NAAC ACCREDITED



तेजस्वि नावधीतमस्तु ISO 9001:2008 & 14001:2004

Institute of Management & Technology

'A' Grade Institute by DHE, Govt. of NCT Delhi and Approved by the Bar Council of India and NCTE

Reference Material for Three Years

Bachelor of Computer Application

Code :020

Semester – I

FIMT Campus, Kapashera, New Delhi-110037, Phones : 011-25063208/09/10/11, 25066256/ 57/58/59/60 Fax : 011-250 63212 Mob. : 09312352942, 09811568155 E-mail : fimtoffice@gmail.com Website : www.fimt-ggsipu.org

DISCLAIMER :FIMT, ND has exercised due care and caution in collecting the data before publishing tis Reference Material. In spite of this ,if any omission,inaccuracy or any other error occurs with regards to the data contained in this reference material, FIMT, ND will not be held responsible or liable. FIMT, ND will be grateful if you could point out any such error or your suggestions which will be of great help for other readers.

INDEX

Three Years

Bachelor of Computer Application

Code : 020

Semester – I

S.NO.	SUBJECTS	CODE	PG.NO.
1	MATHEMATICS-I	101	03-205
2	TECHNICAL COMMUNICATION	103	206-266
3	INTRO. TO PROGRAMMING LANGUAGE	105	267-384
4	INTRO. TO COMPUTER & IT	107	385-586
5	PHYSICS	109	587-747

150 9001:2015 & 14001:2015

MATHEMATICS-I (101)

UNIT - I

DETERMINANTS:

In linear algebra, the **determinant** is a value associated with a square matrix. It can be computed from the entries of the matrix by a specific arithmetic expression, while other ways to determine its value exist as well. The determinant provides important information about a matrix of coefficients of a system of linear equations, or about a matrix that corresponds to a linear transformation of a vector space. In the first case the system has a unique solution exactly when the determinant is nonzero; when the determinant is zero there are either no solutions or many solutions. In the second case the transformation has an inverse operation exactly when the determinant is nonzero. A geometric interpretation can be given to the value of the determinant of a square matrix with real entries: the absolute value of the determinant gives the scale factor by which area or volume (or a higher dimensional analogue) is multiplied under the associated linear transformation, while its sign indicates whether the transformation preserves orientation. Thus a 2×2 matrix with determinant -2, when applied to a region of the plane with finite area, will transform that region into one with twice the area, while reversing its orientation. The determinant of the matrix

$$\begin{bmatrix} a & b & c \\ d & e & f \\ g & h & i \end{bmatrix}$$

is written

$$\begin{vmatrix} a & b & c \\ d & e & f \\ g & h & i \end{vmatrix}$$

and has the value

aei + bfg + cdh - ceg - bdi - afh.

Minor

An element, aij, to the value of the determinant of order n - 1, obtained by deleting the row i and the column j in the matrix is called a minor.

Cofactor

The cofactor of the element aij is its minor prefixing:

The + sign if $\mathbf{i}+\mathbf{j}$ is even.

The - sign if i+j is odd.

1	2	1	Б	
2	5	4	$\rightarrow - \frac{2}{5}$	
3	6	2	p	2

Properties of Determinants:-

The determinant has many properties. Some basic properties of determinants are:

1. $det(I_n) = 1_{Where I_n}$ is the $n \times n$ identity matrix.

ARM

- $\det(A^{\mathrm{T}}) = \det(A).$
- 3.

 $\det(A^{-1}) = \frac{1}{\det(A)}.$

4. For square matrices A and B of equal size, det(AB) = det(A) det(B).

- 5. $\det(cA) = c^n \det(A)_{\text{for an } n \times n \text{ matrix.}}$
- 6. If *A* is a triangular matrix, i.e. $a_{i,j} = 0$ whenever i > j or, alternatively, whenever i < j, then its determinant equals the product of the diagonal entries:

$$\det(A) = a_{1,1}a_{2,2}\cdots a_{n,n} = \prod_{i=1}^{n} a_{i,i}$$

This can be deduced from some of the properties below, but it follows most easily directly from the Leibniz formula (or from the Laplace expansion), in which the identity permutation is the only one that gives a non-zero contribution.

A number of additional properties relate to the effects on the determinant of changing particular rows or columns:

- 7. Viewing an $n \times n$ matrix as being composed of n columns, the determinant is an n-linear function. This means that if one column of a matrix A is written as a sum v + w of two column vectors, and all other columns are left unchanged, then the determinant of A is the sum determinants of the matrices obtained from A by replacing the column by v respectively by w (and a similar relation holds when writing a column as a scalar multiple of a column vector).
- 8. This *n*-linear function is an alternating form. This means that whenever two columns of a matrix are identical, or more generally some column can be expressed as a linear combination of the other columns (i.e. the columns of the matrix form a linearly dependent set), its determinant is 0.

MATRICES:

In mathematics, a **matrix** (plural **matrices**) is a rectangular array of numbers, symbols, or expressions, arranged in *rows* and *columns*.^{[1][2]} The individual items in a matrix are called its *elements* or *entries*. An example of a matrix with 2 rows and 3 columns is

[1	9	-13	-
20	5	-6	·

Types of Matrices

A matrix may be classified by types. It is possible for a matrix to belong to more than one type.

A row matrix is a matrix with only one row.

$$E = (4)$$

E is a row matrix of order 1×1

$$B = (9 -2 5)$$

B is a row matrix of order 1×3

A column matrix is a matrix with only one column.

$$C = (3)$$

C is a column matrix of order 1×1

 $D = \begin{pmatrix} -5 \\ 3 \end{pmatrix}$ NAAC ACCREDITED

D is a column matrix of order 2 × 1

A column matrix of order 2 ×1 is also called a vector matrix.

A zero matrix or a null matrix is a matrix that has all its elements zero.

 $O = \begin{pmatrix} 0 & 0 & 0 \\ 0 & 0 & 0 \end{pmatrix}$

O is a zero matrix of order 2 × 3

A square matrix is a matrix with an equal number of rows and columns.

$$T = \begin{pmatrix} 6 & 3 \\ 0 & 4 \end{pmatrix}$$

T is a square matrix of order 2×2

$$V = \begin{pmatrix} 7 & 1 & 9 \\ 3 & 2 & 5 \\ 2 & 1 & 8 \end{pmatrix}$$

V is a square matrix of order 3×3

A **diagonal matrix** is a square matrix that has all its elements zero except for those in the diagonal from top left to bottom right; which is known as the **leading diagonal** of the matrix.

	(3	0	0)
B =	0	8	0
	lo	0	2)

B is a diagonal matrix

A unit matrix is a diagonal matrix whose elements in the diagonal are all ones.



P is a unit matrix.

Matrix Addition and subtraction:-

Two matrices can only be added or subtracted if they have the same size. Matrix addition and subtraction are done entry-wise, which means that each entry in A+B is the sum of the corresponding entries in A and B.

Here is an example of matrix addition

$$A = \begin{bmatrix} 7 & 5 & 3\\ 4 & 0 & 5 \end{bmatrix} \qquad B = \begin{bmatrix} 1 & 1 & 1\\ -1 & 3 & 2 \end{bmatrix}$$
$$A + B = \begin{bmatrix} 7+1 & 5+1 & 3+1\\ 4-1 & 0+3 & 5+2 \end{bmatrix} = \begin{bmatrix} 8 & 6 & 4\\ 3 & 3 & 7 \end{bmatrix}$$

And an example of subtraction

$$A = \begin{bmatrix} 7 & 5 & 3 \\ 4 & 0 & 5 \end{bmatrix} \qquad B = \begin{bmatrix} 1 & 1 & 1 \\ -1 & 3 & 2 \end{bmatrix}$$
$$A - B = \begin{bmatrix} 7 - 1 & 5 - 1 & 3 - 1 \\ 4 + 1 & 0 - 3 & 5 - 2 \end{bmatrix} = \begin{bmatrix} 6 & 4 & 2 \\ 5 & -3 & 3 \end{bmatrix}$$

Remember you can not add or subtract two matrices of different sizes.

The following rules apply to sums and scalar multiples of matrices. Let A, B, and C be matrices of the same size, and let r and s be scalars.

- A + B = B + A
- (A + B) + C = A + (B + C)
- A + 0 = A
- r(A + B) = rA + rB
- (r+s)A = rA + sA

Multiplication

What is matrix multiplication? You can multiply two matrices if, and only if, the number of columns in the first matrix equals the number of rows in the second matrix.

Otherwise, the product of two matrices is undefined. The product matrix's dimensions are \rightarrow (rows of first matrix) × (columns of the second matrix) In above multiplication, the matrices can be multiplied since the number of columns in the 1st one, matrix A, equals the number of rows in the 2nd, matrix B. The Dimensions of the product matrix. Rows of 1st matrix × Columns of 2nd 4 × 3.

$$\mathbf{A} = \begin{pmatrix} a & b \\ c & d \end{pmatrix}, \quad \mathbf{B} = \begin{pmatrix} \alpha & \beta \\ \gamma & \delta \end{pmatrix},$$

their matrix products are:

$$\mathbf{AB} = \begin{pmatrix} a & b \\ c & d \end{pmatrix} \begin{pmatrix} \alpha & \beta \\ \gamma & \delta \end{pmatrix} = \begin{pmatrix} a\alpha + b\gamma & a\beta + b\delta \\ c\alpha + d\gamma & c\beta + d\delta \end{pmatrix},_{\text{and}}$$
$$\mathbf{BA} = \begin{pmatrix} \alpha & \beta \\ \gamma & \delta \end{pmatrix} \begin{pmatrix} a & b \\ c & d \end{pmatrix} = \begin{pmatrix} \alpha a + \beta c & \alpha b + \beta d \\ \gamma a + \delta c & \gamma b + \delta d \end{pmatrix}.$$

Adjoint of Matrix: - The classical Adjoint of a square matrix A the transpose of the matrix who (i, j) entry is a i j cofactor.

(Adjoint of a Matrix) Let Abe an matrix. The matrix $B = [b_{ij}]$ with $b_{ij} = C_{ji}$, for $1 \le i, j \le n$ is called the Adjoint of A, denoted Adj(A).

EXAMPLE

$$A = \begin{bmatrix} 1 & 2 & 3 \\ 2 & 3 & 1 \\ 1 & 2 & 2 \end{bmatrix}. \qquad Adj(A) = \begin{bmatrix} 4 & 2 & -7 \\ -3 & -1 & 5 \\ 1 & 0 & -1 \end{bmatrix};$$

Let

 $C_{11} = (-1)^{1+1}A_{11} = 4, C_{12} = (-1)^{1+2}A_{12} = -3, C_{13} = (-1)^{1+3}A_{13} = 1,$ as and so on.

Inverse of a Matrix

IAAC ACCREDITED **Definition and Examples**

Recall that functions f and g are inverses if

$$f(g(x)) = g(f(x)) = x$$

We will see later that matrices can be considered as functions from Rⁿ to R^m and that matrix multiplication is composition of these functions. With this knowledge, we have the following:

Let A and B be n x n matrices then A and B are inverses of each other, then

 $AB = BA = I_n$

Example

Consider the matrices

 $B = \begin{pmatrix} 2 & 1 \\ -4 & -2 \\ 3 & 2 \end{pmatrix}$ $A = \begin{bmatrix} 2 & 0 \\ -3 & 0 \end{bmatrix}$ We can check that when we multiply A and B in either order we get the identity matrix.

(Check this.)

Not all square matrices have inverses. If a matrix has an inverse, we call it nonsingular or invertible. Otherwise it is called singular. We will see in the next section how to determine if a matrix is singular or nonsingular.

Determinants and Cramer's Rule

The 2 X 2 system

ax + by = e;

cx + dy = f;

has a unique solution provided Δ = ad- bc is nonzero, in which case the

NHAGE

NAAC ACCREDITE

Solution is given by

x = (de - bf)/(ad - bc); y = (af - ce)/(ad - bc)

This result, called Cramer's Rule for 2 X2 systems.\

Cofactor Expansion

The special subject of cofactor expansions is used to justify Cramer's rule and to provide an alternative method for computation of determinants. There is no claim that cofactor expansion is e_cient, only that it is possible, and di_erent than Sarrus' rule or the use of the four properties.

The Cayley-Hamilton Theorem:-

Presented here is an adjoint formula $F^{-1} = adj(F) = det(F)$ derivation

for the celebrated Cayley-Hamilton formula

$(-A)^{-n} + p_{n-1}(-A)^{n-1} + ___ + p0I = 0.$

The nxn matrix A is given and I is the identity matrix. The coe_cients pk in above are determined by the characteristic polynomial of matrix A, which is de_ned by the determinant expansion formula

Det (A).

Dependence of Vectors:-

A subset **S** of a vector space **V** is called *linearly dependent* if there exist a **finite** number of **distinct** vectors $u_1, u_2, ..., u_n$ in **S** and scalars $a_1, a_2, ..., a_n$, not all zero, such that Note that the zero on the right is the zero vectors, not the number zero.

For any vectors $u_1, u_2, ..., u_n$ we have that

This is called the trivial representation of 0 as a linear combination of u_1 , u_2 ,..., u_n , this motivates a very simple definition of both linear independence and linear dependence, for a set to be linearly dependent, there must exist a non-trivial representation of 0 as a linear combination of vectors in the set.

A subset **S** of a vector space **V** is then said to be *linearly independent* if it is not linearly dependent, in other words, a set is linearly independent if the only representation of 0 as a linear combination its vectors are trivial representations.^[1]

Note that in both definitions we also say that the vectors in the subset S are linearly dependent or linearly independent.

More generally, let **V** be a vector space over a field **K**, and let $\{\mathbf{v}_i \mid i \in I\}$ be a family of elements of **V**. The family is *linearly dependent* over **K** if there exists a family $\{a_j \mid j \in J\}$ of elements of **K**, not all zero, such that where the index set *J* is a nonempty, finite subset of *I*. Unit - II

Lecture 4: Concept of Limit

Definition: We say that the limit of f(x) is L as x approaches a and write this as

 $\lim f(x) = L.$

An alternative notation that we will occasionally use in denoting limits is $f(x) \rightarrow L$ as $x \rightarrow a$

without actually letting x = a.

This means that the definition says that as x gets closer and closer to x = a from both sides of course then f(x) must be getting closer and closer to L or, as we move in towards x = a then f(x) must be moving in towards L. **Definition:** Right-handed limit is denoted by $\lim_{x\to a^+} f(x) = L$ and left-handed limit is denoted by $\lim_{x\to a^-} f(x) = L$.

Given a function f(x) if, $\lim_{x \to a^+} f(x) = L = \lim_{x \to a^-} f(x)$ then the limit will exist and $\lim_{x \to a} f(x) = L$

Likewise, if $\lim_{x \to a} f(x) = L$ then, $\lim_{x \to a^+} f(x) = L = \lim_{x \to a^-} f(x)$

If $\lim_{x \to a^+} f(x) \neq \lim_{x \to a^-} f(x)$ then the limit does not exist.

Example 1: Given the following graph, compute each of the following.

1. <i>f</i> (-2)	$2. \lim_{x \to -2^+} f(x)$	3. <i>f</i> (0)	4. $\lim_{x \to 0^+} f(x)$	5. $\lim_{x \to 0^{-}} f(x)$
6. $\lim_{x \to 0} f(x)$	7. <i>f</i> (2)	8. $\lim_{x \to 2^{-}} f(x)$	9. $\lim_{x \to 2^+} f(x)$	10. $\lim_{x \to 2} f(x) = 2$
11. <i>f</i> (-3)	12. $\lim_{x \to 3^{-}} f(x)$	13. $\lim_{x \to 3^+} f(x)$	14. $\lim_{x \to 3} f(x) = 2$	15. <i>f</i> (4)
$1 \subset 1^{\circ} \subset ($				



Exercise 1: Given the following graph, compute each of the following.

1. <i>f</i> (0)	2. $\lim_{x \to 0+} f(x)$	3. <i>f</i> (1)	$4. \lim_{x \to 1+} f(x)$	5. $\lim_{x \to 1^{-}} f(x)$	6. $\lim_{x \to 1} f(x)$
7. <i>f</i> (2)	8. $\lim_{x \to 2^{-}} f(x)$	9. $\lim_{x \to 2^+} f(x)$	10. $\lim_{x \to 2} f(x) = 2$	11. <i>f</i> (3)	12. $\lim_{x \to 3^{-}} f(x)$
13. $\lim_{x \to 3^+} f(x)$	14. $\lim_{x \to 3} f(x)$	15. <i>f</i> (4)	16. $\lim_{x \to 4^{-}} f(x)$		



Limit Properties

- If $\lim_{x \to c} f(x) = L_1$ and $\lim_{x \to c} g(x) = L_2$, then
- 1. $\lim_{x \to c} (f(x) \pm g(x)) = L_1 \pm L_2$. 2. $\lim_{x \to c} (f(x) \cdot g(x)) = L_1 \cdot L_2$.
- 3. $\lim_{x \to c} \frac{f(x)}{g(x)} = \frac{L_1}{L_2}, L_2 \neq 0.$

 $\lim_{x \to c} (af(x)) = aL_1, a \text{ constant}.$

- 5. $\lim_{x \to c} (f(x))^n = (\lim_{x \to c} f(x))^n = L_1^n, n \in N$.
- 6. $\lim_{x \to c} \sqrt[n]{f(x)} = \sqrt[n]{\lim_{x \to c} f(x)} = \sqrt[n]{L_1}, n \in N, \text{ and for } n \text{ even, we assume that } L_1 > 0.$
- 7. $\lim_{x \to c} x = c$. 8. $\lim_{x \to c} x^n = c^n$

Fact: If $p_n(x)$ is a polynomial of degree n, then $\lim_{x \to c} p_n(x) = p_n(c)$.

COPYRIGHT FIMT 2020

4.

Example 2
$$\lim_{x \to 3} \frac{x^3 - 2x^2}{x^2 + 2} = \frac{\lim_{x \to 3} \left(x^3 - 2x^2\right)}{\lim_{x \to 3} \left(x^2 + 2\right)} = \frac{(3)^3 - 2(3)^2}{(3)^2 + 2} = \frac{9}{11}$$

Example 3:

$$\lim_{x \to -3} \sqrt{\frac{x^2 + 2x + 1}{8 + 2x}} = \sqrt{\lim_{x \to -3} \left(\frac{x^2 + 2x + 1}{8 + 2x}\right)} = \sqrt{\frac{\lim_{x \to -3} (x^2 + 2x + 1)}{\lim_{x \to -3} (8 + 2x)}} = \sqrt{\frac{9 - 6 + 1}{8 - 6}} = \sqrt{\frac{4}{2}} = \sqrt{\frac{4}{2}}$$

Example 4: Given the function, $f(x) = \begin{cases} x^3 + 2x + 1, & x < 1 \\ 3x - 1, & x \ge 1 \end{cases}$. Compute the following

limits.

1. $\lim_{x \to 2} f(x)$ 2. $\lim_{x \to 1} f(x)$

Lecture 5: Computing Limits

Remark: Avoid common mistakes of the form $\frac{0}{0}$. Typically zero in the denominator means it's undefined. However that will only be true if the numerator isn't also zero. Also, zero in the numerator usually means that the fraction is zero, unless the denominator is also zero.

14001:2015

So, there are three cases to compute $\lim_{x \to a} \frac{f(x)}{g(x)}$

1.
$$g(a) \neq 0$$
. In this case $\lim_{x \to a} \frac{f(x)}{g(x)} = \frac{f(a)}{g(a)}$.

2. g(a) = 0 and $f(a) \neq 0$. In this case $\lim_{x \to a} \frac{f(x)}{g(x)}$ does not exist.

3. g(a) = 0 and f(a) = 0. In this case $\lim_{x \to a} \frac{f(x)}{g(x)}$ can be calculated by using algebraic

manipulation .

Case 1: $g(a) \neq 0$

Example 1: Evaluate the following limit.

(1)
$$\lim_{x \to 2} \frac{x^3 + x + 1}{x^2 + 2}$$
 (2) $\lim_{x \to 1} \frac{x^2 - 1}{x^2 + 1}$
Solution: (1) $\lim_{x \to 1} \frac{x^3 + x + 1}{x^2 + 1} = \frac{(2)^3 + 2 + 1}{x^2 + 1} = \frac{11}{x^2 + 1}$ (2) $\lim_{x \to 1} \frac{x^2 - 1}{x^2 - 1} = \frac{0}{x^2 + 1}$

R R P

Solution: (1)
$$\lim_{x \to 2} \frac{x + x + 1}{x^2 + 2} = \frac{(2)^2 + 2 + 1}{(2)^2 + 2} = \frac{11}{6}$$
 (2) $\lim_{x \to 1} \frac{x - 1}{x^2 + 1} = \frac{0}{2} = 0$

Case 2: g(a) = 0 and $f(a) \neq 0$ [Limits that equal infinity]

Definition

We say that $\lim_{x \to a} f(x) = \infty$ if we can make f(x) arbitrarily large for all x sufficiently close to x = a, from both sides, without actually letting x = a. We say that $\lim_{x \to a} f(x) = -\infty$ if we can make f(x) arbitrarily large and negative for all x sufficiently close to x = a, from both sides, without actually letting x = a.

Remark: Concider the limit $\lim_{x \to a} \frac{f(x)}{g(x)}$

1. If $\lim f(x) = \infty$ and $\lim f(x) = -\infty$, then the limit doesn't exist.

- 2. If $\lim_{x \to a^+} f(x) = -\infty$ and $\lim_{x \to a^-} f(x) = \infty$, then the limit doesn't exist.
- 3. If $\lim_{x \to a^+} f(x) = \infty$ and $\lim_{x \to a^-} f(x) = \infty$, then the limit doesn't exist and $\lim_{x \to a} f(x) = \infty$.

4. If
$$\lim_{x\to\infty} f(x) = -\infty$$
 and $\lim_{x\to\infty} f(x) = -\infty$, then the limit doesn't exist and
 $\lim_{x\to\infty} f(x) = -\infty$.
Example 2: Evaluate $\lim_{x\to\infty} \frac{x+3}{x-2}$
Solution:
 $\lim_{x\to\infty^{+}} \frac{x+3}{x-2} = -\infty$ and $\lim_{x\to\infty^{+}} \frac{x+3}{x-2} = \infty$. So $\lim_{x\to\infty^{+}} \frac{x+3}{x-2}$ doesn't exist.
Example 3: Evaluate $\lim_{x\to\infty^{+}} \frac{1}{x-2} = \infty$ and $\lim_{x\to\infty^{+}} \frac{1}{x^2} = \infty$. So $\lim_{x\to\infty^{+}} \frac{1}{x^2} = \infty$.
Exercise 1: (1) Evaluate $\lim_{x\to\infty^{+}} \frac{-x}{\sqrt{4-x^2}}$ (2) $\lim_{x\to0^{+}} \frac{-1}{(x-3)^2}$
(1) $\lim_{x\to\infty^{+}} \frac{-1}{\sqrt{4-x^2}} = -\infty$ and $\lim_{x\to\infty^{+}} \frac{-1}{(x-3)^2} = -\infty$.
(2) $\lim_{x\to\infty^{+}} \frac{-x}{\sqrt{4-x^2}} = -\infty$
Example 4: Evaluate the following limits.
 $\lim_{x\to\infty^{+}} \frac{2-x}{x-2x-8}$

16 | Page

Solution:
$$\lim_{x \to +\infty} \frac{2-x}{x^2-2x-8} = -\infty$$
$$\lim_{x \to +\infty} \frac{2-x}{x^2-2x-8} = \infty$$

So
$$\lim_{x \to +\infty} \frac{2-x}{x^2-2x-8}$$
 doesn't exist.

Definition: A line $x = a$ is called a vertical asymptote of the graph of f if either
$$\lim_{x \to +\infty} f(x) = \pm \infty$$
 or
$$\lim_{x \to +\infty} f(x) = \pm \infty$$

Solution: By analyze the sign of f is in the figure, notice that
$$\lim_{x \to +\infty} f(x) = \infty$$
 and
$$\lim_{x \to +\infty} f(x) = -\infty$$

So the vertical asymptote is the line $x = 2$.

Exercise 2: Find the vertical asymptotes of
 $(1) \quad f(x) = \frac{x^2}{x^2-1}$
 $f(x) = \left\{\frac{x+1}{x-2}, x < 0\\ x^3 - 3x - 1, x \ge 0\right\}$
Solution: Notice that
$$(1) \quad \lim_{x \to +\infty} f(x) = -\infty$$

$$\lim_{x \to +\infty} f(x) = -\infty \text{ and } \quad \lim_{x \to +\infty} f(x) = -\infty$$

$$\lim_{x \to +\infty} f(x) = -\infty \text{ and } \quad \lim_{x \to +\infty} f(x) = \infty$$

So the vertical asymptotes are the lines x = 1 and x = -1.

(2) Notice that f has no vertical asymptote because f is defined for all x.

Example 6: Find the vertical asymptotes of (1) $f(x) = xe^{-x}$ (2) $f(x) = \frac{x}{\sqrt{4+x^2}}$

(3) $f(x) = 4 \tan^{-1} x - 1$

Solution: f has no vertical asymptote because f is defined for all x.

Exercise 3: Find the vertical asymptotes of (1) $f(x) = 3e^{-1/x}$. (2) $f(x) = 6x^{1/3} + 3x^{4/3}$

Solution: (1) $\lim_{x \to 0^+} 3e^{-l/x} = 0$ $\lim_{x \to 0^-} 3e^{-l/x} = \infty$. So the vertical asymptote is x = 0 as $x \to 0^-$.

(2) Notice that f has no vertical asymptote because f is defined for all x.

Exercise 4: Find the vertical asymptotes of
$$f(x) = \begin{cases} \frac{4x}{x-4}, x < 0\\ \frac{x^2}{x-2}, 0 \le x < 4\\ \frac{e^{-x}}{x+1}, x \ge 4 \end{cases}$$
.

Solution: $\lim_{x \to 2^{-}} f(x) = -\infty$ and $\lim_{x \to 2^{+}} f(x) = \infty$. So the vertical asymptote is the line

x = 2. **150 9001:2015 & 14001:2015**

Case 3: g(a) = 0 and f(a) = 0

Example 7: Evaluate $\lim_{x \to -2} \frac{x^3 + 8}{x^2 - 4}$ **?**

Solution: take x = -2 we get, $\lim_{x \to -2} \frac{x^3 + 8}{x^2 - 4} = \begin{pmatrix} 0 \\ 0 \end{pmatrix}$ $\lim_{x \to -2} \frac{x^3 + 8}{x^2 - 4} = \lim_{x \to -2} \frac{(x + 2)(x^2 - 2x + 4)}{(x + 2)(x - 2)} = \lim_{x \to -2} \frac{(x^2 - 2x - 4)}{(x - 2)}$ $\lim_{x \to -2} \frac{x^3 + 8}{x^2 - 4} = \lim_{x \to -2} \frac{(x^2 - 2x - 4)}{(x - 2)} = \frac{(-2)^2 - 2(-2) - 4}{-2 - 2} = \frac{4 + 4 - 4}{-4} = -1$ Example 8: Evaluate $\lim_{x \to 4} \frac{x + \sqrt{x} - 6}{\sqrt{x} - 2}$. Solution

$$\lim_{x \to 4} \frac{x + \sqrt{x} - 6}{\sqrt{x} - 2} \left(\frac{0}{0}\right) = \lim_{x \to 4} \frac{(\sqrt{x} + 3)(\sqrt{x} - 2)}{(\sqrt{x} - 2)} = \lim_{x \to 4} (\sqrt{x} + 3) = \sqrt{4} + 3 = 2 + 3 = 5$$

Exercise 5: Evaluate the following limit. (1)
$$\lim_{x \to 3} \frac{x^3 + 4x - 39}{x^3 - 27}$$
 (2)

-

 $\lim_{x \to 3} \frac{x^2 - x - 6}{x^2 - 2x - 3}$ (3) $\lim_{x \to 1} \frac{x^{\frac{3}{2}} - x}{x^{\frac{1}{2}} - 1}$ (4) $\lim_{x \to 0} \frac{(x + 2)^2 - 4}{x}$? $\lim_{x \to 3} \frac{x^3 + 4x - 39}{x^3 - 27} \left(\frac{0}{0}\right) = \lim_{x \to 3} \frac{(x - 3)(x^2 + 3x + 13)}{(x - 3)(x^2 + 3x + 9)} = \lim_{x \to 3} \frac{(x^2 + 3x + 13)}{(x^2 + 3x + 9)}$ Solution: (1) $= \frac{(3)^2 + 3(3) + 13}{(3)^2 + 3(3) + 9} = \frac{9 + 9 + 13}{9 + 9 + 9} = \frac{31}{27}$

(2)
$$\lim_{x \to 3} \frac{x^2 - x - 6}{x^2 - 2x - 3} \left(\frac{0}{0} \right) = \lim_{x \to 3} \frac{(x - 3)(x + 2)}{(x - 3)(x + 1)} = \lim_{x \to 3} \frac{(x + 2)}{(x + 1)} = \frac{3 + 2}{3 + 1} = \frac{5}{4}$$

(3)
$$\lim_{x \to 1} \frac{x^{\frac{3}{2}} - x}{x^{\frac{1}{2}} - 1} \left(\frac{0}{0}\right) = \lim_{x \to 1} \frac{x \left(x^{\frac{1}{2}} - 1\right)}{x^{\frac{1}{2}} - 1} = \lim_{x \to 1} x = 1$$

19 | Page

$$(4) = \lim_{x \to 0} \frac{x(x+4)}{x} = \lim_{x \to 0} (x+4) = 4 \lim_{x \to 0} \frac{(x+2)^2 - 4}{x} = \lim_{x \to 0} \frac{x^2 + 4x + 4 - 4}{x} = \lim_{x \to 0} \frac{x^2 + 4x}{x}$$

Example 8: Evaluate the following limit.

$$\lim_{h \to 0} \frac{\left(h+1\right)^2 - 1}{h}$$

Solution

$$\lim_{h \to 0} \frac{(h+1)^2 - 1}{h} \left(\frac{0}{0}\right) = \lim_{h \to 0} \frac{h^2 + 2h + 1 - 1}{h} = \lim_{h \to 0} \frac{h(h+2)}{h} = \lim_{h \to 0} (h+2) = 2$$

AAGEMEAN

Example 9: Evaluate the following limit.

$$\lim_{x \to 3} \frac{\sqrt{x+6}-3}{x-3}$$

Solution:
$$\lim_{x \to 3} \frac{\sqrt{x+6}-3}{x-3} = \lim_{x \to 3} \frac{\left(\sqrt{x+6}-3\right)\left(\sqrt{x+6}+3\right)}{\left(x-3\right)\left(\sqrt{x+6}+3\right)}$$
$$= \lim_{x \to 3} \frac{(x+6)-(9)}{(x-3)\left(\sqrt{x+6}+3\right)} = \lim_{x \to 3} \frac{(x-3)}{(x-3)\left(\sqrt{x+6}+3\right)} = \lim_{x \to 3} \frac{1}{\left(\sqrt{x+6}+3\right)} = \frac{1}{6}$$

Exercise 6: Evaluate the following limit.

(1)
$$\lim_{x \to 1} \frac{2 - \sqrt{x+3}}{x^2 + 2x - 3}$$
 (2)
$$\lim_{x \to 2} \frac{\sqrt{5x-1}-3}{8-2x^2}$$
 (3)
$$\lim_{x \to 3} \frac{x^2 - 3x}{x - \sqrt{x+1}-1}$$

(4)
$$\lim_{x \to 4} \frac{3 - \sqrt{2x+1}}{\sqrt{x-2}}$$
 (5)
$$\lim_{x \to 0} \frac{\sqrt[3]{1+2x}-1}{x}$$
 (6)
$$\lim_{x \to 1} \frac{\sqrt[3]{x}-1}{\sqrt[4]{x}-1}$$

(7)
$$\lim_{x \to 1} \frac{(x+1)^3 \sqrt{x} - 8}{x-1}$$
 (8) $\lim_{x \to 1} \frac{\frac{5}{2x-3} + 5}{4x^2 - 4}$ (9) $\lim_{x \to 1^+} \frac{\sqrt{x^2 - 1} + \sqrt{x-1}}{\sqrt{x-1}}$

Example 36: Given the function,

$$f(x) = \begin{cases} \frac{x^3 - 1/x^3}{x - 1/x}, & x \neq 1 \\ 3, & x = 1 \end{cases}$$

Compute the following limit.

$$\lim_{x\to 1} f(x)$$

Solution

In doing limits recall that we must always look at what's happening on both sides of the point in question as we move in towards it. As x approaches 1 from the left and from the right is inside the first interval for the function and so there are values of x on both sides of x = 1 inside this interval. This means that we can just use the fact to evaluate this limit.

ALGEMEN-

a she will be a more star. I want that

$$\lim_{x \to 1} f(x) = \lim_{x \to 1} \frac{x^3 - 1/x^3}{x - 1/x} \left(\frac{0}{0}\right) = \lim_{x \to 1} \frac{(x - \frac{1}{x})(x^2 + x(\frac{1}{x}) + \frac{1}{x^2})}{(x - \frac{1}{x})} = \lim_{x \to 1} (x^2 + 1 + \frac{1}{x^2}) = 1 + 1 + 1 = 3$$

Example 37: Evaluate the following limit.

$$\lim_{x \to 0} \frac{x^2 + x}{\sqrt{x^4 + 2x^2}}$$

Solution
$$\lim_{x \to 0} \frac{x^2 + x}{\sqrt{x^4 + 2x^2}} \left(\frac{0}{0}\right) = \lim_{x \to 0} \frac{x^2 + x}{\sqrt{x^2 (x^2 + 2)}} = \lim_{x \to 0} \frac{x^2 + x}{\sqrt{x^2 \sqrt{x^2 + 2}}} = \lim_{x \to 0} \frac{x^2 + x}{|x|\sqrt{x^2 + 2}}$$

Recall that we must always look at what's happening on both sides of the point in question.

$$\lim_{x \to 0^{+}} \frac{x^{2} + x}{|x|\sqrt{x^{2} + 2}} = \lim_{x \to 0^{+}} \frac{x(x+1)}{x\sqrt{x^{2} + 2}} = \lim_{x \to 0^{+}} \frac{(x+1)}{\sqrt{x^{2} + 2}} = \frac{1}{\sqrt{2}}$$
$$\lim_{x \to 0^{-}} \frac{x^{2} + x}{|x|\sqrt{x^{2} + 2}} = \lim_{x \to 0^{-}} \frac{x(x+1)}{-x\sqrt{x^{2} + 2}} = \lim_{x \to 0^{+}} -\frac{(x+1)}{\sqrt{x^{2} + 2}} = \frac{-1}{\sqrt{2}}$$

Therefore, $\lim_{x \to 0} \frac{x^2 + x}{\sqrt{x^4 + 2x^2}}$ does not exist.

Example 38: Evaluate the following limit. $\lim_{x \to 0} \frac{2x - |x|}{|3x| - 2x|}$

Solution

$$\lim_{x \to 0^{+}} \frac{2x - |x|}{|3x| - 2x} = \lim_{x \to 0^{+}} \frac{2x - x}{3x - 2x} = \lim_{x \to 0^{+}} \frac{x}{x} = \lim_{x \to 0^{+}} 1 = 1$$
$$\lim_{x \to 0^{-}} \frac{2x - |x|}{|3x| - 2x} = \lim_{x \to 0^{-}} \frac{2x - (-x)}{(-3x) - 2x} = \lim_{x \to 0^{-}} -\frac{3x}{5x} = \lim_{x \to 0^{-}} -\frac{3}{5} = -\frac{3}{5}$$

We can see that,

 $\lim_{x \to 0^+} \frac{2x - |x|}{|3x| - 2x} \neq \lim_{x \to 0^-} \frac{2x - |x|}{|3x| - 2x}.$ Therefore, $\lim_{x \to 0} \frac{2x - |x|}{|3x| - 2x}$ doesn't exist.

Lecture 6: Limits at infinity

Definition:

By limits at infinity we mean one of the following two limits.

 $\lim_{x\to\infty}f(x)$

 $\lim_{x\to\infty}f(x)$

Theorem: For n > 0 we have

$$\lim_{x \to \infty} \frac{1}{x^n} = 0 \qquad \qquad \lim_{x \to -\infty} \frac{1}{x^n} = 0$$

This fact should make sense if you think about it. We require n > 0 to make sure the term stays in the denominator and as we increase x then x^n will also increase. So, what we end up with is a constant divided by an increasingly large number so the quotient of the two will become increasingly small. In the limit we will get zero.

Theorem:

$$\lim_{x \to +\infty} (a_n x^n + a_{n-1} x^{n-1} + \dots + a_1 x + a_0) = \lim_{x \to +\infty} a_n x^n, a_n \neq 0$$
$$\lim_{x \to -\infty} (a_n x^n + a_{n-1} x^{n-1} + \dots + a_1 x + a_0) = \lim_{x \to -\infty} a_n x^n, a_n \neq 0$$

Remark: You can avoid common mistakes by giving careful consideration to the forms $\frac{\infty}{\infty}$ and $\infty - \infty$ during the computations of the limit. Initially, many students incorrectly conclude that $\frac{\infty}{\infty}$ is equal to 1, or that the limit does not exist, or is $+\infty$ or $-\infty$. Many also conclude that $\infty - \infty$ is equal to 0. In fact, the forms $\frac{\infty}{\infty}$ and $\infty - \infty$ are examples of indeterminant forms. This simply means that you have not yet determind an answer. Usually, these indeterminate forms can be circumvented by using algebraic manipulation. Such tools as algebraic simplification and conjugates can easily be used to circumvent the forms $\frac{\infty}{\infty}$ and $\infty - \infty$ so that the limit can be calculated.

AL DAVIE AND DAVIE AND

Example 1: Compute (1) $\lim_{x \to \infty} (3x^3 - 1000x^2)$. (2) $\lim_{x \to \infty} (x^4 + 5x^2 + 1)$ (3) $\lim_{x \to \infty} \frac{100}{x^2 + 5}$ Solution: (1) $\lim_{x \to \infty} (3x^3 - 1000x^2) = \lim_{x \to \infty} (3x^3) = \infty$ (2) $\lim_{x \to \infty} (x^4 + 5x^2 + 1) = \lim_{x \to \infty} (x^4) = \infty$ (3) $\lim_{x \to \infty} \frac{100}{x^2 + 5} = \left(\frac{100}{\infty}\right) = 0$

Example 2: Compute (1) $\lim_{x \to -\infty} \frac{x+7}{3x+5}$. (2) $\lim_{x \to \infty} \frac{7x^2 + x - 100}{2x^2 - 5x}$ (3) $\lim_{x \to \infty} \frac{x^2 - 3x + 7}{x^3 + 10x - 4}$

Solution: (1)
$$\lim_{x \to -\infty} \frac{x+7}{3x+5} = \left(\frac{-\infty}{-\infty}\right) = \lim_{x \to -\infty} \frac{\frac{x}{x} + \frac{7}{x}}{\frac{3x}{x} + \frac{5}{x}} = \lim_{x \to -\infty} \frac{1+\frac{7}{x}}{3+\frac{5}{x}} = \frac{1+0}{3+0} = \frac{1}{3}$$

(2) $\lim_{x \to \infty} \frac{7x^2 + x - 100}{2x^2 - 5x} = \left(\frac{\infty}{\infty}\right)$ Circumvent it by dividing each term by x^2 .

So
$$\lim_{x \to \infty} \frac{7x^2 + x - 100}{2x^2 - 5x} = \lim_{x \to \infty} \frac{\frac{7x^2}{x^2} + \frac{x}{x^2} - \frac{100}{x^2}}{\frac{2x^2}{x^2} - \frac{5x}{x^2}} = \lim_{x \to \infty} \frac{7 + \frac{1}{x} - \frac{100}{x^2}}{2 - \frac{5}{x}} = \frac{7 + 0 - 0}{2 - 0} = \frac{7}{2}$$

(3) Circumvent it by dividing each term by x^3 .

6 1 1 1 N

So
$$\lim_{x \to \infty} \frac{x^2 - 3x + 7}{x^3 + 10x - 4} = \lim_{x \to \infty} \frac{\frac{x^2}{x^3} - \frac{3x}{x^3} + \frac{7}{x^3}}{\frac{x^3}{x^3} + \frac{10x}{x^3} - \frac{4}{x^3}} = \lim_{x \to \infty} \frac{\frac{1}{x} - \frac{3}{x^2} + \frac{7}{x^3}}{1 + \frac{10}{x^2} - \frac{4}{x^3}} = \frac{0 - 0 + 0}{1 + 0 - 0} = 0$$

Remark: Dividing by x^2 , the highest power of x in the numerator, also leads to the correct answer.

Example 3: Compute (1)
$$\lim_{x \to \infty} \left(x - \sqrt{x^2 + 7} \right)$$
. (2) $\lim_{x \to -\infty} \left(x - \sqrt{x^2 + 7} \right)$

Solution: (1) $\lim_{x \to \infty} \left(x - \sqrt{x^2 + 7} \right) = \left(\infty - \infty \right)$ $\lim_{x \to \infty} \left(x - \sqrt{x^2 + 7} \right) = \lim_{x \to \infty} \frac{\left(x - \sqrt{x^2 + 7} \right) \left(x + \sqrt{x^2 + 7} \right)}{\left(x + \sqrt{x^2 + 7} \right)} = \lim_{x \to \infty} \frac{\left(x^2 - \left(x^2 + 7 \right) \right)}{\left(x + \sqrt{x^2 + 7} \right)} = \lim_{x \to \infty} \frac{-7}{\left(x + \sqrt{x^2 + 7} \right)} = \frac{-7}{\infty} = 0$ (2) $\lim_{x \to -\infty} \left(x - \sqrt{x^2 + 7} \right) = \left(-\infty - \infty \right).$ This is not an indeterminate form. It means $= -\infty$

Exercise 1: Compute (1)
$$\lim_{x \to \infty} \frac{7x^2 + x + 11}{4 - x}$$
 (2) $\lim_{x \to \infty} \frac{x + 3}{\sqrt{9x^2 - 5x}}$ (3) $\lim_{x \to \infty} \frac{x + 3}{\sqrt{9x^2 - 5x}}$

COPYRIGHT FIMT 2020

24 | Page

$$(4) \lim_{x \to \infty} \left(\sqrt{5x^{2} + x + 3} - \sqrt{5x^{2} + 4x + 7} \right) \quad (5) \lim_{x \to \infty} \left[\frac{\sqrt{5x + 9x^{2}}}{1 + 3x} + 2 \right]$$
Solution
$$(1) \qquad \lim_{x \to \infty} \frac{7x^{2} + x + 11}{4 - x} = \left(\frac{\infty}{\infty}\right) = 1$$

$$\lim_{x \to \infty} \frac{7x^{2} + x + 11}{x^{2} - x^{2}} = \lim_{x \to \infty} \frac{7 + \frac{1}{x} + \frac{11}{x^{2}}}{\frac{4}{x^{2}} - \frac{1}{x}} = \frac{7 + 0 + 0}{0 - 0} = \left(\frac{7}{0}\right) = (\infty)$$

$$(2)$$

$$\lim_{x \to \infty} \frac{x + 3}{\sqrt{9x^{2} - 5x}} = \lim_{x \to \infty} \frac{x + 3}{\sqrt{x^{2}}\left(9 - \frac{5}{x}\right)} = \lim_{x \to \infty} \frac{x + 3}{\sqrt{x^{2}}\sqrt{9 - \frac{5}{x}}} = \lim_{x \to \infty} \frac{x + 3}{|x|\sqrt{9 - \frac{5}{x}}} = \lim_{x \to \infty} \frac{x + 3}{-x\sqrt{9 - \frac{5}{x}}}$$

$$= \lim_{x \to \infty} \frac{x \left(1 + \frac{3}{x}\right)}{-x\sqrt{9 - \frac{5}{x}}} = -\lim_{x \to \infty} \frac{\left(1 + \frac{3}{x}\right)}{\sqrt{9 - \frac{5}{x}}} = -\frac{1}{\sqrt{9}} = -\frac{1}{3}.$$

(4)
$$\lim_{x \to \infty} \left(\sqrt{5x^2 + x + 3} - \sqrt{5x^2 + 4x + 7} \right) (\infty - \infty)$$

$$\lim_{x \to \infty} \left(\sqrt{5x^2 + x + 3} - \sqrt{5x^2 + 4x + 7} \right) = \lim_{x \to \infty} \frac{\left(\sqrt{5x^2 + x + 3} - \sqrt{5x^2 + 4x + 7} \right) \left(\sqrt{5x^2 + x + 3} + \sqrt{5x^2 + 4x + 7} \right)}{\left(\sqrt{5x^2 + x + 3} + \sqrt{5x^2 + 4x + 7} \right)}$$

1

77

$$= \lim_{x \to \infty} \frac{\left(5x^2 + x + 3 - \left(5x^2 + 4x + 7\right)\right)}{\left(\sqrt{5x^2 + x + 3} + \sqrt{5x^2 + 4x + 7}\right)} = \lim_{x \to \infty} \frac{-\left(3x + 4\right)}{\left(\sqrt{5x^2 + x + 3} + \sqrt{5x^2 + 4x + 7}\right)}$$
$$= \lim_{x \to \infty} \frac{-\left(3x + 4\right)}{\left(\left|x\right|\sqrt{5 + \frac{1}{x} + \frac{3}{x^2}} + \left|x\right|\sqrt{5 + \frac{4}{x} + \frac{7}{x^2}}\right)} = \lim_{x \to \infty} \frac{-\left(3x + 4\right)}{\left(x\sqrt{5 + \frac{1}{x} + \frac{3}{x^2}} + x\sqrt{5 + \frac{4}{x} + \frac{7}{x^2}}\right)}$$
$$= \lim_{x \to \infty} \frac{-x\left(3 + \frac{4}{x}\right)}{x\left(\sqrt{5 + \frac{1}{x} + \frac{3}{x^2}} + \sqrt{5 + \frac{4}{x} + \frac{7}{x^2}}\right)} = \lim_{x \to \infty} \frac{-\left(3 + \frac{4}{x}\right)}{\left(\sqrt{5 + \frac{1}{x} + \frac{3}{x^2}} + \sqrt{5 + \frac{4}{x} + \frac{7}{x^2}}\right)} = \frac{-3}{2\sqrt{5}}$$

25 | Page

(5) First

$$\lim_{x \to \infty} \frac{\sqrt{5x + 9x^2}}{1 + 3x} = \lim_{x \to \infty} \frac{\sqrt{x^2 \left(\frac{5}{x} + 9\right)}}{1 + 3x} = \lim_{x \to \infty} \frac{\sqrt{x^2} \sqrt{\frac{5}{x} + 9}}{1 + 3x} = \lim_{x \to \infty} \frac{|x| \sqrt{\frac{5}{x} + 9}}{1 + 3x}$$
$$= \lim_{x \to \infty} \frac{x \sqrt{\frac{5}{x} + 9}}{x \left(\frac{1}{x} + 3\right)} = \lim_{x \to \infty} \frac{\sqrt{\frac{5}{x} + 9}}{\left(\frac{1}{x} + 3\right)} = \frac{3}{3} = 1$$

Second: $\lim_{x \to \infty} 2 = 2$. So $\lim_{x \to \infty} \left[\frac{\sqrt{5x + 9x^2}}{1 + 3x} + 2 \right] = \lim_{x \to \infty} \frac{\sqrt{5x + 9x^2}}{1 + 3x} + \lim_{x \to \infty} 2 = 1 + 2 = 3$

Example 4: Compute (1) $\lim_{x \to \infty} \frac{e^x}{4 + 5e^{3x}}$ (2) $\lim_{x \to \infty} \frac{2^x}{3^x}$ (3) $\lim_{x \to \infty} \frac{5^x}{3^x + 2^x}$

(4)
$$\lim_{x \to \infty} \frac{\ln(2+e^{3x})}{\ln(1+e^{x})}$$
 (5) $\lim_{x \to \infty} e^{2x-e^{2x}}$

Solution: (1)
$$\lim_{x \to -\infty} \frac{e^x}{4+5e^{3x}} = \frac{0}{4+0} = 0$$

(2)
$$\lim_{x \to \infty} \frac{2^x}{3^x} \left(\frac{\infty}{\infty}\right) = \lim_{x \to \infty} \left(\frac{2}{3}\right)^x = 0$$

(3)
$$\lim_{x \to \infty} \frac{5^{x}}{3^{x} + 2^{x}} = \lim_{x \to \infty} \frac{\frac{5^{x}}{3^{x}}}{\frac{3^{x}}{3^{x}} + \frac{2^{x}}{3^{x}}} = \lim_{x \to \infty} \frac{\left(\frac{5}{3}\right)^{x}}{1 + \left(\frac{2}{3}\right)^{x}} = \frac{\infty}{0 + 1} = \infty$$

$$\lim_{x \to \infty} \frac{\ln\left(2 + e^{3x}\right)}{\ln\left(1 + e^{x}\right)} \left(\frac{\infty}{\infty}\right) = \lim_{x \to \infty} \frac{\ln\left(e^{3x}\left(\frac{2}{e^{3x}} + 1\right)\right)}{\ln\left(e^{x}\left(\frac{1}{e^{x}} + 1\right)\right)} = \lim_{x \to \infty} \frac{\ln e^{3x} + \ln\left(\frac{2}{e^{3x}} + 1\right)}{\ln e^{x} + \ln\left(\frac{1}{e^{x}} + 1\right)}$$

. (4)
$$= \lim_{x \to \infty} \frac{3x + \ln\left(\frac{2}{e^{3x}} + 1\right)}{x + \ln\left(\frac{1}{e^{x}} + 1\right)} = \lim_{x \to \infty} \frac{3x}{x} = \lim_{x \to \infty} 3 = 3$$

(5) Notice that
$$\lim_{x \to \infty} (2x - 1) = \infty$$
 and $\lim_{x \to \infty} e^x = \infty$. Combining these two results getting $\lim_{x \to \infty} e^{2x - 1} = \infty$

Exercise 2: Compute (1) $\lim_{x \to \infty} \ln 2x$ (2) $\lim_{x \to 0^+} e^{-2/x}$ (3) $\lim_{x \to 0^+} \tan^{-1}(\ln x)$ (4) $\lim_{x \to 0^+} e^{1/x^2}$ (1) Notice that $\lim_{x \to \infty} 2x = \infty$ and $\lim_{x \to \infty} \ln x = \infty$. So $\lim_{x \to \infty} \ln 2x = \infty$ (2) Notices that $\lim_{x \to 0^+} -\frac{2}{x} = -\infty$ and $\lim_{x \to \infty} e^x = 0$. Thus $\lim_{x \to 0^+} e^{-2/x} = 0$ (3) Notices that $\lim_{x \to 0^+} \ln x = -\infty$ and $\lim_{x \to \infty} \tan^{-1} x = -\frac{\pi}{2}$. So $\lim_{x \to 0^+} \tan^{-1}(\ln x) = -\frac{\pi}{2}$ (4) Notices that $\lim_{x \to 0^+} \frac{1}{x^2} = \infty$ and $\lim_{x \to \infty} e^x = \infty$. Thus $\lim_{x \to 0^+} e^{1/x^2} = \infty$ Exercise 3: If $\lim_{x \to \infty} \frac{f(x)}{g(x)} = 0$, find $\lim_{x \to \infty} \frac{f(x) - g(x)}{f(x) + g(x)}$ Solution $\lim_{x \to \infty} \frac{f(x) - g(x)}{g(x)} = \lim_{x \to \infty} \frac{\frac{f(x)}{g(x)} - \frac{g(x)}{g(x)}}{\frac{f(x)}{g(x)}} = \lim_{x \to \infty} \frac{f(x)}{g(x)} = 1$

Definition: A line y = L is called a horizontal asymptote of the graph of f if either

 $\lim_{x \to \infty} f(x) = L \quad \text{or} \quad \lim_{x \to -\infty} f(x) = L$

Example 5: Find the horizontal asymptotes of (1

g(x) = g(x)

1)
$$f(x) = \frac{x-1}{x-2}$$
 (2) $f(x) = \frac{x^2-1}{x^3}$

(3)
$$f(x) = \begin{cases} \frac{x+1}{x-2}, x < 0\\ x^3 - 3x - 1, x \ge 0 \end{cases}$$
 (4) $f(x) = x \ln x^2$ (5) $f(x) = 6x^{1/3} + 3x^{4/3}$

Solution: (1) $\lim_{x \to \infty} f(x) = 1$ and $\lim_{x \to -\infty} f(x) = 1$. So the horizontal asymptote is the line y = 1.

- (2) $\lim_{x \to \infty} f(x) = 1$ and $\lim_{x \to \infty} f(x) = 1$. So the horizontal asymptote is the line y = 1.
- (3) $\lim_{x \to -\infty} f(x) = \lim_{x \to -\infty} \frac{x+1}{x-2} = 1$ $\lim_{x \to \infty} f(x) = \lim_{x \to \infty} (x^3 3x 1) = \infty.$

So the horizontal asymptote is the line y = 1 as $x \rightarrow -\infty$.

(4) $\lim_{x \to \infty} f(x) = \infty$ and $\lim_{x \to \infty} f(x) = -\infty$. So f has no horizontal asymptote.

(5) $\lim_{x \to \infty} \left[6x^{1/3} + 3x^{4/3} \right] = \infty \qquad \qquad \lim_{x \to \infty} \left[6x^{1/3} + 3x^{4/3} \right] = \infty.$ So f has no horizontal asymptote.

Exercise 4: Find the horizontal asymptotes of (1) $f(x) = x^2 \ln x$ (2) $f(x) = \frac{x}{\sqrt{4+x^2}}$

(1 ...

٦

(3)
$$f(x) = 4 \tan^{-1} x - 1$$
 (4) $f(x) = 3e^{-1/x}$ (5) $f(x) = \begin{cases} \frac{4x}{x-4}, x < 0\\ \frac{x^2}{x-2}, 0 \le x < 4\\ \frac{e^{-x}}{x+1}, x \ge 4 \end{cases}$

Solution: (1) $\lim f(x) = \infty$. So f has no horizontal asymptote.

01:201

5 2 9 9 9

(2)
$$\lim_{x \to \infty} \frac{x}{\sqrt{4 + x^2}} = \lim_{x \to \infty} \frac{x}{|x|\sqrt{\frac{4}{x^2} + 1}} = \lim_{x \to \infty} \frac{x}{x\sqrt{\frac{4}{x^2} + 1}} = \lim_{x \to \infty} \frac{1}{\sqrt{\frac{4}{x^2} + 1}} = \frac{1}{\sqrt{1}} = 1$$
$$\lim_{x \to \infty} \frac{x}{\sqrt{4 + x^2}} = \lim_{x \to \infty} \frac{x}{|x|\sqrt{\frac{4}{x^2} + 1}} = \lim_{x \to \infty} \frac{x}{-x\sqrt{\frac{4}{x^2} + 1}} = \lim_{x \to \infty} -\frac{1}{\sqrt{\frac{4}{x^2} + 1}} = -\frac{1}{\sqrt{1}} = -1$$

58 140

So the horizontal asymptotes are y = 1 as $x \to \infty$ and y = -1 as $x \to -\infty$

(3)
$$\lim_{x \to \infty} 4 \tan^{-1} x - 1 = 4 \left(\frac{\pi}{2} \right) - 1 = 2\pi - 1$$
 $\lim_{x \to -\infty} 4 \tan^{-1} x - 1 = 4 \left(-\frac{\pi}{2} \right) - 1 = -2\pi - 1$

So the horizontal asymptotes are $y = 2\pi - 1$ as $x \to \infty$ and $y = -2\pi - 1$ as $x \to -\infty$.

(4)
$$\lim_{x \to \infty} 3e^{-1/x} = 3$$
 $\lim_{x \to \infty} 3e^{-1/x} = 3$. So the horizontal asymptote is the

line y = 3.

(5)
$$\lim_{x \to \infty} f(x) = \lim_{x \to \infty} \frac{4x}{x - 4} = 4$$
 $\lim_{x \to \infty} f(x) = \lim_{x \to \infty} \frac{e^{-x}}{x + 1} = 0$

So the horizontal asymptotes are the lines y = 0 as $x \to \infty$ and y = 4 as $x \to -\infty$.

Lecture 7: Limits of Trigonometric Functions:

Theorem

$$\lim_{x \to 0} \frac{\sin x}{x} = 1$$

Example 1: Evaluate the following limit.

$$\lim_{x \to 0} \left(\frac{1 - \cos x}{x} \right)$$

$$\lim_{x \to 0} \frac{1 - \cos x}{x} \left(\frac{0}{0}\right) = \lim_{x \to 0} \frac{1 - \cos x}{x} \cdot \frac{1 + \cos x}{1 + \cos x} = \lim_{x \to 0} \frac{1 - \cos^2 x}{x(1 + \cos x)} = \lim_{x \to 0} \frac{\sin^2 x}{x(1 + \cos x)}$$
$$= \lim_{x \to 0} \frac{\sin x}{x} \cdot \lim_{x \to 0} \frac{\sin x}{(1 + \cos x)} = (1)(\frac{0}{1 + 1}) = (1)(0) = 0$$

Example 2: show that.

$$1.\lim_{x \to 0} \frac{\sin ax}{bx} = \lim_{x \to 0} \frac{ax}{\sin bx} = \lim_{x \to 0} \frac{\sin ax}{\sin bx} = \frac{a}{b}$$
$$2.\lim_{x \to 0} \frac{\tan ax}{bx} = \lim_{x \to 0} \frac{ax}{\tan bx} = \lim_{x \to 0} \frac{\tan ax}{\tan bx} = \frac{a}{b}$$
$$3.\lim_{x \to 0} \frac{\sin ax}{\tan bx} = \lim_{x \to 0} \frac{\tan ax}{\sin bx} = \frac{a}{b}$$

Solution

We will show that

$\sin ax a$	tan ax a	$\sin ax = a$
$\lim_{m \to \infty} m = -$	$l_{1m} = -$	$\lim_{n \to \infty} m = m$
$x \rightarrow 0 bx \qquad b$	$x \rightarrow 0 bx b$	$x \to 0 \tan b x b$

ETRANAG

and left the rest as an exercise.

Solution 1

$$\lim_{x \to 0} \frac{\sin ax}{bx} = \lim_{x \to 0} \frac{\frac{\sin ax}{ax}}{\frac{bx}{ax}} = \lim_{x \to 0} \frac{\frac{\sin ax}{ax}}{\frac{b}{a}} = \lim_{x \to 0} \frac{a}{b} \frac{\sin ax}{ax}$$

_.

let
$$y = ax$$
, as $x \to 0, y \to 0$. So

$$\lim_{x \to 0} \frac{\sin ax}{bx} = \frac{a}{b} \lim_{x \to 0} \frac{\sin ax}{ax} = \frac{a}{b} \lim_{y \to 0} \frac{\sin y}{y} = \frac{a}{b} (1) = \frac{a}{b}$$

Solution 2

$$\lim_{x \to 0} \frac{\tan ax}{bx} = \lim_{x \to 0} \left(\frac{\frac{\sin ax}{\cos bx}}{bx} \right) = \lim_{x \to 0} \left(\frac{1}{\cos bx} \cdot \frac{\sin ax}{bx} \right) = \lim_{x \to 0} \left(\frac{1}{\cos bx} \cdot \frac{\sin ax}{bx} \right) = 1 \cdot \frac{a}{b} = \frac{a}{b}$$





Example 5: Evaluate the following limit

 $\lim_{\theta\to 0}\frac{1-\cos\theta}{\theta\sin\theta}$

$$\lim_{\theta \to 0} \frac{1 - \cos \theta}{\theta \sin \theta} \left(\frac{0}{0} \right) = \lim_{\theta \to 0} \frac{1 - \cos \theta}{\theta \sin \theta} \cdot \frac{1 + \cos \theta}{1 + \cos \theta} = \lim_{\theta \to 0} \frac{1 - \cos^2 \theta}{\theta \sin \theta (1 + \cos \theta)} = \lim_{\theta \to 0} \frac{\sin^2 \theta}{\theta \sin \theta (1 + \cos \theta)}$$
$$= \lim_{\theta \to 0} \frac{\sin \theta}{\theta (1 + \cos \theta)} = \lim_{\theta \to 0} \frac{\sin \theta}{\theta} \cdot \lim_{\theta \to 0} \frac{1}{(1 + \cos \theta)} = (1)(\frac{1}{1 + 1}) = \frac{1}{2}$$

Example 6: Evaluate the following limit

K MAMA

 $\lim_{t\to\pi}\frac{\sin t}{\pi-t}$

Solution

$$\lim_{t \to \pi} \frac{\sin t}{\pi - t} = \left(\frac{0}{0}\right). \text{ Let } y = \pi - t \text{ , as } t \to \pi \text{ , } y \to 0. \text{ So}$$

 $\lim_{t \to \pi} \frac{\sin t}{\pi - t} = \lim_{y \to 0} \frac{\sin(\pi - y)}{y} = \lim_{y \to 0} \frac{\sin \pi \cos y - \cos \pi \sin y}{y} = \lim_{y \to 0} \frac{0 - (-1)\sin y}{y} = \lim_{y \to 0} \frac{\sin y}{y} = 1$

Example 7: Evaluate the following limit

```
\lim_{x \to 0} \frac{x \tan x + \cos 2x - 1}{x^2}
```

Solution

$$\lim_{x \to 0} \left[\frac{x \tan x + \cos 2x - 1}{x^2} \right] \left(\frac{0}{0} \right) = \lim_{x \to 0} \left[\frac{x \tan x + (1 - 2\sin^2 x) - 1}{x^2} \right] \text{ use the identity } \left[\cos(2x) = 1 - 2\sin^2 x \right] \\ = \lim_{x \to 0} \left[\frac{x \tan x - 2\sin^2 x}{x^2} \right] = \lim_{x \to 0} \left[\frac{x \tan x}{x^2} - \frac{2\sin^2 x}{x^2} \right] = \lim_{x \to 0} \left[\frac{\tan x}{x} - 2 \cdot \frac{\sin x}{x} \cdot \frac{\sin x}{x} \right] \\ = \lim_{x \to 0} \left[\frac{\tan x}{x} \right] - 2\lim_{x \to 0} \left[\frac{\sin x}{x} \right] \cdot \lim_{x \to 0} \left[\frac{\sin x}{x} \right] = 1 - (2)(1)(1) = -1$$

Example 8: Evaluate the following limit.

$$\lim_{x \to \pi/2} \frac{\cos(x)}{\cos(-x)}$$

$$\lim_{x \to \pi/2} \frac{\cos\left(x\right)}{\cos\left(-x\right)} \left(\frac{0}{0}\right) = \lim_{x \to \pi/2} \frac{\cos\left(x\right)}{\cos\left(x\right)} = \lim_{x \to \pi/2} 1 = 1$$

OF tRAMAR

Example 9: Evaluate the following limit

 $\lim_{x\to 1}\frac{\sin(3x-3)}{1-x^3}$

Solution

$$\lim_{x \to 1} \left[\frac{\sin(3x-3)}{1-x^3} \right] \left(\frac{0}{0} \right) = \lim_{x \to 1} \frac{\sin[3(x-1)]}{-(x-1)(1+x+x^2)} = \lim_{x \to 1} \left[\frac{-\sin[3(x-1)]}{(x-1)} \right] \cdot \lim_{x \to 1} \left[\frac{1}{1+x+x^2} \right]$$

To evaluate
$$\lim_{x \to 1} \frac{\sin 3(x-1)}{(x-1)}$$

Let y=x-1, as $x \to 1, y \to 0$.So

$$\lim_{x \to 1} \frac{\sin 3(x-1)}{(x-1)} = \lim_{y \to 0} \frac{\sin 3y}{y} = 3$$

92

Therefore,

$$\lim_{x \to 1} \left[\frac{\sin(3x-3)}{1-x^3} \right] \left(\frac{0}{0} \right) = \lim_{x \to 1} \frac{\sin[3(x-1)]}{-(x-1)(1+x+x^2)} = 3\left(-\frac{1}{3}\right) = -1$$

Example 10: Evaluate the following limit

$$\lim_{x \to 0} \frac{\tan 2x}{\sqrt{3x+1}-1}$$

$$\lim_{x \to 0} \frac{\tan 2x}{\sqrt{3x+1}-1} \left(\frac{0}{0}\right) = \lim_{x \to 0} \frac{\tan 2x}{\sqrt{3x+1}-1} \cdot \frac{\sqrt{3x+1}+1}{\sqrt{3x+1}+1} = \lim_{x \to 0} \frac{\tan 2x \left(\sqrt{3x+1}+1\right)}{(3x+1)-1}$$
$$= \lim_{x \to 0} \frac{\tan 2x \left(\sqrt{3x+1}+1\right)}{3x} = \lim_{x \to 0} \frac{\tan 2x}{3x} \cdot \lim_{x \to 0} (\sqrt{3x+1}+1) = (\frac{2}{3})(2) = \frac{4}{3}$$

Example 11: Evaluate the following limit

 $\lim_{x\to 1}\frac{\sin^2(\pi x)}{1-2x+x^2}$

Solution

 $\lim_{x \to 1} \frac{(\sin(\pi x))^2}{1 - 2x + x^2} \left(\frac{0}{0}\right) = \lim_{x \to 1} \frac{(\sin(\pi - \pi x))^2}{(1 - x)^2} = \lim_{x \to 1} \frac{(\sin(\pi(1 - x))^2)}{(1 - x)^2} = \lim_{x \to 1} \frac{\sin(\pi(1 - x))}{(1 - x)} \cdot \lim_{x \to 1} \frac{\sin(\pi(1 - x))}{(1 - x)}$ Let y = x - 1, as $x \to 1$, $y \to 0$. So

$$\lim_{x \to 1} \frac{\sin \pi (1-x)}{(1-x)} \cdot \lim_{x \to 1} \frac{\sin \pi (1-x)}{(1-x)} = \lim_{y \to 0} \frac{\sin \pi y}{y} \cdot \lim_{y \to 0} \frac{\sin \pi y}{y} = (\pi)(\pi) = \pi^2$$

Example 12: Find all values of k if $\lim_{x \to 0} \frac{\sin^2(kx)}{x^2} = 4$.

