes.
4. **Retransmission Timer:** This timer is used for handling any lost segments. Its value is twice the Round trip time, i.e.,  $2 \times RTT$ . RTT is time needed for a segment to reach a destination and for an acknowledgement to be received.

## APPLICATION LAYER

An interface between the networks is called application. This section introduces two important concepts:

- Application Layer: The application layer of the OSI model provides the first step of getting data onto the network.

- **Application Software:** Applications are the software programs used by people to communicate over the network. Examples of application software, includes HTTP, FTP, e-mail, and others, used to explain the differences between these two concepts.

In the OSI model, information is passed from one layer to the next, starting at the application layer on the transmitting host and proceeding down the hierarchy to the physical layer, then passing over the communications channel to the destination host, where the information proceeds back up the hierarchy, ending at the application layer.

The following six steps explain the procedure:

1. People create the communication.
2. The application layer prepares human communication for transmission over the data network.
3. Software and hardware converts communication to digital format.
4. Application layer services initiate the data transfer.
5. Each layer plays its role. The OSI layers encapsulate data down the stack. Encapsulated data travels across the media to the destination. OSI layers at the destination unencapsulate the data.
6. The application layer receives data from the network and prepares it for human use.

The application layer, layer 7, is the top layer of both the OSI and TCP/IP models. Layer 7 provides the interface between the application you use to communicate and the underlying network over which your messages are transmitted. Application layer protocols are used to exchange data between programs, running on the source and destination hosts.

### TCP/IP Application Layer Protocol

The most widely known TCP/IP application layer protocols are those that provide the exchange of user information. These protocols specify the format and control information necessary for many of the common internet communication functions. Among these, TCP/IP protocols are the following.

- Domain name system (DNS) is used to resolve internet names to IP addresses.
- Hypertext transfer protocol (HTTP) is used to transfer files that make up the web pages of the world wide web.
- Simple mail transfer protocol (SMTP) is used for the transfer of mail messages and attachments.
- Telnet, a terminal emulation protocol, is used to provide remote access to servers and networking devices.
- File transfer protocol (FTP) is used for interactive file transfers between systems.

### Application Layer Services

Programs such as file transfer or network print spooling, might need the assistance of application layer services to use network resources. Although transparent to

the user, these services have interface with the network and prepares the data for transfer. Different types of data whether it is text, graphics or video require different network services to ensure that it is properly prepared for processing by the functions occurring at the lower layers of OSI model. Application layer services establish an interface to the network and protocols provide the rules and formats that govern how data is treated, a single executable program can use all three components. For example, while discussing “Telnet”, you could be referring to the Telnet application, the Telnet service, or the Telnet protocol.

### Application Layer Protocol Functions

Both the source and destination devices use application layer protocols during a communication session. For the communications to be successful, the application layer protocols implemented on the source and destination host must match.

### Protocols perform the following tasks

- Establish consistent rules for exchanging data between applications and services loaded on the participating devices.
- Specifies how data inside the messages is structured and the types of messages that are sent between source and destination. These messages can be requests for services, acknowledgements, data messages, status messages, or error messages.
- Defines message dialogues, ensuring that a message being sent is met by the expected response and that the correct services are invoked when data transfer occurs.

Applications and services can also use multiple protocols in the course of a single conversation. One protocol might specify how to establish the network connection and another might describe the process for the data transfer when the message is passed to the next lower layer.

A single application can employ many different supporting application layer services. Thus, what appears to the user as one request for a web page might, in fact, amount to dozens of individual requests. For each request, multiple processes can be executed. For example, the FTP requires a client to initiate a control process and a data stream process to a server. Additionally, servers typically have multiple clients requesting information at the same time, as shown in the figure below. For example, a Telnet server can have many clients requesting connections to it. These individual client requests must be handled simultaneously and separately for the network to succeed. The application layer processes and services rely on support from lower layer functions to successfully manage the multiple conversations.

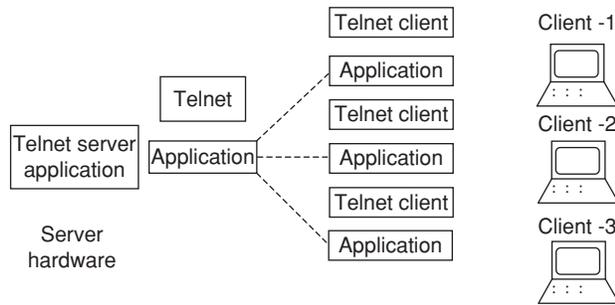


Figure 6 Multiple client's service Requests

## APPLICATION LAYER PROTOCOLS

The transport Layer uses an addressing scheme called a port number. Port numbers identify application layer services that are source and destination of data. Server programs generally use predefined port numbers that are commonly known by clients.

Some of these services are

- Domain Name System (DNS): TCP/UDP Port 53
- HTTP: TCP Port 80
- Simple Mail Transfer Protocol (SMTP): TCP Port 25
- Post office Protocol (POP): UDP Port 110
- Telnet: TCP Port 23
- DHCP: UDP Port 67
- FTP: TCP Ports 20 and 21

## Internet Control Message Protocol (ICMP)

- Used by hosts and gateways to send notification of datagram problems back to the sender.
- Used for error reporting and query messages.
- Helpful in network debugging.
- Uses the services of TCP and UDP with the port number 7 as the ping command which is used for testing, this testing is done from a source which starts at the application layer and reaches network through transport layer.
- ICMP is encapsulated into an IP datagram and then transmitted into the network, if the protocol field in the IP datagram is 1 then the IP datagram is said to be carrying ICMP message.

## Types of messages

### Error reporting

- **Destination Unreachable:** The packet is discarded due to the host not present in the network or the host is not responding to the request.
- **Source Quench:** The packet is discarded due to the congestion in the network.
- **Parameter Problem:** The packet is discarded due to the processing problem observing a change in the header format of the I/P datagram.
- **Time Exceeded:** The packet is discarded because the TTL value is decremented to zero(0).

- **Redirection:** Here the packet is not discarded but redirected to a network as the host doesn't belong to this network.

### Query message

**Router solicitation and router advertisement request and reply:** Router solicitation is a request generated by the source requesting the router's presence in the network.

The response is a router advertisement generated by the router broadcasting its network id and its presence in the network.

**Address mask request and reply:** If by any means the node is unable to identify the network bits in its I/P address then this request is used by the source to a router requesting for the network id, the reply is also unicast in this scenario.

**Time stamp echo request and reply:** This is used to calculate the round trip time of a packet for network diagnose or debugging.

**Echo request and reply:** This is used to see the presence of a host or a router in the network. For example PING.

## SMTP

- SMTP stands for simple mail transfer protocol.
- It uses the services of TCP on port number 25.
- It is a push protocol. Even when the destination is not interested to receive the message this push approach of the SMTP makes the receiver receive the message.
- Components of SMTP:

### 1. User Agent (UA) :

- (i) It provides Graphical User Interface access to the user.

**Example:** Netscape navigation, Mozilla Firefox. It also provides command-driven access in early days.

- (ii) It handles the inbox transactions:

- (a) Composing messages: Helps the user compose the e-mail message to be sent out.
- (b) Reading messages: Helps to read incoming messages by checking the mail in the incoming mail box.
- (c) Replying to messages: Sends the message to the sender or recipients of the copy.
- (d) Forwarding messages: Sends the message to a third party.
- (e) Handling mailboxes: Two mailboxes, an inbox and an outbox are created by the user agent. The inbox keeps all the received e-mails until they are deleted by the user. The outbox keeps all the sent e-mails until the user deletes them.

2. Mail transfer agent (MTA): The actual mail is transferred using MTA.
3. Multipurpose Internet mail extension (MIME): By default SMTP uses ASCII format for transaction. But few languages like Japanese, German etc do not support ASCII format. Hence for carrying non-ASCII form of transactions MIME is used in conjunction with SMTP. Thus, MIME is a set of software functions that transforms non-ASCII data (stream of bits) to ASCII data and vice-versa.
4. Mail access protocol (MAP): MAP is a pull approach where the emails of a client are retrieved from the mail server i.e., it is used to retrieve the clients emails from the mail server.

Two protocol of MAP are

- (i) POP 3 (Post Office Protocol)
- (ii) IMAP4 (Internet MAP)

### POP3

1. It is a pull protocol.
2. It uses the services of TCP on port number 110.
3. POP3 has several drawbacks and hence it is currently not in use.
  - A user cannot have different folders on the server.
  - A user cannot partially check the contents of the mail before downloading.
  - A user cannot search a mail with a keyword.
  - The user is not allowed to organize the mail on the server.
- (4) Modes of POP3
  - (i) Copy mode: The mails are copied from the mail server onto the client.
  - (ii) Delete mode: The mails are transferred from the mail server to the client and deleted at the mail server. By default POP3 uses delete mode.

### IMAP 4

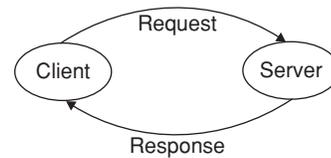
To overcome the drawbacks of POP3, IMAP4 is in current use. It provides the following functions:

1. A user can create, delete or rename mail boxes on the mail server.
2. A user can create a hierarchy of mailboxes in a folder.
3. A user can partially download e-mail.
4. A user can check the e-mail header before downloading and can search the contents of the e-mail for any specific character prior to downloading.

### HTTP

- HTTP stands for Hyper Text Transfer Protocol.
- It uses the services of TCP on well known port 80.
- It is a protocol mainly used to access data on the World Wide Web (www).

- HTTP functions as a combination of FTP and SMTP.
- It uses only one TCP connection, there is no separate control connection, only data is transferred between the client and the server.
- HTTP messages are read and interpreted by the HTTP server and HTTP client (browser).
- It works on two commands request and reply.
- It is a stateless protocol as it does not have any mapping from one transaction onto the other and treats a request and reply as a pair every time.



HTTP1.1 has several request types called methods:

1. GET: Requests a document from the server.
  2. HEAD: Requests information about a document but not the document itself.
  3. POST: Sends some information from the client to the server.
  4. PUT: Sends a document from the server to the client.
  5. TRACE: Echoes the incoming request.
  6. CONNECT: Reserved.
  7. OPTION: Inquires about available options.
- HTTP supports proxy servers. A proxy server is a computer that keeps copies of responses to recent requests. This reduces the load on the original server, decreases traffic and improves latency.
  - HTTP Connections:
    - (i) *Non-persistence*: In this connection approach for every request and reply (response) as a pair, a separate TCP connection is established every time. It suffers from slow start process. This was present in http version 1.0. Two RTTS are required to fetch each object.
    - (ii) *Persistence*: Here a single TCP connection is set on which multiple request and response can be made. This is observed from http version 2.0 onwards (apache http server). For http/1.1 is default. Hence we have reduced network congestion and faster content delivery.

### File Transfer Protocol (FTP)

- FTP uses the services of TCP.
- It needs two TCP connections:
  - Uses well-known port 21 for the control connection.
  - Uses well-known port 20 for the data connection.

- Mode of access:
  - FTP(TCP) – requires username and password.
  - TFTP(UDP) – requires no username and password.
- Types of files supported by FTP:
  - ASCII: By default FTP follows ASCII mode for file transfer. It is composed of 7-bit + 1 parity bit.
  - EBCDIC: If any node supports EBCDIC then this type of technique is used for file transfer. BCDIC supports 8 bits data format and is used in IBM. There is no error control i.e., there is no parity bit.
  - Image file: If the file to be sent is very large then continuous streams of 0s and 1s are sent to the transport layer. This is image file. Here FTP does not care of code, it is done by the lower layers.
- **Transmission mode of FTP:** FTP can transfer a file across the data connection by using one of the following three transmission modes:
  - Stream mode: This is the default mode. Data are delivered from FTP to TCP as a continuous stream of bytes.
  - Block mode: Data is delivered from FTP to TCP in blocks. Each block is preceded by a 3-byte header. The first byte is called the block descriptor, the next two bytes define the size of the block in bytes.
  - Compressed mode: If the file is big then the data is compressed. The compression method which is mostly used is run-length encoding. Consecutive appearances of a data unit are replaced by one occurrence and the number of repetitions. In a binary file, null characters are compressed.

## DNS

- Stands for Domain Name System.
- The DNS is a client/server application that identifies each host on the Internet with a unique user-friendly name i.e., it is used to map an Uniform Resource Locator (URL) to an IP address.
- DNS can use the services of UDP or TCP using the well-known port 53.
- If the size of the response message is more than 512 bytes, it uses the TCP connection.
- When the size of the response message is less than 512 bytes, UDP connection is used. Even though the size of message is not known then also UDP can be used. The UDP server will truncate the message if the message size is more than 512 bytes.
- DNS organizes the namespace in a hierarchical structure to decentralize the responsibilities involved in naming.
- DNS can be pictured as an inverted hierarchical tree structure with one root node at the top and a maximum of 128 levels. Each node in the tree has a domain name.

For example, on the Internet, the domain names, such as `http/www.cisco.com`, are much easier for people to

remember than `198.132.219.25`. Also if, cisco decides to change the numeric address, it is transparent to the user, because the domain name will remain `http/www.cisco.com`. The new address will simply linked to the existing domain name and connectivity is maintained as shown in the figure.

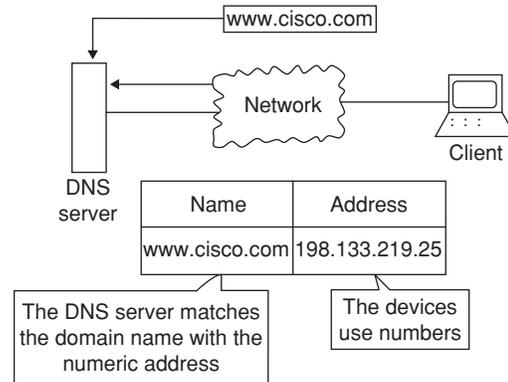


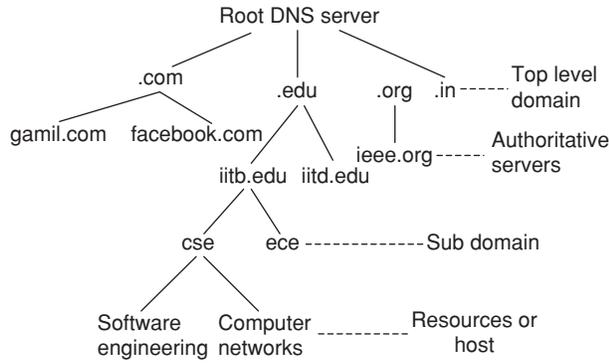
Figure 7 DNS addresses

When networks were small, it was a simple task to maintain the mapping between domain names and the addresses they represent. However, as networks began to grow and the number of devices increased, this manual system became unworkable. DNS was created for domain name to address resolution for these networks. DNS uses a distributed set of servers to resolve the names associated with these numbered addresses.

The DNS protocol defines an automated service that matches resource names with the required numeric network address. It includes the format for queries, responses, and data formats. DNS protocol communications use a single format called a message. This message format is used for all types of client queries and server responses, error messages, and the transfer of resource record information between servers. DNS is a client/server service, however, it differs from the other client/server services. Where as other services use a client that is an application (Web browser, e – mail, client, and so on) the DNS client runs as a service itself. The DNS client, sometimes called the DNS resolver, supports name resolution for the other network applications and other services that need it.

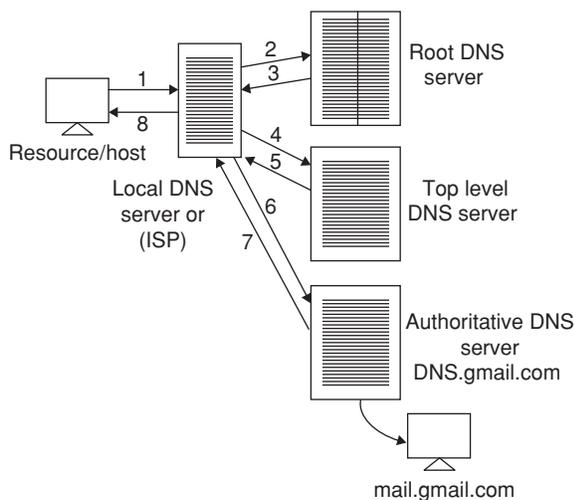
When configuring a network device, you generally provide one or more DNS server addresses that the DNS client can use for name resolution. Usually the Internet Service Provider (ISP) gives you the address to use for the DNS servers. When a user's application requests to connect to a remote device by name, the requesting DNS client queries one of these DNS servers to resolve the name to a numeric address.

- The domain name space consists of a tree of domain names. Each node or leaf in the tree has zero or more resource records, which holds information associated with the domain name. The tree sub-divides into zones beginning at the root zone. A DNS zone consists of a collection of connected nodes authoritatively served by an authoritative name server.



### Components of DNS

1. Root DNS Server : Root name servers keep track of all the authoritative name servers of each of the top level domain (TLD) name servers.
2. Top Level Domain: It provides the information regarding the presence of different zone files like
  - (i) based on geographical location (country domain): us—for United Statesm, in—for India
  - (ii) based on general attributes (generic domain):
    - com—used by commercial organization  
Example, gmail.com
    - .edu—used by educational institutes
    - .org—used by non-profit organizations  
Example, ieee.org
    - .gov—used by government institutions  
Example, nasa.gov
    - .mil—used by military organizations  
Example, army.mil
3. Zones: The TLD and the domains under TLD are divided into smaller units with the help of delegation. The domain is divided into small units, so that it can be managed easily. These small units are zones.
4. Authoritative DNS servers checks whether authoritative name servers are located in the DNS hierarchy.



### Dns Resource Records (RR)

- Every domain, whether it is a TLD, subdomain or single host have a set of resource records associated with it in the DNS distributed data base.
- Resource Records provide the mapping of host name to IP address. When a query is made to the DNS server, the host or server. who sends that query receives a response which is nothing but the resource record associated with it.
- A Resource Record (RR) is a 5 tuple that contains (Name, Time to live, class, Type, Value)
  - (i) Name: It is the domain name to which this RR belongs to. More than one resource records may exists for the same domain.
  - (ii) Time to live: The TTL is measured in seconds and it is a 32-bit integer.
  - (iii) Class: This field contains the value 'IN' which tells whether this record is used by internet or not.
  - (iv) Type: Defines type of RR address, name service, canonical name.
  - (v) Value: This field can be a number, ASCII strings or any domain.

## NETWORKING DEVICES

### Repeater

In digital communication systems, a repeater is a device that receives a digital signal on an electromagnetic or optical transmission medium and regenerates the signal. Repeaters remove the unwanted noise in an incoming signal. Unlike an analog signal, the original digital signal, even if weak or distorted, can be clearly perceived and restored. With analog transmission, signals are re strengthened with *amplifiers* which unfortunately also amplify noise as well as information.

### Hub

A hub is the central part of a wheel where the spokes come together. The term is familiar to frequent fliers who travel through airport “hubs” to make connecting flights from one point to another. In data communications, a hub is a place of convergence where data arrives from one or more directions and is forwarded out in one or more other directions. A hub usually includes a switch of some kind. (And a product that is called a “switch” could usually be considered a hub as well.)

### Switch

In a telecommunications network, a switch is a device that channels incoming data from any of multiple input ports to the specific output port that will take the data towards its intended destination. In the traditional circuit-switched telephone network, one or more switches are used to set up a dedicated though temporary connection or circuit for an exchange between two or more parties.

In the open systems Interconnection (OSI) communications model, a switch performs the Layer 2 or Data-link layer

function. That is, it simply looks at each packet or data unit and determines from a physical address (the “MAC address”) which device a data unit is intended for and switches it out towards that device. However, in wide area networks such as the Internet, the destination address requires a look-up in a routing table by a device known as a router. Some newer switches also perform routing functions (Layer 3 or the Network layer functions in OSI) and are sometimes called IP switches. On larger networks, the trip from one switch point to another in the network is called a hop. The time a switch takes to figure out where to forward a data unit is called its latency. The price paid for having the flexibility that switches provide in a network is this latency. In the simplest networks, a switch is not required for messages that are sent and received within the network. For example, a local area network may be organized in a token ring or bus arrangement in which each possible destination inspects each message and reads any message with its address.

## Bridge

A bridge is a product that connects a local area network (LAN) to another local area network that uses the same protocol (for example, Ethernet or token ring). You can envision a bridge as being a device that decides whether a message from you to someone else is going to the local area network in your building or to someone on the local area network in the building across the street. A bridge examines each message on a LAN, passing those to be within the same LAN and forwarding those known to be on the other interconnected LAN (or LANs).

In bridging networks, computer or node addresses have no specific relationship to location. For this reason, messages are sent out to every address on the network and accepted only by the intended destination node. Bridges learn which addresses are on which network and develops a *learning table* so that subsequent messages can be forwarded to the right network.

Bridging networks are generally always interconnected local area networks since broadcasting every message to all possible destinations would flood a larger network with unnecessary traffic. For this reason, router networks such as the Internet uses a scheme that assigns addresses to nodes so that a message or packet can be forwarded only in one general direction rather than forwarded in all directions. A bridge works at the data-link (physical network) level of a network, copying a data frame from one network to the next network along the communications path. A bridge is sometimes combined with a router in a product called a brouter.

## Routers

Routers operate on the Network layer, which is a higher level in the OSI conceptual model. Routers use a combination of

software and hardware, but it is used to route data from its source to its destination. Routers actually have a sophisticated OS that allows them to configure various connection ports. You can setup a router to route data packets from different network protocol stacks, which include TCP/IP, IPX/SPX and AppleTalk.

Routers are also used to connect remote LANs together using different WAN technologies. But, when a router has become large, the large network is divided into logical segments called subnets. This division of the network is based on the addressing scheme related to a particular subnet is kept local. The router only forwards data that is meant for the subnets on the extended network.

Routers also help to decide how to forward data packets to their destination based on the routing table. The protocols built into the router’s operating system is used to identify neighboring routers and their network addresses. This allows routers to build a routing table.

## Brouter

A brouter is a network bridge and a router combined in a single product. A bridge is a device that connects one local area network (LAN) to another local area network that uses the same protocol (for example, Ethernet or token ring). If a data unit on one LAN is intended for a destination on an interconnected LAN, the bridge forwards the data unit to that LAN; otherwise, it passes it along the same LAN. A bridge usually offers only one path to a given interconnected LAN. A router connects a network to one or more other networks that are usually part of a wide area network and may offer a number of paths out to destinations on those networks. A router therefore needs to have more information than a bridge about the interconnected networks. It consults a routing table for this information. Since a given outgoing data unit from a computer may be intended for an address on the local network, on an interconnected LAN, or the wide area network, it makes sense to have a single unit that examines all data units and forwards them appropriately.

## Gateway

A gateway is a network point that acts as an entrance to another network. On the Internet, a node or stopping point can be either a gateway node or a host (end-point) node.

In the network for an enterprise, a computer server acting as a gateway node is often also acting as a proxy server and a firewall server. A gateway is often associated with a router, which knows where to direct a given packet of data that arrives at the gateway, and a switch, which furnishes the actual path in and out of the gateway for a given packet.

## EXERCISES

## Practice Problems I

**Directions for questions 1 to 15:** Select the correct alternative from the given choices.

- If TCP RTT is currently 40 m/sec and the following acknowledgements come in after 26, 32 and 24 m/sec respectively. What is the new RTT estimate?  $\alpha = 0.9$ .  
(A) 32.69 (B) 24.31 (C) 36.55 (D) 42.23
- If a TCP connection is transferring a file of 5000 bytes. The first byte is numbered 1001. What are the sequence numbers for each segment if data is sent in five segments, each carrying 1000 bytes?  
(A) 1001, 2001, 3001, 4001, 5001  
(B) 1000, 2000, 3000, 4000, 5000  
(C) 5000, 6000, 7000, 8000, 9000  
(D) 5001, 6001, 7001, 8001, 9001
- Which of the below statements hold good with respect to routing done by a bridge?  
(i) they can route packets using IP addresses  
(ii) they use data link layer addresses to do routing  
(iii) the LAN route IPv4, IPv6, Apple Talk, ATM, OSI packets  
(A) (i), (ii) (B) (ii), (iii)  
(C) (i), (iii) (D) (i), (ii), (iii)
- Match the following:  

|     |           |   |                                       |
|-----|-----------|---|---------------------------------------|
| i   | Repeaters | p | connects different nodes of a LAN     |
| ii  | Hub       | q | amplifies the signal between segments |
| iii | Switch    | r | connects different LANs               |
| iv  | Bridge    |   |                                       |

  
 (A) i – q ii – r iii – p iv – r  
 (B) i – r ii – p iii – q iv – q  
 (C) i – q ii – p iii – p iv – r  
 (D) i – p ii – p iii – q iv – r
- Match the following.  

|     |                      |   |                                                                                                         |
|-----|----------------------|---|---------------------------------------------------------------------------------------------------------|
| i   | Retransmission timer | p | goes off when a TCP connection is idle for a long time                                                  |
| ii  | Keep-alive timer     | q | goes off if sender and receiver are waiting for each other                                              |
| iii | Persistence timer    | r | goes off to trigger the delivery of a segment in case acknowledgement is not received for first attempt |

  
 (A) i – r ii – q iii – p  
 (B) i – q ii – r iii – p  
 (C) i – p ii – r iii – q  
 (D) i – p ii – q iii – r
- In T/TCP (Transactional TCP) what does the packet that is sent by client, consist of?  
(i) SYN (ii) REQUEST (iii) FIN  
(A) (i), (ii) (B) (ii), (iii)  
(C) (i), (iii) (D) (i), (ii), (iii)

- Assume TCP uses 32-bit sequence numbers and sequence numbers are given to each byte that gets transmitted. If data is transmitted at 1 Gbps. What is the wraparound time for sequence numbers?  
(A) 14.4 sec (B) 24.24 sec  
(C) 34.36 sec (D) 44.45 sec
- What are the disadvantages of NAT?  
(i) NAT forms link between sender and receiver and then link can be broken irreparably during a connection.  
(ii) NAT violates architectural model of IP.  
(iii) NAT hacks source port field of TCP header which is of limited size.  
(iv) NAT alleviates IP shortage.  
(A) (i), (ii), (iii) (B) (ii), (iii), (iv)  
(C) (i), (iii), (iv) (D) (i), (ii), (iv)
- What is the main protocol in the transport layer?  
(A) TCP (B) UDP  
(C) FTP (D) Both (A) and (B)
- Number of bytes for header in UDP segment and TCP segment are  
(A) 8 bytes, 20 bytes (B) 16 bytes, 16 bytes  
(C) 32-bits, 20-bits (D) None of these
- TCP maintains a variable RTT (Round trip time), for determining the time to reach destination and receiving acknowledgement, the formula for RTT is  
(A)  $RTT = RTT + D$   
(B)  $RTT = 4RTT$   
(C)  $RTT = \alpha RTT + (1 - \alpha) M$  ( $\alpha = 7/8$ )  
(D) None of these.
- Maximum segment size is  
(A) The size of the segment without header.  
(B) The size of the segment with limit.  
(C) The transmission link capacity.  
(D) Less than maximum transfer unit.
- What is meant by silly window syndrome that ruins TCP performance?  
(A) This occurs when sender sends data in large blocks and receiver receives in large blocks.  
(B) This occurs when sender sends data in large blocks and receiver receives in or reads one byte at a time.  
(C) Both (A) and (B)  
(D) None of these

**Common data for questions 14 and 15:** A TCP segment begins with a fixed-format, 20-byte header. The header is followed by reader options. After the options, upto 65,495 bytes of data may follow.

- Number of one bit flags available in the TCP header are  
(A) 5 (B) 6  
(C) 2 (D) None of these
- Which of the flags is used for establishing connections?  
(A) PSH (B) ACK (C) URG (D) SYN

## Practice Problems 2

**Directions for questions 1 to 15:** Select the correct alternative from the given choices.

- Which of the below TCP primitives block a port?  
(i) LISTEN (ii) CONNECT (iii) RECEIVE  
(A) (i), (ii) (B) (ii), (iii)  
(C) (i), (iii) (D) (i), (ii), (iii)
- In the context of TCP sockets how is a symmetric DISCONNECT different from that of an asymmetric one?  
(i) In symmetric DISCONNECT each direction is closed separately.  
(ii) In asymmetric DISCONNECT each direction is closed separately.  
(iii) In asymmetric DISCONNECT transport user can release the connection  
(A) (i), (ii) (B) (ii), (iii)  
(C) (i), (iii) (D) (i), (ii), (iii)
- When does RPC/UDP does not make a good combination?  
(i) When the caller and callee machines are separated by small network distance.  
(ii) When the parameters of the procedures are too huge in size.  
(iii) When the procedure requested cannot be repeated safely as needed.  
(A) (i), (ii) (B) (ii), (iii)  
(C) (i), (iii) (D) (i), (ii), (iii)
- Which of the following statements below are true with reference to RTP (Real Time Transport Protocol)?  
(i) It multiplexes server real time data stream into a single stream of UDP packets.  
(ii) RTP has flow control, error control mechanism.  
(iii) RTP has no mechanism for retransmission.  
(A) (i), (ii) (B) (ii), (iii)  
(C) (i), (iii) (D) (i), (ii), (iii)
- What does RTCP (real time transport control protocol) accomplish?  
(i) Provides feedback on delay, jitter etc to sources.  
(ii) Handles intostream synchronization.  
(iii) Provides a way to name the sources.  
(A) (i), (ii) (B) (ii), (iii)  
(C) (i), (iii) (D) (i), (ii), (iii)
- Which of the following are applicable to TCP?  
(i) Breaks the data coming from upper layers into 64 kbyte size packets and transmits them.  
(ii) Manages the time out and re-uses them.  
(iii) Should reassemble the packets in correct order at receiving end.  
(iv) TCP supports multicasting.  
(A) (i), (ii), (iii) (B) (ii), (iii), (iv)  
(C) (i), (iii), (iv) (D) (i), (ii), (iv)
- Which of the below statements about sockets is/are true?  
(i) For sender and receiver to avail TCP service sockets have to be created.  
(ii) Each socket is a 16 bit number local to that host.  
(iii) Sockets can involve themselves in one connection at a time.  
(iv) Ports below 1024 are reserved.  
(A) (i), (ii), (iii) (B) (ii), (iii), (iv)  
(C) (iii), (iv), (i) (D) (i), (ii), (iv)
- What are the functions of application layer?  
(A) Mail service provides a basis for electronic mails forwarding and storage  
(B) File access transfer and management  
(C) Creates virtual terminal that allows us to log onto remote host  
(D) All the above
- Which of the following application uses UDP?  
(A) Streaming a multimedia  
(B) Client-server interaction  
(C) Internet telephony  
(D) All the above
- What are the reasons for choosing an UDP by an application?  
(A) No connection establishment  
(B) No connection state  
(C) Small packet header  
(D) All the above
- TCP uses multiple timers to do its work, the timers are  
(A) Retransmission timer  
(B) Persistence timer  
(C) Keep alive timer  
(D) All the above
- Which of the following is supported by TCP connections?  
(A) Full-duplex (B) Point-to-point  
(C) Multicasting (D) Both (A) and (B)
- TCP connection is \_\_\_\_\_ stream.  
(A) Byte (B) Message  
(C) Packet (D) None of these.
- If a sender wants to indicate that, it has no data for the receiver, one of the following bits is set.  
(A) PSH (B) RST  
(C) FIN (D) ACK
- If the receiver host is responding by sending a primitive SYN ( $SEQ = y$ ,  $ACK = x + 1$ ) means  
(A) The receiver data sequence number is  $y$ .  
(B) It has received up to  $x + 1$  bytes of data.  
(C) Both (A) and (B)  
(D) None of these

## PREVIOUS YEARS' QUESTIONS

1. The transport layer protocols used for real time multimedia, file transfer, DNS and email respectively are [2013]
  - (A) TCP, UDP, UDP and TCP
  - (B) UDP, TCP, TCP and UDP
  - (C) UDP, TCP, UDP and TCP
  - (D) TCP, UDP, TCP and UDP
2. Which one of the following socket API functions converts an unconnected active TCP socket into a passive socket? [2014]
  - (A) Connect
  - (B) Bind
  - (C) Listen
  - (D) Accept
3. Suppose two hosts use a TCP connection to transfer a large file. Which of the following statements is/are FALSE with respect to the TCP connection? [2015]
  - I. If the sequence number of a segment is  $m$ , then the sequence number of the sub sequent segment is always  $m + 1$ .
  - II. If the estimated round trip time at any given point of time is  $t$  sec, the value of the retransmission timeout is always set to greater than or equal to  $t$  sec.
  - III. The size of the advertised window never changes during the course of the TCP connection.
  - IV. The number of unacknowledged bytes at the sender is always less than or equal to the advertised window.
  - (A) III only
  - (B) I and III only
  - (C) I and IV only
  - (D) II and IV only
4. In one of the pairs of protocols given below, both the protocols can use multiple TCP connections between the same client and the server. Which one is that? [2015]
  - (A) HTTP, FTP
  - (B) HTTP, TELNET
  - (C) FTP, SMTP
  - (D) HTTP, SMTP
5. Assume that the bandwidth for a TCP connection is 1048560 bits/sec. Let  $\alpha$  be the value of RTT in milliseconds (rounded off to the nearest integer) after which the TCP window scale option is needed. Let  $\beta$  be the maximum possible window size with window scale option. Then the values of  $\alpha$  and  $\beta$  are [2015]
  - (A) 63 milliseconds,  $65535 \times 2^{14}$
  - (B) 63 milliseconds,  $65535 \times 2^{16}$
  - (C) 500 milliseconds,  $65535 \times 2^{14}$
  - (D) 500 milliseconds,  $65535 \times 2^{16}$
6. Consider the following statements
  1. TCP connections are full duplex
  2. TCP has no option for selective acknowledgement
  3. TCP connections are message streams
  - (A) Only 1 is correct
  - (B) Only 1 and 3 are correct
  - (C) Only 2 and 3 are correct
  - (D) All of 1, 2 and 3 are correct
7. Which one of the following protocols is **NOT** used to resolve one form of address to another one? [2016]
  - (A) DNS
  - (B) ARP
  - (C) DHCP
  - (D) RARP
8. Which of the following is/are example(s) of stateful application layer protocols? [2016]
  - (i) HTTP
  - (ii) FTP
  - (iii) TCP
  - (iv) POP3
  - (A) (i) and (ii) only
  - (B) (ii) and (iii) only
  - (C) (ii) and (iv) only
  - (D) (iv) only
9. Identify the correct sequence in which the following packets are transmitted on the network by a host when a browser requests a webpage from a remote server, assuming that the host has just been restarted. [2016]
  - (A) **HTTP GET** request, **DNS** query, **TCP SYN**
  - (B) **DNS** query, **HTTP GET** request, **TCP SYN**
  - (C) **DNS** query, **TCP SYN**, **HTTP GET** request
  - (D) **TCP SYN**, **DNS** query, **HTTP GET** request
10. Consider a TCP client and a TCP server running on two different machines. After completing data transfer, the TCP client calls **close** to terminate the connection and a FIN segment is sent to the TCP server. Server-side TCP responds by sending an ACK, which is received by the client-side TCP. As per the TCP connection state diagram (RFC 793), in which state does the client-side TCP connection wait for the FIN from the server-side TCP? [2017]
  - (A) LAST-ACK
  - (B) TIME-WAIT
  - (C) FIN-WAIT-1
  - (D) FIN-WAIT-2
11. Consider socket API on a Linux machine that supports connected UDP sockets. A connected UDP socket is a UDP socket on which **connect** function has already been called. Which of the following statements is/are CORRECT? [2017]
  - I. A connected UDP socket can be used to communicate with multiple peers simultaneously.
  - II. A process can successfully call **connect** function again for an already connected UDP socket.
  - (A) I only
  - (B) II only
  - (C) Both I and II
  - (D) Neither I nor II
12. Consider the following statements regarding the slow start phase of the TCP congestion control algorithm. Note that **cwnd** stands for the TCP congestion window

and MSS denotes the Maximum Segment Size.