Solution

$$\lim_{x \to 0} \frac{\sin^2(kx)}{x^2} = 4 \Longrightarrow \lim_{x \to 0} \left[\frac{\sin(kx)}{x} \cdot \frac{\sin(kx)}{x} \right] = 4 \Longrightarrow k^2 = 4 \Longrightarrow k = -2, 2$$

Example 13: Evaluate the following limit.

$$\lim_{x \to \frac{\pi}{2}} \tan x$$

Solution

Notice that

34 | Page

x=5

 $\lim_{x \to \frac{\pi}{2}} \tan x = \lim_{x \to \frac{\pi}{2}} \frac{\sin x}{\cos x}$

The limit of the numerator is 1, and the limit of the denominator is 0. So the limit of the ratio does not exist. To be more specific than this, we need to analyze the sign of the ratio.

Let's take a look at the left-handed limit first. In this case we are going to be assuming that whatever x is it will be less than $\frac{\pi}{2}$. Therefore, as x gets closer and closer to $x = \frac{\pi}{2}$ the numerator is getting closer and closer to 1 while the denominator is getting closer and closer to 0 and will always be positive since we know that whatever x is it must satisfy $x < \frac{\pi}{2}$ (see the figure).

So, as we get closer and closer to $x = \frac{\pi}{2}$ (from the left) we have a positive, finite number in the numerator divided by an increasingly smaller positive number. This will result in increasing large and positive numbers. In other words,

 $\lim_{x \to \frac{\pi}{2}^{-}} \tan x = \infty$

The right-handed limit is similar. As we move in towards $x = \frac{\pi}{2}$ from the right we will always

have $x > \frac{\pi}{2}$ and so we will have a positive, finite number in the numerator divided by a increasingly smaller negative number and so the whole thing should be getting larger and larger. In this case the right-handed limit is,

 $\lim_{x \to \frac{\pi}{2^+}} \frac{\sin x}{\cos x} = -\infty$

So $\lim_{x \to \frac{\pi}{2}} \frac{\sin x}{\cos x}$ doesn't exist.

Theorem: Squeeze Theorem

Suppose that for all x on [a,b] we have,

 $g(x) \le f(x) \le h(x)$

Also suppose that,

 $\lim_{x \to \infty} g(x) = \lim_{x \to \infty} h(x) = L$

for some $a \leq c \leq b$. Then,

 $\lim_{x \to c} f(x) = L$

The Squeeze theorem is also known as the Sandwich Theorem and the Pinching Theorem.

NAGE

C ACCREDITED

So, how do we use this theorem to help us with limits? Let's take a look at the following example to see the theorem in action.

Example 14: Evaluate the following limit.

 $\lim_{x \to 0} x^2 \cos\left(\frac{1}{x}\right)$

Solution: In this example none of the previous examples can help us. There's no factoring or simplifying to do. We can't rationalize and one-sided limits won't work. There's even a question as to whether this limit will exist since we have division by zero inside the cosine at x = 0.

We know the following about cosine.

 $-1 \le \cos x \le 1$

We don't just have an x in the cosine, but as long as we avoid x = 0 we can say the same thing for our cosine.

$$-1 \le \cos\left(\frac{1}{x}\right) \le 1$$
 [for all $x \ne 0$]
Its okay for us to ignore x = 0 since we are taking a limit and we know that limits don't care about what's actually going on at x = 0 in this case.

Now if we have the above inequality for our cosine we can just multiply everything by an x^2 and get the following.

$$-x^{2} \leq x^{2} \cos\left(\frac{1}{x}\right) \leq x^{2}$$

In other words we've managed to squeeze the function that we where interested in between two other function that are very easy to deal with. So, the limits of the two outer functions are.

$$\lim_{x \to 0} x^{2} = 0 = \lim_{x \to 0} \left(-x^{2} \right)$$

These are the same and so by the Squeeze theorem we must also have,

$$\lim_{x \to 0} x^2 \cos\left(\frac{1}{x}\right) = 0$$

Remark: the squeezing theorem also holds for one sided limits and limits at ∞ and $-\infty$.

Example 15: Evaluate the following limit.

$$\lim_{x\to\infty} \left(e^{-3x} \cos 2x \right)$$

Solution: We know the following about cosine.

$$-1 \le \cos x \le 1$$

We can say the same thing for our cosine.

$$-1 \le \cos(2x) \le 1$$

Now if we have the above inequality for our cosine we can just multiply everything by an e^{-x} and get the following.

$$-e^{-x} \leq e^{-x} \cos(2x) \leq e^{-x}$$

The limits of the two outer functions are.

$$\lim_{x\to\infty} -e^{-x} = 0 = \lim_{x\to\infty} e^{-x}$$

These are the same and so by the Squeeze theorem we must also have,

$$\lim_{x\to\infty} \left(e^{-x} \cos\left(2x\right) \right) = 0$$

Example 16: Evaluate the following limit.

$$\lim_{x \to 0^+} \sqrt{x} \sin\left(x + \frac{1}{x}\right)$$

Solution: We know the following about sine.

$$-1 \le \sin x \le 1$$

We can say the same thing for our sine.

$$-1 \le \sin\left(x + \frac{1}{x}\right) \le 1$$
 [for all $x \ne 0$]

If x > 0, then $\sqrt{x} > 0$. Multiply through by \sqrt{x} , we get

$$-\sqrt{x} \le \sqrt{x} \sin\left(x + \frac{1}{x}\right) \le \sqrt{x}$$

In other words we've managed to squeeze the function that we where interested in between two other function that are very easy to deal with. So, the limits of the two outer functions are.

$$\lim_{x \to 0^+} -\sqrt{x} = 0 = \lim_{x \to 0^+} \sqrt{x}$$

These are the same and so by the Squeeze theorem we must also have,

$$\lim_{x \to 0^+} \sqrt{x} \cos\left(x + \frac{1}{x}\right) = 0$$

Example 17: Evaluate the following limit.

 $\lim_{x \to 0+} x^{1/3} \cos\left(2 + \frac{1}{x}\right)$

Solution: We know the following about cosine.

$$-1 \le \cos x \le 1$$

We can say the same thing for our cosine.

$$-1 \le \cos\left(2 + \frac{1}{x}\right) \le 1$$
 [for all $x \ne 0$

If x > 0, then $x^{1/3} > 0$ and so

$$-x^{1/3} \le x^{1/3} \cos\left(2 + \frac{1}{x}\right) \le x^{1/3}$$

The limits of the two outer functions are.

$$\lim_{x \to 0^+} -x^{1/3} = 0 = \lim_{x \to 0^+} x^{1/3}$$

These are the same and so by the Squeeze theorem we must also have,

$$\lim_{x \to 0+} x^{1/3} \cos\left(2 + \frac{1}{x}\right) = 0$$

Continuity

Definition

A function f(x) is said to be continuous at x = c if

1. f(c) is defined.

2. $\lim_{x \to c} f(x)$ is exist.

3. $\lim_{x \to c} f(x) = f(c)$

COPYRIGHT FIMT 2020

39 | Page

A function is said to be continuous on the interval [a,b] if it is continuous at each point in the interval.

Now, this definition justifies how we've been computing limits for awhile now and so it's good in that sense. However, it doesn't really tell us just what it means for a function to be continuous. Let's take a look at the following example to help us understand just what it means for a function to be continuous.

Example 1: Given the graph of f(x) shown below, determine if f(x) is continuous at x = 1, x = 2, and x = 3.

Solution: To answer the question for each point we'll need to get both the limit at that point and the function value at that point. If they are equal the function is continuous at that point and if they aren't equal the function isn't continuous at that point.



First x = 1.

f(1) = 1 lim f(x) doesn't exist

The function value and the limit aren't the same and so the function is not continuous at this point.

Now x = 2.

f(2): undefined

 $\lim_{x \to 1} f(x) = 2$

The function value and the limit aren't the same and so the function is not continuous at this point.

Finally x = 3.

f(3) = 2 $\lim_{x \to 2^{2}} f(x) = 2$

The function is continuous at this point since the function and limit have the same value.

Example 2: discuss the continuity at x = 3 for the following functions.

$$f(x) = \frac{x^2 - 9}{x - 3}$$

$$g(x) = \begin{cases} \frac{x^2 - 9}{x - 3}, & x \neq 3 \\ 4, & x = 3 \end{cases}$$

$$h(x) = \begin{cases} \frac{x^2 - 9}{x - 3}, & x \neq 3 \\ 6, & x = 3 \end{cases}$$
Solution

1. the functions is undefined at x = 3, and hence is not continuous at that point.

2.

$$f'(3) = 4$$
$$\lim_{x \to 3} f(x) = \lim_{x \to 3} \frac{x^2 - 9}{x - 3} = \lim_{x \to 3} \frac{(x - 3)(x + 3)}{x - 3} = \lim_{x \to 3} (x + 3) = 6$$

The function value and the limit aren't the same and so the function is discontinuous at x = 3.

3.

$$f(3) = 6$$
 $\lim_{x \to 3} f(x) = 6$

100

The function is continuous at x = 3 since the function and limit have the same value.

Example 3: Find values of a, b where the function

$$f(x) = \begin{cases} ax^{2} + 1, & x > 2 \\ -11, & x = 2 \\ x^{3} + b, & x < 2 \end{cases}.$$

Is continuous.

Solution: f is continuous every where. In particular f is continuous at x = 2. This implies that:

1. $\lim_{x \to 2} f(x)$ exist. This implies that:

 $\lim_{x \to 2+} f(x) = \lim_{x \to 2-} f(x)$ $\Rightarrow \lim_{x \to 2+} \left(ax^2 + 1 \right) = \lim_{x \to 2-} \left(x^3 + b \right)$ $\Rightarrow 4a + 1 = 8 + b \qquad (1)$

2. $\lim_{x \to 2^{+}} f(x) = f(2)$. this implies that

 $\lim_{x \to 2^{+}} f(x) = -11$ $\Rightarrow 4a + 1 = -11 \qquad (or \ 8 + b = -11)$ $\Rightarrow a = -3$

Substitute a = -3 in (1), getting

b = -19

Example 4: Find all values of the constant k that make

$$f(x) = \begin{cases} \frac{2x}{\tan kx} + 1, & x < 0\\ 3x + k, & x \ge 0 \end{cases}$$

Is continuous at x = 0.

Solution: To make f continuous at x = 0, we must have $\lim_{x \to 0} f(x)$ exist.

i.e

$$\lim_{x \to 0^+} f(x) = \lim_{x \to 0^-} f(x)$$
$$\Rightarrow \lim_{x \to 0^+} (3x+k) = \lim_{x \to 0^-} \left(\frac{2x}{\tan kx} + 1\right) \Rightarrow k = \left(\frac{2}{k}\right) + 1, k \neq 0$$
$$\Rightarrow k^2 = 2+k \Rightarrow k^2 - k - 2 = 0 \Rightarrow (k-2)(k+1) = 0 \Rightarrow k = 2 \text{ and } k = -1$$

Example 5: Find value(s) for the constant k so that

$$f(x) = \begin{cases} \frac{\sin kx}{x}, & x < 0\\ 3x + 2k^2, & x \ge 0 \end{cases}$$

will be continuous at x = 0.

Solution: To make f continuous at x = 0, we must have $\lim_{x \to 0} f(x)$ exist.

i.e

$$\lim_{x \to 0^+} f(x) = \lim_{x \to 0^-} f(x)$$
$$\Rightarrow \lim_{x \to 0^+} 3x + 2k^2 = \lim_{x \to 0^-} \frac{\sin kx}{x} \Rightarrow 2k^2 = k \Rightarrow 2k^2 - k = 0 \Rightarrow k (2k - 1) = 0 \Rightarrow k = 0, k = \frac{1}{2}$$

Theorem: Polynomials are continuous every where

Theorem: A rational function is continuous every where except at the point where the denominator is 0.

Example 6: discuss the points of discontinuity for the function.

$$f(x) = \frac{x^2 - 9}{x^2 - 5x + 6}$$

Solution: A rational function is continuous every where except at the point where the denominator is 0. Solving this equation

$$x^2 - 5x + 6 = 0$$

Yields two points of discontinuity x = 2 and x = 3.

Theorem: If $\lim g(x) = L$ and if the function f(x) is continuous at L, then

 $\lim_{x \to \infty} f(g(x)) = f(\lim_{x \to \infty} g(x))$

Remark: the above theorem also holds for one sided limits and limits at ∞ and $-\infty$.

With this fact we can now do limits like the following examples

Example 7: Evaluate the following limit.

 $\lim_{x\to 0} e^{\sin x}$

Solution: Since we know that exponentials are continuous we can use the fact above.

$$\lim_{x \to 0} e^{\sin x} = e^{\lim_{x \to 0} \sin x} = e^{0} = 1$$

Example 8: Compute $\lim_{x \to \infty} \log_{10} \left(\frac{x^6 - 500}{x^6 + 500} \right).$

Solution

 $\lim_{x \to \infty} \log_{10} \left(\frac{x^{6} - 500}{x^{6} + 500} \right) = \log_{10} \left[\lim_{x \to \infty} \left(\frac{x^{6} - 500}{x^{6} + 500} \right) \right]$

The previous step is valid because of the continuity of the logarithm function. Note also that the expression $\frac{x^6 - 500}{x^6 + 500}$ leads to the indeterminate form $\frac{\infty}{\infty}$. Circumvent it by dividing each term by x^6 , the highest power of x.

NAAC ACCREDITED

$$= \log_{10} \left[\lim_{x \to \infty} \left(\frac{\frac{x^{6}}{x^{6}} - \frac{500}{x^{6}}}{\frac{x^{6}}{x^{6}} + \frac{500}{x^{6}}} \right) \right] = \log_{10} \left[\lim_{x \to \infty} \left(\frac{1 - \frac{500}{x^{6}}}{1 + \frac{500}{x^{6}}} \right) \right] = \log_{10} \left(\frac{1 - 0}{1 + 0} \right) = \log_{10} 1 = 0$$

The term $\frac{500}{x^6}$ approaches 0 as x approaches ∞ .

Example 9: Evaluate the following limit.

$$\lim_{x \to 3} \sqrt{\frac{3x - 9}{2x^2 - 18}}$$

Solution
$$\lim_{x \to 3} \sqrt{\frac{3x - 9}{2x^2 - 18}} = \sqrt{\lim_{x \to 3} \frac{3x - 9}{2x^2 - 18}} \left(\frac{0}{0}\right)$$

previous step is valid because of the continuity of the square root function.

$$=\sqrt{\lim_{x \to 3} \frac{3(x-3)}{2(x^2-9)}} = \sqrt{\lim_{x \to 3} \frac{3(x-3)}{2(x-3)(x+3)}}$$
$$=\sqrt{\lim_{x \to 3} \frac{3}{2(x+3)}} = \sqrt{\frac{3}{2(3+2)}} = \sqrt{\frac{3}{12}} = \sqrt{\frac{1}{4}} = \frac{1}{2}$$

Example 10: Compute $\lim_{x \to \infty} \sqrt{\frac{x^3 + 7x}{4x^3 + 5}}$

Solution

$$\lim_{x \to \infty} \sqrt{\frac{x^3 + 7x}{4x^3 + 5}} = \sqrt{\lim_{x \to \infty} \frac{x^3 + 7x}{4x^3 + 5}} \left(\frac{\infty}{\infty}\right)$$

The previous step is valid because of the continuity of the square root function.

Inside the square root sign lies an indeterminate form. Circumvent it by dividing each term by x^3 , the highest power of x inside the square root sign.

CCREDITED

$$= \sqrt{\lim_{x \to \infty} \frac{\frac{x^3}{x^3} + \frac{7x}{x^3}}{\frac{4x^3}{x^3} + \frac{5}{x^3}}} = \lim_{x \to \infty} \sqrt{\lim_{x \to \infty} \frac{1 + \frac{7}{x^2}}{4 + \frac{5}{x^3}}} = \sqrt{\frac{1 + 0}{4 + 0}} = \frac{1}{2}$$

Each of the two expressions $\frac{7}{x^2}$ and $\frac{5}{x^3}$ approaches 0 as x approaches ∞ .

Another very nice consequence of continuity is the Intermediate Value Theorem.

Theorem: Intermediate Value Theorem

Suppose that f(x) is continuous on [a,b] and let M be any number between f(a) and f(b).

Then there exists a number $c \in (a, b)$ such that, f(c) = M.

Remark: All the Intermediate Value Theorem is really saying is that a continuous function will take on all values between f(a) and f(b). Below is a graph of a continuous function that illustrates the Intermediate Value Theorem.

As we can see from this image if we pick any value, M, that is between the value of f(a) and the value of f(b) and draw a line straight out from this point the line will hit the graph in at least one point. In other words somewhere between a and b the function will take on the value of M. Note



that the figure also shows that it may take on the value more than one place.

It's also important to note that the Intermediate value theorem only says that the function will take on the value of M somewhere between a and b. It doesn't say just what that value will be. It only says that it exists. It also does not tell us how many times the function may take on this value. It only tells us that it takes the value at least once.

A nice use of the Intermediate Value Theorem is to prove the existence of roots of equations.

Example 11: Show that $f(x) = x^3 + 2x - 15$ has a root somewhere in the interval [2,3].

Solution

What we're really asking here is whether or not the function will take on the value

$$f(x) = 0$$

somewhere between 2 and 3. In other words, we're using M = 0 in the Intermediate value theorem. All we need to show is that 0 is between f(2) and f(3) and we'll be done.

f(2) = -3 f(3) = 18

So by the Intermediate Value Theorem there must be a number c somewhere between 2 and 3 such that

f(c) = 0

Therefore the function does have a root between 2 and 3.

Example 12: If f and g are continuous on [a,b] and f(a) > g(a), f(b) < g(b), then there is at least one solution of the equation f(x) = g(x) in (a,b).

Solution

Let h(x) = f(x) - g(x). Since both f(x) and g(x) are both continuous on [a,b], then h(x) is continuous on [a,b].

Notice that h(a) = f(a) - g(a) > 0 and h(b) = f(b) - g(b) < 0

So by the Intermediate Value Theorem there must be a number c somewhere between a and b such that

h(c) = 0

This implies that f(c) - g(c) = 0 or f(c) = g(c).

Therefore the equation f(x) = g(x) has a solution in (a,b).

UNIT-III

Techniques of Differentiation

I. Notations for the Derivative

The derivative of y = f(x) may be written in any of the following ways:

$$f'(x)$$
, y' , $\frac{dy}{dx}$, $\frac{d}{dx}[f(x)]$, or $D_x[f(x)]$.

- II. Basic Differentiation Rules
 - A. Suppose *c* and *n* are constants, and f *and* g are differentiable functions.

$$(1) \quad f(x) = cg(x)$$

$$f'(x) = \lim_{b \to x} \frac{f(b) - f(x)}{b - x} = \lim_{b \to x} \frac{cg(b) - cg(x)}{b - x} = c \lim_{b \to x} \frac{g(b) - g(x)}{b - x} = cg'(x)$$

(2)
$$f(x) = g(x) \pm k(x)$$

$$f'(x) = \lim_{b \to x} \frac{f(b) - f(x)}{b - x} = \lim_{b \to x} \frac{[g(b) \pm k(b)] - [g(x) \pm k(x)]}{b - x} =$$

$$\lim_{b \to x} \frac{g(b) - g(x)}{b - x} \pm \lim_{b \to x} \frac{k(b) - k(x)}{b - x} = g'(x) \pm k'(x)$$

(3)
$$f(x) = g(x)k(x)$$

$$f'(x) = \lim_{b \to x} \frac{f(b) - f(x)}{b - x} = \lim_{b \to x} \frac{g(b)k(b) - g(x)k(x)}{b - x} =$$
$$\lim_{b \to x} \frac{g(b)k(b) - g(b)k(x) + g(b)k(x) - g(x)k(x)}{b - x} =$$
$$\left[\lim_{b \to x} g(b)\right] \left[\lim_{b \to x} \frac{k(b) - k(x)}{b - x}\right] + \left[\lim_{b \to x} k(x)\right] \left[\lim_{b \to x} \frac{g(b) - g(x)}{b - x}\right] =$$

g(x)k'(x) + k(x)g'(x) (Product Rule)

(4)
$$f(x) = \frac{g(x)}{k(x)} \Rightarrow f(x)k(x) = g(x) \Rightarrow g'(x) = f(x)k'(x) + k(x)f'(x) \Rightarrow$$

$$f'(x) = \frac{g'(x) - f(x)k'(x)}{k(x)} = \frac{g'(x) - \left[\frac{g(x)}{k(x)}\right]k'(x)}{k(x)} = \frac{k(x)g'(x) - g(x)k'(x)}{[k(x)]^2}.$$

This derivative rule is called the **Quotient Rule.** (5) f(x) = c

$$f'(x) = \lim_{b \to x} \frac{f(b) - f(x)}{b - x} = \lim_{b \to x} \frac{c - c}{b - x} = \lim_{b \to x} \frac{0}{b - x} = \lim_{b \to x} 0 = 0$$

(6) f(x) = x

$$f'(x) = \lim_{b \to x} \frac{f(b) - f(x)}{b - x} = \lim_{b \to x} \frac{b - x}{b - x} = \lim_{b \to x} 1 = 1$$

(7) $f(x) = x^n$

$$f'(x) = \lim_{h \to 0} \frac{f(x+h) - f(x)}{h} = \lim_{h \to 0} \frac{(x+h)^n - x^n}{h} =$$

$$\lim_{h \to 0} \frac{\left[x^{n} + nx^{n-1}h + \frac{n(n-1)}{2}x^{n-2}h^{2} + \dots \right] - x^{n}}{h} =$$

$$\lim_{h \to 0} \left[\frac{nx^{n-1}h + h^{2} \left(\frac{n(n-1)}{2}x^{n-2} + \dots \right)}{h} \right] =$$

$$\lim_{h \to 0} \left[nx^{n-1} + h \left(\frac{n(n-1)}{2}x^{n-2} + \dots \right) \right] = nx^{n-1} \text{ (Power Rule)}$$

Example 1: Suppose f and g are differentiable functions such that f(1) = 3,

$$g(1) = 7$$
, $f'(1) = -2$, and $g'(1) = 4$. Find (i) $(f + g)'(1)$, (ii) $(g - f)'(1)$,

(iii)
$$(fg)'(1)$$
, (iv) $\left(\frac{g}{f}\right)'(1)$, and $\left(\frac{f}{g}\right)'(1)$.
(i) $(f+g)'(1) = f'(1) + g'(1) = -2 + 4 = 2$
(ii) $(g-f)'(1) = g'(1) - f'(1) = 4 - (-2) = 6$
(iii) $(fg)'(1) = f(1)g'(1) + g(1)f'(1) = 3(4) + 7(-2) = 12 + (-14) = -2$

(iv)
$$\left(\frac{g}{f}\right)'(1) = \frac{f(1)g'(1) - g(1)f'(1)}{[f(1)]^2} = \frac{3(4) - 7(-2)}{3^2} = \frac{12 + 14}{9} = \frac{26}{9}$$

(v)
$$\left(\frac{f}{g}\right)(1) = \frac{g(1)f'(1) - f(1)g'(1)}{[g(1)]^2} = \frac{7(-2) - 3(4)}{7^2} = \frac{-14 - 12}{49} = \frac{-26}{49}$$

Example 2: If $f(x) = x^4 - 3x^3 + 5x^2 - 7x + 11$, find f'(x).

$$f'(x) = 4x^3 - 3(3x^2) + 5(2x) - 7(1) + 0 = 4x^3 - 9x^2 + 10x - 7$$

Example 3: If $f(x) = 4\sqrt{x} - \frac{3}{\sqrt[3]{x^2}} + \frac{5}{x} - \frac{7}{x^5}$, then find f'(x).

$$f(x) = 4\sqrt{x} - \frac{3}{\sqrt[3]{x^2}} + \frac{5}{x} - \frac{7}{x^5} = 4x^{\frac{1}{2}} - 3x^{-\frac{2}{3}} + 5x^{-1} - 7x^{-5} \Rightarrow$$

$$f'(x) = 4\left(\frac{1}{2}x^{-\frac{1}{2}}\right) - 3\left(-\frac{2}{3}x^{-\frac{5}{3}}\right) + 5\left(-1x^{-2}\right) - 7\left(-5x^{-6}\right) =$$

$$2x^{-\frac{1}{2}} + 2x^{-\frac{5}{3}} - 5x^{-2} + 35x^{-6} = \frac{2}{\sqrt{x}} + \frac{2}{\sqrt[3]{x^5}} - \frac{5}{x^2} + \frac{35}{x^6}$$

Example 4: If
$$f(x) = \frac{x^2 + 2x - 3}{3x - 4}$$
, then find $f'(1)$.

$$f'(x) = \frac{(3x-4)(2x+2) - (x^2 + 2x - 3)(3)}{(3x-4)^2} = \frac{6x^2 - 2x - 8 - 3x^2 - 6x + 9}{(3x-4)^2} =$$

3

$$\frac{3x^2 - 8x + 1}{(3x - 4)^2} \Rightarrow f'(1) = \frac{3(1)^2 - 8(1) + 1}{[3(1) - 4]^2} = \frac{-4}{1} = -4 \text{ or}$$

$$f'(1) = \frac{[3(1) - 4][2(1) + 2] - [1^2 + 2(1) - 3](3)}{[3(1) - 4]^2} = \frac{(-1)(4) - (0)(3)}{(-1)^2} = \frac{-4}{1} = -4$$

2

B. Trigonometric functions

-

(1)
$$f(x) = \sin x$$
$$f'(x) = \lim_{h \to 0} \frac{f(x+h) - f(x)}{h} = \lim_{h \to 0} \frac{\sin(x+h) - \sin x}{h} =$$
$$\lim_{h \to 0} \frac{\sin x \cosh + \cos x \sinh - \sin x}{h} = \lim_{h \to 0} \frac{\sin x (\cosh - 1) + \cos x \sinh - 1}{h} =$$

$$(\sin x)\left[\lim_{h \to 0} \frac{\cosh(-1)}{h}\right] + (\cos x)\left[\lim_{h \to 0} \frac{\sinh(-1)}{h}\right] = (\sin x)(0) + (\cos x)(1) =$$

$$\cos x$$

(2)
$$f(x) = \cos x$$

$$f'(x) = \lim_{h \to 0} \frac{f(x+h) - f(x)}{h} = \lim_{h \to 0} \frac{\cos(x+h) - \cos x}{h} =$$

$$\lim_{h \to 0} \frac{\cos x \cosh - \sin x \sinh - \cos x}{h} = \lim_{h \to 0} \frac{\cos x (\cosh - 1) - \sin x \sinh h}{h} =$$

$$(\cos x) \left[\lim_{h \to 0} \frac{\cosh - 1}{h} \right] - (\sin x) \left[\lim_{h \to 0} \frac{\sinh h}{h} \right] = (\cos x) (0) - (\sin x) (1) =$$

$$-\sin x$$

$$(3) \quad f(x) = \tan x = \frac{\sin x}{\cos x}$$

$$f'(x) = \frac{(\cos x)(\cos x) - (\sin x)(-\sin x)}{(\cos x)^2} = \frac{\cos^2 x + \sin^2 x}{\cos^2 x} = \frac{1}{\cos^2 x} = \sec^2 x$$

$$(4) \quad f(x) = \sec x = \frac{1}{\cos x}$$

$$f'(x) = \frac{(\cos x)(0) - 1(-\sin x)}{(\cos x)^2} = \frac{\sin x}{\cos^2 x} = \frac{1}{\cos x} \cdot \frac{\sin x}{\cos x} = \sec x \tan x$$

$$(5) \quad f(x) = \csc x = \frac{1}{\sin x}$$

$$f'(x) = \frac{(\sin x)(0) - 1(\cos x)}{(\sin x)^2} = \frac{-\cos x}{\sin^2 x} = \frac{-1}{\sin x} \cdot \frac{\cos x}{\sin x} = -\csc x \cot x$$

COPYRIGHT FIMT 2020

52 | Page

(6)
$$f(x) = \cot x = \frac{\cos x}{\sin x}$$

$$f'(x) = \frac{(\sin x)(\sin x) - (\cos x)(\cos x)}{(\sin x)^2} = \frac{-\cos^2 x - \sin^2 x}{\sin^2 x} = \frac{-1}{\sin^2 x} = -\csc^2 x$$

AGEME

C. Composition and the generalized derivative rules

(1)
$$f(x) = (g \circ k)(x) = g(k(x))$$

$$f'(x) = \lim_{b \to x} \frac{f(b) - f(x)}{b - x} = \lim_{b \to x} \frac{g(k(b)) - g(k(x))}{b - x} = \lim_{b \to x} \frac{g(k(b)) - g(k(x))}{b - x}$$

$$\frac{k(b) - k(x)}{k(b) - k(x)} = \lim_{b \to x} \frac{g(k(b)) - g(k(x))}{k(b) - k(x)} \cdot \lim_{b \to x} \frac{k(b) - k(x)}{b - x} =$$

 $\lim_{k(b) \to k(x)} \frac{g(k(b)) - g(k(x))}{k(b) - k(x)} \cdot \lim_{b \to x} \frac{k(b) - k(x)}{b - x} = g'(k(x)) \cdot k'(x).$

This derivative rule for the composition of functions is called the **Chain Rule**.

(2) Suppose that f(x) = g(k(x)) where $g(x) = x^n$. Then $f(x) = [k(x)]^n$.

$$g(x) = x^n \Rightarrow g'(x) = nx^{n-1} \Rightarrow g'(k(x)) = n[k(x)]^{n-1}$$
. Thus, $f'(x) =$

 $g'(k(x)) \cdot k'(x) = n[k(x)]^{n-1} \cdot k'(x)$. This derivative rule for the power of a function is called the **Generalized Power Rule**.

(3) Suppose that f(x) = g(k(x)) where $g(x) = \sin x$. Then $f(x) = \sin[k(x)]$.

$$g(x) = \sin x \Rightarrow g'(x) = \cos x \Rightarrow g'(k(x)) = \cos[k(x)]$$
. Thus, $f'(x) = \cos[k(x)]$.

$$g'(k(x)) \cdot k'(x) = \cos[k(x)] \cdot k'(x).$$

(4) Similarly, if $f(x) = \cos[k(x)]$, then $f'(x) = -\sin[k(x)] \cdot k'(x)$.

(5) If $f(x) = \tan[k(x)]$, then $f'(x) = \sec^2[k(x)] \cdot k'(x)$.

(6) If $f(x) = \sec[k(x)]$, then $f'(x) = \sec[k(x)] \tan[k(x)] \cdot k'(x)$.

(7) If $f(x) = \cot[k(x)]$, then $f'(x) = -\csc^2[k(x)] \cdot k'(x)$.

(8) If $f(x) = \csc[k(x)]$, then $f'(x) = -\csc[k(x)]\cot[k(x)] \cdot k'(x)$

Example 1: Suppose f and g are differentiable functions such that:

f(1) = 9 f(2) = -5 g(1) = 2 g(9) = 3

$$f'(1) = -2$$
 $f'(2) = -6$ $g'(1) = 4$ $g'(9) = 7$

Find each of the following: (i) $(f \circ g)'(1)$; (ii) $(g \circ f)'(1)$; (iii) h'(1) if

$$h(x) = \sqrt{f(x)}$$
; (iv) $j'(1)$ if $j(x) = [g(x)]^5$; (v) $l'(1)$ if $l(x) = \frac{3}{[f(x)]^2}$;

(vi) s'(1) if $s(x) = \sin[f(x)]$; and (vii) m'(1) if $m(x) = \sec[g(x)]$.

(i)
$$(f \circ g)'(1) = f'(g(1)) \cdot g'(1) = f'(2) \cdot g'(1) = (-6)(4) = -24$$

NAAL ALLKEUH

(ii)
$$(g \circ f)(1) = g(f(1)) \cdot f(1) = g(9) \cdot f(1) = 7(-2) = -14$$

(iii)
$$h(x) = \sqrt{f(x)} = [f(x)]^{\frac{1}{2}} \Rightarrow h'(x) = \frac{1}{2} [f(x)]^{-\frac{1}{2}} \cdot f'(x) = \frac{f'(x)}{2\sqrt{f(x)}} \Rightarrow$$

 $h'(1) = \frac{f'(1)}{2\sqrt{f(1)}} = \frac{-2}{2\sqrt{9}} = -\frac{1}{3}$

(iv)
$$j(x) = [g(x)]^5 \Rightarrow j'(x) = 5[g(x)]^4 \cdot g'(x) \Rightarrow j'(1) = 5[g(1)]^4 \cdot g'(1) =$$

 $5(2)^4(4) = 320$

(v)
$$l(x) = \frac{3}{[f(x)]^2} = 3[f(x)]^{-2} \Rightarrow l'(x) = -6[f(x)]^{-3} \cdot f'(x) \Rightarrow l'(1) =$$

$$\frac{-6f'(1)}{[f(1)]^3} = \frac{-6(-2)}{9^3} = \frac{12}{729} = \frac{4}{243}$$

(vi) $s'(x) = \cos[f(x)] \cdot f'(x) \Rightarrow s'(1) = \cos[f(1)] \cdot f'(1) = \cos(9) \cdot (-2) = -2\cos 9$

(vii) $m'(x) = \sec[g(x)] \tan[g(x)] \cdot g'(x) \Longrightarrow m'(1) = \sec[g(1)] \tan[g(1)] \cdot g'(1) =$

 $\sec(2)\tan(2)\cdot 4 = 4\sec 2\tan 2$

10.00

Example 2: If $f(x) = \sqrt[3]{2x^4 - x^2 + 5x + 2}$, then find f'(1).

$$f(x) = \sqrt[3]{2x^4 - x^2 + 5x + 2} = (2x^4 - x^2 + 5x + 2)^{1/3} \implies f'(x) =$$

$$\frac{1}{3}(2x^4 - x^2 + 5x + 2)^{-\frac{2}{3}}(8x^3 - 2x + 5) = \frac{8x^3 - 2x + 5}{3\sqrt[3]{(2x^4 - x^2 + 5x + 2)^2}} \Rightarrow$$

$$f'(1) = \frac{8 - 2 + 5}{3\sqrt[3]{(2 - 1 + 5 + 2)^2}} = \frac{11}{3\sqrt[3]{64}} = \frac{11}{12}$$

Example 3: If $g(x) = \frac{4}{(x^3 + 4)^8}$, then find g'(x).

$$g(x) = \frac{4}{(x^3 + 4)^8} = 4(x^3 + 4)^{-8} \Rightarrow g'(x) = -32(x^3 + 4)^{-9}(3x^2) = \frac{-96x^2}{(x^3 + 4)^9}$$

RELU

Example 4: If
$$h(x) = \sin(\cos x)$$
, then find $h'(x)$.
 $h'(x) = \cos(\cos x) \cdot (-\sin x)$

Example 5: If $j(x) = \tan(2x^2 - 3x + 1)$, then find j'(x).

$$j'(x) = \sec^2(2x^2 - 3x + 1) \cdot (4x - 3)$$

COPYRIGHT FIMT 2020

56 | Page

Example 6: If $k(x) = x^2 \sqrt{3x+4}$, then find k'(x).

$$k(x) = x^2 \sqrt{3x+4} = x^2 (3x+4)^{1/2} \Longrightarrow k'(x) = x^2 \left[\frac{1}{2} (3x+4)^{-1/2} (3) \right] +$$

$$(3x+4)^{\frac{1}{2}}(2x) = \frac{3x^2}{2(3x+4)^{\frac{1}{2}}} + \frac{2x(3x+4)^{\frac{1}{2}}}{1} = \frac{3x^2+4x(3x+4)}{2(3x+4)^{\frac{1}{2}}} =$$

 $\frac{15x^2 + 16x}{2(3x+4)^{1/2}}$

Example 7: If $l(x) = \left(\frac{2x-1}{3x+4}\right)^4$, then find l'(x).

$$l'(x) = 4\left(\frac{2x-1}{3x+4}\right)^3 \left[\frac{(3x+4)(2) - (2x-1)(3)}{(3x+4)^2}\right] = \frac{4(2x-1)^3}{(3x+4)^3} \left[\frac{11}{(3x+4)^2}\right] =$$

RELE

 $\frac{44(2x-1)^3}{(3x+4)^5}.$

Example 8: If $k(x) = \frac{\sin x}{1 + \cos x}$, then find k'(x).

$$k'(x) = \frac{(1+\cos x)(\cos x) - (\sin x)(-\sin x)}{(1+\cos x)^2} = \frac{\cos x + \cos^2 x + \sin^2 x}{(1+\cos x)^2} =$$

 $\frac{\cos x + 1}{(1 + \cos x)^2} = \frac{1}{1 + \cos x}.$

Example 9: If $s(x) = \sin^{3}(x^{2} - 1)$, then find s'(x).

$$s(x) = \sin^3 (x^2 - 1) = [\sin(x^2 - 1)]^3 \Rightarrow s'(x) = 3[\sin(x^2 - 1)]^2 \cdot \cos(x^2 - 1) \cdot 2x = 6x \sin^2 (x^2 - 1) \cos(x^2 - 1).$$

III. Implicit Differentiation

Example 1: Find the slope of the tangent line to the circle $x^2 + y^2 = 25$ at the

AAGE





(0, - 5)

Solution 1 : A circle is not a function. However, $x^2 + y^2 = 25 \Longrightarrow y^2 =$

 $25 - x^2 \Rightarrow y = \pm \sqrt{25 - x^2} \Rightarrow y = \sqrt{25 - x^2}$ is the equation of the upper

half circle and $y = -\sqrt{25 - x^2}$ is the equation of the lower half circle.

Since the point (3, 4) is on the upper half circle, use the function f(x) =

$$\sqrt{25 - x^2} = \left(25 - x^2\right)^{\frac{1}{2}} \Rightarrow f'(x) = \frac{1}{2}\left(25 - x^2\right)^{-\frac{1}{2}}(-2x) = \frac{-x}{\sqrt{25 - x^2}} \Rightarrow$$
$$m = f'(3) = \frac{-3}{\sqrt{25 - 3^3}} = \frac{-3}{\sqrt{25 - 9}} = \frac{-3}{\sqrt{16}} = -\frac{3}{4}.$$

Sometimes, an equation $[x^2 + y^2 = 25]$ in two variables, say x and y, is given, but it is not in the form of y = f(x). In this case, for each value of one of the variables, one or more values of the other variable may exist. Thus, such an equation may

describe one or more functions [$y = \sqrt{25 - x^2}$ and $y = -\sqrt{25 - x^2}$]. Any function

defined in this manner is said to be defined implicitly. For such equations, we may not be able to solve for y explicitly in terms of x [in the example, I was able to solve for yexplicitly in terms of x]. In fact, there are applications where it is not essential to obtain a formula for y in terms of x. Instead, the value of the derivative at certain

points must be obtained. It is possible to accomplish this goal by using a technique

called implicit differentiation. Suppose an equation in two variables, say x and y, is

given and we are told that this equation defines a differentiable function f with y = f(x). Use the following steps to differentiate implicitly:

(1) Simplify the equation if possible. That is, get rid of parentheses by multiplying using the distributive property or by redefining subtraction, and clear fractions by multiplying every term of the equation by a common

denominator for all the fractions; simplify and combine like terms.

(2) Differentiate both sides of the equation with respect to *x*. Use all the relevant differentiation rules, being careful to use the **Chain Rule** when differentiating

expressions involving y.

(3) Solve for $\frac{dy}{dx}$.

Note: It might be helpful to substitute f(x) into the equation for y before

differentiating with respect to x. This will remind you when you must use the

generalized forms of the **Chain Rule**. Since $f'(x) = \frac{dy}{dx}$, you differentiate

with respect to x and substitute y for f(x) and $\frac{dy}{dx}$ for f'(x). Then you can

solve for $\frac{dy}{dx}$.

Solution 2:
$$x^2 + y^2 = 25 \Rightarrow x^2 + [f(x)]^2 = 25 \Rightarrow \frac{d}{dx} \left(x^2 + [f(x)]^2 = 25 \right) \Rightarrow$$

 $2x + 2[f(x)]f'(x) = 0 \Rightarrow f'(x) = \frac{-2x}{2[f(x)]} \Rightarrow \frac{dy}{dx} = \frac{-x}{y} \Rightarrow \frac{dy}{dx} \Big|_{\substack{x=3\\y=4}} = -\frac{3}{4}.$
Example 2: Suppose that the equation $\frac{2}{x} + \frac{3}{y} = x$ defines a function f with $y = f(x)$

Find $\frac{dy}{dx}$ and the slope of the tangent line at the point (2, 3).

Solution 1: Solve for y.
$$xy\left(\frac{2}{x} + \frac{3}{y}\right) = xy(x) \Rightarrow 2y + 3x = x^2y \Rightarrow y = \frac{3x}{x^2 - 2} \Rightarrow$$

 $\frac{dy}{dx} = \frac{(x^2 - 2)(3) - 3x(2x)}{(x^2 - 2)^2} = \frac{-3x^2 - 6}{(x^2 - 2)^2} \Rightarrow \frac{dy}{dx}\Big|_{x=2} = \frac{-18}{4} = -\frac{9}{2}$
Solution 2: Clear fractions $\Rightarrow 2y + 3x = x^2y \Rightarrow \frac{d}{dx}(2y + 3x = x^2y) \Rightarrow$
 $2\frac{dy}{dx} + 3 = x^2\frac{dy}{dx} + 2xy \Rightarrow \frac{dy}{dx} = \frac{3 - 2xy}{x^2 - 2} \Rightarrow \frac{dy}{dx}\Big|_{x=2} = \frac{3 - 12}{2} = -\frac{9}{2}$
Solution 3: $\frac{d}{dx}\left(\frac{2}{x} + \frac{3}{y} = x\right) \Rightarrow \frac{d}{dx}\left(2x^{-1} + 3y^{-1} = x\right) \Rightarrow -2x^{-2} - 3y^{-2}\frac{dy}{dx} = 1 \Rightarrow$
 $\frac{-2}{x^2} - \frac{3}{y^2}\frac{dy}{dx} = 1 \Rightarrow -2y^2 - 3x^2\frac{dy}{dx} = x^2y^2 \Rightarrow \frac{dy}{dx} = \frac{-2y^2 - x^2y^2}{3x^2} \Rightarrow$
 $\frac{dy}{dx}\Big|_{x=2} = \frac{-18 - 36}{12} = \frac{-54}{12} = -\frac{9}{2}$

Example 3: If $\cos(xy) = y$, then find $\frac{dy}{dx}$

$$\frac{d}{dx}(\cos(xy) = y) \Rightarrow -\sin(xy)\left[x\frac{dy}{dx} + y(1)\right] = \frac{dy}{dx} \Rightarrow -x\sin(xy)\frac{dy}{dx} - y\sin(xy) = \frac{dy}{dx} \Rightarrow -y\sin(xy) = \frac{dy}{dx}(1 + x\sin(xy)) \Rightarrow \frac{dy}{dx} = \frac{-y\sin(xy)}{1 + x\sin(xy)}$$

IV. Higher Order Derivatives

A. Notation

(1) 1st derivative (derivative of the original function y = f(x)): $\frac{dy}{dx} = f'(x)$

(2) 2nd derivative (derivative of the 1st derivative):
$$\frac{d^2y}{dx^2} = f''(x)$$

(3) 3rd derivative (derivative of the 2nd derivative): $\frac{d^3y}{dx^3} = f'''(x)$

B. Distance functions

Suppose s(t) is a distance function with respect to time t. Then s'(t) = v(t)is an instantaneous velocity (or velocity) function with respect to time t, and s''(t) = v'(t) = a(t) is an acceleration function with respect to time t.

AAGEMEN

Example 1: If $f(x) = x^2 \sin x$, then find f'(x) and f''(x).

 $f'(x) = x^2 \cos x + 2x \sin x$

 $f''(x) = x^{2}(-\sin x) + 2x\cos x + 2x\cos x + 2\sin x = -x^{2}\sin x + 4x\cos x + 2\sin x$

Example 2: If $g(x) = \frac{2x+3}{4x-5}$, then find g'(x) and g''(x).

150 9001:2015 & 14001:2015

$$g'(x) = \frac{(4x-5)(2) - (2x+3)(4)}{(4x-5)^2} = \frac{8x-10-8x-12}{(4x-5)^2} = \frac{-22}{(4x-5)^2} = -22(4x-5)^{-2}$$

$$g''(x) = 44(4x-5)^{-3}(4) = 176(4x-5)^{-3} = \frac{176}{(4x-5)^3}$$

Example 3: If $x^2 + y^2 = 25$, then find $\frac{dy}{dx}$ and $\frac{d^2y}{dx^2}$.

$$\frac{d}{dx}\left(x^{2} + y^{2} = 25\right) \Rightarrow 2x + 2y\frac{dy}{dx} = 0 \Rightarrow \frac{dy}{dx} = \frac{-2x}{2y} = \frac{-x}{y}$$
$$\frac{d^{2}y}{dx^{2}} = \frac{d}{dx}\left(\frac{dy}{dx}\right) = \frac{d}{dx}\left(\frac{-x}{y}\right) = \frac{y(-1) - (-x)\left(\frac{dy}{dx}\right)}{y^{2}} = \frac{-y + x\left(\frac{-x}{y}\right)}{y^{2}} = \frac{-y^{2} - x^{2}}{y^{3}} = \frac{-(x^{2} + y^{2})}{y^{3}} = \frac{-25}{y^{3}}$$

Practice Sheet - Techniques of Differentiation

I. Find the derivative of each function defined as follows; there is no need to simplify your answers.

 $\frac{x^3}{y^3}$

(1)
$$f(x) = x^4 - 5x^3 + 9x^2 - 7x + 5$$
 (2) $y = \frac{9}{x^2} - \frac{8}{x^3} + \frac{2}{x^4}$

(3)
$$g(x) = 8\sqrt{x} - \frac{6}{\sqrt[3]{x^2}}$$
 (4) $y = \frac{3x^2 - 6x}{x^3}$

(5)
$$h(x) = \frac{3x+2}{x^2+1}$$
 (6) $y = x^2 \cos x$

(7)
$$f(x) = \frac{\sin x}{x}$$
 (8) $y = \sqrt[3]{x^2 - 3x + 4}$

(9)
$$g(x) = \sin(\sqrt{x})$$
 (10) $y = \cos^3 x$

(11)
$$h(x) = \sqrt{\frac{2x+1}{3x-4}}$$
 (12) $y = \frac{\sec x}{1+\tan x}$

HAGEME

- 77

4

(13)
$$k(x) = x\sqrt{9-x^2}$$
 (14) $y = \sin(3x)\cos(4x)$

49

(15)
$$f(x) = \tan^4(x^3)$$
 (16) $y = \frac{\sqrt{x}+1}{\sqrt{x}-1}$

(17)
$$g(x) = x \sec\left(\frac{1}{x}\right)$$
 (18) $y = \sqrt{1 + \sin 2x}$

II. Find
$$\frac{dy}{dx}$$
 by implicit differentiation for each of the following:

(1)
$$-3xy-4y^2 = 2$$
 (2) $8x^2 = 2y^3 + 3xy^2$

(3)
$$\frac{3}{2x} + \frac{1}{y} = y$$
 (4) $3x^2 = \frac{2-y}{2+y}$

(5)
$$x = \tan y$$
 (6) $y = \cos(x - y)$

(7) $x \sin y + y \sin x = 1$ (8) $x = \sec^3(y^2 - 1)$ NAAC ACCREDITED

III. Find the slope of the tangent line at the given point on each curve defined by the DANAG

given equation:

(1) $x^{2} + 3y^{2} = 21;$ (3, -2) (2) $x^{3} + \sqrt[3]{y} = 3;$ (1, 8) (3) $\sqrt{xy} - y = -2;$ (1, 4) (4) $3xy - 2x^{4} = y^{3} - 23;$ (2, -3) (5) $x = \cos y;$ $\left(\frac{1}{2}, \frac{-\pi}{3}\right)$ (6) $\sin(xy) = x;$ $\left(1, \frac{\pi}{2}\right)$

IV. For each of the following functions f(x), find f'(x) and f''(x).

(1)
$$f(x) = 3x^4 - 4x^2 + 7x - 11$$

(2)
$$f(x) = \frac{3x+1}{2x-1}$$

(3)
$$f(x) = x^3 \cos(4x)$$

(4) $f(x) = \sin^4 x$

V. Suppose the distance (in feet) that an object travels in t seconds is given by the

formula $s(t) = 2t^3 + 4t - 5$. Find s(2), v(2), and a(2).

Solution Key for Techniques of Differentiation

1. (1)
$$f'(x) = 4x^3 - 15x^2 - 18x - 7$$

(2) $y = 9x^{-2} - 8x^{-3} + 2x^{-4} \Rightarrow \frac{dy}{dx} = -18x^{-3} + 24x^{-4} - 8x^{-5}$
(3) $g(x) = 8x^{\frac{1}{2}} - 6x^{-\frac{2}{3}} \Rightarrow g'(x) = 4x^{-\frac{1}{2}} + 4x^{-\frac{5}{3}}$
(4) $y = 3x^{-1} - 6x^{-2} \Rightarrow \frac{dy}{dx} = -3x^{-2} + 12x^{-3}$
(5) $h'(x) = \frac{(x^2 + 1)(3) - (3x + 2)(2x)}{(x^2 + 1)^2}$
(6) $\frac{dy}{dx} = x^2(-\sin x) + (\cos x)(2x)$
(7) $f'(x) = \frac{x(\cos x) - (\sin x)(1)}{x^2}$
(8) $y = (x^2 - 3x + 4)^{\frac{1}{3}} \Rightarrow \frac{dy}{dx} = \frac{1}{3}(x^2 - 3x + 4)^{-\frac{2}{3}}(2x - 3)$
(9) $g'(x) = \cos(\sqrt{x}) \cdot (\frac{1}{2}x^{-\frac{1}{2}})$
(10) $y = (\cos x)^3 \Rightarrow \frac{dy}{dx} = 3(\cos x)^2(-\sin x)$
(11) $h(x) = (\frac{2x + 1}{3x - 4})^{\frac{1}{2}} \Rightarrow h'(x) = \frac{1}{2}(\frac{2x + 1}{3x - 4})^{-\frac{1}{2}}[\frac{(3x - 4)(2) - (2x + 1)(3)}{(3x - 4)^2}]$

(11)
$$h(x) = \left(\frac{2x+1}{3x-4}\right)^{1/2} \Rightarrow h'(x) = \frac{1}{2} \left(\frac{2x+1}{3x-4}\right)^{-1/2} \left[\frac{(3x-4)(2) - (2x+1)(3)}{(3x-4)^2}\right]^{-1/2}$$

(12)
$$\frac{dy}{dx} = \frac{(1 + \tan x)(\sec x \tan x) - \sec x(\sec^2 x)}{(1 + \tan x)^2}$$

(13)
$$k(x) = x\left(9 - x^2\right)^{\frac{1}{2}} \Rightarrow k'(x) = x\left[\frac{1}{2}\left(9 - x^2\right)^{-\frac{1}{2}}(-2x)\right] + \left(9 - x^2\right)^{\frac{1}{2}}(1)$$

(14)
$$\frac{dy}{dx} = \sin(3x) \left[-\sin(4x)(4) \right] + \cos(4x) \left[\cos(3x)(3) \right]$$

(15)
$$f(x) = \left[\tan(x^3) \right]^4 \implies f'(x) = 4 \left[\tan(x^3) \right]^3 \left[\sec^2(x^3) \right] (3x^2)$$

(16)
$$\frac{dy}{dx} = \frac{\left(\sqrt{x} - 1\left(\frac{1}{2\sqrt{x}}\right) - \left(\sqrt{x} + 1\left(\frac{1}{2\sqrt{x}}\right)\right)}{\left(\sqrt{x} - 1\right)^2} \qquad \left[\frac{d}{dx}\left(\sqrt{x}\right) = \frac{1}{2\sqrt{x}}\right]$$

(17)
$$g'(x) = x \left[\sec\left(\frac{1}{x}\right) \tan\left(\frac{1}{x}\right) \left(-\frac{1}{x^2}\right) \right] + \sec\left(\frac{1}{x}\right) (1)$$

(18)
$$y = (1 + \sin 2x)^{\frac{1}{2}} \Rightarrow \frac{dy}{dx} = \frac{1}{2}(1 + \sin 2x)^{-\frac{1}{2}}(\cos 2x)(2)$$

II. (1)
$$-3x\frac{dy}{dx} - 3y - 8y\frac{dy}{dx} = 0 \Rightarrow \frac{dy}{dx} = \frac{-3y}{3x + 8y}$$

(2) $16x = 6y^2\frac{dy}{dx} + 6xy\frac{dy}{dx} + 3y^2 \Rightarrow \frac{dy}{dx} = \frac{16x - 3y^2}{6y^2 + 6xy}$
(3) $\frac{3}{2x} + \frac{1}{y} = y \Rightarrow 3y + 2x = 2xy^2 \Rightarrow 3\frac{dy}{dx} + 2 = 4xy\frac{dy}{dx} + 2y^2 \Rightarrow \frac{dy}{dx} = \frac{2 - 2y^2}{4xy - 3}$

67 | P a g e

(4)
$$6x^{2} + 3x^{2}y = 2 - y \Rightarrow 12x + 3x^{2}\frac{dy}{dx} + 6xy = -\frac{dy}{dx} \Rightarrow \frac{dy}{dx} = \frac{-12x - 6xy}{3x^{2} + 1}$$

(5) $1 = (\sec^{2} y)\frac{dy}{dx} \Rightarrow \frac{dy}{dx} = \frac{1}{\sec^{2} y} = \cos^{2} y \text{ or } \frac{dy}{dx} = \frac{1}{1 + \tan^{2} y} = \frac{1}{1 + x^{2}}$
(6) $\frac{dy}{dx} = -\sin(x - y) \left[1 - \frac{dy}{dx} \right] \Rightarrow \frac{dy}{dx} = \frac{-\sin(x - y)}{1 - \sin(x - y)}$
(7) $x(\cos y)\frac{dy}{dx} + \sin y + y\cos x + (\sin x)\frac{dy}{dx} = 0 \Rightarrow \frac{dy}{dx} = \frac{-\sin y - y\cos x}{1 - \sin(x - y)}$

$$\frac{dx}{dx} = \frac{dx}{dx} = \frac{dx$$

(8)
$$1 = 3\sec^2(y^2 - 1)\left[\sec(y^2 - 1)\tan(y^2 - 1)\left(2y\frac{dy}{dx}\right) \Rightarrow \frac{dy}{dx} =$$

$$\frac{1}{6y \sec^3(y^2 - 1)\tan(y^2 - 1)} = \frac{1}{6xy \tan(y^2 - 1)}$$

III. (1)
$$2x + 6y \frac{dy}{dx} = 0 \Rightarrow 2(3) + 6(-2) \frac{dy}{dx} = 0 \Rightarrow 6 - 12 \frac{dy}{dx} = 0 \Rightarrow \frac{dy}{dx} = \frac{1}{2}$$

(2)
$$3x^{2} + \left(\frac{1}{3}y^{-2/3}\right)\frac{dy}{dx} = 0 \Rightarrow 3x^{2} + \left(\frac{1}{3\sqrt[3]{y^{2}}}\right)\frac{dy}{dx} = 0 \Rightarrow 3(1)^{2} + \left(\frac{1}{3\sqrt[3]{8^{2}}}\right)\frac{dy}{dx} = 0 \Rightarrow$$

$$3 + \frac{1}{12}\frac{dy}{dx} = 0 \Longrightarrow \frac{dy}{dx} = -36$$

(3)
$$xy = (y-2)^2 \Rightarrow x\frac{dy}{dx} + y = 2(y-2)\frac{dy}{dx} \Rightarrow 1\frac{dy}{dx} + 4 = 4\frac{dy}{dx} \Rightarrow \frac{dy}{dx} = \frac{4}{3}$$

COPYRIGHT FIMT 2020

68 | Page

(4)
$$3x\frac{dy}{dx} + 3y - 8x^3 = 3y^2\frac{dy}{dx} \Rightarrow 3(2)\frac{dy}{dx} + 3(-3) - 8(2^3) = 3(-3)^2\frac{dy}{dx} \Rightarrow$$

$$6\frac{dy}{dx} - 9 - 64 = 27\frac{dy}{dx} \Rightarrow \frac{dy}{dx} = \frac{-73}{21}$$
(5) $1 = (-\sin y)\frac{dy}{dx} \Rightarrow 1 = \left[-\sin\left(\frac{-\pi}{3}\right)\right]\frac{dy}{dx} \Rightarrow 1 = \frac{\sqrt{3}}{2}\frac{dy}{dx} \Rightarrow \frac{dy}{dx} = \frac{2}{\sqrt{3}}$
(6) $\cos(xy)\left[x\frac{dy}{dx} + y\right] = 1 \Rightarrow \frac{dy}{dx} = \frac{1 - y\cos(xy)}{x\cos(xy)} \Rightarrow \frac{dy}{dx} \cdot x = 1$
 $y = \frac{\pi}{2} = \frac{1 - \frac{\pi}{2}\cos\left(\frac{\pi}{2}\right)}{1\cos\left(\frac{\pi}{2}\right)} = 1$

 $\frac{1}{0} \Rightarrow \frac{dy}{dx}$ does not exist.

IV. (1)
$$f'(x) = 12x^3 - 8x + 7$$
 and $f''(x) = 36x^2 - 8x^3 - 8x$

(2)
$$f'(x) = \frac{(2x-1)(3) - (3x+1)(2)}{(2x-1)^2} = \frac{-5}{(2x-1)^2} = -5(2x-1)^{-2}$$
 and
 $f''(x) = 10(2x-1)^{-3}(2) = \frac{20}{(2x-1)^3}$
(3) $f'(x) = x^3 [-4\sin(4x)] + 3x^2 \cos(4x) = -4x^3 \sin(4x) + 3x^2 \cos(4x)$ and

$$f''(x) = -4x^{3} [4\cos(4x)] - 12x^{2} \sin(4x) + 3x^{2} [-4\sin(4x)] + 6x\cos(4x) =$$
$$-16x^{3} \cos(4x) - 24x^{2} \sin(4x) + 6x\cos(4x)$$

(4)
$$f'(x) = 4(\sin x)^3 \cos x$$
 and $f''(x) = 4(\sin x)^3(-\sin x) + (\cos x)(12(\sin x)^2(\cos x)) =$

$$-4\sin^4 x + 12\sin^2 x \cos^2 x$$

V.
$$s(2) = 2(2^3) + 4(2) - 5 = 16 + 8 - 5 = 19$$
 ft

$$v(t) = s'(t) = 6t^{2} + 4 \Longrightarrow v(2) = 6(2^{2}) + 4 = 28 \text{ ft/sec}$$

$$a(t) = v'(t) = s''(t) = 12t \implies a(2) = 12(2) = 24 \ ft/\sec^2$$

This is one of a series of worksheets designed to help you increase your confidence in handling Mathematics. This worksheet contains both theory and exercises which cover:-

- 1. Exponential functions 2. Logarithmic functions
- 3. Implicit Differentiation 4. Logarithmic Differentiation
- 5. Parametric Equations
- **1. Exponential Functions**

It can be shown that e^x is the function such that $\frac{d(e^x)}{dx}$

Examples

Differentiate the following (i) $y = e^{2x}$ (ii) $y = e^{f(x)}$ (iii) $y = \frac{1+x}{e^{2x}}$

(i) $y = e^{2x}$ is a function of a function

Writing
$$y = e^{u}$$
 where $u = 2x \Rightarrow \frac{dy}{du} = e^{u}$ and $\frac{du}{dx} = 2$

giving
$$\frac{dy}{dx} = \frac{dy}{du} \times \frac{du}{dx} = e^u \times 2 = 2e^{2x}$$

(ii)
$$y = e^{f(x)}$$
 is a function of a function

Writing
$$y = e^u$$
 where $u = f(x) \Rightarrow \frac{dy}{du} = e^u$ and $\frac{du}{dx} = f'(x)$

giving
$$\frac{dy}{dx} = \frac{dy}{du} \times \frac{du}{dx} = e^u \times f'(x) = f'(x)e^{f(x)}$$

(iii)
$$y = \frac{1+x}{e^{2x}}$$
 is a quotient

writing
$$u = 1 + x$$
, $v = e^{2x}$ gives $u' = 1$, $v' = 2e^{2x}$

using
$$\frac{dy}{dx} = \frac{u v - uv}{v^2}$$

gives $\frac{dy}{dx} = \frac{1 \times e^{2x} - (1 + x) \times 2e^{2x}}{(e^{2x})^2}$
 $= \frac{(1 - 2 - 2x)e^{2x}}{e^{4x}} = \frac{-(1 + 2x)}{e^{2x}}$

Exercise 1

Differentiate the following

1.
$$e^{7x}$$
 2. $e^{\cos x}$ 3. $2e^{2x} + 3e^{x^2} - e$ 4. $e^{(\sin x + \cos x)}$
5. $x^2 e^x$ 6. $\frac{e^x + 1}{x^3}$ 7. $\frac{(e^x - 1)}{e^x}$ 8. $\frac{e^{x^2}}{x}$

2. Logarithmic functions

hence

The inverse function of e^x is $\log_e x$ which is usually written as $\ln x$ (shorthand for natural or Napierian logarithms after Napier who developed them). For more information see the logs booklet.

Given $y = \ln x$ then, from the definition of logarithms,

$$x = e^{y}$$
 which gives $\frac{dx}{dy} = e^{y} \Rightarrow \frac{dy}{dx} = \frac{1}{e^{y}} = \frac{1}{x}$
$$\frac{d(\ln x)}{dx} = \frac{1}{x}$$

Extending this to differentiate $y = \ln[f(x)]$ which is a function (In) of the function f(x).

write
$$y = \ln u$$
 where $u = f(x) \Rightarrow \frac{dy}{du} = \frac{1}{u}$ and $\frac{du}{dx} = f'(x)$
using the chain rule $\frac{dy}{dx} = \frac{dy}{du} \times \frac{du}{dx}$
giving $\frac{dy}{dx} = \frac{1}{f(x)} \times f'(x) = \frac{f'(x)}{f(x)}$

Another important result to learn
$$\left| \frac{d}{dx} \left[\ln(f(x)) \right] = \frac{f'(x)}{f(x)} \right|$$

COPYRIGHT FIMT 2020

72 | Page
Examples

Differentiate the following functions

(i)
$$y = \ln(5x^2 - 6)$$
 (ii) $y = \ln(\frac{x+2}{x+3})$ (iii) $y = \ln(\frac{\sin^2 x}{\cos x})$

(i) Using the above

$$y = \ln(5x^2 - 6)$$
 gives $\frac{dy}{dx} = \frac{10x}{5x^2 - 6}$

(ii) Simplifying the expression gives

$$y = \ln\left(\frac{x+2}{x+3}\right) = \ln(x+2) - \ln(x+3)$$

Hence
$$\frac{dy}{dx} = \frac{1}{x+2} - \frac{1}{x+3}$$

= $\frac{(x+3) - (x+2)}{(x+2)(x+3)} = \frac{1}{(x+2)(x+3)}$

Note you could do this without simplifying but it is more difficult!

(iii) Simplifying the expression gives

$$y = \ln\left(\frac{\sin^2 x}{\cos x}\right) = \ln\left(\sin^2 x\right) - \ln(\cos x)$$

= $2\ln(\sin x) - \ln(\cos x)$
 $\frac{dy}{dx} = 2\frac{\cos x}{\sin x} - \frac{-\sin x}{\cos x}$
= $\frac{2\cos^2 x + \sin^2 x}{\sin x\cos x} = \frac{1 + \cos^2 x}{\sin x\cos x} (\text{using } \cos^2 x + \sin^2 x = 1)$

Exercise 2

Differentiate the following

1.
$$\ln 7x$$

2. $-\ln 6x$
3. $\ln(x^2 + 3)$
4. $\ln(x^2 + 2x - 1)$
5. $\ln(\frac{5x}{2x - 3})$
6. $\ln(\frac{x^2}{1 + x})$
7. $\ln(\frac{\sin x}{x \cos x})$
8. $\ln(\frac{(x - 1)^2}{\sqrt{x + 1}})$
9. $\ln[(\frac{x + 4}{\sqrt[3]{5x + 2}})^2]$

3. Implicit Functions

A function such as $y = x^5 + 3x^3 - x + 5$ is called an explicit function as y is explicitly given in terms of x.

A function such as $x^5 + 3x^3y - xy^2 + 5y - 3x = 15$ is called an implicit function as y is not given explicitly in terms of x nor x in terms of y.

An implicit function can be differentiated with respect to x as it stands.