- (i) The cwnd increases by 2 MSS on every successful acknowledgment.
- (ii) The cwnd approximately doubles on every successful acknowledgement.
- (iii) The cwnd increases by 1 MSS every round trip time.
- (iv) The cwnd approximately doubles every round trip time.

Which one of the following is correct? [2018]

- (A) Only (ii) and (iii) are true
- (B) Only (i) and (iii) are true
- (C) Only (iv) is true
- (D) Only (i) and (iv) are true

13. Consider a long-lived TCP session with an end-to-end bandwidth of 1 Gbps ( $= 10^9$  bits-per-second). The session starts with a sequence number of 1234. The minimum time (in seconds, rounded to the closest integer) before this sequence number can be used again is \_\_\_\_\_. [2018]

## ANSWER KEYS

### EXERCISES

#### Practice Problems 1

1. C    2. A    3. B    4. D    5. A    6. D    7. C    8. A    9. D    10. A  
11. C    12. D    13. B    14. B    15. D

#### Practice Problems 2

1. C    2. C    3. B    4. C    5. D    6. A    7. D    8. D    9. D    10. D  
11. D    12. D    13. A    14. C    15. C

#### Previous Years' Questions

1. C    2. C    3. B    4. A    5. C    6. A    7. C    8. C    9. C    10. D  
11. B    12. C    13. 34

# Chapter 4

## IP(v4)

### LEARNING OBJECTIVES

- IP addressing
- Class A
- Class B
- Class C
- Class D
- Subnet mask
- Classless Inter Domain Routing (CIDR)
- Network address
- Network Address Translation (NAT)
- IP-Protocol

### IP ADDRESSING

Every machine on the internet has a unique identification number, called an IP Address. A typical IP address looks like this:

216.27.61.137

To make it easier for humans to remember, IP addresses are normally expressed in decimal format as a “*dotted decimal number*” like the one above. But computers communicate in binary form. Look at the same IP address in binary:

11011000.00011011.00111101.10001001

The four numbers in an IP address are called octets, because they each have eight positions when viewed in binary form. If you add all the positions together, you get 32, which is why IP addresses are considered 32-bit numbers. Since each of the eight positions can have two different states (1 or 0) the total number of possible combinations per octet is  $2^8$  or 256. So each octet can contain any value between 0 and 255. Combine the four octets and you get  $2^{32}$  or a possible 4,294,967,296 unique values.

Out of the almost 4.3 billion possible combinations, certain values are restricted from use as typical IP addresses. For example, the IP address 0.0.0.0 is reserved for the default network and the address 255.255.255.255 is used for broadcasts.

The octets serve a purpose other than simply separating the numbers. They are used to create classes of IP addresses that can be assigned to a particular business, government or other entities based on size and need. The octets are split into two sections: Net and Host. The Net section always contains the first octet. It is used to identify the network that a computer belongs to. Host

(sometimes referred to as Node) identifies the actual computer on the network. The Host section always contains the last octet. There are five IP classes plus certain special addresses. They are

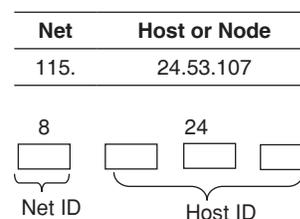
1. Class A
2. Class B
3. Class C
4. Class D
5. Class E

**Default Network:** The IP address of 0.0.0.0 is used for the default network.

### Class A

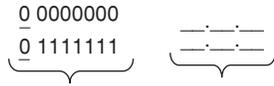
This class is for very large networks, such as, a major international company. IP addresses with a first octet from 1 to 126 are part of this class. The other three octets are used to identify each host. This means that there are 126 Class A networks each with 16,777,214 ( $2^{24} - 2$ ) possible hosts for a total of 2,147,483,648 ( $2^{31}$ ) unique IP addresses. Class A networks account for half of the total available IP addresses. In Class A networks, the high order bit value (the very first binary number) in the first octet is always 0.

**Example:**



$2^8 - 256$  Networks  
 $2^{24} - 1, 67, 77, 216$  hosts

- Used for large organizations, e.g., CISCO.
- x.x.x.x/8 – default mask of 8-bits in Network.

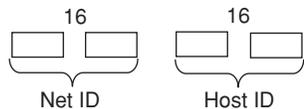


- 1st bit of 1st octet is '0'
- '1' bit is fixed and 2 special addresses are fixed  
 $\therefore 2^7 - 2$  Networks
- All 1's and all 0's are not used in host portion.  
 $\therefore 2^{24} - 2$  hosts

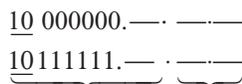
**Loopback** The IP address 127.0.0.1 is used as the loopback address. This means that it is used by the host computer to send a message back to itself. It is commonly used for troubleshooting and network testing.

### Class B

Class B is used for medium-sized networks. A good example is a large college campus. IP addresses with first octet from 128 to 191 are part of this class. Class B addresses also include the second octet as part of the Net identifier. The other two octets are used to identify each host. This means that there are 16,384 ( $2^{14}$ ) Class B networks each with 65,534 ( $2^{16} - 2$ ) possible hosts for a total of 1,073,741,824 ( $2^{30}$ ) unique IP addresses. Class B networks make up a quarter of the total available IP addresses. Class B networks have a first bit value of 1 and a second bit value of 0 in the first octet.



- Used for medium organization eg: universities, medium companies.
- x.x.x.x/16 – default mask of 16 bits in network



- 1st two bits of 1st octet is '10'
- 2 bits fixed. So  $2^{14}$  Networks,  $2^{16} - 2$  hosts.

### Class C

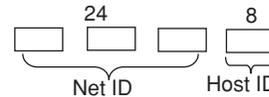
Class C addresses are commonly used for small to mid-size businesses. IP addresses with a first octet from 192 to 223 are part of this class. Class C addresses also include the second and third octets as part of the Net identifier. The last octet is used to identify each host. This means that there

are 2,097,152 ( $2^{21}$ ) Class C networks each with 254 ( $2^8 - 2$ ) possible hosts for a total of 536,870,912 ( $2^{29}$ ) unique IP addresses. Class C networks make up an eighth of the total available IP addresses. Class C networks have a first bit value of 1, second bit value of 1 and a third bit value of 0 in the first octet.

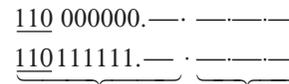
**Example:**

| Net        | Host or Node |
|------------|--------------|
| 195.24.53. | 107          |

- Used for small organizations, e.g., colleges



- x.x.x.x/24 – default mask of 24-bits in network.



- 1st 3-bits of 1st octet is '110', So  $2^{21}$  Networks,  $2^{24} - 2$  hosts.

### Class D

Used for multicasts, Class D is slightly different from the first three classes. It has a first bit value of 1, second bit value of 1, third bit value of 1 and fourth bit value of 0. The other 28-bits are used to identify the group of computers the multicast message is intended for. Class D accounts for 1/16th ( $268,435,456$  or  $2^{28}$ ) of the available IP addresses.

**Example:**

| Net  | Host or Node |
|------|--------------|
| 224. | 24.53.107    |

- Used for multicasting
- No Net ID or Host ID.
- Whole address is used for multicasting.

### Class E

Class E is used for experimental purposes only. Like Class D, it is different from the first three classes. It has a first bit value of 1, second bit value of 1, third bit value of 1 and fourth bit value of 1. The other 28-bits are used to identify the group of computers the multicast message is intended for. Class E accounts for 1/16th ( $268,435,456$  or 228) of the available IP addresses.

| Net  | Host or Node |
|------|--------------|
| 240. | 24.53.107    |

### Broadcast

- Messages that are intended for all Computers are broadcasted using the IP Address 255.255.255.255.

**Notes:**

- If we use class of IP address, we waste many IP's So, we use subnetting to save IP's.
- We can use following private addresses to save IP's.  
10.0.0.0 – 10.255.255.255 – class A network  
172.16.0.0. – 172.31.255.255 – 16 class B network  
192.168.0.0 – 192.168.255.255 – 256 class C network.
- Classful addressing is address with default mask.
- Classless address is a address with any other mask that is not default.

**SUBNET MASK**

**Mask**

The length of the net id and host id is predefined in classful addressing, we can use a mask, also called default mask, which is a 32-bit number made of contiguous 1's followed by contiguous 0's. The masks of classes A, B and C are shown below. Masking is not applicable to class D and class E.

1. Class A mask is  
1111 1111.0000 0000.0000 0000. 0000 0000  
255. 0. 0. 0  
CIDR Representation is '/8'.
2. Class B mask is  
1111 1111.1111 1111.0000 0000.0000 0000  
255. 255. 0. 0  
CIDR Representation is '/16'.
3. Class C mask is  
1111 1111.1111 1111. 1111 1111.0000 0000  
255. 255. 255. 0  
CIDR Representation is '/24'.

The mask will be helpful in finding the netid and the hostid. For Example the mask for class A address has eight 1's, which means the first 8 bits of any address in class A define the netid, the next 24 bits define the hostid.

**Subnetting packet structure** Internet protocol is layer-3 protocol in OSI, It takes data segments from layer-4 and divides it into packets.

IP packet encapsulates data unit received from above layer and add to its own header information.



**Figure 1** IP Encapsulation

The encapsulated data is referred to as IP-payload.

**Subnetting** Each IP-class has its own default mask, which bounds that IP class to have pre fixed number of Networks and prefixed number of Hosts per network.

- CIDR provides the flexibility of borrowing bits of Host part of the IP-address and use them as network within a network, called subnet.
- By using subnetting, one single class A IP-address can be used to have smaller sub networks, that provides better network management.

**Class-A subnets**

To make more subnets in class A, bits from Host part are borrowed and the subnet mask is changed accordingly.

**Table 1** Possible combination of class-A subnets

| Network bits | Subnet Mask     | Bits-borrowed | Subnets | Hosts/subnet |
|--------------|-----------------|---------------|---------|--------------|
| 8            | 255.0.0.0       | 0             | 1       | 16777214     |
| 9            | 255.128.0.0     | 1             | 2       | 8388606      |
| 10           | 255.192.0.0     | 2             | 4       | 4194302      |
| 11           | 255.224.0.0     | 3             | 8       | 2097150      |
| .            | .               | .             | .       | .            |
| .            | .               | .             | .       | .            |
| .            | .               | .             | .       | .            |
| 25           | 255.255.255.128 | 17            | 131072  | 126          |
| 27           | 255.255.255.224 | 19            | 524288  | 30           |
| 30           | 255.255.255.252 | 22            | 4194304 | 2            |

In case of subnetting also, the first and last IP address of every subnet is used for subnet number and subnet Broadcast IP-address. These 2 IP-addresses cannot be assigned to hosts, subnetting cannot be implemented using more than 30-bits as network bits, because that provides less than 2 hosts per subnet.

**Class-B subnets**

Class B IP addresses can be subnetted the same way as class-A addresses, that is by borrowing bits from Host bits.

**Table 2** Possible combination of class-B subnets

| Network bits | Subnet Mask     | Bits-borrowed | Subnets | Hosts/subnet |
|--------------|-----------------|---------------|---------|--------------|
| 16           | 255.255.0.0     | 0             | 1       | 65534        |
| 17           | 255.255.128.0   | 1             | 2       | 32766        |
| 18           | 255.255.192.0   | 2             | 4       | 16382        |
| 19           | 255.255.224.0   | 3             | 8       | 8190         |
| .            | .               | .             | .       | .            |
| .            | .               | .             | .       | .            |
| .            | .               | .             | .       | .            |
| 28           | 255.255.255.240 | 12            | 4096    | 14           |
| 29           | 255.255.255.248 | 13            | 8192    | 6            |
| 30           | 255.255.255.252 | 14            | 16384   | 2            |

**Class-C subnets**

Class C IP-addresses are usually assigned to a very small size network because it can only have 254 Hosts in a network.

**Table 3** Possible combination of class-C subnets

| Network Bits | Subnet Mask     | Bits-borrowed | Subnets | Hosts/subnet |
|--------------|-----------------|---------------|---------|--------------|
| 24           | 255.255.255.0   | 0             | 1       | 254          |
| 25           | 255.255.255.128 | 1             | 2       | 126          |
| 26           | 255.255.255.192 | 2             | 4       | 62           |
| 27           | 25.255.255.224  | 3             | 8       | 30           |
| 28           | 255.255.255.240 | 4             | 16      | 14           |
| 29           | 255.255.255.248 | 5             | 32      | 6            |
| 30           | 255.255.255.252 | 6             | 64      | 2            |

## CLASSLESS INTER DOMAIN ROUTING (CIDR)

The mask shown in the form  $'/n'$  where  $'n'$  can be 8, 16, 24 in classful Addressing, This notation is also called slash notation or class less Inter domain Routing (CIDR) notation. This notation is used in class less addressing.

The Internet is rapidly running out of IP addresses. In particular, the problem is with class 'B' network. For most organizations, a class A network, with 16 million addresses is too big, and a class C network, with 256 address is too small.

A class B network with 65, 536, is just right. In reality a class 'B' address is too large for most organizations. It is known that more than half of all class B networks have fewer than 50 hosts.

If the class B address had split, and allocated 20-bits for network number, another problem would have emerged, i.e., the routing table explosion. From the routers point of view, the IP address space is a 2-level hierarchy, with network numbers and host numbers.

One solution is CIDR. The basic Idea behind CIDR, is to allocate the remaining IP addresses in variable – sized blocks, without regard to the classes. If an organization needs 2000 addresses, it is given a block of 2048 address on a 2048-byte boundary.

Dropping the classes makes forwarding more complicated.

In classful addressing system, the forwarding is carried in the following way:

When a packet arrived at a router, a copy of the IP address was shifted right 28 bits to yield a 4-bit class number. A 16-way branch then sorted packet into A, B, C with eight of the cases for class A, four of the cases for class B, 2 of the cases for class 'C'. The code for each class then masked off the 8-, 16-, 24-bit network number and right aligned it in a 32-bit word. The network number was then looked up in the A, B or C table, by indexing for A and B networks and hashing for C networks. Once the entry was found, the outgoing line could be looked up and the packet forwarded.

The above described algorithm will not work for CIDR, Instead, each routing table entry is extended by giving it a 32-bit mask. There is only a single routing table for all networks. Consisting of an array of (IP-address, subnet mask, outgoing line) triples.

When a packet comes, its destination IP-address is first extracted. Then the routing table is scanned entry by entry, masking the destination address and comparing it to the table entry looking for a match.

It is possible that multiple entries match, in which case the longest mask is used. Thus, if there is a match for a/20 mask and a/24 mask, the/24 entry is used.

Complex algorithms have been devised to speed up the address matching process.

To overcome address depletion and give more organizations access to the internet, classless addressing was

designed and implemented. In this scheme, there are no classes, but the addresses are still granted in blocks.

In classless addressing, when an entity, small or large, needs to be connected to the Internet, it is granted a block of addresses. The size of the block varies based on the nature and size of the entity. For example, household may be given only 2 addresses, a large organization may be given thousands of addresses. An ISP(Internet Service Provider) provides addresses to customers.

To simplify the handling of addresses, the Internet authorities impose 3 restrictions on class less address blocks:

1. The addresses in a block must be contiguous.
2. The number of addresses in a block must be a power of 2 (1, 2, 4, 8, 16, ...).
3. The first address must be evenly divisible by the number of addresses.

A better way to define block of addresses is to select any address in the block and the mask. In classless addressing the mask for a block can take any value from 0 to 32.

In IPv4 addressing, a block of addresses can be defined as

$$x.y.z.w/n$$

in which  $x.y.z.w$  defines one of the addresses and the  $n$  defines the mask.

The address and the  $n$  notation completely define the whole block (First block address, Last address and the number of addresses)

**First Address** The first address in the block can be found by setting the  $(32-n)$  right most bits in the binary notation of the address to Zeros.

### Solved Examples

**Example 1:** A block of addresses is granted to a small organization, We know that one of the addresses is 209.17.38.40/28, what is the first address in the block?

**Solution 1:** The binary representation of the given address is: 209.17.38.40

$$11010001.00010001.00100110.00101000$$

If we set  $(32-n) = 32 - 28 = 4$  right most bits to '0', we get

$$11010001.00010001.00100110.00100000$$

First address:

$$209. \quad 17. \quad 38. \quad 32$$

**Solution 2:** Another way to find the first address is to represent the mask as a 32-bit binary number.  $'/28'$  can be represented as

$$11111111.11111111.11111111.11110000$$

(28 ones and 4 zeros)

The first address can be found by ANDing the given address with the mask.

Address: 11010001.00010001.00100110.00101000

Mask: 11111111.11111111.11111111.11110000

First Address:

11010001.00010001.00100110.00100000

209. 17. 38. 32

**Last address and number of addresses** The last address in the block can be found by setting the  $(32-n)$  right most bits in the binary notation of the address to 1's.

**Example 2:** Find the last address and Number of addresses in the block, for the given CIDR, in the above example?

**Solution 1:** The binary representation of the given address is

11010001.00010001.00100110.00101111

If we set  $(32-n) = (32-28) = 4$  right most bits to 1, we get

11010001.00010001.00100110.00101111

209. 17. 38. 47

**Solution 2:** Another way to find the last address is by ORing the given address with the complement of the mask. The complement of a number is found by changing each 1 to 0 and each 0 to 1.

Address: 11010001.00010001.00100110.00101000

Mask Complement:

00000000.00000000.00000000.00001111

Last Address:

11010001.00010001.00100110.00101111

209. 17. 38. 47

**Number of Addresses** The number of addresses in the block is the difference between the last and first address. It can be found using formula  $2^{32-n}$ . The value of  $n$  is 28, which means that number of addresses is  $2^{32-28} = 2^4 = 16$

The number of addresses can also be found by complementing the mask, interpreting it as a decimal number, and adding 1 to it.

Mask Complement:

00000000.00000000.00000000.00001111

Number of addresses:  $15 + 1 = 16$

**Example 3:** Suppose that a University-1 needs 2048 addresses and is assigned the addresses 184.26.0.0 through 184.26.7.255, along with mask 255.255.248.0. Next University-2 asks for 4096 addresses. Since a block of 4096 must lie on a 4096-byte boundary, they cannot be given addresses starting at 184.26.8.0. Instead, they get 184.26.16.0 through 184.26.31.255 along with subnet mask 255.255.240.0. University-3 asks for 1024 addresses and is assigned addresses 184.26.8.0 through 184.26.11.255 and mask 255.255.252.0. These assignments are summarized in below table:

| University   | First Address | Last Address  | Number of Addresses | Written as     |
|--------------|---------------|---------------|---------------------|----------------|
| University-1 | 184.26.0.0    | 184.26.7.255  | 2048                | 184.26.0.0/21  |
| University-3 | 184.26.8.0    | 184.26.11.255 | 1024                | 184.26.8.0/22  |
| Available    | 184.26.12.0   | 184.26.15.255 | 1024                | 184.26.12.0/22 |
| University-2 | 184.26.16.0   | 184.26.31.255 | 4096                | 184.26.16.0/20 |

What are the masks for three universities?

**Solution:**

University-1 Address 184.26.0.0

10111000.00011010.00000000.00000000

Mask: 11111111.11111111.11111000.00000000

University-2 Address 184.26.16.0

10111000.00011010.00010000.00000000

Mask: 11111111.11111111.11110000.00000000

University-3 Address 184.26.8.0

10111000.00011010.00001000.00000000

Mask: 11111111.11111111.11111100.00000000

What happens when a packet comes in addressed to 184.26.17.4?

**Solution:**

The process has address 184.26.17.4

The binary representation of 184.26.17.4 is

10111000.00011010.00010001.00000100

First it is Boolean ANDed with University-1 mask which gives

10111000.00011010.00010001.00000100

11111111.11111111.11111000.00000000

10111000.00011010.00010000.00000000

184. 26. 16. 0

This value does not match the base address of University-1, so the original address is next ANDed with University-2 mask which gives

10111000.00011010.00010001.00000100

11111111.11111111.11110000.00000000

10111000.00011010.00010000.00000000

184. 26. 16. 0

This value does not match the University-2 base address. Finally original address is ANDed with University-3

$$\begin{array}{r} 10111000.00011010.00010001.00000100 \\ \underline{11111111.11111111.11111100.00000000} \\ \underline{10111000.00011010.00010000.00000000} \\ = 184. \quad 26. \quad 16. \quad 0 \end{array}$$

The University-3 entry is used and the packet is sent along the line named in it.

## NETWORK ADDRESS

When an organization is given a block of addresses, the organization is free to allocate the addresses to the devices that need to be connect to the Internet. The first address in the class is treated as a special address. The first address is called the network address and defines the organization network. First address is the one that is used by routers to direct the message sent to the organization from the outside.

### Two-level Hierarchy Without Subnetting

An IP-address can be only 2 levels of hierarchy when not subnetted. The  $n$  left most bits of the address  $x.y.z.w/n$  define the network (organization network), the  $(32-n)$  right most bits define the particular host (computer or router) to the network. The 2 common terms are prefix and suffix. The part of the address that defines the network is called the prefix; the part that defines the host is called the suffix.

Prefix is common to all addresses in the network; the suffix changes from one device to another.

### Three-level of Hierarchy with Subnetting

An organization that is granted a large block of addresses may want to create clusters of networks called subnets and divide the address between the different subnets. The rest of the world still sees the organization as one entity. Internally there are several subnets. All messages are sent to the router address that connects the organization to the rest of the Internet; the router routes the message to the appropriate subnets.

The organization, needs to create small sub blocks of addresses, each assigned to specific subnets. The organization has its own mask; each subnet must also have its own.

**Example 4:** Suppose an organization is given the block 19.18.50.0/26, which contains 64 addresses. The organization has 3 offices and needs to divide the addresses into 3 sub blocks of 32, 16 and 16 addresses. Find the new masks for all the 3 sub blocks?

**Solution:** First subnet has to be allocated 32 addresses. Suppose the mask for the first subnet is  $n_1$ , then  $2^{32-n_1}$  must be 32, which means  $n_1 = 27$

$$(\therefore 2^{32-n_1} = 2^5)$$

Second subnet has to be allocated 16 addresses suppose the mask for the second subnet is  $n_2$ , then  $2^{32-n_2}$  must be 16,  $n_2 = 28$

$$(\therefore 2^{32-n_2} = 2^4)$$

Suppose the mask for third subnet is  $n_3$ , then  $2^{32-n_3}$  must be 16,  $n_3 = 28$

$$(\therefore 2^{32-n_3} = 2^4)$$

We have the masks 27, 28, 28 with the organization mask being 26.

## Network Address Translation (NAT)

NAT enables a user to have a large set of addresses internally and one address, or a small set of addresses, externally. To separate the addresses used inside the home or business and the ones used for the Internet, the Internet authorities have reserved three sets of addresses as private addresses.

| Range                          | Total    |
|--------------------------------|----------|
| 10.0.0.0 to 10.255.255.255     | $2^{24}$ |
| 172.16.0.0 to 172.31.255.255   | $2^{20}$ |
| 192.168.0.0 to 192.168.255.255 | $2^{16}$ |

Figure 2 Addresses for private networks

Any organization can use an address out of this set without permission from the Internet Authorities. These reserved addresses are for private networks. They are unique inside the organization, but they are not unique globally.

No Router will forward a packet that has one of these addresses as the destination address.

## IP-PROTOCOL

The job of Internet protocol is to provide a best effort to transport datagrams from source to destination, without regard to whether these machines are on the same network or other networks.

- Communication in the Internet works as follows. The transport layer takes data streams and breaks them up into datagrams.
- Datagrams can be upto 64 kbytes each, but in practice they are usually not more than 1500 bytes (so that the datagrams fit in one Ethernet frame). Each datagram is transmitted through the Internet, being fragmented into smaller units as it goes. When all the pieces finally get to the destination machine, they are reassembled by the network layer into the original datagram.
- An IP datagram consists of a header part and a text part. The header has a 20-byte fixed part and a variable length optional part.
- It is transmitted in big endian order: from left to right, with the high-order bit of the version field going first.

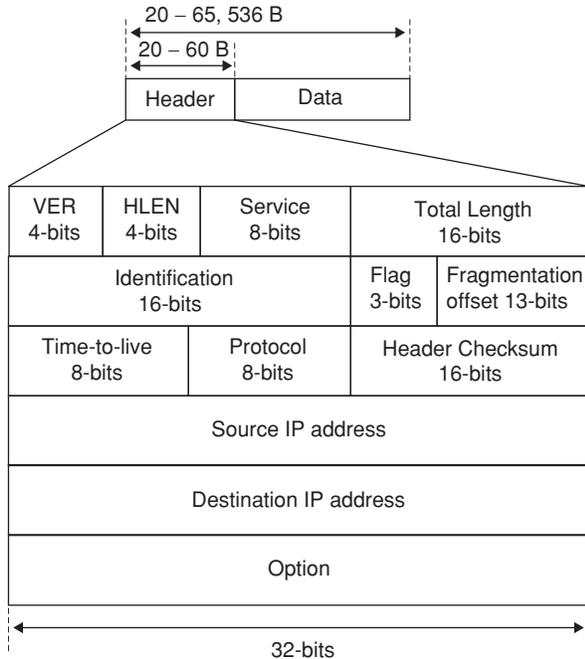


Figure 3 IPv4 datagram format

The Internet protocol version 4 (IPv4) is the delivery mechanism used by the TCP/IP protocols.

- IPv4 is an unreliable and connection less datagram protocol.
- IPv4 provides no error control or flow control (Expect for error detection on the header).
- If reliability is important, IPv4 must be paired with a reliable protocol such as TCP.
- IPv4 is connection less protocol for a packet-switching network that uses the datagram approach. This means that each datagram can follow a different route to the destination. This implies that datagrams sent by the same source to the same destination could arrive out of order. Some of them could be lost or corrupted during transmission. Again, IPv4 relies on a higher-level protocol to take care of all these problems.
- A datagram is a variable-length packet consisting of 2-parts: Header and Data.

The header is 20 to 60 bytes in length and contains information essential to routing and delivery.

The fields of IPv4 are:

### Version (VER)

This 4-bit field defines the version of the IPv4 protocol. Currently the version is 4. This field tells the IPv4 software running in the processing machine that the datagram has the format of version 4. All fields must be interpreted as specified in the fourth version of the protocol. If the machine is using some version of IPv4, the datagram is discarded rather than interpreted incorrectly.

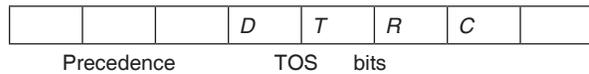
### Header Length: (HLEN)

This 4-bit defines the total length of the datagram header in 4-byte words. This field is required because the length of the header is variable (between 20 and 60 bytes). When there are no options, the header length is 20 bytes, and the value of this field is 5 ( $5 \times 4 = 20$ ). When the options field is at its maximum size, the value of this field is 15 ( $15 \times 4 = 60$ ).

### Services

This 8-bit field, previously called service type, is now called differentiated services. Both Interpretations are given below.

1. Service Type:



- D: Minimize delay
- R: Maximize reliability
- T: Maximize throughput
- C: Minimize cost

In this interpretation, the first 3-bits are called precedence bits. The next 4-bits are called Type-of-service (TOS) bits, and the last bit is not used.

**Precedence** Precedence is a 3-bit sub field ranging from 0(000 binary) to 7(111 binary) precedence is used to give priority to the datagram in issues such as congestion. If a router is congested and needs to discard some datagrams, the datagrams with lowest precedence are discarded first.

**Example:** A datagram used for network management is much more important than a datagram containing optional information for a group.

The precedence sub field was part of version 4, but never used.

**TOS bits** It is a 4-bit sub field with each bit having a special meaning. One and only one of the bits can have the value of 1 in each datagram. The bit pattern interpretation and their services are given below.

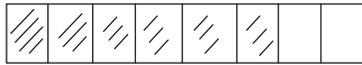
| TOS Bits | Description          |
|----------|----------------------|
| 0000     | Default              |
| 0001     | Minimize cost        |
| 0010     | Maximize reliability |
| 0100     | Maximize throughput  |
| 1000     | Minimize delay       |

Figure 4 Types of service

- Interactive activities, Activities requiring immediate attention, and activities requiring immediate response need minimum delay.
- Activities that send bulk data require maximum throughput.

- Management activities need maximum reliability
- Background activities need minimum cost.

## 2. Differentiated Services:



Code point

The first 6-bits make up the code point sub field, and the last 2-bits are not used. The code point sub field can be used in two different ways.

1. When the 3 rightmost bits are 0's, the 3 leftmost bits are interpreted the same way as the precedence bits in the service type interpretation.
2. When the 3 rightmost bits are not all 0s, the 6 bits define 64 services based on the priority assignment by the Internet Authorities.

## Total length

This is a 16-bit field that defines the total length (header plus data) of the IPv4 datagram in bytes. To find the length of the data coming from the upper layer, subtract the header length from the total length. The header length can be found by multiplying the value in the HLEN field by 4.

∴ Length of data = Total Length – Header Length.

The field length is 16-bits, the total length of the IPv4 datagram is limited to 65,535 ( $2^{16} - 1$ ) bytes, of which 20 to 60 bytes are the header and the rest is data from the upper layer.

- The total length field defines the total length of the datagram including the header.
- If the size of an IPv4 datagram is less than 46 bytes, some padding will be added to meet this requirement. In this case, when a machine decapsulates the datagram, it needs to check the total field to determine how much is really data and how much is padding.

## Fragmentation fields

The fields that are related to fragmentation and reassembly of an IPv4 datagram are

1. Identification
2. Flags
3. Fragmentation offset

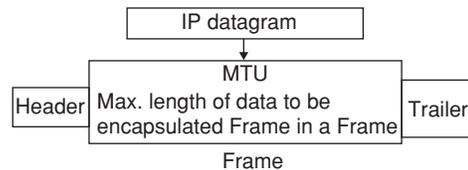
## Fragmentation

A datagram can travel through different networks. Each router decapsulates the IPV4 datagram from the frame it receives, processes it, and then encapsulates it in another frame. The size and format of the received frame depends on the protocol used by the physical network through which the frame has travelled. The format and size of the sent frame depend on the protocol used by the physical network through which the frame is going to travel.

**Example:** If a router connects a LAN to a WAN, it receives a frame in the LAN format and sends a frame in the WAN format.

## Maximum transfer unit (MTU)

When a datagram is encapsulated in a frame, the total size of the datagram must be less than this maximum size. The value of the MTU depends on the physical network protocol.



Following table shows MTU's for some networks.

| Protocol           | MTU    |
|--------------------|--------|
| Hyper Channel      | 65,535 |
| Token Ring(16Mbps) | 17,914 |
| Token Ring(4Mbps)  | 4,464  |
| FDDI               | 4,352  |
| Ethernet           | 1,500  |
| X.25               | 576    |
| PPP                | 296    |

To make the IPv4 protocol independent of the physical network, the designers decided to make the maximum length of the IPV4 datagram equal to 65,535 bytes.

This makes transmission more efficient if we use a protocol with an MTU of this size. But, for other physical network, we must divide the datagram to make it possible to pass through these networks. This is called Fragmentation.

- When a datagram is fragmented, each fragment has its own header with most of the fields repeated.
- A fragmented datagram may itself be fragmented if it encounters a network with an even smaller MTU.
- A datagram can be fragmented several times before it reaches the final destination.
- The host or router that fragments a datagram must change the values of three fields.
  - Flags
  - Fragmentation offset
  - Total length

The rest of the fields must be copied. The value of the checksum must be recalculated regardless of fragmentation.

**Identification** This 16-bit field identifies a datagram originating from the source host. The combination of the identification and source IPv4 address must uniquely define a datagram as it leaves the source host.

To guarantee uniqueness, the IPv4 protocol uses a counter to label the datagrams.

The counter is initialized to a positive number. When the IPv4 protocol sends a datagram, it copies the current value of the counter to the identification field and increments the counter by 1.

As long as the counter is kept in the main memory uniqueness is guaranteed.

When a datagram is fragmented, the value in the identification field is copied to all fragments.

All fragments have the same identification number, the same as the original datagram.

The identification number helps in reassembling the datagram at destination side.

**Flags** This is a 3-bit field, the first bit is reserved.

The second bit is called the ‘do not fragment’ bit. If its value is 1, the machine must not fragment the datagram.

If it cannot pass the datagram through any available physical network, it discards the datagram and sends an ICMP error message to the source host.

If its value is ‘0’, the datagram can be fragmented if necessary.

The 3rd bit is called the more fragment bit. If its value is 1, it means the datagram is not the last fragment; there are more fragments after this one. If its value is 0, it means this is the last or only fragment.



*D*: Do Not Fragment

*M*: More Fragments

**Fragmentation offset** This is a 13-bit field, shows the relative position of this fragment with respect to the whole datagram.

**Example:** A datagram with a data size of 4000 bytes fragmented into three fragments. The bytes in the original datagram are numbered 0 to 3999.