Consider $x^2 + 3y - y^2 + xy - 3x = 15$

Differentiating each term with respect to x we get:

$$\frac{d(x^2)}{dx} + \frac{d(3y)}{dx} - \frac{d(y^2)}{dx} + \frac{d(xy)}{dx} - \frac{d(3x)}{dx} = \frac{d(15)}{dx}$$

To differentiate a function of y with respect to x we need

to use the chain rule
$$\frac{d[f(y)]}{dx} = \frac{d[f(y)]}{dy} \times \frac{dy}{dx}$$

giving
$$\frac{d(3y)}{dx} = \frac{3d(y)}{dy} \times \frac{dy}{dx} = 3\frac{dy}{dx}$$
 and $\frac{d(y^2)}{dx} = \frac{d(y^2)}{dy} \times \frac{dy}{dx} = 2y\frac{dy}{dx}$

using the product formula $\frac{d(xy)}{dx} = \frac{d(x)}{dx} \times \frac{y}{1} + \frac{x}{1} \times \frac{dy}{dx} = y + x\frac{dy}{dx}$

Putting these together we have:

$$\frac{d(x^2)}{dx} + \frac{d(3y)}{dx} - \frac{d(y^2)}{dx} + \frac{d(xy)}{dx} - \frac{d(3x)}{dx} = \frac{d(15)}{dx}$$

$$2x + 3\frac{dy}{dx} - 2y\frac{dy}{dx} + y + x\frac{dy}{dx} - 3 = 0$$

$$(3 - 2y + x)\frac{dy}{dx} = 3 - y - 2x$$

$$\frac{dy}{dx} = \frac{3 - y - 2x}{3 - 2y + x}$$

Example

Find the gradient of the curve $x^2 + y^2 + 2xy - 5x + 3y = 10$ at the points where x = 1

First we need to find the values of y when x = 1

Putting x = 1 we get $1 + y^2 + 2y - 5 + 3y = 10 \Rightarrow y^2 + 5y - 14 = 0$

which gives $(y-2)(y+7) = 0 \Rightarrow y = 2 \text{ or } y = -7$

notice that there are two points to consider (1, 2) and (1, -7)

Differentiating the function $x^2 + y^2 + 2xy - 5x + 3y = 10$

gives
$$\frac{d(x^2)}{dx} + \frac{d(y^2)}{dx} + \frac{d(2xy)}{dx} - \frac{d(5x)}{dx} + \frac{d(3y)}{dx} = \frac{d(10)}{dx}$$
$$2x + 2y\frac{dy}{dx} + 2y + 2x\frac{dy}{dx} - 5 + 3\frac{dy}{dx} = 0$$

giving
$$\frac{dy}{dx} = \frac{5 - 2x - 2y}{2y + 2x + 3}$$

at
$$P = (1, 2)$$
 $\frac{dy}{dx} = \frac{5 - 2 - 4}{4 + 2 + 3} = -\frac{1}{9}$
at $Q = (1, -7)$ $\frac{dy}{dx} = \frac{5 - 2 + 14}{-14 + 2 + 3} = -\frac{17}{9}$



Note you could substitute in and find

the value of $\frac{dy}{dx}$ without making it the

subject.

The sketch of the graph shows the two points *P* and *Q*. From the sketch you can see that the gradient is negative in both cases.

(iv) $2x^3 + 3xy^2 - y^3 = 0$

Exercise 3

- 1. In the following find $\frac{dy}{dx}$ in terms of x and y
 - (i) $x^2 + y^2 = 10$ (ii) $2x^2 + 2y^2 + 3x = 10 + 7y$

(iii)
$$x^2 - y^2 + 3xy = 6$$

- 2. Find the gradient of the curve $x^2 + 6y^2 = 10$ at the points where x = 2.
- 3. Find the gradient of the curve $x^3 + 4xy = y^2 + 15$ at the points where x = 2.

4. Logarithmic Differentiation

The function $y = a^x$ cannot be differentiated by any of the methods developed so far. But taking the natural logarithm of both sides overcomes the problem!

To solve $y = a^x$ take logs $\ln y = \ln(a^x) = x \ln a$

differentiate
$$\frac{d(\ln y)}{dx} = \frac{d(x \ln a)}{dx}$$

By the chain rule the left hand side gives $\frac{d(\ln y)}{dx} = \frac{d(\ln y)}{dy}\frac{dy}{dx} = \frac{1}{y}\frac{dy}{dx}$

the right hand side gives

putting these together gives

$$\frac{1}{y}\frac{dy}{dx} = \ln a$$

 $\frac{d(x\ln a)}{dx} = \ln a \frac{d(x)}{dx} = \ln a$

hence

$$\frac{dy}{dx} = (\ln a)y = (\ln a)a^x$$

This method can simplify differentiation in a number of cases, as shown in the following examples.

Examples (The first two could be differentiated as quotients.)

1. Find $\frac{dy}{dx}$ given the function $y = \frac{\sin x}{\cos x}$ (ie tanx)

Taking logs gives

 $\ln y = \ln \sin x - \ln \cos x$

Differentiate
$$\frac{1}{y}\frac{dy}{dx} = \frac{\cos x}{\sin x} - \frac{-\sin x}{\cos x} = \frac{\cos^2 x + \sin^2 x}{\sin x \cos x} = \frac{1}{\sin x \cos x}$$
$$\frac{dy}{dx} = \frac{1}{\sin x \cos x} \times y = \frac{1}{\sin x \cos x} \times \frac{\sin x}{\cos x} = \frac{1}{\cos^2 x} = \sec^2 x$$

The result should be known $\frac{d}{dx}(\tan x) = \sec^2 x$

2. Find
$$\frac{dy}{dx}$$
 given the function $y = \frac{x \sin x}{(x+1)\cos x}$

$$\ln(y) = \ln\left(\frac{x\sin x}{(x+1)\cos x}\right) = \ln(x\sin x) - \ln[(x+1)\cos x]$$
$$= \ln(x) + \ln(\sin x) - \ln(x+1) - \ln(\cos x)$$

Differentiating gives

$$\frac{1}{y}\frac{dy}{dx} = \frac{1}{x} + \frac{\cos x}{\sin x} - \frac{1}{x+1} + \frac{\sin x}{\cos x}$$
$$\frac{dy}{dx} = y \left[\frac{1}{x} + \frac{\cos x}{\sin x} - \frac{1}{x+1} + \frac{\sin x}{\cos x} \right]$$
$$= \frac{x \sin x}{(x+1)\cos x} \left[\frac{1}{x} + \frac{\cos x}{\sin x} - \frac{1}{x+1} + \frac{\sin x}{\cos x} \right]$$

which is a lot easier than using the quotient method. It could be 'simplified' but this rarely needs to be done.

3. Find
$$\frac{dy}{dx}$$
 given the function $y = x^x$

Take natural logs $\ln(y) = \ln(x^x) = x \ln(x)$

Differentiate $\frac{1}{y}\frac{dy}{dx} = x \times \frac{1}{x} + \ln(x) = 1 + \ln(x)$ (using the product rule)

$$\frac{dy}{dx} = y[1 + \ln(x)] = x^x[1 + \ln(x)]$$

Exercise 4

Use logarithmic differentiation to differentiate the following:

1.
$$r = 2^{\theta}$$

2. $y = x^{x}$
3. $s = \sin^{t} t = (\sin t)^{t}$
4. $v = \sin(u^{u})$
5. $y = \frac{xe^{x} + 1}{e^{x}(x+1)}$
6. $y = \frac{\sin^{2} x}{1 + \cos x}$
7. $y = \frac{(x+1)}{(2x+3)^{2}(x-4)}$

5. Parametric Differentiation

When a function is given in parametric form it means that x and y are given in terms of another variable, the parameter. i.e. x = f(t), y = g(t).

 $x = t^2$, y = 2t are parametric equations. Frequently the parameter can be eliminated.

 $y = 2t \Rightarrow t = \frac{1}{2}y$ but $x = t^2$ hence $x = (\frac{1}{2}y)^2 = \frac{1}{4}y^2$ or $y^2 = 4x$, the equation of a parabola

To find the gradient of such a function in parametric form we need to use the chain rule

		2.0					dy			
dy_	dy	dt	which	canhe	writton	$as \frac{dy}{dy}$	dt	or	<i>y</i> '	
dx	dt	dx	which	canbe	wintten	$\frac{ds}{dx} =$	dx	01	x'	
			-				dt			

Given $x = t^2$, y = 2twe have $\frac{dx}{dt} = 2t$, $\frac{dy}{dt} = 2$ hence $\frac{dy}{dx} = \frac{2}{2t} = \frac{1}{t}$

In this case we can also find the gradient using the Cartesian equations:

Given $y^2 = 4x$ we have $2y \frac{dy}{dx} = 4$ hence $\frac{dy}{dx} = \frac{4}{2y} = \frac{2}{y}$

Comparing the two answers, as y = 2t then $\frac{2}{y} = \frac{1}{t}$ so the two answers are the same (as expected!)

Examples

1. Find the gradient of the curve given by $x = \sin t$, $y = \cos 2t$ when $t = \frac{\pi}{3}$.

Find
$$\frac{dy}{dt}$$
 and $\frac{dx}{dt}$
and use $\frac{dy}{dx} = \frac{dy}{dt} \div \frac{dx}{dt} = \frac{y'}{x'} = \frac{-2\sin 2t}{\cos t}$
when $t = \frac{\pi}{3}, \frac{dy}{dx} = \frac{-2\sin\left(\frac{2\pi}{3}\right)}{\cos\left(\frac{\pi}{3}\right)} = \frac{-2\left(\frac{\sqrt{3}}{2}\right)}{\frac{1}{2}} = -2\sqrt{3}$

Finally substitute for t

Notes a) It would be possible to eliminate t and obtain the Cartesian equation $y = 1 - 2x^2$ which will give the same value for the gradient.

b) By putting $\sin 2x = 2\sin x \cos x$, $\frac{dy}{dx} = \frac{-2\sin 2t}{\cos t}$ can be simplified to

 $\frac{dy}{dx} = \frac{-4\sin t\cos t}{\cos t} = -4\sin t$ if necessary.

2. Find the gradient of the curve given by $x = \theta + \sin \theta$, $y = 1 - \cos \theta$ when $\theta = \frac{\pi}{2}$ and when $\theta = \pi$.

alde

$$x = \theta + \sin\theta \Rightarrow \frac{dx}{d\theta} = 1 + \cos\theta; \quad y = 1 - \cos\theta \Rightarrow \frac{dy}{d\theta} = \sin\theta$$

$$\frac{dy}{dx} = \frac{y}{x'} = \frac{\sin\theta}{1 + \cos\theta}$$

when
$$\theta = \frac{\pi}{2}$$
 $\frac{dy}{dx} = \frac{\sin(\frac{\pi}{2})}{1 + \cos(\frac{\pi}{2})} = \frac{1}{1+0} = 1$, tangent at 45°

when,
$$\theta = \pi \quad \frac{dy}{dx} = \frac{\sin \pi}{1 + \cos \pi} = \frac{0}{1 - 1} = \infty$$
, tangent vertical

(really the value is indeterminate)

Notes a) $x = \theta + \sin \theta$ and $y = 1 - \cos \theta$ cannot be made into a simple Cartesian equation!

b)
$$\frac{dy}{dx} = \frac{\sin\theta}{1 + \cos\theta}$$
 can be simplified by putting $\sin\theta = 2\sin\left(\frac{\theta}{2}\right)\cos\left(\frac{\theta}{2}\right)$

and
$$\cos\theta = 2\cos^2\left(\frac{\theta}{2}\right) - 1$$
 giving $\frac{dy}{dx} = \frac{2\sin\left(\frac{\theta}{2}\right)\cos\left(\frac{\theta}{2}\right)}{1 + 2\cos^2\left(\frac{\theta}{2}\right) - 1} = \frac{2\sin\left(\frac{\theta}{2}\right)\cos\left(\frac{\theta}{2}\right)}{2\cos^2\left(\frac{\theta}{2}\right)} = \tan\left(\frac{\theta}{2}\right)$

Exercise 5 In questions 1 to 8 find $\frac{dy}{dx}$ in terms of the given parameter.

1.	$x = 1 + 3t, y = 3 - t^2$	2. $x = t^2 + 3t$, $y = 2t - t^2$
3.	$x = 1 - 4s, y = 1 + s^2$	4. $x = \sin \phi$, $y = \cos 2\phi + 1$
5.	$x = e^{\theta}, y = 1 + \theta$	6. $x = e^{u} + e^{-u}, y = e^{u} - e^{-u}$
7.	$x = \frac{1}{s}, y = \frac{s+1}{s-1}$	8. $x = \frac{1}{1+t}, y = \frac{t^2}{1+t}$
	150 9001	:2015 & 14001:2015

In questions 9 to 14 find the gradient of the curve at the given point.

9. x = 1 - 3t, $y = 3 + 2t^2$; t = 2 10. $x = t^2 - t$, $y = 2t - t^2$; t = -211. $x = \cos \alpha$, $y = \sin 2\alpha$; $\alpha = \frac{\pi}{3}$ 12. $x = \sin \phi - 1$, $y = \cos \phi + 1$; $\phi = 0$ 13. $x = \ln s$, $y = s \ln(s+1)$; s = 1 14. $x = e^{2r} + e$, $y = 2e^r - r$; r = 0

ANSWERS

Exercise 1

1.
$$7e^{7x}$$
 2. $-\sin x e^{\cos x}$ 3. $4e^{2x} + 6x e^{x^2}$ 4. $(\cos x - \sin x)e^{(\sin x + \cos x)}$
5. $(x^2 + 2x)e^x$ 6. $\frac{xe^x - 3e^x - 3}{x^4}$ 7. $e^{-x} = \frac{1}{e^x}$ 8. $\frac{e^x(2x^2 - 1)}{x^2}$

MAGEMENT

Exercise 2

1.
$$\frac{1}{x}$$
 2. $-\frac{1}{x}$ 3. $\frac{2x}{x^2+3}$ 4. $\frac{2x+2}{x^2+2x-1}$ 5. $\frac{1}{x} - \frac{2}{2x-3} = \frac{-3}{x(2x-3)}$
6. $\frac{2}{x} - \frac{1}{x+1} = \frac{x+2}{x(x+1)}$ 7. $\cot x + \tan x - \frac{1}{x} = \frac{1}{\cos x \sin x} - \frac{1}{x}$

8. Simplify to
$$2\ln(x-1) - \frac{1}{2}\ln(x+1)$$
; answer $\frac{2}{x-1} - \frac{1}{2(x+1)} = \frac{3x+5}{2(x-1)(x+1)}$

9. Hint first simplify as above; answer $\frac{2}{x+4} - \frac{10}{3(5x+2)} = \frac{4(5x-7)}{3(x+4)(5x+2)}$

Exercise 3

S. 10

1. (i)
$$-\frac{x}{y}$$
 (ii) $\frac{3+4x}{7-4y}$ (iii) $\frac{3y+2x}{2y-3x}$ (iv) $\frac{6x^2+3y^2}{3y^2-6xy} = \frac{2x^2+y^2}{y(y-2x)}$

2.
$$x = 2 \Rightarrow y = \pm 1$$
 $\frac{dy}{dx} = \frac{-x}{6y}$ $\text{grad} = \pm \frac{1}{3}$

3.
$$x = 2 \Rightarrow y = 1$$
, 7 $\frac{dy}{dx} = \frac{3x^2 + 4y}{2y - 4x}$ grad $= -\frac{8}{3}, \frac{20}{3}$

Exercise 4

1.
$$(\ln 2)2^{\theta}$$
 2. $(1 + \ln x)x^{x}$ 3. $s = \left(\ln(\sin t) + \frac{t\cos t}{\sin t}\right)(\sin t)^{t}$ 4. $\cos(u^{u})(1 + \ln u)u^{u}$
5. $\frac{e^{x} - x - 2}{e^{x}(x+1)^{2}}$ 6. $\sin x$ 7. $\frac{-4x^{2} + 2x + 1}{(x-4)^{2}(2x+3)^{2}}$

NAAC ACCREDITED

SHAGEME

Exercise 5

Taylor & Maclaurin Series

Notation – If f(x) is continuously differentiable *n* times then we say that $f(x) = f^{(0)}(x)$, $f'(x) = f^{(1)}(x)$, $f''(x) = f^{(2)}(x)$, and that the *i*th derivative of *f* is $f^{(i)}(x)$.

Suppose that a function f has a power series representation $f(x) = c_0 + c_1(x - a) + c_2(x - a)^2 + c_3(x - a)^3 + c_4(x - a)^4 + \dots$ for |x - a| < R. Differentiating term by term find the following:

- f⁽⁰⁾(x)
- $f^{(1)}(x)$
- f⁽²⁾(x)
- $f^{(3)}(x)$
- f⁽ⁱ⁾(x)

a

Find a formula for c_n given the above information (Recall the convention 0! = 1)

Theorem (Taylor Series Representation) – If f has a power series expansion about a (i.e. if

$$f(x) = \sum_{n=0}^{\infty} c_n (x-a)^n$$
) then the power series must come in the form

$$f(x) = \sum_{n=0}^{\infty} \frac{f^{(n)}(a)}{n!} (x-a)^n$$

COPYRIGHT FIMT 2020

2015

Definition – The power series $\sum_{n=0}^{\infty} \frac{f^{(n)}(a)}{n!} (x-a)^n$ is called the **Taylor Series associated with**

the function *f* about *a*. In the special case when a = 0 we call $\sum_{n=0}^{\infty} \frac{f^{(n)}(0)}{n!} x^n$ the Maclaurin series of *f*.

Show that the Maclaurin series associated with $f(x) = e^x$ is given by $\sum_{n=0}^{\infty} \frac{x^n}{n!}$. Show that this power series converges for *all* real values for *x*.

Notice: We do NOT know that $e^x = \sum_{n=0}^{\infty} \frac{x^n}{n!}$ yet.

Question: When is it true that a function f(x) with derivatives of all orders is equal to its Taylor series? Equivalently, when is it true that a function f(x) has a power series?

Notation – Let $T_n(x) = \sum_{i=0}^n \frac{f^{(i)}(a)}{i!} (x-a)^i$ be called the **n**th **degree Taylor Polynomial of** f**about** a. Let $R_n(x) = f(x) - T_n(x)$ be called the **remainder** of the Taylor Series of f about a. Find and graph the 1st 2nd and 3rd degree Taylor Polynomials of $f(x) = e^x$ about 0 along with $f(x) = e^x$.

Notice: If we can show that $\lim_{n\to\infty} R_n(x) = 0$ then

$$0 = \lim_{n \to \infty} R_n(x) = \lim_{n \to \infty} [f(x) - T_n(x)] = f(x) - \lim_{n \to \infty} T_n(x)$$

Hence f(x) must be equal to its Taylor Series $\sum_{n=0}^{\infty} \frac{f^{(n)}(a)}{n!} (x-a)^n$

Theorem (Taylor's Inequality)(A bound on the remainder of a Taylor Series) -

Suppose that f(x) has derivatives of all orders. If $|f^{(n+1)}(x)| \le M$ for all $|x - a| \le d$, then the remainder $R_n(x)$ of the Taylor Series satisifes the following inequality

$$|R_n(x)| \le \frac{M}{(n+1)!} |x-a|^{n+1}$$
 for all $|x-a| \le d$

Proof (case n = 1):

Assume that $|f''(x)| \le M$. Then $f''(x) \le M$ for $a - d \le x \le a + d$ so

 $\int_{a}^{x} f''(t)dt \le \int_{a}^{x} Mdt$

$$f'(x) - f'(a) \le M(x - a) \text{ or } f'(x) \le f'(a) + M(x - a)$$

Integrating both sides again yields

$$\int_{a}^{x} f'(t)dt \leq \int_{a}^{x} \left[f'(a) + M(t-a)\right]dt$$

$$f(x) - f(a) \le f'(a)(x-a) + M \frac{(x-a)^2}{2}$$

But then we have

$$R_1(x) = f(x) - T_1(x) = f(x) - f(a) - f'(a)(x - a) \le M \frac{(x - a)^2}{2}$$

as desired.

Case (n = 2) (write out yourself)

Prove that $e^x = \sum_{n=0}^{\infty} \frac{x^n}{n!}$. Then approximate *e* correct to 5 decimal places.

Prove that $\sin x = \sum_{n=0}^{\infty} (-1)^n \frac{x^{2n+1}}{(2n+1)!}$ and approximate $\sin 2$ correct to five decimal places

Prove that $\cos x = \sum_{n=0}^{\infty} (-1)^n \frac{x^{2n}}{(2n)!}$ and approximate $\cos 2$ correct to five decimal places

Differentiate the Taylor series for sin x term by term. Does the result agree with the Taylor series for $\cos x$?

Write out the Taylor series for e^{-x^2}

Find the Taylor Series for $\int e^{-x^2} dx$

Evaluate $\int_{0}^{1} e^{-x^{2}} dx$ correct to three decimal places.

Show that $\lim_{x\to 0} \frac{e^x - 1 - x}{x^2} = 0.5$ using the Maclaurin series for e^x . Check your result using l'Hôpital.

i nopital.

Evaluate $\int \frac{e^x - 1}{x} dx$ as a power series term by term.

What should your calculator output if you type in taylor(e^x , x, 4, 0). Explain how you knew what to predict.

Limits – Indeterminate Forms and L'Hospital's Rule

I. Indeterminate Form of the Type $\frac{0}{0}$

We have previously studied limits with the indeterminate form $\frac{0}{0}$ as shown in the

following examples:

Example 1:
$$\lim_{x \to 2} \frac{x^2 - 4}{x - 2} = \lim_{x \to 2} \frac{(x + 2)(x - 2)}{x - 2} = \lim_{x \to 2} (x + 2) = 2 + 2 = 4$$

Example 2: $\lim_{x \to 0} \frac{\tan 3x}{\sin 2x} = \lim_{x \to 0} \frac{\frac{\sin 3x}{\cos 3x}}{\sin 2x} = \lim_{x \to 0} \frac{\sin 3x}{1} \cdot \frac{1}{\cos 3x} \cdot \frac{1}{\sin 2x} =$

$$\frac{3}{2} \left(\lim_{3x \to 0} \frac{\sin 3x}{3x} \right) \left(\lim_{x \to 0} \frac{1}{\cos 3x} \right) \left(\lim_{2x \to 0} \frac{2x}{\sin 2x} \right) = \frac{3}{2} (1)(1)(1) = \frac{3}{2}$$

[*Note*: We use the given limit
$$\lim_{\Delta \to 0} \frac{\sin \Delta}{\Delta} = 1$$
.]

Example 3:
$$\lim_{h \to 0} \frac{\sqrt[3]{8+h}-2}{h} = f'(8) = \frac{1}{3\sqrt[3]{8^2}} = \frac{1}{12}$$
. [*Note*: We use the definition

of the derivative $f'(a) = \lim_{h \to 0} \frac{f(a+h) - f(a)}{h}$ where $f(x) = \sqrt[3]{x}$

Example 4: $\lim_{x \to \pi/3} \frac{\cos x - \frac{1}{2}}{x - \pi/3} = f'(\pi/3) = -\sin(\pi/3) = -\sqrt{3}/2$. [*Note*: We use the

definition of the derivative $f'(a) = \lim_{x \to a} \frac{f(x) - f(a)}{x - a}$ where

$$f(x) = \cos x$$
 and $a = \frac{\pi}{3}$.]

However, there is a general, systematic method for determining limits with the

indeterminate form $\frac{0}{0}$. Suppose that *f* and *g* are differentiable functions at *x* = *a*

and that $\lim_{x \to a} \frac{f(x)}{g(x)}$ is an indeterminate form of the type $\frac{0}{0}$; that is, $\lim_{x \to a} f(x) = 0$

and $\lim_{x \to a} g(x) = 0$. Since f and g are differentiable functions at x = a, then f and g

are continuous at
$$x = a$$
; that is, $f(a) = \lim_{x \to a} f(x) = 0$ and $g(a) = \lim_{x \to a} g(x) = 0$.

THE HER

Furthermore, since f and g are differentiable functions at x = a, then f'(a) =

$$\lim_{x \to a} \frac{f(x) - f(a)}{x - a} \text{ and } g'(a) = \lim_{x \to a} \frac{g(x) - g(a)}{x - a}. \text{ Thus, if } g'(a) \neq 0 \text{, then}$$

$$\lim_{x \to a} \frac{f(x)}{g(x)} = \lim_{x \to a} \frac{f(x) - f(a)}{g(x) - g(a)} = \lim_{x \to a} \frac{\frac{f(x) - f(a)}{x - a}}{\frac{g(x) - g(a)}{x - a}} = \frac{f'(a)}{g'(a)} = \lim_{x \to a} \frac{f'(x)}{g'(x)} \text{ if } f' \text{ and}$$

g' are continuous at x = a. This illustrates a special case of the technique known as

L'Hospital's Rule.

L'Hospital's Rule for Form $\frac{0}{0}$ Suppose that f and g are differentiable functions on an open interval containing x = a, except possibly at x = a, and that $\lim_{x \to a} f(x) = 0$ and $\lim_{x \to a} g(x) = 0$. If $\lim_{x \to a} \frac{f'(x)}{g'(x)}$ has a finite limit, or if this limit is $+\infty$ or $-\infty$, then $\lim_{x \to a} \frac{f(x)}{g(x)} = \lim_{x \to a} \frac{f'(x)}{g'(x)}$. Moreover, this statement is also true in the case of a limit as $x \to a^-, x \to a^+, x \to -\infty$, or as $x \to +\infty$.

In the following examples, we will use the following three-step process:

Step 1. Check that the limit of $\frac{f(x)}{g(x)}$ is an indeterminate form of type $\frac{0}{0}$. If it

is not, then L'Hospital's Rule cannot be used.

150 9001:2015 & 14001:2015

Step 2. Differentiate *f* and *g* separately. [*Note*: **Do not differentiate** $\frac{f(x)}{g(x)}$

using the quotient rule!]

Step 3. Find the limit of $\frac{f'(x)}{g'(x)}$. If this limit is finite, $+\infty$, or $-\infty$, then it is

equal to the limit of $\frac{f(x)}{g(x)}$. If the limit is an indeterminate form of type

$$\frac{0}{0}$$
, then simplify $\frac{f'(x)}{g'(x)}$ algebraically and apply **L'Hospital's Rule** again.

Example 1:
$$\lim_{x \to 2} \frac{x^2 - 4}{x - 2} = \lim_{x \to 2} \frac{2x}{1} = 2(2) = 4$$

Example 2: $\lim_{x \to 2} \frac{\tan 3x}{x} = \lim_{x \to 2} \frac{3\sec^2 3x}{1} = \frac{3(1)}{2} = \frac{3}{2}$

Example 2: $\lim_{x \to 0} \frac{\tan 2x}{\sin 2x} = \lim_{x \to 0} \frac{2}{2\cos 2x} = \frac{2}{2(1)} = \frac{1}{2}$

Example 3:
$$\lim_{h \to 0} \frac{\sqrt[3]{8+h}-2}{h} = \lim_{h \to 0} \frac{\frac{1}{3}(8+h)^{-\frac{2}{3}}(1)}{1} = \lim_{h \to 0} \frac{1}{3(8+h)^{\frac{2}{3}}} = \frac{1}{3(8)^{\frac{2}{3}}} = \frac{1}{12}$$

Example 4:
$$\lim_{x \to \pi/3} \frac{\cos x - 1/2}{x - \pi/3} = \lim_{x \to \pi/3} \frac{-\sin x}{1} = -\sin(\pi/3) = -\frac{\sqrt{3}}{2}$$

Example 5:
$$\lim_{x \to 0} \frac{e^x - x - 1}{x^2} = \lim_{x \to 0} \frac{e^x - 1}{2x} = \lim_{x \to 0} \frac{e^x}{2} = \frac{1}{2}$$
 [Use L'Hospital's Rule

twice.]

Example 6:
$$\lim_{x \to +\infty} \frac{\frac{1}{x^2}}{\sin(\frac{1}{x})} = \lim_{x \to +\infty} \frac{\frac{-2}{x^3}}{\cos(\frac{1}{x}) - \frac{1}{x^2}} = \lim_{x \to +\infty} \frac{\frac{2}{x}}{\cos(\frac{1}{x})} = \frac{0}{1} = 0$$
, or

$$\lim_{x \to +\infty} \frac{\frac{1}{x^2}}{\sin(\frac{1}{x})} = \lim_{y \to 0^+} \frac{y^2}{\sin y} = \lim_{y \to 0^+} \frac{2y}{\cos y} = \frac{2(0)}{1} = 0 \text{ where } y = \frac{1}{x}.$$

Example 7: $\lim_{x \to 0} \frac{x}{\ln x} = \lim_{x \to 0} x \left(\frac{1}{\ln x} \right) = 0(0) = 0$ [This limit is **not** an indeterminate

form of the type
$$\displaystyle rac{0}{0}$$
 , so **L'Hospital's Rule** cannot be used.]

II. Indeterminate Form of the Type $\frac{\infty}{\infty}$

We have previously studied limits with the indeterminate form $\frac{\infty}{\infty}$ as shown in the

following examples:

Example 1:
$$\lim_{x \to +\infty} \frac{3x^2 + 5x - 7}{2x^2 - 3x + 1} = \lim_{x \to +\infty} \frac{\frac{3x^2}{x^2} + \frac{5x}{x^2} - \frac{7}{x^2}}{\frac{2x^2}{x^2} - \frac{3x}{x^2} + \frac{1}{x^2}} = \\\lim_{x \to +\infty} \frac{3 + \frac{5}{x} - \frac{7}{x^2}}{2 - \frac{3}{x} + \frac{1}{x^2}} = \lim_{x \to +\infty} \frac{3 + 0 - 0}{2 - 0 + 0} = \frac{3}{2}$$
Example 2:
$$\lim_{x \to -\infty} \frac{3x - 1}{x^2 + 1} = \lim_{x \to -\infty} \frac{\frac{3x}{x^2} - \frac{1}{x^2}}{\frac{x^2}{x^2} + \frac{1}{x^2}} = \lim_{x \to -\infty} \frac{\frac{3}{x} - \frac{1}{x^2}}{1 + \frac{1}{x^2}} = \frac{0 - 0}{1 + 0} = \frac{0}{1} = 0$$
Example 3:
$$\lim_{x \to \infty} \frac{3x^3 - 4}{2x^2 + 1} = \lim_{x \to -\infty} \frac{\frac{3x^3}{x^2} - \frac{4}{x^3}}{\frac{2x^2}{x^2} + \frac{1}{x^2}} = \lim_{x \to -\infty} \frac{3 - \frac{4}{x^3}}{\frac{2}{x} + \frac{1}{x^3}} = \frac{3 - 0}{0 + 0} = \frac{3}{0} \Rightarrow$$
I init does not exist.
Example 4:
$$\lim_{x \to -\infty} \frac{\sqrt{4x^2 + 1}}{x + 1} = \lim_{x \to -\infty} \frac{\sqrt{4x^2 + 1}}{\frac{x + 1}{x}} = \lim_{x \to -\infty} \frac{\sqrt{4x^2 + 1}}{\frac{x + 1}{x}} (\sqrt{x^2} = -x)$$

because
$$x < 0$$
 and thus $x = -\sqrt{x^2}$) = $\lim_{x \to -\infty} \frac{-\sqrt{\frac{4x^2 + 1}{x^2}}}{\frac{x + 1}{x}} =$

$$\lim_{x \to -\infty} \frac{-\sqrt{4 + \frac{1}{x^2}}}{1 + \frac{1}{x^2}} = \frac{-\sqrt{4}}{1} = -2.$$

However, we could use another version of L'Hospital's Rule.

L'Hospital's Rule for Form
$$\frac{\infty}{\infty}$$

Suppose that f and g are differentiable functions on an open interval
containing $x = a$, except possibly at $x = a$, and that $\lim_{x \to a} f(x) = \infty$ and
 $\lim_{x \to a} g(x) = \infty$. If $\lim_{x \to a} \frac{f'(x)}{g'(x)}$ has a finite limit, or if this limit is $+\infty$ or
 $-\infty$, then $\lim_{x \to a} \frac{f(x)}{g(x)} = \lim_{x \to a} \frac{f'(x)}{g'(x)}$. Moreover, this statement is also true
in the case of a limit as $x \to a^-, x \to a^+, x \to -\infty$, or as $x \to +\infty$.
Example 1: $\lim_{x \to \infty} \frac{3x^2 + 5x - 7}{2x^2 - 3x + 1} = \lim_{x \to \infty} \frac{6x + 5}{4x - 3} = \lim_{x \to \infty} \frac{6}{4} = \frac{3}{2}$
Example 2: $\lim_{x \to \infty} \frac{3x^{-1}}{x^2 + 1} = \lim_{x \to \infty} \frac{3}{2x} = \frac{3}{2} \lim_{x \to \infty} \frac{1}{x} = \frac{3}{2}(0) = 0$
Example 3: $\lim_{x \to \infty} \frac{3x^3 - 4}{2x^2 + 1} = \lim_{x \to \infty} \frac{9x^2}{4x} = \lim_{x \to \infty} \frac{18x}{4} = \infty$

COPYRIGHT FIMT 2020

91 | Page

Example 4:
$$\lim_{x \to \infty} \frac{\sqrt{4x^2 + 1}}{x + 1} = \lim_{x \to \infty} \frac{2\sqrt{4x^2 + 1}}{1} = \lim_{x \to \infty} \frac{4x}{\sqrt{4x^2 + 1}} \Rightarrow L'Hospital's$$

8*x*

Rule does not help in this situation. We would find the limit as we

did previously.

Example 5:
$$\lim_{x \to +\infty} \frac{\ln(x^2 + 1)}{\ln(x^3 + 1)} = \lim_{x \to +\infty} \frac{\frac{2x}{x^2 + 1}}{\frac{3x^2}{x^3 + 1}} = \lim_{x \to +\infty} \frac{2x(x^3 + 1)}{3x^2(x^2 + 1)} = \lim_{x \to +\infty} \frac{2x^4 + 2x}{3x^4 + 3x^2} =$$

NAAF AFFDEDITED

$$\lim_{x \to +\infty} \frac{8x^3 + 2}{12x^3 + 6x} = \lim_{x \to +\infty} \frac{24x^2}{36x^2 + 6} = \lim_{x \to +\infty} \frac{48x}{72x} = \frac{48}{72} = \frac{2}{3}$$

Example 6:
$$\lim_{x \to 0^+} \frac{\ln x}{\frac{1}{x^2}} = \lim_{x \to 0^+} \frac{\frac{1}{x}}{-\frac{2}{x^3}} = \lim_{x \to 0^+} \frac{x^3}{-2x} = \lim_{x \to 0^+} \frac{x^2}{-2} = \frac{0^2}{-2} = 0$$

Example 7:
$$\lim_{x \to +\infty} \frac{\arctan x}{x} = \left(\lim_{x \to +\infty} \frac{1}{x}\right) \left(\lim_{x \to +\infty} \arctan x\right) = (0) \left(\frac{\pi}{2}\right) = 0$$
 [This limit is

not an indeterminate form of the type $\frac{\infty}{\infty}$, so **L'Hospital's Rule**

0

cannot be used.]

III. Indeterminate Form of the Type $0 \cdot \infty$

20

Indeterminate forms of the type $0\cdot\infty$ can sometimes be evaluated by rewriting the 150 9001:2015 & 14001-20

forms of type
$$\frac{0}{0}$$
 or $\frac{\infty}{\infty}$.

Example 1:
$$\lim_{x \to 0^+} x \ln x = \lim_{x \to 0^+} \frac{\ln x}{\frac{1}{x}} = \lim_{x \to 0^+} \frac{\frac{1}{x}}{-\frac{1}{x^2}} = \lim_{x \to 0^+} \frac{-x^2}{x} = \lim_{x \to 0^+} (-x) = 0$$

Example 2:
$$\lim_{x \to 0^+} (\sin x) \ln x = \lim_{x \to 0^+} \frac{\ln x}{\csc x} = \lim_{x \to 0^+} \frac{\frac{1}{x}}{-\csc x \cot x} = \lim_{x \to 0^+} \frac{-\sin x \tan x}{x} = \lim_{x \to 0^+} \frac{-\sin x \tan x}{x}$$

$$\left(\lim_{x \to 0^+} \frac{-\sin x}{x}\right) \left(\lim_{x \to 0^+} \tan x\right) = (-1)(0) = 0$$

Example 3:
$$\lim_{x \to +\infty} x \sin\left(\frac{1}{x}\right) = \lim_{x \to +\infty} \frac{\sin\left(\frac{1}{x}\right)}{\frac{1}{x}} = \lim_{y \to 0^+} \frac{\sin y}{y} = 1$$
 [Let $y = \frac{1}{x}$.]

IV. Indeterminate Form of the Type $\infty - \infty$

A limit problem that leads to one of the expressions

$$(+\infty) - (+\infty)$$
, $(-\infty) - (-\infty)$, $(+\infty) + (-\infty)$, $(-\infty) + (+\infty)$

is called an **indeterminate form of type** $\infty - \infty$. Such limits are indeterminate because the two terms exert conflicting influences on the expression; one pushes it in the positive direction and the other pushes it in the negative direction. However, limits problems that lead to one the expressions

$$(+\infty) + (+\infty)$$
, $(+\infty) - (-\infty)$, $(-\infty) + (-\infty)$, $(-\infty) - (+\infty)$

are not indeterminate, since the two terms work together (the first two produce a limit of $+\infty$ and the last two produce a limit of $-\infty$). Indeterminate forms of the type $\infty - \infty$ can sometimes be evaluated by combining the terms and manipulating

the result to produce an indeterminate form of type $\frac{0}{0}$ or $\frac{\infty}{\infty}$.

Example 1:
$$\lim_{x \to 0^+} \left(\frac{1}{x} - \frac{1}{\sin x} \right) = \lim_{x \to 0^+} \left(\frac{\sin x - x}{x \sin x} \right) = \lim_{x \to 0^+} \frac{\cos x - 1}{x \cos x + \sin x} =$$

$$\lim_{x \to 0^+} \frac{-\sin x}{-x\sin x + \cos x + \cos x} = \frac{0}{2} = 0$$

Example 2:
$$\lim_{x \to 0} \left[\ln(1 - \cos x) - \ln(x^2) \right] = \lim_{x \to 0} \left[\ln\left(\frac{1 - \cos x}{x^2}\right) \right] = \ln\left[\lim_{x \to 0} \left(\frac{1 - \cos x}{x^2}\right) \right] = \ln\left[\lim_{x \to 0} \left(\frac{\sin x}{2x}\right) \right] = \ln\left(\frac{1}{2}\right)$$

V. Indeterminate Forms of the Types $\,0^0,\,\infty^0,\,1^\infty$

Limits of the form $\lim_{x\to a} [f(x)]^{g(x)} \left\{ or \lim_{x\to\infty} [f(x)]^{g(x)} \right\}$ frequently give rise to

indeterminate forms of the types $\,0^{0},\,\infty^{0},\,1^{\infty}\,.\,$ These indeterminate forms can

sometimes be evaluated as follows:

 $r \rightarrow 0$

(1)
$$y = [f(x)]^{g(x)}$$

(2) $\ln y = \ln [f(x)]^{g(x)} = g(x) \ln [f(x)]$
(3) $\lim_{x \to a} [\ln y] = \lim_{x \to a} \{g(x) \ln [f(x)]\}$

The limit on the righthand side of the equation will usually be an

indeterminate limit of the type $0 \cdot \infty$. Evaluate

this limit using the

technique previously described. Assume that

 $\lim_{x \to a} \{g(x) \ln [f(x)]\} = L.$ (4)
Finally, $\lim_{x \to a} [\ln y] = L \Rightarrow \ln [\lim_{x \to a} y] = L \Rightarrow \lim_{x \to a} y = e^{L}.$ Example 1: Find $\lim_{x \to a} x^{x}$.

This is an indeterminate form of the type 0^0 . Let $y = x^x \Longrightarrow \ln y = \ln x^x =$

$$x \ln x$$
. $\lim_{x \to 0^+} \ln y = \lim_{x \to 0^+} x \ln x = \lim_{x \to 0^+} \frac{\ln x}{\frac{1}{x}} = \lim_{x \to 0^+} \frac{\frac{1}{x}}{-\frac{1}{x^2}} = \lim_{x \to 0^+} (-x) = 0.$

Thus, $\lim_{x \to 0^+} x^x = e^0 = 1$.

Example 2: Find $\lim_{x \to +\infty} (e^x + 1)^{-\frac{2}{x}}$.

This is an indeterminate form of the type ∞^0 . Let $y = (e^x + 1)^{-2/x} \Rightarrow$

- RGEBER

$$\ln y = \ln\left[\left(e^{x}+1\right)^{-\frac{2}{x}}\right] = \frac{-2\ln(e^{x}+1)}{x} \cdot \lim_{x \to +\infty} \ln y = \lim_{x \to +\infty} \frac{-2\ln(e^{x}+1)}{x} = \lim_{x \to +\infty} \frac{-2\left(\frac{e^{x}}{e^{x}+1}\right)}{1} = \lim_{x \to +\infty} \frac{-2e^{x}}{e^{x}+1} = \lim_{x \to +\infty} \frac{-2e^{x}}{e^{x}} = -2 \cdot \text{Thus, } \lim_{x \to +\infty} \left(e^{x}+1\right)^{-\frac{2}{x}} =$$

Example 3: Find $\lim_{x\to 0^+} (\cos x)^{\frac{1}{x}}$.

 e^{-2} .

This is an indeterminate form of the type 1^{∞} . Let $y = (\cos x)^{\frac{1}{x}} \Rightarrow$

$$\ln y = \ln \left[(\cos x)^{\frac{1}{x}} \right] = \frac{\ln(\cos x)}{x} \cdot \lim_{x \to 0^+} \ln y = \lim_{x \to 0^+} \frac{\ln(\cos x)}{x} =$$

 $\lim_{x \to 0^+} (-\tan x) = 0. \text{ Thus, } \lim_{x \to 0^+} (\cos x)^{\frac{1}{x}} = e^0 = 1.$

Practice Sheet for L'Hospital's Rule

(1) $\lim_{x \to 0} \frac{xe^{3x} - x}{1 - \cos(2x)} =$

(2)
$$\lim_{x \to +\infty} \frac{x}{(\ln x)^3} =$$

(3)
$$\lim_{x \to 0} [\ln(1 - \cos x) - \ln(x^{2})] =$$

(4)
$$\lim_{x \to +\infty} \left(1 - \frac{2}{x} \right)^{3x} =$$

(5)
$$\lim_{x \to +\infty} \frac{\cos\left(\frac{1}{x}\right) - 1}{\cos\left(\frac{2}{x}\right) - 1} =$$

(6)
$$\lim_{x \to 0} \frac{\sqrt{1 - x} - \sqrt{1 - x^{2}}}{x} =$$

(7)
$$\lim_{x \to 0} (\cos x)^{1/x^{2}} =$$

(8)
$$\lim_{x \to 1} \frac{5x^{4} - 7x^{3} + x^{2} - x + 2}{3x^{4} - 8x^{3} + 6x^{2} - 1} =$$

(9)
$$\lim_{x \to 0} \frac{9 - \sqrt{81 - 5x}}{x} =$$

(10)
$$\lim_{x \to +\infty} \frac{\ln(x^{3} + 2)}{\ln(5x^{3} - 1)} =$$

(11)
$$\lim_{x \to +\infty} (e^{x} + 1)^{-\frac{2}{x}} =$$

(12)
$$\lim_{x \to 0} \frac{\sin(4x) - 2\sin(2x)}{x^3} =$$
(13)
$$\lim_{x \to \infty} \left[\frac{1}{\sin x} - \frac{1}{x} \right] =$$
(14)
$$\lim_{x \to +\infty} x \left(e^{\frac{1}{x}} - 1 \right) =$$
(15)
$$\lim_{x \to 0^+} \sqrt[3]{x} \ln x =$$
(16)
$$\lim_{x \to 0} \frac{\ln\left(\frac{2x+1}{5x+1}\right)}{x} =$$
(17)
$$\lim_{x \to +\infty} \left(1 + \frac{e}{x} \right)^{\frac{3}{2}} =$$
(18)
$$\lim_{x \to 0^+} \frac{\arctan(\sin 3x)}{x} =$$
(19)
$$\lim_{x \to 0^+} \frac{\int_{x \to 0}^{0} \frac{e^{2x} + x}{x^3} =$$
(20)
$$\lim_{x \to 0^+} \frac{e^{2x} + x}{x} =$$
(21)
$$\lim_{x \to +\infty} \frac{\arctan x}{x} =$$
(22)
$$\lim_{x \to 0^+} \frac{\arctan(\sin (3x))}{x^2} =$$
(23)
$$\lim_{x \to 0^+} \frac{\ln(\cos x)}{x^2} =$$

COPYRIGHT FIMT 2020

97 | Page

$$(24) \quad \lim_{x \to +\infty} \left(1 + \frac{1}{2x}\right)^x =$$

- (25) $\lim_{x \to +\infty} (\ln x)^{\frac{1}{x}} =$
- (26) $\lim_{x \to +\infty} \left(\ln(e^x + 1) x \right) =$

Solution Key for L'Hospital's Rule

(1)
$$\lim_{x \to 0} \frac{xe^{3x} - x}{1 - \cos(2x)} = \lim_{x \to 0} \frac{3xe^{3x} + e^{3x} - 1}{2\sin 2x} = \lim_{x \to 0} \frac{9xe^{3x} + 3e^{3x} + 3e^{3x}}{4\cos 2x} = \frac{6}{4} = \frac{3}{2}$$

10.000

(2)
$$\lim_{x \to +\infty} \frac{x}{(\ln x)^3} = \lim_{x \to +\infty} \frac{1}{3(\ln x)^2 (\frac{1}{x})} = \lim_{x \to +\infty} \frac{x}{3(\ln x)^2} = \lim_{x \to +\infty} \frac{1}{6\ln x (\frac{1}{x})} = \lim_{x \to +\infty} \frac{x}{6\ln x} = \lim_{x \to +\infty} \frac{1}{6\ln x} = \lim_{x \to +\infty} \frac{1}{6$$

$$\ln\left\{\lim_{x\to 0} \left(\frac{\sin x}{2x}\right)\right\} = \ln\left(\frac{1}{2}\right) = -\ln 2$$

(4) Let
$$y = \frac{1}{x} \Rightarrow \lim_{x \to +\infty} \left(1 - \frac{2}{x}\right)^{3x} = \lim_{y \to 0^+} (1 - 2y)^{3/y}$$
. Now, let $z = (1 - 2y)^{3/y} \Rightarrow \ln z = -6$

 $\ln(1-2y)^{3/y} = \frac{3\ln(1-2y)}{y} \Longrightarrow \lim_{y \to 0^+} \ln z = \lim_{y \to 0^+} \frac{3\ln(1-2y)}{y} = \lim_{y \to 0^+} \frac{1-2y}{1} = -6.$ Thus,

$$\lim_{y \to 0^+} \ln z = -6 \Longrightarrow \ln\left(\lim_{y \to 0^+} z\right) = -6 \Longrightarrow \lim_{y \to 0^+} z = e^{-6} \Longrightarrow \lim_{x \to +\infty} \left(1 - \frac{2}{x}\right)^{3x} = \lim_{y \to 0^+} \left(1 - 2y\right)^{3/y} = 1$$

 $\lim z = e^{-6}$. $y \rightarrow 0^+$

(5) Let
$$y = \frac{1}{x} \Rightarrow \lim_{x \to \infty} \frac{\cos\left(\frac{1}{x}\right) - 1}{\cos\left(\frac{2}{x}\right) - 1} = \lim_{y \to 0^+} \frac{\cos(y) - 1}{\cos(2y) - 1} = \lim_{y \to 0^+} \frac{-\sin(y)}{-2\sin(2y)} =$$

$$\lim_{y \to 0^+} \frac{\sin y}{4\sin y \cos y} = \lim_{x \to 0^+} \frac{1}{4\cos y} = \frac{1}{4}.$$
(6)
$$\lim_{x \to 0^+} \frac{\sqrt{1 - x} - \sqrt{1 - x^2}}{x} = \lim_{x \to 0^+} \frac{\sqrt{1 - x} - \sqrt{1 - x^2}}{x} \cdot \frac{\sqrt{1 - x} + \sqrt{1 - x^2}}{\sqrt{1 - x} + \sqrt{1 - x^2}} =$$

$$\lim_{x \to 0^+} \frac{(1 - x) - (1 - x^2)}{x(\sqrt{1 - x} + \sqrt{1 - x^2})} = \lim_{x \to 0^+} \frac{x^2 - x}{x(\sqrt{1 - x} + \sqrt{1 - x^2})} = \lim_{x \to 0^+} \frac{x - 1}{\sqrt{1 - x} + \sqrt{1 - x^2}} = -\frac{1}{2}.$$
(7) Let $y = (\cos x)^{\frac{1}{2}x^2} \Rightarrow \ln y = \ln(\cos x)^{\frac{1}{2}x^2} = \frac{\ln(\cos x)}{x^2} \Rightarrow \lim_{x \to 0^+} (\ln y) = \lim_{x \to 0^+} \frac{\ln(\cos x)}{x^2} =$

$$\lim_{x \to 0^+} \frac{-\sin x}{2x} = \lim_{x \to 0^+} \frac{-\sin x}{2x\cos x} = \lim_{x \to 0^+} \left(\frac{\sin x}{x}\right) \left(\frac{-1}{2\cos x}\right) = -\frac{1}{2}.$$
Thus, $\lim_{x \to 0^+} (\ln y) = -\frac{1}{2} \Rightarrow$

$$\ln\left(\lim_{x \to 0^+} y\right) = -\frac{1}{2} \Rightarrow \lim_{x \to 0^+} 2x\cos x = \lim_{x \to 0^+} (\cos x)^{\frac{1}{2}x^2} = \lim_{x \to 0^+} y = e^{-\frac{1}{2}}.$$
(8) $\lim_{x \to 0^+} \frac{5x^4 - 7x^3 + x^2 - x + 2}{x} = \lim_{x \to 0^+} \frac{20x^3 - 21x^2 + 2x - 1}{12x^3 - 24x^2 + 12x} = \lim_{x \to 0^+} \frac{60x^2 - 42x + 2}{36x^2 - 48x + 12} = \frac{20}{0} \Rightarrow \lim_{x \to 0^+} \frac{9 - \sqrt{81 - 5x}}{x} = \frac{9 + \sqrt{81 - 5x}}{y + \sqrt{81 - 5x}} = \lim_{x \to 0^+} \frac{81 - (81 - 5x)}{x(9 + \sqrt{81 - 5x})} =$

$$\lim_{x \to 0^+} \frac{5x}{x(9 + \sqrt{81 - 5x})} = \lim_{x \to 0^+} \frac{5}{9 + \sqrt{81 - 5x}} = \frac{5}{18}.$$

COPYRIGHT FIMT 2020

99 | Page

$$(10) \lim_{x \to \infty} \frac{\ln(x^{3} + 2)}{\ln(5x^{3} - 1)} = \lim_{x \to \infty} \frac{\frac{3x^{3}}{15x^{3}}}{\frac{5x^{3}}{5x^{3} - 1}} = \lim_{x \to \infty} \frac{3(5x^{3} - 1)}{15(x^{3} + 2)} = \lim_{x \to \infty} \frac{15x^{3} - 3}{15x^{3} + 30} = 1$$

$$(11) \text{ Let } y = (e^{x} + 1)^{-\frac{7}{x}} \Rightarrow \ln y = \ln(e^{x} + 1)^{-\frac{7}{x}} = \frac{-2\ln(e^{x} + 1)}{x} \Rightarrow \lim_{x \to \infty} \ln y =$$

$$\lim_{x \to \infty} \frac{-2\ln(e^{x} + 1)}{x} = \lim_{x \to \infty} \frac{-2e^{x}}{1} = \lim_{x \to \infty} \frac{-2e^{x}}{e^{x} + 1} = \lim_{x \to \infty} \frac{-2e^{x}}{e^{x}} = -2. \text{ Thus, } \lim_{x \to \infty} \ln y =$$

$$-2 \Rightarrow \ln\left(\lim_{x \to \infty} y\right) = -2 \Rightarrow \lim_{x \to \infty} y = e^{-2} \Rightarrow \lim_{x \to \infty} (e^{x} + 1)^{-\frac{7}{x}} = \lim_{x \to \infty} \frac{-2e^{x}}{6^{x}} = -2. \text{ Thus, } \lim_{x \to \infty} \ln y =$$

$$-2 \Rightarrow \ln\left(\lim_{x \to \infty} y\right) = -2 \Rightarrow \lim_{x \to \infty} y = e^{-2} \Rightarrow \lim_{x \to \infty} (e^{x} + 1)^{-\frac{7}{x}} = \lim_{x \to \infty} y = e^{-2}.$$

$$(12) \lim_{x \to 0} \frac{\sin(4x) - 2\sin(2x)}{x^{3}} = \lim_{x \to 0} \frac{4\cos(4x) - 4\cos(2x)}{3x^{2}} = \lim_{x \to \infty} \frac{-16\sin(4x) + 8\sin(2x)}{6x} =$$

$$\lim_{x \to 0} \frac{-64\cos(4x) + 16\cos(2x)}{6} = \frac{-48}{6} = -8.$$

$$(13) \lim_{x \to 0} \left[\frac{1}{\sin x} - \frac{1}{x}\right] = \lim_{x \to 0} \left(\frac{x - \sin x}{x\sin x}\right) = \lim_{x \to 0} \frac{1 - \cos x}{x\cos x + \sin x} =$$

$$\lim_{x \to 0} \frac{\sin x}{x\sin x + \cos x + \cos x} = \frac{0}{2} = 0.$$

$$(14) \lim_{x \to \infty} x(e^{1/x} - 1) = \lim_{x \to 0} \frac{e^{\frac{1}{x}} - 1}{1/x} = \lim_{x \to \infty} \frac{e^{\frac{1}{x}} (-\frac{1}{x})^{\frac{1}{x}}}{-\frac{1}{x}^{\frac{1}{x}}} = \lim_{x \to 0} \frac{-3x^{\frac{1}{3}}}{x - 0} = 1.$$

$$(15) \lim_{x \to 0} \frac{3\sqrt{x} \ln x}{x} = \lim_{x \to 0} \frac{\ln x}{x^{\frac{1}{x}}} = \lim_{x \to 0} \frac{-\frac{1}{2x^{\frac{1}{x}}} = \lim_{x \to 0} \frac{-3x^{\frac{1}{3}}}{x - 0} = 1.$$

$$(16) \lim_{x \to 0} \frac{\ln\left(\frac{2x+1}{x}\right)}{x - 0} = \lim_{x \to 0} \frac{(\frac{5x+1}{2x+1})\left(\frac{2(5x+1)-5(2x+1)}{(5x+1)^{2}}\right)}{1} = \lim_{x \to 0} \frac{-3}{(2x+1)(5x+1)} = -3.$$

COPYRIGHT FIMT 2020

100 | Page

(17) Let
$$y = \frac{1}{x} \Rightarrow \lim_{x \to +\infty} \left(1 + \frac{e}{x} \right)^{\frac{x}{2}} = \lim_{y \to 0^+} \left(1 + ey \right)^{\frac{1}{2}y}$$
. Next, let $z = (1 + ey)^{\frac{1}{2}y} \Rightarrow \ln z =$

$$\ln(1+ey)^{\frac{1}{2}y} = \frac{\ln(1+ey)}{2y} \Longrightarrow \lim_{y \to 0^+} \ln z = \lim_{y \to 0^+} \frac{\ln(1+ey)}{2y} = \lim_{y \to 0^+} \frac{\frac{e}{1+ey}}{2} = \frac{e}{2}.$$
 Thus,

13

$$\lim_{y \to 0^+} \ln z = \frac{e}{2} \Longrightarrow \ln\left(\lim_{y \to 0^+} z\right) = \frac{e}{2} \Longrightarrow \lim_{y \to 0^+} z = e^{\frac{e}{2}} \Longrightarrow \lim_{x \to +\infty} \left(1 + \frac{e}{x}\right)^{\frac{x}{2}} = \lim_{y \to 0^+} \left(1 + ey\right)^{\frac{1}{2}y} = \lim_{y \to 0^+} \left(1 + e^y\right)^{\frac{1}{2}y} = \lim_$$

$$\lim_{y \to 0^+} z = e^{\frac{e}{2}} \cdot y = \left(e^{2x} + x\right)^{\frac{1}{x}} \Rightarrow \ln y = \ln\left(e^{2x} + x\right)^{\frac{1}{x}} = \frac{\ln\left(e^{2x} + x\right)}{x} \Rightarrow \lim_{x \to 0} \ln y =$$

۰

(18)
$$\lim_{x \to 0} \frac{\arctan(\sin 3x)}{x} = \lim_{x \to 0} \frac{\frac{3\cos 3x}{1+\sin^2 3x}}{1} = 3.$$

(19)
$$\lim_{x \to 0^+} \frac{\int_{0}^{\sin(t^2)} dt}{x^3} = \lim_{x \to 0^+} \frac{\sin(x^2)}{3x^2} = \lim_{x \to 0^+} \frac{2x\cos(x^2)}{6x} = \lim_{x \to 0^+} \frac{\cos(x^2)}{3} = \frac{1}{3}.$$

(20) Let
$$y = (e^{2x} + x)^{1/x} \Rightarrow \ln y = \ln(e^{2x} + x)^{1/x} = \frac{\ln(e^{2x} + x)}{x} \Rightarrow \lim_{x \to 0} \ln y =$$

$$\lim_{x \to 0} \frac{\ln(e^{2x} + x)}{x} = \lim_{x \to 0} \frac{\frac{2e^{2x} + 1}{e^{2x} + x}}{1} = 3. \text{ Thus } \lim_{x \to 0} \ln y = 3 \Longrightarrow \ln\left(\lim_{x \to 0} y\right) = 3 \Longrightarrow$$

$$\lim_{x\to 0} y = e^3 \Longrightarrow \lim_{x\to 0} \left(e^{2x} + x \right)^{\frac{1}{x}} = \lim_{x\to 0} y = e^3.$$

COPYRIGHT FIMT 2020

х

$$(21) \lim_{x \to \infty} \frac{\arctan x}{x} = \frac{\pi/2}{+\infty} = 0.$$

$$(22) \lim_{x \to 0} \frac{\arctan(\sin(3x))}{\arctan(2\tan x)} = \lim_{x \to 0} \frac{\frac{3\cos 3x}{1+\sin^2 3x}}{\sqrt{1-4\tan^2 x}} = \frac{3/4}{2/4} = \frac{3}{2}.$$

$$(23) \lim_{x \to 0} \frac{\ln(\cos x)}{x^2} = \lim_{x \to 0} \frac{-\sin x}{2x} = \lim_{x \to 0} \left(\frac{\sin x}{x}\right) \left(\frac{-1}{2\cos x}\right) = -\frac{1}{2}.$$

$$(24) \text{ Let } y = \frac{1}{x} \Rightarrow \lim_{x \to \infty} \left(1 + \frac{1}{2x}\right)^x = \lim_{y \to 0^+} \left(1 + \frac{1}{2}y\right)^{\frac{1}{2}}. \text{ Let } z = \left(1 + \frac{1}{2}y\right)^{\frac{1}{2}} \Rightarrow \ln z =$$

$$\ln\left(1 + \frac{1}{2}y\right)^{\frac{1}{2}} = \frac{\ln\left(1 + \frac{1}{2}y\right)}{y} \Rightarrow \lim_{y \to 0^+} \ln z = \lim_{y \to 0^+} \frac{\ln\left(1 + \frac{1}{2}y\right)}{y} = \lim_{x \to \infty} \frac{\frac{1}{2}}{\frac{1+\frac{1}{2}y}{1}} = \frac{1}{2}.$$

$$(25) \text{ Let } y = (\ln x)^{\frac{1}{4}} \Rightarrow \ln y = \ln(\ln x)^{\frac{1}{4}} = \frac{\ln(\ln x)}{x} \Rightarrow \lim_{x \to \infty} \ln y = \lim_{x \to \infty} \frac{\ln(n x)}{x} =$$

$$\lim_{x \to \infty} \frac{1}{x} \frac{1}{x} = \lim_{x \to \infty} \frac{1}{x} \ln x = 0. \text{ Thus, } \lim_{x \to \infty} \ln y = 0 \Rightarrow \ln\left(\lim_{x \to \infty} y\right) = 0 \Rightarrow \lim_{x \to \infty} y =$$

$$e^0 = 1 \Rightarrow \lim_{x \to \infty} (\ln x)^{\frac{1}{4}} = \lim_{x \to \infty} y = 1.$$

(26)
$$\lim_{x \to +\infty} \left(\ln(e^x + 1) - x \right) = \lim_{x \to +\infty} \left(\ln\left(e^x + 1\right) - \ln e^x \right) = \lim_{x \to +\infty} \left(\ln\left(\frac{e^x + 1}{e^x}\right) \right) = \ln\left(\lim_{x \to +\infty} \left(\frac{e^x + 1}{e^x}\right) \right) = \ln\left(\lim_{x \to +\infty} \left(\frac{e^x}{e^x}\right) \right) = \ln 1 = 0.$$

NAAC ACCREDITED

The Mean-value Theorem

Consider the quantity Q defined by the equation

$$\frac{f(b) - f(a)}{b - a} = Q,$$

or

$$f(b) - f(a) - (b - a)Q = 0$$

F(x)

$$F(x) = f(x) - f(a) - (x - a)Q.$$

 $F(b) = 0 \qquad F(a) = 0$ From (13.2), , and from (13.3), ; therefore, by Rolle's Theorem (see F'(x) substituting F'(x) must be zero for at least one value of x between a and b, say for x. But by differentiating (13.3) we get

F'(x) = f'(x) - Q.

Therefore, since $F'(x_1)=0$, then also $f'(x_1)-Q=0$, and $Q=f'(x_1)$. Substituting this value of Q in (<u>13.1</u>), we get the Theorem of Mean Value^{13<u>.1</sub>},</sup></u>

$$\frac{f(b) - f(a)}{b - a} = f'(x_1), \ a < x_1 < b$$

NAAC ACCREDITED

ABGEME

where in general all we know about $__$ is that it lies between a and b.

The Theorem of Mean Value interpreted Geometrically.

y = f(x)Let the curve in the figure be the locus of



Figure : Geometric illustration of the Mean value theorem.

Take
$$OC = a$$
 and $OD = b$; then $f(a) = CA$ $f(b) = DB$ $AE = b - a$ and $EB = f(b) - f(a)$.
Therefore the slope of the chord AB is $\tan EAB = \frac{EB}{AE} = \frac{f(b) - f(a)}{b - a}$.

1.201C P 14/0 11-21 There is at least one point on the curve between A and B (as P) where the tangent (or curve)

is parallel to the chord AB. If the abscissa of P is x_1 the slope at P is

 $\tan t = f'(x_1) = \tan EAB.$

Equating these last two equations, we get

$$\frac{f(b) - f(a)}{b - a} = f'(x_1),$$

which is the Theorem of Mean Value.

The student should draw curves to show that there may be more than one such point in the interval; and curves to illustrate, on the other hand, that the theorem may not be true if f'(x)

becomes discontinuous for any value of x between a and b or if becomes discontinuous

Clearing of fractions, we may also write the theorem in the form

$$f(b) = f(a) + (b - a)f'(x_1)$$

Let $b = a + \Delta a$; then $b - a = \Delta a$, and since x_1 is a number lying between a and b, we may write

 $x_1 = a + \theta \cdot \Delta a,$

where θ is a positive proper fraction. Substituting in we get another form of the Theorem of Mean Value.

$$f(a + \Delta a) - f(a) = \Delta a f'(a + \theta \cdot \Delta a), \quad 0 < \theta < 1.$$

Rolle's Theorem

A differentiable and continuous function, which attains equal values at two points, must have a point somewhere between them where the slope of the tangent line to the graph of the function is zero.



The application of the Mean Value Theorem applies to Rolle's Theorem but only where the

f'(c) = 0 =slope of the tangent is equal to zero. \Rightarrow MANA

MAX. & MIN

Introduction

A great many practical problems occur where we have to deal with functions of such a nature that they have a greatest (maximum) value or a least (minimum) valueand it is very important to know what particular value of the variable gives such a value of the function.

Example For instance, suppose that it is required to find the dimensions of the rectangle of greatest area that can be inscribed in a circle of radius 5 inches. Consider the circle in Figure:



Figure: A rectangle with circumscribed circle

$$DE = \sqrt{100 - x^2}$$

Inscribe any rectangle, as BCDE. Let CD = x; then rectangle is evidently

, and the area of the

$$A = A(x) = x\sqrt{100 - x^2}.$$

That a rectangle of maximum area must exist may be seen as follows: Let the base $\ CD$ ($=\sqrt{100-x^2}$

= x) increase to 10 inches (the diameter); then the altitude DE() will decrease to zero and the area will become zero. Now let the base decrease to zero; then the altitude will increase to 10 inches and the area will again become zero. It is therefore intuitionally evident that there exists a greatest rectangle. By a careful study of the figure we might suspect that when the rectangle becomes a square its area would be the greatest, but this would at best be mere guesswork. A better way would evidently be to plot the graph of the function A = A(x) and note its behavior. To aid us in drawing the graph of A(x)

, we observe that

(a)

from the nature of the problem it is evident that x and A must both be positive; and

(b)

the values of x range from zero to 10 inclusive.

Now construct a table of values and draw the graph. What do we learn from the graph?



Figure: The area of a rectangle with fixed circumscribed circle.

(a) If the rectangle is carefully drawn, we may find quite accurately the area of the rectangle corresponding to any value x by measuring the length of the corresponding ordinate. Thus,

when x = OM = 3 inches, then A = MP = 28.6 square inches; and when $A = NQ \approx 35.8$ inches, then sq. in. (found by measurement).

(b) There is one horizontal tangent (RS). The ordinate TH from its point of contact T is greater than any other ordinate. Hence this discovery: One of the inscribed rectangles has evidently a greater area than any of the others. In other words, we may infer from this that

the function defined by A = A(x) has a maximum value. We cannot find this value (= HT) exactly by measurement, but it is very easy to find, using Calculus methods. We observed that at T the tangent was horizontal; hence the slope will be zero at that point . To find the

abscissa of T we then find the first derivative of $egin{aligned} A(x) \ & ,$

$$\begin{array}{rcl} & = x\sqrt{100 - x^2}, \\ \frac{d_{11}}{dx} & = \frac{100 - 2x^2}{\sqrt{100 - x^2}}, \\ \frac{100 - 2x^2}{\sqrt{100 - x^2}} & = 0. \end{array}$$

 $x = 5\sqrt{2}$ $DE = \sqrt{100 - x^2} = 5\sqrt{2}$ Solving, . Substituting back, we get . Hence the rectangle of maximum area inscribed in the circle is a square of area $A = CD \times DE = 5\sqrt{2} \times 5\sqrt{2} = 50$ square inches. The length of HT is therefore [32].

Example A wooden box is to be built to contain 108 cu. ft. It is to have an open top and a square base. What must be its dimensions in order that the amount of material required shall be a minimum; that is, what dimensions will make the cost the least?


Hence

$$M = M(x) = x^2 + \frac{432}{x}$$

is a formula giving the number of square feet required in any such box having a capacity of

M(x) 108cu. ft. Draw a graph of

150 9001:2015 & 14001:2015



What do we learn from the graph?

(a) If the box is carefully drawn, we may measure the ordinate corresponding to any length (= x) of the side of the square base and so determine the number of square feet of lumber required.

(b) There is one horizontal tangent (RS). The ordinate from its point of contact T is less than any other ordinate. Hence this discovery: One of the boxes evidently takes less lumber than any of the others. In other words, we may infer that the function defined by a minimum value. Let us find this point on the graph exactly, using our Calculus. M(x)Differentiating to get the slope at any point, we have

$$\frac{dM}{dx} = 2x - \frac{432}{x^2}.$$

At the lowest point T the slope will be zero. Hence

$$2x - \frac{432}{x^2} = 0;$$

that is, when x = 6 the least amount of lumber will be needed.

Substituting in ${M(x) \over }$, we see that this is M=108 sq. ft.

The fact that a least value of M exists is also shown by the following reasoning. Let the base increase from a very small square to a very large one. In the former case the height must be very great and therefore the amount of lumber required will be large. In the latter case, while the height is small, the base will take a great deal of lumber. Hence M varies from a large value, grows less, then increases again to another large value. It follows, then, that the graph must have a ``lowest'' point corresponding to the dimensions which require the least amount of lumber, and therefore would involve the least cost.

Here is how to compute the critical points in SAGE:

[fontsize=\small,fontfamily=courier,fontshape=tt,frame=single,label=\sage]

sage: x = var("x")

sage: $f = x^2 + 432/x$

sage: solve(f.diff(x)==0,x)

[x = 3*sqrt(3)*I - 3, x = -3*sqrt(3)*I - 3, x = -6]

 $(x^2 + 432/x)' = 0$ This says that above at x = 6.

We will now proceed to the treatment in detail of the subject of maxima and minima.

Increasing and decreasing functions

A function is said to be *increasing* when it increases as the variable increases and decreases as the variable decreases. A function is said to be *decreasing* when it decreases as the variable increases and increases as the variable decreases.

The graph of a function indicates plainly whether it is increasing or decreasing.

Example (1) For instance, consider the function $\underline{a^x}$ whose graph is the locus of the equation $y = a^x \quad a > 1$ $\lambda^x \quad -1 < x < 1$ Figure: SAGE plot of As we move along the curve from left to right the curve is rising; that is, as x increases the $\frac{a}{2}$) always increases. Therefore $\underline{a^x}$ is an increasing function for all values of x. function ($(a - x)^{3}$ whose graph is the locus of the (2) On the other hand, consider the function $y = (a - x)^3$ equation y = (2Figure : SAGE plot of Now as we move along the curve from left to right the curve is falling; that is, as x $(a - x)^3$ is a decreasing function for all values of x. (3) That a function may be sometimes increasing and sometimes decreasing is shown by the

graph of

 $y = 2x^3 - 9x^2 + 12x - 3.$



As we move along the curve from left to right the curve rises until we reach the point A when x = 1, then it falls from A to the point B when x = 2, and to the right of B it is always rising. Hence

(a)

from <u>to</u> x = 1 the function is increasing;

(b)

from x = 1 to x = 2 the function is decreasing;

(c)

from x = 2 to $x = +\infty$ the function is increasing.

The student should study the curve carefully in order to note the behavior of the function when x = 1 and x = 2. Evidently A and B are turning points. At A the function ceases to increase and commences to decrease; at B, the reverse is true. At A and B the tangent (or curve) is evidently parallel to the x-axis, and therefore the slope is zero.

Tests for determining when a function is increasing or decreasing

It is evident from that at a point where a function

$$y = f(x)$$

is increasing, the tangent in general makes an acute angle with the $\, x$ -axis; hence

 $= an au = rac{dy}{dx} = f'(x) =$ a positive number.

Similarly, at a point where a function is decreasing, the tangent in general makes an obtuse angle with the x-axis; therefore

$$= an au = rac{dy}{dx} = f'(x) =$$
 a negative number.

In order, then, that the function shall change from an increasing to a decreasing function, or vice versa, it is a necessary and sufficient condition that the first derivative shall change sign. But this can only happen for a continuous derivative by passing through the value zero. Thus in Figure as we pass along the curve the derivative (= slope) changes sign at the points where x = 1 and x = 2. In general, then, we have at ``turning points,"

$$\frac{dy}{dx} = f'(x) = 0.$$

y = f(x)A value of satisfying this condition is called a *critical point* of the function f. The derivative is continuous in nearly all our important applications, but it is interesting to note the case when the derivative (= slope) changes sign by passing through ∞ . This would evidently happen at the points one a curve where the tangents (and curve) are perpendicular to the *x*-axis. At such exceptional critical points

 $\frac{dy}{dx} = f'(x) = \inf;$

or, what amounts to the same thing,

$$\frac{1}{f'(x)} = 0.$$

Maximum and minimum values of a function

A *maximum value* of a function is one that is greater than any values immediately preceding or following. A *minimum value* of a function is one that is less than any values immediately preceding or following.

For example, it is clear that the function has a maximum value (----) when x = 1, and a minimum value (y = l) when x = 2.

The student should observe that a maximum value is not necessarily the greatest possible value of a function nor a minimum value the least. It is seen that the function (--) has values to the right of x = 1 that are greater than the maximum 2, and values to the left of x = 1 that are less than the minimum 1.



Figure : A continuous function.

A function may have several maximum and minimum values. Suppose that represents the

graph of a function f(x)

At B, F the function is at a local maximum, and at D, G a minimum. That some particular minimum value of a function may be greater than some particular maximum value is shown in the figure, the minimum value at D being greater than the maximum value at G.

At the ordinary critical points D, F, H the tangent (or curve) is parallel to the x-axis; therefore

$$slope = \frac{dy}{dx} = f'(x) = 0.$$

At the exceptional critical points A, B, G the tangent (or curve) is perpendicular to the x-axis, giving

$$slope = \frac{dy}{dx} = f'(x) = \infty.$$

One of these two conditions is then necessary in order that the function shall have a maximum or a minimum value. But such a condition is not sufficient; for at H the slope is zero and at A it is infinite, and yet the function has neither a maximum nor a minimum value at either point. It is necessary for us to know, in addition, how the function behaves in the neighborhood of each point. Thus at the points of maximum value, B, F, the function changes from an increasing to a decreasing function, and at the points of minimum value, D, G, the function changes from a decreasing to an increasing function. It therefore follows from that at maximum points

 $=rac{dy}{dx}=f'(x)$ must change from + to -,

and at minimum points

$$=rac{dy}{dx}=f'(x)$$
 slope must change from - to -

when we move along the curve from left to right.

At such points as A and H where the slope is zero or infinite, but which are neither maximum nor minimum points,

 $= \frac{dy}{dx} = f'(x)$ slope does not change sign.

We may then state the conditions in general for maximum and minimum values of for certain values of the variable as follows:

f(x) is a maximum if f'(x) = 0, and f'(x) changes from + to -. (8.1)

f(x) is a minimum if f'(x) = 0, and f'(x) changes from - to -. (8.2)

NAAC ACCREDITED

The values of the variable at the turning points of a function are called *critical values*; thus x = 1 and x = 2 are the critical values of the variable for the function whose graph is shown in Figure. The critical values at turning points where the tangent is parallel to the x-axis are evidently found by placing the first derivative equal to zero and solving for real values of x, just as under. (Similarly, if we wish to examine a function at exceptional turning points where the tangent is perpendicular to the x-axis, we set the reciprocal of the first derivative equal to zero and solve to find critical values.)

To determine the sign of the first derivative at points near a particular turning point, substitute in it, first, a value of the variable just a little less than the corresponding critical value, and then one a little greaterIf the first gives $-\frac{y}{(as at L, Figure)}$ and the second - (as at M), then the function ($-\frac{y}{(as at P)}$) has a maximum value in that interval (as at I). If the first gives $-\frac{y}{(as at P)}$ and the second $-\frac{y}{(as at N)}$, then the function ($-\frac{y}{(as at C)}$) has a minimum value in that interval (as at C).

If the sign is the same in both cases (as at Q and R), then the function (_____) has neither a maximum nor a minimum value in that interval (as at F).

We shall now summarize our results into a compact working rule.

Examining a function for extremal values: first method

Working rule:

- FIRST STEP. Find the first derivative of the function.
- SECOND STEP. Set the first derivative equal to zero^{8.7} and solve the resulting equation for real roots in order to find the critical values of the variable.
- THIRD STEP. Write the derivative in factor form; if it is algebraic, write it in linear form.
- FOURTH STEP. Considering one critical value at a time, test the first derivative, first for a value a trifle less and then for a value a trifle greater than the critical value. If

the sign of the derivative is first - and then -, the function has a maximum value for that particular critical value of the variable; but if the reverse is true, then it has a minimum value. If the sign does not change, the function has neither.

Example In the problem worked out, we showed by means of the graph of the function

$$A = x\sqrt{100 - x^2}$$

that the rectangle of maximum area inscribed in a circle of radius **5** inches contained **5** square inches. This may now be proved analytically as follows by applying the above rule.

$$f(x) = x\sqrt{100 - x^2}$$

Solution.

$$f'(x) = \frac{100 - 2x^2}{\sqrt{100 - x^2}}$$

First step.

$$\frac{100-2x^2}{\sqrt{100-x^2}} = 0 \qquad x = 5\sqrt{2}$$
Second step. _____implies _____, which is the critical value. Only the positive sign of the radical is taken, since, from the nature of the problem, the negative sign has no meaning.

$$f'(x) = \frac{2(5\sqrt{2}-x)(5\sqrt{2}+x)}{\sqrt{(10-x)(10+x)}}$$

Third step.

$$x < 5\sqrt{2}$$
 $f'(x) = \frac{2(+)(+)}{\sqrt{(+)(+)}} = +$ $x > 5\sqrt{2}$ When

Fourth step. When

$$f'(x) = \frac{2(+)(+)}{\sqrt{(-)(+)}} = -$$

 $x = 5\sqrt{2}$ Since the sign of the first derivative changes from ____to ___at , the function has a maximum value

 $f(5\sqrt{2}) = 5\sqrt{2} \cdot 5\sqrt{2} = 50.$

is a critical point, at which the area is 50 square inches and at This tells us that which the area changes from increasing to decreasing. This implies that the area is a ARMAGE maximum at this point

Problems

It is desired to make an open-top box of greatest possible volume from a square piece of tin whose side is a, by cutting equal squares out of the corners and then folding up the tin to form the sides. What should be the length of a side of the squares cut out?

a-2xSolution. Let x = side of small square = depth of box; then ____ ____ = side of square forming

 $V = (a - 2x)^2 x$, which is the function to be made a bottom of box, and volume is maximum by varying x. Applying rule:

$$\frac{dV}{dx} = (a - 2x)^2 - 4x(a - 2x) = a^2 - 8ax + 12x^2$$

First step.

 $-8ax + 12x^2 = 0$ $x=rac{a}{2}$ and $rac{a}{6}$ Second step. Solving gives critical values

It is evident that must give a minimum, for then all the tin would be cut away, leaving x =no material out of which to make a box. By the usual test, is found to give a maximum volume _____. Hence the side of the square to be cut out is one sixth of the side of the given square.

The drawing of the graph of the function in this and the following problems is left to the student.

Assuming that the strength of a beam with rectangular cross section varies directly as the breadth and as the square of the depth, what are the dimensions of the strongest beam that can be sawed out of a round log whose diameter is d?

Solution. If x = breadth and $\stackrel{y}{=}$ depth, then the beam will have maximum strength xy^2

when the function is a maximum. From the construction and the Pythagorean

theorem, ; hence we should test the function

$$f(x) = x(d^2 - x^2)$$

 $f'(x) = -2x^2 + d^2 - x^2 = d^2 - 3x^2$ First step.

 $d^2 - 3x^2 = 0$ Second step. Therefore, $x = \frac{d}{\sqrt{3}} =$ critical value which gives a maximum.

Therefore, if the beam is cut so that depth = _____of diameter of log, and breadth =

v ³ _____of diameter of log, the beam will have maximum strength.

What is the width of the rectangle of maximum area that can be inscribed in a given



HINT. If OC = h, BC = h - x and PP' = 2y; therefore the area of rectangle 2(h - x)y.

 $y^2=2px$ But since P lies on the parabola , the function to be tested is $2(h-x)\sqrt{2px}$

C ACCREDIT

AGEMEN

Ans. Width = $\frac{\frac{2}{3}h}{\frac{2}{3}h}$

Find the altitude of the cone of maximum volume that can be inscribed in a sphere of radius r



Figure An inscribed cone, height <u>and base radius</u> x, in a sphere.

HINT. Volume of cone = $\frac{\frac{1}{3}\pi^2}{2y}$ $x^2 = BC \times CD = y(2r - y)$; therefore the $f(y) = \frac{\pi}{3}y^2(2r - y)$ function to be tested is .

Find the altitude of the cylinder of maximum volume that can be inscribed in a given right cone



Figure: An inscribed cylinder in a cone.

HINT. Let AU = r and BC = h. Volume of cylinder = $\frac{\pi x^2 y}{h}$. But from similar triangles ABC and DBG, r/x = h/(h-y), so $x = \frac{r(h-y)}{h}$. Hence the function to be tested is

 $f(y) = \frac{r^2}{h^2} y(h-y)^2$

Ans. Altitude =

ivide *a* into two parts such that their product is a maximum.

Ans. Each part $=\frac{\pi}{2}$

Divide 10 into two such parts that the sum of the double of one and square of the other may be a minimum.

0 9001:2015 & 14001:2

Ans. 9and 1.

Find the number that exceeds its square by the greatest possible quantity.

Ans.

What number added to its reciprocal gives the least possible sum?

Ans. 1.

Assuming that the stiffness of a beam of rectangular cross section varies directly as the breadth and the cube of the depth, what must be the breadth of the stiffest beam that can be cut from a log 16 inches in diameter?

Ans. Breadth = 8 inches.

A water tank is to be constructed with a square base and open top, and is to hold **64** cubic yards. If the cost of the sides is \$ 1 a square yard, and of the bottom \$ 2 a square yard, what are the dimensions when the cost is a minimum? What is the minimum cost?

Ans. Side of base = 4yd., height = 4yd., cost \$ 96.

A rectangular tract of land is to be bought for the purpose of laying out a quarter-mile track with straightaway sides and semicircular ends. In addition a strip **35** yards wide along each straightaway is to be bought for grand stands, training quarters, etc. If the land costs \$ 200 an acre, what will be the maximum cost of the land required?

Ans. \$ 856.

A torpedo boat is anchored 9miles from the nearest point of a beach, and it is desired to send a messenger in the shortest possible time to a military camp situated 15miles from that point along the shore. If he can walk 5miles an hour but row only 4miles an hour, required the place he must land.

Ans. 3 miles from the camp.

A gas holder is a cylindrical vessel closed at the top and open at the bottom, where it sinks into the water. What should be its proportions for a given volume to require the least material (this would also give least weight)?

Ans. Diameter = double the height.

What should be the dimensions and weight of a gas holder of $\frac{8,000,000}{\text{cubic feet}}$ capacity, built in the most economical manner out of sheet iron $\frac{1}{16}$ of an inch thick and weighing $\frac{5}{2}$ lb. per sq. ft.?

Ans. Height = 137 ft., diameter = 273 ft., weight = 220 tons.

A sheet of paper is to contain 18 sq. in. of printed matter. The margins at the top and bottom are to be 2 inches each and at the sides 1 inch each. Determine the dimensions of the sheet which will require the least amount of paper.

Ans. 5in. by 10in.

A paper-box manufacturer has in stock a quantity of strawboard 30 inches by 14 inches. Out of this material he wishes to make open-top boxes by cutting equal squares out of each corner and then folding up to form the sides. Find the side of the square that should be cut out in order to give the boxes maximum volume.

Ans. 3 inches.

A roofer wishes to make an open gutter of maximum capacity whose bottom and sides are each 4 inches wide and whose sides have the same slope. What should be the width across the top?

Ans. 8inches. 4

Assuming that the energy expended in driving a steamboat through the water varies as the cube of her velocity, find her most economical rate per hour when steaming against a current running *c*miles per hour.

HINT. Let v = most economical speed; then $av^3 =$ energy expended each hour, abeing a constant depending upon the particular conditions, and $\frac{v-c}{av^3} =$ actual $\frac{av^3}{av^3}$

distance advanced per hour. Hence ______is the energy expended per mile of distance advanced, and it is therefore the function whose minimum is wanted.

Prove that a conical tent of a given capacity will require the least amount of canvas when $\sqrt{2}$

the height is times the radius of the base. Show that when the canvas is laid out flat it will be a circle with a sector of $152^{0}9' = 2.65555...$ cut out. A bell tent 10ft. high should then have a base of diameter 14ft. and would require 272sq. ft. of canvas.

A cylindrical steam boiler is to be constructed having a capacity of 1000 cu. ft. The material for the side costs \$ 2 a square foot, and for the ends \$ 3 a square foot. Find radius when the cost is the least.

Ans.
$$\frac{1}{\sqrt[3]{3\pi}}$$
ft.

In the corner of a field bounded by two perpendicular roads a spring is situated **6**rods from one road and **8**rods from the other.

- (a) How should a straight road be run by this spring and across the corner so as to cut off as little of the field as possible?
- (b) What would be the length of the shortest road that could be run across?

($6^{\frac{2}{3}} + 8^{\frac{2}{3}}$) $\frac{3}{2}$ Ans. (a) 12and 16rods from corner. (b) _____rods.

Show that a square is the rectangle of maximum perimeter that can be inscribed in a given circle.

Two poles of height a and b feet are standing upright and are *c*feet apart. Find the point on the line joining their bases such that the sum of the squares of the distances from this point to the tops of the poles is a minimum. (Ans. Midway between the poles.) When will the sum of these distances be a minimum?

A conical tank with open top is to be built to contain V cubic feet. Determine the shape if the material used is a minimum.

An isosceles triangle has a base 12 in. long and altitude 10 in. Find the rectangle of maximum area that can be inscribed in it, one side of the rectangle coinciding with the base of the triangle.

ivide the number 4 into two such parts that the sum of the cube of one part and three times the square of the other shall have a maximum value.

Divide the number **a** into two parts such that the product of one part by the fourth power of the other part shall be a maximum.

can buoy in the form of a double cone is to be made from two equal circular iron plates of radius r. Find the radius of the base of the cone when the buoy has the greatest displacement (maximum volume).

Into a full conical wineglass of depth *a* and generating angle *a* there is carefully dropped a sphere of such size as to cause the greatest overflow. Show that the radius of the sphere is

ANAGE

 $\frac{\alpha \sin \alpha}{\sin \alpha \cos 2\alpha}$

A wall 27ft. high is 8ft. from a house. Find the length of the shortest ladder that will reach the house if one end rests on the ground outside of the wall.

13√13 Ans.

Here's how to solve this using SAGE: Let *h* be the height above ground at which the ladder hits the house and let *d* be the distance from the wall that the ladder hits the ground the other side of the triangles, on wall. By similar $h/27 = (8+d)/d = 1 + \frac{8}{d}$ d + 8 =. The length of the ladder is, by the

$$f(h) = \sqrt{h^2 + (8+d)^2} = \sqrt{h^2 + (8\frac{h}{h-27})^2}$$

Pythagorean theorem,

[fontsize=\small,fontfamily=courier,fontshape=tt,frame=single,label=\sage]

sage: h = var("h")

sage: f(h) = sqrt(h^2+(8*h/(h-27))^2)

sage: f1(h) = diff(f(h),h)

sage: f2(h) = diff(f(h),h,2)

sage: crit_pts = solve(f1(h) == 0,h); crit_pts

```
[h == 21 - 6^* \operatorname{sqrt}(3)^* I, h == 6^* \operatorname{sqrt}(3)^* I + 21, h == 39, h == 0]
```

```
sage: h0 = crit_pts[2].rhs(); h0
```

39

sage: f(h0)

IAAC ACCREDITED 13*sqrt(13) sage: f2(h0) RMAG

3/(4*sqrt(13))

f(h) has four critical points, but only one of which is meaningful, $h_0 = 39$ This savs At this point, is a minimum.

A vessel is anchored 3 miles offshore, and opposite a point 5 miles further along the shore another vessel is anchored 9 miles from the shore. A boat from the first vessel is to land a passenger on the shore and then proceed to the other vessel. What is the shortest course of the boat?

Ans. 13 miles.

A steel girder 25 ft. long is moved on rollers along a passageway 12.8 ft. wide and into a corridor at right angles to the passageway. Neglecting the width of the girder, how wide must the corridor be?

Ans 5.4ft

A miner wishes to dig a tunnel from a point A to a point B 300 feet below and 500 feet to the east of A. Below the level of A it is bed rock and above A is soft earth. If the cost of tunneling through earth is \$ 1 and through rock \$ 3 per linear foot, find the minimum cost of a tunnel.

Ans. \$ 1348.53.

A carpenter has 108 sq. ft. of lumber with which to build a box with a square base and open top. Find the dimensions of the largest possible box he can make.

6 × 6 × 3 Ans. _____

Find the right triangle of maximum area that can be constructed on a line of length h as hypotenuse.

Ans. $\frac{h}{\sqrt{2}}$ = length of both legs.

What is the isosceles triangle of maximum area that can be inscribed in a given circle?

Ans. An equilateral triangle.

Find the altitude of the maximum rectangle that can be inscribed in a right triangle with base b and altitude h.

Ans. Altitude = $\overline{2}$

Find the dimensions of the rectangle of maximum area that can be inscribed in the $b^2x^2 + a^2y^2 = a^2b^2$

ellipse

 $a\sqrt{2} \times b\sqrt{2}$ Ans. ; area = 2*ab*.

Find the altitude of the right cylinder of maximum volume that can be inscribed in a sphere of radius r.

Ans. Altitude of cylinder = $\frac{\frac{2r}{\sqrt{3}}}{.}$

Find the altitude of the right cylinder of maximum convex (curved) surface that can be inscribed in a given sphere.

Ans. Altitude of cylinder =

What are the dimensions of the right hexagonal prism of minimum surface whose volume is **36**cubic feet?

Ans. Altitude = $2\sqrt{3}$; side of hexagon = 2.

Find the altitude of the right cone of minimum volume circumscribed about a given sphere.

Ans. Altitude = 4r, and volume = -vol. of sphere.

A right cone of maximum volume is inscribed in a given right cone, the vertex of the inside cone being at the center of the base of the given cone. Show that the altitude of the inside cone is one third the altitude of the given cone.

$$y^2 = 2px$$

Given a point on the axis of the parabola at a distance **a** from the vertex; find the abscissa of the point of the curve nearest to it.

Ans. x = a - p

What is the length of the shortest line that can be drawn tangent to the ellipse $b^2x^2 + a^2y^2 = a^2b^2$

and meeting the coordinate axes?

Ans.
$$a +$$

Ь

A Norman window consists of a rectangle surmounted by a semicircle. Given the perimeter, required the height and breadth of the window when the quantity of light admitted is a maximum.

Ans. Radius of circle = height of rectangle.

A tapestry 7feet in height is hung on a wall so that its lower edge is 9feet above an observer's eye. At what distance from the wall should he stand in order to obtain the most favorable view? (HINT. The vertical angle subtended by the tapestry in the eye of the observer must be at a maximum.)

Ans. 12feet.

What are the most economical proportions of a tin can which shall have a given capacity, making allowance for waste? (HINT. There is no waste in cutting out tin for the side of the can, but for top and bottom a hexagon of tin circumscribing the circular pieces required is used up. NOTE 1. If no allowance is made for waste, then height = diameter. NOTE 2. We know that the shape of a bee cell is hexagonal, giving a certain capacity for honey with the greatest possible economy of wax.)

$$\frac{2\sqrt{3}}{\pi}$$
 ×

Ans. Height = _____diameter of base.

An open cylindrical trough is constructed by bending a given sheet of tin at breadth 2a. Find the radius of the cylinder of which the trough forms a part when the capacity of the trough is a maximum.

Ans. Rad. = $\frac{\pi}{1}$; i.e. it must be bent in the form of a semicircle.

A weight W is to be raised by means of a lever with the force F at one end and the point of support at the other. If the weight is suspended from a point at a distance a from the point of support, and the weight of the beam is w pounds per linear foot, what should be the length of the lever in order that the force required to lift it shall be a minimum?



An electric arc light is to be placed directly over the center of a circular plot of grass 100 feet in diameter. Assuming that the intensity of light varies directly as the sine of the angle under which it strikes an illuminated surface, and inversely as the square of its distance from the surface, how high should the light he hung in order that the best possible light shall fall on a walk along the circumference of the plot?



The lower corner of a leaf, whose width is **a**, is folded over so as just to reach the inner edge of the page.



Figure: A leafed page of width **a**.

(a) Find the width of the part folded over when the length of the crease is a minimum.

(b) Find the width when the area folded over is a minimum.

Ans. (a) $\frac{\frac{3}{4}a}{;}$ (b)

A rectangular stockade is to be built which must have a certain area. If a stone wall already constructed is available for one of the sides, find the dimensions which would make the cost of construction the least.

Ans. Side parallel to wall = twice the length of each end.

When the resistance of air is taken into account, the inclination of a pendulum to the $heta=ae^{-kt}\cos{(nt+\eta)}$ vertical may be given by the formula

elongations occur at equal intervals ⁿ of time.

It is required to measure a certain unknown magnitude x with precision. Suppose that nequally careful observations of the magnitude are made, giving the results $a_1, a_2, a_3, \ldots, a_n$. The errors of these observations are evidently $x - a_1, x - a_2, x - a_3, \cdots, x - a_n$, some of which are positive and some negative. It

has been agreed that the most probable value of x is such that it renders the sum of the

 $(x-a_1)^2 + (x-a_2)^2 + (x-a_3)^2 + \dots + (x-a_n)^2$

squares of the errors, namely

minimum. Show that this gives the arithmetical mean of the observations as the most probable value of \boldsymbol{x} .

(This is related to the method of least squares, discovered by Gauss, a commonly used technique in statistical applications.)

The bending moment at x of a beam of length ℓ , uniformly loaded, is given by the

formula $M = \frac{1}{2}w\ell x - \frac{1}{2}wx^2$, where w= load per unit length. Show that the maximum bending moment is at the center of the beam.

If the total waste per mile in an electric conductor is _____, where c= current in amperes (a constant), r= resistance in ohms per mile, and t= a constant depending on the interest on the investment and the depreciation of the plant, what is the relation between c, r, and t when the waste is a minimum?

 $W = c^2 r + \frac{c^2}{2}$

Ans. cr = t.

A submarine telegraph cable consists of a core of copper wires with a covering made of nonconducting material. If x denote the ratio of the radius of the core to the thickness of the covering, it is known that the speed of signaling varies as

$$x^2 \log \frac{1}{x}$$

Show that the greatest speed is attained when

Assuming that the power given out by a voltaic cell is given by the formula

$$P = \frac{E^2 R}{(r+R)^2},$$

when E= constant electromotive force, r= constant internal resistance, R= external resistance, prove that P is a maximum when r = R.

The force exerted by a circular electric current of radius **a** on a small magnet whose axis coincides with the axis of the circle varies as

$$\frac{x}{(a^2 + x^2)^{\frac{5}{2}}}.$$

where x = distance of magnet from plane of circle. Prove that the force is a maximum when $x = \frac{a}{2}$.

We have two sources of heat at A and B, which we visualize on the real line (with B to the right or A), with intensities **a** and **b** respectively. The total intensity of heat at a point P between A and B at a distance of **x** from A is given by the formula $I = \frac{a}{x^2} + \frac{b}{(d-x)^2}$ Show that the temperature at P will be the lowest when

that is, the distances BP and AP have the same ratio as the cube roots of the

corresponding heat intensities. The distance of P from A is

The range of a projectile in a vacuum is given by the formula $R = \frac{v_0^2 \sin 2\phi}{g}$, where $\underline{v_0}$ = initial velocity, $\underline{v_0}$ = acceleration due to gravity, ϕ = angle of projection with the

horizontal. Find the angle of projection which gives the greatest range for a given initial velocity.



The total time of flight of the projectile in the last problem is given by the formula $T = \frac{2v_0 \sin \phi}{2}$

. At what angle should it be projected in order to make the time of flight a maximum?

$$\phi=90^o=\pi/2$$
 Ans.

Fourth step. Examine first for critical value x = 1.

f'(x) = 5(+)(+)2(+) = +. Therefore, when x = 1 the function has a minimum Examine now for the critical value $x=rac{1}{5}$. When $x < \frac{1}{5}$ value $f'(x) = 5(-)(+)^2(-) = +$ $x > \frac{1}{5}$ $f'(x) = 5(-)(+)^2(+) = -$. When . Therefore, $x = \frac{1}{5}$ $f(rac{1}{5})=1.11$ the function has a maximum value . Examine lastly for the when

critical value x = -1 When x < -1 f'(x) = 5(-)(-)2(-) = +

f'(x) = 5(-)(+)2(+) = -

 $\phi = 45^o = \pi/4$

$$(x-1)^2(x+1)^3$$

Examine the function for maximum and minimum values. Use the first method.

$$f(x) = (x-1)^2(x+1)^3$$

Solution.

Second step.

Third step.

When

Ans.

First

$$f'(x) = 2(x-1)(x+1)^3 + 3(x-1)^2(x+1)^2 = (x-1)(x+1)^2(5x-1)$$

$$(x-1)(x+1)^2(5x-1) = 0$$
 $x = 1, -1, \frac{1}{5}$

, which are critical values.

When

$$T = \left[\frac{2}{g \sin 2\phi} \right]$$

x-axis) is given by the formula . Neglecting friction, etc., what must be the value of ϕ to make the quickest descent?

 ϕ The ti the

COPYRIGHT FIMT 2020

. When

$$f'(x) = 5(x-1)(x+1)^2(x-\frac{1}{5})$$

step.

$$f'(x) = 5(-)(+)2(-) = +$$

. Therefore, when $x = -1$ the function has neither a

maximum nor a minimum value.

,

Examine the following functions for maximum and minimum values:

$$(x - 3)^{2}(x - 2)$$
Ans. $x = \frac{7}{3}$, gives max. $= \frac{4}{27}$; $x = 3$, gives min. $= 0$.
 $(x - 1)^{3}(x - 2)^{2}$
Ans. $\frac{x = \frac{8}{5}}{5}$, gives max. $= 0.03456$; $x = 2$, gives min. $= 0$; $x = 1$, gives neither.
 $(x - 4)^{5}(x + 2)^{4}$
Ans. $\frac{x = -2}{-3}$, gives max.; $x = \frac{2}{3}$ gives min; $x = 4$, gives neither.
 $(x - 2)^{5}(2x + 1)^{4}$
Ans. $x = -\frac{1}{2}$, gives max.; $x = \frac{11}{18}$, gives min.; $x = 2$, gives neither.
 $\frac{(x + 1)^{\frac{2}{3}}(x - 5)^{2}}{-4}$

Figure: SAGE plot of

$$\begin{array}{c}
y = (x+1)^{\frac{3}{2}}(x-5)^{\frac{1}{2}} \\
\text{Ans.} \quad x = \frac{1}{2} \\
\text{Ans.} \quad x = \frac{1}{2} \\
x = -1 \\
\text{Ans.} \quad x = \frac{2a}{3} \\
x = \frac{2a}{3} \\
x = 1 \\
x$$

Ans. x = 4, gives max. x = 16, gives min.

$$\frac{(a-x)^3}{a-2x}$$

Ans.
$$x = \frac{a}{4}$$
, gives min.

$$\frac{1-\frac{a}{2}+\frac{a^{2}}{2}}{\frac{1}{2}+\frac{a}{2}+\frac{a}{2}}$$
Ans.
$$x = \frac{1}{2}$$
, gives min.

$$\frac{x^{2}-\frac{3x+2}}{x^{2}+\frac{3x+2}}$$

$$x = \sqrt{2}$$
, gives min. = $(2\sqrt{2}-17, x = -\sqrt{2}, gives max. = (-12\sqrt{2}-17, x)$, gives max. = $(-12\sqrt{2}-17, x)$, gives max. = $(-1, -2)$, give neither.

$$\frac{x}{x^{2}} = \frac{a^{2}}{a+b}$$
, gives max. = $(\frac{(a-b)^{2}}{a+b}$, gives max.

$$\frac{x^{2}}{x^{2}} = \frac{a^{2}}{a-b}$$
, gives max. = $(\frac{(a-b)^{2}}{a+b}$, gives max.

$$\frac{x^{2}}{x^{2}} = \frac{a^{2}}{a-b}$$
, gives min.; $(x = \frac{a^{2}}{a+b})$, gives max.

$$\frac{x^{3}}{x^{2}} = 3x^{2} - 9x + 5$$
.
Examine $x^{3} - 3x^{2} - 9x + 5$.
Solution.

$$f(x) = x^{3} - 3x^{2} - 9x + 5$$
.
Solution.

$$f(x) = x^{3} - 3x^{2} - 9x + 5$$
.
Second step.

$$f'(x) = 3x^{2} - 6x - 9$$
.
Second step.

$$f'(x) = 6x - 6$$
.
Third step.

Fourth step.

$$f''(-1) = -12$$
Fourth step.

$$f(-1) = 10 = \max \text{ maximum value.} \qquad f''(3) = +12$$
Therefore,

$$f(3) = -22 = \min \text{ minimum value.}$$
Examine sti ² *x* cos *x* for maximum and minimum values.

$$f(x) = \sin^2 x \cos x$$
Solution.

$$f'(x) = 2\sin x \cos^2 x - \sin^3 x$$
First step.

$$2\sin x \cos^2 x - \sin^3 x = 0$$
Second step.

$$2\sin x \cos^2 x - \sin^3 x = 0$$
Second step.

$$f''(x) = \cos x (2\cos^2 x - 7\sin^2 x)$$
Third step.

$$f''(0) = + \qquad \text{Therefore,} \qquad f(0) = 0 \qquad \text{minimum value.}$$
Fourth step.

$$f''(0) = + \qquad \text{Therefore,} \qquad f(0) = - \qquad \text{Therefore,} \qquad f(\alpha) = -$$
Therefore,
$$f(\alpha) = - \qquad \text{Therefore,} \qquad f(\alpha) = - \qquad f(\alpha) = -$$

$$\frac{x^3}{3} - 21x^2 + 3x + 1$$

Ans.
$$oldsymbol{x}=1$$
 , gives max. = $\ddot{\overline{5}}$; $oldsymbol{x}=oldsymbol{3}$, gives min. = $\ 1$.

 $2x^3 - 15x^2 + 36x + 10$

Ans. x = 2, gives max. = 38; x = 3, gives min. = 37.

D.G.E

$$x^3 - 9x^2 + 15x - 3$$

Ans. x = 1, gives max. = 4; x = 5, gives min. =

$$x^3 - 3x^2 + 6x + 10$$

Ans. No max. or min.

0.00

$$x^5 - 5x^4 + 5x^3 + 1$$

. $x = 1$, gives max. = 2; $x = 3$, gives min. = $\frac{-26}{-26}$; $x = 0$,

gives neither.

$$3x^5 - 125x^2 + 2160x$$

$$x = -4$$

and 3, give max.; $x = -3$
and 4, give min.
$$2x^{3} - 3x^{2} - 12x + 4$$

$$2x^{3} - 21x^{2} + 36x - 20$$

$$x^{4} - 2x^{2} + 10$$

COPYRIGHT FIMT 2020

.

$$x^{3}-8$$

$$4-z^{6} \sin x(1+\cos x)$$

$$x = 2n\pi + \frac{\pi}{3}, \text{ give max.} = \frac{\pi}{3}\sqrt{3}, x = 2n\pi - \frac{\pi}{3}, \text{ give min.} = \frac{\pi}{3}\sqrt{3}, x = n\pi, \text{ give min.} = \frac{\pi}{3}\sqrt{3}, x = n\pi, \text{ give min.} = \frac{\pi}{3}\sqrt{3}, x = n\pi, \text{ give min.} = \frac{\pi}{3}\sqrt{3}$$
Ans. $x = e, \text{ gives min.} = e; x = 1, \text{ gives neither.}$
Ans. $x = e, \text{ gives max.}$

$$ae^{kx} + be^{-kx}$$

$$x = \frac{1}{k}\log\sqrt{\frac{5}{a}}, \text{ gives min.} = 2\sqrt{ab}$$
Ans. $\frac{x^{2}}{2}, \frac{\pi}{3}, \frac{\pi}{2}, \frac{\pi}{3}, \frac{\pi}{3}$

Ans.
$$x = \frac{\pi}{6}$$
, gives max.; $x = -\frac{\pi}{6}$, gives min.
 $\frac{x + \tan x}{4}$
Ans. No max. or min.
 $\sin^3 x \cos x$.
Ans. $x = n\pi + \frac{\pi}{3}$, gives max. $= \frac{3}{16}\sqrt{3}$, $x = n\pi - \frac{\pi}{3}$, gives min. $= -\frac{3}{16}\sqrt{3}$; $x = n\pi$
, gives neither.
 $x \cos x$.
 $y = x \cos(x)$
Ans. x such that $\frac{x \sin x}{x} = \frac{x}{2}$, gives max/min.
 $\sin x + \cos 2x$
Ans. $x = \frac{\pi}{4}$, gives max.; $x = \frac{\pi}{2}$, gives min.
 $2 \tan x - \tan^2 x$
Ans. $\frac{x = \frac{\pi}{4}}{x}$, gives max.

Ans.
$$x = \frac{\pi}{4}$$
, gives max.
 $\frac{x}{1+x \tan x}$.
 $x = \cos x$, gives max.; $\frac{x = -\cos x}{x}$, gives min.

Definition of successive derivatives

We have seen that the derivative of a function of \boldsymbol{x} is in general also a function of \boldsymbol{x} . This new function may also be differentiable, in which case the derivative of the first derivative is called the second derivative of the original function. Similarly, the derivative of the second derivative is called the third derivative; and so on to the \boldsymbol{n} -th derivative. Thus, if

NAAC ACCREDITED

$$y = 3x^{4},$$

$$\frac{dy}{dx} = 12x^{3},$$

$$\frac{d}{dx}\left(\frac{dy}{dx}\right) = 36x^{2},$$

$$\frac{d}{dx}\left[\frac{d}{dx}\left(\frac{dy}{dx}\right)\right] = 72x,$$

Notation

The symbols for the successive derivatives are usually abbreviated as follows:

$$\frac{d}{dx} \left(\frac{dy}{dx} \right) = \frac{d^2 y}{dx^2},$$

$$\frac{d}{dx} \left[\frac{d}{dx} \left(\frac{dy}{dx} \right) \right] = \frac{d}{dx} \left(\frac{d^2 y}{dx^2} \right) = \frac{d^3 y}{dx^3},$$

$$\dots \dots$$

$$\frac{d}{dx} \left(\frac{d^{n-1} y}{dx^{n-1}} \right) = \frac{d^n y}{dx^n}.$$

y = f(x)If , the successive derivatives are also denoted by

$$f'(x), f''(x), f'''(x), f^{(4)}(x), ..., f^{(n)}(x);$$

or

$$y', y'', y''', y^{(4)}, ..., y^{(n)};$$

or,

$$rac{d}{dx}f(x), \; rac{d^2}{dx^2}f(x), \; rac{d^3}{dx^3}f(x), \; rac{d^n}{dx^n}f(x), \; ..., \; rac{d^n}{dx^n}f(x).$$

The *n*-th derivative

For certain functions a general expression involving n may be found for the n-th derivative. The usual plan is to find a number of the first successive derivatives, as many as may be necessary to discover their law of formation, and then by induction write down the n-th derivative.

 $\frac{d^n y}{dx^n} = a^n e^{ax}$

Example Given $y = e^{ax}$, find $\frac{d^n y}{dx^n}$

Solution.

Example Given $y = \log x$, find $\frac{d^n y}{dx^n}$

 $rac{dy}{dx}=ae^{ax}$ $rac{d^2y}{dx^2}=a^2e^{ax}$

 $\frac{dy}{dx} = \frac{1}{x} \quad \frac{d^2y}{dx^2} = -\frac{1}{x^2} \quad \frac{d^3y}{dx^3} = \frac{1\cdot 2}{x^3} \quad \frac{d^4y}{dx^4} = \frac{1\cdot 2\cdot 3}{x^3} \quad \frac{d^ny}{dx^n} = (-1)^{n-1} \frac{(n-1)!}{x^n}$ Solution.

Example Given
$$\frac{y = \sin x}{dx^n}$$

 $\frac{d^n y}{dx^n}$.
 $\frac{dy}{dx} = \cos x = \sin \left(x + \frac{\pi}{2}\right)$
Solution.
 $\frac{d^2 y}{dx^2} = \frac{d}{dx} \sin \left(x + \frac{\pi}{2}\right) = \cos \left(x + \frac{\pi}{2}\right) = \sin \left(x + \frac{2\pi}{2}\right),$
 $\frac{d^3 y}{dx^3} = \frac{d}{dx} \sin \left(x + \frac{2\pi}{2}\right) = \cos \left(x + \frac{2\pi}{2}\right) = \sin \left(x + \frac{3\pi}{2}\right)$

$$\frac{d^n y}{dx^n} = \sin\left(x + \frac{n\pi}{2}\right).$$

Leibnitz's Formula for the n-th derivative of a product

This formula expresses the n-th derivative of the product of two variables in terms of the variables themselves and their successive derivatives.

If u and v are functions of x, we have, from equation (V) in above,

 $\frac{d}{dx}(uv) = \frac{du}{dx}v + u\frac{dv}{dx}.$

Differentiating again with respect to x,

$$\frac{d^2}{dx^2}(uv) = \frac{d^2u}{dx^2}v + \frac{du}{dx}\frac{dv}{dx} + \frac{du}{dx}\frac{dv}{dx} + u\frac{d^2v}{dx^2} = \frac{d^2u}{dx^2}v + 2\frac{du}{dx}\frac{dv}{dx} + u\frac{d^2v}{dx^2}$$

Similarly,

$$\frac{d^3}{dx^3}(uv) = \frac{d^3u}{dx^3} + \frac{d^2u}{dx^2}\frac{dv}{dx} + 2\frac{d^2u}{dx^2}\frac{dv}{dx} + 2\frac{du}{dx}\frac{d^2v}{dx^2} + \frac{du}{dx}\frac{d^2v}{dx^2} + u\frac{d^3v}{dx^3} \\ = \frac{d^3u}{dx^3}v + 3\frac{d^2u}{dx^2}\frac{dv}{dx} + 3\frac{du}{dx}\frac{d^2v}{dx^2} + u\frac{d^3v}{dx^3}.$$

However far this process may be continued, it will be seen that the numerical coefficients follow the same law as those of the Binomial Theorem, and the indices of the derivatives correspond^{7.1} to the exponents of the Binomial Theorem. Reasoning then by mathematical (m+1) induction from the m-th to the m-th to the m-st derivative of the product, we can prove Leibnitz's Formula

$$\frac{d^{n}}{dx^{n}}(uv) = \frac{d^{n}u}{dx^{n}}v + n\frac{d^{n-1}u}{dx^{n-1}}\frac{dv}{dx} + \frac{n(n-1)}{2!}\frac{d^{n-2}u}{dx^{n-2}}\frac{d^{2}v}{dx^{2}} + \dots + n\frac{du}{dx}\frac{d^{n-1}v}{dx^{n-1}} + u\frac{d^{n}v}{dx^{n}}$$
(7.1)

COPYRIGHT FIMT 2020

144 | Page
$y=e^x\log x$, find $rac{d^3y}{dx^3}$ by Leibnitz's Formula. Example Given

Solution. Let
$$\underline{\vec{w}} = e^x$$
, and $v = \log x$; then $\frac{du}{dx} = e^x$, $\frac{dv}{dx} = \frac{1}{x}$, $\frac{d^2u}{dx^2} = e^x$, $\frac{d^2v}{dx^2} = -\frac{1}{x^2}$, $\frac{d^3u}{dx^3} = e^x$, $\frac{d^3v}{dx^3} = \frac{2}{x^3}$
Substituting in we get

$$\frac{d^3y}{dy^3} = e^x \log x + \frac{3x^2}{x} - \frac{3e^x}{x^2} = e^x \left(\log x + \frac{3}{x} - \frac{3}{x^2} + \frac{2}{x^3}\right).$$

This can be verified using the SAGE commands:

[fontsize=\small,fontfamily=courier,fontshape=tt,frame=single,label=\sage]

sage: x = var("x")

sage: f = exp(x)*log(x)

sage: diff(f,x,3)

e^x*log(x) + 3*e^x/x - 3*e^x/x^2 + 2*e^x/x^3

 $y = x^2 e^{ax}$ Example Given by Leibnitz's Formula. , find

Solution. Let
$$u = x^2$$
, and $\underline{v = e^{ax}}$; then $\frac{du}{dx} = 2x$, $\frac{dv}{dx} = ae^{ax}$, $\frac{d^2u}{dx^2} = 2x$, $\frac{d^2v}{dx^2} = a^2e^{ax}$, $\frac{d^3u}{dx^2} = 0$, $\frac{d^3v}{dx^3} = 0$, $\frac{d^3v}{dx^3} = a^3e^{ax}$, $\frac{d^nu}{dx^n} = 0$, $\frac{d^nv}{dx^n} = a^ne^{ax}$, Substituting in (7.1), we get

$$\frac{d^{n}y}{dx^{n}} = x^{2}a^{n}e^{ax} + 2na^{n-1}xe^{ax} + n(n-1)a^{n-2}e^{ax} = a^{n-2}e^{ax}[x^{2}a^{2} + 2nax + n(n-1)].$$

Successive differentiation of implicit functions

To illustrate the process we shall find from the equation of the hyperbola

 $b^2 x^2 - a^2 y^2 = a^2 b^2.$

Differentiating with respect to x, as ,

 $2b^2x - 2a^2y\frac{dy}{dx} = 0,$

MANA

or,

$$\frac{dy}{dx} = \frac{b^2x}{a^2y}.$$

Differentiating again, remembering that $_$ is a function of x,

$$rac{d^2y}{dx^2} = rac{a^2yb^2 - b^2\pi a^2rac{dy}{dx}}{a^4y^2}.$$

Substituting for $\frac{d\bar{x}}{dx}$ its value from

$$\frac{d^2y}{dx^2} = \frac{a^2b^2y - a^2b^2x\left(\frac{b^2y}{a^2y}\right)}{a^4y^2} = -\frac{b^2(b^2x^2 - a^2y^2)}{a^4y^3}.$$

 $b^3 x^2 - a^2 y^2 = a^2 b^2$, therefore gives,

$$\frac{d^2y}{dx^2}=-\frac{b^4}{a^2y^3}.$$

This basically says

$$y' = \frac{dy}{dx} = \frac{b^2x}{a^2y},$$

and

$$y'' = \frac{d^2y}{dx^2} = -\frac{b^2 - a^2(y')^2}{a^2y}.$$

Exercises

Verify the following derivatives:

$$y = 4x^{3} - 6x^{2} + 4x + 7$$

$$f(x) = \frac{12(2x - 1)}{1 - x}$$

$$f(x) = \frac{x^{3}}{1 - x}$$

$$f(x) = \frac{1}{1 - x}$$
Ans.
$$f(y) = y^{6}$$

$$f^{(6)}(y) = 6!$$
Ans.
$$y = x^3 \log x$$

$$\frac{d^4y}{dx^4} = \frac{6}{x}$$
Ans.

$$y = \frac{c}{x^n} \quad y'' = \frac{n(n+1)c}{x^{n+2}}$$

$$y = (x - 3)e^{2x} + 4xe^x + x$$

.

Ans.

$$y'' = 4e^{x}[(x-2)e^{x} + x + 2]$$

$$y = \frac{a}{2}(e^{\frac{x}{a}} + e^{-\frac{x}{a}})$$

$$y'' = \frac{1}{2a}(e^{\frac{x}{a}} + e^{-\frac{x}{a}}) = \frac{y}{a^{x}}$$
Ans.

$$f(x) = ax^{2} + bx + c$$
Ans.

$$f''(x) = 0$$
Ans.

$$f(x) = \log(x+1)$$
Ans.

$$\frac{f'(4)(x) = -\frac{6}{(x+1)^{2}}}{f(x) = \log(e^{x} + e^{-x})}$$
Ans.

$$r = \sin a\theta.$$

$$\frac{\frac{d^{4}x}{d\theta^{4}} = a^{4} \sin a\theta = a^{4}r}{Ans.}$$

$$r = \tan \phi$$
Ans.

$$\frac{\frac{d^{4}x}{d\theta^{4}} = 6 \sec^{6} \phi - 4 \sec^{2} \phi}{Ans.}$$

$$r''' = 2 \cot \phi \csc^{2} \phi$$
Ans.

$$f(t) = e^{-t} \cos t$$

$$f^{(4)}(t) = -4e^{-t}\cos t = -4f(t)$$

Ans.

 $f(\theta) = \sqrt{\sec 2\theta}$

$$f''(heta)=3[f(heta)]5-f(heta)$$

Ans.

$$p = (q^2 + a^2) \arctan \frac{q}{a}$$

$$\frac{d^3p}{dq^3} = \frac{4a^3}{(a^2+q^2)^2}.$$

Ans.

$$y = a^x$$

$$\frac{d^n y}{dx^n} = (\log a)^n a^x$$

Ans.

 $y = \log(1+x)$

$$\frac{d^n y}{dx^n} = (-1)^{n-1} \frac{(n-1)!}{(1+x)^n}$$

.

Ans.

$$\frac{d^{n}y}{dx^{n}} = a^{n}\cos\left(ax + \frac{n\pi}{2}\right)$$
Ans.
$$y = x^{n-1}\log x$$

$$\frac{d^{n}y}{dx^{n}} = \frac{(n-1)!}{2}$$

$$\frac{d^n y}{dx^n} = \frac{(n-1)}{x}$$
Ans.

$$y = \tfrac{1-x}{1+x}$$

$$rac{d^n y}{dx^n} = 2(-1)^n rac{n!}{(1+x)^{n+1}}$$

А

Hint: Reduce fraction to form
$$\begin{array}{l} -1 + \frac{2}{1+x} \\ \text{before differentiating.} \end{array}$$

$$\begin{array}{l} y = e^x \sin x \\ \frac{d^2y}{dx^2} - 2\frac{dy}{dx} + 2y = 0 \\ \text{If } & \\ \end{array}$$

$$\begin{array}{l} y = a \cos(\log x) + b \sin(\log x) \\ \text{If } & \\ \end{array}$$

$$\begin{array}{l} x^2 \frac{d^2y}{dx^2} + x \frac{dy}{dx} + y = 0 \\ \end{array}$$

Use Leibnitz's Formula in the next four examples:

 $y = x^2 a^x$

$$\frac{d^n y}{dx^n} = a^x (\log a)^{n-2} [(x \log a + n)^2 - n]$$

Ans.

 $y = xe^x$

$$rac{d^n y}{dx^n} = (x+n)e^x$$
 Ans.

 $f(x) = e^x \sin x$

$$f^{(n)}(x) = (\sqrt{2})^n e^x \sin\left(x + \frac{n\pi}{4}\right)$$

$$f(\theta) = \cos a\theta \cos b\theta$$

$$f^n(\theta) = \frac{(a+b)^n}{2} \cos\left[(a+b)\theta + \frac{n\pi}{2}\right] + \frac{(a-b)^n}{2} \cos\left[(a-b)\theta + \frac{n\pi}{2}\right]$$

Ans.

4

$$a = rac{d^2s}{dt^2}$$
 $a_x = rac{d^2x}{dt^2}$

0

Show that the formulas for acceleration, , may be written _____

$$a_y = \frac{d^2y}{dt^2}$$

$$y^2 = 4ax$$

$$\frac{d^2y}{dx^2} = -\frac{4a^2}{y^3}$$
.

Ans.

 $b^2 x^2 + a^2 y^2 = a^2 b^2$

$$rac{d^2 y}{dx^2} = -rac{b^4}{a^2 y^3} \quad rac{d^3 y}{dx^2} = -rac{3N^2 \pi}{\pi^3 y^3}$$

Ans.

$$x^{2} + y^{2} = r^{2} \quad \frac{d^{2}y}{dx^{2}} = -\frac{r^{2}}{y^{3}}$$

$$y^2 + y = x^2$$

$$\frac{d^3y}{dx^3} = -\frac{24x}{(1+2y)^5}$$

Ans.

$$ax^2 + 2hxy + by^2 = 1$$

$$\frac{d^2y}{dx^2} = \frac{h^2 - ab}{(hx + by)^3}$$

Ans.

$$y^2 - 2xy = a^2$$

$$\frac{d^2y}{dx^2} = \frac{a^2}{(y-x)^3}; \frac{d^3y}{dx^3} = -\frac{3a^2x}{(y-x)^5}$$
Ans.
sec $\phi \cos \theta = c$

$$\frac{d^2\theta}{d\phi^2} = \frac{\tan^2\theta - \tan^2\phi}{\tan^3\theta}.$$

Ans.

 $heta = an(\phi + heta)$

$$\frac{d^3\theta}{d\phi^3} = -\frac{2(5+8\theta^2+3\theta^4)}{\theta^8}$$
Ans.

Find the second derivative in the following:

(a)	$\log(u+v) = u - v.$	(e)	$y^3 - x^3 - 3axy = 0.$
(b)	$e^u + u = e^v + v.$	(f)	$y^2 - 2mxy + x^2 - a = 0.$
(c)	$s = 1 + te^s$.	(g)	$y = \sin(x + y).$
(d)	$e^s + st - e = 0.$	(h)	$\langle x^{x^{(1+y)}} = xy.$

ARHAG

UNIT-IV

Integration

What is integration?

The dictionary definition of *integration* is combining parts so that they work together or form a whole. Mathematically, integration stands for finding the area under a curve from one point to another. It is represented by

$$\int_{a}^{b} f(x) dx$$

where the symbol \int is an integral sign, and a and b are the lower and upper limits of integration, respectively, the function f is the integrand of the integral, and x is the variable of integration. Figure 1 represents a graphical demonstration of the concept.

Riemann Sum

Let f be defined on the closed interval [a,b], and let Δ be an arbitrary partition of [a,b] such as: $a = x_0 < x_1 < x_2 < \dots < x_{n-1} < x_n = b$, where Δx_i is the length of the i^{th} subinterval (Figure 2).

If c_i is any point in the i^{th} subinterval, then the sum

9001-2015

$$\sum_{i=1}^{n} f(c_i) \Delta x_i, x_{i-1} \le c_i \le x_i$$

is called a Riemann sum of the function f for the partition Δ on the interval [a,b]. For a given partition Δ , the length of the longest subinterval is called the norm of the partition. It is denoted by $\|\Delta\|$ (the norm of Δ). The following limit is used to define the



Figure The definite integral as the area of a region under the curve, $Area = \int_{a}^{b} f(x) dx$.

If c_i is any point in the i^{th} subinterval, then the sum



Figure Division of interval into *n* segments.

is called a Riemann sum of the function f for the partition Δ on the interval [a,b]. For a given partition Δ , the length of the longest subinterval is called the norm of the partition. It is denoted by $\|\Delta\|$ (the norm of Δ). The following limit is used to define the definite integral.

$$\lim_{\|\Delta\| \to 0} \sum_{i=1}^{n} f(c_i) \Delta x_i = I$$

NAAC ACCREDITED

This limit exists if and only if for any positive number ε , there exists a positive number δ such that for every partition Δ of [a,b] with $\|\Delta\| < \delta$, it follows that

$$\left|I - \sum_{i=1}^{n} f(c_i) \Delta x_i\right| < \varepsilon$$

for any choice of c_i in the i^{th} subinterval of Δ .

If the limit of a Riemann sum of f exists, then the function f is said to be integrable over [a,b] and the Riemann sum of f on [a,b] approaches the number I.

$$\lim_{\|\Delta\|\to 0}\sum_{i=1}^n f(c_i)\Delta x_i = I$$

where

$$I = \int_{a}^{b} f(x) dx$$

Example

Find the area of the region between the parabola $y = x^2$ and the *x*-axis on the interval [0,4.5]. Use Riemann's sum with four partitions.

Solution

We evaluate the integral for the area as a limit of Riemann sums. We sketch the region (Figure 3), and partition [0,4.5] into four subintervals of length



Figure Graph of the function $y = x^2$.

The points of partition are

$$x_0 = 0, x_1 = 1.125, x_2 = 2.25, x_3 = 3.375, x_4 = 4.5$$

Let's choose c_i 's to be right hand endpoint of its subinterval. Thus,

$$c_1 = x_1 = 1.125, c_2 = x_2 = 2.25, c_3 = x_3 = 3.375, c_4 = x_4 = 4.5$$

The rectangles defined by these choices have the following areas:

 $f(c_1)\Delta x = f(1.125) \times (1.125) = (1.125)^2 (1.125) = 1.4238$

$$f(c_2)\Delta x = f(2.25) \times (1.125) = (2.25)^2 (1.125) = 5.6953$$

$$f(c_3)\Delta x = f(3.375) \times (1.125) = (3.375)^2 (1.125) = 12.814$$

$$f(c_4)\Delta x = f(4.5) \times (1.125) = (4.5)^2 (1.125) = 22.781$$

The sum of the areas then is

$$\int_{0}^{4.5} x^{2} dx \approx \sum_{i=1}^{4} f(c_{i}) \Delta x,$$

= 1.4238+5.6953+12.814+22.781
= 42.715

How does this compare with the exact value of the integral $\int x^2 dx$?

Example

Find the exact area of the region between the parabola $y = x^2$ and the x-axis on the interval [0,b]. Use Riemann's sum.

Solution

Note that in Example 1 for $y = x^2$ that

$$f(c_i)\Delta x = i^2 (\Delta x)^3$$

Thus, the sum of these areas, if the interval is divided into n equal segments is

$$S_n = \sum_{i=1}^n f(c_i) \Delta x$$
$$= \sum_{i=1}^n i^2 (\Delta x)^3$$
$$= (\Delta x)^3 \sum_{i=1}^n i^2$$

Since

$$\Delta x = \frac{b}{n}$$
, and
 $\sum_{i=1}^{n} i^2 = \frac{n(n+1)(2n+1)}{6}$

then

$$S_{n} = \frac{b^{3}}{n^{3}} \frac{n(n+1)(2n+1)}{6}$$
$$= \frac{b^{3}}{6} \frac{2n^{2} + n + 2n + 1}{n^{2}}$$
$$= \frac{b^{3}}{6} \left(2 + \frac{3}{n} + \frac{1}{n^{2}}\right)$$

The definition of a definite integral can now be used

$$\int_{a}^{b} f(x) dx = \lim_{\|\Delta x\| \to 0} \sum_{i=1}^{n} f(c_i) \Delta x$$

To find the area under the parabola from x = 0 to x = b, we have

$$\int_{0}^{b} x^{2} dx = \lim_{|\Delta| \to 0} \sum_{i=1}^{n} f(c_{i}) \Delta x$$
$$= \lim_{n \to \infty} S_{n}$$
$$= \lim_{n \to \infty} \frac{b^{3}}{6} \left(2 + \frac{3}{n} + \frac{1}{n^{2}} \right)$$
$$= \frac{b^{3}}{6} \left(2 + 0 + 0 \right)$$

$$=\frac{b^3}{3}$$

For the value of b = 4.5 as given in Example 1,

$$\int_{0}^{4.5} x^2 dx = \frac{4.5^3}{3}$$

= 30.375

The Mean Value Theorem for Integrals

The area of a region under a curve is usually greater than the area of an inscribed rectangle and less than the area of a circumscribed rectangle. The mean value theorem for integrals states that somewhere between these two rectangles, there exists a rectangle whose area is exactly equal to the area of the region under the curve, as shown in Figure 4. Another variation states that if a function f is continuous between a and b, then there is at least one point in [a,b] where the function equals the average value of the function f over [a,b]

Theorem: If the function f is continuous on the closed interval [a,b], then there exists a number c in [a,b] such that:

$$f(c) = \frac{1}{b-a} \int_{a}^{b} f(x) dx$$

Example

Graph the function $f(x) = (x-1)^2$, and find its average value over the interval [0,3]. At what point in the given interval does the function assume its average value?



Figure Mean value rectangle.

RAMAG

Solution

Average(f) =
$$\frac{1}{b-a} \int_{a}^{b} f(x) dx$$

= $\frac{1}{3-0} \int_{0}^{3} (x-1)^{2} dx$
= $\frac{1}{3} \int_{0}^{3} (x^{2} - 2x + 1) dx$

= 1

$$=\frac{1}{3}\left[\left(\frac{1}{3}\times27-9+3\right)-0\right]$$

The average value of the function f over the interval [0,3] is 1. Thus, the function assumes its average value at

$$f(c) = 1$$

 $(c-1)^2 = 1$
 $c = 0, 2$

The connection between integrals and area can be exploited in two ways. When a formula for the area of the region between the x-axis and the graph of a continuous function is known, it can be used to evaluate the integral of the function. However, if the area of region is not known, the integral of the function can be used to define and calculate the area. Table 1 lists a number of standard indefinite integral forms.



Figure The function $f(x) = (x-1)^2$.

Example 4

Find the area of the region between the circle $x^2 + y^2 = 1$ and the *x*-axis on the interval

150 9001:2015 & 14001:20

[0,1] (the shaded region) in two different ways.

Solution



The first and easy way to solve this problem is by recognizing that it is a quarter circle. Hence the area of the shaded area is

$$A = \frac{1}{4}\pi r^{2}$$
$$= \frac{1}{4}\pi (1)^{2}$$
$$= \frac{\pi}{4}$$

The second way is to use the integrals and the trigonometric functions. First, let's simplify the function $x^2 + y^2 = 1$.

 $x^{2} + y^{2} = 1$ $y^{2} = 1 - x^{2}$ $y = \sqrt{1 - x^{2}}$

The area of the shaded region is the equal to

$$A = \int_0^1 \sqrt{1 - x^2} \, dx$$

We set $x = \sin \theta$, $dx = \cos \theta d\theta$

$$A = \int_{0}^{\pi/2} \sqrt{1 - x^2} dx$$

$$= \int_{0}^{\pi/2} \sqrt{(1 - \sin^2 \theta)} \cos \theta \, d\theta$$

$$= \int_{0}^{\pi/2} \sqrt{(\cos^2 \theta)} \cos \theta \, d\theta$$

$$= \int_{0}^{\pi/2} \cos^2 \theta \, d\theta$$
By using the following formula
$$\cos^2 \theta = \frac{1 + \cos 2\theta}{2},$$
we have

$$\cos^2\theta = \frac{1+\cos 2\theta}{2}$$

we have

$$A = \int_{0}^{\pi/2} \frac{1 + \cos 2\theta}{2} d\theta$$
$$= \int_{0}^{\pi/2} \left(\frac{1}{2} + \frac{\cos 2\theta}{2}\right) d\theta$$
$$= \left[\frac{1}{2}\theta + \frac{\sin 2\theta}{4}\right]_{0}^{\pi/2}$$
$$= \left(\frac{\pi}{4} + 0\right) - (0 + 0)$$

$$=\frac{\pi}{4}$$

The following are some more examples of exact integration. You can use the brief table of integrals given in Table 1.