1. The first fragment carries bytes 0 to 1399. The offset for this datagram is  $0/8 = 0$ .
2. The second fragment carries bytes 1400 to 2799; the offset values for this fragment is  $1400/8 = 175$ .
3. The third fragment carries bytes 2800 to 3999. The offset value for this fragment is  $2800/8 = 350$ 
  - All fragment datagrams follow different paths to reach destination.
  - Even though each fragment follows a different path and arrives out of order, the final destination host can reassemble the original datagram from the fragments received by using the following strategy:
    - The first fragment has an offset field Value of zero
    - Divide the length of the first fragment by 8. The second fragment has an offset value equal to that result.
    - Divide the total length of the first and second fragments by 8. The third fragment has an offset value equal to that result.
    - The last fragment has a more bit value of 0.

### Time-to-live

This field was originally designed to hold a time stamp, which was decremented by each visited router. The datagram was discarded when the value became zero.

This field is used mostly to control the maximum number of hops (routers) visited by the datagram, When a source host sends the datagram, it stores a number in this field. This value is approximately 2 times the maximum number of routes between any 2 hosts.

Each router that processes the datagram decrements this number by 1. If this value becomes zero, the router discards the datagram.

This field limits the lifetime of a datagram and avoids loops (A datagram may travel between 2 or more routers for a long time without ever getting delivered to the destination host)

### Protocol

This 8-bit field defines the higher-level protocol that uses the services of the IPv4 layer.

- An IPv4 datagram can encapsulate data from several higher-level protocols such as TCP, UDP, ICMP and IGMP.
- The field specifies the final destination protocol to which the IPv4 datagram is delivered.

**Source address** This 32-bit field defines the IPv4 address of the source. This field must remain unchanged during the time the IPv4 datagram travels from the source host to the destination host.

### Destination address

This 32-bit field defines the IPv4 address of the destination. This field must remain unchanged during the time the IPv4 datagram travels from the source host to the destination host.

**Example 5:** An IPv4 packet has arrived with the first 8-bits as shown,

01000100

Will the packet be discarded?

**Solution:** There is an error in this packet.

The 4 left most bits (0100) show the version, which is correct. The next 4 bits (0100) show an invalid header length ( $4 \times 4 = 16$ )

The minimum number of bytes in the header must be 20. The receiver discards the packet.

**Example 6:** In an IPv4 packet, the value of HLEN is 1100 in binary. How many bytes of options are being carried by this packet?

**Solution:** The HLEN value is 12, which means the total number of bytes in the header is  $12 \times 4 = 48$  bytes. The first 20 bytes are the base header, the next 28 bytes are the options.

**Example 7:** In an IPv4 packet, the value of HLEN is 5 and the value of the total length field is  $0 \times 0038$ . How many bytes of data are being carried by this packet?

**Solution:** Then HLEN value is 5, which means the total number of bytes in the header is  $5 \times 4 = 20$  bytes (no options). The total length is  $(16^1 \times 3 + 16^0 \times 8) = 48 + 8 = 56$ .

Which means the packet is carrying 36 bytes of data.  
 $\therefore (56 - 20) = 36$

**Option:** The header of the IPv4 datagram is made of 2 parts:

1. Fixed part
2. Variable part

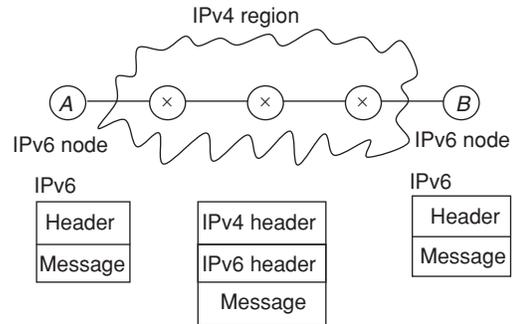
The fixed part is 20 bytes long. The variable part comprises the options that can be a maximum of 40 bytes.

Options, as the name implies, are not required for a datagram. They can be used for network testing and debugging.

## Tunneling

If two computers are using IPv6, and want to communicate across the region using IPv4, then Tunneling concept is used.

To pass a packet using IPv6 across a region using IPv4, packet should have IPv4 address. So, IPv6 packet is encapsulated in IPv4 packet, once the packet leaves region it leaves its capsule.



Tunneling cannot be used for translation of IPv6 address to IPv4 address.

## Types of tunneling

There are two types of tunneling methods (Corresponding to VPN]

1. End to End tunneling
2. Node to Node tunneling

End to end tunneling acts an interface between LAN and Internet. It is used in Remote access VPN connection.

Node to node tunneling acts as an interface between nodes which are present at an edge of a private network. It is mostly used in the site to site VPN connection.

The other types of tunneling are

1. Layer 2 tunneling (Data link layer)
3. Layer 3 tunneling (Network layer)

Based on the tunneling protocol used for data encapsulation we have different tunneling methods.

Layer 2 tunneling protocol uses frames for encapsulating message it is mostly used in point to point tunnels between client and VPN server.

Layer 3 tunneling protocol adds a new IP header to the packet, it is mostly used for connecting two or more private networks.

## EXERCISES

### Practice Problems I

**Directions for questions 1 to 15:** Select the correct alternative from the given choices.

1. Consider the given IP address, 156.216.24.65 with a subnet mask of 7-bits, what are the number of hosts and subnets?  
 (A) 512, 128 (B) 510, 126  
 (C) 511, 127 (D) 509, 125
2. IP address is 198.250.144.23 and mask is 255.255.255.240, find the class, network mask length and broadcast address?  
 (A) 128,198.250.144.31  
 (B) 26,198.250.144.63  
 (C) 30,198.250.144.3  
 (D) 32,198.250.144.0

3. If subnet addresses are 129.253.4.0, 129.253.8.0, 129.253.12.0, 129.253.16.0  
 What is subnet mask?  
 (A) 129.253.7.0 (B) 129.253.31.0  
 (C) 129.253.192.0 (D) 129.253.252.0

4. For a given source IP 192.16.9.10 and destination network 10.0.0.0, match the following:

|                                 |                   |
|---------------------------------|-------------------|
| i Network address               | P 127.0.0.5       |
| ii Directed broadcast address   | Q 0.0.0.5         |
| iii Limited broadcast address   | R 0.0.0.0         |
| iv This host on this network    | S 255.255.255.255 |
| v Specific host on this network | T 10.255.255.255  |
| vi Loop back address            | U 192.16.9.0      |

8.62 | Unit 8 • Networks, Information Systems, Software Engineering and Web Technology

Find out all related addresses given above?

- (A) i – U, ii – R, iii – Q, iv – S, v – T, vi – P  
 (B) i – R, ii – Q, iii – P, iv – T, v – S, vi – U

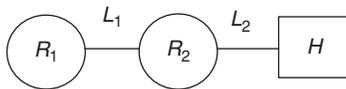
- (C) i – U, ii – T, iii – S, iv – R, v – Q, vi – P  
 (D) i – U, ii – T, iii – S, iv – Q, v – R, vi – P

**Common data questions 5 and 6:** Consider the routing table below:

| Destination    | Gateway         | Mask            | Flags | Interface   |
|----------------|-----------------|-----------------|-------|-------------|
| 165.230.198.64 | 165.230.198.119 | 255.255.255.192 | U     | eth0        |
| 192.168.1.0    | 192.168.1.1     | 255.255.255.0   | U     | eth1        |
| 127.0.0.0      | 127.0.0.1       | 255.0.0.0       | U     | Loop back 0 |
| Default        | 165.230.198.65  | 255.255.255.255 | UG    | eth0        |

5. How many local subnets is this machine attached to?  
 (A) 0 (B) 1  
 (C) 2 (D) 3
6. How many IP addresses can this machine reach to (excluding the loop back route)?  
 (A) 64 (B) 256  
 (C) 320 (D) 192
7. If a packet is of size 1000 bytes and network span 15 km and speed of propagation is 70% of speed of light. What is the throughput of the system?  
 (A) 58 Mbps (B) 60 Mbps  
 (C) 56 Mbps (D) 54 Mbps

**Common data for questions 8 and 9:** The diagram shows router  $R_1$  sending a datagram to host  $H$  through Router  $R_2$ .



Link  $L_1$  permits a maximum transfer unit of 1500 bytes. Link  $L_2$  only permits a maximum transfer unit of 1100 bytes.  $A$  is an IP datagram which

- (i) Has size 4000 bytes (The size of the datagram includes the header of 20 bytes).  
 (ii) Is not using any of the option field in the header.  $A$  must be fragmented as it is sent from  $R_1$  to  $H$ . Assume that all datagrams are received successfully.
8. The sizes of IP datagram  $A$  is fragmented in sending it from  $R_1$  to  $R_2$  over  $L_1$  are  
 (A) 1500, 1500, 1500  
 (B) 1500, 1500, 1060  
 (C) 1500, 1500, 1100  
 (D) 1500, 1500, 1040
9. The number of IP datagrams received by  $H$  are  
 (A) 6 (B) 5  
 (C) 4 (D) 3

**Common data for questions 10 and 11:** In a network, system packet size is 2 kB, propagation time is 30 milliseconds and channel capacity is  $10^6$  bits/sec.

10. What will be the transmission time?  
 (A) 21 microsecond (B) 16.3 microsecond  
 (C) 18.3 millisecond (D) 16.3 millisecond
11. What is the channel utilization of sender?  
 (A) 21% (B) 12%  
 (C) 16% (D) 30%
12. In an IPv4 header fragment offset is set to a size of 13 bits. If the maximum size of datagram is 64 kB, what is the maximum number of fragments possible?  
 (A) 8191 (B) 8192  
 (C) 13 (D) 12
13. A packet arriving at main router is addressed to 95.80.15.6. The subnet mask used is 255.255.252.0/22. What is the resultant address?  
 (A) 95.80.255.0 (B) 95.80.24.0  
 (C) 95.80.12.0 (D) 95.80.6.6
14. What does 'record route' option signify in an IP header?  
 (A) Routes that processed the packet, stores the packet details in local memory.  
 (B) Follow the path already available without further use of any routing algorithm.  
 (C) Make each router append its IP address to the packet in transit.  
 (D) Make routers inform the source about the path taken.
15. Which of the following address is used or reserved for loop back testing?  
 (A) 0. 0. 0. 0  
 (B) 1. 1. 1. 1  
 (C) 127.xx.yy.zz  
 (D) None of these

## Practice Problems 2

**Directions for questions 1 to 15:** Select the correct alternative from the given choices.

- The internet layer of TCP/IP model is similar to \_\_\_\_\_ layer of OSI model  
(A) Transport layer (B) Session layer  
(C) Presentation layer (D) Network layer
- Consider the following statements:  
 $S_1$ : A system can have multiple IP addresses.  
 $S_2$ : A system can have multiple physical addresses.  
Which one of the following is correct?  
(A) Both  $S_1$  and  $S_2$  are true  
(B) Both  $S_1$  and  $S_2$  are false  
(C)  $S_1$  is true,  $S_2$  is false  
(D)  $S_2$  is true,  $S_1$  is false
- Time to live (TTL) Field in IP header is used  
(A) To avoid infinite loops.  
(B) To Fragment the packets in the subnet.  
(C) To calculate the shortest path from source and destination.  
(D) None of these
- Match the following IP header fields to their functionalities.

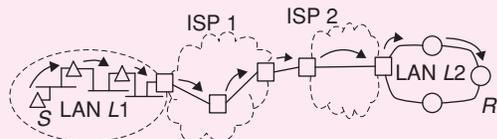
|                        |                                                                     |
|------------------------|---------------------------------------------------------------------|
| i DF                   | p Indicates if there are more fragments                             |
| ii MF                  | q Indicates the transport process to which packets need to be given |
| iii Protocol           | r If this bit is set routers are not supposed to fragment           |
| iv Header Checksum     | s Needed to relate all the fragments of a datagram                  |
| v Identification value | t Can vary from hop to hop                                          |

- (A) i – q, ii – p, iii – t, iv – r, v – s  
(B) i – p, ii – q, iii – s, iv – t, v – r  
(C) i – s, ii – p, iii – r, iv – q, v – t  
(D) i – r, ii – p, iii – q, iv – t, v – s
- The header part of IP contains \_\_\_\_\_ bytes fixed part.  
(A) 20 (B) 24  
(C) 16 (D) 60
  - The header checksum in the IP header is used to verify  
(A) Only header  
(B) Only data  
(C) Both (A) and (B)  
(D) None of these

- The source address length in IPv4 is  
(A) 8 bytes (B) 16 bytes  
(C) 32-bits (D) 16-bits
- The highest IP address in digital notation is  
(A) 255 . 0 . 0 . 0  
(B) 255 . 255 . 0 . 0  
(C) 255 . 255 . 255 . 0  
(D) 255 . 255 . 255 . 255
- Which type address class is used for multicast address?  
(A) Class A (B) Class B  
(C) Class C (D) Class D
- In the leaky bucket algorithm, the leaky bucket means \_\_\_\_\_  
(A) Infinite buffer.  
(B) Finite internal queue.  
(C) Constant service time.  
(D) Both (B) and (C)
- The mechanism of leaky bucket algorithm  
(A) Reduces congestion.  
(B) Turns uneven flow of packet into even flow.  
(C) Smoothens out bursts.  
(D) All the above.
- In the IP protocol header we have two one bit flags DF and MF.  
What are the uses of DF bit flag?  
(A) DF (means don't fragment) orders router not to fragment the packet.  
(B) DF is set when the destination is incapable of putting the pieces back together again.  
(C) Both (A) and (B)  
(D) None of the above
- In the IPv4 header, what is the maximum value of total length field?  
(A) 60 bytes  
(B) 255 bytes  
(C) 576 bytes  
(D) 65535 bytes
- What are the characteristics of a flow specification input?  
(A) Maximum packet size (bytes).  
(B) Token bucket rate (bytes/sec).  
(C) Token bucket size (bytes)  
(D) All the above
- Standard protocols like HTTP, SMTP, NNTP are part of  
(A) Presentation layer  
(B) Application layer  
(C) Session layer  
(D) Not part of any layer

## PREVIOUS YEARS' QUESTIONS

- Suppose computers A and B have IP addresses 10.105.1.113 and 10.105.1.91 respectively and they both use the same net mask  $N$ . Which of the values of  $N$  given below should not be used if A and B should belong to the same network? [2010]  
(A) 255.255.255.0 (B) 255.255.255.128  
(C) 255.255.255.192 (D) 255.255.255.224
- In the IPv4 addressing format, the number of networks allowed under class C addresses is [2012]  
(A)  $2^{14}$  (B)  $2^7$   
(C)  $2^{21}$  (D)  $2^{24}$
- In an IPv4 datagram, the  $M$ -bit is 0, the value of HLEN is 10, the value of total length is 400 and the fragment offset value is 300. The position of the datagram, the sequence numbers of the first and the last bytes of the payload, respectively are [2013]  
(A) Last fragment, 2400 and 2789  
(B) First fragment, 2400 and 2759  
(C) Last fragment, 2400 and 2759  
(D) Middle fragment, 300 and 689
- In the diagram shown below,  $L1$  is an Ethernet LAN and  $L2$  is a Token-Ring LAN. An IP packet originates from sender  $S$  and traverses to  $R$ , as shown. The links within each ISP and across the two ISPs, are all point-to-point optical links. The initial value of the TTL field is 32. The maximum possible value of the TTL field when  $R$  receives the datagram is \_\_\_\_\_. [2014]



- Host  $A$  (on TCP/IPv4 network  $A$ ) sends an IP datagram  $D$  to host  $B$  (also on TCP/IPv4 network  $B$ ). Assume that no error occurred during the transmission of  $D$ . When  $D$  reaches  $B$ , which of the following IP header field(s) may be different from that of the original datagram  $D$ ? [2014]  
(i) TTL (ii) Checksum  
(iii) Fragment Offset  
(A) (i) only (B) (i) and (ii) only  
(C) (ii) and (iii) only (D) (i), (ii), and (iii)
- An IP router implementing Classless Inter-Domain Routing (CIDR) receives a packet with address 131.23.151.76. The router's routing table has the following entries:

| Prefix        | Out Interface Identifier |
|---------------|--------------------------|
| 131.16.0.0/12 | 3                        |
| 131.28.0.0/14 | 5                        |
| 131.19.0.0/16 | 2                        |
| 131.22.0.0/15 | 1                        |

The identifier of the output interface on which this packet will be forwarded is \_\_\_\_\_. [2014]

- Every host in an IPv4 network has a 1-second resolution real-time clock with battery backup, each host needs to generate up to 1000 unique identifiers per second. Assume that each host has a globally unique IPv4 address. Design a 50-bit globally unique ID for this purpose. After what period (in seconds) will the identifiers generated by a host wrap around? [2014]
- Which one of the following fields of an IP header is NOT modified by a typical IP router? [2015]  
(A) Checksum (B) Source address  
(C) Time to Live (TTL) (D) Length
- Consider the following routing table at an IP router:

| Network No.  | Net Mask      | Next Hop    |
|--------------|---------------|-------------|
| 128.96.170.0 | 255.255.254.0 | Interface 0 |
| 128.96.168.0 | 255.255.254.0 | Interface 1 |
| 128.96.166.0 | 255.255.254.0 | R2          |
| 128.96.164.0 | 255.255.252.0 | R3          |
| 0.0.0.0      | Default       | R4          |

For each IP address in Group I identify the correct choice of the next hop from Group II using the entries from the routing table above. [2015]

| Group I             | Group II       |
|---------------------|----------------|
| i) 128.96.171.92    | a) Interface 0 |
| ii) 128.96.167.151  | b) Interface 1 |
| iii) 128.96.163.151 | c) R2          |
| iv) 128.96.165.121  | d) R3          |
|                     | e) R4          |

- (A) i-a, ii-c, iii-e, iv-d (B) i-a, ii-d, iii-b, iv-e  
(C) i-b, ii-c, iii-d, iv-e (D) i-b, ii-c, iii-e, iv-d
- Host  $A$  sends a UDP datagram containing 8880 bytes of user data to host  $B$  over an Ethernet LAN. Ethernet frames may carry data upto 1500 bytes (i.e., MTU = 1500 bytes). Size of UDP header is 8 bytes and size of IP header is 20 bytes. There is no option field in IP header. How many total number of IP fragments will be transmitted and what will be the contents of offset field in the last fragment?  
(A) 6 and 925 (B) 6 and 7400  
(C) 7 and 1110 (D) 7 and 8880
- In the network 200.10.11.144/27, the fourth octet (in decimal) of the last IP address of the network which can be assigned to a host is \_\_\_\_\_. [2015]
- An IP datagram of size 1000 bytes arrives at a router. The router has to forward this packet on a link whose MTU (maximum transmission unit) is 100 bytes.

Assume that the size of the IP header is 20 bytes.

The number of fragments that the IP datagram will be divided into for transmission is \_\_\_\_\_. **[2016]**

13. For the IEEE 802.11 MAC protocol for wireless communication, which of the following statements is/ are **TRUE**? **[2016]**
- I. At least three non-overlapping channels are available for transmissions.
  - II. The RTS-CTS mechanism is used for collision detection.
  - III. Unicast frames are ACKed.
- (A) All I, II and III                      (B) I and III only  
(C) II and III only                        (D) II only
14. The maximum number of IPv4 router addresses that can be listed in the record route (RR) option field of an IPv4 header is \_\_\_\_\_. **[2017]**
15. Match the following:

|    | Field                    |     | Length in bits |
|----|--------------------------|-----|----------------|
| P. | UDP Header's Port Number | I.  | 48             |
| Q. | Ethernet MAC Address     | II. | 8              |

|    | Field                        |      | Length in bits |
|----|------------------------------|------|----------------|
| R. | IPv6 Next Header             | III. | 32             |
| S. | TCP Header's Sequence Number | IV.  | 16             |

**[2018]**

- (A) P-III, Q-IV, R-II, S-I  
(B) P-II, Q-I, R-IV, S-III  
(C) P-IV, Q-I, R-II, S-III  
(D) P-IV, Q-I, R-III, S-II
16. Consider an IP packet with a length of 4,500 bytes that includes a 20-byte IPv4 header and a 40-byte TCP header. The packet is forwarded to an IPv4 router that supports a Maximum Transmission Unit (MTU) of 600 bytes. Assume that the length of the IP header in all the outgoing fragments of this packet is 20 bytes. Assume that the fragmentation offset value stored in the first fragment is 0.
- The fragmentation offset value stored in the third fragment is \_\_\_\_\_. **[2018]**

## ANSWER KEYS

### EXERCISES

#### Practice Problems 1

1. B      2. A      3. D      4. A      5. C      6. C      7. C      8. D      9. B      10. D  
11. A      12. B      13. C      14. C      15. C

#### Practice Problems 2

1. D      2. C      3. A      4. D      5. A      6. A      7. C      8. D      9. D      10. B  
11. D      12. C      13. D      14. D      15. B

#### Previous Years' Questions

1. D      2. C      3. C      4. 26      5. D      6. 1      7. 256      8. B      9. A      10. C  
11. 158      12. 13      13. B      14. 9      15. C      16. 144

# Chapter 5

## Network Security

### LEARNING OBJECTIVES

- ▣ Network security basics
- ▣ Terminologies
- ▣ Cryptographic techniques
- ▣ Encryptions
- ▣ Types of keys
- ▣ Traditional cipher algorithms
- ▣ Substitution cipher
- ▣ Traditional cipher
- ▣ Symmetric key encryption
- ▣ Asymmetric key encryption
- ▣ Diffie-hellman
- ▣ Digital signatures and certificates

### NETWORK SECURITY BASICS

It is necessary to define some fundamental terms relating to network security and are the elements used to measure the security of a network. These terms are used to measure the security of a network. To be considered sufficiently advanced along the spectrum of security, a system must adequately address identification, integrity, accountability, non-repudiation, authentication, availability, confidentiality each of which is defined in the following sections:

#### Identification

Identification is simply the process of identifying one's self to another entity or determining the identity of the individual or entity, with whom you are communicating.

#### Authentication

Authentication serves as proof that you are who you say you are or what you claim to be. Authentication is critical if there is to be any trust between parties. Authentication is required when communicating over a network or logging into a network. When communicating over a network you should ask yourself two questions.

1. With whom am I communicating?
2. Why do I believe this person or entity is who he claims to be?

#### Access Control (Authorization)

This refers to the ability to control the level of access that individuals or entities have to a network or system and how much information they can receive. Level of authorization basically determines what you're allowed to do once you are authenticated and allowed access to a network, system or some other resource such as data

or information. Access control is the determination of the level of authorization to a system, network or information.

#### Availability

This refers to whether the network, system, hardware and software are reliable and can recover quickly and completely in the event of an interruption in service. Ideally, these elements should not be susceptible to denial of service attacks.

#### Confidentiality

This is also be called privacy or secrecy to the protection of information from unauthorized disclosure. Usually achieved either by restricting access to the information or by encrypting the information so that it is not meaningful to unauthorized individuals or entities.

#### Integrity

This can be thought of as accuracy, this refers to the ability to protect information, data, or transmissions from unauthorized, uncontrolled, or accidental alterations.

#### Accountability

This refers to the ability to track or audit what an individual or entity is doing on a network or system.

#### Non-repudiation

The ability to prevent individuals or entities from denying (repudiating) that information, data or files were sent or received or that information or files were accessed or altered, when in fact they were. This capability is crucial in e-commerce, without if an individual or

entity can deny that he, she or it is responsible for a transaction and that he, she or it is, therefore, not financially liable.

### Threats

A threat is anything that can disrupt the operation, functioning, integrity, or availability of a network or system. This can take any form and can be malevolent, accidental, or simply an act of nature.

### Vulnerabilities

A vulnerability is an inherent weakness in the design, configuration, implementation, or management of a network or system that renders it susceptible to a threat. Vulnerabilities are what make networks susceptible to information loss and downtime. Every network and system has some kind of vulnerability.

### Attacks

An attack is a specific technique used to exploit a vulnerability. For example, a threat could be a denial of service. A vulnerability is in the design of the operating system, and an attack could be a ‘Ping of death’. There are two general categories of attacks:

1. Passive
2. Active

**Passive attacks** These are very difficult to detect because there is no overt activity that can be monitored or detected.

Examples of passive attacks would be packet sniffing or traffic analysis.

These types of attacks are designed to monitor and record traffic on the network. They are usually employed for gathering information that can be used later in active attacks.

**Active attacks** These employ more overt actions on the network or system. As a result, they can be easier to detect, but at the same time they can be much more devastating to a network.

Examples of this type of attack would be a denial-of-service attack or active probing of systems and networks.

### Viruses

A virus, a parasitic program that cannot function independently, is a program or code fragment that is self propagating. It is called a virus, because like its biological counterpart, it requires a ‘host’ to function. In the case of a computer virus the host is some other program to which the virus attaches itself. A virus is usually spread by executing an infected program or by sending an infected file to someone else, usually in the form of an e-mail attachment.

### Worm

A worm is a self-contained and independent program that is usually designed to propagate or spawn itself on infected

systems and to seek other systems via available networks. The difference between a virus and a Worm is that a virus is not an independent program.

### Trojan horses

A trojan horse is a program or code fragment that hides inside a program and performs a disguised function. A trojan horse program hides within another program or disguises itself as a legitimate program. This can be accomplished by modifying the existing program or by simply replacing the existing program with a new one. The Trojan horse program functions much the same way as the legitimate program, but usually it also performs some other function, such a recording sensitive information or providing a trap door. An example would be a ‘password grabber’.

### Logic bombs

A logic bomb is a program or subsection of a program designed with malevolent intent. It is referred to as a logic bomb, because the program is triggered when certain logical conditions are met. This type of attack is almost always perpetrated by an insider with privileged access to the network. The perpetrator could be a programmer or a vendor that supplies software.

### Denial of service (DOS)

Denial of service attacks are designed to shut down or render inoperable a system or network. The goal of the denial-of-service attack is not to gain access or information but to make a network or system unavailable for use by other users. It is called denial-of-service attack, because the end result is to deny legitimate users access to network services.

### Protection against network threats

Network threats may cause a massive harm to the system, as the network users are increasing, there is a good chance to attack a system protection against threats should be done.

To protect system form virus and worms, a security suite should be installed.

Similarly, to protect a system from Trojan horse, internet security suite prevents from downloading Trojan horse.

SPAM filters should be used to stop SPAM, this is available within the mail servers by default.

A strong encryption should be used to protect against packet sniffers, so that packets become unreadable making packet sniffers useless.

## CRYPTOGRAPHIC TECHNIQUES

For the exchange of information and commerce to be secure on any network, a system or process must be put in place that satisfies requirements for confidentiality, access control, authentication, integrity, and non-repudiation. The key

to the securing information on a network is cryptography. Cryptography can be used as a tool to provide privacy.

Traditionally, cryptography conjures up thoughts of spies and secret codes. In reality, cryptography and encryption have found broad applications in society. Every time you use an ATM machine to get cash or a point-of-sale machine to make a purchase, you are using encryption.

### Encryption

Encryption is the process of scrambling the contents of a file or message to make it unintelligible to anyone not in possession of the ‘key’ required to unscramble it.

|   |   |   |   |   |   |   |   |   |   |   |   |   |   |   |   |   |   |   |   |   |   |   |   |   |   |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| A | B | C | D | E | F | G | H | I | J | K | L | M | N | O | P | Q | R | S | T | U | V | W | X | Y | Z |
| 1 | 2 | 3 | 4 | 5 | 6 | A | B | C | D | E | F | G | H | I | J | K | L | M | N | O | P | Q | R | S | T |

To illustrate how this works see the following where the cipher is used to scramble the message:

‘Little green apples’

Cipher text: FCNNF5 AL55H 1JF5M

Clear text: LITTLE GREEN APPLES

This cipher would not be effective at keeping a message secret for long. It does not comply with one of the qualities of a truly effective cipher. Ciphers usually fall into one to two categories:

1. Block Ciphers
2. Stream Cipher

### Stream ciphers

Stream cipher algorithms process plaintext to produce a stream of cipher text. The cipher inputs the plaintext in a stream and outputs a stream of cipher text.

#### Example:

Plaintext: LET US TALK ONE TO ONE

Cipher text: F5N OM NLFE ITS NI ITS

Stream cipher have several weaknesses. The most crucial short coming of stream ciphers is the fact that patterns in the plain text can be reflected in the cipher text. Knowing that certain words repeat makes breaking the code easier. In addition, certain words in the English language appear with predictable regularity. Letters of the alphabet also appear in predictable regularity. The most commonly used letters of the alphabet in the English language are E, T, A, O, N and I. The least commonly used letters are J, K, X, Q and Z. The most common combination of letters in the English language is ‘th’. As a result, if a code breaker is able to find a ‘t’ in a code, it doesn’t take long to find an ‘h’.

### Block ciphers

Block ciphers differ from stream ciphers in that they encrypt and decrypt information in fixed size blocks rather than

A cryptosystem or algorithm is the process or procedure to turn plain text into crypto text. A crypto algorithm is also known as a ‘cipher’. Theoretically, all algorithms can be broken by one method or another. However, an algorithm should not contain an inherent weakness that an attacker can easily exploit

**Example:** Below is an example of a cipher, to scramble a message with this cipher, simply match each letter in a message to the first row and convert it into the number or letter in the second row. To unscramble a message, match each letter or number in a message to the corresponding number or letter in the second row and convert it into the letter in the first row.

encrypting and decrypting each letter or word individually. A block cipher passes a block of data or plaintext through its algorithm to generate a block of cipher text. Another requirement of block cipher is that the cipher texts should contain no detectable pattern.

### Types of keys

We deal with three types of keys in cryptography:

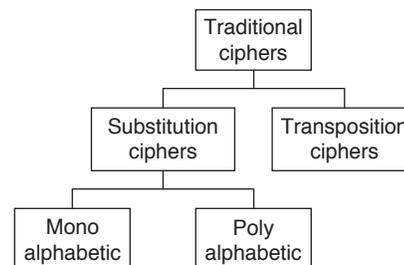
1. Secret key
2. Public key
3. Private Key

- The secret key, is the shared key used in symmetric-key cryptography.
- Public and Private keys are used in asymmetric-key cryptography.
- In symmetric-key cryptography, the same key locks and unlocks the box.
- In asymmetric-key cryptography, one key locks the box, but another key is needed to unlock it.

## TRADITIONAL CIPHER ALGORITHMS

Traditional ciphers are character oriented, these ciphers can be divided into two broad categories:

1. Substitution ciphers
2. Transposition ciphers.



### Substitution Cipher

A substitution cipher substitutes one symbol with another. If the symbols in the plain text are alphabetic characters, we replace one character with another. Substitution ciphers can be categorized as either mono-alphabetic or poly-alphabetic ciphers.

- In a mono-alphabetic cipher, a character or symbol in the plaintext is always changed to the same character or symbol in the cipher text regardless of its position in the text. For example if the algorithm says that character 'A' in the plain text is changed to character 'E', every character 'A' is changed to character 'E'.
- The relationship between characters in the plain text and the cipher text is a one-to-one relationship.
- In a poly-alphabetic cipher, each occurrence of a character can have a different substitute. The relationship between a character in the plain text to a character in the cipher text is a one-to-many relationship.
- To achieve this goal, we need to divide the text into groups of characters and use a set of keys.
- In substitution cipher, if 'a' becomes D, 'b' becomes 'E' then the word 'corrupt' becomes ETUUXSW, plain text will be given in lower case, and cipher text in upper case.
- A slight generalization of the ceasar cipher allows the cipher text alphabet to be shifted by 'K' letters, instead of always '3'.
- The next improvement is to have each of the symbols in the plain text, say, the 26 letter for simplicity, map onto some other letter.

**Example:**

|             |   |   |   |   |   |   |   |   |   |   |
|-------------|---|---|---|---|---|---|---|---|---|---|
| Plain Text  | a | b | c | d | e | f | g | h | i | j |
| Cipher Text | L | N | O | B | R | M | S | U | V | Z |

|             |   |   |   |   |   |   |   |   |   |   |   |
|-------------|---|---|---|---|---|---|---|---|---|---|---|
| Plain text  | k | l | m | n | o | p | q | r | s | t | u |
| Cipher Text | P | A | K | C | L | H | W | Q | X | Y | J |

|             |   |   |   |   |   |
|-------------|---|---|---|---|---|
| Plain Text  | v | w | x | y | z |
| Cipher Text | E | F | D | G | J |

|             |         |
|-------------|---------|
| Plain Text  | corrupt |
| Cipher Text | OIQQJHY |

- In this method, if a small cipher is given it can be broken easily. The basic attack takes advantage of the statistical properties of natural languages. For example, In English, 'e' is the most common letter followed by *t, o, a, n, i* etc.
- The most common 2 letter combinations, are *th, in, er, re* and *an*.
- The most common three-letter combinations are are, the, ing, and, and ion.
- By making guesses at common letters, digrams and trigrams and knowing about likely patterns of vowels and consonants, the cryptanalyst builds up a tentative plain-text, letter by letter.

### Transposition Ciphers

Substitution ciphers preserve the order of the plaintext symbols but disguise them.

Transposition ciphers, in contrast, reorder the letters but do not disguise them. Following figure depicts a common transposition cipher, the columnar transposition.

- The cipher is keyed by a word or phrase not containing any repeated letters.

**Example:** 'NETWORKS' is the key.

Plaintext: Transfer ten million dollars to my account.  
What is the cipher text using transposition cipher?

**Solution:** Key: NETWORKS

|   |   |   |   |   |   |   |   |
|---|---|---|---|---|---|---|---|
| N | E | T | W | O | R | K | S |
| 3 | 1 | 7 | 8 | 4 | 5 | 2 | 6 |
| T | r | a | n | s | f | e | r |
| t | e | n | m | i | l | l | i |
| o | n | d | o | l | l | a | r |
| s | t | o | m | y | a | c | c |
| o | u | n | t | a | b | c | d |

The purpose of key is to number the columns, column 1 being under the key letter closest to the start of the alphabet, and so on.

The plain text is written horizontally in rows, padding is required to fill the matrix, if it is not complete'. The cipher text is read out by columns, starting with the column whose key letter is the lowest.

Plain text: Transfer ten million dollars to my account  
Cipher Text: rentue laccttososilyfllabircdandonnmomt.

### SYMMETRIC KEY ENCRYPTION

Symmetric key, also referred to as private key or secret key, is based on a single key and algorithm being shared between the parties who are exchanging encrypted information. The same key both encrypts and decrypts messages.

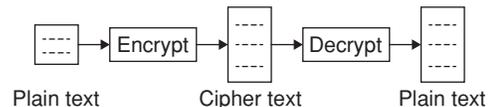


Figure 1 Symmetric key encryption

The strength of the scheme is largely dependent on the size of the key and on keeping it secret. Generally the larger the key, the more secure the scheme. In addition, symmetric key encryption is relatively fast. Private key cryptosystems are not well suited for spontaneous communication over open and unsecured networks. Symmetric key provides on process for authentication or non-repudiation.