TableA brief table of integrals $\int dx = x + C$	$\int \sin x dx = -\cos x + C$
$\int a f(x) dx = a \int f(x) dx + C$	$\int \cos x dx = \sin x + C$
$\int \left[u(x) \pm v(x) \right] dx = \int u(x) dx \pm \int v(x) dx + C$	$\int \tan x dx = -\ln \cos x + C = \ln \sec x + C$
$\int x^n dx = \frac{x^{n+1}}{n+1} + C$	$\int \sec(ax)dx = \frac{1}{a}\ln\left \sec(ax) + \tan(ax)\right + C$
$\int u dv = u v - \int v du + C$	$\int \cot x dx = -\ln \csc x + C = \ln \sin x + C$
$\int \frac{dx}{ax+b} = \frac{1}{a} \ln \left ax+b \right + C$	$\int \sec^2 ax dx = \frac{1}{a} \tan(ax) + C$
$\int a^x dx = \frac{a^x}{\ln a} + C$	$\int \sec(x) \tan(x) dx = \sec(x) + C$
$\int e^{ax} dx = \frac{e^{ax}}{a} + C$	$\int \csc(x)\cot(x)dx = -\csc(x) + C$
तजारव ना	वधातमस्तु

Example 5

Evaluate the following integral

 $\int_{0}^{1} 2x e^{-x^2} dx$

Solution

Let
$$u = -x^{2}$$
, $du = -2xdx$
At $x = 0$, $u = -(0)^{2} = 0$
At $x = 1$, $u = -(1)^{2} = -1$
 $\int_{0}^{1} 2xe^{-x^{2}}dx = \int_{0}^{1} (-e^{-x^{2}})(-2xdx)$
 $= \int_{0}^{-1} (-e^{u})(du)$
 $= [-e^{u}]_{0}^{1}$
 $= -e^{-1} - (-e^{0})$
 $= 0.6321$
Example 6
Evaluate
 $\int_{0}^{\pi/4} \frac{1 + \sin x}{\cos^{2} x} dx$
Solution
 $\int_{0}^{\pi/4} \frac{1 + \sin x}{\cos^{2} x} dx = \int_{0}^{\pi/4} \left(\frac{1}{(\cos^{2} x} + \frac{\sin x}{\cos^{2} x}) dx - \frac{1}{(\cos^{2} x)(\tan x)} dx \right)$

COPYRIGHT FIMT 2020

164 | Page

$$= [\tan x]_0^{\pi/4} + [\sec x]_0^{\pi/4}$$
$$= (1-0) + (\sqrt{2}-1)$$
$$= \sqrt{2}$$

MANAR

Example 7 Evaluate $\int x \sec^2 x \, dx$

Solution

We use the formula

 $\int u dv = uv - \int v du$

Let u = x, du = dx, and $dv = \sec^2 x \, dx$, $v = \tan x$

So the new integral is

$$\int x \sec^2 x \, dx = x \tan x - \int \tan x \, dx$$

 $= x \tan x + \ln \left| \cos x \right| + C$

Example 8

Evaluate

Solution

Let
$$u = \ln x$$
, $du = \frac{1}{x} dx$ and $dv = x dx$, $v = \frac{x^2}{2}$

Using the formula $\int u dv = uv - \int v du$, the new integral is

$$\int_{1}^{2} (x)(\ln x) dx = \left[\ln x \times \frac{x^2}{2} \right]_{1}^{2} - \int_{1}^{2} \left(\frac{x^2}{2} \right) \left(\frac{1}{x} dx \right)$$

$$= \left[\ln x \times \frac{x^2}{2} \right]_{1}^{2} - \int_{1}^{2} \frac{x}{2} dx$$

$$= \left[\ln x \times \frac{x^2}{2} \right]_{1}^{2} - \left[\frac{x^2}{4} \right]_{1}^{2}$$

$$= \left[\left(\ln 2 \times \frac{2^2}{2} \right) - \left(\ln 1 \times \frac{1^2}{2} \right) \right] - \left[\left(\frac{2^2}{4} \right) - \left(\frac{1^2}{4} \right) \right]$$

$$= \left[(2\ln 2) - \left(\frac{1}{2} \ln 1 \right) \right] - \left[\left(\frac{4}{4} \right) - \left(\frac{1}{4} \right) \right]$$

$$= \left[(2\ln 2) - \left(\frac{1}{2} \times 0 \right) \right] - \left[1 - \frac{1}{4} \right]$$

$$= 0.6362$$

Example 9

Evaluate

$$\int_{0}^{1} \frac{5x}{(4+x^{2})^{2}} dx$$

Solution

We use the formula $\int_{a}^{b} f(g(x))g'(x)dx = \int_{g(a)}^{g(b)} f(u)du$, by substituting u = g(x), du = g'(x)dx

नावध

then integrating from g(a) to g(b).

Let

$$u = g(x) = 4 + x^2,$$

SO

$$g(0) = 4, g(1) = 5$$
, and

$$du = (2x)dx$$

The new integral is

 $\int_{0}^{1} \frac{5x}{(4+x^{2})^{2}} dx = \int_{0}^{1} \frac{1}{(4+x^{2})^{2}} \times \frac{5}{2} \times (2x) dx$ TED $=\frac{5}{2}\int_{4}^{5}\frac{1}{u^{2}}du$ $=\frac{5}{2}\left[-\frac{1}{u}\right]_{4}^{5}$ $=\frac{5}{2}\left[(-\frac{1}{5})-(-\frac{1}{4})\right]$ $=\frac{5}{2}\times\frac{1}{20}$ = 0.125 Example 10 Evaluate $\int_{0}^{4} |2x-1| dx$ व नाव Solution First, let's analyze the expression |2x-1|. $|2x-1| = -(2x-1), x < \frac{1}{2}$ $=(2x-1), x \ge \frac{1}{2}$

$$\int_{0}^{4} |2x - 1| dx = \int_{0}^{1/2} -(2x - 1) dx + \int_{1/2}^{4} (2x - 1) dx$$
$$= \left[-x^{2} + x \right]_{0}^{1/2} + \left[x^{2} - x \right]_{1/2}^{4}$$
$$= \left[\left(-\frac{1}{4} + \frac{1}{2} \right) - 0 \right] + \left[(16 - 4) - \left(\frac{1}{4} - \frac{1}{2} \right) \right]$$
$$= 12.5$$

Example 11

Evaluate

$$\int_{-\infty}^{-2} \frac{2}{x^2 - 1} dx$$

Solution

$$\int_{-\infty}^{2} \frac{2}{x^{2}-1} dx = \int_{-\infty}^{2} \frac{2}{(x-1)\times(x+1)} dx$$

$$= \int_{-\infty}^{2} \frac{(x+1)-(x-1)}{(x-1)\times(x+1)} dx$$

$$= \int_{-\infty}^{2} \frac{x+1}{(x-1)\times(x+1)} - \frac{x-1}{(x-1)\times(x+1)} dx$$

$$= \int_{-\infty}^{2} \frac{1}{x-1} dx - \int_{-\infty}^{2} \frac{1}{x+1} dx$$

$$= \lim_{b \to \infty} \left[\ln |x-1| \right]_{b}^{b} - \lim_{b \to -\infty} \left[\ln |x+1| \right]_{b}^{b}$$

$$= \lim_{b \to \infty} \left[\ln \left| \frac{x-1}{x+1} \right| \right]_{-2}^{b}$$

$$= \ln (3) - \ln \left(\lim_{b \to -\infty} \left| \frac{b-1}{b+1} \right| \right)$$
$$= \ln (3) - \ln (1)$$
$$= \ln (3)$$
$$= 1.0986$$

Example 12

Graph the function $y = \frac{1}{3}(x^2 + 2)^{3/2}$, and find the length of the curve from x = 0 to x = 3.

Solution

We use the equation

$$L = \int_{a}^{b} \sqrt{1 + \left(\frac{dy}{dx}\right)^2} \, dx$$

We have:

$$y = \frac{1}{3}(x^2 + 2)^{3/2}$$

So,

$$\frac{dy}{dx} = \left(\frac{1}{3}\right) \times \left(\frac{3}{2}\right) \times (x^2 + 2)^{3/2 - 1} \times (2x)$$
$$= x\sqrt{x^2 + 2}$$
$$L = \int_{0}^{3} \sqrt{1 + \left(x\sqrt{x^2 + 2}\right)^2} dx$$



ARGE

Figure 7 Graph of the function $y = \frac{1}{3}(x^2 + 2)^{3/2}$

$$= \int_{0}^{3} \sqrt{1 + x^{2}(x^{2} + 2)} dx$$

= $\int_{0}^{3} \sqrt{1 + x^{4} + 2x^{2}} dx$
= $\int_{0}^{3} \sqrt{(x^{2} + 1)^{2}} dx$
= $\int_{0}^{3} (x^{2} + 1) dx$
= $\left[\frac{x^{3}}{3} + x\right]_{0}^{3}$
= 12

Example 13

Find the area of the shaded region given in Figure 8.





Figure 8 Graph of the function $\cos^2 x$.

Solution

For the sketch given,

$$a=rac{\pi}{2}, b=\pi$$
 , and

$$f(x) - g(x) = 1 - \cos^2 x = \sin^2 x$$

$$A = \int_{\pi/2}^{\pi} \sin^{2}(x) dx$$

= $\int_{\pi/2}^{\pi} \frac{1 - \cos 2x}{2} dx$
= $\int_{\pi/2}^{\pi} \left[\frac{1}{2} - \frac{\cos 2x}{2} \right] dx$

$$= \left\lfloor \frac{x}{2} - \frac{\sin 2x}{4} \right\rfloor_{\pi/2}^{\pi}$$

$$= \left[\left(\frac{\pi}{2} - \frac{\sin(2\pi)}{4} \right) - \left(\frac{\pi}{4} - \frac{\sin 2\left(\frac{\pi}{2}\right)}{4} \right) \right]$$
$$= \left[\left(\frac{\pi}{2} - 0 \right) - \left(\frac{\pi}{4} - 0 \right) \right]$$
$$= \frac{\pi}{4}$$

Example 14

Find the volume of the solid generated by revolving the shaded region in Figure 9 about the *y*-axis.





Solution

We use the formula $V = \int_{a}^{b} \pi (radius)^{2} dy$

Let

$$u = \frac{\pi}{4} y, \ du = \frac{\pi}{4} dy.$$

Therefore, at y = 0, u = 0

$$y = 1, u = \frac{\pi}{4}$$

$$V = \int_{0}^{1} \pi [R(y)]^{2} dy$$

$$= \pi \int_{0}^{1} \left[\tan\left(\frac{\pi}{4}y\right) \right]^{2} dy$$

$$= \pi \times \frac{4}{\pi} \int_{0}^{1} \left[\tan\left(\frac{\pi}{4}y\right) \right]^{2} \frac{\pi}{4} dy$$

$$= 4 \int_{0}^{\pi/4} (\tan u)^{2} du \quad (\text{Choosing } u = \frac{\pi}{4}y)$$

$$= 4 \int_{0}^{\pi/4} (-1 + \sec^{2} u) du$$

$$= 4 \left[-u + \tan u \right]_{0}^{\pi/4}$$

$$= 4 \left[\left(-\frac{\pi}{4} + \tan \frac{\pi}{4} \right) - (0 + \tan 0) \right]$$

$$= 4 \left[\left(-\frac{\pi}{4} + 1 \right) - (0 + 0) \right]$$

$$= 0.8584$$

Partial Fractions:-

Partial Fractions provides a way to integrate all rational functions.

Rational functions= $\frac{P(x)}{Q(x)}$ when P and Q are polynomials

This is the technique to find $\int \frac{P(x)}{Q(x)} dx$

Rule 1: The degree of the numerator must be less than the degree of the denominator. If this is not the case we first must divide the numerator into the denominator.

Step 1: If Q has a quadratic factor $ax^2 + bx + c$ which corresponds to a complex root of order

k, then the partial fraction expansion of $\frac{P}{Q}$ contains a term of the form

 $\frac{B_1 x + C_1}{(ax^2 + bx + c)} + \frac{B_2 x + C_2}{(ax^2 + bx + c)^2} + \dots + \frac{B_k + C_k}{(ax^2 + bx + c)^k}$

Where $B_1, B_2, ..., B_k$ and $C_1, C_2, ..., C_k$ are unknown constants.

Step 2: Set the sum of the terms of equal to the partial fraction expansion

Example:
$$\frac{1}{(x-2)(x-5)} = \frac{A}{x-2} + \frac{B}{x-5}$$

Step 3: When then multiply both sides by Q to get some expression that is equal to P

Example: 1 = A(x-5) + B(x-2)

1= (A+B)x-5A-2B

Step 4: Use the theory that 2 polynomials are equal if and only if the corresponding coefficients are equal

Example: 5A-2B=1 and A+B=0

Step 5: Solve for A, B, and C

Example: A= -1/3 B= 1/3

Step 6: Express integral of $\frac{P}{Q}$ as the sum of the integrals of the terms of partial fraction expansion.

Example:
$$\int \frac{1}{(x-2)(x-5)} dx = \int \frac{\frac{-1}{3}}{(x-2)} dx + \int \frac{\frac{1}{3}}{(x-5)} dx$$
$$= \frac{-1}{3} \ln|x-2| + \frac{1}{3} \ln|x-5| + C$$

AMAGEMEN

Example 2:

Find
$$\int \frac{x^4 - 2x^2 + 4x + 1}{x^3 - x^2 - x + 1} dx$$

 $\frac{x^4 - 2x^2 + 4x + 1}{x^3 - x^2 - x + 1} = x + 1 + \frac{4x}{x^3 - x^2 - x + 1}$ Note: long division

$$\frac{4x}{(x-1)^2(x+1)}$$
 Note: Factor Q(x)= x³ - x² - x +1

 $\frac{A}{(x-1)} + \frac{B}{(x-1)^2} + \frac{C}{x+1}$ Note: Partial fraction decomposition since (x-1)²'s factor is linear.

There is a constant on top for the and power and first power

 $4x = A(x-1)(x+1) + B(x+1) + C(x-1)^{2}$

Note: multiply by Least common denominator

-7015

(x-1)² (x+1)

$$= (A+C)x^{2} + (B-2C)x+(-A+B+C)$$
A+C = 0
B-2C = 4
A+B+C = 0
Note: Equate equations
A=1 B=2 C=-1
Note: Solve for coefficients
$$\int (x+1)dx + \int \frac{1}{x-1}dx + \int \frac{2}{(x-1)^{2}}dx - \int \frac{1}{x+1}dx$$

$$= \frac{x^{2}}{2} + x + \ln|x-1| - \frac{2}{x-1} - \ln|x+1| + C$$

$$= \frac{x^{2}}{2} + x - \frac{2}{x-1} + \ln\left|\frac{x-1}{x+1}\right| + C$$
Example 3:
$$Ion \int \frac{2x^{2} - x + 4}{x^{2} + 4x} dx$$

$$= \frac{2x^{2} - x + 4}{x^{2} + 4x} = \frac{A}{x} + \frac{Bx + C}{B^{2} + 4}$$
Note: x²+4 is quadratic
$$2x^{2} - x + 4 = A(x^{2}+4) + (Bx+C)x$$
Note: multiplying x(x²+4)
$$= (A+B)x^{2} - Cx + A$$

$$A+B=2 C=-1 \quad AA=4$$
Note: Equating coefficients
$$A=1 \quad B=1 \quad C=-1$$

$$\int \frac{2x^2 - x + 4}{x^3 + 4x} dx = \int \frac{1}{x} dx + \int \frac{x - 1}{x^2 + 4} dx = \int \frac{1}{x} dx + \int \frac{x}{x^2 + 4} dx - \int \frac{1}{x^2 + 4} dx$$
$$= \ln|x| + \frac{1}{2} \ln|x^2 + 4| - \frac{1}{2} \tan^{-1}\left(\frac{x}{2}\right) + C$$

INTEGRATION BY PART

This is a method to evaluate integrals that cannot be evaluated by eye or by u-substitution. It is usually applied to expressions with varied functions within each other, or multiplied by each other. A good rule is: if the expression has a chain of functions (f(g(x))) or if the expression has a product of functions x(f(x)), integration by parts will be necessary. Here are some examples of problems that would be solved with integration by parts:

 $x^3 \ln(x)$

 $x \cdot \sec^2(x)$

Let's start with:

 $x^3 \ln(x)$

In integration by parts, you separate the expression into two parts: u, and $\partial(v)$. The u should be easy to differentiate, and the d(v) should be easy to integrate. Once you have chosen a u and a d(v), set up a chart like this:

u = ln(x)
$$d(v) = \frac{\ln(x) \cdot x^4}{4} - \frac{x^4}{16}$$

COPYRIGHT FIMT 2020

 $\theta^2 \sin \theta$

$$d(u) = x^{-1}$$
 $v = \frac{x^4}{4}$

Now, the formula to solve this is:

$$uv - \int (d(u)v)$$

so here, the equation to solve is:

$$\frac{\ln(x)\cdot x^4}{4} - \int \left(x^{-1}\cdot \frac{x^4}{4}\right) d(x)$$

which simplifies to:

$$\frac{\ln(x)\cdot x^4}{4} - \int \left(\frac{x^3}{4}\right) d(x)$$

and solve the integral to get:

$$\frac{\ln(x)\cdot x^4}{4} - \frac{x^4}{16}$$

Unfortunately, it is not always so simple. Sometimes, you must use u-substitution, or even integration by parts again within the solution. Take, for example:

$$\theta^2 \sin \theta$$

To solve this, you would set up a chart again.

$$\theta^2 \sin \theta$$

$$u = \theta^2 \qquad \qquad d(v) = \sin \theta$$

$$d(u) = 2\theta \qquad \qquad \mathsf{v} = -\cos\theta$$

With this chart, you can set up the solution using the $uv - \int (d(u)v)$ formula:

$$= -\theta^2 \cos \theta + \int (2\theta \cos \theta) d(\theta)$$

But the second part of this, $\int (2\theta \cos \theta) d(\theta)$ cannot be solved by eye. You must set up a second chart:

$$\int (2\theta \cos \theta) d(\theta)$$

u = 2 θ $d(v) = \cos \theta$

 $v = \sin \theta$

This gives us:

$$= 2\theta\sin\theta + \int (2\sin\theta)d(\theta)$$

d(u

Which can be simplified to:

$$=2\theta\sin\theta+2\cos\theta$$

Now, you can substitute it into the original solution, in the place of $\int (2\theta \cos \theta) d(\theta)$, giving you:

$$= -\theta^2 \cos\theta + 2\theta \sin\theta + 2\cos\theta$$

Trigonometric Integrals(REDUCTION FORMULA)

- I. Integrating Powers of the Sine and Cosine Functions
 - A. Useful trigonometric identities
 - 1. $\sin^2 x + \cos^2 x = 1$
 - 2. $\sin 2x = 2\sin x \cos x$

3.
$$\cos 2x = \cos^2 x - \sin^2 x = 2\cos^2 x - 1 = 1 - 2\sin^2 x$$

4. $\sin^2 x = \frac{1 - \cos 2x}{2}$
5. $\cos^2 x = \frac{1 + \cos 2x}{2}$
6. $\sin x \cos y = \frac{1}{2} [\sin(x - y) + \sin(x + y)]$
7. $\sin x \sin y = \frac{1}{2} [\cos(x - y) - \cos(x + y)]$
8. $\cos x \cos y = \frac{1}{2} [\cos(x - y) + \cos(x + y)]$

B. Reduction formulas

24

1.
$$\int \sin^n x \, dx = -\frac{1}{n} \sin^{n-1} x \cos x + \frac{n-1}{n} \int \sin^{n-2} x \, dx$$

2.
$$\int \cos^{n} x \, dx = \frac{1}{n} \cos^{n-1} x \sin x + \frac{n-1}{n} \int \cos^{n-2} x \, dx$$

C. Examples
1. Find
$$\int \sin^{2} x \, dx$$
.

Method 1(Integration by parts): $\int \sin^2 x \, dx = \int \sin x \, (\sin x \, dx)$. Let
$$u = \sin x$$
 and $dv = \sin x \, dx \Rightarrow du = \cos x \, dx$ and $v = \int \sin x \, dx = \int \sin x \, dx$

1

$$-\cos x. \text{ Thus, } \int \sin^2 x \, dx = (\sin x)(-\cos x) + \int \cos^2 x \, dx = -\sin x \cos x + \int (1 - \sin^2 x) \, dx = -\sin x \cos x + \int 1 \, dx - \int \sin^2 x \, dx = -\sin x \cos x + x - \int \sin^2 x \, dx = 2 \int \sin^2 x \, dx = -\sin x \cos x + x \Rightarrow \int \sin^2 x \, dx = -\frac{1}{2} \sin x \cos x + \frac{1}{2} x + C.$$

Method 2(Trig identity): $\int \sin^2 x \, dx = \frac{1}{2} \int (1 - \cos 2x) \, dx = \frac{1}{2} x - \frac{1}{4} \sin 2x + C$.

-

Method 3(Reduction formula): $\int \sin^2 x \, dx = -\frac{1}{2} \sin x \cos x + \frac{1}{2} \int 1 dx =$

$$-\frac{1}{2}\sin x\cos x + \frac{1}{2}x + C$$

2. Find $\cos^3 x \, dx$.

Use the reduction formula: $\int \cos^3 x \, dx = \frac{1}{3} \cos^2 x \sin x + \frac{2}{3} \int \cos x \, dx =$

$$\frac{1}{3}\cos^2 x \sin x + \frac{2}{3}\sin x + C = \frac{1}{3}\sin x(1 - \sin^2 x) + \frac{2}{3}\sin x + C =$$

$$\sin x - \frac{1}{3}\sin^3 x + C.$$

3. Find $\int \sin^3 x \cos^2 x \, dx$.

$$\int \sin^3 x \cos^2 x \, dx = \int \sin^2 x \sin x \cos^2 x \, dx = \int (1 - \cos^2 x) \cos^2 x \sin x dx =$$
$$\int (\cos^2 x - \cos^4 x) (\sin x \, dx). \text{ Let } u = \cos x \Longrightarrow du = -\sin x \, dx. \text{ Thus,}$$

1

$$\int (\cos^2 x - \cos^4 x)(\sin x \, dx) = -\int (u^2 - u^4) \, du = -\frac{1}{3}u^3 + \frac{1}{5}u^5 + C =$$

ARGEMEN

 $-\frac{1}{3}\cos^3 x + \frac{1}{5}\cos^5 x + C.$

2

4. Find
$$\int \sin^2 x \cos^2 x \, dx$$
.

$$\int \sin^2 x \cos^2 x \, dx = \int \left(\frac{1-\cos 2x}{2}\right) \left(\frac{1+\cos 2x}{2}\right) \, dx = \frac{1}{4} \int (1-\cos^2 2x) \, dx =$$
$$\frac{1}{4} \int \sin^2 2x \, dx = \frac{1}{4} \int \left(\frac{1-\cos 4x}{2}\right) \, dx = \frac{1}{8} \int 1 \, dx - \frac{1}{8} \int \cos 4x \, dx =$$
$$\frac{1}{8} x - \frac{1}{32} \sin 4x + C \, .$$

5. Find $\int \sin 4x \cos 3x \, dx$.

Method 1(Integration by parts): Let $u = \sin 4x$ and $dv = \cos 3x \, dx \Rightarrow du =$

$$4\cos 4x \ dx \ and \ v = \frac{1}{3}\sin 3x \ . \text{Thus}, \int \sin 4x\cos 3x \ dx =$$

$$(\sin 4x) \left(\frac{1}{3}\sin 3x\right) - \frac{4}{3} \int \cos 4x \sin 3x \ dx = \frac{1}{3}\sin 4x \sin 3x -$$

$$\frac{4}{3} \int \cos 4x \sin 3x \ dx \ . \text{Find} \ \int \cos 4x \sin 3x \ dx \ . \text{Let} \ u = \cos 4x \ \text{and} \ dv =$$

$$\sin 3x \ dx \Rightarrow du = -4\sin 4x \ dx \ \text{and} \ v = -\frac{1}{3}\cos 3x \ . \text{Thus},$$

$$\int \cos 4x \sin 3x \ dx = -\frac{1}{3}\cos 4x \cos 3x - \frac{4}{3} \int \sin 4x \cos 3x \ dx \ . \text{Returning to}$$

$$\text{the original integral,} \ \int \sin 4x \cos 3x \ dx = \frac{1}{3}\sin 4x \sin 3x -$$

$$\frac{4}{3} \left\{ -\frac{1}{3}\cos 4x \cos 3x - \frac{4}{3} \int \sin 4x \cos 3x \ dx \right\} = \frac{1}{3}\sin 4x \sin 3x +$$

$$\frac{4}{9}\cos 4x \cos 3x + \frac{16}{9} \int \sin 4x \cos 3x \ dx \Rightarrow -\frac{7}{9} \int \sin 4x \cos 3x \ dx =$$

$$\frac{1}{3}\sin 4x \sin 3x + \frac{4}{9}\cos 4x \cos 3x \Rightarrow \int \sin 4x \cos 3x \ dx =$$

$$-\frac{3}{7}\sin 4x \sin 3x - \frac{4}{7}\cos 4x \cos 3x \ dx = \frac{1}{2} \int (\sin x + \sin 7x) \ dx =$$

$$-\frac{1}{2}\cos x - \frac{1}{14}\cos 7x + C .$$

3

II. Integrating Powers of the Tangent and Secant Functions

A. Useful trigonometric identity: $\tan^2 x + 1 = \sec^2 x$

- B. Useful integrals
 - 1. $\int \sec x \tan x \, dx = \sec x + C$
 - 2. $\int \sec^2 x \, dx = \tan x + C$

3.
$$\int \tan x \, dx = \ln |\sec x| + C = -\ln |\cos x| + C$$

4.
$$\int \sec x \, dx = \ln \left| \sec x + \tan x \right| + C$$

C. Reduction formulas

1.
$$\int \sec^{n} x \, dx = \frac{\sec^{n-2} x \tan x}{n-1} + \frac{n-2}{n-1} \int \sec^{n-2} x \, dx$$

2.
$$\int \tan^n x \, dx = \frac{\tan^{n-1} x}{n-1} - \int \tan^{n-2} x \, dx$$

D. Examples
1. Find
$$\int \tan^2 x \, dx$$
.
 $\int \tan^2 x \, dx = \int (\sec^2 x - 1) \, dx = \int \sec^2 x \, dx - \int 1 \, dx = \tan x - x + C$.
2. Find $\int \tan^3 x \, dx$.

$$\int \tan^3 x dx = \frac{\tan^2 x}{2} - \int \tan x dx = \frac{1}{2} \tan^2 x - \ln|\sec x| + C.$$
3. Find $\int \sec^3 x dx.$

$$\int \sec^3 x dx = \frac{\sec x \tan x}{2} + \frac{1}{2} \int \sec x dx = \frac{1}{2} \sec x \tan x + \frac{1}{2} \ln|\sec x + \tan x| + C.$$
4. Find $\int \tan x \sec^2 x dx.$
Let $u = \tan x \Rightarrow du = \sec^2 x dx \Rightarrow \int \tan x \sec^2 x dx = \int u du = \frac{1}{2} u^2 + C = \frac{1}{2} \tan^2 x + C.$
5. Find $\int \tan x \sec^4 x dx.$

$$\int \tan x \sec^4 x dx = \int \tan x \sec^2 x dx = \int \tan x (1 + \tan^2 x) \sec^2 x dx = \int \tan x \sec^2 x dx = \int \tan x \sec^2 x dx.$$

$$\int \tan x \sec^4 x dx = \int \tan^2 x \sec^2 x dx = \int \tan x (1 + \tan^2 x) \sec^2 x dx = \int \tan x \sec^2 x dx = \int \tan^2 x + \frac{1}{4} \tan^2 x + C.$$

6. Find
$$\int \tan x \sec^3 x \, dx$$
.

$$\int \tan x \sec^3 x \, dx = \int \sec^2 x (\sec x \tan x \, dx). \text{ Let } u = \sec x \Longrightarrow du = \sec x \tan x \, dx.$$

- Thus, $\int \tan x \sec^3 x \, dx = \int u^2 \, du = \frac{1}{3}u^3 + C = \frac{1}{3}\sec^3 x + C$.
- 7. Find $\int \tan^2 x \sec^3 x \, dx$.

$$\int \tan^2 x \sec^3 x \, dx = \int (\sec^2 x - 1) \sec^3 x \, dx = \int \sec^5 x \, dx - \int \sec^3 x \, dx.$$
 Using
the reduction formula, $\int \sec^5 x \, dx = \frac{1}{2} \sec^3 \tan x + \frac{3}{2} \int \sec^3 x \, dx.$ Thus,

$$\int \tan^2 x \sec^3 x \, dx = \int \sec^5 x \, dx - \int \sec^3 x \, dx = \frac{1}{4} \sec^3 x \tan x + \frac{3}{4} \int \sec^3 x \, dx - \frac{1}{4} \sec^3 x \tan x + \frac{3}{4} \int \sec^3 x \, dx = \frac{1}{4} \sec^3 x \tan^2 x + \frac{3}{4} \int \sec^3 x \, dx = \frac{1}{4} \sec^3 x \tan^2 x + \frac{3}{4} \int \sec^3 x \, dx = \frac{1}{4} \sec^3 x \tan^2 x + \frac{3}{4} \int \sec^3 x \, dx = \frac{1}{4} \sec^3 x \tan^2 x + \frac{3}{4} \int \sec^3 x \, dx = \frac{1}{4} \sec^3 x \tan^2 x + \frac{3}{4} \int \sec^3 x \, dx = \frac{1}{4} \sec^3 x \tan^2 x + \frac{3}{4} \int \sec^3 x \, dx = \frac{1}{4} \sec^3 x \, dx = \frac{1}{4} \sec^3 x \tan^2 x + \frac{3}{4} \int \sec^3 x \, dx = \frac{1}{4} \sec^3 x \tan^2 x + \frac{3}{4} \int \sec^3 x \, dx = \frac{1}{4} \sec^3 x \, dx = \frac{1}{4} \sec^3 x \tan^2 x + \frac{3}{4} \int \sec^3 x \, dx = \frac{1}{4} \sec^3$$

$$\int \sec^3 x \, dx = \frac{1}{4} \sec^3 x \tan x - \frac{1}{4} \int \sec^3 x \, dx = \frac{1}{4} \sec^3 x \tan x - \frac{1}{8} \sec x \tan x + \frac{1}{8} \sec^3 x \tan^2 x + \frac{1}{8} \sec^2 x \tan^2 x + \frac{1}{8} \sec^$$

 $\frac{1}{8}\ln|\sec x + \tan x| + C.$

8. Find $\sqrt{\tan x} \sec^4 x \, dx$.

$$\sqrt{\tan x} \sec^4 x \, dx = \int \sqrt{\tan x} \sec^2 x \sec^2 x \, dx = \int \sqrt{\tan x} (1 + \tan^2 x) \sec^2 x \, dx.$$

Let
$$u = \tan x \Rightarrow du = \sec^2 x \, dx \Rightarrow \int \sqrt{\tan x} \sec^4 x \, dx = \int \sqrt{\tan x} \sec^2 x \, dx + \int \sqrt{\tan x} \tan^2 x \sec^2 x \, dx = \int u^{\frac{1}{2}} du + \int u^{\frac{5}{2}} du = \frac{2}{3} u^{\frac{3}{2}} + \frac{2}{7} u^{\frac{7}{2}} + C = \frac{2}{3} (\tan x)^{\frac{3}{2}} + \frac{2}{7} (\tan x)^{\frac{7}{2}} + C.$$

9. Find $\int \sqrt{\sec x} \tan x \, dx$.

Let
$$u = \sqrt{\sec x} \Rightarrow u^2 = \sec x \Rightarrow 2udu = \sec x \tan x dx = u^2 \tan x dx \Rightarrow$$

 $\tan x dx = \frac{2udu}{u^2} = \frac{2}{u} du$. Thus, $\int \sqrt{\sec x} \tan x dx = \int u \left(\frac{2}{u} du\right) = 2 \int 1 du =$
 $2u + C = 2\sqrt{\sec x} + C$.

Practice Sheet for Trigonometric Integrals

(1) Prove the reduction formula:
$$\int \sin^n x \, dx = -\frac{1}{n} \sin^{n-1} x \cos x + \frac{n-1}{n} \int \sin^{n-2} x \, dx$$

E.A.

(2) Prove the reduction formula: $\int \cos^n x \, dx = \frac{1}{n} \cos^{n-1} x \sin x + \frac{n-1}{n} \int \cos^{n-2} x \, dx$

(3) Prove the reduction formula:
$$\int \sec^n x \, dx = \frac{\sec^{n-2} x \tan x}{n-1} + \frac{n-2}{n-1} \int \sec^{n-2} x \, dx$$

(4) Prove the reduction formula:
$$\int \tan^n x \, dx = \frac{\tan^{n-1} x}{n-1} - \int \tan^{n-2} x \, dx$$

(s)
$$\int_{0}^{\pi/2} \cos^{2}(2x) dx =$$
(c)
$$\int_{0}^{\pi/2} \cos^{2}(2x) dx =$$
(d)
$$\int_{0}^{\pi/2} \sin(5x)\cos(3x) dx =$$
(e)
$$\int_{0}^{\pi/2} \sin^{3}x \sec^{3}x dx =$$
(f)
$$\int_{0}^{\sqrt{\sin x} \cos^{3}x} dx =$$
(g)
$$\int_{0}^{\sqrt{\sin x} \cos^{3}x} dx =$$
(h)
$$\int_{0}^{\sqrt{\sin x} \cos^{3}x} dx =$$
(h)
$$\int_{0}^{\pi/2} \cos^{3}x \sin^{2}x dx =$$
(h)
$$\int_{0}^{\pi/2} \sin^{3}x dx =$$
(h)
$$\int_{0}^{\pi/2} \sin^{3}x dx =$$

$$(12) \int \sin^2 x \cos^2 x dx =$$

(13)
$$\int \tan^5 x \sec x \, dx =$$

NAAC ACCREDITED

Solution Key for Trigonometric Integrals

(1) $\int \sin^n x \, dx = \int \sin^{n-1} x \sin x \, dx$. Use integration by parts with $u = \sin^{n-1} x$ and

 $dv = \sin x \, dx \Rightarrow du = (n-1)\sin^{n-2} x \cos x \, dx$ and $v = \int \sin x \, dx = -\cos x \Rightarrow$

$$\int \sin^{n} x \, dx = \int \sin^{n-1} x \sin x \, dx = -\sin^{n-1} x \cos x + (n-1) \int \sin^{n-2} x \cos^{2} x \, dx =$$
$$-\sin^{n-1} x \cos x + (n-1) \int \sin^{n-2} x \left(1 - \sin^{2} x\right) dx = -\sin^{n-1} x \cos x +$$

RELE

$$(n-1)\int \sin^{n-2}x\,dx - (n-1)\int \sin^n x\,dx \Longrightarrow n\int \sin^n x\,dx = -\sin^{n-1}x\cos x +$$

$$(n-1)\int \sin^{n-2} x \, dx \Longrightarrow \int \sin^n x \, dx = -\frac{1}{n} \sin^{n-1} x \cos x + \frac{n-1}{n} \int \sin^{n-2} x \, dx \, .$$

(2) $\int \cos^n x \, dx = \int \cos^{n-1} x \cos x \, dx$. Use integration by parts with $u = \cos^{n-1} x$ and

 $dv = \cos x \, dx \Longrightarrow du = (n-1)\cos^{n-2} x (-\sin x) \, dx$ and $v = \int \cos x \, dx = \sin x \Longrightarrow$

 $\int \cos^n x \, dx = \int \cos^{n-1} x \cos x \, dx = \cos^{n-1} x \sin x + (n-1) \int \cos^{n-2} x \sin^2 x \, dx =$

 $\cos^{n-1} x \sin x + (n-1) \int \cos^{n-2} x (1 - \cos^2 x) dx = \cos^{n-1} x \sin x + (n-1) \int \cos^{n-2} x (1 - \cos^2 x) dx = \cos^{n-1} x \sin x + (n-1) \int \cos^{n-2} x (1 - \cos^2 x) dx = \cos^{n-1} x \sin^{n-1} x \cos$

$$(n-1)\int \cos^{n-2}x\,dx - (n-1)\int \cos^n x\,dx \Rightarrow n\int \cos^n x\,dx = \cos^{n-1}x\sin x + \frac{1}{2}\sin^n x\,dx$$

$$(n-1)\int \cos^{n-2} x \, dx \Rightarrow \int \cos^n x \, dx = \frac{1}{n} \cos^{n-1} x \sin x + \frac{n-1}{n} \int \cos^{n-2} x \, dx$$

400

(3) $\int \sec^n x \, dx = \int \sec^{n-2} x \sec^2 x \, dx$. Use integration by parts with $u = \sec^{n-2} x$ and

$$dv = \sec^2 x \, dx \Rightarrow du = (n-2) \sec^{n-3} x (\sec x \tan x \, dx) \text{ and } v = \int \sec^2 x \, dx = \tan x \Rightarrow$$

$$\int \sec^{n} x \, dx = \int \sec^{n-2} x \sec^2 x \, dx = \sec^{n-2} x \tan x - (n-2) \int \sec^{n-2} x \tan^2 x \, dx =$$
$$\sec^{n-2} x \tan x - (n-2) \int \sec^{n-2} x \left(\sec^2 x - 1\right) dx = \sec^{n-2} x \tan x - (n-2) \int \sec^n x \, dx =$$

$$(n-2)\int \sec^{n-2} x \, dx \Longrightarrow (n-1)\int \sec^n x \, dx = \sec^{n-2} x \tan x + (n-2)\int \sec^{n-2} x \, dx \Longrightarrow$$

$$\int \sec^{n} x \, dx = \frac{\sec^{n-2} x \tan x}{n-1} + \frac{n-2}{n-1} \int \sec^{n-2} x \, dx \, .$$

(4)
$$\int \tan^n x \, dx = \int \tan^{n-2} x \tan^2 x \, dx = \int \tan^{n-2} x \left(\sec^2 x - 1 \right) dx = \int \tan^{n-2} x \sec^2 x \, dx - \frac{1}{2} \int \tan^{n-2} x \sec^2 x \, dx = \int \tan^{n-2} x \, dx = \int \tan^{n$$

$$\int \tan^{n-2} x \, dx = \frac{\tan^{n-1} x}{n-1} - \int \tan^{n-2} x \, dx \, .$$

(5) Let
$$u = 3x \Rightarrow du = 3 dx \Rightarrow \int \tan^3(3x) dx = \frac{1}{3} \int \tan^3(3x) 3 dx = \frac{1}{3} \int \tan^3 u \, du$$
. Use

reduction formula #4 above to get $\frac{1}{3}\int \tan^3 u \, du = \frac{1}{3}\left(\frac{\tan^2 u}{2}\right) - \frac{1}{3}\int \tan u \, du =$

$$\frac{1}{6}\tan^2 u - \frac{1}{3}\ln|\sec u| \Rightarrow \int_{0}^{\frac{\pi}{4}} \tan^3(3x) \, dx = \left\{\frac{1}{6}\tan^2(3x) - \frac{1}{3}\ln|\sec(3x)|\right\}_{0}^{\frac{\pi}{4}} =$$

$$\left\{\frac{1}{6}\tan^{2}\left(\frac{3\pi}{4}\right) - \frac{1}{3}\ln\left|\sec\left(\frac{3\pi}{4}\right)\right|\right\} - \left\{\frac{1}{6}\tan^{2}(0) - \frac{1}{3}\ln\left|\sec(0)\right|\right\} = \frac{1}{6}(-1)^{2} - \frac{1}{3}\ln\left|-\sqrt{2}\right| - \frac{1}{6}\ln\left|\frac{1}{6}\left(-\frac{1}{6}\right)^{2}\right| - \frac{1}{6$$

$$\frac{1}{6}(0)^2 + \frac{1}{3}\ln 1 = \frac{1}{6} - \frac{1}{3}\ln\left(\sqrt{2}\right).$$

(6) Use the trigonometric identity $\cos^2 \Delta = \frac{1 + \cos 2\Delta}{2}$ to get $\int \cos^2(2x) dx =$

$$\int \frac{1 + \cos(4x)}{2} dx = \frac{1}{2} \int 1 dx + \frac{1}{2} \int \cos(4x) dx = \frac{1}{2} x + \frac{1}{8} \sin(4x) \Rightarrow \int \cos^2(2x) dx = \frac{\pi}{4}$$

$$\left\{\frac{1}{2}\left(\frac{\pi}{4}\right) + \frac{1}{8}\sin\pi\right\} - \left\{\frac{1}{2}(0) + \frac{1}{8}\sin(0) = \frac{\pi}{8}\right\}$$

(7) Use the trigonometric identity $\sin x \cos y = \frac{1}{2} [\sin(x-y) + \sin(x+y)]$ to get

$$\int \sin(5x)\cos(3x)\,dx = \frac{1}{2}\int \sin(2x)\,dx + \frac{1}{2}\int \sin(8x)\,dx = -\frac{1}{4}\cos(2x) - \frac{1}{16}\cos(8x) \Rightarrow$$

$$\int_{0}^{\frac{\pi}{8}} \sin(5x)\cos(3x)\,dx = \left\{-\frac{1}{4}\cos\left(\frac{\pi}{4}\right) - \frac{1}{16}\cos(\pi)\right\} - \left\{-\frac{1}{4}\cos(0 - \frac{1}{16}\cos(0)\right\} =$$

$$-\frac{1}{4}\left(\frac{\sqrt{2}}{2}\right) + \frac{1}{16} + \frac{1}{4} + \frac{1}{16} = \frac{3-\sqrt{2}}{8}$$
(8) $\int \tan^3 x \sec^3 x \, dx = \int \tan^2 x \sec^2 x (\sec x \tan x \, dx) =$

$$\int (\sec^2 x - 1) \sec^2 x (\sec x \tan x \, dx) = \int \sec^4 x (\sec x \tan x \, dx) -$$

$$\int \sec^2 x (\sec x \tan x \, dx) = \frac{1}{5} \sec^5 x - \frac{1}{3} \sec^3 x + C.$$
(9) $\int \sqrt{\sin x} \cos^3 x \, dx = \int \sqrt{\sin x} \left(\cos^2 x\right) (\cos x \, dx) = \int (\sin x)^{\frac{1}{2}} (1 - \sin^2 x) \cos x \, dx =$

$$\int (\sin x)^{\frac{1}{2}} \cos x \, dx - \int (\sin x)^{\frac{5}{2}} \cos x \, dx = \frac{2}{3} (\sin x)^{\frac{3}{2}} - \frac{2}{7} (\sin x)^{\frac{7}{2}} + C.$$

(10) $\int \cos^3 x \sin^2 x \, dx = \int \cos^2 x \sin^2 x (\cos x \, dx) = \int (1 - \sin^2 x) (\sin^2 x) \cos x \, dx =$

 $\int \sin^2 x (\cos x \, dx) - \int \sin^4 x (\cos x \, dx) = \frac{1}{3} \sin^3 x - \frac{1}{5} \sin^5 x + C.$

(11)
$$\int \frac{\sin^3 x}{\sqrt{\cos x}} \, dx = \int (\cos x)^{-\frac{1}{2}} \sin^2 x (\sin x \, dx = \int (\cos x)^{-\frac{1}{2}} (1 - \cos^2 x) \sin x \, dx =$$

$$\int (\cos x)^{-\frac{1}{2}} (\sin x \, dx) - \int (\cos x)^{\frac{3}{2}} (\sin x \, dx) = -2(\cos x)^{\frac{1}{2}} + \frac{2}{5}(\cos x)^{\frac{5}{2}} \Rightarrow$$

$$\int_{0}^{\pi/2} \frac{\sin^3 x}{\sqrt{\cos x}} \, dx = \left\{ -2\cos\left(\frac{\pi}{2}\right) + \frac{2}{5}\left(\cos\left(\frac{\pi}{2}\right)\right)^{5/2} \right\} - \left\{ -2\cos\left(\frac{\pi}{2}\right)^{5/2} \right\} = \frac{8}{5}$$

(12) Use the trigonometric identities $\cos^2 \Delta = \frac{1 + \cos 2\Delta}{2}$ and $\sin^2 \Delta = \frac{1 - \cos 2\Delta}{2}$.

$$\int \sin^{2} x \cos^{2} x dx = \int \left(\frac{1-\cos 2x}{2}\right) \left(\frac{1+\cos 2x}{2}\right) dx = \frac{1}{4} \int (1-\cos^{2} 2x) dx = \frac{1}{4} \int 1 dx - \frac{1}{4} \int \cos^{2} 2x dx = \frac{1}{4} x - \frac{1}{4} \int \left(\frac{1+\cos 4x}{2}\right) dx = \frac{1}{4} x - \frac{1}{8} \int 1 dx - \frac{1}{8} \int \cos 4x dx = \frac{1}{4} x - \frac{1}{8} x - \frac{1}{32} \sin 4x + C = \frac{1}{8} x - \frac{1}{32} \sin 4x + C.$$
(13)
$$\int \tan^{5} x \sec x dx = \int \tan^{4} x \tan x \sec x dx = \int (\tan^{2} x)^{2} \tan x \sec x dx = \int (\sec^{2} x - 1)^{2} \sec x \tan x dx = \int (\sec^{4} x - 2\sec^{2} x + 1) \sec x \tan x dx = \int \sec^{4} x (\sec x \tan x dx) - 2 \int \sec^{2} x (\sec x \tan x dx) + \int \sec x \tan x dx = \frac{1}{5} \sec^{5} x - \frac{2}{3} \sec^{3} x + \sec x + C.$$

Gamma Function



The (complete) gamma function $\Gamma(n)$ is defined to be an extension of the factorial to complex and real number arguments. It is related to the factorial by

 $\Gamma(n) = (n-1)!,$

a slightly unfortunate notation due to Legendre which is now universally used instead of Gauss's simpler $\Pi(n) = n!$ (Gauss 1812; Edwards 2001, p. 8).

It is analytic everywhere except at z = 0, -1, -2, ..., and the residue at z = -k is

$$\operatorname{Res}_{z=-k} \Gamma(z) = \frac{(-1)^k}{k!}$$

There are no points zat which $\Gamma(z) = 0$.

The gamma function is implemented in *Mathematica* as Gamma[z].

There are a number of notational conventions in common use for indication of a power of a gamma functions. While authors such as Watson (1939) use $\Gamma^n(z)$ (i.e., using a trigonometric function-like convention), it is also common to write $[\Gamma(z)]^n$.

The gamma function can be defined as a definite integral for R[z] > 0 (Euler's integral form)

$$\Gamma(z) \equiv \int_0^\infty t^{z-1} e^{-t} dt$$
$$= 2 \int_0^\infty e^{-t^2} t^{2z-1} dt,$$

-	
0	r
-	•

$$\Gamma(z) \equiv \int_0^1 \left[\ln\left(\frac{1}{t}\right) \right]^{z-1} dt.$$

Plots of the real and imaginary parts of $\Gamma(z)$ in the complex plane are illustrated above. Integrating equation (3) by parts for a real argument, it can be seen that

$$\Gamma(x) = \int_{0}^{\infty} t^{x-1} e^{-t} dt$$

$$= [-t^{x-1} e^{-t}]_{0}^{\infty} + \int_{0}^{\infty} (x-1) t^{x-2} e^{-t} dt$$

$$= (x-1) \int_{0}^{\infty} t^{x-2} e^{-t} dt$$

$$= (x-1) \Gamma(x-1).$$
If *x* is an integer *n* = 1, 2, 3, ..., then
$$\Gamma(n) = (n-1) \Gamma(n-1) = (n-1)(n-2) \Gamma(n-2)$$

$$= (n-1) (n-2) \cdots 1 = (n-1)!,$$

so the gamma function reduces to the factorial for a positive integer argument.

A beautiful relationship between $\Gamma(z)$ and the Riemann zeta function $\zeta(z)$ is given by

$$\zeta(z)\,\Gamma(z) = \int_0^\infty \frac{u^{z-1}}{e^u - 1}\,dz$$

for **R**[z] > 1(Havil 2003, p. 60).

The gamma function can also be defined by an infinite product form (Weierstrass form)

$$\Gamma(z) \equiv \left[z \, e^{\gamma z} \prod_{r=1}^{\infty} \left(1 + \frac{z}{r} \right) e^{-z/r} \right]^{-1},$$

where γ is the Euler-Mascheroni constant (Krantz 1999, p. 157; Havil 2003, p. 57). This can be written

$$\Gamma(z) = \frac{1}{z} \exp\left[\sum_{k=1}^{\infty} \frac{(-1)^k s_k}{k} z^k\right],$$

where

$$s_1 \equiv \gamma$$

 $s_k \equiv \zeta(k$

for $k \ge 2$, where $\zeta(z)$ is the Riemann zeta function (Finch 2003). Taking the logarithm of both sides of (\diamondsuit),

$$-\ln\left[\Gamma\left(z\right)\right] = \ln z + \gamma z + \sum_{n=1}^{\infty} \left[\ln\left(1 + \frac{z}{n}\right) - \frac{z}{n}\right].$$

Differentiating,

$$-\frac{\Gamma'(z)}{\Gamma(z)} = \frac{1}{z} + \gamma + \sum_{n=1}^{\infty} \left(\frac{\frac{1}{n}}{1 + \frac{z}{n}} - \frac{1}{n}\right)$$
$$= \frac{1}{z} + \gamma + \sum_{n=1}^{\infty} \left(\frac{1}{n+z} - \frac{1}{n}\right)$$

$$\Gamma'(z) \qquad = \qquad -\Gamma(z)\left[\frac{1}{z}+\gamma+\sum_{n=1}^{\infty}\left(\frac{1}{n+z}-\frac{1}{n}\right)\right]$$

$$\equiv \Gamma(z) \Psi(z) = \Gamma(z) \psi_0(z)$$

$$\Gamma'(1) = -\Gamma(1)\left\{1 + \gamma + \left[\left(\frac{1}{2} - 1\right) + \left(\frac{1}{3} - \frac{1}{2}\right) + \dots + \left(\frac{1}{n+1} - \frac{1}{n}\right) + \dots\right]\right\}$$

$$= -(1 + \gamma - 1) = -\gamma$$

$$\Gamma'(n) = -\Gamma(n)\left\{\frac{1}{n} + \gamma + \left[\left(\frac{1}{1+n} - 1\right) + \left(\frac{1}{2+n} - \frac{1}{2}\right) + \left(\frac{1}{3+n} - \frac{1}{3}\right) + \dots\right]\right\}$$

$$= -(n-1)!\left(\frac{1}{n} + \gamma - \sum_{k=1}^{n} \frac{1}{k}\right),$$

where $\Psi(z)$ is the digamma function and $\psi_0(z)$ is the polygamma function. *n*th derivatives are given in terms of the polygamma functions ψ_n , ψ_{n-1} , ..., ψ_0 .

COPYRIGHT FIMT 2020

197 | Page

The minimum value x_0 of $\Gamma(x)$ for real positive $x = x_0$ is achieved when

$$\Gamma'(x_0) = \Gamma(x_0)\psi_0(x_0) = 0$$

 $\psi_0(x_0)=0.$

This can be solved numerically to give $x_0 = 1.46163$...(Sloane's A030169; Wrench 1968), which has continued fraction [1, 2, 6, 63, 135, 1, 1, 1, 1, 4, 1, 38, ...] (Sloane's A030170). At x_0 , $\Gamma(x_0)$ achieves the value 0.8856031944... (Sloane's A030171), which has continued fraction [0, 1, 7, 1, 2, 1, 6, 1, 1, ...] (Sloane's A030172).

The Euler limit form is

$$\Gamma(z) = \frac{1}{z} \prod_{n=1}^{\infty} \left[\left(1 + \frac{1}{n} \right)^{z} \left(1 + \frac{z}{n} \right)^{-1} \right],$$

so

$$\Gamma(z) \qquad = \qquad \lim_{n \to \infty} \frac{(n+1)^z}{z (1+z) \left(1+\frac{z}{2}\right) \left(1+\frac{z}{3}\right) \cdots \left(1+\frac{z}{n}\right)}$$

 $= \lim_{n \to \infty} \frac{(n+1)^{z} n!}{z (z+1) (z+2) (z+3) \cdots (z+n)}$

$$\lim_{n\to\infty}\frac{n!}{(z)_{n+1}}(n+1)^z$$

$$= \lim_{n \to \infty} \frac{1}{(z)_{n+1}} n^{2}$$

(Krantz 1999, p. 156).

=

One over the gamma function $1/\Gamma(z)$ is an entire function and can be expressed as

$$\frac{1}{\Gamma(z)} = z \exp\left[\gamma z - \sum_{k=2}^{\infty} \frac{(-1)^k \zeta(k) z^k}{k}\right],$$

where γ is the Euler-Mascheroni constant and $\zeta(z)$ is the Riemann zeta function (Wrench 1968). An asymptotic series for $1/\Gamma(z)$ is given by

$$\frac{1}{\Gamma(z)} \sim z + \gamma z^2 + \frac{1}{12} \left(6 \gamma^2 - \pi^2 \right) z^3 + \frac{1}{12} \left[2 \gamma^3 - \gamma \pi^2 + 4 \zeta(3) \right] z^4 + \dots$$

Writing

$$\frac{1}{\Gamma(z)} = \sum_{k=1}^{\infty} a_k \, z^k,$$

the *a*_ksatisfy

$$a_n = n a_1 a_n - a_2 a_{n-1} + \sum_{k=2}^n (-1)^k \zeta(k) a_{n-k}$$

(Bourguet 1883, Davis 1933, Isaacson and Salzer 1943, Wrench 1968). Wrench (1968) numerically computed the coefficients for the series expansion about 0 of

$$\frac{1}{z(1+z)\Gamma(z)} = 1 + (\gamma-1)z + \left[1 + \frac{1}{2}(\gamma-2)\gamma - \frac{1}{12}\pi^2\right]z^2 + \dots$$

The Lanczos approximation gives a series expansion for $\Gamma(z+1)$ for z > 0 in terms of an arbitrary constant σ such that $R[z + \sigma + 1/2] > 0$.

The gamma function satisfies the functional equations

- $\Gamma(1+z) = z \Gamma(z)$
- $\Gamma(1-z) = -z \Gamma(-z).$

Additional identities are

$$\Gamma(x) \Gamma(-x) = -\frac{\pi}{x \sin(\pi x)}$$

$$\Gamma(x) \Gamma(1-x) = \frac{\pi}{\sin(\pi x)}$$

$$|(ix)!|^2 = \frac{\pi x}{\sinh(\pi x)}$$

$$|(n+ix)!| = \sqrt{\frac{\pi x}{\sinh(\pi x)}} \prod_{s=1}^n \sqrt{s^2 + x^2}.$$

Using (41), the gamma function $\Gamma(r)$ of a rational number r can be reduced to a constant times $\Gamma(\text{frac}(r))$ or $1/\Gamma(\text{frac}(r))$. For example,

$$\Gamma\left(\frac{2}{3}\right) = \frac{2\pi}{\sqrt{3}\,\Gamma\left(\frac{1}{3}\right)}$$

 $\Gamma\left(\frac{3}{4}\right) = \frac{\sqrt{2}\pi}{\Gamma\left(\frac{1}{4}\right)}$

$$\Gamma\left(\frac{3}{5}\right) = \sqrt{2 - \frac{2}{\sqrt{5}}} \frac{\pi}{\Gamma\left(\frac{2}{5}\right)}$$

$$\Gamma\left(\frac{4}{5}\right) \qquad = \qquad \sqrt{2 + \frac{2}{\sqrt{5}}} \frac{\pi}{\Gamma\left(\frac{1}{5}\right)}.$$

For $R[x] = -\frac{1}{2}$,

$$\left|\left(-\frac{1}{2}+iy\right)!\right|^2 = \frac{\pi}{\cosh(\pi y)}.$$

Gamma functions of argument 2zcan be expressed using the Legendre duplication formula

$$\Gamma(2z) = (2\pi)^{-1/2} 2^{2z-1/2} \Gamma(z) \Gamma(z+\frac{1}{2}).$$

Gamma functions of argument ³zcan be expressed using a triplication formula

$$\Gamma(3 z) = (2\pi)^{-1} 3^{3 z - 1/2} \Gamma(z) \Gamma(z + \frac{1}{3}) \Gamma(z + \frac{2}{3}).$$

The general result is the Gauss multiplication formula

$$\Gamma(z) \Gamma\left(z+\frac{1}{n}\right) \cdots \Gamma\left(z+\frac{n-1}{n}\right) = (2\pi)^{(n-1)/2} n^{1/2-nz} \Gamma(nz).$$

The gamma function is also related to the Riemann zeta function $\zeta(z)$ by \Box

$$\Gamma\left(\frac{s}{2}\right)\pi^{-s/2}\zeta(s)=\Gamma\left(\frac{1-s}{2}\right)\pi^{-(1-s)/2}\zeta(1-s).$$

For integer n = 1, 2, ..., the first few values of $\Gamma(n)$ are 1, 1, 2, 6, 24, 120, 720, 5040, 40320, 362880, ... (Sloane's A000142). For half-integer arguments, $\Gamma(n/2)$ has the special form

$$\Gamma\left(\frac{1}{2}n\right) = \frac{(n-2)!!\sqrt{\pi}}{2^{(n-1)/2}},$$

where n!! is a double factorial. The first few values for n = 1, 3, 5, ... are therefore

$$\Gamma(\frac{1}{2}) = \sqrt{\pi}$$

$$\Gamma\left(\frac{3}{2}\right) = \frac{1}{2}\sqrt{\pi}$$

$$\Gamma\left(\frac{5}{2}\right) = \frac{3}{4}\sqrt{\pi},$$

 $15\sqrt{\pi}/8$, $105\sqrt{\pi}/16$, ... (Sloane's A001147 and A000079; Wells 1986, p. 40). In general, for *n*a positive integer n = 1, 2, ...

$$\Gamma\left(\frac{1}{2}+n\right) = \frac{1\cdot 3\cdot 5\cdots (2\,n-1)}{2^n}\,\sqrt{\pi}$$

$$= \frac{(2\,n-1)\,!!}{2^n}\,\sqrt{\pi}$$

$$\Gamma\left(\frac{1}{2}-n\right) = \frac{(-1)^n\,2^n}{1\cdot 3\cdot 5\cdots (2\,n-1)}\,\sqrt{\pi}$$

$$= \frac{(-1)^n\,2^n}{(2\,n-1)\,!!}\,\sqrt{\pi}.$$

Simple closed-form expressions of this type do not appear to exist for $\Gamma(1/n)$ for n a positive integer n > 2. However, Borwein and Zucker (1992) give a variety of identities relating gamma functions to square roots and elliptic integral singular values k_n , i.e., elliptic moduli k_n such that

$$\frac{K'(k_n)}{K(k_n)} = \sqrt{n} ,$$

where K(k) is a complete elliptic integral of the first kind and $K'(k) = K(k') = K(\sqrt{1-k^2})_{is}$ the complementary integral. M. Trott (pers. comm.) has developed an algorithm for automatically generating hundreds of such identities.

$$\Gamma\left(\frac{1}{3}\right) = 2^{7/9} 3^{-1/12} \pi^{1/3} [K(k_3)]^{1/3}$$

$\Gamma(\frac{1}{4})$	=	$2\pi^{1/4} \left[K(k_1) \right]^{1/2}$
$\Gamma\left(\frac{1}{6}\right)$	=	$2^{-1/3} 3^{1/2} \pi^{-1/2} \left[\Gamma \left(\frac{1}{3} \right) \right]^2$
$\Gamma\left(\frac{1}{8}\right)\Gamma\left(\frac{3}{8}\right)$	=	$\left(\sqrt{2} - 1\right)^{1/2} 2^{13/4} \pi^{1/2} K\left(k_2\right)$
$\frac{\Gamma\left(\frac{1}{8}\right)}{\Gamma\left(\frac{3}{8}\right)}$	-AC	$2\left(\sqrt{2}+1\right)^{1/2}\pi^{-1/4}\left[K\left(k_{1}\right)\right]^{1/2}$
$\Gamma\left(\frac{1}{12}\right)$	and	$2^{-1/4} 3^{3/8} \left(\sqrt{3} + 1\right)^{1/2} \pi^{-1/2} \Gamma\left(\frac{1}{4}\right) \Gamma\left(\frac{1}{3}\right)$
$\Gamma\left(\frac{5}{12}\right)$	J.	$2^{1/4} 3^{-1/8} \left(\sqrt{3} - 1\right)^{1/2} \pi^{1/2} \frac{\Gamma\left(\frac{1}{4}\right)}{\Gamma\left(\frac{1}{3}\right)}$
$\frac{\Gamma\left(\frac{1}{24}\right)\Gamma\left(\frac{11}{24}\right)}{\Gamma\left(\frac{5}{24}\right)\Gamma\left(\frac{7}{24}\right)}$	-	$\sqrt{3}\sqrt{2+\sqrt{3}}$
$\frac{\Gamma\left(\frac{1}{24}\right)\Gamma\left(\frac{5}{24}\right)}{\Gamma\left(\frac{7}{24}\right)\Gamma\left(\frac{11}{24}\right)}$	_	$4 \cdot 3^{1/4} \left(\sqrt{3} + \sqrt{2} \right) \pi^{-1/2} K(k_1)$
$\frac{\Gamma\left(\frac{1}{24}\right)\Gamma\left(\frac{7}{24}\right)}{\Gamma\left(\frac{5}{24}\right)\Gamma\left(\frac{11}{24}\right)}$	-11	$2^{25/18} 3^{1/3} \left(\sqrt{2} + 1\right) \pi^{-1/3} \left[K\left(k_3\right)\right]^{2/3}$
$\Gamma\left(\frac{1}{24}\right)\Gamma\left(\frac{5}{24}\right)\Gamma\left(\frac{7}{24}\right)\Gamma\left(\frac{11}{24}\right)$	= R	$384\left(\sqrt{2}+1\right)\left(\sqrt{3}-\sqrt{2}\right)\left(2-\sqrt{3}\right)\pi\left[K(k_{6})\right]^{2}$
$\Gamma(\frac{1}{10})$	=	$2^{-7/10} 5^{1/4} \left(\sqrt{5} + 1\right)^{1/2} \pi^{-1/2} \Gamma\left(\frac{1}{5}\right) \Gamma\left(\frac{2}{5}\right)$
$\Gamma\left(\frac{3}{10}\right)$	0	$2^{-3/5} \left(\sqrt{5} - 1\right) \pi^{1/2} \frac{\Gamma\left(\frac{1}{5}\right)}{\Gamma\left(\frac{2}{5}\right)}$
$\frac{\Gamma\left(\frac{1}{15}\right)\Gamma\left(\frac{4}{15}\right)\Gamma\left(\frac{7}{15}\right)}{\Gamma\left(\frac{2}{15}\right)}$:201	$2 \cdot 3^{1/2} 5^{1/6} \sin\left(\frac{2}{15}\pi\right) \left[\Gamma\left(\frac{1}{3}\right)\right]^2$
$\frac{\Gamma\left(\frac{1}{15}\right)\Gamma\left(\frac{2}{15}\right)\Gamma\left(\frac{7}{15}\right)}{\Gamma\left(\frac{4}{15}\right)}$	=	$2^2 \cdot 3^{2/5} \sin\left(\frac{1}{5}\pi\right) \sin\left(\frac{4}{15}\pi\right) \left[\Gamma\left(\frac{1}{5}\right)\right]^2$

$\frac{\Gamma\left(\frac{2}{15}\right)\Gamma\left(\frac{4}{15}\right)\Gamma\left(\frac{7}{15}\right)}{\Gamma\left(\frac{1}{15}\right)}$	=	$\frac{2^{-3/2} 3^{-1/5} 5^{1/4} \left(\sqrt{5} -1\right)^{1/2} \left[\Gamma\left(\frac{2}{5}\right)\right]^2}{\sin\left(\frac{4}{15} \pi\right)}$
$\frac{\Gamma\left(\frac{1}{15}\right)\Gamma\left(\frac{2}{15}\right)\Gamma\left(\frac{4}{15}\right)}{\Gamma\left(\frac{7}{15}\right)}$	=	$60\left(\sqrt{5}-1\right)\sin\left(\frac{7}{15}\pi\right)\left[K\left(k_{15}\right)\right]^{2}$
$\frac{\Gamma\left(\frac{1}{20}\right)\Gamma\left(\frac{9}{20}\right)}{\Gamma\left(\frac{3}{20}\right)\Gamma\left(\frac{7}{20}\right)}$	-A($2^{-1} 5^{1/4} (\sqrt{5} + 1)$
$\frac{\Gamma\left(\frac{1}{20}\right)\Gamma\left(\frac{3}{20}\right)}{\Gamma\left(\frac{7}{20}\right)\Gamma\left(\frac{9}{20}\right)}$	EL P.C	$2^{4/5} \left(10 - 2\sqrt{5}\right)^{1/2} \pi^{-1} \sin\left(\frac{7}{20}\pi\right) \sin\left(\frac{9}{20}\pi\right) \left[\Gamma\left(\frac{1}{5}\right)\right]^2$
$\frac{\Gamma\left(\frac{1}{20}\right)\Gamma\left(\frac{7}{20}\right)}{\Gamma\left(\frac{3}{20}\right)\Gamma\left(\frac{9}{20}\right)}$	-	$2^{3/5} \left(10 + 2\sqrt{5}\right)^{1/2} \pi^{-1} \sin\left(\frac{3}{20}\pi\right) \sin\left(\frac{9}{20}\pi\right) \left[\Gamma\left(\frac{2}{5}\right)\right]^2$
$\Gamma\left(\frac{1}{20} ight)\Gamma\left(\frac{3}{20} ight)\Gamma\left(\frac{7}{20} ight)\Gamma\left(\frac{9}{20} ight)$	=	$160\left(\sqrt{5}-2\right)^{1/2}\pi\left[K\left(k_{5}\right)\right]^{2}.$

Several of these are also given in Campbell (1966, p. 31).