### Data Encryption Standard: (DES)

DES consists of an algorithm and a key. The key is a sequence of eight bytes, each containing eight bits for a 64 bit key. Since each byte contains one parity bit, the key is actually 56 bits in length. DES is widely used in automated teller machine (ATM) and point-of-sale (POS) networks, so if you use an ATM or debit card you are using DES.

### ASYMMETRIC KEY ENCRYPTION

Asymmetric cryptography is also known as public key cryptography, public key cryptography uses two keys one is public key and the other is private key. The key names describe their function. One key is kept private, and the other key is made public. Knowing the public key doesn't reveal the private key. A message encrypted by the private key can only be decrypted by the corresponding public key. Conversely, a message encrypted by the public key can only be decrypted by the private key.

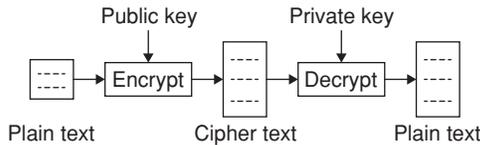


Figure 2 Asymmetric key encryption

With the aid of public key cryptography, it is possible to establish secure communications with any individual or entity when using a compatible software or hardware device.

There are three public key algorithms in wide use today:

1. Diffie–Hellman
2. RSA
3. Digital Signature Algorithm (DSA)

### Diffie–Hellman

It was the first usable public key algorithm. Diffie–Hellman is based on the difficulty of computing discrete logarithms. It can be used to establish a shared secret key that can be used by two parties for symmetric encryption. Diffie–Hellman is often used for IPsec key management protocols. For spontaneous communications with Diffie–Hellman, two communicating entities would each generate a random number that is used as their private keys. They exchange public keys they each apply their private keys to the other's. public key to compute identical values (shared secret key). They then use the shared secret key to encrypt and exchange information.

### Diffie–Hellman key exchange

The protocol that allows strangers to establish a shared secret key is called the Diffie–Hellman key exchange and works as follows:

- Ana and Brat have to agree on 2 large numbers, 'n' and 'g', where 'n' is a prime.
- (n - 1)/2 is also a prime and certain conditions apply to 'g'.

- These numbers may be public, so either one of them can just pick 'n' and 'g' and tell the other openly.
- Now Ana picks a large number (suppose 512-bit) 'x', and keeps it secret. Similarly Brat picks a large secret number, 'y'.
- Ana initiates the key exchange protocol by sending Brat a message containing (n, g, g<sup>x</sup> mod n)
- Brat responds by sending Ana a message containing (g<sup>y</sup> mod n)
- Now Ana raises the number Brat sent her to the xth power modulo 'n' to get [(g<sup>y</sup> mod n)<sup>x</sup> mod n]
- Brat performs a similar operation to get [(g<sup>x</sup> mod n)<sup>y</sup> mod n], Both the calculations yield (g<sup>xy</sup> mod n).

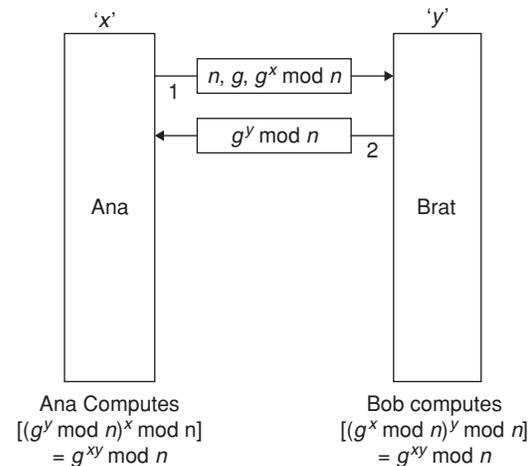


Figure 3 Diffie-Hellman key exchange

### RSA (Rivest, Shamir, Adelman)

RSA multiplies large prime numbers together to generate keys. It's strength lies in the fact that it is extremely difficult to factor the product of large prime numbers. This algorithm is the one, most often associated with public key encryption. The RSA algorithm also provides digital signature capabilities.

#### Example:

- Select two large primes = p, q p = 17, q = 11
- n = p × q = 17 × 11 = 187
- calculate φ = (p - 1) (q - 1) = 16 × 10 = 160
- select e, such that LCD (φ, e) = 1, 0 < e < φ say, e = 7
- calculate d such that d mod φ = 1
- 160k + 1 = 161, 321, 481, 641,
- Check which of these is divisible by 7
- 161 is divisible by 7 giving d = 161/7 = 23
- Key 1 = {7, 187}, key 2 = {23, 187}

### Digital Signatures

A digital signature allows a receiver to authenticate (to a limited extent) the identity of the sender and to verify the integrity of the message for the authentication process, you

must already know the sender's public key, either from prior knowledge or from some trusted third party. Digital signatures are used to ensure message integrity and authentication. In its simplest form, a digital signature is created by using the sender's private key to hash the entire contents

of the message being sent to create a message digest. The recipient uses the sender's public key to verify the integrity of the message by recreating the message digest. By this process you ensure the integrity of the message and authenticate the sender.

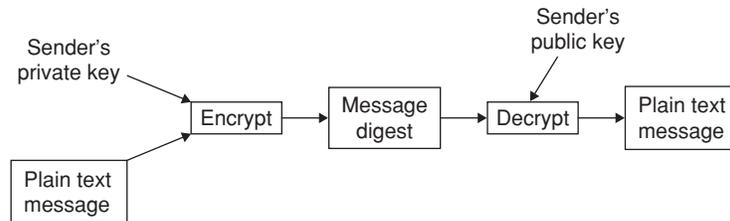


Figure 4 Digital signature

To sign a message, senders usually append their digital signature to the end of a message and encrypt it using the recipient's public key. Recipients decrypt the message using their own private key and verify the sender's identity and the message integrity by decrypting the sender's digital signature using the sender's public key. The strength of digital signatures are that they are almost impossible to counterfeit and they are easily verified.

### Digital certificate

Digital signatures can be used to verify that a message has been delivered unaltered and to verify the identity of the sender by public key. The problem with authenticating a digital signature, however, is that you must be able to verify that a public key does in fact belong to the individual or entity that claims to have sent it and that the individual or entity is in fact who or what it claims to be.

A digital certificate issued by a certification authority (CA) utilizing a hierarchical public key infrastructure (PKI) can be used to authenticate a sender's identity for spontaneous, first-time contacts. Digital certificates provide a means for secure first time spontaneous communication. A digital certificate provides a high level of confidence in the identity of the individual.

A digital certificate is issued by a trusted/unknown third party (CA) to bind an individual or entity to a public key. The digital certificate is digitally signed by the CA with the CA's private key. This provides independent confirmation that an individual or entity is in fact who it claims to be. The CA issued digital certificates that certify for the identities of those to whom the certificates were issued.

### Firewalls

Firewall is a control link between internet and organization intranet. It protects network premises from internet based attacks by providing single choke point. All the network traffic is forced to travel through this fire wall. Firewall allows only authorized traffic to pass through.

The different types of firewalls are:

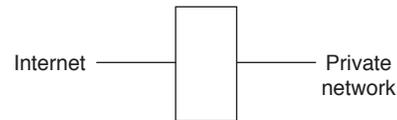
1. Packet – filtering router
2. Application level gateways

3. Circuit level gateways
4. Bastion host

### Packet filtering router

It filters packets with incoming and outgoing interfaces, and permits or denies certain services. It uses the information of transport layer like IP sources, ICMP message etc.

The drawbacks are IP address spoofing, tiny fragment attack and source routing attacks.



### Application level gateway

It provides proxies for each service, when user requests service, it validates the request as legal one and return results to the user.

Application level gateway is more secure than the packet filter.

The drawback of this gateway is processing overhead at each connection.

### Circuit-level gateway

It is application level gateway functionality for certain applications. It does not allow end-end TCP connection, rather it maintains two connections, one with the inner host and the other with the outer host. Once the connections are established TCP segment is allowed without examining contents. It only checks the incoming data.

### Bastion host

It provides a platform for the application gateway (or) circuit level gateway, it is a critical strong point in network security.

An additional authentication is required for the user who want access to proxy services. Even proxy service authenticates itself before granting the access to user.

Only essential services are installed in the Bastion host which are decided by admin.

## EXERCISES

## Practice Problems I

**Directions for questions 1 to 15:** Select the correct alternative from the given choices.

- In an encryption scheme that uses RSA, values, for  $p$  and  $q$  are selected to be 5 and 7 respectively what could be the value of  $d$ ?  
(A) 12 (B) 3 (C) 11 (D) 9
- A person  $x$  is supposed to send a document with digitized signature to another person  $y$  using public key Cryptography.  $p$  is the message.  $D_x, D_y$  are private keys of  $x$  and  $y$  respectively.  $E_x, E_y$  are public keys of  $x, y$  respectively. Select the best possible sequence of events from below:
  - $D_x(p)$
  - $D_y(p)$
  - $E_y(D_x(p))$
  - $D_y(D_x(p))$
  - $D_y(E_y(p))$
  - $D_y(E_y(D_x(p)))$
  - $E_x(D_x(p))$
  - $E_y(p)$
  - $E_x(D_y(p))$
  - $D_x(E_y(p))$

(A) (ii), (ix), (viii), (v) (B) (viii), (x), (v), (i)  
(C) (i), (iii), (v), (vii) (D) (vii), (v), (iii), (i)
- Select correct statements about PGP:
  - Uses existing cryptographic algorithms that have been quite successful.
  - Support text compression, digital signatures.
  - Takes plaintext as feed and generates base-64 text.
  - No key management capability is provided.

(A) (i), (ii), (iii) (B) (ii), (iii), (iv)  
(C) (i), (iii), (iv) (D) (i), (ii), (iv)

## Linked answer questions 4 and 5:

- Using mono alphabetic substitution a string a b b a c a a b c d is transformed to one of the below strings. Select the most appropriate option:  
(A) p q q p r p p s r s (B) j t t x j j i t x t x  
(C) u s s u a u s a b (D) d c c d b b b c b a
- Using the mapping obtained above, encrypt the phrase 'bad cab' using same method: Assume space is not encrypted.  
(A) q p s r p q (B) t j z x j t  
(C) s u b a u s (D) c d a b d c
- Select the correct statements with regard to packet filters of a firewall:
  - They are usually driven by a table with information in regards to acceptable sources and destinations.
  - Default rules about what needs to be done in regards to packets coming from or going to other machines.
  - Can block TCP ports.

- (A) (i), (ii) (B) (ii), (iii)  
(C) (i), (iii) (D) (i), (ii), (iii)

- What is meant by non-repudiation in the area of digital signatures?
  - Receiver verifying the signature of the sender.
  - Receiver concocting the message.
  - Sender denying having signed digitally.
  - Receiver changing the contents after receiving the signed document.
- Which of the following statements about DES is/are true?
  - DES is public key algorithm.
  - DES has 19 distinct stages.
  - In the 16 iterations of DES, different keys are used.

(A) (i), (ii) (B) (ii), (iii)  
(C) (i), (iii) (D) (i), (ii), (iii)
- Which of the below represents Triple encryption using DES? ( $P$  is the unencrypted input, 'C' is encrypted output,  $k_1, k_2, k_3$  are keys used in encryption and decryption,  $E$  stands for encryption and  $D$  stands for decryption).
  - 
  - 
  - 
  -
- Which of the below statements are applied for cipher block chaining?
  - Each plaintext block is XOR'ed with previous block before encryption.
  - Encryption is a mono alphabetic substitution cipher.
  - Cipher block chaining can result in same plaintext blocks encrypted to different cipher text blocks.

(A) (i), (ii) (B) (ii), (iii)  
(C) (i), (iii) (D) (i), (ii), (iii)
- Which of the below statements are applied to RSA algorithm?
  - RSA is a relatively slow algorithm when encrypting large data.
  - Mainly used where key is to be distributed.
  - The strength of the algorithm lies in the fact that determining the key can take exceedingly long time by brute force.

- (A) (i), (ii)                      (B) (ii), (iii)  
 (C) (i), (iii)                      (D) (i), (ii), (iii)
12. The security and usefulness of a digital signature depends on  
 (A) A public hash function  
 (B) A two-way hash function  
 (C) Protection of user's private key  
 (D) Protection of user's public key
13. Let ' $M$ ' be the message to be encrypted,  $E$  be Encryption key and  $N$  be the product of two random prime numbers, then what is the cipher text using RSA algorithm?  
 (A)  $C = E^m \text{ mod } N$               (B)  $C = M^E \text{ mod } N$   
 (C)  $C = N^E \text{ mod } M$               (D)  $C = E^N \text{ mod } M$
14. Which of the following best describes the decryption in Triple DES?

- (A) Plain text =  $D_{K_1}(E_{K_2}(D_{K_1}(\text{cipher text})))$   
 (B) Plain text =  $D_{K_1}(E_{K_2}(D_{K_3}(\text{cipher text})))$   
 (C) Plain text =  $E_{K_1}(D_{K_2}(E_{K_1}(\text{cipher text})))$   
 (D) Plain text =  $E_{K_1}(D_{K_2}(E_{K_1}(\text{cipher text})))$
15. In which cipher mode, all cipher blocks will be chained so that if one is modified the cipher text cannot be decrypted correctly?  
 (A) Electronic Code Book  
 (B) Cipher Block Chaining  
 (C) Cipher Feedback Mode  
 (D) Counter Mode

## Practice Problems 2

**Directions for questions 1 to 15:** Select the correct alternative from the given choices.

1. 'All algorithms must be public only the keys are secret' is  
 (A) Rijndael Principle  
 (B) Kerckhoff's principle  
 (C) Rivest shamir Adleman principle  
 (D) None of these
2. Pretty Good Privacy encrypts data by using a block cipher called  
 (A) RSA                              (B) MD5  
 (C) IDEA                              (D) DES
3. E-mail security package is related to  
 (A) Pretty Good Privacy  
 (B) DNS spoofing  
 (C) Secure Socket Layer  
 (D) Transport Layer Security
4. Which of the following protocols will be proxy, on an application firewall?  
 (A) IPX                              (B) FTP  
 (C) POP                              (D) SMS
5. A good recommendation is that if a private key is \_\_\_\_\_ or longer, the key is thought to be secure.  
 (A) 40 bits                              (B) 60 bits  
 (C) 70 bits                              (D) 80 bits
6. Which issue is related to server side security?  
 (A) Protection of the server from legitimate web access  
 (B) Security of the information stored on server  
 (C) Security of the customer's physical credit card  
 (D) Security of the customer's computer
7. Which of the following is not an active attack?  
 (A) Denial of service              (B) Traffic Analysis  
 (C) Replay                              (D) Masquerade
8. Verifying the true identity of the sender of a message recipient is known as \_\_\_\_\_.

- (A) Authentication              (B) fabrication  
 (C) Cryptography              (D) availability
9. In which of the following techniques, letters are arranged in a different order?  
 (A) Transposition  
 (B) Substitution  
 (C) Private key Encryption  
 (D) None of the above
10. In which type of attack, Algorithm, cipher text, chosen plaintext and cipher text are known?  
 (A) Cipher text only  
 (B) Known plain text  
 (C) Chosen cipher text  
 (D) Chosen text
11. In which type of ciphers the encryption depends on current state?  
 (A) Link cipher  
 (B) Block cipher  
 (C) Stream cipher  
 (D) Current cipher
12. Traffic Analysis can be counted using  
 (A) Encryption                      (B) Decryption  
 (C) Replay                              (D) Data padding
13. DES Algorithm is vulnerable to  
 (A) Masquerade attack  
 (B) Replay attack  
 (C) Denial of service  
 (D) Brute Force attack
14. What is the size of key in Triple DES?  
 (A) 168 bits                              (B) 112 bits  
 (C) 56 bits                              (D) Either (A) or (B) or (C)
15. Direct digital signature involves  
 (A) Source only  
 (B) Destination only  
 (C) Communicating parties, sender and receiver.  
 (D) Everyone including communicating parties.

## PREVIOUS YEARS' QUESTIONS

1. Suppose that everyone in a group of  $N$  people wants to communicate secretly with the  $N - 1$  others, using symmetric key cryptographic system. The communication between any two persons should not be decodable by the others in the group. The number of keys required in the system as a whole to satisfy the confidentiality requirement is [2015]
- (A)  $2N$  (B)  $N(N - 1)$   
 (C)  $N(N - 1)/2$  (D)  $(N - 1)^2$
2. Consider that  $B$  wants to send a message  $m$  that is digitally signed to  $A$ . Let the pair of private and public keys for  $A$  and  $B$  be denoted by  $K_x^-$  and  $K_x^+$  for  $x = A, B$ , respectively. Let  $K_x(m)$  represent the operation of encrypting  $m$  with a key  $K_x$  and  $H(m)$  represent the message digest. Which one of the following indicates the CORRECT way of sending the message  $m$  along with the digital signature to  $A$ ? [2016]
- (A)  $\{m, K_B^+(H(m))\}$  (B)  $\{m, K_B^-(H(m))\}$   
 (C)  $\{m, K_A^-(H(m))\}$  (D)  $\{m, K_A^+(m)\}$
3. Anarkali digitally signs a message and sends it to Salim. Verification of the signature by Salim requires [2016]
- (A) Anarkali's public key.  
 (B) Salim's public key.  
 (C) Salim's private key.  
 (D) Anarkali's private key.
4. A sender  $S$  sends a message  $m$  to receiver  $R$ , which is digitally signed by  $S$  with its private key. In this scenario, one or more of the following security violations can take place.
- (I)  $S$  can launch a birthday attack to replace  $m$  with a fraudulent message.  
 (II) A third party attacker can launch a birthday attack to replace  $m$  with a fraudulent message.  
 (III)  $R$  can launch a birthday attack to replace  $m$  with a fraudulent message.
- Which of the following are possible security violations? [2017]
- (A) (I) and (II) only (B) (I) only  
 (C) (II) only (D) (II) and (III) only
5. In a RSA cryptosystem, a participant  $A$  uses two prime numbers  $p = 13$  and  $q = 17$  to generate her public and private keys. If the public key of  $A$  is 35, then the private key of  $A$  is \_\_\_\_\_. [2017]

## ANSWER KEYS

## EXERCISES

## Practice Problems 1

1. C    2. C    3. A    4. C    5. C    6. D    7. C    8. B    9. C    10. C  
 11. D    12. C    13. B    14. B    15. B

## Practice Problems 2

1. B    2. C    3. A    4. B    5. C    6. B    7. B    8. A    9. A    10. D  
 11. C    12. D    13. D    14. D    15. C

## Previous Years' Questions

1. C    2. B    3. A    4. B    5. 11

## TEST

## COMPUTER NETWORKS

Time: 60 min.

**Directions for questions 1 to 30:** Select the correct alternative from the given choices

- What is the Hamming distance between 000, 011?  
(A) 0 (B) 1  
(C) 2 (D) 3
- Consider the given data:
 

| Dataword | Codeword |
|----------|----------|
| 00       | 00000    |
| 01       | 01011    |
| 10       | 10101    |
| 11       | 11110    |

Find the minimum hamming distance?  
(A) 2 (B) 3  
(C) 4 (D) 5
- In Go-back- $n$ , what should be the Window size?  
(A)  $2^m$  (B)  $2^{m-1}$   
(C)  $2^{m-2}$  (D)  $2^{2m}$
- If there are 16 sequence numbers, what are the sender and receiver window sizes in go-back- $n$  and selective repeat respectively?  
(A) (15, 1) (8, 8) (B) (14, 2) (8, 8)  
(C) (15, 1) (7, 8) (D) (15, 1) (8, 7)
- A code needs to be designed with 8 data bits and  $r$  check bits. What is the minimum value of  $r$  in order to correct single bit errors?  
(A) 1 (B) 2  
(C) 3 (D) 4
- A code has hamming distance of 6. What is the maximum number of bit errors that can be corrected?  
(A) 1 (B) 2  
(C) 3 (D) 4
- In the above case what is the number of errors that can be detected?  
(A) 3 (B) 4  
(C) 5 (D) 6
- CRC is being used to do error detection and correction. The frame with data 101001001 needs to be sent and the generator polynomial being used is  $x^4 + x + 1$ . What is the final transmitted frame?  
(A) 1010010011110 (B) 1010010010010  
(C) 1010010011010 (D) 1010010010000
- OSI model seven layer is based on which of the following principles:  
(A) A layer should be created where a different level of abstraction is needed  
(B) Each layer should perform a well defined function

- (C) The layer boundaries should be chosen to minimize the information flow across the interfaces  
(D) All the above
- Which of the following is/are the tasks of physical layer?  
(A) How to link two or more devices physically  
(B) What type of data flow is needed between two devices  
(C) Type of topology required  
(D) All the above
  - The functions of the data link layer are  
(A) It provides services to network layer and accepts services from physical layer  
(B) It is responsible for error control and detection within the network.  
(C) It regulates the amount of data that can be transmitted on one line  
(D) All the above
  - Which one of the following layers deals with problems that arise when packet travels from one network to another?  
(A) Transport layer (B) Physical layer  
(C) Data link layer (D) Network layer
  - What is the main function of the network layer?  
(A) Routing  
(B) Congestion control  
(C) Both (A) and (B)  
(D) None of these.
  - Which layer ensures interoperability among the communicating devices, and also computers to communicate even if their internal representation is different?  
(A) Session layer (B) Transport layer  
(C) Presentation layer (D) Application layer
  - Which of the following is not a layer in TCP/IP reference model?  
(A) Application layer (B) Transport layer  
(C) Data link layer (D) Host to Network layer
  - Suppose we want to transmit a character 'C', the binary value is 1000011 if we pass through an even parity generator then the output is  
(A) 10000110 (B) 10000111  
(C) 1000011 (D) 1000010
  - What type of frames can be recognized by stop and wait protocols?  
(A) Damaged frames  
(B) Lost frames  
(C) Lost of acknowledgement frames  
(D) All the above

18. IEEE project 802 divides the data link layer into two sub layers. What is the upper sublayer?  
 (A) LLC (B) MAC  
 (C) PDU (D) HDLC

**Common data for questions 19 and 20:** When a data frame, arrives at the receiver, instead of sending an acknowledgement separately the receiver rests itself and waits until the network layer passes it the next packet. The acknowledgement is attached to the outgoing data frame.

19. The technique of temporarily delaying outgoing ACK so that they can be hooked onto the next outgoing data frame is called\_\_\_\_\_
- (A) Pipelining (B) Piggybacking  
 (C) Flooding (D) None
20. Which layer implements technique of piggybacking?  
 (A) Physical layer (B) Data link layer  
 (C) Transport layer (D) Session layer
21. What is the protocol used in one bit sliding window protocol?  
 (A) Unrestricted simplex  
 (B) Simplex stop and wait  
 (C) Simplex protocol for noisy channel  
 (D) Restricted duplex.
22. The technique of keeping the sender window appropriately in such a way, that it can continuously transmit frames for a time equal to the round trip time, so that acknowledgement of first frame will arrive just after transmitting the last frame, is called  
 (A) Flooding (B) Piggy backing  
 (C) Pipelining (D) Selective repeat
23. Pick the incorrect statement from the following  
 (A) Go-Back- $N$  method requires more storage at the receiving side.  
 (B) Selective repeat involves complex logic than Go-back- $N$   
 (C) Go-back- $N$  has better line utilization  
 (D) Selective repeat has better line utilization
24. In stop and wait flow control, to send 'n' data packets how many acknowledgements are needed.

- (A)  $n$  (B)  $2n$   
 (C)  $n - 1$  (D)  $n + 1$

25. In sliding window flow control, if the window size is 64 what is the range of sequence numbers?  
 (A) 0 to 63 (B) 0 to 64  
 (C) 1 to 63 (D) 1 to 64
26. In Go-Back- $N$  Automatic Repeat Request (ARR), if frames 4, 5, 6 are received successfully, the receiver will send which ACK number to the sender?  
 (A) 5 (B) 6  
 (C) 7 (D) 4
27. Which of the following are the responsibilities of a token ring monitor station?  
 (A) Check to see that token is not lost  
 (B) Taking action when ring breaks  
 (C) Clearing the ring when garbled frames appear  
 (D) All the above

**Common data for questions 28 and 29:**

28. In ISO-OSI reference model the layer that provides necessary translation of different control codes, character set and graphic character and it ensures interoperability among communicating devices. The above explanation is about which of the following layers?  
 (A) Session layer  
 (B) Data link layer  
 (C) Presentation layer  
 (D) Application layer
29. What are the other tasks performed by the above layer?  
 (A) Encryption and compression  
 (B) Token management and synchronization  
 (C) Error detection and error correction  
 (D) None of these
30. The layer that takes a raw transmission and transforms it into a line that appears free of undetected transmission errors and it takes care of traffic regulation to keep fast transmitter from drowning slow receiver. The layer that provides these services is  
 (A) Physical layer (B) Transport layer  
 (C) Data link layer (D) Application layer

### ANSWER KEYS

1. C 2. B 3. A 4. A 5. D 6. B 7. C 8. D 9. D 10. D  
 11. D 12. D 13. C 14. C 15. C 16. B 17. D 18. A 19. B 20. B  
 21. B 22. C 23. C 24. A 25. A 26. C 27. D 28. C 29. A 30. C

## Part B Information Systems

**Chapter 1:** Process Life Cycle 8.79

**Chapter 2:** Project Management  
and Maintenance 8.89

U

N

I

T

8

*This page is intentionally left blank*

# Chapter 1

## Process Life Cycle

### LEARNING OBJECTIVES

- Introduction
- Process vs program
- Software component and elements
- Information gathering
- Requirement analysis
- Feasibility study
- Data flow diagram
- Process specification
- Input/output design
- Software process life cycle
- Software process model

### INTRODUCTION

A system can be defined as an orderly grouping of interdependent components linked together according to a plan to achieve a specific objective.

**Example:** Telephone system, transportation system, accounting system, etc.

### PROCESS VERSUS PROGRAM

A software process gives all steps used to create a software application, from the customer's requirements to the finished product.

- The software process determines the organization and flexibility of the project.
- There are several different software processes and each describes their own solution to develop a valid software.
- Software programs are written programs or rules with associated documentation pertaining to the operation of a computer system.

### SOFTWARE COMPONENTS AND ELEMENTS

#### Software Component

It is a software element that can be independently deployed and composed without modification according to a composition standard.

- A component model implementation is the dedicated set of executable software elements required to support the execution of components.
- A component has clearly defined interfaces.

- An interface standard is the mandatory requirement enforced to enable software elements to directly interact with other software elements.
- An interface standard declares, when an interface comprises.

**Standard** An object or measure serving as a basis to which others should conform, by which the quality of others is judged.

#### Software Element

A sequence of abstract program statements that describe computations, which has to be performed by a machine.

#### Interface

It describes the behavior of a component that is obtained by considering only the interactions of that interface and by hiding all other interactions.

- An abstraction of the behavior consists of subset of the interactions of one component together with a set of constraints.

#### Interaction

It is defined as action between 2 or more software elements.

**Composition** It is a combination of 2 or more software components, the newly formed component, behaviour will be at a different level of abstraction.

The characteristics of new component is determined the components combined and the way in which they are combined.

## INFORMATION GATHERING

Complete and accurate information is essential in building computer-based systems. Information about the organization, the staff who uses the system and the workflow should be gathered.

Information about the organization's policies, goals, objectives and structure explains the kind of environment the computer-based system should produce.

Information about the people who run the present system, their job functions and information requirements, the relationships of their jobs to the existing system and the interpersonal network that holds the user groups together are required for determining the importance of the existing system for the organization and also for planning the proposed system.

Workflow focuses on what happens to the data through various points in a system and can be shown by a data flow diagram or a system flow chart.

Information can be gathered by studying documents, forms and files of existing system. Onsite observation of the system is also an effective method for gathering information. It is the process of recognizing and noting people, objects and occurrences to obtain information. Interview is one of the most often and oldest method for gathering information. Interview has the advantage of identifying relations or verifying information and also capture information face to face with the concerned person. Questionnaire is another method for gathering information and is an inexpensive mean for gathering data which can be tabulated and analyzed quickly. Visiting companies that have developed similar systems, reading journals and other computer related books, which specify how others have solved similar problems is also another means of information gathering.

## REQUIREMENT ANALYSIS

Requirement analysis results in specification of software's operational characteristics, indicates software's interface with other system elements and establishes constraints that the software must meet.

During requirement analysis, the primary focus should be on *what* not *how*. It should define what user interaction occurs in a particular circumstance, what objects does the system manipulate, what functions must the system perform, what behaviours does the system exhibit, what interfaces are defined and what constraints applied.

The requirement analysis model must achieve three primary objectives:

1. To describe what the customer requires.
2. To establish a basis for the creation of a software design
3. To define a set of requirements that can be validated once the software is built.

## Requirement Negotiation

Requirement negotiation is required to have a win-win result. The customer should get product which satisfies most of his/her needs, and software team should develop a product within a budget, working in real-time environment and within deadlines.

Boehm has defined negotiation activities at the beginning of each software process iteration.

- Identify the key stake holder's system or subsystem.
- Determine the 'win conditions' of the stake holder.
- Negotiate the 'win conditions' of stake holder and establish them into win-win condition.

## Requirement Elicitation

Requirement elicitation is gathering the requirements from stakeholders, customers, etc. The question and answer format is suitable for the first encounter with users and the remaining phases are replaced with requirement elicitation. As we can't get all the requirements by having questions and answers session, requirement elicitation practices should be implemented, which includes interviews, workshops, user scenarios, etc.

The approaches that are followed for eliciting the requirements are:

1. Collaborative requirement gathering
2. Quality function deployment
3. User scenarios
4. Elicitation work products

## Functional Requirements

Functional requirements are primary actions that must take place in software in accepting and processing the input and in processing and generating the output.

Functional requirements capture the intended behaviour of the system, which could be expressed as service task (or) functions of the system.

These are core functionalities of the system. It also includes exact sequence of operations, input validation, mapping of outputs to inputs and error handling and recovery. These requirements are implemented in system design.

## Non-functional Requirements

Non-functional requirements are expected requirements of a user. These define operational constraints based on the user characteristics.

Non-functional requirements are product, business and external-based. These requirements define how a software system has to be. It also defines the quality of product, type of reliability and usability of the system. Implementation requirements depend on organization, and delivery requirements are defined in the non-functional requirements.

Non-functional requirements are implemented in system-architecture.

## Measuring Requirements

The requirements are the major component of project. The metrics for the requirement activities are:

1. Product size
2. Requirement quality
3. Requirement status
4. Requirement change (request for changes)
5. Effort

**Product size** refers to the count of functional and non-functional requirements. It tracks, whether these requirements are implemented as a function of time.

**Requirement quality** refers to the inspection of specification of requirements, counting the defects, missing of requirements incompleteness, ambiguities, etc.

**Requirements status** is monitoring of requirements over time, gives out project status. The status could be proposed, approved, implemented, verified deferred, deleted, rejected.

**Requirement management** handles the addition, modify and deletion of requirement, track the change of requirement which affects multiple requirements of different level.

**Effort** is the time taken to record requirements activities which includes development and management of requirements.

## Types of Requirements

These are the services that a software system has to provide and constraints under which it must operate.

### User requirements

- Written for customers.
- These statements will be given in natural language and diagrams of the services that the system provides and operational constraints.

### System requirements

- Written as a contract between contractor and client.
- A structured document with detailed descriptions of the system services.

### Software specification

- Written for developers.
- A detailed description of software that can act as basis for a design (or) implementation.

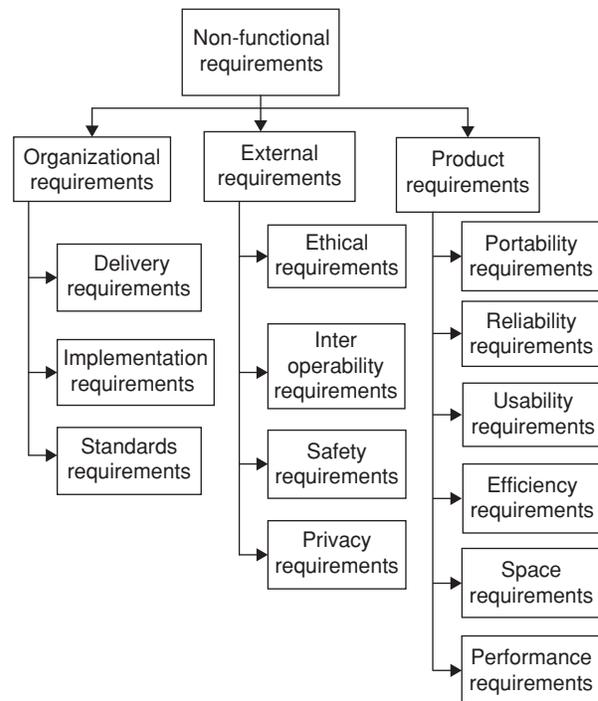
### Functional requirements

The system should provide statements of services, how the system should react to particular inputs and how the system should behave in particular situations.

## Non-functional requirements

Functions offered by the system such as timing constraints, constraints on the development process, standards, etc.

1. **Product requirements:** It specifies, that the delivered product must behave in a particular way.  
Example: Execution speed, reliability
2. **Organizational requirements:** These are consequences of organizational policies and procedures.  
Example: Process standards, implementation requirements
3. **External requirements:** These arise from factors which are external to the system and its development process.  
Example: Interoperability requirements



## Domain requirements

These come from the application domain of the system that reflects the characteristics of the domain.

- These could be functional or non-functional.

## FEASIBILITY ANALYSIS

Feasibility study is a test of a system proposal according to its workability, impact on the organization, ability to meet user needs, and effective use of resources.

The objective of a feasibility study is not to solve the problem but to acquire a sense of its scope. Costs and benefits are estimated with greater accuracy at this stage. Feasibility analysis helps to identify the best solution to the end user.

The key considerations involved in feasibility analysis are:

1. Economic feasibility
2. Technical feasibility
3. Behavioral feasibility

Economic analysis is the most frequently used method for evaluating the effectiveness of a candidate system. Also known as cost/benefit analysis, economic analysis determines the benefits and savings that are expected from the candidate system and compares them with the costs. If benefits outweigh costs, then decision is made to design and implement the system, else alterations are made if it has a chance of being approved.

Technical feasibility is concerned with hardware and software requirements to implement the system. Technical analysis centres around the existing computer system (hardware, software, etc.) identifies, to what extent it can support the proposed addition. Additional hardware and software (OS, databases) requirements are identified and checks whether financial considerations/constraints can accommodate these technical enhancements.

Behavioural analysis makes an estimate of how strong a user staff is likely to react towards the development of a computerized system. Computer installations usually changes employee job status, and also there may be transfer, training period, etc. Thus the introduction of a new system requires special effort to educate, sell and train the staff on new ways of conducting business.

## DATA FLOW DIAGRAMS (DFD)

DFD also called bubble chart, clarifies system requirements and identifies major transformations that will become programs in system design. It functionally decomposes the requirements specification down to the lowest level of detail.

The four DFD symbols are:

1.  (or)  Source/Destination of data
2.  (or)  Data flow
3.  (or)  (or)  Process
4.  (or)  (or)  Data store

The first symbol defines a source or destination of system data. The second symbol specifies data flow direction. It can be considered as a pipeline through which the information flows. The third symbol represents a process that transforms incoming data flows into outgoing data flows, and the fourth symbol is used to represent storage of data.

In short, DFD takes an input-process-output view of a system. That is, data objects flow into the software is transformed by processing elements and resultant data objects flow out of the software. Data objects are represented by labelled arrows and transformations are represented by circles (also called

bubbles). The DFD is presented in a hierarchical fashion. That is, the first data flow model sometimes called level 0 DFD or context diagram represents the system as a whole. Subsequent data flow diagrams refine the context diagram, providing increase in detail with each subsequent level.

## PROCESS SPECIFICATION (PSPEC)

The process specification (PSPEC) is used to describe all flow model processes that appear at the final level of refinement. The content of the process specification can include narrative text, a program design language (PDL) description of the process algorithm, mathematical equations, tables or UML activity diagrams.

By providing a PSPEC to accompany each transformation (bubble) in the flow model, a 'mini-spec' can be created that serves as a guide for design of the software component that will implement the transformation.

## INPUT/OUTPUT DESIGN

### Input Design

The most common cause of errors in data processing is inaccurate input data. Errors occurred during data entry can be controlled by input design.

Input design is the process of converting user-originated inputs to a computer-based format.

Input data is collected and organized into groups of similar data. The goal of designing input data is to make data entry easy, logical and free from errors.

Source data is captured initially on original paper or a source document. A source document should be logical and easy to understand. Each area in the form should be clearly identified and should specify to the user what to write and where to write.

Source documents may enter into the system from punch cards, diskettes, optical character recognition (OCR) reader, Magnetic ink character recognition (MICR) reader, barcode reader, etc. Touch screen or voice input can be used for online data entry, for example, ATM.

There are three major approaches for entering data into the computer—menus, formatted forms and prompts.

A menu is a selection list that simplifies computer data access or entry. The user can choose what to enter from a list of options. Though a menu limits a user choice of responses, it reduces the chances of errors in data entry.

A formatted form is a preprinted form or a template that requests the user to enter data in appropriate locations (fill-in the blank type form). The form is displayed on the screen and the user can fill information by positioning the cursor in appropriate text boxes.

In prompt, the system displays one inquiry at a time, asking the user for a response, for example, asking for user-id and password.

## Output Design

Computer output is the most important and direct source of information to the user. Efficient and intelligible output design will improve system's relationships with the user and helps in decision making.

The devices available for providing computer-based output are printer, CRT screen display, audio response (speaker), plotters, etc.

The task of output preparation is very critical, regaining skill and ability to align user requirements with the capabilities of the system in operation.

## SOFTWARE PROCESS LIFE CYCLE

A software process can be defined as a frame work of the activities, actions and tasks that are required to build quality software.

All these activities, actions and tasks reside within a frame work or model that defines their relationship with the process and with one another.

A generic process frame work for software engineering encompasses five activities:

**Communication** Proper communication and collaboration with the customer is made in this activity to understand the objectives for the project and also to gather requirements that help to define software features and functions.

**Planning** This activity develops a software project plan which defines the software engineering work by specifying the technical tasks to be conducted, the risks that may occur, the resources that will require, the work products to produce and the work schedule.

**Modelling** A software engineer creates models to better understand software requirements and the design that will achieve those requirements.

**Construction** This activity combines code generation and testing required to uncover errors in the code.

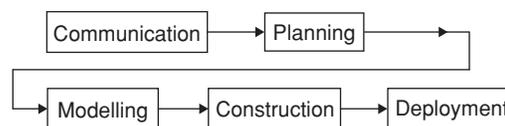
**Deployment** In this activity, the software (as a complete product or as a partial increment) is delivered to the customer. The customer evaluates the delivered product and provides feedback based on evaluation.

Another important aspect of the software process called process flow describes how the frame work activities and the actions and tasks that occur within each framework activity are organized with respect to sequence and time.

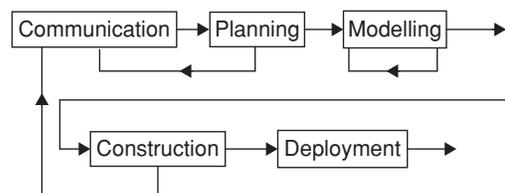
## Process Quality and Improvement

Quality refers to characteristic or attribute of something. Process quality factors are portability, usability, reusability, correctness and maintainability. The process quality is the implementation of the following steps firstly initiates the process and design the solutions, implement these solutions with the impact demonstration.

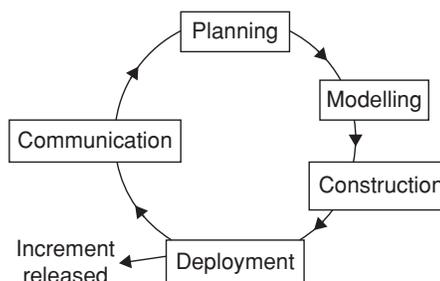
**Linear process flow** executes each of the five framework activities in sequence, beginning with communication and ends with deployment.



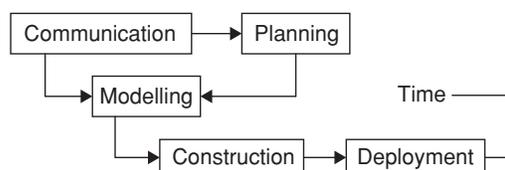
**Iterative process flow** repeats one or more of the activities before proceeding to the next activity.



**Evolutionary process flow** executes the activities in a circular manner. Each circuit through the five activities leads to a more complete version of the software.



**Parallel process flow** executes one or more activities in parallel with other activities.



## The Unified Process

The unified process (UP) is an attempt to draw on the best features and characteristics of conventional software process models. It recognizes the importance of customer communication and streamlined methods for describing the customer's view of a system. It helps the architect focus on the right goals, such as understandability, reliance to future changes, and reuse. It suggests a process flow that is iterative and incremental, providing the evolutionary feel that is essential in modern software development.

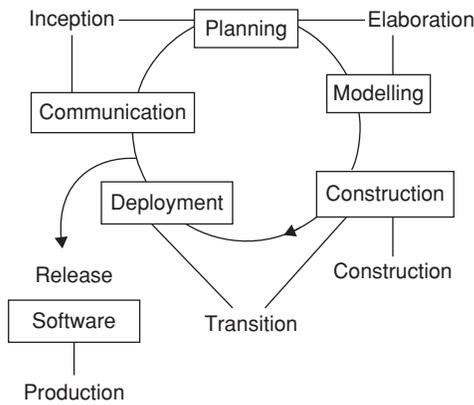


Figure 1 Phases of the unified process

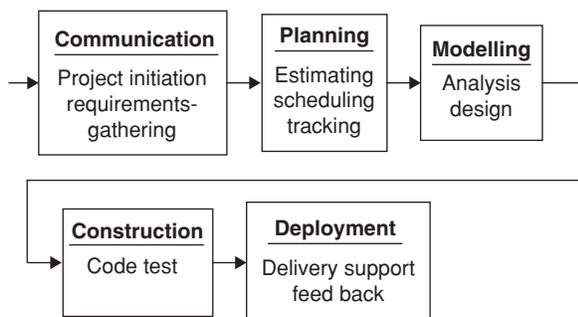
The unified process is an incremental model in which five phases are defined:

1. Inception phase: Encompasses both customer communication and planning activities and emphasizes the development and refinement of use cases as a primary model.
2. Elaboration Phase: Encompasses the customer communication and modelling activities focusing on the creation of analysis and design models with an emphasis on class definitions and architectural representations.
3. Construction phase: Refines and translates the design model into implemented software components.
4. Transition phase: Transfers the software from the developer to the end user for beta testing and acceptance.
5. Production phase: Ongoing monitoring and support are conducted. Defect reports and requests for changes are also submitted and reevaluated.

## SOFTWARE PROCESS MODELS

### The Waterfall Model

The waterfall model also called classic life cycle, follows a systematic sequential approach to software development.



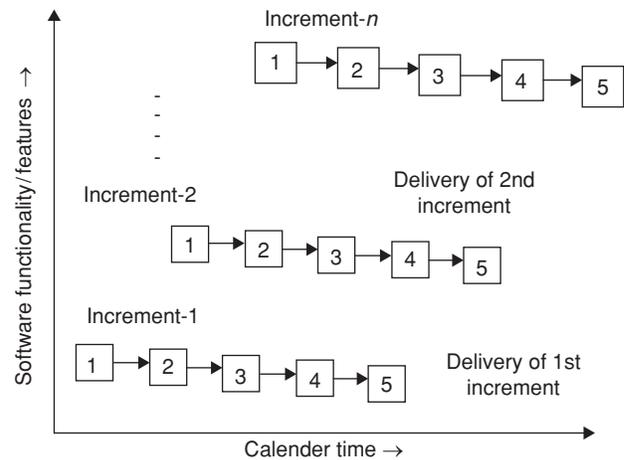
It begins with customer specification of requirements and progresses through planning, modelling, construction and deployment, culminating in on-going support of the completed software.

The waterfall model is the oldest paradigm for software engineering. The problems encountered when this model is applied are:

1. Real projects rarely follow the sequential flow that the model proposes.
2. This model requires the requirements explicitly which the customer cannot state all the requirements as it is difficult.
3. A working version of the program will not be available until late in the project time span. If a major blunder is undetected until the working program is reviewed, it can be disastrous.

### Incremental Process Model

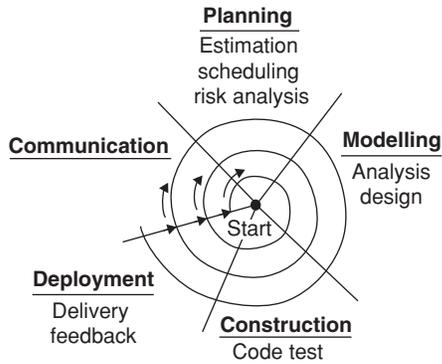
1. Communication
2. Planning
3. Modelling (analysis, design)
4. Construction (code, test)
5. Deployment (delivery, feedback)



Incremental development is particularly useful when staffing is unavailable for a complete implementation by the business deadline that has been established for the project.

### Spiral Model

Spiral model is an evolutionary software process model. Using spiral model, software is developed in a series of evolutionary releases. During early iterations, the release might be a model or prototype. Later iterations produce more complete versions of the system.



The spiral development model is a risk-driven process model generator that is used to guide multi stakeholder concurrent engineering of software intensive systems.

The spiral model is a realistic approach to the development of large scale systems and software. It uses prototyping as a risk reduction mechanism but, more importantly, enables the developer to apply the prototyping approach at any stage in the evolution of the product. At demands a direct consideration of technical risks at all stages of the project.

**Conceptual Modelling**

Conceptual modelling refers to abstraction of a model which fits for the purpose. The purpose of this modelling is to make model valid credible, feasible and useful.

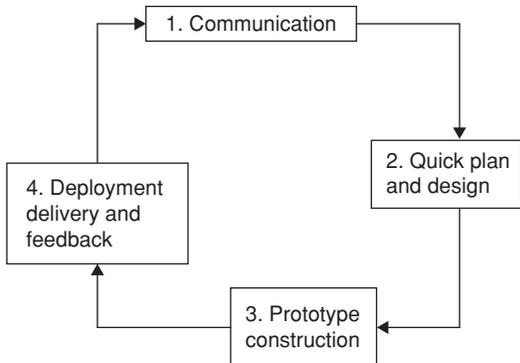
The main objective of conceptual modelling is improvising the understanding of an individual with respect to the system, an approach which will convey the system details among the stakeholders.

For the extraction of system specifications when a software is developed, some of the failures could occur in future due to lack of requirements [unclear requirements (or) changing requirements] This could be traced with the help of conceptual modelling.

**Prototyping Model**

Prototyping model is used when the user is not sure about the addition of requirements in the product. It is also implemented when the developer is not sure about the algorithm efficiency, operation system adaptability, etc. Prototyping paradigm provides the approaches.

The prototyping model is implemented as follows:



Prototyping model starts with communication in which software objectives is defined requirement, identification are done. In quick design, all the software aspects are represented quick design leads to prototype construction.

This prototype model is deployed, in which requirements are evaluated and refined by customer.

The iteration is done until customer gets satisfied with the needs at the same time developer will come to know what are the needs to be done.

**Disadvantages**

1. Developer may compromise at implementation, as prototyping works quickly. Un-ideal implementation issues may become an integral part of the system.
2. Customer just sees the working version of the software, he could not able to consider the quality of software and long-term maintenance.

Though there are some problems with prototyping, but it is effective paradigm for the software engineering when a software is developed using prototyping, both developer and customer should agree on the prototype.

It is more advantageous when the customer and user are not sure what they want it maintains a template of the older software.

**Role of metrics and measurement in software development**

The software attributes that were present in process, project and product levels are called measurement.

The metric refers to the attributes that are included in the project.

A software engineer gets the measurements and develops the metrics.

Measurement is done in two ways:

1. Direct measure
2. Indirect measure

Direct measure includes the lines of code, (least, moderate, worst) execution speed and size of the memory.

Indirect measure is done with the help of functional points. It measures quality, maintainability, efficiency and reliability.

Metric is used to control the cost, schedule, project quality. It means metric provides information for the control of process development.

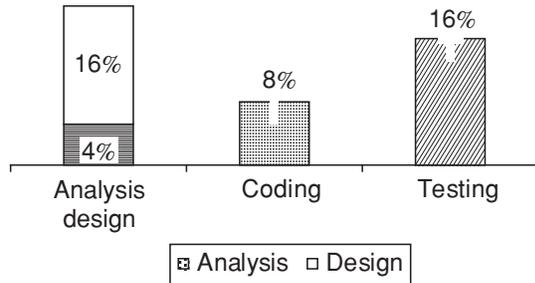
**Effort distribution with phases**

Software development is done in phases. It includes analysis, design, coding and testing.

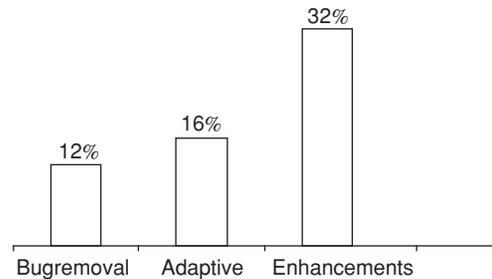
Design and testing plays major role in development, while coding is having least preference.

40% of the efforts were done on development and 60% of efforts are on the maintenance.

Distribution of the efforts on the development is shown below:



Maintenance includes removal of bug and corrective maintenance, adaptive maintenance and enhancement. Distribution of efforts in maintenance is shown below:



## EXERCISES

### Practice Problems I

**Directions for questions 1 to 15:** Select the correct alternative from the given choices.

- Which of the following statements is true?
  - The first step to the system study project is to announce the study project.
  - During the system study analysis determine manager's information needs by asking questions.
  - During the system study, flowcharts are drawn using general symbols.
  - All the above
- Which of the following statement(s) is true regarding the spiral model of software development?
  - In the spiral model of software development, the primary determinant in selecting activities in each interaction is risk.
  - The spiral model is a risk driven process model generator that is to guide multi-stakeholder. Concurrent engineering of software intensive systems.
  - Using the spiral model, software is developed in a series of evolutionary releases.
  - All the above
- Which of the following is a step in feasibility analysis?
  - Form a project team and appoint a project head.
  - Determine and evaluate performance and cost effectiveness of each candidate system.
  - Weigh system performance and cost data.
  - All the above
- Which of the following statement(s) is true?
  - The risk driven nature of the spiral model allows it to accommodate any mixture of specification oriented or some other approach.
  - Each cycle of spiral is completed by review which covers all the products developed during that cycle, including plans for the next cycle.
  - Spiral model works for development as well as enhancement project.
  - All the above
- Data flow diagram, regular expression and transition table can be combined to provide
  - decision table for functional specification of system software.
  - finite state automata for functional specification of system software.
  - event table for functional specification of system software.
  - None of these
- Which of the following statements are true about software configuration management tool?
  - It keeps track of the schedule based on the mile stones reached.
  - It manages man power distribution by changing the structure of the project.
  - It maintains different versions of the configurable items.
  - All the above
- The cost incurred on a project was ₹250,000 and benefits were ₹30,000 per month. The payback period using simple pay back method is
 

|               |                |
|---------------|----------------|
| (A) 8 months  | (B) 8.3 months |
| (C) 12 months | (D) 1.2 months |
- Which of the following phase has the maximum effort distribution?
 

|                          |                           |
|--------------------------|---------------------------|
| (A) Testing              | (B) Information gathering |
| (C) Requirement analysis | (D) Coding                |
- Which of the following statement is true regarding cost benefit analysis?
  - It evaluates tangible and non-tangible factors.
  - It estimates the hardware and software costs.
  - It compares the cost with the benefits of introducing a computer-based system.
  - All the statements are true.
- A project is considered economically feasible if the following factor holds good.
  - Return on investment (ROI)
  - Total cost of ownership (TCO)

- (C) Gross domestic product (GDP)  
 (D) Net present value (NPV)
11. At the end of the feasibility study the system analyst  
 (A) meets the users for a discussion.  
 (B) gives system proposal to management.  
 (C) gives a feasibility report to management.  
 (D) gives a software requirement specification (SRS).
12. In a data flow diagram, data flows cannot take place between  
 (A) two data stores  
 (B) two external entities  
 (C) a data store and an external entity  
 (D) Both (A) and (B)
13. Consider the decision table shown below. It is

|    | R <sub>1</sub> | R <sub>2</sub> | R <sub>3</sub> | R <sub>4</sub> |
|----|----------------|----------------|----------------|----------------|
| C1 | Y              | Y              | Y              | N              |
| C2 | Y              | N              | Y              | Y              |
| C3 |                |                | Y              | Y              |
| A1 | X              |                | X              |                |
| A2 |                | X              |                | X              |

- (A) an ambiguous decision table.  
 (B) a complete decision table.  
 (C) an incomplete decision table.  
 (D) Both (A) and (B)
14. Which of the following requirement specifications can be validated?  
 (S1): If the system fails during any operation, there should not be any loss of data.  
 (S2): Checking the hardware compatibility.  
 (S3): Defining a data interface.  
 (S4): Specification of response time for various functions.  
 (A) S1 and S2  
 (B) S2, S3 and S4  
 (C) S1, S3 and S4  
 (D) S1 and S4
15. Which of the following are true?  
 (i) A DFD should have loops.  
 (ii) A DFD should not have crossing lines.  
 (iii) Leveled DFD is easier to understand.  
 (iv) Context diagrams are not used in DFDs.  
 (A) (ii) and (i)  
 (B) (i) and (iv)  
 (C) (ii) and (iii)  
 (D) (iii) and (iv)

## Practice Problems 2

**Directions for questions 1 to 15:** Select the correct alternative from the given choices.

- Questionnaire consists of  
 (A) Forms  
 (B) Documents  
 (C) Qualitative data  
 (D) Quantitative data
- The method to obtain qualitative information is  
 (A) Background information  
 (B) Questionnaires  
 (C) Interviewing technique  
 (D) Journals and reports on similar systems
- Which among the following is a functional requirement?  
 (A) Description of all input data and their sources  
 (B) Capacity requirements  
 (C) Operating system available on the system  
 (D) Maintaining a log of activities
- The advantage of use case during requirement analysis phase is, it  
 (A) focuses on external behaviour only.  
 (B) focuses on internal behaviour only.  
 (C) focuses on additional behaviour.  
 (D) focuses on internal and external behaviour.
- Operational feasibility refers to  
 (A) technology needed is available and if available whether it is usable  
 (B) the proposed solution can fit in with existing operations  
 (C) the money spent is recovered by savings  
 (D) superior quality of products
- Software engineering is the application of  
 (A) Systematic approach of the development  
 (B) Quantifiable approach of the development  
 (C) Discipline approach of the development  
 (D) All of these
- The data flow model of an application mainly shows:  
 (A) The underlying data and the relationship among them  
 (B) Processing requirement and the flow of data  
 (C) Decision and control information  
 (D) Communication network structure
- DFD completeness is  
 (A) The process of discovering discrepancies between two or more sets of DFDs or discrepancies within a single DFD.  
 (B) The extent to which all necessary components of a data flow diagram have been included and fully decomposed.  
 (C) The conversation of inputs and outputs to a DFD process when that process is decomposed to a lower level.  
 (D) An iterative process of breaking the description of a system down into a finer and finer details, which creates a set of charts in which one on a given chart is explained in greater detail on another chart.
- The requirement analysis is performed in  
 (A) System design phase  
 (B) System development phase  
 (C) System analysis phase  
 (D) System investigation phase

10. In data flow diagram, an originator or receiver of data is usually designed by  
 (A) square box (B) circle  
 (C) rectangle (D) arrow
11. A feasibility document should contain all the following except  
 (A) project name  
 (B) problem description  
 (C) feasible alternative  
 (D) data flow diagrams
12. SRS document is \_\_\_\_\_ between customers and developers.  
 (A) legal contract  
 (B) standard  
 (C) request proposal  
 (D) None of the above
13. According to Brooks, adding more people to an already late software project makes it  
 (A) late  
 (B) fast  
 (C) does not impact schedule  
 (D) None of the above
14. The following is a quality metric:  
 (A) Correctness  
 (B) Maintainability  
 (C) Usability  
 (D) All of the above
15. Feasibility study should focus on  
 (A) Technical feasibility  
 (B) Economic feasibility  
 (C) Operational feasibility  
 (D) All of the above

**PREVIOUS YEARS' QUESTIONS**

1. What is the appropriate pairing of items in the two columns listing various activities encountered in a software life cycle? **[2010]**

|   |                      |   |                                      |
|---|----------------------|---|--------------------------------------|
| P | Requirements capture | 1 | Module development and integration   |
| Q | Design               | 2 | Domain analysis                      |
| R | Implementation       | 3 | Structural and behavioural modelling |
| S | Maintenance          | 4 | Performance tuning                   |

- (A) P-3, Q-2, R-4, S-1  
 (B) P-2, Q-3, R-1, S-4  
 (C) P-3, Q-2, R-1, S-4  
 (D) P-2, Q-3, R-4, S-1

2. Which one of the following is NOT desired in a good SRS document? **[2011]**  
 (A) Functional requirements  
 (B) Non-functional requirements  
 (C) Goals of implementation  
 (D) Algorithms for software implementation

**ANSWER KEYS**

**EXERCISES**

**Practice Problems 1**

1. D    2. D    3. D    4. D    5. B    6. C    7. B    8. A    9. B    10. A  
 11. C    12. D    13. C    14. B    15. C

**Practice Problems 2**

1. D    2. C    3. A    4. A    5. B    6. D    7. B    8. B    9. C    10. A  
 11. D    12. A    13. A    14. D    15. D

**Previous Years' Questions**

1. B    2. D

# Chapter 2

## Project Management and Maintenance

### LEARNING OBJECTIVES

- Project management
- Software design
- Modeling component level design
- SRS
- Software testing
- White-box testing
- Black box testing
- Implementation maintenance
- Software quality assurance
- Software Re-engineering
- COCOMO MODEL

### PROJECT MANAGEMENT

Project management is a technique used to ensure successful completion of a project by the project managers.

The functions included in project management are:

- Estimating resource requirements
- Scheduling tasks and events
- Providing training and site preparation
- Selecting qualified staff and supervising their work
- Monitoring the projects program
- Documenting
- Periodic evaluation
- Contingency planning

Project management involves planning, organization and control projects. It uses tools and software packages for planning and managing projects.

Project planning involves plotting project activities against time frame.

### PROJECT PLANNING TOOLS

- Tools used during software planning
- Helps the top level managers to take critical decisions during planning stage

### Gantt Charts

This activity scheduling method introduced in 1914 by Henry L. Gantt, uses horizontal bars to show the duration of actions or tasks.

The left end marks the beginning of the task and the right end its finish. Earlier tasks appear in the upper left and later ones in the lower right.

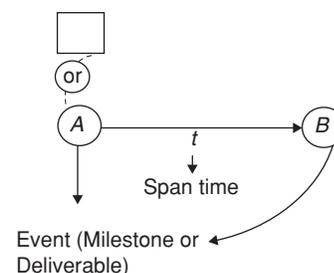
In real-life applications, an allowance for contingencies is provided. This is called **slack time**. Each project allows between 5 to 25 percent slack time for completion.

### Program Evaluation and Review Technique (Pert)

Gantt charts do not show precedence relationships among the tasks and milestones of a project.

A PERT chart is a project management tool used to schedule, organize and coordinate tasks within a project.

A PERT chart presents a graphic illustration of a project as a network diagram consisting of numbered nodes (either circles or rectangles) representing events, or milestones in the project linked by labelled vectors (directional lines) representing tasks in the project. The direction of the arrows on the lines indicates the sequence of tasks.

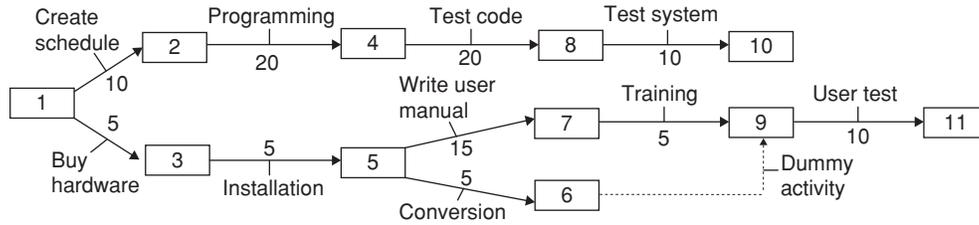


In the diagram, shown below the tasks between nodes 1, 2, 4, 8 and 10 must be completed in sequence and are called dependent or serial tasks. The tasks between nodes 1 and 2 and nodes 1 and 3 are not dependent on the completion of one to start the other and can be undertaken simultaneously. These tasks are called parallel or concurrent tasks. Tasks that must be completed in sequence but don't require resources or completion time are represented by dotted lines

with arrows and are called dummy activities (Example: dashed arrow linking 6 and 9).

Numbers on the opposite sides of the vectors indicate the time allotted for the task.

The PERT chart is preferred over Gantt chart since it clearly illustrates task dependencies. But on complex projects, PERT chart may be much more difficult to interpret.



Thus in short,

Dependency diagrams can be defined as a formal notation to help in the construction and analysis of complex schedules. Dependency diagrams are drawn as a connected graph of nodes and arrows. Dependency diagrams consists of three elements:

- Event—A significant occurrence in the life of a project.
- Activity—Amount of work required to move from one event to the next.
- Span time—Actual calendar time required to complete an activity.

## SOFTWARE DESIGN

Software design is the process in which requirements are translated into a blue print for constructing the software.

Once software requirements have been analyzed and modelled, software design is the last software engineering action within the modelling activity and sets the stage for construction (code generation and testing).

Architectural design defines the relationship between major structural elements of the software, the architectural

styles and design patterns, that can be used to achieve the requirements defined for the system, and the constraints that affect the way in which architecture can be implemented.

The interface design describes how the software communicates with systems that interoperate with it, and with humans who use it.

The component-level design transforms structural elements of the software architecture into a procedural description of software components.

The major goals of the design process are:

- The design must implement all of the explicit requirements contained in the requirements model, and it must accommodate all the implicit requirements, desired by stakeholders.
- The design must be a readable, understandable guide for those who test and subsequently support the software.
- The design should provide a complete picture of the software, addressing the data, functional and behavioral domains from an implementation perspective.

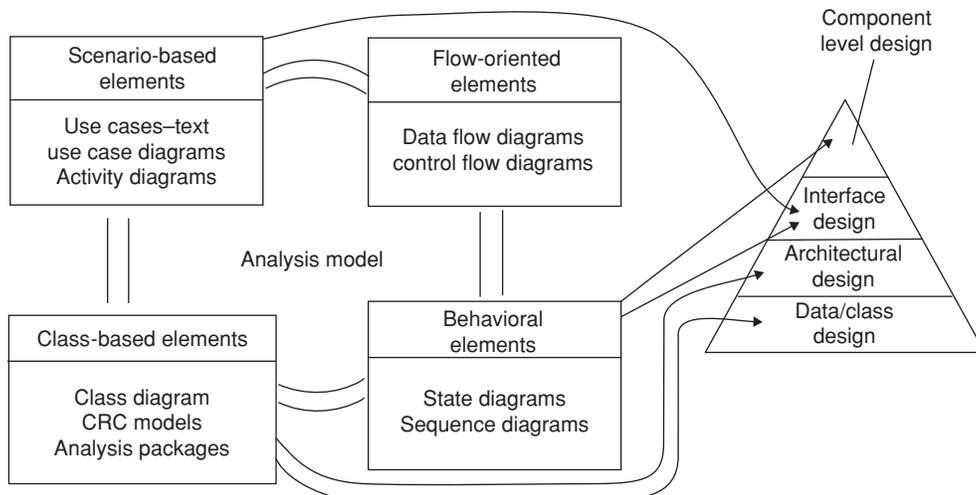


Figure 1 Design model.

Software design sits at the technical kernel of software engineering and is applied regardless of the software process model that is used. In the beginning, once the software requirements have been analyzed and modeled, software design is the last software engineering action within the modeling activity and sets the stage for construction (code generation and testing).

- The data/class design transforms analysis–class models into design class realizations and the requisite data structures required to implement the software.
- The architectural design defines the relationship between major structural elements of the software, the architectural styles and design patterns that can be used to achieve the requirements defined for the system.
- The interface design describes how the software communicates with systems that interoperate with it, and with humans who use it. An interface implies a flow of information (data/control) and a specific type of behavior.
- The component-level design transforms structural elements of the software architecture into a procedural description of software components.

## Design Concepts

Important software design concepts:

**Abstraction** Many levels of abstraction can be posed while considering a modular solution to any problem. At highest level of abstraction, a solution is stated in broad terms and at lower levels, a more detailed description of the solution is provided. At the lowest level of abstraction, the solution is stated in a manner that can be directly implemented.

**Architecture** Architecture is the structure or organization of program components (modules), the manner in which these components interact, and the structure of data that are used by components.

**Patterns** The intent of each design pattern is to provide a description that enables a designer to determine:

1. whether the pattern is applicable to current work.
2. whether the pattern can be reused.
3. whether the pattern can serve as a guide for developing a similar, but functionally or structurally different pattern.

**Separation of concerns** Separation of concerns is a design concept that suggests that any complex problem can be more easily handled if it is subdivided into pieces that can be solved and/or optimized independently.

A concern is a feature or behaviour that is specified as part of the requirement model for the software.

**Modularity** common manifestation of separation of concerns. Software is divided into separately named and addressable components (modules) that are integrated to satisfy problem requirements.

**Information hiding** Modules should be specified and designed so that information (algorithm and data) contained within a module is inaccessible to other modules that have no need for such information.

**Functional independence** Software should be designed in such a way that each module addresses a specific subset of requirements and has a simple interface when viewed from other parts of the program structure.

Functional independence is achieved by developing modules, which can perform a single function.

**Refinement** Refinement is a process of elaboration, begins with a statement or function defined at a high level of abstraction and then elaborates the original statement, providing more and more details as each successive refinement (elaboration) occurs.

**Refactoring** Refactoring is the process of changing a software system in such a way that it does not alter the external behaviour of the code (design), yet improves its internal structure.

When software is refactored, the existing design is examined for redundancy, unused design elements, inefficient or unnecessary algorithms, poorly constructed or inappropriate data structures or any other design failures that can be corrected to yield a better design.

## Modeling Component Level Design

Component level design occurs after the first iteration of architectural design has been completed. At this stage the overall data and program structure of the software has been established.

**Component** A component is a modular building block for computer software.

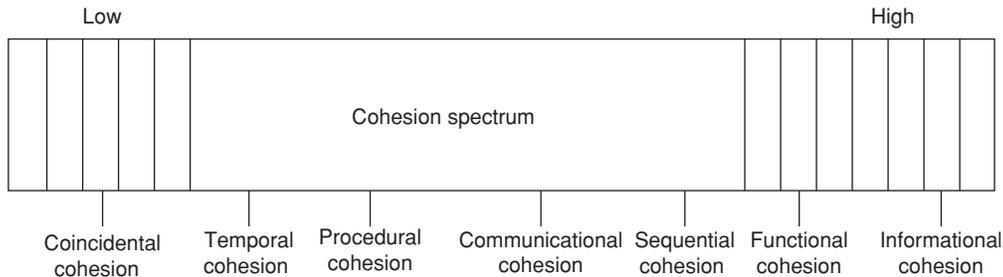
**Cohesion** Cohesion implies that a component or class encapsulates only attributes and operations that are closely related to one another and to the class or component itself. Cohesion is a measure of internal relative strength of a module. It should be more. Different types of cohesion are:

1. **Coincidental cohesion:** If elements of a module are unrelated, then it is coincidental cohesive.
2. **Logical cohesion:** If elements of a module are related, then it is logical cohesion.
3. **Temporal cohesion:** If the elements of a module are related and the elements are confined to initialization or time, it is temporal cohesion.
4. **Procedural cohesion:** If the elements are confined to one name and if they perform a set of operations, then the module is said to be procedural cohesive.
5. **Communicational cohesion:** If the elements in a module interact through data declared in it, then the module is said to be communicational cohesion.

- 6. **Sequential cohesion:** If the elements are related and if they perform a set of operations in which the output of one operation is the input for another operation.
- 7. **Functional cohesion:** If the elements are related and if they are confined to one name and if they perform

one and only one task, the module is functional cohesive.

- 8. **Informational cohesion:** If the elements of a module are confined to abstraction, it is informational cohesion.