A few curious identities include

$$\prod_{n=1}^{2} \Gamma\left(\frac{1}{3}n\right) = \frac{2\pi}{\sqrt{3}}$$

$$\prod_{n=1}^{3} \Gamma\left(\frac{1}{3}n\right) = \frac{2\pi}{\sqrt{3}}$$

$$\prod_{n=1}^{4} \Gamma\left(\frac{1}{3}n\right) = \frac{2\pi\Gamma\left(\frac{1}{3}\right)}{3\sqrt{3}}$$

$$\prod_{n=1}^{5} \Gamma\left(\frac{1}{3}n\right) = \frac{8}{27}\pi^{2}$$

$$\prod_{n=1}^{6} \Gamma\left(\frac{1}{3}n\right) = \frac{8}{27}\pi^{2}$$

$$\prod_{n=1}^{7} \Gamma\left(\frac{1}{3}n\right) = \frac{32}{243} \pi^{2} \Gamma\left(\frac{1}{3}\right)$$
$$\prod_{n=1}^{8} \Gamma\left(\frac{1}{3}n\right) = \frac{640 \pi^{3}}{2187 \sqrt{3}},$$

of which Magnus and Oberhettinger 1949, p. 1 give only the last case,

$$\frac{\left[\Gamma\left(\frac{1}{4}\right)\right]^4}{16\pi^2} = \frac{3^2}{3^2 - 1} \frac{5^2 - 1}{5^2} \frac{7^2}{7^2 - 1} \cdots,$$
$$\frac{\Gamma'(1)}{\Gamma(1)} - \frac{\Gamma'\left(\frac{1}{2}\right)}{\Gamma\left(\frac{1}{2}\right)} = 2\ln 2$$

and

$$\frac{\Gamma'(1)}{\Gamma(1)} - \frac{\Gamma'\left(\frac{1}{2}\right)}{\Gamma\left(\frac{1}{2}\right)} = 2\ln 2$$

(Magnus and Oberhettinger 1949, p. 1). Ramanujan also gave a number of fascinating identities:

$$\frac{\Gamma^2(n+1)}{\Gamma(n+x\,i+1)\,\Gamma(n-x\,i+1)} = \prod_{k=1}^{\infty} \left[1 + \frac{x^2}{(n+k)^2} \right]$$

 $\phi(m, n) \phi(n, m) = \frac{\Gamma^3(m+1)\Gamma^3(n+1)}{\Gamma(2m+n+1)\Gamma(2n+m+1)} \frac{\cosh\left[\pi(m+n)\sqrt{3}\right] - \cos\left[\pi(m-n)\right]}{2\pi^2(m^2+mn+n^2)},$

where

$$\phi(m, n) \equiv \prod_{k=1}^{\infty} \left[1 + \left(\frac{m+n}{k+m}\right)^3 \right],$$

$$\prod_{k=1}^{\infty} \left[1 + \left(\frac{n}{k}\right)^3 \right] \prod_{k=1}^{\infty} \left[1 + 3\left(\frac{n}{n+2k}\right)^2 \right] = \frac{\Gamma\left(\frac{1}{2}n\right)}{\Gamma\left[\frac{1}{2}(n+1)\right]} \frac{\cosh\left(\pi n\sqrt{3}\right) - \cos\left(\pi n\right)}{2^{n+2}\pi^{3/2}n}$$

(Berndt 1994).

Ramanujan gave the infinite sums

$$\begin{split} &\sum_{k=0}^{\infty} (8\ k+1) \left[\frac{\Gamma\left(k+\frac{1}{4}\right)}{k!\,\Gamma\left(\frac{1}{4}\right)} \right]^4 \\ &= 1+9\left(\frac{1}{4}\right)^4 + 17\left(\frac{1\cdot 5}{4\cdot 8}\right)^4 + 25\left(\frac{1\cdot 5\cdot 9}{4\cdot 8\cdot 12}\right)^4 + \dots \\ &= \frac{2^{3/2}}{\sqrt{\pi} \left[\Gamma\left(\frac{3}{4}\right)\right]^2} \end{split}$$

and

NAAC ACCREDITED

$$\sum_{k=0}^{\infty} (-1)^k (4k+1) \left[\frac{(2k-1)!!}{(2k)!!} \right]^5$$

= $1 - 5 \left(\frac{1}{2} \right)^5 + 9 \left(\frac{1 \cdot 3}{2 \cdot 4} \right)^5 - 13 \left(\frac{1 \cdot 3 \cdot 5}{2 \cdot 4 \cdot 6} \right)^5 + \dots$
= $\frac{2}{\left[\Gamma \left(\frac{3}{4} \right) \right]^4}$

(The following asymptotic series is occasionally useful in probability theory (e.g., the onedimensional random walk):

$$\frac{\Gamma\left(J+\frac{1}{2}\right)}{\Gamma\left(J\right)} = \sqrt{J}\left(1-\frac{1}{8\,J}+\frac{1}{128\,J^2}+\frac{5}{1024\,J^3}-\frac{21}{32\,768\,J^4}+\dots\right)$$

(Graham *et al.* 1994). This series also gives a nice asymptotic generalization of Stirling numbers of the first kind to fractional values.

It has long been known that $\Gamma(\frac{1}{4})\pi^{-1/4}$ is transcendental (Davis 1959), as is $\Gamma(\frac{1}{3})$ (Le Lionnais 1983; Borwein and Bailey 2003, p. 138), and Chudnovsky has apparently recently proved that $\Gamma(\frac{1}{4})$ is itself transcendental (Borwein and Bailey 2003, p. 138).

There exist efficient iterative algorithms for $\Gamma(k/24)$ for all integers k (Borwein and Bailey 2003, p. 137). For example, a quadratically converging iteration for $\Gamma(1/4) = 3.6256099 \dots$ (Sloane's A068466) is given by defining

$$x_n = \frac{1}{2} \left(x_{n-1}^{1/2} + x_{n-1}^{-1/2} \right)$$

$$= \frac{y_{n-1} x_{n-1}^{1/2} + x_{n-1}^{-1/2}}{y_{n-1} + 1}$$

COPYRIGHT FIMT 2020

 y_n

TECHNICAL COMMUNICATION (103)

UNIT – 1 CONCEPTS AND FUNDAMENTALS

TECHNICAL COMMUNICATION:

Technical communication is playing an increasingly important role in business, industry, government, and even education. Many students take their first technical communication course

not really knowing much about the realities of the field. The tremendous expansion of the Internet has given anyone with a computer and a phone line access to vast quantities of information. Workers in every field are being called on to evaluate and interpret this information in various ways, using progress reports, news releases, product descriptions and specifications, instructions or procedures, funding proposals, feasibility studies, policy statements etc. In other words, most employees are being asked to become technical communicators: people who collect, organize, and interpret complex information and present it clearly and efficiently.

MEANING AND DEFINITION OF COMMUNICATION:

The word "communication" derived from the Latin word 'communicare' that means to impart, to

Participate, to share or to make common. It is a process of exchange of facts, ideas, and opinions and as a means that individual or organization share meaning and understanding with one another. Communication adds meaning to human life. It helps to build relationship and fosters love and understanding. It enriches our knowledge of the universe and makes living worthwhile. American Management Association defines, 'Communication is any behaviour that results in an

exchange of meaning'. Peter little defines communication as, 'Communication is the process by which information is transmitted between individuals and/or organizations so that an understanding response result'. Newman and Summer Jr. state that, 'Communication is an exchange of facts, ideas, opinions or emotions by two or more persons'.

IMPORTANCE OF COMMUNICATION:

- 1. For instruction
- 2. For integration
- 3. For information
- 4. For evaluation
- 5. For direction
- 6. For teaching
- 7. For influencing
- 8. For image building

COMMUNICATION SCOPE:

1. **Communication in personal life:** The entire life form birth to death of any person involves communication. No one can spend any moment without communication. A man interacts with his family members, fellow friends or colleagues involve communication. Even when he enjoys a T.V. program or reads newspaper, he is engaged in communication.

2. Communication in social life: Man is a social being. So, people live in a group in the society. To live in a society a man is to take part in the social functions and to maintain relations with the other members of the society.

3. **Communication in organizational life:** Communication is most important in business organization. No organization exists without communication. Communication is used in the following organization activities.

4. Communication in political life: Political parties give special emphasis on communication. Success of any political party depends on mass communication system.

5. **Communication in state affairs:** Various Government department used communication to run the administration and to inform people about development programs and other issues.

TYPES OF COMMUNICATION

Types of communication based on the communication channels used are:

1. Verbal Communication

Verbal communication refers to the form of communication in which message is transmitted verbally; communication is done by word of mouth and a piece of writing. Objective of every communication is to have people understand what we are trying to convey. In verbal communication remember the acronym KISS (keep it short and simple).

Verbal Communication is further divided into:

a) Oral Communication

b) Written Communication

Oral Communication

In oral communication, Spoken words are used. It includes face-to-face conversations, speech, telephonic conversation, video, radio, television, voice over internet. In oral communication, communication is influence by pitch, volume, speed and clarity of speaking.

Written Communication

In written communication, written signs or symbols are used to communicate. A written message may be printed or hand written. In written communication message can be transmitted via email, letter, report, memo etc. Message, in written communication, is influenced by the vocabulary & grammar used, writing style, precision and clarity of the language used.

2. Nonverbal Communication

Nonverbal communication is the sending or receiving of wordless messages. We can say that communication other than oral and written, such as gesture, body language, posture, tone of voice or facial expressions, is called nonverbal communication. Nonverbal communication is all

about the body language of speaker. Nonverbal communication helps receiver in interpreting the message received. Often, nonverbal signals reflect the situation more accurately than verbal messages. Sometimes nonverbal responses contradict verbal communication and hence affect the effectiveness of message.

Nonverbal communication has the following three elements:

Appearance: Speaker: clothing, hairstyle, neatness, use of cosmetics Surrounding: room size, lighting, decorations, furnishings

Body Language: Facial expressions, gestures, postures

Sounds: Voice Tone, Volume, Speech rate

THE COMMUNICATION PROCESS/CYCLE:



PROCESS OF COMMUNICATION

1. Sender or transmitter: The person who desires to convey the message is known as sender. Sender initiates the message and changes the behaviour of the receiver.

2. Message: It is a subject matter of any communication. It may involve any fact, idea, opinion or information. It must exist in the mind of the sender if communication is to take place.

3. **Encoding:** The communicator of the information organises his idea into series of symbols (words, signs, etc.) which, he feels will communicate to the intended receiver or receivers.

4. **Communication channel:** The sender has to select the channel for sending the information.

Communication channel is the media through which the message passes. It is the link that connects the sender and the receiver.

5. Receiver: The person who receives the message is called receiver or receiver is the person to

whom the particular message is sent by the transmitter. The communication process is incomplete without the existence of receiver of the message. It is a receiver who receives and tries to understand the message.

6. Decoding: Decoding is the process of interpretation of an encoded message into the understandable meaning. Decoding helps the receiver to drive meaning from the message.

7. Feedback: Communication is an exchange process. For the exchange to be complete the information must go back to whom from where it started (or sender), so that he can know the reaction of the receiver. The reaction or response of the receiver is known as feedback.



THE SHANNON-WEAVER MATHEMATICAL MODEL, 1949

Background

Claude Shannon, an engineer for the Bell Telephone Company, designed the most influential of all early communication models. His goal was to formulate a theory to guide the efforts of engineers in finding the most efficient way of transmitting electrical signals from one location to another. Later Shannon introduced a mechanism in the receiver which corrected for differences between the transmitted and received signal; this monitoring or correcting mechanism was the forerunner of the now widely used concept of feedback.

The Shannon-Weaver Mathematical Model, 1949



BERLO'S S-M-C-R, 1960

Background

GI Ehninger, Gronbeck and Monroe: "The simplest and most influential message-centred model of

our time came from David Berlo (Simplified from David K. Berlo, The Process of Communication (New York: Holt, Rinehart, and Winston, 1960):"

The model recognized that receivers were important to communication, for they were the targets.

The idea of "encoding" and "decoding" emphasized the problems we all have (psycholinguistically) in translating our own thoughts into words or other symbols and in deciphering the words or symbols of others into terms we ourselves can understand.

Berlo's Model of Com



A Source encode S-M-C-R Model. who

SCHRAMM'S INTERACTIVE MODEL, 1954 Background

Wilbur Schramm (1954) was one of the first to alter the mathematical model of Shannon and Weaver. He conceived of decoding and encoding as activities maintained simultaneously by sender and receiver; he also made provisions for a two-way interchange of messages. Notice also

the inclusion of an "interpreter" as an abstract representation of the problem of meaning.



Schramm's Model of Communication, 1954

NON-LINEAR MODELS DANCE'S HELICAL SPIRAL, 1967

Background Depicts communication as a dynamic process. Mortensen: "The helix represents the way communication evolves in an individual from his birth to the existing moment." Dance: "At any and all times, the helix gives geometrical testimony to the concept that communication while moving forward is at the same moment coming back upon itself and being affected by its past behaviour, for the coming curve of the helix is fundamentally affected by the curve from which it emerges. Yet, even though slowly, the helix can gradually free itself from its lower-level distortions. The communication process, like the helix, is constantly moving forward

and yet is always to some degree dependent upon the past, which informs the present and the future. The helical communication model offers a flexible communication process".



THEORIES OF COMMUNICATION

Bull's Eye Theory:

Bull's Eye Theory Action view is the basis for the theory of communication. The whole process

of communication is based on one-way action doing something to someone. The sender plays an

important role who encodes the message with the help of arbitrary symbols. The demonstration

or doing skills of the sender is for the purpose to change the behaviour of receiver.

Ping-Pong Theory:

Ping-Pong Theory this theory is also called interaction or interpersonal view. Ping-Pong is the game of table tennis, represents the interaction theory of communication. In communication process, the turns take place between the sender and receiver. In this theory, there is linear cause and effect.

Spiral Theory:

Spiral Theory the spiral theory of communication is also called as transactions view of communication. It recognizes more than one interaction between sender and the receiver. A transaction implies independence, mutual and reciprocal causality. Communication is not static but dynamic and life time experience.

ESSENTIALS OF GOOD COMMUNICATION - THE SEVEN CS OF COMMUNICATION

1. **Completeness -** The communication must be complete. It should convey all facts required by the audience. The sender of the message must take into consideration the receiver's mind set and convey the message accordingly.

2. Conciseness - Conciseness means wordiness, i.e., communicating what you want to convey in least possible words without forgoing the other C's of communication. Conciseness is a necessity for effective communication.

3.Consideration - Consideration implies "stepping into the shoes of others". Effective communication must take the audience into consideration, i.e., the audience's view points, background, mind-set, education level, etc. Make an attempt to envisage your audience, their requirements, emotions as well as problems.

4.Clarity - Clarity implies emphasizing on a specific message or goal at a time, rather than trying to achieve too much at once.

5. Concreteness - Concrete communication implies being particular and clear rather than fuzzy and general. Concreteness strengthens the confidence.

6. Courtesy - Courtesy in message implies the message should show the sender's expression as well as should respect the receiver. The sender of the message should be sincerely polite, judicious, reflective and enthusiastic.

7. Correctness - Correctness in communication implies that there are no grammatical errors in communication.

FACTORS RESPONSIBLE FOR GROWING IMPORTANCE OF COMMUNICATION

(1) Growth in the size and multiple locations of organizations

(2) Growth of trade unions: Over the last so many decades, trade unions have been growing strong. No management can be successful without taking the trade unions into confidence. Effective communication will create relationship between the management and the workers.

(3) Growing importance of human relations

(4) Public relations: Every organization has a social responsibility towards customers, government, suppliers and the public at large. Communication is the only way an organization

can project a positive image of itself.

(5) Advances in behavioural sciences: Modern management is deeply influenced by exciting discoveries made in behavioural sciences like psychology, sociology, transactional analysis etc. All of them throw light on suitable aspects of human nature and help in developing a positive attitude towards life and building up meaningful relationship. This is possible only through

communication.

(6) Technological advancement

CHANNELS OF COMMUNICATION

We divide the different types of communication medium into two different categories:

Physical media
 Mechanical media
 Physical media

These are the channels where the person who is talking can be seen and heard by the audience. In this not only hear the messages but to see the body language and feel the climate in the room is also important. This does not need to be two-way channels. In certain situations the receiver expects physical communication. This is the case especially when dealing with high concern messages. If a message is perceived as important to the receiver they expect to hear it live from their manager.

- Large meetings, town hall meetings
- Department meetings (weekly meetings)
- Up close and personal (exclusive meetings)
- Video conferences
- Viral communication or word of mouth
- Large meetings

Mechanical media

The second of the two types of communication medium is mechanical media. With mechanical media we mean written or electronic channels. These channels can be used as archives for messages or for giving the big picture and a deeper knowledge. But they can also be very fast. Typically though, because it is written, it is always interpret by the reader based on his or her mental condition. Irony or even humour rarely travels well in mechanical channels.

IAAC ACCREDITE

- E-mail
- Weekly letters or newsletters
- Personal letters
- Intranet
- Magazines or papers
- SMS
- Social media
- E-mail

VERBAL AND NON-VERBAL COMMUNICATION

BRHAG

Communication is at the heart of any relationship, be it familial, business, romantic, or friendly. While there have been significant advances in how we understand body language and other forms of communication, verbal communication continues to be the most important aspect of our interaction with other people. It's important to understand both the benefits and shortcomings of this most basic communication.

Advantages of Verbal Communication

- This can greatly increase both the speed and accuracy of communication.
- Verbal communication is far more precise than non-verbal cues
- Verbal communication is most effective when combined with other forms of communication like body language and gestures to help cue the intensity of the verbiage.

• Verbal communication is also the most effective way of explaining intangible concepts, as problem areas can be readily addressed and explained.

• Verbal communication also does not use natural resources in the way that technological methods or printing can.

Disadvantages of Verbal Communication

- There is a much smaller chance of an objective record.
- Verbal communication can also be quickly forgotten, especially if there are multiple points to consider.
- There is always the possibility of miscommunications leading to angry responses.

Non Verbal Communication

Nonverbal communication is the process of communication through sending and receiving wordless (mostly visual) cues between people. Messages can be communicated through gestures and touch, body language or posture, physical distance, facial expression and eye contact, which are all types of nonverbal communication.

Speech contains nonverbal elements known as paralanguage, including voice quality, rate, pitch,

volume, and speaking style, as well as prosodic features such as rhythm, intonation, and stress. However, much of the study of nonverbal communication has focused on face-to-face interaction, where it can be classified into three principal areas: environmental conditions where communication takes place, physical characteristics of the communicators, and behaviours of

1. Facial Expression: for happiness, sadness, anger and fear are similar throughout the world.

2. Gestures: Deliberate movements and signals are an important way to communicate meaning without words. Common gestures include waving, pointing, and using fingers to indicate numeric amounts.

3. Paralinguistic: refers to vocal communication that is separate from actual language. This includes factors such as tone of voice, loudness, inflection and pitch.

4. Body Language and Posture: Posture and movement can also convey a great deal on information. While these nonverbal behaviours can indicate feelings and attitudes

FORMAL AND INFORMAL COMMUNICATION

Formal Communication

In formal communication, certain rules, conventions and principles are followed while communicating message. Formal communication occurs in formal and official style. Usually professional settings, corporate meetings, conferences undergoes in formal pattern. In formal communication, use of slang and foul language is avoided and correct pronunciation is required. Authority lines are needed to be followed in formal communication.

£ 14001:201

Informal Communication

Informal communication is done using channels that are in contrast with formal communication channels. It's just a casual talk. It is established for societal affiliations of members in an organization and face-to-face discussions. It happens among friends and family. In informal communication use of slang words, foul language is not restricted. Usually, informal communication is done orally and using gestures. Informal communication, unlike formal communication, doesn't follow authority lines. In an organization, it helps in finding out staff grievances as people express more when talking informally. Informal communications, unlike in building relationships.

BARRIERS OF COMMUNICATION 1. Physiological Barrier

Physiological barriers to communication are related with the limitations of the human body and the human mind (memory, attention, and perception). Physiological barriers may result from individuals' personal discomfort, caused by ill-health, poor eye sight, or hearing difficulties.

- a) Poor Listening Skills
- b) Information Overload
- c) Inattention
- d) Emotions
- e) Poor Retention

2. Psychological Barrier

Psychological factors such as misperception, filtering, distrust, unhappy emotions, and people's state of mind can jeopardize the process of communication. We all tend to feel happier and more

MANAG

receptive to information when the sun shines. Similarly, if someone has personal problems such as worries and stress about a chronic illness, it may impinge his/her communication with others.

3. Social Barriers

Social barriers to communication include the social psychological phenomenon of conformity, a process in which the norms, values, and behaviours of an individual begin to follow those of the wider group. Social factors such as age, gender, socioeconomic status, and marital status may act

as a barrier to communication in certain situations.

4. Cultural Barriers

Culture shapes the way we think and behave. It can be seen as both shaping and being shaped by our established patterns of communication. Cultural barrier to communication often arises when individuals in one social group have developed different norms, values, or behaviours to individuals associated with another group. Cultural difference leads to difference in interest, knowledge, value, and tradition.

5. Semantic Barrier

Language, jargon, slang, etc., are some of the semantic barriers. Different languages across different regions represent a national barrier to communication. Use of jargon and slang also act as barrier to communication

6. Linguistic Barriers

Individual linguistic ability may sometimes become a barrier to communication. The use of difficult or inappropriate words in communication can prevent the people from understanding
the message. Poorly explained or misunderstood messages can also result in confusion. The linguistic differences between the people can also lead to communication breakdown. The same word may mean differently to different individuals.

7. Organizational Barriers

Unclear planning, structure, information overload, timing, technology, and status difference are

the organizational factors that may act as barriers to communication.

AIDS TO COMMUNICATION

Audio/video aids in business communication and training help in numerous ways. Each individual understands and retains information differently, which is why professional training organizations use multiple audio and visual tools during presentation sessions. This type of session is commonly known as a "multimedia presentation," which can include written, visual, auditory and sometimes interactive methods. Using visual aids can save your business time, especially if the subject contains information that may be too lengthy for written or oral communication. Pie charts, graphs, diagrams, photographs, video shorts and animation can often help explain subject matter quickly, and in a manner that is more easily absorbed by the learner. Visual aids can include projectors, flip charts, models, white boards or any combination thereof.

Auditory Aids

The type of auditory aid used in your multimedia presentation is most likely based on your budget, but the effectiveness of each type must be considered.

Handouts

Handouts serve to reinforce oral and visual components, and can go into further detail if desired or warranted. Handouts allow participants to follow along with presented information, make notes, formulate questions or refer to key points in the future as a "refresher course."

UNIT-II

WRITTEN COMMUNICATION

Written communication has great significance in today's business world. It is an innovative activity of the mind. Effective written communication is essential for preparing worthy promotional materials for business development. Speech came before writing. But writing is more unique and formal than speech. Effective writing involves careful choice of words, their organization in correct order in sentences formation as well as cohesive composition of sentences. Also, writing is more valid and reliable than speech. But while speech is spontaneous,

writing causes delay and takes time as feedback is not immediate.

OBJECTIVES OF WRITTEN COMMUNICATION

Written communication aims to inform someone of something in a way that they are able to read

and understand the message, with an intention of responding to it. (Story telling, narrating) If a form of written communication cannot be understood by the recipient then the message may well as not exist. In order to write a piece of information that can be understood clearly you need to have the correct spelling, punctuation and grammar. In addition, depending on the form of written communication you need to make sure you use the right format. For example, if you are writing a letter you need to ensure you are using the appropriate format.

MEDIA OF WRITTEN COMMUNICATION

- Letters
- E-mails
- Books
- Pamphlet
- Memorandum
- Notices
- Circulars
- Magazines
- The Internet or via other media, etc.

MERITS AND DEMERITS OF WRITTEN COMMUNICATION

th P. H. P.

Merits:

1. It is a permanent means of communication. Thus, it is useful where record maintenance is required.

- 2. It assists in proper delegation of responsibilities.
- 3. Written communication is more precise and explicit.
- 4. Effective written communication develops and enhances an organization's image.
- 5. It provides ready records and references.

Demerits:

1) Written communication does not save upon the costs. It costs huge in terms of stationery and the manpower employed in writing/typing and delivering letters.

2) Also, if the receivers of the written message are separated by distance and if they need to clear their doubts, the response is not spontaneous.

3) Written communication is time-consuming as the feedback is not immediate. The encoding and sending of message takes time.

4) Too much paper work and e-mails burden is involved.

PLANNING AND PREPARING OF EFFECTIVE BUSINESS MESSAGES.

While preparing a written or an oral business message, planning a business message is a must do task, you need to plan, organize, compose, edit and revise it. The message must also be proofread and corrected before it is mailed. Apart from the steps mentioned above the writer must take care of seven C qualities and also of legal aspect. Careful preparation of communication is important, even if the writer / speaker has the modern technology. These are the steps for planning a business message.

PLANNING A BUSINESS MESSAGE STEPS

- 1. Define the purpose of the message.
- 2. Analyze your audience readers or listeners.
- 3. Choose the ideas to include.
- 4. Collect all the facts to back up these ideas.5. Outline organize your message.

PERSUASIVE WRITING

Persuasive writing, known as creative writing or an argument, is a piece of writing in which the writer uses words to convince the reader that the writer's opinion is correct with regard to an issue. Persuasive writing sometimes involves convincing the reader to perform an action, or it may simply consist of an argument or several arguments to align the reader with the writer's point of view. Persuasive writing is one of the most commonly used writing types in the world. This type of writing is often used for advertising copy. A well-written persuasive piece is supported with a series of facts which help the author argue his or her point.

OVERVIEW OF TECHNICAL RESEARCH AND REPORT WRITING TECHNICAL WRITING

It is a communication in any field whose primary aim is to convey a particular piece of information for a particular purpose to particular readers.

_ It is objective, clear and accurate, concise and unemotional in its presentation of facts.

_ Special techniques that often uses are definitions, descriptions of mechanism, and descriptions of processes, classifications and interpretations.

- _ Deals with topics of : technical nature, Science, engineering, and technology
- _ Deals with an object, process, system, or abstract idea
- _ It used to Stress on the accuracy rather than style
- It focus on the Technical content and Not the author's feelings about it

NATURE OF TECHNICAL WRITING

1. Technical writing is exposition about scientific subjects and about various technical subjects associated with the sciences.

2. Technical writing is characterized by certain formal elements such as its scientific and technical vocabulary, its use of graphical aids and its use of conventional report forms.

3. Technical writing is ideally characterized by the maintenance of an attitude of impartially and objectivity, by the extreme care to convey information accurately and concisely, and by the absence of any attempt to arouse emotion 4. Technical writing is writing in which there is a relatively high concentration of certain complex and important writing techniques in particular, definition, description of mechanism, and description of a process, classification and interpretation.

PURPOSE OF TECHNICAL WRITING

1. It gives information in decision making and task accomplishments.

- 2. It analyzes events and their implications, the failure of systems.
- 3. It persuades and influences decision making.

CHARACTERISTICS OF TECHNICAL WRITING

1. Technical writing information flow easily and clearly.

2. Technical writing emphasizes objective reporting with no room for different interpretations, sentences structure and paragraph organization, declarative sentences with third-person pronouns.

3. Technical writing emphasizes factual data, statistics and measurable elements.

BASIC PRINCIPLES OF GOOD TECHNICAL WRITING

1. Writers should always have in mind a specific reader, real or imaginary, when writing their report and always assume that they are intelligent but uninformed.

2. They should decide on their exact purpose in writing.

- 3. They should use simple, concrete and familiar language.
- 4. They should check or review their writing from time to time.
- 5. They should make the paper as neat and as attractive as possible.

PROPERTIES OF TECHNICAL WRITING

Accuracy

One of the essential characteristics of technical writing is maintaining accuracy.

Clarity

Write the technical document in a layman's tone so that the customers are also able to understand what the product is all about. Try cutting down on the use of jargons because again, this is going to confuse the customers.

Descriptiveness

Be as descriptive in distinguishing the technical product as you can. More than half of the customers come from a non technical background and they need to have sufficient details otherwise they will not be able to picture the product correctly.

Correctness

Technical writing requires that you use correct grammar and sentence structure. Write down the key features in the form of headings, sub headings or bullet points as this will make the manual easy for the customers to read.

TECHNICAL WRITING PROCESS

1) planning
2) writing
3) delivery
4) archiving

ROLE OF TECHNICAL WRITING

_ Technical writers take complex information and communicate it clearly, concisely and accurately without relying on technical or corporate jargon to explain what they're trying to say.

_ Technical communications are created and distributed by most employees in service organizations today, especially by professional staff and management.

_ Writing well is difficult and time-consuming, and writing in a technical way and about technical subjects compounds the difficulties. The entire point of communications is to disseminate useful information. To be useful, information must be understood and acted upon.

_ Effective communications require quality content, language, format, and more. To present the appropriate content, it is imperative to understand one's audience and writing purpose.

_ If a document does not communicate the information that the writer intends and what he or she wants the reader to understand, then the communication is meaningless.

THE HOLISTIC GUIDE TO TECHNICAL WRITING

Technical writing is simplifying complex task, situation or tool in simple concise, easy to understand form. Technical writing is an art to delivering technical information to technical or non-technical users in a simple and easy to understandable form. Technical Writing is done for a purpose. Technical writing is creating documents that help someone install, uninstall, configure or use a product or a tool or a service. It results in the creation of things such as user manuals, admin guides, instruction booklets and help systems, installation manual, but it just not restricted to this you can also write brochure, PowerPoint slides, etc.

• Writing Ability: It is not required for a technical writer be an expert in any technology but a technical writer must have flair of writing simple, customized, concise, and error less documents. A technical writer must be grammatically well sound. If you are not grammatically sound then you cannot be a successful technical writer. A technical writer must have ability to write on diverse topics.

• Analytical Nature: A technical writer must have the capability to analysis the given content and

to create good technical documents on that analysis.

• Information Gathering: Information gathering skills of a technical writer makes him/her successful in his/her the field of technical writing. If you have flair of writing but don't have the

capability to gather the relevant information for a technical document then you are a waste for an organization. Technical writing revolves around the information gathering.

• Quick Learning and Adaptive Capability: A technical writer must be a quick learner and have capability to grab the new concepts in a short period. With the growth of economy, various new products are being launched in the market.

ABC's of Technical Writing

1. Accuracy - must be tactful in the recording of data, statement of calculating mathematical figures.

2. Brevity - it's easier to grasp the main idea of the report written if you have a brief report.

3. Confidence - a writer has to be decisive or sure of what is he writing about.

4. **Dignity** - ring of authority - all grammatical constructions must be complete, no flowery words. Ideas/info must be well-organized, simplified, summarized & expressed in straightforward manner.

5. Emphasis - stressing the major points & subordinating them.

6. Facility - devices used by the writer.

a) Pacing - technical/unfamiliar info should be presented from small to large quantity.

b) Sequence - familiar to unfamiliar; simple to complex.

c) **Arrangement** - significant details should be stressed & balanced to show their proper relationship.

d) Continuity - thought should be clearly established, illustrated, or stated.

7. Grammatical Correctness - reflects the communicative competence of the writer.

8. Honesty - if the writer borrowed some statements, ideas, or quotations, he has to acknowledge

them either in footnotes, end-notes, or cite the author or sources.

OUTPUTS OR END PRODUCTS OF TECHNICAL WRITING

Technical writing produces Help topics, User Manuals, Reference Manuals and other documentation generally produced by the manufacturer of technology.

Some results of Technical Writing are:

Technical Report, Abstract ,Feasibility Report, Business Letter, Brochure, Contract, Instructional

Manual, Proposal, Progress Report, Policy, Article for a Technical Journal, Monograph, Memorandum, Graphic Aids, Specification, Printed Action Memo, Survey Report, Letter Report,

Laboratory Report, Technical Paper

WRITING A PROPOSAL

Writing a proposal is usually the result of one of two scenarios; you are responding to a Request for Proposal (RFP) which asks for a solution to a problem or; you have identified an opportunity to provide a product or service to satisfy a requirement a potential client may not even know that they need.

IAAC ACCREDITED

MAN

Components of a Proposal

- Title Page
- Cover Letter
- Proprietary Notice
- Introduction
- Technical Approach
- Project Team
- Relevant Experience
- Project Price/Budget
- Schedule
- Certifications
- Appendices
- Final Review

WRITING PROPOSAL

The general purpose of any proposal is to persuade the readers to do something, whether it is to persuade a potential customer to purchase goods and/or services, or to persuade your employer to fund a project or to implement a program that you would like to launch. Any proposal offers a plan to fill a need, and your reader will evaluate your plan according to how well your written presentation answers the questions of *WHAT* you are proposing, *HOW* you plan to do it, *WHEN* you plan to do it, and *HOW MUCH* it is going to cost. The most basic composition of a proposal, as with any other written document, is simple; it needs a *beginning* (the Introduction), a *middle* (the Body of material to be presented) and an *end* (the Conclusion/Recommendation).

• The INTRODUCTION presents and summarizes the problem you intend to solve and your solution to that problem, including the benefits the reader/group will receive from the solution and the cost of that solution.

• The BODY of the proposal should explain the complete details of the solution: how the job will be done, broken into separate tasks; what method will be used to do it, including the equipment, material, and personnel that would be required; when the work will begin; and, when the job will be completed. It should also present a detailed cost breakdown for the entire job.

• The CONCLUSION should emphasize the benefits that the reader will realize from your solution to the problem and should urge the reader to action. It should be encouraging, confident and assertive in tone.

FORMAT OF PROPOSAL

- Front Matter
- Letter of transmittal
- Title Page
- Project Summary (approx. 200 word abstract)
- Introduction
- Body

• Project Proposal: (Includes Statement of the Problem, Proposed Solution(s), Program of Implementation, Conclusions/Recommendations)

- Conclusion/Recommendations
- Back Matter
- Bibliography and/or Works Cited
- Qualifications (of writer(s) and/or project implementers)
- Budget

Example: format of proposal for event:

	Event Proposal
	[Event Title]
To,	
	[name of the sponsor to whom the proposal is being submitted] [proposed date of the event in DD/MM/YYYY format]
Event outline	-
	[name of the event]
	[topic on which the event is based] [estimated duration of the event]
	Main Activities of Event
	Budget Outlines
Total:	[total cost]
	[Signature] Event Proposal Template
15	0 9001:2015 & 14001:2015

1

WRITING LETTERS

BUSINESS LETTERS: A business letter is usually used when writing from one company to another, or for correspondence between such organizations and their customers, clients and other external parties. The overall style of letter will depend on the relationship between the parties concerned. There are many reasons to write a business letter. It could be to request direct information or action from another party, to order supplies from a supplier, to identify

a mistake that was committed, to reply directly to a request, to apologize for a wrong or simply to convey goodwill.

SIGNIFICANCE

Assist in sustaining business relationship to convey complex information serve as permanent record to reach a large and geographically diverse audience

PURPOSES

- to informto congratulate
- to enquire
- to order
- to request
- to collect dues
- to complain

BUSINESS LETTERS STRUCTURE

th PHA A

- _ Heading
- _ Inside name & address
- _ Salutation
- _ Subject
- _ body of letter
- _ Complementary close
- _ signature
- _ Additional Elements
- Enclosure notation
- _ Postscript/ Identification mark



COPYRIGHT FIMT 2020

OFFICE MEMORANDUM

Memorandum: A memo usually has a smaller demographic of audience and is usually more exclusive. Memos are less public and normally targeted at a more exclusive, smaller audience. They are often used as a way of reminding someone of something that needs to be done, or to pass on a proposal or idea of some kind. They are most commonly typed in today's technological era, however they can also be hand written.

GOOD NEWS AND BAD NEWS LETTERS

If you've experienced something positive in your life such as a promotion, the birth of a child, or any award, it is a good idea to notify your friends, family and well-wishers through a good news letter. This is a simple, friendly document, usually informal in tone and style. It should inform

the addressee about the good news before going into details of the news. If the news is about a personal achievement, you should remember to be brief and humble. The letter is an occasion to inform, not to brag about your achievement. You should end the letter by thanking the addressee for his / her support.

Example:

{Date}

{Address of Business}

Dear {Name},

We have received your letter concerning {issue}, and we would like to extend to you first and foremost our sincerest apology, as well as let you know that we can help you! We would like to take immediate action to rectify this situation, and as such we will be immediately {explain what you will be doing to fix the problem}. We hope that this solution is to your satisfaction. In order for us to render you this service, we require you to fill out the enclosed form and follow the instructions carefully. Upon receipt of your form, we will promptly take the promised actions. Thank you for your patience and your business. We appreciate this opportunity to correct the problem.

Sincerely,

{Name}

PERSUASIVE LETTERS

111 Main Street Fallsington 19054 April 25, 2005 Mr. Steven Rogel Veotrhauser Company Federal Way, WA 98063-9

I am a student at Fallsington Elementary School in Fallsington, PA. I in thing this letter to ask you to not cut down too much of the rainforest. I hat rainforests once covered 14% of the earth and that they now only cover % of the earth?

If you decide to out down the rainforest we will have fewer medicines and you have a set of the out down too much of the rainforest or there will be no oxygen. I hope that you will try to preserve the animala, plants and trees that can only live in the rainforest. Instead of outling down the rainforest you could make paper from recycled paper.

oril 22, 200

SALES LETTERS

_ A sales letter is also referred as Letter of Sale, Marketing Sales Letter and Business Sales Letter.

_ It is a type of business letter; meant for generating business.

_ Letter written to publicise, advertise and ultimately sell a product or a service to the consumers.

_ These letters enable a businessman to approach present & potential customers easily & at low cost.

OBJECTIVES OF SALES LETTERS

- _ To promote sales of product, a service or an idea.
- _ To introduce new products in the market quickly, effectively & at low cost.
- _ To introduced the salesman to the potential customer.
- _ To widen the market for existing products.
- _ To remind customers about the product/ service.
- _ To keep customers in regular touch with the company & its products & services.

ELEMENTS OF SALES LETTER

- _ Appealing & persuasive
- _ Attractive
- _ Creative in nature
- _ Complete, explain the product or service in detail
- _ A brochure/ pamphlet, etc. may be attached with it.

PURPOSE OF SALES LETTERS

To persuade the readers to "buy" a product, service, idea, or point of view

- Grab the reader's attention
- Highlight the product's appeal
- Show the product's use
- Conclude with a request for action (buy it!)
- To make direct sales
- To announce and test the reaction to new services and products

Headline

Header:

- Sender's Letterhead or Sender's Name and Address
- _ The Recipient's Name (specific official, person and organisation)

- 7

and Address

- _ Date (can be placed after senders address)
- _ Reference or Subject (optional)

_ Salutation - Dear Sir/Madam/ Mr./Ms.

Introduction

Introductory lines regarding the product or service

Body

D 1 C	1 . /	• •	1 • 1	1.0
Relevance of t	nroduct /s	service 1	n daily	lite
	product / s		n uany	me

- _ Assistive information towards the purchase process
- _ Compliments and offer of assistance

Closing

- _ Complimentary Gesture Thanks, Thank you etc.
- _ Valediction Ex. Sincerely
- _ Signature or Signature Line
- _ Your Typed Name
- ENCL (optional) stands for 'Enclosure'

GEST

Logo or Motto of your Organisation

LETTER STYLES / LAYOUT

Forms of layout in letter writing are as follows:

- 1. Indented Form (traditional form)
- 2. Hanging Indention
- 3. Block Form (more modern form)

4. Semi Block

These different forms of layout are shown below:

Layout 1: Indented Form

SENDER'S NAME AND ADDRESS

Tel No: Ref:

E Mail: Date:

Dear

		Complimentary close	
Layout 2: Hanging Indentation	1 Form	and Signature	
Tel	SENDER'S NAME AND ADDRESS	Ref	
E-mail:		Date:	
Mr R. Shermani, Munager, Bank of Baroda, Hyderabad.			
Dear Sir,			
		Complimentary close	
Layout 3: Block Form		and englished to	
Tel	SENDER'S NAME AND ADDRESS	Ref	
E-mail:		Date:	
Mr R. Shermani,			
Bank of Baroda,			
Hyderabad.			
Dear Sir,			

COPYRIGHT FIMT 2020

REPORT WRITING

A report is a presentation of facts and findings, usually as a basis for recommendations; written for a specific readership, and probably intended to be kept as a record. Report writing is an essential skill for professionals. A report aims to inform, as clearly and succinctly as possible.

An effective report can be written going through the following steps-

1. Determine the objective of the report, i.e., identify the problem.

2. Collect the required material (facts) for the report.

- 3. Study and examine the facts gathered.
- 4. Plan the facts for the report.
- 5. Prepare an outline for the report, i.e., draft the report.
- 6. Edit the drafted report.

7. Distribute the draft report to the advisory team and ask for feedback and recommendations.

A report should generally include the following sections.(Sections marked with an asterisk (*) are essential: others are optional depending on the type, length and purpose of the report.)

- Letter of transmittal
- Title page*
- Table of contents
- List of abbreviations and/or glossary
- Executive summary/abstract
- Introduction*
- Body*
- Conclusion*
- Recommendations
- Bibliography
- Appendices

The essentials of good/effective report writing are as follows-

1) Know your objective, i.e., be focused.

2) Analyze the niche audience, i.e., make an analysis of the target audience, the purpose for which audience requires the report, kind of data audience is looking for in the report, the implications of report reading, etc.

3) Decide the length of report.

4) Disclose correct and true information in a report.

5) Discuss all sides of the problem reasonably and impartially. Include all relevant facts in a report.

6) Concentrate on the report structure and matter. Pre-decide the report writing style. Use vivid structure of sentences.

7) The report should be neatly presented and should be carefully documented.

8) Highlight and recap the main message in a report.

9) Encourage feedback on the report from the critics. The feedback, if negative, might be useful if properly supported with reasons by the critics. The report can be modified based on such feedback.

TYPES OF REPORTS

Business Reports

Business reports are a type of assignment in which you analyse a situation (either a real situation

or a case study) and apply business theories to produce a range of suggestions for improvement.

Business reports are typically assigned to enable you to:

- _ Examine available and potential solutions to a problem, situation, or issue.
- _ Apply business and management theory to a practical situation.
- _ Demonstrate your analytical, reasoning, and evaluation skills in identifying and

weighing-up possible solutions and outcomes.

- _ Reach conclusions about a problem or issue.
- _ Provide recommendations for future action.
- _ Show concise and clear communication skills.

Academic Reports

It is a research study on various aspects of the subjects. It generally takes form of research report which covers a wide variety of subjects & its coverage is also quite expensive Contents are similar in both the reports but some additional contents are there in academic reports which are as follows:-

- 1) Statement of the problem
- 2) Overview of literature
- 3) The conceptual framework
- 4) Research questions/hypothesis
- 5) Coverage
- 6) Data collection
- 7) Data processing

FORMAT OF REPORT

• Title Section - If the report is short, the front cover can include any information that you feel is necessary including the author(s) and the date prepared. In a longer report, you may want to include a table of contents and definitions of terms.

• Summary - There needs to be a summary of the major points, conclusions, and recommendations. It needs to be short as it is a general overview of the report. Some people will read the summary and only skim the report, so make sure you include all the relevant information. It would be best to write this last so you will include everything, even the points that might be added at the last minute.

• Introduction - The first page of the report needs to have an introduction. You will explain the problem and show the reader why the report is being made. You need to give a definition of terms if you did not include these in the title section, and explain how the details of the report are arranged.

• Body - This is the main section of the report. This section can include jargon from your industry. There needs to be several sections, with each having a subtitle. Information is usually arranged in order of importance with the most important information coming first. If you wish, a "Discussion" section can be included at the end of the Body to go over your findings and their significance.

• Conclusion - This is where everything comes together. Keep this section free of jargon as most people will read the Summary and Conclusion.

• Recommendations - This is what needs to be done. In plain English, explain your recommendations, putting them in order of priority.

• Appendices - This includes information that the experts in the field will read. It has all the technical details that support your conclusions.

This report writing format will make it easier for the reader to find what he is looking for. Remember to write all the sections in plain English, except for the Body. Also remember that the information needs to be organized logically with the most important information coming first.

LAYOUT /DRAFTING OF THE REPORT

- 1. Title Page
- 2. Figures and Tables
- 3. Equations and Formulae
- 4. Chapter Numbering System
- 5. Font
- 6. Appendices

ESSENTIAL REQUIREMENT OF GOOD REPORT WRITING

All reports need to be clear, concise and well structured. The key to writing an effective report is to allocate time for planning and preparation. With careful planning, the writing of a report will be made much easier. The essential stages of successful report writing are described below.

Consider how long each stage is likely to take and divide the time before the deadline between the different stages. Be sure to leave time for final proof reading and checking.

Stage One: Understanding the report brief

This first stage is the most important. You need to be confident that you understand the purpose of your report as described in your report brief or instructions.

Stage Two: Gathering and selecting information

Once you are clear about the purpose of your report, you need to begin to gather relevant information. Your information may come from a variety of sources, but how much information you will need will depend on how much detail is required in the report.

Stage Three: Organizing your material

Once you have gathered information you need to decide what will be included and in what sequence it should be presented. Begin by grouping together points that are related. These may form sections or chapters.

Stage Four: Analyzing your material

Before you begin to write your first draft of the report, take time to consider and make notes on the points you will make using the facts and evidence you have gathered. What conclusions can be drawn from the material? What are the limitations or flaws in the evidence?

Stage Five: Writing the report

Having organized your material into appropriate sections and headings you can begin to write the first draft of your report. You may find it easier to write the summary and contents page at the end when you know exactly what will be included. Aim for a writing style that is direct and precise. Avoid waffle and make your points clearly and concisely.

Stage Six: Reviewing and redrafting

Ideally, you should leave time to take a break before you review your first draft. Be prepared to rearrange or rewrite sections in the light of your review. Writing on a word processor makes it easier to rewrite and rearrange sections or paragraphs in your first draft. If you write your first draft by hand, try writing each section on a separate piece of paper to make redrafting easier.

Stage Seven: Presentation

Once you are satisfied with the content and structure of your redrafted report, you can turn your attention to the presentation. Check that the wording of each chapter/section/subheading is clear and accurate. Check that you have adhered to the instructions in your report brief regarding format and presentation. Check for consistency in numbering of chapters, sections and appendices. Make sure that all your sources are acknowledged and correctly referenced. **JOB APPLICATION**

An application for employment, job application, or application form is a form or collection

of forms that an individual seeking employment, called an applicant, must fill out as part of the

process of informing an employer of the applicant's availability and desire to be employed, and

persuading the employer to offer the applicant employment

RNAC

TYPES OF JOB APPLICATION

1. Handwritten

2. Electronic / Online

CONTENT OF JOB APPLICATION

Your name You're Address Your city, state, code Your phone number Your Email Id Your contact name Dear Sir.

DRAFTING THE APPLICATION

Job Application

William Mathura (Your Name) Model Village (Your Address) North Point, Hong Kong (Your Address) 22 July 2011 (Your Address) Mr. Hanukah Chan (Recipient's Name) Personnel Manager (Recipient's Designation) Wong & Lim Consultants (Recipient's Address) P.O. Box 583 (Recipient's Address) Kwai Chung, Kowloon (Recipient's Address) Dear Mr. Chan (Salutation)

Sub: Application for the Post of MTO

I am writing to apply for the post of MTO (Management Training Officer), which was advertised on the Oriental Daily Newspaper of the Hong Kong and Student Board of Polytechnic University on 21st July 2011. I have one year working experience at King City Garment Manufactory Limited. This experience plays a pivot role in improving my leadership skills, communication skills and ability to work in a team environment. I can fluently speak and write English. I also have fluency in speaking and writing Mandarin, and can therefore work in mainland China.

Currently I am studying a M.B.A. in Management at the Hong Kong Polytechnic University, graduating in 2012. I am studying subjects relevant to the post of MTO including Operations Management, Human Resources Management, Accounting, Marketing and Strategic Management. My final year project is entitled "Research and Knowledge Management Practices" in HK. Execution of this project will surely improve my communication skills, my leadership skills and my ability to lead and supervise subordinates effectively. I have also learned how to run a project from the planning stage to its completion.

Working for Wong & Lim Consultants appeals to me because it has a good reputation and it provides excellent training. Your organization produces a high-quality service, and I can contribute to this with my leadership skills and my ability to work under pressure. I am available for interview at any time. I can be contacted most easily on the cell phone number given below. I look forward to meet you soon.

ACCREDI

Yours sincerely, William Mathura Phone: 24862893 Mobile: 95427415 E-mail: abc@hkinternet.com Encl: Resume

(Your resume is enclosed with this Application Letter)

(Closing)

PREPARATION OF RESUME

A resume is a summary of the qualities and qualifications of a person. It is a informative and inspiring piece of written communication.

ARGEME

5 Rules for Building a Great Resume

Your resume has one job: To convince the reader that you're a candidate worth interviewing. Here are five rules to help you write a resume that does its job:

- 1. Summarize Your Unique Value
- 2. Communicate with Confidence
- 3. Watch Your Language
- 4. Key in on Keywords
- 5. Keep it Concise

RESUME FORMAT

A resume's "format" is based on the headings you use; the order in which they appear; and the dates of employment for each position. Each format serves a particular purpose.

1. CHRONOLOGICAL

This is widely used resume format emphasizes your career progression by focusing on the dates and job titles you've held, followed by your education. Contrary to its name, a chronological resume actually lists your work history in reverse chronological order, starting with your current or most recent position and going back through each position you've held for the past 10 or 15 years.

2. FUNCTIONAL

It highlights *what* you can do, rather than *when* you did it and for whom. In other words, it defines your value by focusing on skills, not job history. A functional resume calls attention to your specific areas of expertise and lists them under such headings as "Accounting Skills," "Marketing Skills," or "IT Skills."

3. COMBINATION/HYBRID

This format combines elements from both the chronological and functional formats. It balances the focus on your skills and accomplishments with your work history, including employment dates and job titles.

UNIT-III

ORAL COMMUNICATION

Oral communication implies communication through mouth. It includes individuals conversing with each other, be it direct conversation or telephonic conversation. Speeches, presentations, discussions are all forms of oral communication. Oral communication is generally recommended when the communication matter is of temporary kind or where a direct interaction is required. Face to face communication (meetings, lectures, conferences, interviews, etc.) is significant so as to build a rapport and trust.

PRINCIPLES OF EFFECTIVE ORAL COMMUNICATION 1) KNOW YOUR LISTENERS AND ADAPT YOUR MESSAGE TO THEM

_ Think about your audience's demographics—age, gender, occupation, race or ethnicity, religion, cultural heritage, etc.

_ Consider what your audience already knows about your topic, how familiar they are with the terminology, how closely their views match yours, and how committed they are to existing attitudes and beliefs.

_ The best communicators are those who understand their listeners and adjust their message in order to "reach them where they are."

2) SPEAKING IS FUNDAMENTALLY DIFFERENT FROM WRITING BECAUSE LISTENING IS FUNDAMENTALLY DIFFERENT FROM READING

_ A reader chooses when and where to focus attention; a speaker must focus a listener's attention on what he or she is saying at this moment.

_ A reader controls how fast he or she will move through a text; a speaker controls how fast listeners will move through an oral presentation.

_ Readers have the option of going back and re-reading; listeners must grasp material as the speaker presents it.

3) UNDERSTAND YOUR NERVOUSNESS

_ It's normal: 3 out of 4 people say they feel nervous about speaking in public. It's like getting up for an athletic contest: you want to do well, you've prepared, and you're ready to go!

_ Your performance is important, but *it's not the main thing*. The main thing is *sharing your message*—the ideas, feelings, information. It's about learning together.

_ Nobody expects perfection. If you mess up something, just fix it and go on. Your audience is your partner: they want to learn from you; they want you to succeed.

_ Some nervousness is a good thing. Heightened activation can energize your presentation, enhance your alertness and animation, and boost audience engagement.

ADVANTAGES OF ORAL COMMUNICATION

• There is high level of understanding and transparency in oral communication as it is interpersonal.

• There is no element of rigidity in oral communication. There is flexibility for allowing changes in the decisions previously taken.

• The feedback is spontaneous in case of oral communication. Thus, decisions can be made quickly without any delay.

• Oral communication is not only time saving, but it also saves upon money and efforts.

• Oral communication is best in case of problem resolution. The conflicts, disputes and many issues/differences can be put to an end by talking them over.

DISADVANTAGES OF ORAL COMMUNICATION

• Relying only on oral communication may not be sufficient as business communication is formal and very organized.

• Oral communication is less authentic than written communication as they are informal and not as organized as written communication.

• Oral communication is time-saving as far as daily interactions are concerned, but in case of meetings, long speeches consume lot of time and are unproductive at times.

• Oral communications are not easy to maintain and thus they are unsteady.

• There may be misunderstandings as the information is not complete and may lack essentials.

MEDIA OF ORAL COMMUNICATION

• Face to Face communication

• Teleconferencing

• Telephone

In **face to face communication**, we have all the cues available to us: words, facial expression, gestures, body language, tone of voice, room temperature, room noise, and other people in the room that might be present. The message will be more complete and clear when all cues are present.

Videoconferencing and web conferencing are almost as effective as face to face communication. The only cue that is missing in video and web conferencing is the shared presence or surroundings that may give people additional information about the message meaning.

Telephone communication lacks nonverbal cues. When we are having a phone conversation, we don't have facial expressions or body language to help us decode messages, so we must focus on every word being said, and the tone of voice that is being used. We compensate for the absence of nonverbal cues by adding more weight to the words being said and the tone of voice being used.

STYLES OF ORAL COMMUNICATION

One-on-One Conversations (Face to face) Group Discussions Presentations Client Interaction

INTERVIEWS

Interview is another medium of communication. It is formal meeting & discussion with someone on a particular subject. Usually it is a means of getting information it involves:

• Giving information that will help the applicant make up his mind about the company.

• Giving advice that may serve to change the mental or emotional attitude to the interviewee Interviewing the candidates is an important aspect of selection procedure the final selection is partly based on the performance of the candidate in different tests and partly on his performance

in the final interview.

In interview, the candidate has to appear before the interview board or a group of interviewers. The overall personality of the candidate is judged by the interview which may last for 10-20 minutes or even more. Various questions are asked to the candidate in order to the candidate in order to judge his ability, knowledge, capacity and so on.

PURPOSE OF INTERVIEW

Main purpose for interviewee:-

- Communicate information about yourself, your experience and your abilities.
- Seek further information about the position and the organization
- Evaluate the match between your needs and what the job offers.
- Main purpose for interviewer:-
- To gather relevant information about the candidate.

i. Interview preparation:-Interest in and knowledge of the industry, the position and the organization.

ii. Communication skills:-oral presentation skill and the ability to interact with others

iii. Qualifications:- academic, work, volunteer & other experiences

iv. Leadership potential & teamwork:- Demonstrate ability to work with others and to get other to work together

v. Clear & realistic career goals:- future plans and awareness of career paths.

vi. Work ethic:- acceptance of responsibility, ability to keep commitments and attitude of the importance of work.

ART OF INTERVIEWING

According to **S.G. Ginsburg**,"the interviewer's questions must explore viewpoints as well as experiences; they must be as tough as the problems that will face the person who gets the job" As in any other profession, interviewing is an art which demands training and experience. He makes such decision after talking with an applicant for 20-30 minutes. Evaluate the appearance, general manners and relevant experience and training of an

applicant. In today's, most complicated problems of the business is people. Technical processes may be mastered; plans and offices may be built to exacting specifications for performing work with accuracy. But if the human element is disregarded, it automatically leads to trouble so, there is need to encourage the population to grow and aware about the importance of the human factor in industry. Successful interviewing should always be based on fair dealing with people drawing them out analyzing evaluating their strength and weakness.

TYPES OF INTERVIEW

Interviews have been categorized on the basis of various characteristics and qualities.

1. On the basis of an objective: this is done to ascertain weakness in the candidate and making attempts to remove them, or for collecting information. These are of four types:

a) Clinical interview: it is used in medical profession and is done to learn the cause of certain psychological abnormalities. Once the cause is found remedial measures are taken.

b) Selection interview: this is done to select a person on the basis of certain qualities.

c) **Diagnostic interview:** when the objective of the interview is confined to investigating an issue or problem it is called diagnostic interview

d) **Research interview:** interview conducted for the purpose of data collection or hypothesis building in a research is called research interview.

2. On the basis of number of respondents:

a) Group interview: this type of interview is conducted for a group. It can last 1-2 hours and has 10-12 members with one moderator. This method is used in marketing research to collect information on a product type such as detergent.

b) Individual interview: this is an interview where a single person is interviewed.

3. On the basis of from:

a) Structured interview: the pattern of this type of interview is pre-decided. The questions, their wording and their sequence are fixed. An interviewer may be allowed some liberty in asking questions.

b) Unstructured interview: this type of interview is flexible and open, and the questions structure is not pre-decided.

4. On the basis of formation:

a) **Panel interview:** such an interview has the advantage of bringing in the experience of a number of people. Those who may have to work with the candidate get a chance to voice their opinions about him.

b) Two interviews: this type of interview is a small panel interview where the interviewer may decide to adopt opposing notes, the one sympathetic and the other confrontational. The method is supposed to reveal the candidate's probable reaction to pressure in the workplace.

c) **One-to-one interview:** it is the most preferred type of Interview by a candidate. It is more conversational and easy to handle.

INTERVIEW STYLES

Interview styles means,"the degree/level of patterning of the interactions between interviewer and interviewee. The style can be informal or any modification or combining of the two the degree of formality is dependent largely on the relationship of the interviewer and interviewee"

On the basis of practices styles are as follows:-

1) **Direct interview**: - it is face to face observational method. In this; one measures the attitude, knowledge& suitability of interviewee with the help of questions and answers.

2) **Indirect method**:-it is not straight forward questions and answers method. Interviewee is given an opportunity & conducive atmosphere to feel free to talk. Interviewee plays a role of speaking on a particular issue & the interviewer plays a mainly a listening role.

3) **Patterned interview**: - questions to the interviewee are standardized in advance and ask according to pattern.

4) **Depth interview**:- Number of questions on a particular area are put to the interviewee an answer of any one questions does not cover full information A number of follow up questions are put by the interviewer.

5) **Stress interview**:- worry/pressure experienced by the interviewee in a particular circumstances or anxiety caused by the stress created deliberately by the interviewer.

6) **Board interview**:-when a group of people propose to interview respondents, it is called "panel" or "board" interview.

7) **Group interview**: - in a group interview a group of respondents or interviewees are allowed together to interact and exchange each other. Interviewer plays the role of observational and listener to appraise the qualities of respondents in a group.

ESSENTIAL FEATURES OF INTERVIEW

• IT is a face to face interaction between two or more persons.

• It is carried out with a definite objective to either know a person and his capability or views or ideas

• It is a person to person interaction in a controlled setting.

INTERVIEW STRUCTURE

This template offers interview panels a structure that will save you time and ensure you provide information that helps the organisation, the interviewee and the panel.

The structure comprises four main segments:

- Introduction
- Evidence gathering
- Applicant's questions and comments
- Close

For each segment, you need to decide who will do and say what.

GUIDELINES FOR INTERVIEWER

1) He/She Should Be Of A Certain Status, Standing & Experience. They Should Possess The Working Knowledge Of Topics.

2) Skilful interviewing is an art & like all other arts, it requires training & experience it is learnt better by practice than by reading a book

3) Interviewer should not begin the interview without the thorough study of the relevant data contained in the candidate's application

4) An interviewer should know the traits that need to be assessed during interview, namely, intelligence, ability to present ideas, emotional balance, readiness and response

5) Time of interview should not to consume in collecting routine information from records & documents from school institutions etc.

GUIDELINES FOR INTERVIEWEE

1) Be comfortable discussing everything on your resume some interviewers may use it as their only guide for the interview

2) Dress appropriately a positive first impression gets the interview off to a good start.

3) Listen attentively to the interviewer if you do not understand a question ask to have it restated.

4) Get directly to the point ask the interviewer if he would like you to go into great detail before you do so.

5) Do not open yourself to areas of questioning that could pose difficulties for you.

MEETINGS

The word meeting denotes an arrangement to come face to face with advance, plan for a purpose. A meeting is a gathering of two or more people that has been convened for the purpose of achieving a common goal through verbal interaction, such as sharing information or reaching agreement. Meetings may occur face to face or virtually, as mediated by communications technology, such as a telephone conference call, a Skype or a videoconference. Thus, a meeting may be distinguished from other gatherings, such as a chance encounter (not convened), a sports game or a concert and a demonstration (whose common goal is achieved mainly through the number of demonstrator's present, not verbal interaction).

KIND OF MEETINGS

Meetings of Members: These are meetings where the members / shareholders of the company

meet and discuss various matters. Member's meetings are of the following types:-

A. Statutory Meeting:

A public company limited by shares or a guarantee company having share capital is required to hold a statutory meeting. Such a statutory meeting is held only once in the lifetime of the company. Such a meeting must be held within a period of not less than one month or within a period not more than six months from the date on which it is entitled to commence business i.e. it obtains certificate of commencement of business. In a statutory meeting, the following matters only can be discussed:-

- a. Floatation of shares / debentures by the company
- b. Modification to contracts mentioned in the prospectus

B. Annual General Meeting

Must be held by every type of company, public or private, limited by shares or by guarantee, with or without share capital or unlimited company, once a year. Every company must in each year hold an annual general meeting. Not more than 15 months must elapse between two annual

general meetings. However, a company may hold its first annual general meeting within 18 months from the date of its incorporation. In such a case, it need not hold any annual general meeting in the year of its incorporation as well as in the following year only. In the case there is any difficulty in holding any annual general meeting (except the first annual meeting), the Registrar may, for any special reasons shown, grant an extension of time for holding the meeting by a period not exceeding 3 months provided the application for the purpose is made before the due date of the annual general meeting.

C. Extraordinary General Meeting

Such meeting is usually called by the Board of Directors for some urgent business which cannot wait to be decided till the next AGM. Every business transacted at such a meeting is special business. An explanatory statement of the special business must also accompany the notice calling the meeting. The notice must also give the nature and extent of the interest of the directors or manager in the special business, as also the extent of the shareholding interest in the company of every such person. In case approval of any document has to be done by the members at the meeting, the notice must also state that the document would be available for inspection at the Registered Office of the company during the specified dates and timings.

D. Class Meeting

Class meetings are meetings which are held by holders of a particular class of shares, e.g., preference shareholders. Such meetings are normally called when it is proposed to vary the rights of that particular class of shares. At such meetings, these members discuss the pros and cons of the proposal and vote accordingly. (See provisions on variations of shareholder's rights). Class meetings are held to pass resolution which will bind only the members of the class concerned, and only members of that class can attend and vote.

II. Meetings of the Board of Directors

- Meeting of the Board of Directors
- Meeting of a Committee of the Board

III. Other MeetingsA. Meeting of debenture holdersB. Meeting of creditors

Advantages of meetings/ committees

Information Sharing
Encourages Teamwork
Disadvantage of meetings/ committees
Time
Lack of Leader

PLANNING AND ORGANIZATION OF MEETINGS

Effective meeting planning and organization guidelines

Good meeting planning is a necessary prerequisite for any effective business meeting. Here are the key considerations and practical tips to guide you through the meeting organization process. The first and foremost question to ask before you start planning any meeting is "What are the desired outcomes from that meeting?" The second critical question is "What is the best tool to reach that outcome?" Remember that a meeting is just one of the tools of interpersonal communication. The next important question of meeting planning and organization is "Who are the right people to be at the meeting?" Decide on the appropriate meeting format. Given the meeting purpose and participants, would it work better as formal or informal, Public or private? More like a conference with a number of presentations?

Who will chair the meeting? The chair person, whether it is you or somebody else, needs enough authority and ability to keep meeting running smoothly. An effective chairperson keeps the meeting on track, maintains constructive and positive atmosphere, ensures that nobody hijacks or sabotages the meeting, and helps all participants contribute most effectively. The next meeting planning step is to decide on time and place. Of course, you want to schedule time when all the participants, or at least the majority, are available.

PROJECT PRESENTATIONS

A presentation means speaking on a topic before a select audience. It is a form of oral communication with an audience on some formal occasion. The formal presentation of information is divided into two broad categories: Presentation Skills and Personal Presentation. These two aspects are interwoven and can be described as the preparation, presentation and practice of verbal and non-verbal communication. This article is an overview of how to prepare and structure a presentation, and how to manage notes and/or illustrations at any speaking event.

Many people feel terrified when asked to make their first public talk. Some of these initial fears can be reduced by good preparation which will also lay the groundwork for making an effective presentation.

A Presentation Is...

A presentation is a means of communication which can be adapted to various speaking situations, such as talking to a group, addressing a meeting or briefing a team. To be effective, step-by-step preparation and the method and means of presenting the information should be carefully considered. A presentation concerns getting a message across to the listeners and may often contain a 'persuasive' element.

ADVANTAGES

- **1.** Topic will become interesting.
- 2. Easy to understand
- 3. Diagrams are inbuilt in presentation.

DISADVANTAGES

- 1. Technical knowledge required
- 2. It is a costly process of teaching

EXECUTIVE SUMMARY

An executive summary, sometimes known as a management summary, is a short document or section of a document, produced for business purposes, that summarizes a longer report or proposal or a group of related reports in such a way that readers can rapidly become acquainted with a large body of material without having to read it all. It will usually contain a brief statement of the problem or proposal covered in the major document(s), background information, concise analysis and main conclusions. It is intended as an aid to decision making by managers and has been described as possibly the most important part of a business plan. They must be short and to the point.

ACCREDITED

CHARTS

A **chart** is a graphical representation of data, in which "the data is represented by symbols, such as bars in a bar chart, lines in a line chart, or slices in a pie chart". A chart can represent tabular numeric data, functions or some kinds of qualitative structure and provides different information. The term "chart" as a graphical representation of data has multiple meanings:

• A data chart is a type of diagram or graph that organizes and represents a set of numerical or qualitative data.

• Maps that are adorned with extra information for some specific purpose are often known as charts, such as a nautical chart or aeronautical.

Common charts





Histogram Bar chart Pie chart Line chart DISTRIBUTION OF TIME

Equal distribution of time while presenting a presentation so that there will be time for the query session of the audience after summing up or conclusion of the presentation.

- _ Presentation
- _ Questions & answers
- _ Summing up

VISUAL PRESENTATION

You should only use visual aids if they are necessary to maintain interest and assist comprehension in your presentation. Do not use visual aids just to demonstrate your technological competence - doing so may compromise the main point of your presentation - getting your messages across clearly and concisely.

If visual aids are used well they will enhance a presentation by adding impact and strengthening audience involvement, yet if they are managed badly they can ruin a presentation. Most visual aids will need advance preparation and should be operated with efficiency. This page gives details of the following common visual aids:

9001:2015 & 140

- · Whiteboards and Interactive Whiteboards
- Flip chart
- Over-head projector (OHP)
- Slides
- Video
- PowerPoint or other presentation software
- Handouts

GUIDELINES FOR USING VISUAL AIDS

- _ Prepare Visual Aids in Advance
- _ Keep Visual Aids Simple
- _ Make Sure Visual Aids Are Large Enough
- _ Use Fonts That Are Easy to Read

- _ Use a Limited Number of Fonts
- _ Use Color Effectively

GUIDELINES FOR PRESENTING VISUAL AIDS

- Avoid Using the Chalkboard
- _ Display Visual Aids Where Listeners Can See Them
- _ Avoid Passing Visual Aids Among the Audience
- _ Display Visual Aids Only While Discussing Them
- _ Explain Visual Aids Clearly and Concisely Talk to Your Audience, Not to Your Visual Aid
- Practice with Your Visual Aids

ELECTRONIC MEDIA (POWER-POINT PRESENTATION)

Electronic media are media that use electronics or electromechanical energy for the end-user (audience) to access the content. This is in contrast to static media (mainly print media), which today are most often created electronically, but don't require electronics to be accessed by the end-user in the printed form. The primary electronic media sources familiar to the general public are better known as video recordings, audio recordings, multimedia presentations, slide presentations, CD-ROM and online content. Most new media are in the form of digital media. However, electronic media may be in either analog or digital format.

LISTENING SKILL

Listening may be defined as the process of hearing, understanding & interpreting the spoken words.

Good listening for improved communications

Effective listening requires both deliberate efforts and a keen mind. Effective listeners appreciate flow of new ideas and information. Organizations that follow the principles of effective listening are always informed timely, updated with the changes and implementations, and are always out of crisis situation. Effective listening promotes organizational relationships, encourages product delivery and innovation, as well as helps organization to deal with the diversity in employees and customers it serves.

01:2015 & 14001:2015

Effective Listening Skills

- 1. Discover your interests' field.
- 2. Grasp and understand the matter/content.
- 3. Remain calm. Do not lose your temper. Anger hampers and inhibits communication.
- Angry people jam their minds to the words of others.
- 4. Be open to accept new ideas and information.
- 5. Jot down and take a note of important points.
- 6. Work upon listening. Analyze and evaluate the speech in spare time.
- 7. Rephrase and summarize the speaker's ideas.

8. Keep on asking questions. This demonstrates that how well you understand the speaker's ideas and also that you are listening.

9. Avoid distractions.

ART OF LISTENING

The ability and need to communicate touches every area of our lives. Everything we do in life requires communication with others. Just try to not communicate at work for a day or in your business transactions and see what happens. Refuse to communicate in your personal relationships and see what kind of interesting results you'll create.

~~ ~

Much of communication theory focuses on how to speak to others and how to convey your message. But, communication is really a two-way process. It is an activity, not a one-time event. The listener's role is as central to the communication process as the speaker's role. Real communication and connection occur when the speaker and listener participate in the process. Exercise active listening skills. Try asking more questions. If you need clarification ask the speaker to say more, give an example or to explain further. Give feedback or paraphrase what you've heard: "Are you saying such and such? What I heard you say is this. Is this what you meant?" Try nodding your head to show interest. Or ask a question of interest to demonstrate that you are really listening to what is being said. Add the occasional "uh-huh". Try making eye contact with the speaker. Even though you are sitting and listening quietly, this may not be enough for the speaker to feel that s/he is truly being understood.

Become aware of your personal filters and triggers. Each of us is a product of our upbringing, culture, life experiences and anything and everything that makes us unique as human beings. Our uniqueness can sometimes be an obstacle to being an effective listener. As you listen, try to remain open to what you are hearing and withhold evaluation or judgment. Become aware of what your triggers are in the communication process and what shuts your listening down.

NATURE OF LISTENING

Besides the division of the skills as 'receptive' and 'productive', another subdivision focuses on 'one way reception' and 'interactive reception' in this age of active learning. Reading and writing are one-way skills where learners don't get direct feedback. But in speaking and listening, learners may have their understanding and reproduction checked instantly. Thus active and self learning takes place Moreover, there is a traditional labelling for reading and listening as= "passive" skills. But linguists believe that a listener is involved in guessing, anticipating, checking, interpreting, interacting and organizing by associating and accommodating their prior knowledge of meaning and form.

PROCESS OF LISTENING

1- Sound Recognition

First of all the listener hears sound. He tries to percept the words or speech. He discriminate the herd sounds. He also recognizes the tone of the speaker. Technically speaking he recognizes stress pattern of the speech along with intonation pauses and breaks in the speech.

2- Word Recognition

The second step in the listening process is the word recognition. The listener recognizes the sound pattern as words. Now sound becomes meaningful and takes the form of words. He also locates the word in the word list; regain lexical, grammatical and semantic information about the word.

3- Parsing or analysis

Thirdly the listener learns sentence processing or parsing. He detects sentence constituents or parts of sentence and builds a structure or frame of the sentence.

4- Construction of Meanings

The fourth step in the listening process is to construct the literal meaning of the sentence. If the sentence has some ambiguous or difficult words he tries to select the relevant meanings of the sentence.

5- Creates Short Term memory

When the listener become efficient in making sense of the sentence he keeps back the information he has attained from the sentence and hold it the in short-term memory.

6- Discourse Devices

Afterwards he becomes familiar with interrelated devices in talk.

7- Infer Implied Meanings

When he has learned the discourse devices being used by the speaker he tries to infer the implied meaning and intention of the speaker. It is an elevated step of listening where the listener is near to reach the perfection of the listening.

8- Make Prediction

After a lot of such purposeful and effective listening the listener is able to predict. He can easily tell the next contour of the discourse. He predicts forehand what is to be said.

9- Deciding Response

The last step in the process is the selection of the response on the part of listener. After completely understanding the speech he decides himself how to respond to this speech talk or discourse.

TYPES OF LISTENING Passive listening

- _ No active participation be the listener,
- _ Not absorbed in the memory for future use
- _ Constrained to physiological and psychological factors
- Lack of interest, fatigue, ill health, disregard
- Lack of wavelength
- _ And it is perceived as it has been understood.

Selective Listening

- _ Listening what they want to _ Attention is not focused
- _ message is not thoroughly processed

_ Not in a position to concentrate, thinks that the speaker is not well informed or he thinks that he is better.

Active listening

- _ Concrete effort by both the ends
- _ Actively spoken and heard
- _ Retained for future use
- _ Good participation be both the sides

IMPORTANCE OF LISTENING

Listening skills are cultivated by

- _ Stake holders need to be at good level at all times and grow
- _ Businesses thrive on listening
- _ Knowledge up gradation
- Listen to customer complaints

DIFFERENCES BETWEEN LISTENING AND HEARING

- **_ Listening** is active; **hearing** is passive
- _ Listening is emotional; hearing is passionless
- Listening is deep in experience; hearing is never beyond glossy exteriors
- _ Listening is timing; hearing is passive listening at times
- _ Listening is caring; hearing is passive listening

PRINCIPLES OF GOOD LISTENING

1-Maintain eye contact. This is the first and foremost principle of good listening. It helps the listener to concentrate on the speaker's words. It saves him from distracting his attention from

the speaker.

2-Focus on content, not delivery. A Good listener never focus on the delivery, he always become all ears for the contents. He does not allow his attention to go astray from the words of

the speaker.

3-Avoid emotional involvement. When you are too emotionally involved in listening, you tend

to hear what you want to hear--not what is actually being said. Try to remain objective and open minded.

4-Avoid diversion and distraction. A good and active listener focus on the contents he is listening not the nearby voices or sounds.

5-Consider the listening an inspiring intellectual job. Listening to an educational lecture is not

a passive work. It is always edifying and useful for listener that is why he should take listening

as an inspiring job.

6- Ask questions in your mind. Active listening keeps you alert. Always ask yourself some questions about the contents of the speaker. Keep evaluating his ideas during listening.

7- Keep your mind focused on the Speech.

Use the rate of speech and your rate of thought to anticipate what the speaker will say next. In this way you will be able to keep your mind from straying.

8- Try to infer the main idea .The main ideas are the most important points the speaker wants

to convey. Always try to infer the theme of the contents of the speaker.

9- Express that you are listening and understanding. Try to look at the speaker and express

your attention by nodding now and then. Appropriate feedback at different points with smile, frown or laugh shows that you are actively involved in listening.

10-Remain Objective .Avoid emotional involvement because it will deviate your attention from

:2015 & 14001:2015

the real contents of the speech.

BARRIERS IN LISTENING

- _ Physical barriers
- _ Defective microphones, speakers
- _ Voices and noises
- _ Speaker is close to microphone
- _ Interruptions
- _ Transmission failures

People related barriers

- 1. Shrill voice
- 2. Accent (rapid)
- 3. Receiver doesn't thinks that the speaker is well informed.
- 4. No focus of listener
- 5. Not having adequate authority(thought by listener)

NEGOTIATION SKILLS

Negotiation skills is defined as "a process by which two parties interact to resolve a conflict jointly" According to **J.L.Graham**, "negotiation is a face to face decision –making process between parties concerning a specific product"

FEATURES OF NEGOTIATION SKILLS

- Minimum two parties present
- Both parties have pre determined goals
- There is an outcome
- Both parties believe that the outcome of the negotiation may be satisfactory
- Parties understand the purpose of negotiation

FACTORS THAT CAN INFLUENCE NEGOTIATION

1) Place:- familiarity with surrounding help in boosting confidence

2) Time:- time should be adequate for smooth exchange of ideas & securing agreement before it is to late

3) Attitude:- Attitude of both parties should be positive. That is willingness to make an agreement or deal

4) Subjective factors: - Like relation of two parties involved, status difference, information and expertise.

What skills do we need to negotiate?

- Preparation & planning skills
- Knowledge of the subject
- Ability to think clearly & rapidly under pressure & uncertainty
- Ability to express thoughts to verbally listening skills
- Patience
- General problem solving & analytical skills

NEGOTIATION PROCESS

1) PREPARING:-preparing for a negotiation involves preparing background information before negotiating an outcome you want. it involves:-

• Setting objectives

- Determine what you'll do if the negotiation fails
- Listening ranking and valuing issues
- Analyzing other parties
- Researching information
- Researching the negotiation

2) Proposal:- proposal regarding about the issue to negotiate.

- 3) Planning:-
- Planning & presenting your agenda
- Planning your critical first offer
- Identifying strategies to help overcome negotiating challenges
- 4) Discussions:-
- Discuss the agenda calmly
- Listen the view points of both the parties
- 5) Bargaining:- two parties bargains with each other to reach on the decision
- 6) Agreement:- preparation of the two parties which indicate that both the parties are

satisfied with the decision taken by the negotiation

7) Implementation:-implement the decision taken in a negotiation process.

STRATEGIES TO IMPROVE ORAL COMMUNICATION

To improve oral communication

- Think before you speak.
- Use direct and concise language.
- Vary your vocal tone.
- Pronounce words completely and correctly.
- Master your nonverbal communication skills.

STRATEGIES TO IMPROVE PRESENTATION

To improve Presentation

- Be well versed in the subject matter you are presenting.
- Know who your audience is.
- Stay relaxed but focused.
- Keep the presentation content based.
- Speak as if you were conversing with one person in relaxed plain language.
- Visual aids can be a plus. Just don't put your audience to sleep.
- Never read from a script. It's a good way to get lost.
- If you are new at this, never look directly at the audience, but slightly over their heads.

5 & 14001:201

- Make sure your appearance is appropriate. First impressions are critical.
- When you speak, do not shout, or whisper. Know the acoustics of the room.

STRATEGIES TO IMPROVE SPEAKING AND LISTENING SKILLS

To improve speaking skills

- _ Using minimal responses
- _ Recognizing scripts
- _ Using language to talk about language
- To improve listening skills
- _ Desire to listen
- _ Resistance distractions
- _ Focus on message
- _ Delay evaluation & Premature conclusions
- _ Taking notes

UNIT – 4 SOFT SKILLS

Non- Verbal Communication

- _ Communication without using words.
- _ Simple & limited
- _ used: maps, charts, graphs etc
- _ Methods may be visual or auditory

KINESICS (BODY LANGUAGE)

_ Study of the role of body movements such as winking, shrugging, etc in communication.

ARMA

Elements

- _ Personal appearance
- _ Facial expression
- _ Head
- _ Posture / body position
- _ Eye contact
- _ Gestures
- _ Body shape
- _ Smell & touch

Make effective use of it

- _ Mind the body talk
- _ Select the proper sitting posture
- _ Be careful with the handshake
- _Est. good eye contact
- _ Maintain your self esteem

PROXEMICS (SPACE LANGUAGE)

_ Subject that deals with the way people use physical space to communicate

_ Zone / territory constructed - doesn't allow to be invaded during communication unless

relation b/w speaker & listener is intimate.

Acc to Edward T Hall – 4 kinds of distance

- _ Intimate physical contact to 18 inches
- _ Personal 18 inches 4 feet
- $_$ Social -4 feet -12 feet
- _ Public 12feet range of eyesight & hearing

PARALANGUAGE (LIKE LANGUAGE)

_ Para refers to like

_ Non verbal factors like tone of voice, emphasis given, the breaks in the sentences, the speed of delivery, the degree of loudness or softness, & the pitch of voice, which affect the spoken words.
Why it is used?

- _ Completion of message
- Personal & educational background of the sender
- _ Regional or national background
- _ Mental state of the communicator
- _ Learning exercise

- Limitations _ Lack of reliability
- _ Chance of misguiding / misleading
- Requires attention for understanding
- _ Lack of uniformity
- _ Requires patience

INTER PERSONAL SKILLS

_ Interpersonal skills are the life skills we use every day to communicate and interact with other people, individually and in groups. It includes not only how we communicate with others, but also our confidence and our ability to listen and understand. Problem solving, decision making and personal stress management are also considered interpersonal skills.

3 Key Interpersonal Skills You Need At Work Today

- **_ Empower** joint problem solving
- **Encourage** better connections with others
- _ Engage people to want to help you

CORPORATE COMMUNICATION SKILLS

- _ Builds strong business relationship
- _ Internal and external coordination
- _ Build & maintain brand image of company
- Gives competitive adv to company.

BUSINESS ETIOUETTES

_ Etiquettes – means conventional rules of social behavior. Rules generally unwritten & are passed on from 1 generation to another 001-2015

FUNDAMENTAL RULES OF BUSINESS ETIQUETTES

- _ Impact
- _ I- integrity
- _ M- manners
- _ P- personality
- A- appearance
- _ C- consideration
- T- tact

LANGUAGE SKILLS IMPROVING COMMAND IN ENGLISH

The importance of the English language cannot be overemphasized. Comfort with English is almost a prerequisite for success in the world today. Regardless of the industry, proficiency in English is an important factor in both hiring and promotion decisions. A lot of us have studied English in school and are fairly comfortable with reading and writing. However, we hesitate while speaking because we feel that we lack the fluency and may make grammatical mistakes. We are afraid of speaking English in formal situations and we are quick to switch to our native language once we are in the company of our family and friends. Here are English some tips for success in achieving proficiency and fluency in English:

1. **Do not hesitate.** Talk to whoever you can. Decide among your circle of friends that you will only talk in English with each other. This way you can get rid of hesitation and also have your friends correct you when you are wrong.

2. **Start a conversation with strangers in English**. Since you do not know them personally, you will feel less conscious about what they would feel about you.

3. Maintaining a diary to record the events of your day is a great way to practice your writing skills. Take your time to use new words and phrases when you write in your diary.

4. **Read the newspaper**. Read it aloud when you can. Concentrate on each word. Note down the words you don't understand and learn their meanings. Try to use these words in your own sentences.

5. Watch English movies and English shows on television. Initially, you can read the subtitles to follow the conversation. As you practice more, you will realize that you are able to follow the conversation without needing to read the sub-titles.

6. Set aside an hour every day to watch English news channels. This is one of the most effective ways of improving your comprehension.

IMPROVING VOCABULARY

1. **Make reading the newspaper a daily ritual.** You may be comfortable reading a particular section but make an effort to read different articles on every page. The editorial page is highly recommended not only for vocabulary but also for structuring and presenting thought.

2. Make it a habit to read a new book every week. It is not surprising that those who read a lot develop a good vocabulary. You can consider becoming a member of the local library. Make a list of words that are new to you and look up their meanings in the dictionary.

3. Watching English movies and television shows is important for improving English and learning new English words. The best part about watching English videos is that you can learn the correct pronunciation as well.

4. Use vocabulary cards. Vocabulary cards are used by students who are trying to learn many words in a short period of time. You can make your own cards by writing the word on one side and the meaning on the other side of a square piece of paper. It is a convenient tool to learn new words in your free time

5. Use the internet. The internet is an unlimited resource for reading material. Pick up a topic of your choice and search for articles about it. You will come across plenty of material to read, which you might find interesting, and importantly, will also introduce you to new words. Be sure to look them up in a dictionary.

6. **Don't forget the new words.** The best way to ensure that you never forget the new words you learn is to start using them in your day to day conversation. Do not try to force them into a conversation but do use them if you think they are appropriate.

7. Learn pronunciation. Most dictionaries provide us with pronunciations of words using phonetic symbols. It is important to learn the sounds that correspond to these phonetic symbols, in order to become comfortable pronouncing new words.

CHOICE OF WORDS

In order to communicate our message to our listeners, we need to choose the right words that clearly express out thoughts. Very often, in both informal conversation and public speaking, we make statements that are not very clear. For example:

• Unclear: Let's go to a place where they sell those things we need for the office.

• Clear: Let's go to the bookstore that sells books and school supplies.

To develop the skill in choosing the right words, we can use:

1. Simple words – our ideas will be easier to understand if we use simple words and phrases rather than the complex or difficult ones. The English language gives us a variety of words that can be used to express various shades of meaning.

Simple Words Difficult Words Dislike Abhor Increase Abound Betray Circumvent Belief Credence

2. Precise words – these express our thoughts and feelings accurately. The use of vague words confuses the listeners and does not clearly express our intended meaning. To choose precise words that clearly and accurately communicate our point, we need to increase our

speaking vocabulary. This will give us choices and will make us sensitive to differences in meaning.

- Vague: We had a bad meeting yesterday.
- Precise: We had a disorganized meeting yesterday.
- Vague: My supervisor looked through the monthly report.
- Precise: My supervisor examined the monthly report.

3. Specific words – these identify items within a category while general words refer to an entire category. Almost any concept can be made more general or more specific by the words we use to express it. Specific words help our listeners to form a picture in their minds of the exact images we want them to see.

• General Word: The purchasing officer bought a lot of THINGS at the bookstore.

• Specific Word: The purchasing officer bought pencils, brown envelopes, white board pens, and bond paper at the bookstore.

4. Concrete words – these name things that can be perceived by one or more of the five senses.

In contrast, abstract words name ideas or beliefs that cannot be perceived by the senses.

COMMON ERRORS

COMMON ERRORS IN BUSINESS WRITING

Sentence structure

- Incomplete sentences e.g., because he wasn't at work that day.
- Run-on sentences e.g., the meeting was adjourned we all left right after that.
- Comma splices e.g., Julie presented the layout to the clients, they liked it.

Spelling and punctuation

- Errors using capitals e.g., She works in the ford building.
- Incorrect punctuation e.g., doesn't use the photocopier.
- Spelling mistakes e.g., Punctuality is very important.

Grammar

- Incorrect use of verb forms e.g., I seen the client yesterday.
- Incorrect subject-verb agreement e.g., every manager and employee in the company agree with the decision.
- Improper use of transitions e.g., the meeting was long although the team talked a lot.

• Pronoun references that are unclear e.g., Sam went to the manager's office to pick up his report.

• Improper use of articles e.g., they discussed advantages and disadvantages of the proposal.

• Incorrect use of subject/object pronouns e.g., the boss gave a raise to Julia and I.

Choice

• Incorrect use of similar words e.g., their very concerned about how the layoffs will affect their morale.

Improper tone, style or level of formality

• Use of informal style or casual language e.g., He's really ticked off that people continue that is not appropriate in business writing to come late. E.g. He left

COMMON ERRORS WITH: VERBS

1. Verb Tense Error

1. Present tense is used to indicate that an action is occurring at the present time. In the first sentence below, the verb is in present perfect tense. However, the action of the sentence is occurring in the present. You can correct this error by changing the verb to present tense. For example, Error: This example has explained [PRP] verb tense problems.

Correct: This example explains [PR] verb tense problems.

2. Present perfect tense is used to indicate that an action occurring in the present began in the past. In the first sentence below, the verb is in past tense. However, the action of the sentence is not finished: it occurred in the past and is still occurring in the present. You can correct this error by changing the verb to present perfect tense. For example, Error: Environmental pollution alarmed [PA] many people.

Correct: Environmental pollution has alarmed [PRP] many people.

3. Past tense is used to indicate that an action was completed at some unspecified past time. In the first sentence below, the verb is in present tense. However, the adverbial phase *two weeks ago* indicates that the action of the sentence occurred in the past. You can correct this error by changing the verb to past tense. For example,

Error: Jane buys [PR] a new car two weeks ago.

Correct: Jane bought [PA] a new car two weeks ago.

4. Past perfect tense is used to indicate that an action was completed at some specific past time. In the first sentence below, the verb *finished* is in past tense. However, the preposition *after* indicates that the action was already completed by the time another action occurred (*Jim had dessert*). In other words, the time relationship expressed by the tense of the verb *finish* depends on two actions. You can correct this error by changing the verb to past perfect tense. For example,

Error: Jim had [PA] dessert after he finished [PA] the three main courses. Correct: Jim had [PA] dessert after he had finished [PAP] the three main courses.

5. Future tense is used to indicate that an action will be completed at some unspecified future time. In the first sentence below, the verb is in present tense. However, the prepositional phase *in two years* indicates that the action of the sentence will occur in the future. You can correct this error by changing the verb to future tense. For example, Error: In two years, Jane buys [PR] a new car.

Correct: In two years, Jane will buy [FU] a new car.

6. Future perfect tense is used to indicate that an action will be completed at some specific future time. In the first sentence below, the verb *will elect* is in future tense. However, that action will occur before the time another action takes place (*by January 2005*). In other words, the time relationship expressed by the tense of the verb *elect* depends on two actions. You can correct this error by changing the verb to future perfect tense. For example,

Error: By January 2005, we will elect [FU] a new president.

Correct: By January 2005, we will have elected [FUP] a new president.

2. Verb Tense Shift

Verb tenses should be used consistently. A verb tense shift error occurs when verb tenses shift arbitrarily. In the first sentence below, the tenses shift from present to past although the actions occurred at the same time. Specifically, the verb *study* is in present tense but the verb *reviewed* is in past tense. You can correct this error by making the tenses consistent. For example,

Error: For the test, they study [PR] the book and reviewed [PA] their notes.

Correct: For the test, they study [PR] the book and review [PR] their notes.

Correct: For the test, they studied [PA] the book and reviewed [PA] their notes.

In the first sentence below, there is a time relationship between the two actions. In other words, the second action, *investigates*, occurs before the first action, *buys*. However, both verbs are in present tense. The first verb may remain in present tense, as it will occur as soon as the second action is complete (as is indicated by the preposition *before*). However, the second verb should be in future tense, as it has not occurred yet.

3. Verb Mood Problems

As you learned in Unit 2, verbs have three moods: indicative, imperative, and subjunctive. Errors with verb mood typically occur with the subjunctive. The subjunctive mood is used to express wishes, requests, or conditions contrary to fact. It is also used in dependent clauses that contain an order or a recommendation.

1. The present subjunctive form of to be is be for all persons.

In the first sentence below, the dependent clause (*that*...) contains an order; therefore, the mood should be subjunctive. Therefore, the verb should be *been told* rather than *are told*. For example,

Error: It is necessary that they are told what happened. Correct: It is necessary that they be told what happened.

2. The present subjunctive forms for verbs other than *to be* is identical to the base form of the verb for all persons. In the first sentence below, the dependent clause (*that*...) contains an order; therefore, the mood should be subjunctive. Therefore, the verb should be *known* rather than *knows*. For example, Error: It is necessary that she knows what happened. Correct: It is necessary that she know what happened.

4. Verb Mood Shift

Verb moods should be used consistently. A verb mood shift error occurs when verb moods shift arbitrarily. In the first sentence below, the mood shifts from indicative to imperative. You can correct this error by making the moods consistent. As the author is reporting on the manual, the mood of the second sentence should be indicative rather than imperative. For example, Error: The manual gives explicit directions. [IND] Open the application first. [IMP] Correct: The manual gives explicit directions. [IND] It recommends that you open the application first.

ADJECTIVES

An adjective modifies a noun or a pronoun by describing, identifying, or quantifying words. An adjective often precedes the noun or the pronoun which it modifies. 'Big', 'boring', 'purple', 'quick' and 'obvious' are all adjectives

ADVERBS

An **adverb** is a word that modifies verbs, adjectives and other adverbs. It is a word that describes or modifies a verb. Ex: carefully, quickly, wisely. Also sometimes modifies an adjective. ("She was very tall." 'Very' is an adverb modifying 'tall,' which in turn an adjective is modifying

'she'?) Adverbs usually, but not always, end in "-ly".

Common Adjective and Adverb Errors

Some writers make mistakes in their choices of adjectives and adverbs. They may use an Adjective where an adverb is correct or vice versa. These adjective and adverb mistakes are easy to make because the incorrect versions are often used in informal speech. For example, you might say to a friend, "That's real sad," but "real" is incorrect. You should use the adverb "really" because it modifies "sad," which is an adjective. Be especially careful with these adjectives - adverb pairs:

Adjective	Adverb
good	well
bad	badly
real	really
slow	slowly
terrible	terribly
quick	Quickly

PRONOUNS

Pronouns are words that substitute a noun or another pronoun. Examples of pronouns are *he*,

she, who, themselves...
In the example:
Mike likes his daughter.
Mike and his daughter can be replaced by him and her:
He likes her

PRONOUN AGREEMENT ERRORS

Pronoun errors--pronouns that don't represent their antecedents correctly--are among the most common

Errors of grammar in college level writing. There are three ways in which pronouns may be incorrect, and they correspond to the three ways in which pronouns are identified: in number, in person, and in case.

Pronoun Case Errors

Conversational English makes many allowances for pronoun case errors because, in the spontaneous flow of dialogue and everyday interaction, no one stops to demand correct grammar. Writing, on the other hand, is not spontaneous; it can be reread or read slowly, and these mistakes are more glaring indications of the writer's poor command of English. The most common pronoun case errors confuse the subject case with object case, or misuse the reflexive

case.

1. Subject-Object Case Error

I made him promise to keep this delicate matter between him and I. Me and Terrence camped in Yosemite over the weekend. Yours and your friend's ties are exactly alike.

2. Reflexive Case Error

This application should be filled out by yourself only. Samantha, Peter and myself are now roommates. Her friends respect herself more than she does.

COMMON MISTAKES IN TENSES

1. Incorrect: It is raining since yesterday. Correct: It has been raining since yesterday.

2. Incorrect: I reached the station before the train had arrived. Correct: I had reached the station before the train arrived.

3. Incorrect: Yesterday we had seen the panther. Correct: **Yesterday we saw the panther.**

4. Incorrect: I had been playing a match.Correct: I had been playing a match yesterday, so couldn't go out with my friends.

CONJUNCTIONS

To conjugate a verb is to state the form the verb takes for each person. For example, to conjugate the verb "to have" (in the present tense) you say "I have, you have, he/she/it has, we have, y'all have, they have."

A **conjunction** is merely a connecting word. It has no other function in the sentence. In most languages of European origin, clauses are joined together by conjunctions in similar ways. However, students who speak non-European-type languages often experience some problems in using English conjunctions correctly.

Correct use of some conjunctions

Unless

Unless means **if not**, so it will be superfluous to introduce another **not** into the following clause.

• Unless you give the keys of the safe, you will be shot.

OR

• If you do not give the keys of the safe, you will be shot.

(NOT unless you do not give the keys of the safe, you will be shot.)

Lest

Lest means that... not, and, therefore, it will be wrong to add another not in the following clause. Moreover, it should be noted that the only auxiliary verb that can be used after lest it should.

• Take care lest you fall. (NOT Take care lest you do not fall.)

• Take care lest you should fall.

• Book your tickets early **lest you should** miss this chance.

As...as...

As is used in comparisons of equality'. Than and that are not used in this way.

• My hands were **as** cold **as** ice.

• Your eyes are **the same** colour **as** mine.

PUNCTUATIONS

Improving Punctuation

Why do we need punctuation? Are commas and colons required only to give new writers a hard time? Punctuation allows us to clarify the meaning of words when voice or "body language" cues are removed. Punctuation tells the reader how to make sense of words alone. Poorly or wrongly used punctuation contributes to awkward writing and reader confusion. For example, how would you speak these words:

• I'm sorry I still love you (Notice the difference punctuation makes)

• I'm sorry. I still love you.

• I'm sorry I still love you!

When we speak, we emphasize certain words to make our meaning clear. When we write, however, we confuse the reader if we don't punctuate well.

A. END PUNCTUATION (Period, question mark, exclamation points, and ellipsis)

Periods (.): Ordinary sentences end with periods.

Question Marks (?): Don't forget to use a question mark when your sentence asks a question. Do you understand?

Exclamation Points (!): Reserve exclamation points for direct orders and commands, and for genuine exclamations. Use only one at a time!

B. **JOINING AND LISTING PUNCTUATION** (Comma, Semicolon, Colon, and Dash) When we start joining sentences--or parts of sentences--together, we need punctuation to insure smooth splicing.

SIMPLE SENTENCE + SIMPLE SENTENCE = COMPOUND SENTENCE

Betty loves bologna. + Bob won't eat meat. Betty loves bologna, but Bob won't eat any meat.

PREFIX

A prefix is a group of letters that is added to the beginning of a word to make a new word. The spelling rule for prefixes is quite straightforward. Put simply, the spelling of the root word *does not change* when a prefix is added to it. However, errors do occur. Let's take a look at some of the most common spelling mistakes that result from the use of prefixes:

Formation	NAAC ACC	Correct
Incorrect	DUMPAN PAN	a R S for LP I I for LP
dis + similar	<i>Dis</i> similar	disimiliar
im + mature	Immature	imature
in + numerable	Innumerable	inumerable
ir + rational	Irrational	irational
mis + spell	Misspell	mispell
under + rated	Underrated	underated
Formation	Correct	Incorrect
dis + illusion	Disillusion	Disillusion
dis + appearance	Disappearance	Disappearance
in + oculate	Inoculate	Inoculate
un + exceptional	Unexceptional	Unexceptional

There are no exceptions to the prefix rule, so the simplest thing to keep in mind is that when a prefix is added to a word, the spelling of a base word should *not* change. However, if you are ever uncertain about the spelling of a word, it is best to consult the dictionary.

SUFFIX

A suffix is a short word or "word fragment" that sits at the end of a word, and modifies the word's meaning; similar to a prefix at the front of a word. Here are some of the more common suffixes. Different suffixes apply to different classes of words, as you'll see below: Common Suffixes — verbs to nouns

SUFFIX	VERB – NOUN
-AL condition, quality	Arrive - arrival,
	Approve – approval
-ANCE / ENCE action,	Attend – attendance
state, condition or quality	Accept - Acceptance
ATION / TION action or	Educate – education
resulting state	Inform - information
-SION action or resulting	confuse – confusion
State	decide – decision

-URE action or resulting	depart – departure
State	erase – erasure
-MENT state, act, condition	agree- agreement
	pay – payment
-AGE action, state, process	break – breakage
	post – postage

IDIOMATIC USE OF PREPOSITIONS

The use of prepositions can vary greatly between languages, even between two variants of a single language such as American English and British English. When a word phrase or expression is peculiar to a given language and cannot be understood from the individual meanings of its elements, it is called 'idiomatic.' Because idioms (idiomatic word patterns) cannot be deduced from a general knowledge of the words and their meaning, we need to simply memorize them. For native speakers of the language, this process usually happens unconsciously: certain word patterns just sound right. Non-native speakers may have to work at mastering idioms. Here are some common prepositional idioms of American English:

• abide	• by a rule
• abide	• in a place or state
• accords	• with
• according	• to
• Accuse	• Of a crime
• Adapt	• From a source
• Adapt	• to a situation
• Afraid	• Of

SENTENCES AND PARAGRAPH CONSTRUCTION

I Developing Effective Sentence:

(A)Emphasis on Short Sentences

- You can write short, simple sentences in two basic ways:
- 1) By limiting sentence content and
- 2) By using words economically

(B) Determine Emphasis in Sentence Design

(C) Give the Sentence Unity

Good sentences have unity. For a sentence to have unity, all parts of a sentence should concern

one thought. In other words, all the things put in a sentence should have a good reason for being

together. Violations of unity in sentence construction fall into three categories: 1) unrelated ideas, 2) excessive detail, and 3) illogical constructions.

(D) Avoid Sentence Faults.

Fragments, comma splices, and run-on sentences are the three typical sentence faults. A sentence fragment lacks either a subject (Actor) or verb (Action).

II Effective Paragraph Development

A paragraph is a cluster of sentences all related to the same general topic. It is a unit of thought.

A series of paragraphs make up an entire composition. Each paragraph is an important part of the

whole, a key link in the train of thought. Designing paragraphs requires the ability to organize and relate information.

(1) Elements of a Paragraph.

Paragraphs vary widely in length and form. You can communicate effectively in one short paragraph or in pages of lengthy paragraphs, depending on your purpose, your audience, and your message. The typical paragraph contains three basic elements: a topic sentence, related sentences that develop the topic, and transitional words and phrases.

(2) Ways to Develop a Paragraph

IMPROVE SPELLINGS

English spelling and grammar is extremely challenging. Whether you found yourself drifting off

during English class or you are learning English as a second language, spelling and grammar can

be difficult to grasp. That said, if you want to make a good impression on the job, with friends or

at a dinner party, spelling and grammar are critical.

Instructions

- 1. Read lots.
- 2. Look it up.
- 3. Write.

4. Etymology. This isn't the same as entomology, which is the study of insects. Etymology is the study of words and their origins. By understanding the "root" of the word, it is often much easier to remember how to spell versions of it.

व नावध

5. Take classes.

INTRODUCTION TO BUSINESS ENGLISH

Business English is English language especially related to international trade. It is a part of English for Specific Purposes and can be considered a specialism within English language learning and teaching. Business English means different things to different people. For some, it focuses on vocabulary and topics used in the worlds of business, trade, finance, and

international relations. For others it refers to the communication skills used in the workplace, and focuses on the language and skills needed for typical business communication such as presentations, negotiations, meetings, small talk, socializing, correspondence, report writing, and so on. In both of these cases it can be taught to native speakers of English, for example, high school students preparing to enter the job market. It can also be a form of international English. It is possible to study Business English at college and university; institutes around the world have on offer courses (modules) in BE, which can even lead to a degree in the subject.





COPYRIGHT FIMT 2020

DATA STRUCTURE USING "C" (105)

Data Structure and Algorithms

Data Structures are the programmatic way of storing data so that data can be used efficiently. Almost every enterprise application uses various types of data structures in one or the other way. Data Structure is a systematic way to organize data in order to use it efficiently.

Following terms are the foundation terms of a data structure.

- Interface Each data structure has an interface. Interface represents the set of operations that a data structure supports. An interface only provides the list of supported operations, type of parameters they can accept and return type of these operations.
- Implementation Implementation provides the internal representation of a data structure. Implementation also provides the definition of the algorithms used in the operations of the data structure.

Characteristics of a Data Structure

- Correctness Data structure implementation should implement its interface correctly.
- **Time Complexity** Running time or the execution time of operations of data structure must be as small as possible.
- Space Complexity Memory usage of a data structure operation should be as little as possible.

Need for Data Structure

As applications are getting complex and data rich, there are three common problems that applications face now-a-days.

- Data Search Consider an inventory of 1 million(10⁶) items of a store. If the application is to search an item, it has to search an item in 1 million(10⁶) items every time slowing down the search. As data grows, search will become slower.
- **Processor speed** Processor speed although being very high, falls limited if the data grows to billion records.
- **Multiple requests** As thousands of users can search data simultaneously on a web server, even the fast server fails while searching the data.

To solve the above-mentioned problems, data structures come to rescue. Data can be organized in a data structure in such a way that all items may not be required to be searched, and the required data can be searched almost instantly.

Execution Time Cases

There are three cases which are usually used to compare various data structure's execution time in a relative manner.

- Worst Case This is the scenario where a particular data structure operation takes maximum time it can take. If an operation's worst case time is f(n) then this operation will not take more than f(n) time where f(n) represents function of n.
- Average Case This is the scenario depicting the average execution time of an operation of a data structure. If an operation takes f(n) time in execution, then m operations will take mf(n) time.
- Best Case This is the scenario depicting the least possible execution time of an operation of a data structure. If an operation takes f(n) time in execution, then the actual operation may take time as the random number which would be maximum as f(n).

Basic Terminology

- **Data** Data are values or set of values.
- Data Item Data item refers to single unit of values.
- Group Items Data items that are divided into sub items are called as Group Items.
- Elementary Items Data items that cannot be divided are called as Elementary Items.
- Attribute and Entity An entity is that which contains certain attributes or properties, which may be assigned values.
- Entity Set Entities of similar attributes form an entity set.
- Field Field is a single elementary unit of information representing an attribute of an entity.
- Record Record is a collection of field values of a given entity.
- File File is a collection of records of the entities in a given entity set

Data Structures - Algorithms Basics

Algorithm is a step-by-step procedure, which defines a set of instructions to be executed in a certain order to get the desired output. Algorithms are generally created independent of

underlying languages, i.e. an algorithm can be implemented in more than one programming language.

From the data structure point of view, following are some important categories of algorithms

—

- Search Algorithm to search an item in a data structure.
- Sort Algorithm to sort items in a certain order.
- Insert Algorithm to insert item in a data structure.
- Update Algorithm to update an existing item in a data structure.
- **Delete** Algorithm to delete an existing item from a data structure.

Characteristics of an Algorithm

Not all procedures can be called an algorithm. An algorithm should have the following characteristics –

MAGEMEN

- Unambiguous Algorithm should be clear and unambiguous. Each of its steps (or phases), and their inputs/outputs should be clear and must lead to only one meaning.
- Input An algorithm should have 0 or more well-defined inputs.
- **Output** An algorithm should have 1 or more well-defined outputs, and should match the desired output.
- Finiteness Algorithms must terminate after a finite number of steps.
- Feasibility Should be feasible with the available resources.
- **Independent** An algorithm should have step-by-step directions, which should be independent of any programming code.

How to Write an Algorithm?

There are no well-defined standards for writing algorithms. Rather, it is problem and resource dependent. Algorithms are never written to support a particular programming code. As we know that all programming languages share basic code constructs like loops (do, for, while), flow-control (if-else), etc. These common constructs can be used to write an algorithm.

We write algorithms in a step-by-step manner, but it is not always the case. Algorithm writing is a process and is executed after the problem domain is well-defined. That is, we should know the problem domain, for which we are designing a solution.

Example

Let's try to learn algorithm-writing by using an example.

Problem – Design an algorithm to add two numbers and display the result.

step 1 - START

step 2 – declare three integers a, b & c

step 3 – define values of a & b

step 4 – add values of a & b

step 5 – store output of step 4 to c

step 6 - print c

step 7 – STOP

Algorithms tell the programmers how to code the program. Alternatively, the algorithm can be written as –

step 1 – START ADD

step 2 – get values of a & b

step $3 - c \leftarrow a + b$

step 4 – display c

step 5 - STOP

In design and analysis of algorithms, usually the second method is used to describe an algorithm. It makes it easy for the analyst to analyze the algorithm ignoring all unwanted definitions. He can observe what operations are being used and how the process is flowing. Writing **step numbers**, is optional.

We design an algorithm to get a solution of a given problem. A problem can be solved in more than one ways.



Hence, many solution algorithms can be derived for a given problem. The next step is to analyze those proposed solution algorithms and implement the best suitable solution.

Algorithm Analysis

Efficiency of an algorithm can be analyzed at two different stages, before implementation and after implementation. They are the following -

- A Priori Analysis This is a theoretical analysis of an algorithm. Efficiency of an algorithm is measured by assuming that all other factors, for example, processor speed, are constant and have no effect on the implementation.
- A Posterior Analysis This is an empirical analysis of an algorithm. The selected algorithm is implemented using programming language. This is then executed on target computer machine. In this analysis, actual statistics like running time and space required, are collected.

We shall learn about a priori algorithm analysis. Algorithm analysis deals with the execution or running time of various operations involved. The running time of an operation can be defined as the number of computer instructions executed per operation.

Algorithm Complexity

Suppose X is an algorithm and n is the size of input data, the time and space used by the algorithm X are the two main factors, which decide the efficiency of X.

- **Time Factor** Time is measured by counting the number of key operations such as comparisons in the sorting algorithm.
- **Space Factor** Space is measured by counting the maximum memory space required by the algorithm.

The complexity of an algorithm f(n) gives the running time and/or the storage space required by the algorithm in terms of n as the size of input data.

Space Complexity

Space complexity of an algorithm represents the amount of memory space required by the algorithm in its life cycle. The space required by an algorithm is equal to the sum of the following two components –

- A fixed part that is a space required to store certain data and variables, that are independent of the size of the problem. For example, simple variables and constants used, program size, etc.
- A variable part is a space required by variables, whose size depends on the size of the problem. For example, dynamic memory allocation, recursion stack space, etc.

Space complexity S(P) of any algorithm P is S(P) = C + SP(I), where C is the fixed part and S(I) is the variable part of the algorithm, which depends on instance characteristic I. Following is a simple example that tries to explain the concept –

Algorithm: SUM(A, B) Step 1 - START Step 2 - C \leftarrow A + B + 10 Step 3 - Stop

Here we have three variables A, B, and C and one constant. Hence S(P) = 1 + 3. Now, space depends on data types of given variables and constant types and it will be multiplied accordingly.

Time Complexity

Time complexity of an algorithm represents the amount of time required by the algorithm to run to completion. Time requirements can be defined as a numerical function T(n), where T(n) can be measured as the number of steps, provided each step consumes constant time. For example, addition of two n-bit integers takes **n** steps. Consequently, the total computational time is T(n) = c * n, where c is the time taken for the addition of two bits. Here, we observe that T(n) grows linearly as the input size increases.

Data Structures - Asymptotic Analysis

Asymptotic analysis of an algorithm refers to defining the mathematical boundation/framing of its run-time performance. Using asymptotic analysis, we can very well conclude the best case, average case, and worst case scenario of an algorithm.

Asymptotic analysis is input bound i.e., if there's no input to the algorithm, it is concluded to work in a constant time. Other than the "input" all other factors are considered constant.

Asymptotic analysis refers to computing the running time of any operation in mathematical units of computation. For example, the running time of one operation is computed as f(n) and may be for another operation it is computed as $g(n^2)$. This means the first operation running time will increase linearly with the increase in **n** and the running time of the second operation will increase exponentially when **n** increases. Similarly, the running time of both operations will be nearly the same if **n** is significantly small.

Usually, the time required by an algorithm falls under three types -

- Best Case Minimum time required for program execution.
- Average Case Average time required for program execution.
- Worst Case Maximum time required for program execution.

Asymptotic Notations

Following are the commonly used asymptotic notations to calculate the running time complexity of an algorithm.

- O Notation
- Ω Notation
- θ Notation

Big Oh Notation, O

The notation O(n) is the formal way to express the upper bound of an algorithm's running time. It measures the worst case time complexity or the longest amount of time an algorithm can possibly take to complete.

For example, for a function $f(\mathbf{n})$

 $O(f(n)) = \{ g(n) : \text{there exists } c > 0 \text{ and } n_0 \text{ such that } f(n) \le c.g(n) \text{ for all } n > n_0. \}$

Omega Notation, Ω

The notation $\Omega(n)$ is the formal way to express the lower bound of an algorithm's running time. It measures the best case time complexity or the best amount of time an algorithm can possibly take to complete.



For example, for a function *f*(**n**)

 $\Omega(f(n)) \ge \{ g(n) : \text{there exists } c > 0 \text{ and } n_0 \text{ such that } g(n) \le c.f(n) \text{ for all } n > n_0. \}$

Theta Notation, θ

COPYRIGHT FIMT 2020

The notation $\theta(n)$ is the formal way to express both the lower bound and the upper bound of an algorithm's running time. It is represented as follows –



 $\theta(f(n)) = \{ g(n) \text{ if and only if } g(n) = O(f(n)) \text{ and } g(n) = \Omega(f(n)) \text{ for all } n > n_0. \}$

Common Asymptotic Notations

Following is a list of some common asymptotic notations -

constant	-RA	O(1)
logarithmic	-	O(log n)
linear	- 1	O(n)
n log n	1.11	O(n log n)
quadratic	-	O(n ²)
cubic	-	$O(n^3)$
polynomial		n ^{O(1)}
exponential	_	2 ^{O(n)}

Data Structures - Greedy Algorithms

An algorithm is designed to achieve optimum solution for a given problem. In greedy algorithm approach, decisions are made from the given solution domain. As being greedy, the closest solution that seems to provide an optimum solution is chosen.

Greedy algorithms try to find a localized optimum solution, which may eventually lead to globally optimized solutions. However, generally greedy algorithms do not provide globally optimized solutions.

Counting Coins

This problem is to count to a desired value by choosing the least possible coins and the greedy approach forces the algorithm to pick the largest possible coin. If we are provided coins of \gtrless 1, 2, 5 and 10 and we are asked to count \gtrless 18 then the greedy procedure will be –

- 1 -Select one $\gtrless 10$ coin, the remaining count is 8
- 2 Then select one \gtrless 5 coin, the remaining count is 3
- **3** Then select one \gtrless 2 coin, the remaining count is 1
- 4 And finally, the selection of one \gtrless 1 coins solves the problem

Though, it seems to be working fine, for this count we need to pick only 4 coins. But if we slightly change the problem then the same approach may not be able to produce the same optimum result. For the currency system, where we have coins of 1, 7, 10 value, counting coins for value 18 will be absolutely optimum but for count like 15, it may use more coins than necessary. For example, the greedy approach will use 10 + 1 + 1 + 1 + 1 + 1, total 6 coins. Whereas the same problem could be solved by using only 3 coins (7 + 7 + 1) Hence, we may conclude that the greedy approach picks an immediate optimized solution

Examples

Most networking algorithms use the greedy approach. Here is a list of few of them -

- Travelling Salesman Problem
- Prim's Minimal Spanning Tree Algorithm
- Kruskal's Minimal Spanning Tree Algorithm

and may fail where global optimization is a major concern.

- Dijkstra's Minimal Spanning Tree Algorithm
- Graph Map Coloring
- Graph Vertex Cover
- Knapsack Problem
- Job Scheduling Problem

There are lots of similar problems that uses the greedy approach to find an optimum solution.

Data Structures - Divide and Conquer

In divide and conquer approach, the problem in hand, is divided into smaller sub-problems and then each problem is solved independently. When we keep on dividing the subproblems into even smaller sub-problems, we may eventually reach a stage where no more division is possible. Those "atomic" smallest possible sub-problem (fractions) are solved. The solution of all sub-problems is finally merged in order to obtain the solution of an original problem.



Broadly, we can understand divide-and-conquer approach in a three-step process. AGEM

Divide/Break

This step involves breaking the problem into smaller sub-problems. Sub-problems should represent a part of the original problem. This step generally takes a recursive approach to divide the problem until no sub-problem is further divisible. At this stage, sub-problems become atomic in nature but still represent some part of the actual problem.

Conquer/Solve

This step receives a lot of smaller sub-problems to be solved. Generally, at this level, the problems are considered 'solved' on their own.

Merge/Combine

When the smaller sub-problems are solved, this stage recursively combines them until they formulate a solution of the original problem. This algorithmic approach works recursively and conquer & merge steps works so close that they appear as one.

Examples

The following computer algorithms are based on divide-and-conquer programming approach

:2015 & 14001:201

- Merge Sort
- **Quick Sort** •
- Binary Search •
- Strassen's Matrix Multiplication
- Closest pair (points) •

There are various ways available to solve any computer problem, but the mentioned are a good example of divide and conquer approach.

Data Structures - Dynamic Programming

Dynamic programming approach is similar to divide and conquer in breaking down the problem into smaller and yet smaller possible sub-problems. But unlike, divide and conquer, these sub-problems are not solved independently. Rather, results of these smaller sub-problems are remembered and used for similar or overlapping sub-problems.

Dynamic programming is used where we have problems, which can be divided into similar sub-problems, so that their results can be re-used. Mostly, these algorithms are used for optimization. Before solving the in-hand sub-problem, dynamic algorithm will try to examine the results of the previously solved sub-problems. The solutions of sub-problems are combined in order to achieve the best solution.

So we can say that -

- The problem should be able to be divided into smaller overlapping sub-problem.
- An optimum solution can be achieved by using an optimum solution of smaller subproblems.
- Dynamic algorithms use memorization.

Comparison

In contrast to greedy algorithms, where local optimization is addressed, dynamic algorithms are motivated for an overall optimization of the problem.

In contrast to divide and conquer algorithms, where solutions are combined to achieve an overall solution, dynamic algorithms use the output of a smaller sub-problem and then try to optimize a bigger sub-problem. Dynamic algorithms use memorization to remember the output of already solved sub-problems.

Example

The following computer problems can be solved using dynamic programming approach -

15 & 14001:2015

- Fibonacci number series
- Knapsack problem
- Tower of Hanoi
- All pair shortest path by Floyd-Warshall
- Shortest path by Dijkstra
- Project scheduling

Dynamic programming can be used in both top-down and bottom-up manner. And of course, most of the times, referring to the previous solution output is cheaper than recomputing in terms of CPU cycles.

Data Structures & Algorithm Basic Concepts

This chapter explains the basic terms related to data structure.

Data Definition

Data Definition defines a particular data with the following characteristics.

- Atomic Definition should define a single concept.
- Traceable Definition should be able to be mapped to some data element.
- Accurate Definition should be unambiguous.
- Clear and Concise Definition should be understandable.

Data Object

Data Object represents an object having a data.

Data Type

Data type is a way to classify various types of data such as integer, string, etc. which determines the values that can be used with the corresponding type of data, the type of operations that can be performed on the corresponding type of data. There are two data types

- Built-in Data Type
- Derived Data Type

Built-in Data Type

Those data types for which a language has built-in support are known as Built-in Data types. For example, most of the languages provide the following built-in data types.

5 & 14001-2015

- Integers
- Boolean (true, false)
- Floating (Decimal numbers)
- Character and Strings

Derived Data Type

Those data types which are implementation independent as they can be implemented in one or the other way are known as derived data types. These data types are normally built by the combination of primary or built-in data types and associated operations on them. For example –

- List
- Array
- Stack
- Queue

Basic Operations

The data in the data structures are processed by certain operations. The particular data structure chosen largely depends on the frequency of the operation that needs to be performed on the data structure.

NAAC ACCREDITED

ARGEME

- Traversing
- Searching
- Insertion
- Deletion
- Sorting
- Merging

Data Structures and Algorithms - Arrays

Array is a container which can hold a fix number of items and these items should be of the same type. Most of the data structures make use of arrays to implement their algorithms. Following are the important terms to understand the concept of Array.

- Element Each item stored in an array is called an element.
- Index Each location of an element in an array has a numerical index, which is used to identify the element.

Array Representation

Arrays can be declared in various ways in different languages. For illustration, let's take C array declaration.

Nam	e					Eleme	ents				
int arra	y [10]	= { :	35, 3	3, 42	, 10,	14, 1	9, 27	, 44,	26,	31	}
Туре	Size										
elements	35	33	42	10	14	19	27	44	26		31
index	0	1	2	з	4	5	6	7	8		9
					Siz	e:10					

As per the above illustration, following are the important points to be considered.

- Index starts with 0.
- Array length is 10 which means it can store 10 elements.
- Each element can be accessed via its index. For example, we can fetch an element at index 6 as 9.

1001

Basic Operations

Following are the basic operations supported by an array.

- **Traverse** print all the array elements one by one.
- Insertion Adds an element at the given index.
- **Deletion** Deletes an element at the given index.
- Search Searches an element using the given index or by the value.
- Update Updates an element at the given index.

In C, when an array is initialized with size, then it assigns defaults values to its elements in following order.

Data Type	Default Value
bool	false
char	O du du du
int	0
float	0.0
double	0.0f
void	
wchar_t	0

Insertion Operation

Insert operation is to insert one or more data elements into an array. Based on the requirement, a new element can be added at the beginning, end, or any given index of array. Here, we see a practical implementation of insertion operation, where we add data at the end of the array -

Algorithm

Let Array be a linear unordered array of MAX elements.

TRANAGE

Example

Result

Let **LA** be a Linear Array (unordered) with **N** elements and **K** is a positive integer such that $K \le N$. Following is the algorithm where ITEM is inserted into the Kth position of LA –

- 1. Start
- 2. Set J = N
- 3. Set N = N+1
- 4. Repeat steps 5 and 6 while $J \ge K$
- 5. Set LA[J+1] = LA[J]
- 6. Set J = J 1
- 7. Set LA[K] = ITEM
- 8. Stop

Example

Following is the implementation of the above algorithm -

```
#include <stdio.h>
main() {
    int LA[] = {1,3,5,7,8};
    int item = 10, k = 3, n = 5;
    int i = 0, j = n;
```

```
printf("The original array elements are :\
for(i = 0; i<n; i++) {
   printf("LA[%d] = %d \n", i, LA[i]);
 }
             NAAC ACCREDITED
  n = n + 1;
 while (j \ge k) {
                        ERREAG
   LA[j+1] = LA[j];
   j = j - 1;
 LA[k] = item;
        printf("The array elements after insertion :\n"
 for(i = 0; i<n; i++) {
   printf("LA[%d] = %d \n", i, LA[i]);
 }
}
When we compile and execute the above program, it produces the following result -
Output
The original array elements are :
LA[0] = 1
LA[1] = 3
LA[2] = 5
LA[3] = 7
LA[4] = 8
The array elements after insertion :
```

LA[0] = 1 LA[1] = 3 LA[2] = 5 LA[3] = 10 LA[4] = 7LA[5] = 8

For other variations of array insertion operation <u>click here</u>

Deletion Operation

Deletion refers to removing an existing element from the array and re-organizing all elements of an array.

Algorithm

Consider LA is a linear array with N elements and K is a positive integer such that $K \le N$. Following is the algorithm to delete an element available at the Kth position of LA.

1. Start

```
2. Set J = K
```

3. Repeat steps 4 and 5 while J < N

```
4. Set LA[J] = LA[J + 1]
```

- 5. Set J = J+1
- 6. Set N = N-1
- 7. Stop

Example

Following is the implementation of the above algorithm

```
#include <stdio.h>
main() {
    int LA[] = {1,3,5,7,8};
    int k = 3, n = 5;
    int i, j;
```

```
printf("The original array elements are :\n");
 for(i = 0; i<n; i++) {
   printf("LA[%d] = %d \n", i, LA[i]);
 }
               NAAC ACCREDITED
   j = k;
 while (j < n) {
                        ERRHAG
   LA[j-1] = LA[j];
                     S.
   j = j + 1;
 }
 n = n - 1;
   printf("The array elements after deletion :\n");
          for(i = 0; i<n; i++) {
   printf("LA[%d] = %d n", i, LA[i]);
 }
}
When we compile and execute the above program, it produces the following result -
Output
The original array elements are :
LA[0] = 1
LA[1] = 3
LA[2] = 5
LA[3] = 7
LA[4] = 8
The array elements after deletion :
```

LA[0] = 1LA[1] = 3LA[2] = 7LA[3] = 8

Search Operation

You can perform a search for an array element based on its value or its index.

VAAC ACCKEDII

Algorithm

Consider LA is a linear array with N elements and K is a positive integer such that $K \le N$. Following is the algorithm to find an element with a value of ITEM using sequential search.

1. Start

```
2. Set J = 0
```

- 3. Repeat steps 4 and 5 while J < N
- 4. IF LA[J] is equal ITEM THEN GOTO STEP 6
- 5. Set J = J + 1
- 6. PRINT J, ITEM
- 7. Stop

Example

Following is the implementation of the above algorithm -



COPYRIGHT FIMT 2020

Output

The original array elements are : LA[0] = 1 LA[1] = 3 LA[2] = 5 LA[3] = 7 LA[4] = 8Found element 5 at position 3

Update Operation

Update operation refers to updating an existing element from the array at a given index.

Algorithm

Consider LA is a linear array with N elements and K is a positive integer such that $K \le N$. Following is the algorithm to update an element available at the Kth position of LA.

- 1. Start
- 2. Set LA[K-1] = ITEM
- 3. Stop

Example

Following is the implementation of the above algorithm -

```
#include <stdio.h>
                          ARMAR
main() {
 int LA[] = \{1,3,5,7,8\};
 int k = 3, n = 5, item = 10;
 int i, j;
 printf("The original array elements are :\n");
 for(i = 0; i<n; i++) {
   printf("LA[%d] = %d n", i, LA[i]).
  }
 LA[k-1] = item;
 printf("The array elements after updation :\n");
 for(i = 0; i<n; i++) {
                          01:2015 & 14001:20
   printf("LA[%d] = %d \n", i, LA[i]);
  }
}
When we compile and execute the above program, it produces the following result -
```

Output

The original array elements are :

LA[0] = 1 LA[1] = 3 LA[2] = 5 LA[3] = 7 LA[4] = 8The array elements after updation : LA[0] = 1 LA[1] = 3 LA[2] = 10 LA[3] = 7LA[4] = 8

Data Structure and Algorithms - Linked List

A linked list is a sequence of data structures, which are connected together via links.

Linked List is a sequence of links which contains items. Each link contains a connection to another link. Linked list is the second most-used data structure after array. Following are the important terms to understand the concept of Linked List.

- Link Each link of a linked list can store a data called an element.
- Next Each link of a linked list contains a link to the next link called Next.
- LinkedList A Linked List contains the connection link to the first link called First.

Linked List Representation

Linked list can be visualized as a chain of nodes, where every node points to the next node.



As per the above illustration, following are the important points to be considered.

- Linked List contains a link element called first.
- Each link carries a data field(s) and a link field called next.
- Each link is linked with its next link using its next link.
- Last link carries a link as null to mark the end of the list.

Types of Linked List

Following are the various types of linked list.

- Simple Linked List Item navigation is forward only.
- **Doubly Linked List** Items can be navigated forward and backward.
- **Circular Linked List** Last item contains link of the first element as next and the first element has a link to the last element as previous.

Basic Operations

Following are the basic operations supported by a list.

- Insertion Adds an element at the beginning of the list.
- Deletion Deletes an element at the beginning of the list.
- **Display** Displays the complete list.
- Search Searches an element using the given key.
- **Delete** Deletes an element using the given key.

Insertion Operation

Adding a new node in linked list is a more than one step activity. We shall learn this with diagrams here. First, create a node using the same structure and find the location where it has to be inserted.

	NODE				NODE	
ad	Data Items	Next			→ Data Items	Next
			Data Items	Next		
			New NOE	DE		

Imagine that we are inserting a node **B** (NewNode), between **A** (LeftNode) and **C** (RightNode). Then point B.next to C -



It should look like this -

	Data Items	Next
	L	<u>lo</u>
Vext		
4	ext	ext

New NODE

Now, the next node at the left should point to the new node.

```
LeftNode.next -> NewNode;
```


This will put the new node in the middle of the two. The new list should look like this -



Similar steps should be taken if the node is being inserted at the beginning of the list. While inserting it at the end, the second last node of the list should point to the new node and the new node will point to NULL.

Deletion Operation

Deletion is also a more than one step process. We shall learn with pictorial representation. First, locate the target node to be removed, by using searching algorithms.



The left (previous) node of the target node now should point to the next node of the target node –

LeftNode.next -> TargetNode.next; NODE NODE NODE Data Items Next Data Items Next Data Items Next Target NODE

This will remove the link that was pointing to the target node. Now, using the following code, we will remove what the target node is pointing at.

TargetNode.next -> NULL;



We need to use the deleted node. We can keep that in memory otherwise we can simply deallocate memory and wipe off the target node completely.



Reverse Operation

This operation is a thorough one. We need to make the last node to be pointed by the head node and reverse the whole linked list.



First, we traverse to the end of the list. It should be pointing to NULL. Now, we shall make it point to its previous node –



We have to make sure that the last node is not the lost node. So we'll have some temp node, which looks like the head node pointing to the last node. Now, we shall make all left side nodes point to their previous nodes one by one.



Except the node (first node) pointed by the head node, all nodes should point to their predecessor, making them their new successor. The first node will point to NULL.



We'll make the head node point to the new first node by using the temp node.



The linked list is now reversed.

Data Structure - Doubly Linked List

Doubly Linked List is a variation of Linked list in which navigation is possible in both ways, either forward and backward easily as compared to Single Linked List. Following are the important terms to understand the concept of doubly linked list.

- Link Each link of a linked list can store a data called an element.
- Next Each link of a linked list contains a link to the next link called Next.
- **Prev** Each link of a linked list contains a link to the previous link called Prev.
- LinkedList A Linked List contains the connection link to the first link called First and to the last link called Last.

Doubly Linked List Representation



As per the above illustration, following are the important points to be considered.

- Doubly Linked List contains a link element called first and last.
- Each link carries a data field(s) and two link fields called next and prev.
- Each link is linked with its next link using its next link.
- Each link is linked with its previous link using its previous link.
- The last link carries a link as null to mark the end of the list.

Basic Operations

Following are the basic operations supported by a list.

NULL

- Insertion Adds an element at the beginning of the list.
- **Deletion** Deletes an element at the beginning of the list.
- Insert Last Adds an element at the end of the list.
- **Delete Last** Deletes an element from the end of the list.
- Insert After Adds an element after an item of the list.
- **Delete** Deletes an element from the list using the key.
- **Display forward** Displays the complete list in a forward manner.
- Display backward Displays the complete list in a backward manner.

Insertion Operation

Following code demonstrates the insertion operation at the beginning of a doubly linked list.

Example

```
//insert link at the first location
```

void insertFirst(int key, int data) {

//create a link

struct node *link = (struct node*) malloc(sizeof(struct node));

link->key = key;

link->data = data;

```
if(isEmpty()) {
```

//make it the last link

last = link;

} else {

//update first prev link

head->prev = link;

}

//point it to old first link

```
link->next = head;
```

//point first to new first link

```
head = link;
```

}

Deletion Operation

Following code demonstrates the deletion operation at the beginning of a doubly linked list.

SPARGEMENT .

Example

//delete first item

struct node* deleteFirst() {

//save reference to first link

```
struct node *tempLink = head;
```

//if only one link

```
if(head->next == NULL) {
```

```
last = NULL;
```

} else {

head->next->prev = NULL;

}

```
head = head->next;
```

//return the deleted link

return tempLink;

}

Insertion at the End of an Operation

Following code demonstrates the insertion operation at the last position of a doubly linked list.

Example

//insert link at the last location

void insertLast(int key, int data) {

//create a link

struct node *link = (struct node*) malloc(sizeof(struct node));

link->key = key;

link->data = data;

if(isEmpty()) {

//make it the last link

last = link;

} else {

//make link a new last link

last->next = link;

//mark old last node as prev of new link

```
link->prev = last;
```

}

//point last to new last node

last = link;

}

Data Structure - Circular Linked List

Circular Linked List is a variation of Linked list in which the first element points to the last element and the last element points to the first element. Both Singly Linked List and Doubly Linked List can be made into a circular linked list.

Singly Linked List as Circular

In singly linked list, the next pointer of the last node points to the first node.



Doubly Linked List as Circular

In doubly linked list, the next pointer of the last node points to the first node and the previous pointer of the first node points to the last node making the circular in both directions.



As per the above illustration, following are the important points to be considered.

- The last link's next points to the first link of the list in both cases of singly as well as doubly linked list.
- The first link's previous points to the last of the list in case of doubly linked list.

Basic Operations

Following are the important operations supported by a circular list.

- insert Inserts an element at the start of the list.
- delete Deletes an element from the start of the list.
- **display** Displays the list.

Insertion Operation

Following code demonstrates the insertion operation in a circular linked list based on single linked list.

14001:20

Example

//insert link at the first location

```
void insertFirst(int key, int data) {
```

```
//create a link
```

struct node *link = (struct node*) malloc(sizeof(struct node));

link->key = key;

link->data= data;

if (isEmpty()) {

head = link;

head->next = head;

} else {

//point it to old first node

link->next = head;

//point first to new first node

```
head = link;
```

```
}
```

```
}
```

Deletion Operation

Following code demonstrates the deletion operation in a circular linked list based on single linked list.