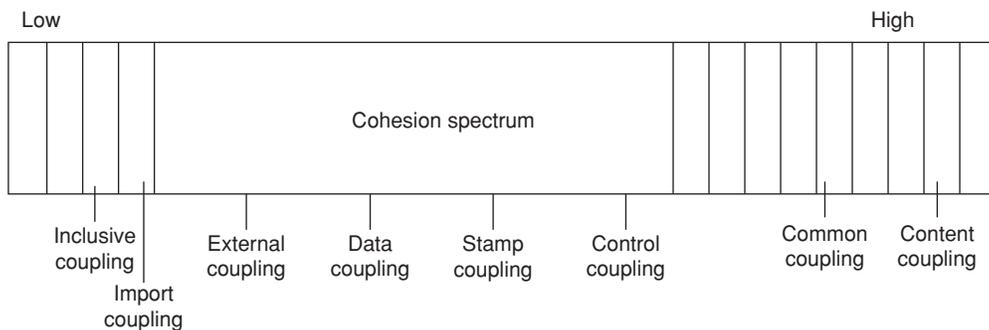


**Note:** Cohesion metric should be high.

**Coupling** Coupling is a qualitative measure of the degree to which classes are connected to one another. As classes and components become more interdependent, coupling increases. In component-level design coupling is to be kept as low as possible. It includes:

- 1. **Procedural or routine call coupling:** A form of coupling in which modules interact nominally more or less they are almost independent.
- 2. **Low coupling:** Form of coupling in which modules interact minimally. In extreme case there is no coupling between them.
- 3. **Inclusion coupling:** A coupling in which source code of one module is included into another module.
- 4. **Import coupling:** A coupling in which one module is declared in another module for its functionality.

- 5. **External coupling:** A coupling in which modules interact with modules written by some third party, which may include specific hardware or software.
- 6. **Data coupling:** Occurs when operations pass long strings of data arguments.
- 7. **Stamp coupling:** Occurs when a class is declared as a type for an argument of an operation of another class.
- 8. **Control coupling:** Coupling in which one module controls the order of execution of other module by using flags.
- 9. **Common coupling:** If the components make use of a global variable, it can lead to uncontrolled error propagation and unforeseen side effects when changes are made.
- 10. **Content coupling:** Type of coupling in when one module refers to other module, in extreme case, it changes internal structure of other modules for its functionality.



**Note:** Coupling metric should be low.

### CODING

Coding may be

- 1. The direct creation of programming language source code (e.g., Java, C).

- 2. The automatic generation of source code using an intermediate design like representation of the component to be built or
- 3. The automatic generation of executable code using a 'fourth generation programming language' (e.g., VC++).

The principles that guide the coding task are closely aligned with programming style, programming languages and programming methods.

The fundamental principles are:

- Understand the problem you are trying to solve.
- Understand basic design principles and concepts.
- Pick a programming language that meets the needs of the software to be built and the environment in which it will operate.
- Select a programming environment that provides tools that will make the work easier.

Create a set of unit tests that will be applied once the component code is completed.

### Characteristics of Good Srs

The characteristics of good SRS are

1. **Correctness:** The requirements specified in the software should meet, then the SRS is correct.
2. **Unambiguous:** The SRS is said to be unambiguous if every specified requirement can be interpreted in only one way.
3. **Completed:** The SRS is said to be complete, if and only if it has all significant requirements, definition of software responses to input data and labels and references to tables, figures and diagrams.
4. **Consistent:** The SRS is said to be consistent if the individual requirements are not defined in a conflict way and the SRS should be a high level document.
5. **Stability:** The SRS is said to be stable (or) ranked for the importance if each requirement has a preference. All the requirements may not have same importance; identify the requirements which are essential and requirements having least preference.
6. **Verifiable:** If each requirement is verifiable then the SRS is said to be verifiable.
7. **Modifiable:** The SRS is said to be modifiable, if the changes to the requirements can be made easily, consistent.
8. **Traceable:** Requirements should be clear so that each requirement can be referenced for enhancement, (or) future developments, which makes the SRS traceable.

**Validation of SRS** Validation of SRS is done to check whether the SRS is reflection of actual requirements and also to check the SRS documents is of good quality.

### Testing

Testing is the process of executing a program with the intent of finding an error.

A good test case is one that has a high probability of finding an as-yet-undiscovered error. A successful test is one that uncovers an as-yet-undiscovered error.

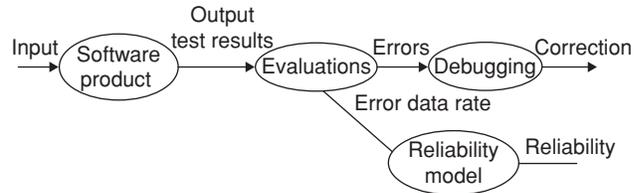


Figure 2 Formal technical review committee (FTR)

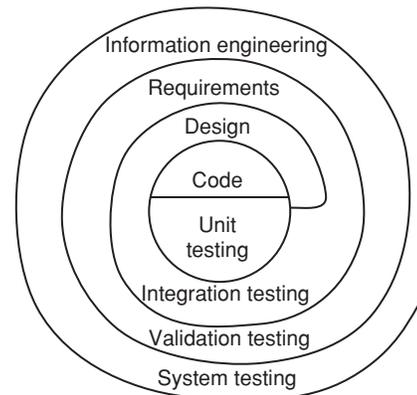
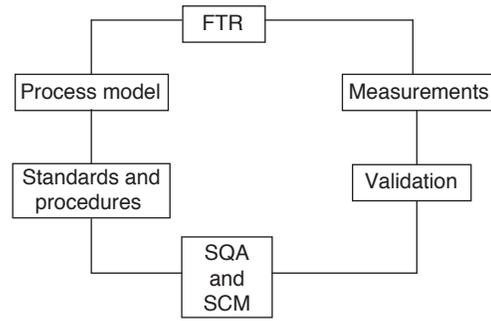


Figure 3 Verification

**There are four software testing strategies:**

1. Unit testing
2. Integration testing
3. Validation testing
4. System testing

### Unit testing

Unit testing concentrates on each unit (e.g., class, component, etc). Unit test focuses on the internal processing logic and data structures within the boundaries of a component. Important control paths are tested to uncover errors within the boundary of the module, using component-level design description as a guide.

### Integration testing

Integration testing focuses on design and construction of the software architecture. Integration testing is a systematic technique for constructing the software architecture while at the same time conducting tests to uncover errors associated with interfacing. The objective is to take unit-tested components and build a program structure that has been dictated by design.

- **Top-down integration** Modules are integrated by moving downward through the control hierarchy, beginning with the main control module (main program). Modules subordinate to the main control module are incorporated into the structure in either a depth-first or breadth-first manner.
- **Bottom-up Integration** begins construction and testing with the components at the lowest levels in the program structure.
- **Regression testing** in the context of an integration test strategy, regression testing is the re-execution of some subset of tests that have already been conducted to ensure that changes have not propagated unintended side effects.
- **Smoke testing** is an integration testing approach that is designed as a pacing mechanism for time-critical projects, allowing the software team to assess the project on a frequent basis.

### Validation testing

Validation succeeds when software functions in a manner that can be reasonably expected by the customer.

In validation testing, the requirements established as part of requirements modeling are validated against the software that has been constructed.

Software validation is achieved through a series of tests that demonstrate conformity with requirements.

Alpha and beta testing can be used to uncover errors that occur only at the end user.

The alpha test is conducted at the developer’s site by a representative group of end users. The software is used in a natural setting by end users in the presence of the developer and the developer records usage problems.

The beta test is conducted at one or more end user sites in the absence of developer. Therefore, beta test is a ‘live’ application of the software in an environment that cannot be controlled by the developer. The customer records all problems and reports to developer.

### System testing

In system testing, the software and other system elements are tested as a whole.

System testing is a series of different tests whose primary purpose is to fully exercise the computer-based system. The types of system tests used for software-based systems are:

- **Recovery testing** is a system test that forces the software to fail in a variety of ways and verifies that recovery is properly performed.
- **Security testing** attempts to verify that protection mechanisms built into a system will protect it from improper penetration.
- **Stress testing** executes a system in a manner that demands resources in abnormal quantity, frequency or volume. A variation of stress testing called sensitivity testing attempts to uncover data combinations within valid input classes that may cause instability or improper processing.

- **Performance testing** is designed to test the run-time performance of software within the context of an integrated system.
- **Deployment testing** also called configuration testing exercises the software in each environment in which it is to operate. It also examines all installation procedures and specialized installation software that will be used by customers, and all documentation that will be used to introduce the software to end users.

## SOFTWARE TESTING

The goal of testing is to find errors and a good test is one that has a high probability of finding an error.

The two ways of testing a software:

1. White-box testing (Internal testing)
2. Black-box testing (External testing)

### White-box Testing

In white-box testing (also called glass-box testing) of software, tests are conducted to ensure that all internal operations are performed according to specifications and all internal components have been adequately exercised.

White-box testing methods should guarantee that:

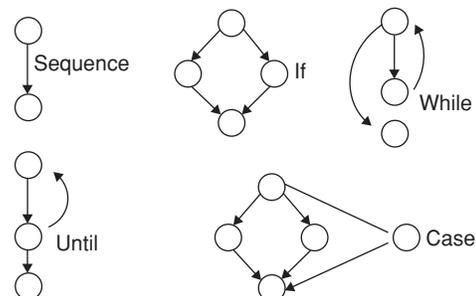
1. All independent paths, within a module are exercised at least once.
2. Exercise all logical decisions on their true or false sides.
3. Execute all loops at their boundaries and within their operational bounds and
4. Exercise internal data structures to ensure their validity.

### Basis path testing

Basis path testing is a white-box testing technique. This method enables the test case designer to derive a logical complexity measure of a procedural design and uses this measure as a guide for defining a basis set of execution paths. Test cases derived are guaranteed to execute every statement in the program at least one time during testing.

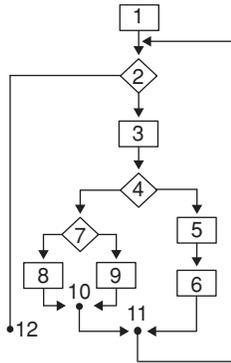
Flow graphs can be used for better understanding the control flow and thus helps basis path testing to execute every statement in the program at least once.

The flow graph symbols are:

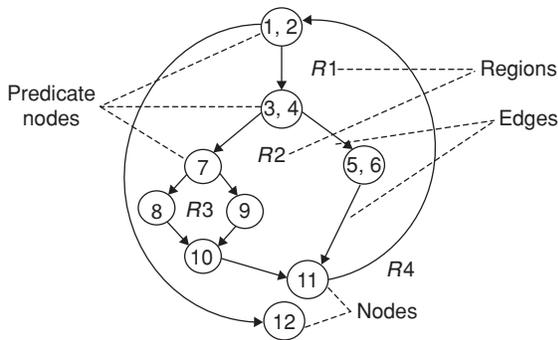


Each circle represents one or more non-branching PDL (Program Design Language) or source code statements.

**Example:  
Flowchart**



Corresponding flow graph is



Each node that contains a condition is called a predicate node. Independent paths (any path through the program that introduces at least one new set of processing statements or a new condition) in the above example are:

- Path 1: 1-2-12
- Path 2: 1-2-3-4-5-6-11-2-12
- Path 3: 1-2-3-4-7-8-10-11-2-12
- Path 4: 1-2-3-4-7-9-10-11-2-12

Thus if tests can be designed to force execution of these paths (a basis set), every statement in the program will have been guaranteed to be executed at least one time, and every condition will have been executed on its true and false sides.

Cyclomatic complexity is a software metric that provides a quantitative measure of the logical complexity of a program. When used in the context of basis path testing method, the value computed for cyclomatic complexity defines the number of independent paths in the basis set of a program and provides an upper bound for the number of tests that must be conducted to ensure that all statements have been executed at least once.

Complexity is calculated in one of the three ways:

1. The number of regions of the flow graph corresponds to the cyclomatic complexity. (i.e., four Regions R1, R2, R3, R4 in the above case)
2. Cyclomatic complexity  $V(G)$  for a flow graph  $G$  is defined as  $V(G) = E - N + 2$ , when  $E$  is the number of flow graph edges and  $N$  is the number of flow graph nodes (i.e., in the above case, there are 11 edges and 9 nodes. Thus  $V(G) = 11 - 9 + 2 = 4$ )

3. Cyclomatic complexity  $V(G)$  for a flow graph  $G$  is also defined as  $V(G) = P + 1$ , where  $P$  is the number of predicate nodes contained in the flow Graph  $G$ . In the above flow graph, there are 3 predicate nodes.

$$\therefore V(G) = 3 + 1 = 4$$

**Control structure testing**

Some of the variations on control structure testing to improve the quality of white-box testing are:

**Condition testing**

Condition testing is a test-case design method that exercises the logical conditions contained in a program module. This method focuses on testing each condition in the program to ensure that it does not contain errors.

**Control Structure Testing**

**Condition testing**

A simple condition is a Boolean variable or a relational expression, possibly preceded with one NOT ( $\neg$ ) operator. A compound condition is composed of two or more simple conditions, Boolean operators and parentheses. The possible types of elements in a condition include a Boolean operator, a Boolean variable, a pair of parentheses (surrounding a simple or compound Boolean condition), a relational operator, or an arithmetic expression.

**Dataflow testing**

This method selects test paths of a program according to the locations of definitions and use of variables in the program.

**Loop testing**

Loop testing is a white-box testing technique that focuses exclusively on the validity of loop constructs. Four classes of loop can be defined as:

**Simple loops**

The following set of tests can be applied to simple loops, where  $n$  is the maximum number of allowable passes through the loop.

1. Skip the loop entirely.
2. Only one pass through the loop.
3. Two passes through the loop
4.  $m$  passes through the loop where  $m < n$
5.  $n - 1, n, n + 1$  passes through the loop.

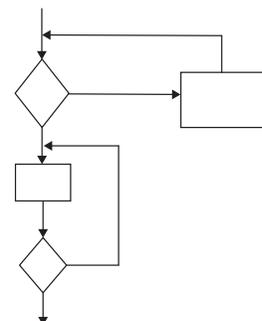


Figure 4 Simple loop.

### Nested loops

Here the number of possible tests grows geometrically as the level of nesting increases. This results in an impractical number of tests.

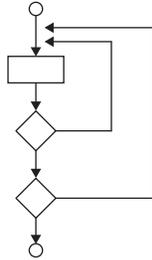


Figure 5 Nested loops.

### Concatenated loops

Concatenated loops can be tested using approach of simple loops, if each of the loops is independent of the other. If two loops are concatenated and the loop counter for loop 1 is used as the initial value for loop 2, then the loops are not independent.

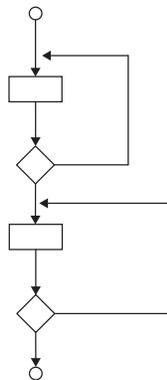


Figure 6 Concatenated loops.

### Unstructured loops

Whenever possible, this class of loops should be redesigned to reflect the use of the structured programming constructs.

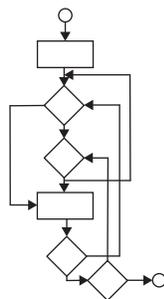


Figure 7 Unstructured loop.

### Black-box Testing

Black-box testing, also called behavioral testing, focuses on the functional requirements of the software.

Black-box testing attempts to find errors in the following categories:

1. Incorrect or missing functions
2. Interface errors
3. Errors in data structures or external database access
4. Behaviour or performance errors and
5. Initialization and termination errors

By applying black-box techniques, we derive a set of test cases that satisfy the following criteria:

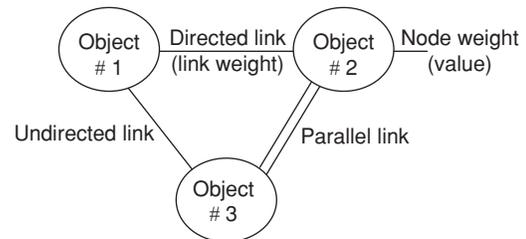
1. Test cases that reduce, by a count that is greater than one, the number of additional test cases that must be designed to achieve reasonable testing.
2. Test cases that tell something about the presence or absence of classes of errors, rather than an error associated only with the specific test at hand.

In graph-based black-box testing methods, software testing begins by creating a graph of important objects and their relationships and then devising a series of tests that will cover the graph so that each object and relationship is exercised and errors are uncovered.

### Graph-based testing methods

To accomplish these steps, the software engineer begins by creating a graph – a collection of nodes that represent objects; links that represent the relationships between objects; node weights that describe the properties of a node and link weights that describe some characteristic of a link.

The symbolic representation of a graph is as shown in the figure.



- Nodes are represented as circles connected by links that take a number of different forms.
- A directed link indicates that a relationship moves in only one direction.
- A bidirectional link (symmetric link) implies that the relationship applies in both directions.
- Parallel links are used when a number of different relationships are established between graph nodes.

### Equivalence partitioning

is a black-box testing method that divides the input domain of a program into classes of data from which test cases can be derived.

Equivalence partitioning strives to define a test case that uncovers classes of errors, thereby reducing the total number of test cases that must be developed.

Test case design for equivalence partitioning is based on an evaluation of equivalence classes for an input condition. An equivalence class represents a set of valid or invalid states for input conditions.

### Boundary value analysis (BVA)

It is developed as a testing technique used to test bounding values since a greater number of error occurring at the boundaries of the input domain than at the centre.

Boundary value analysis is a test case design technique that complements equivalence partitioning. Rather than selecting any element of an equivalence class, BVA leads to the selection of test cases at the ‘edges’ of the class. BVA derives test cases from the input conditions as well as from the output domain.

### Orthogonal array testing

The orthogonal array testing method is useful in finding region faults; an error category associated with faulty logic within a software component.

Orthogonal array testing can be applied to problems in which the input domain is relatively small.

When orthogonal array testing occurs, an Lg orthogonal array of test cases is created. This array has a ‘balancing property’, i.e., test cases are dispersed uniformly throughout the test domain.

### Model-based testing (MBT)

It is a black-box testing technique that uses information contained in the requirements model as the basis for the generation of test cases.

White-box testing is usually performed at the early stages of testing process, while black-box testing tends to be applied during later stages of testing.

## IMPLEMENTATION AND MAINTENANCE

### System Implementation

Implementation is the process of converting a new or a revised system design into an operational one. Major aspects of implementation are conversion, post-implementation review and software maintenance.

There are three types of implementations:

1. Implementation of a computer system to replace a manual system.
2. Implementation of a new computer system to replace an existing one.
3. Implementation of a modified application to replace an existing one using the same computer.

### Conversion

Conversion means changing from one system to another. The objective of conversion is to put the tested system into operation, while holding into costs, risks and personal irritation to a minimum.

It involves:

1. Creating computer-compatible files
2. Training the operating staff
3. Installing terminals and hardware

A very important aspect of conversion is not disrupting the functioning of the organization.

File conversion involves capturing data and creating a computer file from existing files.

### Post implementation review

Every system requires periodic evaluation after implementation. A post-implementation review measures the system’s performance against predefined requirements.

Unlike system testing, which determines where the system fails so that the necessary adjustments can be made, a post-implementation review determines how well the system continues to meet performance specifications. Post-implementation review is done after design and conversion are completed.

### Software Project Estimation

Software is the most expensive element of virtually all computer-based systems. For complex, custom systems, a large cost estimation error can make the difference between profit and loss.

Software project estimation is a form of problem solving, and in most cases, the problem to be solved (i.e., developing a cost and effort estimate for a software project) is too complex to be considered in one piece. For this reason, we decompose the problem recharacterizing it as a set of smaller problems.

### Problem-based estimation

Lines of code (LOC) and function point (FP) are used in two ways during software project estimation.

1. As an estimation variable to ‘size’ each element of the software.
2. As baseline metrics collected from past projects and used in conjunction with estimated variables to develop cost and effort projections.

The project planner begins by estimating a range of values of each information domain value. Using the historical data, the planner estimates an optimistic, most likely, and pessimistic size value for each function or count for each information domain value.

The expected value for the estimation variables is computed as

$$S = \frac{\text{optimistic} + 4 * \text{Most likely} + \text{pessimistic}}{6}$$

### Empirical estimation models

An estimation model for computer software uses empirically derived formulas to predict effort as a function of LOC

or FP. The model should be tested by applying data collected from completed projects, plugging the data into the model and then comparing actual to predicted results.

Some of the LOC-oriented estimation models are

|                                              |                                  |
|----------------------------------------------|----------------------------------|
| $E = 5.2 \times (\text{KLOC})^{0.91}$        | Walston-Felix model              |
| $E = 5.5 + 0.73 \times (\text{KLOC})^{1.16}$ | Bailey-Basili model              |
| $E = 3.2 \times (\text{KLOC})^{1.05}$        | Boehm simple model               |
| $E = 5.288 \times (\text{KLOC})^{1.047}$     | Doty model for $\text{KLOC} > 9$ |

**The software equation** The software equation is a multi-variable model that assumes a specific distribution of effort over the life of a software development project.

$$E = [\text{LOC} \times B^{0.333}/P]^3 \times (1/t^4)$$

where

$E$  = effort in person – months or person – years

$t$  = project duration in months or years

$B$  = Special spills factor

$P$  = Productivity parameter that reflects overall process maturity and management practices, the extent to which good software engineering practices are used, the level of programming languages used, the state of software environment, the skills and experience of the software team, and the complexity of the application.

**Note:**  $B$  increases slowly as ‘the need for integration, testing, quality assurance, and documentation and management skills grows’. For small programs  $\text{KLOC} = 5$  to  $15$ ,  $B = 0.16$ .

For programs greater than 70 KLOC,  $B = 0.39$

Putnam and Myers suggest a set of equations derived from the software equation.

Minimum development time is defined as

$$t_{\min} = 8.14 (\text{LOC}/P)^{0.43} \text{ in months for } t_{\min} > 6 \text{ months}$$

$$E = 180 B t^3 \text{ in person – months for } E \geq 20 \text{ person – months}$$

## Software Maintenance

Maintenance means restoring something to its original condition.

Maintenance is actually the implementation of the post-implementation review plan.

Maintenance is classified into corrective, adaptive or perfective maintenance.

Corrective maintenance repairs processing or performance failures or make changes because of previously uncorrected problems or false assumptions.

Adaptive maintenance means changing the program function.

Perfective maintenance enhances the performance or modify the programs to respond to the user’s additional or changing needs.

About 50–80% of the total system development cost accounts for maintenance. Analysts and programmers spend far more time maintaining programs than they do writing them.

A manufacturer wants to minimize the variation among the products that are produced by maintaining the quality.

User satisfaction = compliant product + good quality + delivery within budget and schedule.

## Software Quality Assurance (SQA)

Software Quality is defined as conformance to explicitly stated functional and performance requirements, explicitly documented development standards, and implicit characteristics that are points regarding quality is expected of all professionally developed software. In addition to the above definition some important

1. Software requirements are the foundation from which quality is measured. Lack of conformance to requirements is lack of quality.
2. If software conforms to its explicit requirements but fails to meet implicit requirements, software quality is suspect.

Activities performed by SQA group:

1. Prepares an SQA plan for a project.
2. Participates in the development of the project’s software process description
3. Reviews software engineering activities to verify compliance with the defined software process.
4. Ensures that deviations in software work and work products are documented and handled according to a documented procedure.
5. Records any non-compliance and reports to senior management.

## Software Reliability

Software reliability is defined as the probability of failure-free operation of a computer program in a specified environment for a specified time.

Measures of reliability and availability:

- A simple measure of reliability is mean-time-between-failure (MTBF).

$$\text{MTBF} = \text{MTTF} + \text{MTTR}$$

where

MTTF = mean-time-to-failure

MTTR = mean-time-to-repair

Although debugging (and related corrections) may be required as a consequence of failure, in many cases the software will work properly after a restart with no other change.

- In addition to a reliability measure, we must develop a measure of availability. Software availability is the probability that a program is operating according to requirements at a given point in time and is defined as

$$\text{Availability} = [\text{MTTF}/(\text{MTTF} + \text{MTTR})] \times 100\%$$

## Software Safety

- Software safety is a software quality assurance activity that focuses on the identification and assessment of potential hazards that may affect software negatively and cause an entire system to fail.

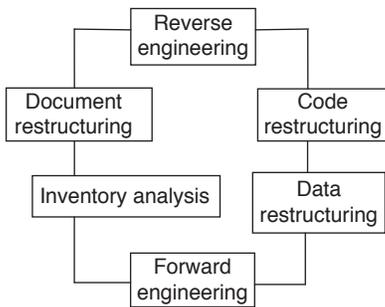
- Software safety examines the ways in which failures result in conditions that can lead to a mishap. That is the failures are evaluated in the context of an entire computer-based system and its environment.

### Software Reengineering

Cost of redevelopment is very high compared to development.

The maintenance of existing software can account for over 60% of all effort expended by a development organization, and the percentage continues to rise as more software is produced.

A reengineering process model is shown below:



- Reengineering takes time, costs significant amount of money and absorbs resources that might be otherwise occupied on immediate concerns.
- Reengineering of information systems is an activity that will absorb information technology resources for many years.
- Inventory analysis : The inventory can be nothing more than a spreadsheet model containing information that provides a detailed description of every active application. It should be revisited on a regular cycle.
- Document restructuring : It creates a framework of documentation that is necessary for the long-term support of an application.
- Code restructuring : The source code is analyzed using a restructuring tool. The restricted code is reviewed and tested to ensure that no anomalies have been introduced.
- Data restructuring : It is a full-scale reengineering activity. Current data architecture is dissected and necessary data models are defined.
- Forward engineering : Also called renovation or reclamation, covers design information from existing software and uses this information to alter or reconstitute the existing system in an effort to improve its overall quality.
- Reverse engineering : It is the process of analyzing a program in an effort to extract data, architectural, and procedural design information.

The abstraction level of a reverse engineering process and the tools used to affect it refers to the sophistication of the design information that can be extracted from source code.

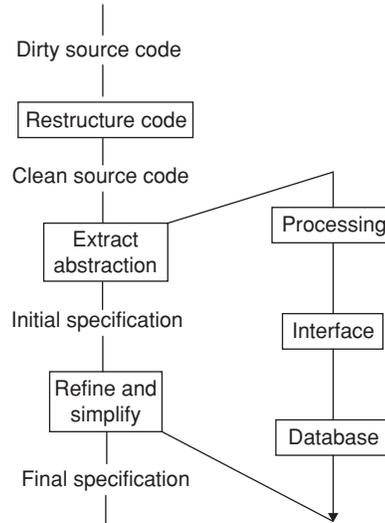


Figure 8 The reverse engineering process.

### COCOMO Model

One of the famous model structures used to estimate the software effort is the constructive cost model, which is often called as COCOMO model. COCOMO was developed by Boehm. The model helps in defining the mathematical relationship between the software development time, the effort in man-months and the maintenance effort.

Basic COCOMO is defined as computers software development effort (and cost) as a function of program size. Program size is expressed in estimated thousand lines of code (KLOC) COCOMO is applied to three classes of software projects:

1. Organic projects
2. Semi-detached projects
3. Embedded projects

#### Organic projects

Organic projects are projects that are having small teams with good working experience with less than rigid requirements.

#### Semi-detached projects

Semi-detached projects are projects with medium teams having mixed working experience with a mix of rigid and less than rigid requirements

#### Embedded projects

Project that are developed within a set of tight constraints (hardware, software, operational...)

The general formula of the basic COCOMO model is

$$E = a(s)^b$$

where

$E \rightarrow$  Represents effort in person-months

$S \rightarrow$  Size of the software development in KLOC

' $a$ ' and ' $b$ ' 5 Values dependent on the development mode

| Development Mode | Value of <i>a</i> | Value of <i>b</i> |
|------------------|-------------------|-------------------|
| Organic          | 2.4               | 1.05              |
| Semi-detached    | 3.0               | 1.12              |
| Embedded         | 3.6               | 1.20              |

Development time  $D = C(E)^d$

People required  $(P) = \frac{E}{D}$  [count]

| Development Mode | Value of <i>c</i> | Value of <i>d</i> |
|------------------|-------------------|-------------------|
| Organic          | 2.5               | 0.38              |
| Semi-detached    | 2.5               | 0.35              |
| Embedded         | 2.5               | 0.32              |

For intermediate COCOMO model, the value of coefficient *Q* and the exponent *b* are given in the table below:

| Development Mode | Value of <i>a</i> | Value of <i>b</i> |
|------------------|-------------------|-------------------|
| Organic          | 3.2               | 1.05              |
| Semi-detached    | 3.0               | 1.12              |
| Embedded         | 2.8               | 1.20              |

**EXERCISES**

**Practice Problems I**

**Directions for questions 1 to 15:** Select the correct alternative from the given choices.

**Common data for questions 1 and 2:** Consider the following payroll program that prints a file of employees and a file of information (transaction file) for the current month and for each employee.

In addition, the program updates the employee file, and produces an earnings report, a deduction report and analysis report. The application is capable of interactive command to print an individually requested pay slip. It also processes a file containing details of payment. This program can give printout of pay slips when they are requested individually. The weight table is shown below:

|                   | Simple | Average | Complex |
|-------------------|--------|---------|---------|
| No. of inputs     | 3      | 4       | 6       |
| No. of outputs    | 4      | 5       | 7       |
| No. of enquiries  | 3      | 4       | 6       |
| No. of files      | 7      | 10      | 15      |
| No. of interfaces | 5      | 7       | 10      |

- What is the unadjusted function point for the given payroll program?  
(A) 60 (B) 62  
(C) 68 (D) 72
- From the above problem, find adjusted function point where  $F4 = 4, F5 = 3, F12 = 2, F14 = 5$ ?  
(A) 49 (B) 62  
(C) 82 (D) 90

**Common data for questions 3 and 4:** The size estimated for software of a certain project is 45,000 lines of code. The average salary paid per engineer is ₹20,000 per month.

- Calculate the effort required if the software is of organic type.  
(A) 100 pm (B) 120 pm  
(C) 130 pm (D) 140 pm
- Calculate the cost required if the software is of semi-detached type.

- (A) 113000 (B) 213000  
(C) 315000 (D) 326515
- A 40 KDSI embedded program for teleprocessing is to be developed. Estimate the time required for the project using basic COCOMO model.  
(A) 12 pm (B) 14 pm  
(C) 16 pm (D) 18 pm

6. Consider the following code:

```
begin
If (x ≤ 0) then x = 0 - x;
a = x;
end
```

Lata wants to test the program with test data. What are the sufficient values to execute both branches of the decision box?

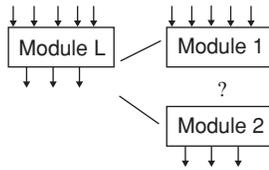
- (A)  $x = 0, 4$  (B)  $x = 0, -4$   
(C)  $x = 1, 4$  (D)  $x = 0, -1$
- What is the maintainability of a software with average number of days of repairing code is 10, adapting code is 20 and for enhancing code is 10?  
(A) 6.3 (B) 12.5  
(C) 32.6 (D) 40
  - Consider a Java program and the SLOC is given as 1000.

```
Class A
{
int x(int a);
int y(int b);
int z(int c);
}
```

What is the modularity?

- (A) 0.001 (B) 0.002  
(C) 0.003 (D) 0.004
- Raj has written a program to add two numbers. Assuming a 32-bit representation for an integer, to exhaustively test his program, the number of test cases required are  
(A)  $2^8$  (B)  $2^{16}$   
(C)  $2^{32}$  (D)  $2^{64}$

10. The module of the length 'L' is split up in two sub modules, module 1 and module 2, each of length  $\frac{L}{2}$ . How many links between the sub modules are allowed so that we maintain the value of information flow metric at same level?



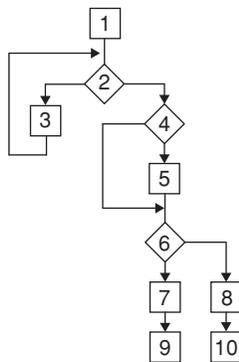
- (A) 2.4 (B) 3.6  
(C) 4.8 (D) 1.2
11. The three estimates of the code size for a particular application for geometric analysis were most optimistic 4600, most likely 6900, most pessimistic is 8600. The value of estimated size that should be taken is
- (A) 4600 (B) 6800  
(C) 6900 (D) 8600

12. For an application of developing new operating system the KLOC is 34.5. What is the number of person-month (effort) best estimated using the intermediate COCOMO model?
- (A) 126 (B) 130  
(C) 158 (D) 196
13. For a real-time software systems the KLOC is 28.2. What is the effort in person-month calculated by using basic COCOMO model?
- (A) 146 (B) 198  
(C) 220 (D) 248
14. For inventory management system the KLOC is 25.5, what is the effort in person-month, using basic COCOMO model?
- (A) 110 (B) 113  
(C) 120 (D) 140
15. For the above, what is the estimated project duration in months?
- (A) 6 (B) 8  
(C) 10 (D) 13

**Practice Problems 2**

*Directions for questions 1 to 16:* Select the correct alternative from the given choices.

**Common data for questions 1 and 2:** Consider the below flow graph:

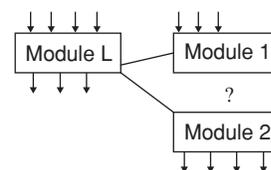


1. What is the number of paths to node 9?
- (A) 2 (B) 3  
(C) 4 (D) 5
2. What is the reachability measure?
- (A) 1.8 (B) 2.8  
(C) 2.4 (D) 2.1

**Common data for questions 3 and 4:** For a software project the estimation is carried out by the Delphi method. Below table shows 5 experts with estimates:

| Estimate | Pessimistic | Most likely | Optimistic |
|----------|-------------|-------------|------------|
| Expert 1 | 30          | 50          | 60         |
| Expert 2 | 10          | 55          | 75         |
| Expert 3 | 20          | 50          | 70         |
| Expert 4 | 30          | 60          | 70         |
| Expert 5 | 25          | 40          | 75         |

3. What is the average estimate?
- (A) 48.3 (B) 49.4  
(C) 50.8 (D) 56.7
4. What is the average variance?
- (A) 5.0 (B) 6.7  
(C) 7.8 (D) 8.3
5. The module of length L is split up into two sub modules (module-1 and module-2) each of length  $\frac{L}{2}$ . How many links between the sub modules exists so that we maintain the value of the information flow metric at the same level as found in the original module?



- (A) 3 (B) 4  
(C) 5 (D) 6

6. Constructive cost model is used to estimate
- Effort in man-month.
  - Effort and schedule based on the size of the software.
  - Size and duration based on the effort of the software.
  - None of these
7. The theoretic concept that will be useful in software testing is
- Hamiltonian circuit
  - Cyclomatic number
  - Eulerian cycle
  - None of these
8. Testing method that is normally used as the acceptance test for a software system is
- Regression testing
  - Integration testing
  - Unit testing
  - None of these
9. Acceptance testing is
- The manner in which each component functions with other component of the system are tested.
  - Running the system with given data by the actual user.
  - The process of testing the changes in a new system or an existing system.
  - None of these
10. Which of the following statements is true?
- Use of independent path testing criterion guarantees execution of each loop in a program under test more than once.
  - Validation is the process of evaluating software at the end of the software development to ensure compliance with the software requirements.
  - Statement coverage cannot guarantee execution of loops in a program under test.
  - None of these
11. The size estimated for software of a certain project is 40,000 lines of code. The average salary paid per engineer is ₹15,000 per month. Calculate the cost required if the software is of organic type.
- 1,60,000
  - 2,20,000
  - 7,90,000
  - 2,25,000
12. The size estimated for a software project is 35 Kloc. The average salary paid per engineer is ₹25,000 per month. Calculate the cost required if the software is of semi-detached type.
- 3,07,500
  - 3,17,500
  - 3,69,952
  - 2,45,000
13. Which of the following statements is false?
- The cyclomatic complexity of a module is the number of decisions in the module plus one where a decision is effectively any conditional statement in the module.
  - A direct flow of control in flow chart representing the lowest cyclomatic complexity.
  - The reasonable limit of the cyclomatic complexity measure is 10.
  - The cyclomatic complexity depends on the number of statements in the flowchart.
14. Which of the following is true regarding software testing?
- Software testing techniques are most effective if applied immediately after requirement specification.
  - Software testing techniques are most effective if applied immediately after design.
  - Software testing techniques are most effective if applied after coding.
  - Software testing methods are most effective if applied after integration.
- Common data for questions 15 and 16:** A software project involves execution of 4 activities  $A_1$ ,  $A_2$ ,  $A_3$ , and  $A_4$ , of duration 11, 7, 8 and 3 days respectively.  $A_1$  is the first one and needs to be completed before any other activity can commence. Activity  $A_2$  and  $A_3$  can be executed in parallel. Activity  $A_4$  cannot commence until both  $A_2$  and  $A_3$  are completed.
15. Find the critical path of the above project.
- $A_1 - A_2 - A_4$
  - $A_1 - A_3 - A_4$
  - $A_1 - A_2 - A_3 - A_4$
  - None of these
16. Find the slack time of the project.
- 0
  - 1
  - 12
  - 13

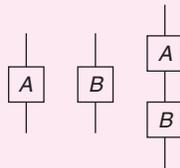
**PREVIOUS YEARS' QUESTIONS**

1. The coupling between different modules of a software is categorized as follows:
- I. Content coupling
  - II. Common coupling
  - III. Control coupling
  - IV. Stamp coupling
  - V. Data coupling

Coupling between modules can be ranked in the order of strongest (least desirable) to weakest (most desirable) as follows:

[2009]

- (A) I-II-III-IV-V                      (B) V-IV-III-II-I  
 (C) I-III-V -II-IV                    (D) IV-II-V -III-I
2. The cyclomatic complexity of each of the modules A and B shown below is 10. What is the cyclomatic complexity of the sequential integration shown below? [2010]



- (A) 19                                      (B) 21  
 (C) 20                                      (D) 10
3. A company needs to develop digital signal processing software for one of its newest inventions. The software is expected to have 40000 lines of code. The company needs to determine the effort in person-months needed to develop this software using the basic COCOMO model. The multiplicative factor for this model is given as 2.8 for the software development on embedded systems, while the exponentiation factor is given as 1.20. What is the estimated effort in person months? [2011]

- (A) 234.25                                (B) 932.50  
 (C) 287.80                                (D) 122.40

4. The following is the comment written for a C function.  
 /\* This function computes the roots of a quadratic equation  $ax^2 + bx + c = 0$ . The function stores two real roots in \* root 1 and \* root 2 and returns the status of validity of roots. It handles four different kinds of cases.
- (i) When coefficient 'a' is zero irrespective of discriminant.
  - (ii) When discriminant is positive.
  - (iii) When discriminant is zero.
  - (iv) When discriminant is negative.

Only in case (ii) and (iii), the stored roots are valid. Otherwise 0 is stored in the roots. The function returns 0 when the roots are valid and -1 otherwise.

The function also ensures  $root1 \geq root2$

```
int get_QuadRoots (float a, float b, float c, float *
root1, float * root 2); /*
```

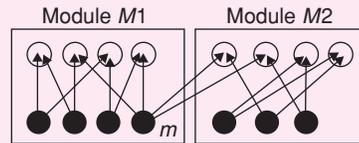
A software test engineer is assigned the job of doing black box testing. He comes up with the following test cases, many of which are redundant.

| Test Case | Input Set |       |      | Expected Output Set |       |              |
|-----------|-----------|-------|------|---------------------|-------|--------------|
|           | a         | b     | c    | root1               | root2 | Return value |
| T1        | 0.0       | 0.0   | 7.0  | 0.0                 | 0.0   | -1           |
| T2        | 0.0       | 1.0   | 3.0  | 0.0                 | 0.0   | -1           |
| T3        | 1.0       | 2.0   | 1.0  | -1.0                | -1.0  | 0            |
| T4        | 4.0       | -12.0 | 9.0  | 1.5                 | 1.5   | 0            |
| T5        | 1.0       | -2.0  | -3.0 | 3.0                 | -1.0  | 0            |
| T6        | 1.0       | 1.0   | 4.0  | 0.0                 | 0.0   | -1           |

Which one of the following options provide the set of non-redundant tests using equivalence class partitioning approach from input perspective for black-box testing? [2011]

- (A) T1, T2, T3, T6                      (B) T1, T3, T4, T5  
 (C) T2, T4, T5, T6                      (D) T2, T3, T4, T5

5. The following figure represents access graphs of two modules M1 and M2. The filled circles represent methods and the unfilled circles represent attributes. If method m is moved to module M2 keeping the attributes where they are, what can we say about the average cohesion and coupling between modules in the system of two modules? [2013]



- (A) There is no change  
 (B) Average cohesion goes up but coupling is reduced  
 (C) Average cohesion goes down and coupling also reduces.  
 (D) Average cohesion and coupling increase.

**Common data for questions 6 and 7:** The procedure given below is required to find and replace certain characters inside an input character string supplied in array A. The characters to be replaced are supplied in array 'oldc', while their respective replacement characters are supplied in array 'newc'. Array A has fixed length of five characters, while arrays 'oldc' and 'newc' contain three characters each.

However, the procedure is flawed.

```
void find_and_replace (char *A, char
*oldc, char *newc) {
for (int i = 0; i <5; i++)
for (int j=0; j<3; j++)
if (A[i] == oldc[j]) A[i] = newc[j];
}
```

The procedure is tested with the following four test cases.

1. oldc = "abc", newc = "dab"
  2. oldc = "cde", newc = "bcd"
  3. oldc = "bca", newc = "cda"
  4. oldc = "abc", newc = "bac"
6. If array *A* is made to hold the string "abcde", which of the above four test cases will be successful in exposing the flaw in this procedure? **[2013]**

- (A) None (B) 2 only  
(C) 3 and 4 only (D) 4 only

7. The tester now tests the program on all input strings of length five consisting of characters 'a', 'b', 'c', 'd' and 'e' with duplicates allowed. If the tester carries out this testing with the four test cases given above, how many test cases will be able to capture the flaw? **[2013]**

- (A) Only one (B) Only two  
(C) Only three (D) All four

8. In the context of modular software design, which one of the following combinations is desirable? **[2014]**

- (A) High cohesion and high coupling  
(B) High cohesion and low coupling  
(C) Low cohesion and high coupling  
(D) Low cohesion and low coupling

## ANSWER KEYS

### EXERCISES

#### Practice Problems 1

1. B    2. A    3. C    4. D    5. C    6. A    7. B    8. C    9. D    10. B  
11. B    12. C    13. B    14. B    15. D

#### Practice Problems 2

1. C    2. B    3. B    4. C    5. B    6. A    7. B    8. D    9. B    10. B  
11. D    12. C    13. D    14. B    15. B    16. B

#### Previous Years' Questions

1. A    2. A    3. A    4. C    5. A    6. C    7. B    8. B

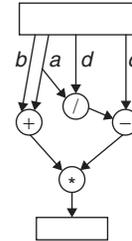
TEST

INFORMATION SYSTEM, SOFTWARE ENGINEERING

Time: 60 min.

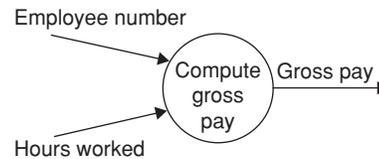
Directions for questions 1 to 30: Select the correct alternative from the given choices

- For a COCOMO model, organic projects are:
  - Projects having small teams with good experience, working with less than rigid requirements.
  - Projects having medium teams with mixed experience, working with more rigid requirements.
  - Projects developed with a set of tight constraints.
  - None of these
- Which of the following statements is true?
  - Basic COCOMO is good for quick estimate of software cost.
  - COCOMO applies to three classes of software projects; organic, semi-detached and embedded.
  - COCOMO does not account for differences in hardware constants, personal quality and experiences, etc.
  - All the above
- The first step in system analysis is
  - software requirement analysis.
  - software requirement specification.
  - system design.
  - information gathering.
- Questionnaire consists of
  - qualitative data.
  - quantitative data.
  - Either (A) or (B)
  - forms and documents.
- The assessment of an intangible benefit is
  - directly measurable.
  - done by discussion amongst users of information system.
  - irrelevant.
  - done by discussion amongst the developers.
- External entities in a DFD may be a
  - source of input data only.
  - destination of results only.
  - source of input data and destination of results.
  - data store.
- A context diagram
  - is a DFD which gives an overview of the system.
  - is a DFD that gives details of the system.
  - is not used in DFDs.
  - do not allow levelling of DFDs.
- Consider the DFD below; derive an expression from the given data:



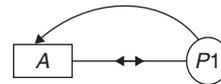
- $c - (a + b) * (a/d)$
- $a + b/d - (c * b)$
- $(a + b) * ((a/d) - c)$
- $(a + b) * (c + (a/d))$

9. Consider the following DFD:



- It calculates the gross pay.
- The process is specified incorrectly.
- Insufficient data flow.
- Data flow diagrams are not used to specify these kind of computations.

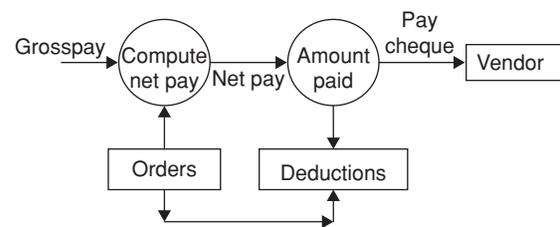
10.



Which of the following is correct for above DFD?

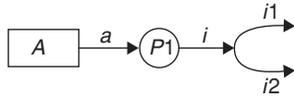
- The given DFD is correct.
- A DFD cannot have arrows pointing in opposite directions.
- Data cannot flow from external entity to a data store.
- Data cannot flow from a data store to an external entity.

11. Consider the given DFD. What is the mistake in the DFD?



- A data flow cannot connect two processing steps.
- A data flow cannot connect two distinct data stores.
- Data stores cannot communicate with a process.
- Data flow cannot connect two distinct external entities.

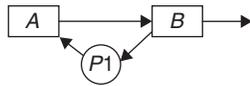
12.



The mistake in above DFD is

- (A) a data flow cannot be given two names.
- (B) a data flow that has crossing lines.
- (C) a DFD which forms loop.
- (D) there are no mistakes in the DFD.

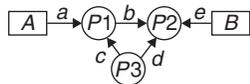
13.



Consider the above DFD. What is the mistake?

- (A) The DFD forms a loop here.
- (B) The DFD is correct.
- (C) DFD does not allow communication among two external entities.
- (D) DFD does not allow data flow among two data stores.

14. Consider the below DFD. What is the mistake?



- (A) Process P2 is not designed properly.
- (B) Process P1 is not designed properly.
- (C) Process P3 is not designed properly.
- (D) The external entities are not properly defined.

15. A good data flow diagram should have the following:

- (A) A process which is a pure decision
- (B) A DFD must be developed bottom up with higher levels giving more details
- (C) Data flow should not act as signals to activate or initiate process
- (D) All the above

16. The first phase of software development is

- (A) Requirements Analysis
- (B) Design
- (C) Coding
- (D) Testing

17. The lowest level of decomposition for a data flow diagram is

- (A) primitive DFD
- (B) unit DFD
- (C) context DFD
- (D) level 0 DFD

18. What is an important information while writing an SRS?

- (A) Nature of SRS
- (B) Characteristics of SRS
- (C) Environment of SRS
- (D) All of these

19. Which of the following is not an estimation metric for project size?

- (A) LOC
- (B) Function Point
- (C) Feature Point
- (D) None of the above

20. Human effort for developing a software project is measured in

- (A) Dollars
- (B) Person-Month
- (C) Refects
- (D) KLOC

21. Flight control software belongs to the following mode (as in basic COCOMO model):

- (A) Organic mode
- (B) Semi-detached mode
- (C) Embedded mode
- (D) None of the above

22. A transaction processing system with fixed requirements for terminal hardware and database software belongs to one of the following modes (in basic COCOMO model):

- (A) Organic mode
- (B) Semi-detached mode
- (C) Embedded mode
- (D) None of the above

23. In a software project, COCOMO is used to estimate

- (A) effort and duration based on the size of the software.
- (B) size, effort and duration based on the cost of the software.
- (C) size and duration based on the effort of the software.
- (D) effort and cost based on the duration of the software.

24. The maximum effort distribution in phases of software development is

- (A) Requirement analysis
- (B) Design phase
- (C) Coding
- (D) Testing

25. The minimum error distribution in the period of software development is in

- (A) Requirement analysis
- (B) Design phase
- (C) Coding
- (D) Testing

26. Basic Relation of COCOMO model is

- (A)  $E = (a * b) * (KLOC)$
- (B)  $E = a * (KLOC^b)$
- (C)  $E = a * (KDL) * b$
- (D)  $E = a / KLOC^b$

27. The extent to which the software can continue to operate correctly despite the input of invalid data is called as:

- (A) Reliability
- (B) Robustness
- (C) Fault-tolerance
- (D) Portability

28. Which of the following statement is false?
- (A) The data flow diagram is presented in hierarchical fashion.
  - (B) Data flow modeling is a core modeling activity in structured analysis.
  - (C) Data flow diagram is formal part of UML.
  - (D) Data flow modeling depicts control flow.
29. For which of the following practices does requirements engineering provide appropriate mechanisms and tools?
- (A) Analyzing need and validating the specification.
  - (B) Ambiguous specification of the solution.
  - (C) Risk Assessment.
  - (D) Implementing the system.
30. Which of the following is a common method of requirements elicitation?
- (A) Transactional Analysis
  - (B) Observation
  - (C) Practical considerations
  - (D) Web accessibility

**ANSWER KEYS**

- |       |       |       |       |       |       |       |       |       |       |
|-------|-------|-------|-------|-------|-------|-------|-------|-------|-------|
| 1. A  | 2. D  | 3. D  | 4. C  | 5. B  | 6. C  | 7. A  | 8. C  | 9. C  | 10. B |
| 11. B | 12. A | 13. C | 14. A | 15. C | 16. A | 17. A | 18. D | 19. D | 20. B |
| 21. C | 22. B | 23. B | 24. D | 25. A | 26. B | 27. B | 28. C | 29. A | 30. B |

*This page is intentionally left blank*

**Part C Software Engineering  
and Web Technology**

**Chapter 1:** Markup Languages

8.111

U

N

I

T

8

*This page is intentionally left blank*

# Chapter 1

## Markup Languages

### LEARNING OBJECTIVES

- HTML
- Structure of an HTML document
- Tags
- Attributes
- External linking
- Hyperlinks
- Breaks
- Tables
- Images and color
- Frames
- CSS
- XML

### HYPERTEXT MARKUP LANGUAGE (HTML)

HyperText markup language is a specialized markup language to create a web page. The language consists of ordinary text and special commands called tags.

HTML is not a formatting language. Rather it defines the parts of a document such as titles, headings, body text and block quotations. These parts are called elements. To define an element, tags are used. Tags give the browsers what they want to display on the web page.

#### Structure of an HTML Document

All HTML documents follow the same structure—a head which contains control information used by the browser and server and a body.

The body contains the content that displays on the screen and tags which control how that content is formatted by the browser.

```
<html>
 <head>
 <title> HTML document </title>
 </head>
 <body>
 <h1> Largest heading </h1>
 <p> A sample paragraph </p>
 <hr>
 </body>
</html>
```

- The entire document is surrounded by `<html> ..... </html>` which tell the software that it is now processing HTML.

- Format of our content should be according to the W3C recommendations.
- `<head>.....</head>` and `<body>.....</body>` tags are compulsory in all HTML documents.
- Programming languages include a mechanism called the comment that lets developers write plain text inside their code files. Comment tags start `<! .....>`. Each comment can contain as many lines of text as you like.
- If the comment runs over a number of lines, each must start and end with `--` and must not contain `--` within its body.

**Example:** `<! .....--`  
`-- ..... --`  
`-- ..... -->`

Comments can be placed in either the head or body of the document.

### TAGS

Tags are instructions in HTML that are embedded directly into the text of an HTML document. Tags, their attributes and values are enclosed between angular brackets '`<`' '`>`'. Tags that come in pairs have a start tag and an end tag. The slash mark is used to denote the end tag. All the text within the start tag and end tag is to be considered part of the element that the tag defines.

Tags are of two types—empty tag and container tag.

A formatted text document is composed of a set of elements such as paragraphs, headings and lists.

- A tag is a format name surrounded by angle brackets, end tags which switch a format off also contain a forward slash.
- Tags are delimited by angled brackets `<h1>`.

- Tags are not case sensitive. The following tags are equivalent: <HEAD>, <head> and <hEad>
- Styles must be switched off by an end tag.
- Some characters have to be replaced in the text by escape sequences if ‘<’ was not escaped, the software would attempt to process anything that followed it as part of a tag.
- White space, tabs and new lines are ignored by the browser; they can be used to make the HTML source more readable without affecting the way the page is displayed.
- Multiple white spaces are replaced by a single space while new lines and tabs are treated as spaces.
- If a browser doesn’t understand a tag it will ignore it.

### Container tags

Tags specified in pairs, delimiting text that will have some type of formatting is called a container tag. A container along with a companion tag, encloses the text to be formatted. The effect of a container tag is applied only to the text they contain. A container tag is also called a paired tag since they always appear as a pair. The general form of a container tag can be represented as:

```
<tag>
Text to be formatted
</tag>
```

The ‘<tag>’ is often called the opening tag and the ‘</tag>’ is called the closing tag. The closing tag will always have a slash ‘/’ to indicate the end of a tag. The opening tag activates the effect and the closing tag turns the effect off.

**Example:** <B> text </B>; the text will appear bold in the browser.

### Empty tag

An empty tag is a single tag representing some formatting commands in HTML. It will not have a companion tag and are hence called stand alone or singular tag.

**Example:** <BR>

The tag will insert a line break at the specified position.

### Attributes

HTML tags sometimes require additional information to be supplied to them. The additional information supplied to an HTML tag is known as attributes of a tag. Attributes are written immediately following the tag, separated by a space. Multiple attributes can be associated with a tag, also separated by a space.

**Example:** <FONT Face = “Arial” size = 12>

Welcome </FONT>

The face and size are attributes of the FONT tag.

- The document body encloses all the page formatting commands. The tags used to indicate the start and end of the main body of textual information are <BODY> and </BODY>. Page defaults like background colour, text colour, font size, etc. can be specified as attributes of the <BODY> tag.

### Attributes of BODY tag

Attributes	Descriptions
Bg colour	Changes the default background colour to the colour specified with this tag. The colour can be specified by name or equivalent hexa decimal number. Example: Bg colour = RED
Text	Changes the body text colour from its default value to the colour specified with this attribute. Example: text = green
Background	Specifies the name of the ‘GIF’ file that will be used as the back ground of the document. This tiles up across the page to give a back ground. Background = “br. gif”

### Text formatting tags

Tags	Descriptions
<P>	Paragraph break: The browser, moves onto a new line skipping one line between the previous line and the new line.
 	A line break is required when the text needs to start from a new line and not continues on the same line. It is an empty tag used to simply instruct the browser to start displaying the remaining text in a new line.
<CENTER>	Center tags are used to centre not only text but anything found between them, like texts, lists, rules, tables, etc.

### HR element attributes

The HR element has no ending tag.

- **ALIGN:**  
ALIGN = “\_\_\_\_\_”, sets the alignment of the line on the page to LEFT, RIGHT or CENTER  
The default is CENTER.  
The alignment has no purpose if the line width is 100%.
- **SIZE:**  
Sets the thickness or size of the line in pixels.
- **WIDTH:**  
Sets the width of the line across the page as a % (or) in pixels.

### Linking

Links are elements in a web page which can be selected by clicking on it. Linking is one of the most important features of HTML. Link allows to connect a text or an image to another web page or section of a web page. A link will

be displayed in a special way in a browser. Links will be highlighted with colours or underlines to indicate that it is a hyperlink. The target of a hyperlink can be another web page, another location on the current page, an image or any other computer file available in the server. Links can be classified into two:

1. External links
2. Internal links

### External linking

External linking refers to linking two documents. When a link in a web page is clicked a new document to which the hyperlink is linked will be opened. An external link points to another HTML document located anywhere in the www.

### Internal linking

Internal linking refers to linking different sections of the same document. When a link is clicked a different section in the same document will be displayed in the browser window.

### Text

The text on an HTML page can be altered in many ways. The actual font used can be changed to attempt to force the browser to use a specific font and the look of the text can be changed for emphasis.

- `<base font size = "n">`  
We can specify the minimum font size for basic text but not for headings. The size argument takes an integer from 1 to 7.
- `<font size = "[+/-]" color = "# rrggbb">` The colour of the text is set with the colour argument. This takes a hex value which represents the amounts of red, green and blue in the chosen colour.

#### Example:

```
<html>
 <head >
 <title> changing font sizes </title >
 </head>
 <body>
<h1> Font sizes </h1>
<base font size = "3">
<p> Here is some text in size 3
<p> Here is some
 Larger
 t
 e
 x
 t
</base font>
</body>
</html>
```

- Other Alternates are
  - `<b> ..... </b>` → Bold
  - `<i> ..... </i>` → Italic
  - `<strong > ..... </Strong >` used as a form of emphasis
  - `<tt > ..... </tt>` mono spaced font
  - `<sub> ..... </sub>` subscript
  - `<sup>..... </sup >` super script

### Break

`<br>`

Forces a line break within a passage of text where a paragraph is not desirable. On complex pages it is sometimes useful to put a `<br>` before and after tables, lists

- To display Escape sequences we need to use the following replacement sequences which always start with an ampersand ‘&’ and are terminated with a semicolon.
  - `& amp;` → &
  - `& nbsp;` → (white space)
  - `& lt ;` → <
  - `& gt;` → >
  - `& quot ;` → “
  - `& copy ;` → ©

### Hyperlinks

The benefit of hypertext is that it lets us create links within a document.

- Links should be used within documents where they either add to the understanding of the work or can be used to reduce download times.
- It is better to have many links to medium sized documents containing about a screenful of information rather than forcing readers to download a single massive document.

```

```

The link tag has 3 sections:

1. The address of the referenced document
  2. A piece of text to display as the link
  3. Closing tag
- The link text can be formatted using any of the text formatting options. Hypertext references, the ‘href’ part of the tag, can be
    1. links to documents or services at other internet sites
    2. links to documents within the same website
    3. links to a specific part of either the current page or another page.

#### Example:

```
 Next page
```

Links to another page in the same directory. The browser displays ‘Next page’ on the screen and highlights it so that readers know it is a hyperlink

#### Example:

```

some site
```

links to another website. This time some sight is displayed and highlighted.

## Lists

One of the most effective ways of structuring a website or its contents is to use lists. For Example, a commercial website may use pictures of its products instead of text in hyperlinks. These can be built as nested lists to provide an interesting graphical interface to the site.

HTML provides 3 types of lists:

1. Basic bulleted list
2. A numbered list
3. A definition list

Each has a different use but generally the definition list is the most flexible of the three lists.

### • Ordered and unordered lists:

```
....
```

The ordered and unordered lists are each made up of sets of list items. Elements of a list may be formatted with any of the usual text formatting tags and may be images or hyperlinks.

The closing tag is not part of HTML.

- ```
<ul [type = "disc"/>"square"/>"circle"] [compact] > ....</ul>
```


The basic unordered list has a bullet in front of each list item.

Everything between the tags must be encapsulated within

```
<li>...</li>
```

 tags.

- To minimize the amount of space that a list uses, we have to add the compact attribute
- ```
<ol [type = "1"|"a"|"A"|"I"|"i"]
```
- ```
[start = "n"] [compact] > .... </ol>
```

An ordered list has a number instead of a bullet in front of each list item. Different numbering schemes can be specified depending upon preference.

A list can number from any value that you desire. The starting value is given by the "start" attribute. All items in an ordered list must be enclosed within

```
<li>..... </li>
```

 tags

Tables

Tables have two uses:

1. Structuring pieces of information
2. Structuring the whole web page
 - Alternatively we can structure a page using frames or images.
 - A table is a grid of information such as, we might have seen in a ledger or spreadsheet.
Unlike a table from a spreadsheet the data items in an HTML table do not need to have any kind of relationship.
 - Most browsers struggle to process complex tables. The browsers are not optimized for tables and where tables are deeply nested on a page the browser may have difficulty displaying the page.
 - Web browsers have a layout engine which arranges the pieces before the web page is displayed.

- It is more difficult if table consists images where the size attribute of the image have not been set.
- ```
<table [align = "center"/>"left"/>"right"]
```

```
[border [= "n"]]
```

```
[cell padding = "n"] [width = "nn%"]
```

```
[cell spacing = "n"]> </table >
```
- Everything between 

```
<table>.....</table>
```

 tags will be part of table.  
These attributes control the formatting of the table as a whole, not that of the items in each cell.
- Tables can be aligned on the screen.
- A table can have a border, which includes a border between the cells. If the border attribute is not set, the table has no border.
- When the border attribute is set but a valid value is not given, a single pixel wide default border is drawn "cell padding" determines how much space is there between the contents of a cell and its border in pixels. Cell spacing sets the amount of white space between cells.
  - The "Width" attribute sets the amount of the screen that the table will use.
- ```
<tr [align = "left"/>"center"/>"right"]
```



```
[valign = "top"/>"center"/>"bottom"]>.... </tr>
```


Each row of the table has to be delimited by these tags. The row can be aligned horizontally and vertically within the table.
 - ```
<th [align = "left"/>"center"/>"right"]
```

```
[valign = "top" | "center" / "bottom"]
```

```
[nowrap] [colspan = "n"] [rowspan = "n"]> </th>
```
- These are table cells which are to be used for headings.
- The contents of the cell can be aligned vertically and horizontally within their row.
- If "nowrap" is set, the contents of the cell will not be automatically wrapped as the table is formatted for the screen.
- The "colspan" and "rowspan" attributes allow individual cells to be larger than a one by one grid.
- ```
<td [align = "left"/>"center"/>"right"]
```



```
[valign = " top"| "center"| "bottom"] [nowrap]
```



```
[colspan = n ] [ rowspan = n] > .... </td>
```


These describe the basic data cell.

Table elements

1.

```
<caption> string </caption>
```


This optional element is used to provide a string which describes the contents of the table. If used it must immediately follow the table element.
 - ```
<thead> </thead>
```

```
<tfoot> </tfoot>
```

```
<tbody> </tbody>
```

  
The rows in a table can be grouped into one of the three divisions.

- The idea is that the browsers will be able to scroll the tbody section of the table without moving either the thead or tfoot sections.
- When long tables extend over more than one page the information in thead and tfoot can be automatically replicated on each page.
- `<colgroup [span = 'n'] [width = 'n']> ..... </col group>`  
Columns within a table can be logically grouped together. Each group of columns can be assigned a default width which will apply to all columns.
- The span indicates the number of columns in the group.
- `<col [span = "n"] [width = "n"]> .... </col>`  
The attributes of individual columns are set using the "col" elements. The 'span' and 'width' attributes work in the same way as the 'colgroup' element.

**Example:**

```

<html>
<head>
 <title> A table </title>
</head>
<body>
<h1> A small table </h1>
<table align = "center" width = "75%"
 border = "1">
<caption> small table </caption >
<colgroup width = "30%" span = "2">
</colgroup>
<colgroup span = "3"> </col group>
<thead>
<tr> <td colspan = "5"> The table header
</td> </tr>
 </thead>
 <tbody>
 <tr>
<td>First </td>
<td> Second </td>
<td> Third </td>
<td> Fourth </td>
<td> Fifth </td>
 </tr>
 <tr>
<td>First </td>
<td> Second </td>
<td> Third </td>
<td> Fourth </td>
<td> Fifth </td>
 </tr>
 </tbody>
</table>
</body>
</html>

```

**Images and Colour**

- Colour can be used in a number of places on a web page; the background can be coloured, individual elements can be altered, and links which are already coloured can have their colours adjusted.
- To change the colours of links or the page background hexadecimal values are placed in the <body> tag.

```

<body bg color = "# nnnnnn" text = "# nnnnnn"
link = "#nnnnnn" vlink = "#nnnnnn"
alink = "#nnnnnn">

```

- The 'Vlink' attribute sets the colour of links visited recently, 'a link' the colour of a currently active link.
- The six figure hexadecimal values must be enclosed in double quotes and preceded by a hash (#).
- The colours of page elements can be altered by using the colour modifier. To change the colour of an individual heading we can use

```

<h2 color = "#a b a b a b"> Heading </h2> and within a
table the table headers could be coloured by:
<th bgcolor = "#a b a b a b">

```

- Images: If we want high quality, good compression and lots of colours use JPG, GIFs are more common as they tend to be smaller files and can be animated.
- `<body background = "URL"> ..... </body>`
- Sets the background of your page to use the given image. Images are tiled (repeated) to fill the available space by default.
- If we want to use a single image across the width of a page make it 1281 pixels wide then it cannot be tiled horizontally.

```

<img src = "URL" | "name" height = "n" width = "n" [alt
= "string"] [align = top] | "center"/ "bottom"]
usemap = "URL">

```

- Displays an inline image, that is an image which appears in the body of the text rather than on a page of its own or in a spawned viewer program.
- The height and width of the image, in pixels tell the browser how much space to allocate to an image when displaying a page.
- Some browsers use these to shrink/stretch images to fit.
- By default any text which follows an image will be aligned alongside its bottom edge. We can alter this so that the first line of text displays alongside the centre or top of the image.
- If we want a block of text shown next to an image we must use a table. To display an image without text, make it into paragraph.

```

<p align = "center"> <img src = "/mygif.gif"
alt = " mine"> </p>

```

The 'usemap' attribute is used in image mapping.

1. `<a href = "URL"> text message </a>`
2. `<a href = "URL"><img src = "filename"> </a>`

The first case uses an ordinary hypertext link but the URL should point to the image file, giving its name and type.

- In the second case we are using an image as the link to another image. This can be useful if we want to display a page of thumbnail images and allow the reader to choose which ones to view full-size. This is one way of speeding up the loading times of graphically intensive sites.
- An **image map** is a large picture which has areas that the reader can click with a mouse.

Each clickable area provides a hypertext link. The image map has 2 parts:

1. Image
2. Map

```

```

It tells the browser to display the source image and to map the second URL, the image map, onto it.

```
<area shape = "circle" | "rect" | "poly" | "default"
href = "URL" Coords = "string" alt = "string">
```

creates a clickable area on an image map. The 'alt' text in this case is displayed by the browser as an indicator for the reader of where the link goes.

- If we do not supply an 'alt', our image map is invalid and may not be displayed.
- The meaning of href should be clear, it is the destination of the link. The clickable area can have one of four shapes. Each shape is defined by coordinates, pairs of integers which give locations on the image in pixels.
- The default location does not require coordinates and is used to indicate what happens if the user clicks outside of the mapped areas.

Each image map can have only one default

- A **'rect'** has four coordinates which are paired. The first pair defines the top left corner and the second pair the bottom right corner of the area.
- A **'circle'** is defined by its centre and its radius centre is given by a pair of values and the radius by a single value. This requires just three values in the coordinate string.
- A **'polygon'** is made from a set of coordinates with the last pair listed being joined to the first to complete the shape.

The following example shows an image map with the mapping in the same file as the image link

```
<img src = "/mappicture.gif" usemap = "# main - map"
height = 30 width = 50>

<map name = "main - map">
<area shape = "rect" href = "/images/ img1.jpg"
alt = "Image one" cords = " 0,0,25,25">
<area shape = "rect" href = "/page1.html"
alt = "page one" coords = " 26,26,50,50">
<area shape = default href = "/page26.html"
alt = "page 26">
</map>

```

## Frames

If we want to represent a complex page structure and not confident about using a table to create it, then we can use frames. Frames are part of the HTML 4 specification. When we talk about frames what we refer to is a 'frameset' which is a special type of web page.

- The frameset page contains a set of references to HTML files, each of which is displayed inside a separate frame.
- All of the pages within a frameset are displayed inside the same browser window and can actually be made to appear to be a single page.
- Frame-based sites display more than one page at the same time, they can be complex to set up.
- Frame-based page is actually made from a set of documents, each displayed in its own frame. Each sub-documents can have its own scrollbars and can be loaded, reloaded, and printed.
- The tags that are needed are

```
<frameset [cols="%, %"] [rows = "%,%"]> ... </frame
set>
```

- This tag determines how the screen will be divided between the various frames we can have as many frames either vertically or horizontally as we want.
- Each frame has to be allocated a percentage of the screen.
- We can also nest framesets so that individual rows or columns can themselves be broken up into frames.

```
<frame [name = "name"] src = "filename"
[scrolling = "yes" | "auto" | "no"]
[frame border = "zero" | "1"]>
```

The src attribute works like an image source or a hyper-link address. It should point to a valid HTML file or image which can be displayed within the frame

- by setting the frame border attribute to "zero" stops it being displayed.
- <a href = " source" target = "n">
- to ensure that pages display in the correct frame we need to extend the basic address tag.

### Example:

```
<html>
<head>
<title> TIME pvt. Limited </title>
<frameset rows = "25% , 50%">
<frame name = "A" src = " ./ company. html">
< frame name = "B" src = " ./ orders.html" scrolling = "no">
</frameset>
</html>
```

## Forms

Forms are used to add an element of interactivity to a website. They are used to let the reader send information back to the server but can also be used to simplify navigation on complex websites.

```
<form action = "URL" method = "post" | "get" > ... </form>
```

- A form can contain virtually all other markup tags but cannot be nested within another form..
- The action attribute specifies the name and location of a CGI script that will be used to process the data. Data can be sent in 2 ways:
  1. post
  2. get
- We should use 'get' to retrieve information from a server and 'post' to send information to a server. The choice of approach is made by the 'method' attribute.
- 'post' is secured than 'get'.
- 'post' is capable of sending a wide variety of character sets but 'get' can only return ASCII data.
- 'post' is used to get data written in non-English languages.

```
<input type = "text" | "password" | "checkbox" | "radio" | "password" | "submit" | "reset" | "button" | "image" name = "string" [value = "string"] [checked] [size = "n"] [maxlength = "n"] [src = "URL"] [align = "top" | "bottom" | "middle" | "left" | "right"] >
```

Following are several types of input widgets:

- **Text** creates an input device up to size characters long and is able to accept up to max length characters as input.
- If value is set, that string will be used as the default text. These fields support only a single line of text. If we want to enter larger amount of text then use a "text area".
- **Password** works exactly like text but the input is not displayed to the screen. Each character is replaced by '\*'. The password is not encoded but is sent to the server as plain text.
- **Radio** creates a radio button. These are always grouped; buttons within a group should have the same name but different values.
 

The CGI script differentiates them by name + value.
- **Checkbox** produces a simple checkbox. It will be returned to the server as name = on if checked at submission.
- **Submit** creates a button which displays the value attribute as its text. It is used to send the data to the server.
- **Reset** also creates a button but this one is used to clear the form.
- **Image** can be used to place a picture on the page instead of a button.
- ```
<select name = "string"> ... </select>
```

It often very useful to have a list of items from which the user can choose.

The tag encloses a set of options and, when sent to the server, the name of the particular select tag and the name of the chosen option are returned.

- ```
<option value = "string" [selected]> ... </option>
```

 the 'select' statement will have several options from which the user can choose. The values will be displayed as the user moves through the list and the chosen one returned to the server.

- ```
<text area name = "string" rows = "n" cols = "n"> ... </text area>
```

 creates a free format plain text area into which the user can enter anything. The area will be sized at 'rows' by cols but will support automatic scrolling.

Example:

```
<html>
  <head>
    <title> my company </title>
  </head>
  <body>
    <h2 align = "center"> Feedback form </h2>
    <hr width = "60%">
    <form action = "http:// www.My company.Com / cgi - bin/ feed back.cgi" Method = "post">
      <p align = "left" > name : <input type = "text" Max length = "32" size = "16">
      <p align = "left"> Email Address:
      <input type = "text" max length = "32" size = "16">
      <p align = "left"> Location:
      <select name = "city" size = "1">
        <option value = "Hyderabad" selected>
          Hyderabad
        <option value = "Chennai"> chennai
        <option value = "Banglore"> Bangalore
        <option value = "Pune"> Pune
      </select>
      <p> comments:
      <br> <textarea name = "comments" rows = "6" cols = "40">
      </text area>
      <P align = "center"> <input type = "submit" name = "feed-back" value = "submit details">
    </form>
    <hr width = "60%">
  </body>
</html>
```

CASCADING STYLE SHEETS

One of the most important aspects of HTML is the capability to separate presentation and content. HTML does not have the facilities that are needed to cope with this diversity, but style sheets provide them.

- A style is simply a set of formatting instructions that can be applied to a piece of text.
- There are 3 mechanisms by which we can apply styles to our HTML documents:
 1. The style can be defined within the basic HTML tag.
 2. Styles can be defined in the <head> section and applied to the whole document.
 3. Styles can be defined in external files called style sheets which can then be used in any document by including the style sheet via a URL.

Styles can be cascaded. This means that formats override any, which were defined or included earlier in the document. We may include an external style sheet which redefines the h_1 tag, then write an h_1 style in the head of your page before finally redefining h_1 in the body of page. The browser will use the last of these definitions when showing the content in the following example the $\langle h_1 \rangle$ tag is redefined.

```
<html>
  <head>
    <title> simple style sheet </title>
    <style>
      <!--
      h1{
        color : black;
        border : thin groove;
        text-align : center ;
      }
      -->
    </style>
  </head>
  <body>
    <h1> simple style sheet </h1>
  </body>
</html>
```

Rules

A style rule has 2 parts, a selector and a set of declarations. The selector is used to create a link between the rule and the HTML tag.

The declaration has 2 parts, a property and a value.

- Selectors can be placed into classes so that a tag can be formatted in a variety of ways.
- Declarations must be separated using colons and terminated using semicolons.

Selector {property: value; property: value ...}

Classes

If we want to apply a style to some paragraphs we have to use classes.

Selector: classname {property : value ; property : value }
<selector class = classname>

- In the style sheet itself the rule is slightly modified by giving the style a unique name which is appended to the selector using a dot.
- In the HTML document when we want to use a named style the tag is extended by including

Class = and the unique name

```
h1 f {
  color :# a b a b a b ;
  background-color : # d 9a b 29;
  font-family : "Book Antiqua", Times, Serif;
  border : thin groove # 9 a b a b a ;
}
<h1 class = "f" > simple heading </h1>
```

Anonymous classes

Some times we want to apply a piece of formatting to many different elements within a page but not necessarily to the entire page.

- Cascading style sheets provides a way of defining styles within reusable classes.

Example:

```
<html>
  <head>
    <title> Anonymous classes </title>
    </style>
    <!--
      f {
        color : # a b a b a b ;
        background-color :# d9a b29;
        font-family : "Book Antiqua", Times, serif;
        border : thin groove #9ab aba;
      }
    -->
  </style>
</head>
<body>
  <h1 class = "f"> A simple Heading </h1>
  <P class = "f"> Applying the style f to a paragraph of text
</P>
</body>
</html>
```

Including style sheets:

```
< link rel = "stylesheet" href = "url" >
```

Type = "text / css" media = "screen">

- The href is a hyperlink to your style sheet, 'rel' tells the browser what type of link you are using.
- We have to tell the browser what type of document we are including, the type statement gives the relevant MIME type.

EXTENSIBLE MARKUP LANGUAGE (XML)

Extensible markup language (XML) is a way to apply structure to a web page. It provides a standard open format and mechanisms for structuring a document so that it can be exchanged and manipulated. Like HTML, XML uses tags to 'markup' data content. Unlike HTML, in XML you define your own tags that meet the exact needs of your document. The custom tags make data easier to organize and search. XML will not change the way your web page look, but it will change the way the documents are read and the way documents are filed and stored.

XML is a markup language. The term 'markup' is used to identify anything put within a document which either adds or provides special meaning. A mark up language is the set of rules. It declares what constitutes markup in a document and defines exactly what the markup means. It also provides a description of document layout and logical structure.

There are three types of markup:

1. Stylistic – How a document is presented.
2. Structural – How the document is to be structured.
3. Semantic – Tells about the content of the data.

In XML, the only type of markup that we are concerned with is structural.

An XML document must begin with the XML declaration statement. This statement alerts the browser or other processing tools that the document contains XML tags.

The declaration is-

```
<? XML version = "1.0" ?>
```

This is the first line of an XML document.

Tags

Tags carry the smallest unit of meaning signifying structure, format or style of the data. They are always enclosed with angular brackets, '<' and '>'. Tags are case-sensitive. This means that the tags <fruit>, <Fruit>, <FRUIT> carry different meanings and cannot be used interchangeably. All the tags must be paired so that they have a start tag and an end tag. For example, <fruit> and </fruit>.

Elements

Tags combined with data form elements. Elements are the building block of a document. An element consists of a start tag, an end tag and the content between them:

```
<fruit> orange is a fruit </fruit>
```

Attributes

An attribute gives information about an element. Attributes are embedded in the element start tag. An attribute consists of an attribute name and attribute value. The attribute name precedes its value and they are separated by an equal sign. Also the attribute value is enclosed in quotes to delimit multiple attributes in the same element.

Example: <fruit number = "6"> orange </fruit>

Document Type Definition

- Well formed XML documents are those which have tags, elements and attributes in a correct nesting structure without really providing further definitions.
- Valid XML documents are documents which follow a more formal structure.
- The main difference between well-formed XML and valid XML is the document type definition.

The document type definition (DTD) is a set of rules that define the elements that may be used and where they may be applied in relation to each other.

XML Parser

The process of taking a file and breaking it into components is called parsing.

The components are defined by a grammar, the rules of the language. Although this may be implied by the file structure rather than formally specified.

XML is not a simple data structure and cannot be handled with regular expressions for parsing. There are 4 parameters which can be used to categorize parsers:

1. Validating
2. Non-validating
3. Stream-based
4. Tree-based

Validating parser A validating parser uses both an XML file and a DTD to check that the XML adheres to the rules of the application.

If the XML breaks the rules by straying from the DTD then the parser will create an error and stop processing the files.

Non-validating parser These parsers will only use the XML document and are quite content if it is well-formed.

Stream-based parser A stream-based parser must read the entire document each time that an operation is requested and send a message to the controlling application when specific events occur.

Example: SAX

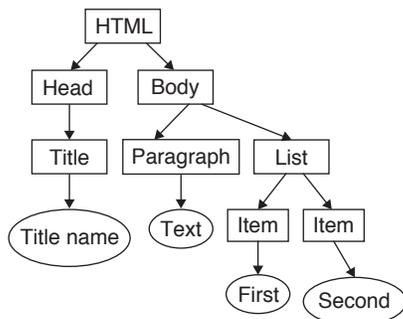
Tree-based parser This type of parser builds a static representation of the document which corresponds to the structure of the original XML.

Example: DOM

Document object model (DOM)

- The DOM is an application program interface (API) for XML documents.
- API is a set of data items and operations which can be used by developers of application programs.
- DOM API specifies the logical structure of XML documents and the ways in which they can be accessed and manipulated.
- DOM API is just a specification.
- DOM-Compliant applications include all the functionality needed to handle XML documents.
- They can build static documents, navigate and search through them, add new elements, delete elements and modify the content of existing elements
- The DOM views XML documents as trees, but this is very much a logical view of the document.
- Each node of the tree, each XML element, is modeled as an object.
- Each node encompasses both data and behavior and that the whole document can be seen as a single complex object.

- Sample DOM is shown below:



- DOM exposes the whole document to applications.

Namespaces

- A namespace is a way of keeping the names used by applications separate from each other.
- Within a particular namespace no duplicate names can exist.
- Applications may use many different namespaces at the same time.
- The implementation of namespaces is system dependent.
- XML developers can specify their own namespaces which can be used in many applications.
- A namespace can be included in the same way as a DTD.

Example:

```
<?xml version = "1.0"?>
<!DOCTYPE items SYSTEM "items.dtd" <
<!xml : namespace ns = "http : //URL/namespaces/jam"
prefix = "jam">
<?!xml : namespace ns = "http : //URL/namespaces/bread
prefix = "bread">
<items>
<item1>
<jam : name> kissan </jam : name>
</item1>
<item>
<bread : name> Roasted </bread : name>
</item>
</items>
```

Each item1 of items has a name element, But a namespace have been declared, so there is no chance of an application to get confused with the two names.

Attributes

- Attributes are important and useful when we are handling complexity.
- Some XML elements need to hold more than one piece of information.
- Some of these pieces are used to control the behaviour of the application. These are included as attributes.

Example:

```
<Quantity amount = "800" unit = "MC"> water </Quantity>
```

Here amount and unit are attributes of Quantity.

- The attributes of XML elements needs to be included in the DTD.
- Associated with the element declaration is an ATTLIST which may contain:
 - The name of the element
 - The name of each attribute
 - The data type of the attribute
 - Any value which will be used as a default if the attribute is omitted from the XML source.
 - Control information about the use of the element.

ATTLIST declaration

```
Example: <!ATTLIST Quantity amount CDATA #
REQUIRED unit CDATA "g">
```

The declaration shows an element with two attributes. The first one is 'amount', which is of CDATA type, means that it holds plain text which will not be passed through XML parser.

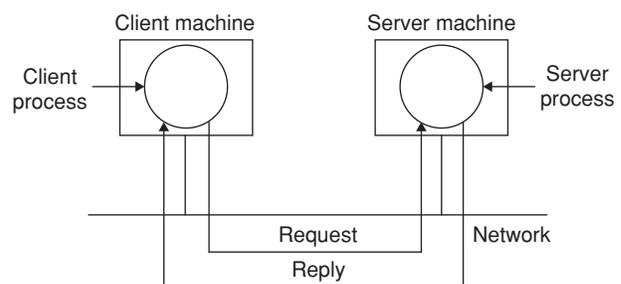
The attribute is REQUIRED which means that it must be included when the element is used. Failure of this will result in an error raised by the parser.

The second element 'unit' is optional which has a default value 'g'. If the value of this attribute is omitted then default value will be used.

CLIENT-SERVER COMPUTING

Client-server networking grew when personal computers (PC's) became the common alternative to older mainframe computers.

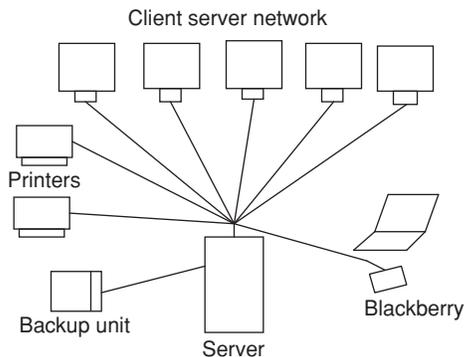
In client-server network, communication generally takes the form of a request message from the client to the server asking for some work to be done. The server does the work and sends back the reply, as shown below:



Usually, there are many clients using a small number of servers. Client devices are typically PC's with network software applications installed, that request and receive information over the network. Mobile devices as well as desktop computers can both function as clients.

A server device typically stores files and databases including more complex applications like websites. Server devices often feature higher powered central processors, more memory and larger disk drives than clients. One server

generally supports numerous clients and multiple servers can be networked together in a pool to handle the increased processing load as the number of clients grows.



Some of the most popular applications on the internet including email, FTP and web services follow client-server model. Each of these clients features a user interface (either graphic or text-based) and a client application that allows

the user to connect to servers. In the case of email and FTP, users enter a computer name (or an IP Address) into the interface to setup connections to the server.

The client-server model was developed to allow more users to share access to database applications. Compared to the mainframe approach, client-server offers improved scalability because connections can be made as needed rather than being fixed. The client-server model also supports modular applications (software applications are divided into modules) that can make creation of software easier.

Example: Users accessing banking services from their computer uses a web browser client to send a request to a web server at a bank. That program may in turn forward the request to its own database client program that sends a request to a database server at another bank computer at bank to retrieve the account information. The balance is returned to the bank database client, which in turn serves it back to the web browser client displaying the results to the user.

EXERCISES

Practice Problems I

Directions for questions 1 to 15: Select the correct alternative from the given choices.

- Which of the following statement is false?
 - HTML is a markup language for hypertext.
 - VML is used for freehand drawing in web page.
 - WML is wireless markup language used for micro procure of mobiles and palmtops.
 - None of these
- The web browser request goes to the server in
 - Hex form
 - ASCII form
 - Binary form
 - Text form
- The tag that contains information about the document including its title, scripts used, style definitions and documentation description is
 - <HTML>, </HTML>
 - <HEAD>, </HEAD>
 - <BODY>, </BODY>
 - <TITLE>, </TITLE>
- Tag that is considered to be illegal in XML is
 - <.document>
 - <document>
 - <\document>
 - None of these
- ___ specifies that a click in this area will not link anywhere.
 - NOHREF
 - HREF = 0
 - NULLHREF
 - HREF
- A _____ specifies the layout for frames, including the locations and characteristics of the frame.
 - frameset
 - border layout
 - table
 - frame border
- What is the correct syntax of the declaration, which defines the XML version?

- <xml version = "1.0"/>
- <?xml version = "1.0"?>
- <?xml version = "1.0"/>
- None of the above

- Which of the following are predefined attributes?
 - xml : lang
 - xml : space
 - both (A) and (B)
 - None of these
- Which of the following XML documents are well formed?
 - ```
<firstElement>
 TIME Hyderabad
</secondElement>
 Head Office
</secondElement>
</firstElement>
```
  - ```
<firstElement>
  TIME Hyderabad
</firstElement>
  Head Office
</secondElement>
<secondElement>
```
 - ```
<firstElement>
 TIME Hyderabad
<secondElement>
 Head Office
</firstElement>
</second Element
```
  - ```
</firstElement>
  TIME Hyderabad
</secondElement>
  Head Office
<secondElement>
</firstElement>
```

10. Which of the following XML fragments are well formed?
- (A) `<myExam mycity = "Hyderabad"/>`
 (B) `<myExam mycity = 'Hyderabad'/>`
 (C) `<myExam mycity = "Hyderabad/">`
 (D) `<myExam mycity = 'Hyderabad/>`
11. In the following anchor tag ` time ` which one is attribute?
- (A) A (B) HREF
 (C) time (D) http
12. A link to the document is like this:
` document `
 Then the link to proposal section will look like (from within same document)
- (A) ``
 (B) ``
 (C) ``
 (D) ``
13. DTD includes the specifications about the markup that can be used within the document, the specifications consists of all EXCEPT
- (A) the browser name
 (B) the size of element name
 (C) entity declarations
 (D) element declarations
14. Every website has a server process listening to TCP port 80 for incoming connections from clients (normally browsers). After a connection has been established the client sends one request and server sends one reply. Then the connection is released. The protocol that defines the legal requests and replies is called
- (A) TFTP (B) FTP
 (C) gopher (D) HTTP
15. Given below are several usages of the anchor tag in HTML:
- I. ` Hello `
 II. ` Hello `
 III. `Hello`
 IV. `Hello`
- Which of the above are valid?
- (A) I and II only (B) I, II, III and IV
 (C) I and III only (D) I, II and III only

Practice Problems 2

Directions for questions 1 to 15: Select the correct alternative from the given choices.

1. Which of the following is not case sensitive?
- (A) VML (B) XHTML
 (C) XML (D) HTML
2. Which of the following statement is false?
- (A) W3C stands for World Wide Web consortium.
 (B) W3C implements the `<layer>` tag.
 (C) W3C sets the HTML standards.
 (D) None of these
3. Which of the following requires a closing tag?
- (A) `<H1>` (B) `<ABBR>`
 (C) `` (D) All of these
4. Which of the following does not require a closing tag?
- (A) `` (B) ``
 (C) `
` (D) `<ABR>`
5. The smallest heading tag in HTML is
- (A) `<H0>` (B) `<H1>`
 (C) `<H6>` (D) `<H8>`
6. The largest size among the heading tags is
- (A) H_6 (B) H_5
 (C) H_7 (D) H_1
7. The ____ tag is effective for formatting program code or similar information, usually in a fixed font with ample space between words and lines.
- (A) `<Pre>` (B) `<Address>`
 (C) `<Blockquote>` (D) ``
8. ____ sets the text characteristics for the document.
- (A) ``
 (B) `<size>`
 (C) `<color>`
 (D) `<basefont>`
9. The tag used for creating a row in a HTML table is
- (A) `<TR>` (B) `<TD>`
 (C) `<Table row>` (D) `<TH>`
10. The SRC attribute is used to point to a ____ of the image.
- (A) folder (B) file
 (C) URL (D) pixel
11. How the position of files will be displayed in browser for the following code?
- ```
<frameset col's = "50%, 50%">
 <frameset rows = "50%, 50%">
 <frame src = "file1.html">
 <frame src = "file3.html">
 </frame set>
 <frame set rows = "50%, 50%">
 <frame src = "file2.html">
 <frame src = "file4.html">
 </frameset?>
</frameset>
```

- (A) 

50	50	
File 1	File 2	30
File 3	File 4	70
- (B) 

50	50	
File 1	File 3	30
File 2	File 4	70
- (C) 

70	30	
File 1	File 4	50
File 2	File 3	50
- (D) 

30	70	
File 1	File 4	50
File 3	File 2	50

12. 'We may have standalone attributes in XML'. This statement is
- (A) True  
 (B) False  
 (C) True if it is well-formed  
 (D) True if it defined in DTD

13. Standalone is one of the possible attributes in the XML declaration. We can set this to \_\_\_ if the document does not refer to any external entity.
- (A) Yes (B) No  
 (C) Not required to set it (D) None of these
14. Which of the following is not the difference between HTML and Java script?
- (A) HTML is used to create web pages, java script is used to customize the web pages.  
 (B) HTML provides security where as java script doesn't provide security.  
 (C) HTML is more preferable by the clients or users where as java script is not more preferable by the users.  
 (D) HTML is less efficient than java script.
15. <HMTL> and </HMTL> tags indicates the beginning and ending of the document which are compulsory because these indicate that the software is \_\_\_\_.
- (A) processing HTML (B) processing XML  
 (C) processing URL (D) deprocessing HTML

### PREVIOUS YEARS' QUESTIONS

1. Match the following: [2015]

(P) Condition coverage	(i) Black-box testing
(Q) Equivalence class partitioning	(ii) System testing
(R) Volume testing	(iii) White-box testing
(S) Alpha testing	(iv) Performance testing

- (A) P-ii, Q-iii, R-i, S-iv  
 (B) P-iii, Q-iv, R-ii, S-i  
 (C) P-iii, Q-i, R-iv, S-ii  
 (D) P-iii, Q-i, R-ii, S-iv

2. Which of the following statements is/are FALSE? [2015]

- I. XML overcomes the limitations in HTML to support a structured way of organizing content.  
 II. XML specification is not case sensitive while HTML specification is case sensitive.  
 III. XML supports user defined tags while HTML uses pre-defined tags.  
 IV. XML tags need not be closed while HTML tags must be closed.

- (A) II only (B) I only  
 (C) II and IV only (D) III and IV only

3. Consider the following C program segment.

```
while(first <= last)
{
 if (array[middle] < search)
 first = middle + 1;
 else if (array[middle] == search)
 found = TRUE;
 else last = middle - 1;
 middle = (first + last)/2;
}
if (first > last) notPresent = TRUE;
```

The cyclomatic complexity of the program segment is \_\_\_\_\_.

4. A software requirements specification (SRS) document should avoid discussing which one of the following? [2015]
- (A) User interface issues  
 (B) Non-functional requirements  
 (C) Design specification  
 (D) Interfaces with third party software
5. Consider the basic COCOMO model where  $E$  is the effort applied in person-months,  $D$  is the development time in chronological months, KLOC is the estimated number of delivered lines of code (in thousands) and  $a_b, b_b, c_b, d_b$  have their usual meanings. The basic COCOMO equations are of the form [2015]
- (A)  $E = a_b(\text{KLOC}) \exp(b_b), D = c_b(E) \exp(d_b)$   
 (B)  $D = a_b(\text{KLOC}) \exp(b_b), E = c_b(D) \exp(d_b)$   
 (C)  $E = a_b \exp(b_b), D = c_b(\text{KLOC}) \exp(d_b)$   
 (D)  $E = a_b \exp(d_b), D = c_b(\text{KLOC}) \exp(b_b)$
6. Which one of the following statements is NOT correct about HTTP cookies? [2015]
- (A) A cookie is a piece of code that has the potential to compromise the security of an Internet user.  
 (B) A cookie gains entry to the user's work area through an HTTP header.  
 (C) A cookie has an expiry date and time.  
 (D) Cookies can be used to track the browsing pattern of a user at a particular site.
7. Which one of the following assertions concerning code inspection and code walkthrough is true? [2015]

- (A) Code inspection is carried out once the code has been unit tested.
- (B) Code inspection and code walkthrough are synonyms
- (C) Adherence to coding standards is checked during code inspection
- (D) Code walkthrough is usually carried out by independent test team

8. Consider a software project with the following information domain characteristics for calculation of function point metric.

Number of external inputs ( $I$ ) = 30

Number of external outputs ( $O$ ) = 60

Number of external inquiries ( $E$ ) = 23

Number of files ( $F$ ) = 08

Number of external interfaces ( $N$ ) = 02

It is given that the complexity weighting factors for  $I$ ,  $O$ ,  $E$ ,  $F$  and  $N$  are 4, 5, 4, 10 and 7, respectively. It is also given that, out of fourteen value adjustment factors that influence the development effort, four factors are not applicable, each of the other four factors have value 3, and each of the remaining factors have value 4. The computed value of function point metric is \_\_\_\_\_ [2015]

9. In a web server, ten WebPages are stored with the URLs of the form `http://www.yourname.com/var.html`; where, var is a different number from 1 to 10 for each Webpage. Suppose, the client stores the Webpage with  $var = 1$  (say  $W_1$ ) in local machine, edits and then tests. Rest of the WebPages remains on the web server.  $W_1$  contains several relative URLs of the form 'var.html' referring to the other WebPages. Which one of the following statements needs to be added in  $W_1$ , so that all the relative URLs in  $W_1$  refer to the appropriate WebPages on the web server? [2015]

- (A) `<a href: http://www.yourname.com/"', href: "var.html">`
- (B) `<base href: http://www.yourname.com/">`
- (C) `<a href: http://www.yourname.com/">`

(D) `<base href: http://www.yourname.com/"', range: "var.html">`

10. Consider a software program that is artificially seeded with 100 faults. While testing this program, 159 faults are detected, out of which 75 faults are from those artificially seeded faults. Assuming that both real and seeded faults are of same nature and have same distribution, the estimated number of undetected real faults is \_\_\_\_\_ [2015]

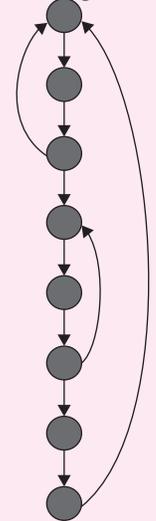
11. Consider three software items: Program-X, Control Flow Diagram of Program-Y and Control Flow Diagram of Program-Z as shown below [2015]

**Program-X:**

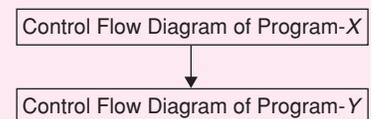
```
Sumcal (int maxint, int value)
{
 int result=0, i=0;
 if (value < 0)
 {
 Value = -value;
 }
 While ((i<value) AND (result <=
 maxint)
 {
 i=i+1;
 result = result + 1;
 }
 if (result <= maxint)
 {
 printf (result);
 }
 else
 {
 printf ("large")
 }
 printf ("end of program");
}
```

**Control Flow**

**Diagram f Program-Y:**



**Control Flow Diagram of Program-Z:**



The value of McCabe's Cyclomatic complexity of program-X, Program-Y, and Program-Z respectively are

- (A) 4, 4, 7
- (B) 3, 4, 7
- (C) 4, 4, 8
- (D) 4, 3, 8

## ANSWER KEYS

### EXERCISES

#### Practice Problems 1

- 1. D
- 2. B
- 3. B
- 4. A
- 5. A
- 6. A
- 7. B
- 8. C
- 9. A
- 10. A
- 11. B
- 12. A
- 13. A
- 14. D
- 15. D

#### Practice Problems 2

- 1. D
- 2. B
- 3. D
- 4. C
- 5. C
- 6. D
- 7. A
- 8. D
- 9. A
- 10. C
- 11. A
- 12. B
- 13. A
- 14. C
- 15. A

#### Previous Years' Questions

- 1. C
- 2. C
- 3. 5
- 4. C
- 5. A
- 6. A
- 7. C
- 8. 612 to 613
- 9. B
- 10. 28
- 11. A

## TEST

## WEB TECHNOLOGIES

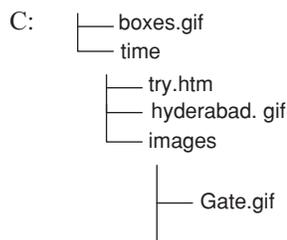
Time: 60 min.

**Directions for questions 1 to 30:** Select the correct alternative from the given choices.

- The \_\_\_\_\_ tag contains information about the document including its title, scripts used, style definitions and documentation descriptions.
  - <Head>
  - <Body>
  - <HTML>
  - <Title>
- Is it easier to process XML than HTML?
  - Yes
  - No
  - Sometimes
  - Can't say
- Which of the following tags is the smallest heading tag?
  - <H1>
  - <H6>
  - <H0>
  - <H8>
- In XML
  - the internal DTD subset is read before the external DTD
  - the external DTD subset is read before the internal DTD
  - there is no external type of DTD
  - there is no internal type of DTD
- Attribute standalone = "no" should be included in XML declaration if a document
  - is linked to an external XSL style sheet
  - has external general references
  - has processing instructions
  - has an external DTD
- Which of the following XML fragments are well-formed?
  - <myEle myAtt = 'val1 = val2' />
  - <myEle myAtt = 'val1>val2' />
  - <myEle myAtt = 'val1>val2' />
  - None of the above
- Parameter entities can appear in
  - XML file
  - DTD file
  - XSL file
  - Both (A) and (B)
- The use of a DTD in XML document is
  - required when validating XML document
  - no longer necessary after the XML editor has been customized
  - used to direct conversion using an XSLT processor
  - a good guide to populating a template to be filled in when generating an XML document automatically.
- To add the attribute named 'Branch' to the <Time> tag the syntax will be
  - <Time attribute Branch = "DSNR">
  - <Time Branch attribute = "DSNR">
  - <Time Branch = "DSNR">
  - None of these
- The syntax for parameter entity is
  - <!ELEMENT % NAME DEFINITION>
  - <!ENTITY % NAME DEFINITION>
  - <!ENTITY \$ NAME DEFINITION>
  - <ENTITY % NAME DEFINITION>
- A schema can be named using the name attribute like
  - <schema attribute = "schema1">
  - <schema nameattribute = "schema1">
  - <schema nameattri = "schema1">
  - <schema name = "schema1">
- A schema describes
  - grammar
  - vocabulary
  - structure
  - data type of XML document
  - (i) and (ii)
  - (iii) and (iv)
  - Both (A) and (B)
  - None of the above
- The XML DOM object is
  - an attribute
  - entity reference
  - comment reference
  - comment data
- The default model for COMPLEXTYPE, in XML schema for element is
  - text only
  - element only
  - no default type
  - Both (A) and (B)
- To create a choice in XML schemas, we use
  - <xsd: select> element
  - <xsd: multi> element
  - <xsd: choice> element
  - <xsd: single> element
- To bind the HTML element <INPUT> Type in text with the data source "dsoCustomer" we use
  - <INPUT TYPE = "TEXT" DATAFIELD = "#dsoCustomer">
  - <INPUT TYPE = "TEXT" DATASRC = "dsoCustomer">
  - <INPUT TYPE = "TEXT" DATASRC = "#dsoCustomer">
  - <INPUT TYPE = "TEXT" DATAFLD = "\*dsoCustomer">
- The attribute used to define a new namespace is
  - XMLN Space
  - Xml Name Space
  - Xmlns
  - XmlNs
- Identify the most accurate statement about the application of XML:
  - XML must be used to produce XML and HTML output.
  - XML can not specify or contain presentation information.

- (C) XML is used to describe hierarchically organized information.
  - (D) XML performs the conversion of information between different e-business applications.
19. What is an XML namespace?
- (A) A set of names applied to specific spaces within an XML document, such as the head and body.
  - (B) A set of names representing a specific XML vocabulary.
  - (C) Both (A) and (B)
  - (D) None of these
20. The standard model for network application is
- (A) producer consumer model
  - (B) node-node model
  - (C) system-system model
  - (D) client-server model
21. Which is/are example of service(s) that a server can provide?
- (A) Return the time-of-day to the client.
  - (B) Print a file for the client
  - (C) Execute a command for the client on the server's system
  - (D) All of the above
22. Print a file, read or write a file for client are handled in a
- (A) iterative fashion
  - (B) concurrent fashion
  - (C) Both (A) and (B)
  - (D) None of these

**Common data for questions 23 and 24:** Consider working with the file try.htm and provided the following directory structure:



23. The tag specification to use boxes.gif file as back-ground will be
- (A) <Body BACKGROUND = “/boxes.gif”>
  - (B) BODY BACKGROUND = “/boxes.gif
  - (C) <BODY BACKGROUND = “boxes.gif”>
  - (D) <BODY BACKGROUND = “#boxes.gif”>
24. The tag specification to use gate.gif file as back ground will be
- (A) <BODY BACKGROUND = “../gate.gif”>
  - (B) <BODY BACKGROUND = “.gate.gif”>
  - (C) <BODY BACKGROUND = “images/gate.gif”>
  - (D) None of these.
25. Consider below HTML code:
- (1) <HTML>
  - (2) <FRAMESET ROWS = “30%, \*”>
  - (3) <FRAMESET COLS = “50%, 50%”>
  - (4) <FRAME Src = “file1. html”>
  - (5) <FRAME Src = “File2.html”>
  - (6) </FRAMESET>
  - (7) <FRAMESET COLS = “50%, 50%”>
  - (8) <FRAME Src = “File3.html”>
  - (9) <FRAME Src = “File4.HTML”>
  - (10) </FRAMESET>
  - (11) </FRAMESET>
  - (12) </HTML>
- We can embed <BODY> tag at line number.
- (A) 2
  - (B) 11
  - (C) Anywhere
  - (D) No where
26. What is the purpose of <area> tag?
- (A) It defines an area inside a table.
  - (B) It defines an area inside an image map.
  - (C) It defines the area of total HTML code.
  - (D) Both (A) and (B)
27. Which of the following tag is used to define a section in a document?
- (A) <frame>
  - (B) <li>
  - (C) <div>
  - (D) <wbr>
28. Which of the following is not a syntax rule of an XML document?
- (A) All XML tags must have a closing tag.
  - (B) All XML tags are case-sensitive.
  - (C) XML elements must be properly nested.
  - (D) XML documents may or may not have a root tag.
29. XML documents follow which structure?
- (A) Graph
  - (B) Tree
  - (C) Stack
  - (D) Queue
30. Which of the following is an advantage of client-server computing?
- (A) Client-server computing provides cost-effective user interface.
  - (B) Client-server computing provides storage for data.
  - (C) Client-server computing provides application services.
  - (D) All the above

### ANSWER KEYS

- |       |       |       |       |       |       |       |       |       |       |
|-------|-------|-------|-------|-------|-------|-------|-------|-------|-------|
| 1. A  | 2. A  | 3. B  | 4. A  | 5. D  | 6. C  | 7. B  | 8. A  | 9. C  | 10. B |
| 11. D | 12. C | 13. B | 14. B | 15. C | 16. C | 17. C | 18. C | 19. B | 20. D |
| 21. D | 22. B | 23. A | 24. C | 25. D | 26. B | 27. C | 28. D | 29. B | 30. D |