100 0001.0010 P 14001.0014

```
//delete first item
```

struct node * deleteFirst() {

//save reference to first link

```
struct node *tempLink = head;
```

```
if(head->next == head) {
```

```
head = NULL;
```

return tempLink;

}

//mark next to first link as first

```
head = head->next;
```

//return the deleted link

return tempLink;

```
}
```

Display List Operation

Following code demonstrates the display list operation in a circular linked list.

```
//display the list
```

void printList() {

```
struct node *ptr = head;
```

```
printf("\n[ ");
```

//start from the beginning

```
if(head != NULL) {
```

```
while(ptr->next != ptr) {
```

printf("(%d,%d) ",ptr->key,ptr->data);

```
ptr = ptr->next;
```

}

}

```
printf(" ]");
```

1

Data Structure and Algorithms - Stack

A stack is an Abstract Data Type (ADT), commonly used in most programming languages. It is named stack as it behaves like a real-world stack, for example – a deck of cards or a pile of plates, etc.



A real-world stack allows operations at one end only. For example, we can place or remove a card or plate from the top of the stack only. Likewise, Stack ADT allows all data operations at one end only. At any given time, we can only access the top element of a stack. This feature makes it LIFO data structure. LIFO stands for Last-in-first-out. Here, the element which is placed (inserted or added) last, is accessed first. In stack terminology, insertion operation is called **PUSH** operation and removal operation is called **POP** operation.

Stack Representation

The following diagram depicts a stack and its operations -



A stack can be implemented by means of Array, Structure, Pointer, and Linked List. Stack can either be a fixed size one or it may have a sense of dynamic resizing. Here, we are going to implement stack using arrays, which makes it a fixed size stack implementation.

Basic Operations

Stack operations may involve initializing the stack, using it and then de-initializing it. Apart from these basic stuffs, a stack is used for the following two primary operations –

- **push**() Pushing (storing) an element on the stack.
- **pop**() Removing (accessing) an element from the stack.

When data is PUSHed onto stack.

To use a stack efficiently, we need to check the status of stack as well. For the same purpose, the following functionality is added to stacks –

- peek() get the top data element of the stack, without removing it.
- **isFull**() check if stack is full.
- **isEmpty**() check if stack is empty.

At all times, we maintain a pointer to the last PUSHed data on the stack. As this pointer always represents the top of the stack, hence named **top**. The **top**pointer provides top value of the stack without actually removing it.

First we should learn about procedures to support stack functions -

peek()

Algorithm of peek() function -

begin procedure peek

return stack[top]

end procedure

Implementation of peek() function in C programming language -

Example

```
int peek() {
```

return stack[top];

}

isfull()

Algorithm of isfull() function -

begin procedure isfull

if top equals to MAXSIZE

return true

else

return false

endif

end procedure

Implementation of isfull() function in C programming language -

Example

bool isfull() {

if(top == MAXSIZE)

return true;

else

return false;

}

isempty()

Algorithm of isempty() function -

begin procedure isempty

if top less than 1

return true

else

return false

endif

end procedure

Implementation of isempty() function in C programming language is slightly different. We initialize top at -1, as the index in array starts from 0. So we check if the top is below zero or -1 to determine if the stack is empty. Here's the code –

Example

```
bool isempty() {
    if(top == -1)
    return true;
    else
    return false;
}
Push Operation
```

The process of putting a new data element onto stack is known as a Push Operation. Push operation involves a series of steps –

- Step 1 Checks if the stack is full.
- Step 2 If the stack is full, produces an error and exit.
- Step 3 If the stack is not full, increments top to point next empty space.
- Step 4 Adds data element to the stack location, where top is pointing.
- Step 5 Returns success.



If the linked list is used to implement the stack, then in step 3, we need to allocate space dynamically.

AAGEMENT

Algorithm for PUSH Operation

A simple algorithm for Push operation can be derived as follows -

begin procedure push: stack, data

if stack is full

return null

endif

 $top \leftarrow top + 1$

stack[top] ← data

end procedure

Implementation of this algorithm in C, is very easy. See the following code -

Example

void push(int data) {

if(!isFull()) {

top = top + 1;

stack[top] = data;

} else {

printf("Could not insert data, Stack is full.\n");

Pop Operation

}

}

Accessing the content while removing it from the stack, is known as a Pop Operation. In an array implementation of pop() operation, the data element is not actually removed, instead **top** is decremented to a lower position in the stack to point to the next value. But in linked-list implementation, pop() actually removes data element and deallocates memory space.

A Pop operation may involve the following steps -

- Step 1 Checks if the stack is empty.
- Step 2 If the stack is empty, produces an error and exit.
- Step 3 If the stack is not empty, accesses the data element at which top is pointing.
- Step 4 Decreases the value of top by 1.
- Step 5 Returns success.



A simple algorithm for Pop operation can be derived as follows –

begin procedure pop: stack

if stack is empty

return null

endif

```
data ← stack[top]
```

 $top \leftarrow top - 1$

return data

end procedure

Implementation of this algorithm in C, is as follows -

Example

```
int pop(int data) {
```

if(!isempty()) {

data = stack[top];

top = top - 1;

return data;

} else {

printf("Could not retrieve data, Stack is empty.\n");

}

}

Data Structure - Expression Parsing

The way to write arithmetic expression is known as a **notation**. An arithmetic expression can be written in three different but equivalent notations, i.e., without changing the essence or output of an expression. These notations are -

व नावधातन

ANA GEMENT

- Infix Notation
- Prefix (Polish) Notation
- Postfix (Reverse-Polish) Notation

These notations are named as how they use operator in expression. We shall learn the same here in this chapter.

Infix Notation

We write expression in **infix** notation, e.g. a - b + c, where operators are used **in**-between operands. It is easy for us humans to read, write, and speak in infix notation but the same does not go well with computing devices. An algorithm to process infix notation could be difficult and costly in terms of time and space consumption.

Prefix Notation

In this notation, operator is **prefixed** to operands, i.e. operator is written ahead of operands. For example, +ab. This is equivalent to its infix notation a + b. Prefix notation is also known as **Polish Notation**.

Postfix Notation

This notation style is known as **Reversed Polish Notation**. In this notation style, the operator is **postfixed** to the operands i.e., the operator is written after the operands. For example, ab+. This is equivalent to its infix notation a + b.

Sr. No.	Infix Notation	Prefix Notation	Postfix Notation
1	a + b	+ a b	a b +
2	(a + b) * c	* + a b c	a b + c *
3	a * (b + c)	* a + b c	a b c + *
4	a / b + c / d	+ / a b / c d	a b / c d / +
5	(a + b) * (c + d)	* + a b + c d	a b + c d + *
6	((a + b) * c) - d	- * + a b c d	a b + c * d -

The following table briefly tries to show the difference in all three notations -

Parsing Expressions

As we have discussed, it is not a very efficient way to design an algorithm or program to parse infix notations. Instead, these infix notations are first converted into either postfix or prefix notations and then computed.

To parse any arithmetic expression, we need to take care of operator precedence and associativity also.

Precedence

When an operand is in between two different operators, which operator will take the operand first, is decided by the precedence of an operator over others. For example –

a+b*c 📥 a+(b*c)

As multiplication operation has precedence over addition, b * c will be evaluated first. A table of operator precedence is provided later.

ARGEME

Associativity

Associativity describes the rule where operators with the same precedence appear in an expression. For example, in expression a + b - c, both + and – have the same precedence, then which part of the expression will be evaluated first, is determined by associativity of those operators. Here, both + and – are left associative, so the expression will be evaluated as (a + b) - c.

Precedence and associativity determines the order of evaluation of an expression. Following is an operator precedence and associativity table (highest to lowest) –

Sr. No.	Operator	Precedence	Associativity
1	Exponentiation ^	Highest	Right Associative
2	Multiplication (*) & Division (/)	Second Highest	Left Associative
3	Addition (+) & Subtraction (-)	Lowest	Left Associative

The above table shows the default behavior of operators. At any point of time in expression evaluation, the order can be altered by using parenthesis. For example –

In $\mathbf{a} + \mathbf{b} \cdot \mathbf{c}$, the expression part $\mathbf{b} \cdot \mathbf{c}$ will be evaluated first, with multiplication as precedence over addition. We here use parenthesis for $\mathbf{a} + \mathbf{b}$ to be evaluated first, like $(\mathbf{a} + \mathbf{b}) \cdot \mathbf{c}$.

Postfix Evaluation Algorithm

We shall now look at the algorithm on how to evaluate postfix notation -

Step 1 – scan the expression from left to right

Step 2 – if it is an operand push it to stack

Step 3 – if it is an operator pull operand from stack and perform operation

Step 4 – store the output of step 3, back to stack

Step 5 - scan the expression until all operands are consumed Step 6 – pop the stack and perform operation

Data Structure and Algorithms - Queue

Queue is an abstract data structure, somewhat similar to Stacks. Unlike stacks, a queue is open at both its ends. One end is always used to insert data (enqueue) and the other is used to remove data (dequeue). Queue follows First-In-First-Out methodology, i.e., the data item stored first will be accessed first.



A real-world example of queue can be a single-lane one-way road, where the vehicle enters first, exits first. More real-world examples can be seen as queues at the ticket windows and bus-stops.

Queue Representation

As we now understand that in queue, we access both ends for different reasons. The following diagram given below tries to explain queue representation as data structure –



Last In Last Out

First In First Out

Queue

As in stacks, a queue can also be implemented using Arrays, Linked-lists, Pointers and Structures. For the sake of simplicity, we shall implement queues using one-dimensional array. :2015 & 14001:2015

Basic Operations

Queue operations may involve initializing or defining the queue, utilizing it, and then completely erasing it from the memory. Here we shall try to understand the basic operations associated with queues -

- enqueue() add (store) an item to the queue.
- **dequeue**() remove (access) an item from the queue.

Few more functions are required to make the above-mentioned queue operation efficient. These are –

- peek() Gets the element at the front of the queue without removing it.
- **isfull**() Checks if the queue is full.
- **isempty**() Checks if the queue is empty.

In queue, we always dequeue (or access) data, pointed by **front** pointer and while enqueing (or storing) data in the queue we take help of **rear** pointer. Let's first learn about supportive functions of a queue –

peek()

This function helps to see the data at the **front** of the queue. The algorithm of peek() function is as follows –

AAGEMEN

Algorithm

begin procedure peek

```
return queue[front]
```

end procedure

Implementation of peek() function in C programming language -

Example

```
int peek() {
```

```
return queue[front];
```

}

isfull()

150 9001:2015 & 14001:2015

As we are using single dimension array to implement queue, we just check for the rear pointer to reach at MAXSIZE to determine that the queue is full. In case we maintain the queue in a circular linked-list, the algorithm will differ. Algorithm of isfull() function –

Algorithm

begin procedure isfull

if rear equals to MAXSIZE

return true

else

return false

endif

end procedure

Implementation of isfull() function in C programming language -

Example

```
bool isfull() {
```

```
if(rear == MAXSIZE - 1)
```

return true;

else

return false;

}

```
isempty()
```



begin procedure isempty

if front is less than MIN OR front is greater than rear

return true

else

return false

endif

end procedure

If the value of **front** is less than MIN or 0, it tells that the queue is not yet initialized, hence empty.

NAAC ACCREDITED

Here's the C programming code -

Example

```
bool isempty() {
```

```
if(front < 0 || front > rear)
```

return true;

else

return false;

}

Enqueue Operation

Queues maintain two data pointers, **front** and **rear**. Therefore, its operations are comparatively difficult to implement than that of stacks.

The following steps should be taken to enqueue (insert) data into a queue -

- Step 1 Check if the queue is full.
- Step 2 If the queue is full, produce overflow error and exit.
- Step 3 If the queue is not full, increment rear pointer to point the next empty space.
- Step 4 Add data element to the queue location, where the rear is pointing.
- Step 5 return success.



Queue Enqueue

Sometimes, we also check to see if a queue is initialized or not, to handle any unforeseen situations.

Algorithm for enqueue operation

procedure enqueue(data)

if queue is full

return overflow

endif

rear \leftarrow rear + 1

queue[rear] ← data

return true

end procedure

Implementation of enqueue() in C programming language -

150 9001:2015 & 14001:2015

Example

int enqueue(int data)

if(isfull())

return 0;

rear = rear + 1;

queue[rear] = data;

return 1:

end procedure

Dequeue Operation

Accessing data from the queue is a process of two tasks - access the data where front is pointing and remove the data after access. The following steps are taken to perform dequeue operation -

- Step 1 Check if the queue is empty.
- Step 2 If the queue is empty, produce underflow error and exit.
- Step 3 If the queue is not empty, access the data where **front** is pointing.
- Step 4 Increment front pointer to point to the next available data element.



• Step 5 – Return success.

return true

end procedure

Implementation of dequeue() in C programming language -

Example

```
int dequeue() {
    if(isempty())
        return 0;
    int data = queue[front];
    front = front + 1;
    return data;
}
```

Data Structure and Algorithms Linear Search

Linear search is a very simple search algorithm. In this type of search, a sequential search is made over all items one by one. Every item is checked and if a match is found then that particular item is returned, otherwise the search continues till the end of the data collection.

Linear Search



Step 1: Set i to 1
Step 2: if i > n then go to step 7
Step 3: if A[i] = x then go to step 6
Step 4: Set i to i + 1

Step 5: Go to Step 2Step 6: Print Element x Found at index i and go to step 8Step 7: Print element not foundStep 8: Exit

Pseudocode



procedure linear_search (list, value)

for each item in the list

if match item == value

reurn the item's location

end if

end for

end procedure

Data Structure and Algorithms Binary Search

Binary search is a fast search algorithm with run-time complexity of O(log n). This search algorithm works on the principle of divide and conquer. For this algorithm to work properly, the data collection should be in the sorted form. Binary search looks for a particular item by comparing the middle most item of the collection. If a match occurs, then the index of item is returned. If the middle item is greater than the item, then the item is searched in the sub-array to the left of the middle item. Otherwise, the item is searched for in the sub-array to the right of the middle item. This process continues on the sub-array as well until the size of the subarray reduces to zero.

RELEL

How Binary Search Works?

For a binary search to work, it is mandatory for the target array to be sorted. We shall learn the process of binary search with a pictorial example. The following is our sorted array and let us assume that we need to search the location of value 31 using binary search.



First, we shall determine half of the array by using this formula –

mid = low + (high - low) / 2

Here it is, 0 + (9 - 0) / 2 = 4 (integer value of 4.5). So, 4 is the mid of the array.



Now we compare the value stored at location 4, with the value being searched, i.e. 31. We find that the value at location 4 is 27, which is not a match. As the value is greater than 27 and we have a sorted array, so we also know that the target value must be in the upper portion of the array.



We change our low to mid + 1 and find the new mid value again.

```
low = mid + 1
mid = low + (high - low) / 2
```

Our new mid is 7 now. We compare the value stored at location 7 with our target value 31.



The value stored at location 7 is not a match, rather it is more than what we are looking for. So, the value must be in the lower part from this location.



Hence, we calculate the mid again. This time it is 5.



We compare the value stored at location 5 with our target value. We find that it is a match.

We conclude that the target value 31 is stored at location 5.

Binary search halves the searchable items and thus reduces the count of comparisons to be made to very less numbers.

Pseudocode

The pseudocode of binary search algorithms should look like this -

Procedure binary_search

 $A \leftarrow \text{sorted array}$

 $n \leftarrow size of array$

 $x \leftarrow$ value to be searched

Set lowerBound = 1

Set upperBound = n

while x not found

if upperBound < lowerBound

EXIT: x does not exists.

set midPoint = lowerBound + (upperBound - lowerBound) / 2

if A[midPoint] < x</pre>

set lowerBound = midPoint + 1

if A[midPoint] > x

set upperBound = midPoint - 1

if A[midPoint] = x

EXIT: x found at location midPoint

end while

end procedure

Data Structure - Interpolation Search

Interpolation search is an improved variant of binary search. This search algorithm works on the probing position of the required value. For this algorithm to work properly, the data collection should be in a sorted form and equally distributed.

Binary search has a huge advantage of time complexity over linear search. Linear search has worst-case complexity of O(n) whereas binary search has $O(\log n)$.

There are cases where the location of target data may be known in advance. For example, in case of a telephone directory, if we want to search the telephone number of Morphius. Here, linear search and even binary search will seem slow as we can directly jump to memory space where the names start from 'M' are stored.

Positioning in Binary Search

In binary search, if the desired data is not found then the rest of the list is divided in two parts, lower and higher. The search is carried out in either of them.

1	2	3	4	5	6	7	8	9	10
1	2	3	4	5	6	7	8	9	10
-									
1	2	3	4	5	6	7	8	9	10
	\sim								
1	2	3	4	5	6	7	8	9	10
	+								

Even when the data is sorted, binary search does not take advantage to probe the position of the desired data.

Position Probing in Interpolation Search

Interpolation search finds a particular item by computing the probe position. Initially, the probe position is the position of the middle most item of the collection.

1	2	3	4	5	6	7	8	9	10
_	~						-		
1	2	3	4	5	6	7	8	9	10
	+								

If a match occurs, then the index of the item is returned. To split the list into two parts, we use the following method –

mid = Lo + ((Hi - Lo) / (A[Hi] - A[Lo])) * (X - A[Lo])

where -

A = list

Lo = Lowest index of the list

Hi = Highest index of the list

A[n] = Value stored at index n in the list

If the middle item is greater than the item, then the probe position is again calculated in the sub-array to the right of the middle item. Otherwise, the item is searched in the subarray to the left of the middle item. This process continues on the sub-array as well until the size of subarray reduces to zero.

Runtime complexity of interpolation search algorithm is $O(\log (\log n))$ as compared to $O(\log n)$ of BST in favorable situations.

Algorithm

As it is an improvisation of the existing BST algorithm, we are mentioning the steps to search the 'target' data value index, using position probing –

Step 1 – Start searching **data** from middle of the list.

Step 2 – If it is a match, return the index of the item, and exit.

Step 3 – If it is not a match, probe position.

Step 4 – Divide the list using probing formula and find the new midle.

Step 5 – If data is greater than middle, search in higher sub-list.

Step 6 – If data is smaller than middle, search in lower sub-list.

Step 7 – Repeat until match.

Pseudocode

 $A \rightarrow Array list$

 $N \rightarrow Size \text{ of } A$

 $X \rightarrow Target Value$

Procedure Interpolation_Search()

Set Lo $\rightarrow 0$

Set Mid \rightarrow -1

Set Hi \rightarrow N-1

While X does not match

if Lo equals to Hi OR A[Lo] equals to A[Hi]

EXIT: Failure, Target not found

end if

Set Mid = Lo + ((Hi - Lo) / (A[Hi] - A[Lo])) * (X - A[Lo])

if A[Mid] = X

EXIT: Success, Target found at Mid

else

if A[Mid] < X

Set Lo to Mid+1

else if A[Mid] > X

Set Hi to Mid-1

end if

end if

End While

End Procedure

Data Structure and Algorithms - Hash Table

Hash Table is a data structure which stores data in an associative manner. In a hash table, data is stored in an array format, where each data value has its own unique index value. Access of data becomes very fast if we know the index of the desired data.

Thus, it becomes a data structure in which insertion and search operations are very fast irrespective of the size of the data. Hash Table uses an array as a storage medium and uses hash technique to generate an index where an element is to be inserted or is to be located from.

Hashing

Hashing is a technique to convert a range of key values into a range of indexes of an array. We're going to use modulo operator to get a range of key values. Consider an example of hash table of size 20, and the following items are to be stored. Item are in the (key,value) format.

	Index	Value
$ \rightarrow \times$	• 0	value_1
key_2	ash 1	value_2
	2	value_3
key 3	3	value 4

200

- (1,20)
- (2,70)(42,80)
- (4,25)
- (12,44)
- (14,32)
- (17,11)
- (13,78)
- (37,98)

Sr. No.	Key	Hash	Array Index
1	1	1 % 20 = 1	1
2	2	2 % 20 = 2	2
3	42	42 % 20 = 2	2
4 150	4	4 % 20 = 4	4
5	12	12 % 20 = 12	12
6	14	14 % 20 = 14	14
7	17	17 % 20 = 17	17

COPYRIGHT FIMT 2020

321 | Page

8	13	13 % 20 = 13	13
9	37	37 % 20 = 17	17

Linear Probing

As we can see, it may happen that the hashing technique is used to create an already used index of the array. In such a case, we can search the next empty location in the array by looking into the next cell until we find an empty cell. This technique is called linear probing.

Sr. No.	Key	Hash	Array Index	After Linear Probing, Array Index
1	1	1 % 20 = 1	1 1 2 2	1
2	2	2 % 20 = 2	2	2
3	42	42 % 20 = 2	2	3
4	4	4 % 20 = 4	4	4
5	12	12 % 20 = 12	12	12
6	14	14 % 20 = 14	14	14
7	17	17 % 20 = 17	17	17
8	13	13 % 20 = 13	13	13
9	37	37 % 20 = 17	17	18

Basic Operations

Following are the basic primary operations of a hash table.

- Search Searches an element in a hash table.
- Insert inserts an element in a hash table.
- **delete** Deletes an element from a hash table.

DataItem

Define a data item having some data and key, based on which the search is to be conducted in a hash table.

```
struct DataItem {
int data;
int key;
```

};

Hash Method

Define a hashing method to compute the hash code of the key of the data item.

int hashCode(int key){

return key % SIZE;

}

Search Operation

Whenever an element is to be searched, compute the hash code of the key passed and locate the element using that hash code as index in the array. Use linear probing to get the element ahead if the element is not found at the computed hash code.

Example

struct DataItem *search(int key) {

//get the hash

int hashIndex = hashCode(key);

//move in array until an empty

while(hashArray[hashIndex] != NULL) {

if(hashArray[hashIndex]->key == key)

return hashArray[hashIndex];

```
//go to next cell
++hashIndex;
//wrap around the table
hashIndex %= SIZE;
}
return NULL;
```

```
}
```

Insert Operation

Whenever an element is to be inserted, compute the hash code of the key passed and locate the index using that hash code as an index in the array. Use linear probing for empty location, if an element is found at the computed hash code.

REIE

-

Example

```
void insert(int key,int data) {
```

struct DataItem *item = (struct DataItem*) malloc(sizeof(struct DataItem));

item->data = data;

item->key = key;

//get the hash
```
int hashIndex = hashCode(key);
```

```
//move in array until an empty or deleted cell
```

while(hashArray[hashIndex] != NULL && hashArray[hashIndex]->key != -1) {

//go to next cell

++hashIndex;

//wrap around the table

hashIndex %= SIZE;

}

hashArray[hashIndex] = item;

}

Delete Operation

Whenever an element is to be deleted, compute the hash code of the key passed and locate the index using that hash code as an index in the array. Use linear probing to get the element ahead if an element is not found at the computed hash code. When found, store a dummy item there to keep the performance of the hash table intact.

0.01.0

Example

struct DataItem* delete(struct DataItem* item) {

000

int key = item->key;

//get the hash

int hashIndex = hashCode(key);

//move i array until an empty

```
while(hashArray[hashIndex] !=NULL) {
 if(hashArray[hashIndex]->key == key) {
   struct DataItem* temp = hashArray[hashIndex];
   //assign a dummy item at deleted position
   hashArray[hashIndex] = dummyItem;
   return temp;
  }
 //go to next cell
 ++hashIndex;
 //wrap around the table
 hashIndex \% = SIZE;
}
return NULL;
```

Data Structure - Sorting Techniques

}

Sorting refers to arranging data in a particular format. Sorting algorithm specifies the way to arrange data in a particular order. Most common orders are in numerical or lexicographical order.

The importance of sorting lies in the fact that data searching can be optimized to a very high level, if data is stored in a sorted manner. Sorting is also used to represent data in more readable formats. Following are some of the examples of sorting in real-life scenarios –

- **Telephone Directory** The telephone directory stores the telephone numbers of people sorted by their names, so that the names can be searched easily.
- **Dictionary** The dictionary stores words in an alphabetical order so that searching of any word becomes easy.

In-place Sorting and Not-in-place Sorting

Sorting algorithms may require some extra space for comparison and temporary storage of few data elements. These algorithms do not require any extra space and sorting is said to happen in-place, or for example, within the array itself. This is called **in-place sorting**. Bubble sort is an example of in-place sorting.

However, in some sorting algorithms, the program requires space which is more than or equal to the elements being sorted. Sorting which uses equal or more space is called **not-in-place sorting**. Merge-sort is an example of not-in-place sorting.

Stable and Not Stable Sorting

If a sorting algorithm, after sorting the contents, does not change the sequence of similar content in which they appear, it is called **stable sorting**.



If a sorting algorithm, after sorting the contents, changes the sequence of similar content in which they appear, it is called **unstable sorting**.



Stability of an algorithm matters when we wish to maintain the sequence of original elements, like in a tuple for example.

Adaptive and Non-Adaptive Sorting Algorithm

A sorting algorithm is said to be adaptive, if it takes advantage of already 'sorted' elements in the list that is to be sorted. That is, while sorting if the source list has some element already sorted, adaptive algorithms will take this into account and will try not to re-order them.

A non-adaptive algorithm is one which does not take into account the elements which are already sorted. They try to force every single element to be re-ordered to confirm their sortedness.

Important Terms

Some terms are generally coined while discussing sorting techniques, here is a brief introduction to them –

Increasing Order

A sequence of values is said to be in **increasing order**, if the successive element is greater than the previous one. For example, 1, 3, 4, 6, 8, 9 are in increasing order, as every next element is greater than the previous element.

Decreasing Order

A sequence of values is said to be in **decreasing order**, if the successive element is less than the current one. For example, 9, 8, 6, 4, 3, 1 are in decreasing order, as every next element is less than the previous element.

Non-Increasing Order

A sequence of values is said to be in **non-increasing order**, if the successive element is less than or equal to its previous element in the sequence. This order occurs when the sequence contains duplicate values. For example, 9, 8, 6, 3, 3, 1 are in non-increasing order, as every next element is less than or equal to (in case of 3) but not greater than any previous element.

Non-Decreasing Order

A sequence of values is said to be in **non-decreasing order**, if the successive element is greater than or equal to its previous element in the sequence. This order occurs when the sequence contains duplicate values. For example, 1, 3, 3, 6, 8, 9 are in non-decreasing order,

as every next element is greater than or equal to (in case of 3) but not less than the previous one.

Data Structure - Bubble Sort Algorithm

Bubble sort is a simple sorting algorithm. This sorting algorithm is comparison-based algorithm in which each pair of adjacent elements is compared and the elements are swapped if they are not in order. This algorithm is not suitable for large data sets as its average and worst case complexity are of $O(n^2)$ where **n** is the number of items.

How Bubble Sort Works?

We take an unsorted array for our example. Bubble sort takes $O(n^2)$ time so we're keeping it short and precise.



Bubble sort starts with very first two elements, comparing them to check which one is greater.



In this case, value 33 is greater than 14, so it is already in sorted locations. Next, we compare 33 with 27.



We find that 27 is smaller than 33 and these two values must be swapped.



Next we compare 33 and 35. We find that both are in already sorted positions.



Then we move to the next two values, 35 and 10.



We know then that 10 is smaller 35. Hence they are not sorted.



We swap these values. We find that we have reached the end of the array. After one iteration, the array should look like this –



To be precise, we are now showing how an array should look like after each iteration. After the second iteration, it should look like this –



Notice that after each iteration, at least one value moves at the end.

4.4	10	07	00	00
14	10	21	33	30

And when there's no swap required, bubble sorts learns that an array is completely sorted.



Now we should look into some practical aspects of bubble sort.

Algorithm

We assume **list** is an array of **n** elements. We further assume that **swap**function swaps the values of the given array elements.

```
begin BubbleSort(list)
```

for all elements of list

if list[i] > list[i+1]

swap(list[i], list[i+1])

end if

end for

return list

end BubbleSort

Pseudocode

We observe in algorithm that Bubble Sort compares each pair of array element unless the whole array is completely sorted in an ascending order. This may cause a few complexity issues like what if the array needs no more swapping as all the elements are already ascending.

To ease-out the issue, we use one flag variable **swapped** which will help us see if any swap has happened or not. If no swap has occurred, i.e. the array requires no more processing to be sorted, it will come out of the loop.

Pseudocode of BubbleSort algorithm can be written as follows -

```
procedure bubbleSort( list : array of items )
```

loop = list.count;

for i = 0 to loop-1 do:

swapped = false

for j = 0 to loop-1 do:

/* compare the adjacent elements */

if list[j] > list[j+1] then

/* swap them */

swap(list[j], list[j+1])

COPYRIGHT FIMT 2020

331 | Page

```
swapped = true
```

end if

end for

/*if no number was swapped that means

array is sorted now, break the loop.*/

if(not swapped) then

break

end if

end for

end procedure return list

Implementation

One more issue we did not address in our original algorithm and its improvised pseudocode, is that, after every iteration the highest values settles down at the end of the array. Hence, the next iteration need not include already sorted elements. For this purpose, in our implementation, we restrict the inner loop to avoid already sorted values.

Data Structure and Algorithms Insertion Sort

This is an in-place comparison-based sorting algorithm. Here, a sub-list is maintained which is always sorted. For example, the lower part of an array is maintained to be sorted. An element which is to be 'insert'ed in this sorted sub-list, has to find its appropriate place and then it has to be inserted there. Hence the name, **insertion sort**.

The array is searched sequentially and unsorted items are moved and inserted into the sorted sub-list (in the same array). This algorithm is not suitable for large data sets as its average and worst case complexity are of $O(n^2)$, where **n** is the number of items.

How Insertion Sort Works?

We take an unsorted array for our example.



Insertion sort compares the first two elements.



It finds that both 14 and 33 are already in ascending order. For now, 14 is in sorted sub-list.



Insertion sort moves ahead and compares 33 with 27.

14	33	27	10	35	19	42	44
	\square						

And finds that 33 is not in the correct position.

14	33	27	10	35	19	42	44
	\square						

It swaps 33 with 27. It also checks with all the elements of sorted sub-list. Here we see that the sorted sub-list has only one element 14, and 27 is greater than 14. Hence, the sorted sub-list remains sorted after swapping.



By now we have 14 and 27 in the sorted sub-list. Next, it compares 33 with 10.

14 27 33 10 35 19 42 44	14	27	33	10	35	19	42	44
-------------------------	----	----	----	----	----	----	----	----

These values are not in a sorted order.



So we swap them.

1	14	27	10	33	35	19	42	44
Ų				\square	\square	\square	\square	

However, swapping makes 27 and 10 unsorted.



Hence, we swap them too.



Again we find 14 and 10 in an unsorted order.



We swap them again. By the end of third iteration, we have a sorted sub-list of 4 items.



This process goes on until all the unsorted values are covered in a sorted sub-list. Now we shall see some programming aspects of insertion sort.

Algorithm

Now we have a bigger picture of how this sorting technique works, so we can derive simple steps by which we can achieve insertion sort.

Step 1 – If it is the first element, it is already sorted. return 1;

Step 2 – Pick next element

- Step 3 Compare with all elements in the sorted sub-list
- Step 4 Shift all the elements in the sorted sub-list that is greater than the value to be sorted

Step 5 – Insert the value

Step 6 – Repeat until list is sorted

Pseudocode

procedure insertionSort(A : array of items)

int holePosition

int valueToInsert

for i = 1 to length(A) inclusive do:

/* select value to be inserted */

```
valueToInsert = A[i]
```

holePosition = i

/*locate hole position for the element to be inserted */

while holePosition > 0 and A[holePosition-1] > valueToInsert do:

A[holePosition] = A[holePosition-1]

holePosition = holePosition -1

end while

/* insert the number at hole position */

A[holePosition] = valueToInsert

end for

end procedure

Data Structure and Algorithms Selection Sort

Selection sort is a simple sorting algorithm. This sorting algorithm is an in-place comparison-based algorithm in which the list is divided into two parts, the sorted part at the left end and the unsorted part at the right end. Initially, the sorted part is empty and the unsorted part is the entire list.

The smallest element is selected from the unsorted array and swapped with the leftmost element, and that element becomes a part of the sorted array. This process continues moving unsorted array boundary by one element to the right.

This algorithm is not suitable for large data sets as its average and worst case complexities are of $O(n^2)$, where **n** is the number of items.

How Selection Sort Works?

Consider the following depicted array as an example.



For the first position in the sorted list, the whole list is scanned sequentially. The first position where 14 is stored presently, we search the whole list and find that 10 is the lowest value.



So we replace 14 with 10. After one iteration 10, which happens to be the minimum value in the list, appears in the first position of the sorted list.



For the second position, where 33 is residing, we start scanning the rest of the list in a linear manner.



We find that 14 is the second lowest value in the list and it should appear at the second place. We swap these values.



After two iterations, two least values are positioned at the beginning in a sorted manner.



The same process is applied to the rest of the items in the array.

Following is a pictorial depiction of the entire sorting process -



Now, let us learn some programming aspects of selection sort.

Algorithm

Step 1 – Set MIN to location 0

Step 2 – Search the minimum element in the list

Step 3 – Swap with value at location MIN

Step 4 – Increment MIN to point to next element

Step 5 – Repeat until list is sorted

Pseudocode

procedure selection sort list : array of items n : size of list for i = 1 to n - 1/* set current element as minimum*/ min = i/* check the element to be minimum */ for j = i+1 to n if list[j] < list[min] then min = j;end if end for /* swap the minimum element with the current element*/ if indexMin != i then

```
swap list min and list i
  end if
end for
       end procedure
```

Data Structures - Merge Sort Algorithm

AGEMEN Merge sort is a sorting technique based on divide and conquer technique. With worst-case time complexity being $O(n \log n)$, it is one of the most respected algorithms.

Merge sort first divides the array into equal halves and then combines them in a sorted manner.

How Merge Sort Works?

To understand merge sort, we take an unsorted array as the following



We know that merge sort first divides the whole array iteratively into equal halves unless the atomic values are achieved. We see here that an array of 8 items is divided into two arrays of size 4.



This does not change the sequence of appearance of items in the original. Now we divide these two arrays into halves.



We further divide these arrays and we achieve atomic value which can no more be divided.



Now, we combine them in exactly the same manner as they were broken down. Please note the color codes given to these lists.

We first compare the element for each list and then combine them into another list in a sorted manner. We see that 14 and 33 are in sorted positions. We compare 27 and 10 and in the target list of 2 values we put 10 first, followed by 27. We change the order of 19 and 35 whereas 42 and 44 are placed sequentially.



In the next iteration of the combining phase, we compare lists of two data values, and merge them into a list of found data values placing all in a sorted order.

After the final merging, the list should look like this

\square							
10	14	19	27	33	35	42	44
	\square				\square		

Now we should learn some programming aspects of merge sorting.

Algorithm

Merge sort keeps on dividing the list into equal halves until it can no more be divided. By definition, if it is only one element in the list, it is sorted. Then, merge sort combines the smaller sorted lists keeping the new list sorted too.

Step 1 -if it is only one element in the list it is already sorted, return.

Step 2 – divide the list recursively into two halves until it can no more be divided.

Step 3 – merge the smaller lists into new list in sorted order.

Pseudocode

We shall now see the pseudocodes for merge sort functions. As our algorithms point out two main functions – divide & merge.

Merge sort works with recursion and we shall see our implementation in the same way.

procedure mergesort(var a as array)

if (n == 1) return a

```
var 11 as array = a[0] \dots a[n/2]
```

```
var l2 as array = a[n/2+1] \dots a[n]
```

```
l1 = mergesort(l1)
```

```
l2 = mergesort(l2)
```

return merge(11,12)

end procedure

procedure merge(var a as array, var b as array)

var c as array

while (a and b have elements)

if (a[0] > b[0])

add b[0] to the end of c

```
remove b[0] from b
```

else

add a[0] to the end of c

remove a[0] from a

end if

end while

while (a has elements)

```
add a[0] to the end of c
```

```
remove a[0] from a
```

end while

while (b has elements)

add b[0] to the end of c

remove b[0] from b

end while

return c

end procedure

Data Structure and Algorithms-Shell Sort

Shell sort is a highly efficient sorting algorithm and is based on insertion sort algorithm. This algorithm avoids large shifts as in case of insertion sort, if the smaller value is to the far right and has to be moved to the far left.

This algorithm uses insertion sort on a widely spread elements, first to sort them and then sorts the less widely spaced elements. This spacing is termed as **interval**. This interval is calculated based on Knuth's formula as -

Knuth's Formula

h = h * 3 + 1

where -

h is interval with initial value 1

This algorithm is quite efficient for medium-sized data sets as its average and worst case complexity are of O(n), where **n** is the number of items.

How Shell Sort Works?

Let us consider the following example to have an idea of how shell sort works. We take the same array we have used in our previous examples. For our example and ease of understanding, we take the interval of 4. Make a virtual sub-list of all values located at the interval of 4 positions. Here these values are {35, 14}, {33, 19}, {42, 27} and {10, 44}



We compare values in each sub-list and swap them (if necessary) in the original array. After this step, the new array should look like this –



Then, we take interval of 2 and this gap generates two sub-lists - $\{14, 27, 35, 42\}$, $\{19, 10, 33, 44\}$



We compare and swap the values, if required, in the original array. After this step, the array should look like this -



Finally, we sort the rest of the array using interval of value 1. Shell sort uses insertion sort to sort the array.

Following	is	the		step-by-step		depiction	—
14	19][27][10][35][33	42	44
14	19	27][10][35][33	42	44
14	19	27	10][35][33	42	44
14	19	27	10	35	33	42	44
14	19	10	27	35	33	42	44
14	10	19	27	35	33	42	44
10	14	19	27	35	33	42	44
10	14	19	27	35	33	42	44
10	14	19	27	33	35	42	44
10	14	19	27] 33]	35	42	44

We see that it required only four swaps to sort the rest of the array.

Algorithm

Following is the algorithm for shell sort.

Step 1 – Initialize the value of *h*

Step 2 – Divide the list into smaller sub-list of equal interval h

Step 3 – Sort these sub-lists using insertion sort

Step 3 – Repeat until complete list is sorted

Pseudocode

Following is the pseudocode for shell sort.

procedure shellSort()

A : array of items

/* calculate interval*/

while interval < A.length /3 do:

interval = interval * 3 + 1

end while

while interval > 0 do:

for outer = interval; outer < A.length; outer ++ do:

/* seect value to be inserted */

valueToInsert = A[outer]

inner = outer;

/*shift element towards right*/

while inner > interval -1 && A[inner - interval] >= valueToInsert do:

A[inner] = A[inner - interval]

inner = inner - interval

end while

/* insert the number at hole position */

A[inner] = valueToInsert

end for

/* calculate interval*/

interval = (interval -1) /3;

end while

end procedure

Data Structure and Algorithms-Quick Sort

Quick sort is a highly efficient sorting algorithm and is based on partitioning of array of data into smaller arrays. A large array is partitioned into two arrays one of which holds values smaller than the specified value, say pivot, based on which the partition is made and another array holds values greater than the pivot value.

Quick sort partitions an array and then calls itself recursively twice to sort the two resulting subarrays. This algorithm is quite efficient for large-sized data sets as its average and worst case complexity are of $O(n^2)$, where **n** is the number of items. ADATE

Partition in Quick Sort

Following animated representation explains how to find the pivot value in an array.

Unsorted Array

42] [10] [14] [19] [35 33 27 44 26 || 31

The pivot value divides the list into two parts. And recursively, we find the pivot for each sub-lists until all lists contains only one element.

Quick Sort Pivot Algorithm

Based on our understanding of partitioning in quick sort, we will now try to write an algorithm for it, which is as follows.

- Step 1 Choose the highest index value has pivot
- **Step 2** Take two variables to point left and right of the list excluding pivot
- **Step 3** left points to the low index
- **Step 4** right points to the high
- **Step 5** while value at left is less than pivot move right
- **Step 6** while value at right is greater than pivot move left
- Step 7 if both step 5 and step 6 does not match swap left and right
- **Step 8** if left \geq right, the point where they met is new pivot

Quick Sort Pivot Pseudocode

The pseudocode for the above algorithm can be derived as -

function partitionFunc(left, right, pivot)

leftPointer = left

rightPointer = right - 1

while True do

while A[++leftPointer] < pivot do

//do-nothing

end while

while rightPointer > 0 && A[--rightPointer] > pivot do

//do-nothing

end while

if leftPointer >= rightPointer

break

else

swap leftPointer,rightPointer

end if

end while

swap leftPointer,right

return leftPointer

end function

Quick Sort Algorithm

Using pivot algorithm recursively, we end up with smaller possible partitions. Each partition is then processed for quick sort. We define recursive algorithm for quicksort as follows –

Step 1 – Make the right-most index value pivot

Step 2 – partition the array using pivot value

Step 3 – quicksort left partition recursively

Step 4 – quicksort right partition recursively

Quick Sort Pseudocode

To get more into it, let see the pseudocode for quick sort algorithm -

procedure quickSort(left, right)

if right-left <= 0

return

else

pivot = A[right]

partition = partitionFunc(left, right, pivot)

quickSort(left,partition-1)

```
quickSort(partition+1,right)
```

end if

end procedure

Data Structure - Graph Data Structure

A graph is a pictorial representation of a set of objects where some pairs of objects are connected by links. The interconnected objects are represented by points termed as **vertices**, and the links that connect the vertices are called **edges**.

Formally, a graph is a pair of sets (V, E), where V is the set of vertices and E is the set of edges, connecting the pairs of vertices. Take a look at the following graph -



In the above graph,

 $V = \{a, b, c, d, e\}$

 $E = \{ab, ac, bd, cd, de\}$

Graph Data Structure

Mathematical graphs can be represented in data structure. We can represent a graph using an array of vertices and a two-dimensional array of edges. Before we proceed further, let's familiarize ourselves with some important terms –

- Vertex Each node of the graph is represented as a vertex. In the following example, the labeled circle represents vertices. Thus, A to G are vertices. We can represent them using an array as shown in the following image. Here A can be identified by index 0. B can be identified using index 1 and so on.
- Edge Edge represents a path between two vertices or a line between two vertices. In the following example, the lines from A to B, B to C, and so on represents edges. We can use a two-dimensional array to represent an array as shown in the following image. Here AB can be represented as 1 at row 0, column 1, BC as 1 at row 1, column 2 and so on, keeping other combinations as 0.
- Adjacency Two node or vertices are adjacent if they are connected to each other through an edge. In the following example, B is adjacent to A, C is adjacent to B, and so on.
- **Path** Path represents a sequence of edges between the two vertices. In the following example, ABCD represents a path from A to D.



Basic Operations

Following are basic primary operations of a Graph –

- Add Vertex Adds a vertex to the graph.
- Add Edge Adds an edge between the two vertices of the graph.
- **Display Vertex** Displays a vertex of the graph.

Data Structure - Depth First Traversal

Depth First Search (DFS) algorithm traverses a graph in a depthward motion and uses a stack to remember to get the next vertex to start a search, when a dead end occurs in any



iteration.

As in the example given above, DFS algorithm traverses from S to A to D to G to E to B first, then to F and lastly to G. It employs the following rules.

- Rule 1 Visit the adjacent unvisited vertex. Mark it as visited. Display it. Push it in a stack.
- **Rule 2** If no adjacent vertex is found, pop up a vertex from the stack. (It will pop up all the vertices from the stack, which do not have adjacent vertices.)
- **Rule 3** Repeat Rule 1 and Rule 2 until the stack is empty.

	Ste	Traversal	Description
--	-----	-----------	-------------





COPYRIGHT FIMT 2020

352 | Page

we

are

As **C** does not have any unvisited adjacent node so we keep popping the stack until we find a node that has an unvisited adjacent node. In this case, there's none and we keep popping until the stack is empty.

Data Structure - Breadth First Traversal

Breadth First Search (BFS) algorithm traverses a graph in a breadthward motion and uses a queue to remember to get the next vertex to start a search, when a dead end occurs in any iteration.



As in the example given above, BFS algorithm traverses from A to B to E to F first then to C and G lastly to D. It employs the following rules.

- Rule 1 Visit the adjacent unvisited vertex. Mark it as visited. Display it. Insert it in a queue.
- Rule 2 If no adjacent vertex is found, remove the first vertex from the queue.
- Rule 3 Repeat Rule 1 and Rule 2 until the queue is empty.







At this stage, we are left with no unmarked (unvisited) nodes. But as per the algorithm we keep on dequeuing in order to get all unvisited nodes. When the queue gets emptied, the program is over.

Data Structure and Algorithms-Tree

Tree represents the nodes connected by edges. We will discuss binary tree or binary search tree specifically.

Binary Tree is a special datastructure used for data storage purposes. A binary tree has a special condition that each node can have a maximum of two children. A binary tree has the benefits of both an ordered array and a linked list as search is as quick as in a sorted array and insertion or deletion operation are as fast as in linked list.



Important Terms

Following are the important terms with respect to tree.

- Path Path refers to the sequence of nodes along the edges of a tree.
- **Root** The node at the top of the tree is called root. There is only one root per tree and one path from the root node to any node.
- **Parent** Any node except the root node has one edge upward to a node called parent.
- Child The node below a given node connected by its edge downward is called its child node.
- Leaf The node which does not have any child node is called the leaf node.
- Subtree Subtree represents the descendants of a node.
- Visiting Visiting refers to checking the value of a node when control is on the node.
- Traversing Traversing means passing through nodes in a specific order.
- Levels Level of a node represents the generation of a node. If the root node is at level 0, then its next child node is at level 1, its grandchild is at level 2, and so on.
- **keys** Key represents a value of a node based on which a search operation is to be carried out for a node.

Binary Search Tree Representation

Binary Search tree exhibits a special behavior. A node's left child must have a value less than its parent's value and the node's right child must have a value greater than its parent value.



We're going to implement tree using node object and connecting them through references.

Tree Node

The code to write a tree node would be similar to what is given below. It has a data part and references to its left and right child nodes.

struct node {

int data;

struct node *leftChild;

struct node *rightChild;

};

In a tree, all nodes share common construct.

BST Basic Operations

The basic operations that can be performed on a binary search tree data structure, are the following –

- Insert Inserts an element in a tree/create a tree.
- Search Searches an element in a tree.
- Preorder Traversal Traverses a tree in a pre-order manner.
- Inorder Traversal Traverses a tree in an in-order manner.
- **Postorder Traversal** Traverses a tree in a post-order manner.

We shall learn creating (inserting into) a tree structure and searching a data item in a tree in this chapter. We shall learn about tree traversing methods in the coming chapter.

Insert Operation

The very first insertion creates the tree. Afterwards, whenever an element is to be inserted, first locate its proper location. Start searching from the root node, then if the data is less than the key value, search for the empty location in the left subtree and insert the data. Otherwise, search for the empty location in the right subtree and insert the data.

Algorithm
If root is NULL
then create root node
return
If root exists then

compare the data with node.data

while until insertion position is located

If data is greater than node.data

goto right subtree

else

goto left subtree

endwhile

insert data

end If

Implementation

The implementation of insert function should look like this -

void insert(int data) {

struct node *tempNode = (struct node*) malloc(sizeof(struct node));

struct node *current;

struct node *parent;

tempNode->data = data;

tempNode->leftChild = NULL;

tempNode->rightChild = NULL;

COPYRIGHT FIMT 2020

358 | Page

//if tree is empty, create root node

```
if(root == NULL) {
```

root = tempNode;

} else {

current = root;

parent = NULL;

while(1) {

parent = current;

//go to left of the tree

```
if(data < parent->data) {
```

```
current = current->leftChild;
```

//insert to the left

if(current == NULL) {

parent->leftChild = tempNode;

return;

}

}

```
//go to right of the tree
     else {
       current = current->rightChild;
               //insert to the right
       if(current == NULL) {
         parent->rightChild = tempNode;
         return;
        }
    }
 }
}
```

Search Operation

Whenever an element is to be searched, start searching from the root node, then if the data is less than the key value, search for the element in the left subtree. Otherwise, search for the element in the right subtree. Follow the same algorithm for each node.

Algorithm

If root.data is equal to search.data

return root
else

while data not found

If data is greater than node.data

goto right subtree

else

goto left subtree

If data found

return node

endwhile

return data not found

end if



The implementation of this algorithm should look like this.

×.

struct node* search(int data) {

ULL.

struct node *current = root;

printf("Visiting elements: ");

while(current->data != data) {

```
if(current != NULL)
```

```
printf("%d ",current->data);
```

//go to left tree

if(current->data > data) {

current = current->leftChild;

}

//else go to right tree

else {

current = current->rightChild;

```
}
```

//not found

if(current == NULL) {

return NULL;

```
}
```

return current;

}

}

Data Structure & Algorithms-Tree Traversal

Traversal is a process to visit all the nodes of a tree and may print their values too. Because, all nodes are connected via edges (links) we always start from the root (head) node. That is, we cannot randomly access a node in a tree. There are three ways which we use to traverse a tree -

- In-order Traversal
- Pre-order Traversal
- Post-order Traversal

Generally, we traverse a tree to search or locate a given item or key in the tree or to print all the values it contains.

AGEMEN

In-order Traversal

In this traversal method, the left subtree is visited first, then the root and later the right subtree. We should always remember that every node may represent a subtree itself.

If a binary tree is traversed **in-order**, the output will produce sorted key values in an ascending order.



We start from **A**, and following in-order traversal, we move to its left subtree **B**. **B** is also traversed in-order. The process goes on until all the nodes are visited. The output of inorder traversal of this tree will be -

 $D \to B \to E \to A \to F \to C \to G$

Algorithm

Until all nodes are traversed -

Step 1 – Recursively traverse left subtree.

Step 2 – Visit root node.

Step 3 – Recursively traverse right subtree.

Pre-order Traversal

In this traversal method, the root node is visited first, then the left subtree and finally the



We start from **A**, and following pre-order traversal, we first visit **A** itself and then move to its left subtree **B**. **B** is also traversed pre-order. The process goes on until all the nodes are visited. The output of pre-order traversal of this tree will be -

 $A \to B \to D \to E \to C \to F \to G$

Algorithm

Until all nodes are traversed -

Step 1 – Visit root node.

Step 2 – Recursively traverse left subtree.

Step 3 – Recursively traverse right subtree.

Post-order Traversal

In this traversal method, the root node is visited last, hence the name. First we traverse the left subtree, then the right subtree and finally the root node.



We start from **A**, and following Post-order traversal, we first visit the left subtree **B**. **B** is also traversed post-order. The process goes on until all the nodes are visited. The output of post-order traversal of this tree will be -

$$D \to E \to B \to F \to G \to C \to A$$

Algorithm

Until all nodes are traversed -

Step 1 – Recursively traverse left subtree.

Step 2 – Recursively traverse right subtree.

Step 3 – Visit root node.

Data Structure - Binary Search Tree

A Binary Search Tree (BST) is a tree in which all the nodes follow the below-mentioned properties –

- The left sub-tree of a node has a key less than or equal to its parent node's key.
- The right sub-tree of a node has a key greater than to its parent node's key.

Thus, BST divides all its sub-trees into two segments; the left sub-tree and the right sub-tree and can be defined as –

 $left_subtree (keys) \le node (key) \le right_subtree (keys)$

Representation

BST is a collection of nodes arranged in a way where they maintain BST properties. Each node has a key and an associated value. While searching, the desired key is compared to the keys in BST and if found, the associated value is retrieved.

Following is a pictorial representation of BST -



We observe that the root node key (27) has all less-valued keys on the left sub-tree and the higher valued keys on the right sub-tree.

Basic Operations

Following are the basic operations of a tree -

- Search Searches an element in a tree.
- Insert Inserts an element in a tree.
- Pre-order Traversal Traverses a tree in a pre-order manner.
- In-order Traversal Traverses a tree in an in-order manner.
- **Post-order Traversal** Traverses a tree in a post-order manner.

Node

Define a node having some data, references to its left and right child nodes.

struct node {

int data;

struct node *leftChild;

struct node *rightChild;

};

Search Operation

Whenever an element is to be searched, start searching from the root node. Then if the data

is less than the key value, search for the element in the left subtree. Otherwise, search for the element in the right subtree. Follow the same algorithm for each node.

Algorithm

struct node* search(int data){

struct node *current = root;

printf("Visiting elements: ");

while(current->data != data){

if(current != NULL) {

printf("%d ",current->data);

//go to left tree

if(current->data > data){

current = current->leftChild;

}//else go to right tree

else {

current = current->rightChild;

}

//not found

if(current == NULL){

return NULL;

}

```
}
}
return current;
}
```

Insert Operation

Whenever an element is to be inserted, first locate its proper location. Start searching from the root node, then if the data is less than the key value, search for the empty location in the left subtree and insert the data. Otherwise, search for the empty location in the right subtree and insert the data.

Algorithm

void insert(int data) {

struct node *tempNode = (struct node*) malloc(sizeof(struct node));

struct node *current;

struct node *parent;

tempNode->data = data;

tempNode->leftChild = NULL;

tempNode->rightChild = NULL;

//if tree is empty

if(root == NULL) {

root = tempNode;

} else {

current = root;

parent = NULL;

while(1) {

parent = current;

//go to left of the tree

if(data < parent->data) {

current = current->leftChild;

//insert to the left

if(current == NULL) {

parent->leftChild = tempNode;

return;

}

}//go to right of the tree

else {

current = current->rightChild;



What if the input to binary search tree comes in a sorted (ascending or descending) manner? It will then look like this –



If input 'appears' non-increasing manner

If input 'appears' in non-decreasing manner

It is observed that BST's worst-case performance is closest to linear search algorithms, that is O(n). In real-time data, we cannot predict data pattern and their frequencies. So, a need arises to balance out the existing BST.

Named after their inventor **Adelson**, **Velski & Landis**, **AVL trees** are height balancing binary search tree. AVL tree checks the height of the left and the right sub-trees and assures that the difference is not more than 1. This difference is called the **Balance Factor**. Here we see that the first tree is balanced and the next two trees are not balanced –



In the second tree, the left subtree of C has height 2 and the right subtree has height 0, so the difference is 2. In the third tree, the right subtree of A has height 2 and the left is missing, so it is 0, and the difference is 2 again. AVL tree permits difference (balance factor) to be only 1.

BalanceFactor = height(left-sutree) - height(right-sutree)

If the difference in the height of left and right sub-trees is more than 1, the tree is balanced using some rotation techniques.

AVL Rotations

To balance itself, an AVL tree may perform the following four kinds of rotations -

- Left rotation
- Right rotation
- Left-Right rotation
- Right-Left rotation

The first two rotations are single rotations and the next two rotations are double rotations. To have an unbalanced tree, we at least need a tree of height 2. With this simple tree, let's understand them one by one.

Left Rotation

If a tree becomes unbalanced, when a node is inserted into the right subtree of the right subtree, then we perform a single left rotation -



In our example, node **A** has become unbalanced as a node is inserted in the right subtree of A's right subtree. We perform the left rotation by making Athe left-subtree of B.

Right Rotation

AVL tree may become unbalanced, if a node is inserted in the left subtree of the left subtree. The tree then needs a right rotation.



Left unbalanced Tree

Right Rotation

Balanced Tree

As depicted, the unbalanced node becomes the right child of its left child by performing a right rotation.

Left-Right Rotation

Double rotations are slightly complex version of already explained versions of rotations. To understand them better, we should take note of each action performed while rotation. Let's first check how to perform Left-Right rotation. A left-right rotation is a combination of left rotation followed by right rotation.

State	Action
2 9001: 1 A B	A node has been inserted into the right subtree of the left subtree. This makes C an unbalanced node. These scenarios cause AVL tree to perform left-right rotation.



Right-Left Rotation

The second type of double rotation is Right-Left Rotation. It is a combination of right rotation followed by left rotation.

State	Action
-------	--------



Data Structure & Algorithms-Spanning Tree

A spanning tree is a subset of Graph G, which has all the vertices covered with minimum possible number of edges. Hence, a spanning tree does not have cycles and it cannot be disconnected..

By this definition, we can draw a conclusion that every connected and undirected Graph G has at least one spanning tree. A disconnected graph does not have any spanning tree, as it cannot be spanned to all its vertices.



We found three spanning trees off one complete graph. A complete undirected graph can have maximum n^{n-2} number of spanning trees, where **n** is the number of nodes. In the above addressed example, $3^{3-2} = 3$ spanning trees are possible.

General Properties of Spanning Tree

We now understand that one graph can have more than one spanning tree. Following are a few properties of the spanning tree connected to graph G –

- A connected graph G can have more than one spanning tree.
- All possible spanning trees of graph G, have the same number of edges and vertices.
- The spanning tree does not have any cycle (loops).
- Removing one edge from the spanning tree will make the graph disconnected, i.e. the spanning tree is **minimally connected**.
- Adding one edge to the spanning tree will create a circuit or loop, i.e. the spanning tree is **maximally acyclic**.

Mathematical Properties of Spanning Tree

- Spanning tree has **n-1** edges, where **n** is the number of nodes (vertices).
- From a complete graph, by removing maximum e n + 1 edges, we can construct a spanning tree.
- A complete graph can have maximum n^{n-2} number of spanning trees.

Thus, we can conclude that spanning trees are a subset of connected Graph G and disconnected graphs do not have spanning tree.

Application of Spanning Tree

Spanning tree is basically used to find a minimum path to connect all nodes in a graph. Common application of spanning trees are –

- Civil Network Planning
- Computer Network Routing Protocol
- Cluster Analysis

Let us understand this through a small example. Consider, city network as a huge graph and now plans to deploy telephone lines in such a way that in minimum lines we can connect to all city nodes. This is where the spanning tree comes into picture.

FUITED

Minimum Spanning Tree (MST)

In a weighted graph, a minimum spanning tree is a spanning tree that has minimum weight than all other spanning trees of the same graph. In real-world situations, this weight can be measured as distance, congestion, traffic load or any arbitrary value denoted to the edges.

Minimum Spanning-Tree Algorithm

We shall learn about two most important spanning tree algorithms here -

- Kruskal's Algorithm
- Prim's Algorithm

Both are greedy algorithms

Heap Data Structures

Heap is a special case of balanced binary tree data structure where the root-node key is compared with its children and arranged accordingly. If α has child node β then –

$key(\alpha) \ge key(\beta)$

As the value of parent is greater than that of child, this property generates Max Heap. Based on this criteria, a heap can be of two types -

For Input \rightarrow 35 33 42 10 14 19 27 44 26 31

Min-Heap – Where the value of the root node is less than or equal to either of its children.



Max-Heap – Where the value of the root node is greater than or equal to either of its children.



Both trees are constructed using the same input and order of arrival.

Max Heap Construction Algorithm

We shall use the same example to demonstrate how a Max Heap is created. The procedure to create Min Heap is similar but we go for min values instead of max values.

We are going to derive an algorithm for max heap by inserting one element at a time. At any point of time, heap must maintain its property. While insertion, we also assume that we are inserting a node in an already heapified tree.

- **Step 1** Create a new node at the end of heap.
- **Step 2** Assign new value to the node.
- Step 3 Compare the value of this child node with its parent.
- Step 4 If value of parent is less than child, then swap them.
- **Step 5** Repeat step 3 & 4 until Heap property holds.

Note – In Min Heap construction algorithm, we expect the value of the parent node to be less than that of the child node.

Let's understand Max Heap construction by an animated illustration. We consider the same input sample that we used earlier.

Max Heap Deletion Algorithm

Let us derive an algorithm to delete from max heap. Deletion in Max (or Min) Heap always happens at the root to remove the Maximum (or minimum) value.

- Step 1 Remove root node.
- **Step 2** Move the last element of last level to root.
- **Step 3** Compare the value of this child node with its parent.
- **Step 4** If value of parent is less than child, then swap them.
- **Step 5** Repeat step 3 & 4 until Heap property holds.



Data Structure-Recursion Basics

Some computer programming languages allow a module or function to call itself. This technique is known as recursion. In recursion, a function α either calls itself directly or calls a function β that in turn calls the original function α . The function α is called recursive function.

Example – a function calling itself.

```
int function(int value) {
    if(value < 1)
    return;</pre>
```

```
function(value - 1);
```

printf("%d ",value);

}

Example – a function that calls another function which in turn calls it again.

AAC ACCREI

```
int function(int value) {
    if(value < 1)</pre>
```

return;

```
function(value - 1);
```

printf("%d ",value);

}

Properties

A recursive function can go infinite like a loop. To avoid infinite running of recursive function, there are two properties that a recursive function must have –

- **Base criteria** There must be at least one base criteria or condition, such that, when this condition is met the function stops calling itself recursively.
- **Progressive approach** The recursive calls should progress in such a way that each time a recursive call is made it comes closer to the base criteria.

Implementation

Many programming languages implement recursion by means of **stacks**. Generally, whenever a function (**caller**) calls another function (**callee**) or itself as callee, the caller function transfers execution control to the callee. This transfer process may also involve some data to be passed from the caller to the callee. This implies, the caller function has to suspend its execution temporarily and resume later when the execution control returns from the callee function. Here, the caller function needs to start exactly from the point of execution where it puts itself on hold. It also needs the exact same data values it was

working on. For this purpose, an activation record (or stack frame) is created for the caller function.



This activation record keeps the information about local variables, formal parameters, return address and all information passed to the caller function.

Analysis of Recursion

One may argue why to use recursion, as the same task can be done with iteration. The first reason is, recursion makes a program more readable and because of latest enhanced CPU systems, recursion is more efficient than iterations.

Time Complexity

In case of iterations, we take number of iterations to count the time complexity. Likewise, in case of recursion, assuming everything is constant, we try to figure out the number of times a recursive call is being made. A call made to a function is O(1), hence the (n) number of times a recursive call is made makes the recursive function O(n).

Space Complexity

Space complexity is counted as what amount of extra space is required for a module to execute. In case of iterations, the compiler hardly requires any extra space. The compiler keeps updating the values of variables used in the iterations. But in case of recursion, the system needs to store activation record each time a recursive call is made. Hence, it is considered that space complexity of recursive function may go higher than that of a function with iteration.

Data Structure & Algorithms-Tower of Hanoi

Tower of Hanoi, is a mathematical puzzle which consists of three towers (pegs) and more than one rings is as depicted –



These rings are of different sizes and stacked upon in an ascending order, i.e. the smaller one sits over the larger one. There are other variations of the puzzle where the number of disks increase, but the tower count remains the same.

Rules

The mission is to move all the disks to some another tower without violating the sequence of arrangement. A few rules to be followed for Tower of Hanoi are –

AGEMEA

- Only one disk can be moved among the towers at any given time.
- Only the "top" disk can be removed.
- No large disk can sit over a small disk.

Following is an animated representation of solving a Tower of Hanoi puzzle with three disks.

Step: 0



Tower of Hanoi puzzle with n disks can be solved in minimum 2^n-1 steps. This presentation shows that a puzzle with 3 disks has taken $2^3 - 1 = 7$ steps.

Algorithm

To write an algorithm for Tower of Hanoi, first we need to learn how to solve this problem with lesser amount of disks, say $\rightarrow 1$ or 2. We mark three towers with name, **source**, **destination** and **aux** (only to help moving the disks). If we have only one disk, then it can easily be moved from source to destination peg.

If we have 2 disks -

- First, we move the smaller (top) disk to aux peg.
- Then, we move the larger (bottom) disk to destination peg.

• And finally, we move the smaller disk from aux to destination peg.



So now, we are in a position to design an algorithm for Tower of Hanoi with more than two disks. We divide the stack of disks in two parts. The largest disk $(n^{th} disk)$ is in one part and all other (n-1) disks are in the second part.

Our ultimate aim is to move disk **n** from source to destination and then put all other (n1) disks onto it. We can imagine to apply the same in a recursive way for all given set of disks.

The steps to follow are –

Step 1 – Move n-1 disks from source to aux

Step 2 – Move n^{th} disk from **source** to **dest**

Step 3 – Move n-1 disks from aux to dest

A recursive algorithm for Tower of Hanoi can be driven as follows –

START		
Procedure Hanoi(disk, source, dest, aux)		
IF disk == 1, THEN		
move disk from source to dest		
ELSE		
Hanoi(disk - 1, source, aux, dest)	// Step 1	
move disk from source to dest	// Step 2	
Hanoi(disk - 1, aux, dest, source)	// Step 3	
END IF		
END Procedure		
STOP		

Data Structure & Algorithms Fibonacci Series

Fibonacci series generates the subsequent number by adding two previous numbers. Fibonacci series starts from two numbers $-\mathbf{F_0} \& \mathbf{F_1}$. The initial values of $F_0 \& F_1$ can be taken 0, 1 or 1, 1 respectively. Fibonacci series satisfies the following conditions -

 $F_n = F_{n\text{-}1} + F_{n\text{-}2}$

Hence, a Fibonacci series can look like this -

 $F_8 = 0 \ 1 \ 1 \ 2 \ 3 \ 5 \ 8 \ 13$

or, this -

 $F_8 = 1 \ 1 \ 2 \ 3 \ 5 \ 8 \ 13 \ 21$

For illustration purpose, Fibonacci of F₈ is displayed as –

Fibonacci Iterative Algorithm

First we try to draft the iterative algorithm for Fibonacci series.

ALGEMEN

Procedure Fibonacci(n)

declare f₀, f₁, fib, loop

set f_0 to 0

set f_1 to 1

display f₀, f₁

for loop $\leftarrow 1$ to n

 $fib \leftarrow f_0 + f_1$ $f_0 \leftarrow f_1$

 $f_1 \leftarrow fib$

display fib

end for

end procedure

To know about the implementation of the above algorithm in C programming language, <u>click here</u>.

Fibonacci Recursive Algorithm

Let us learn how to create a recursive algorithm Fibonacci series. The base criteria of recursion.

START

Procedure Fibonacci(n)

declare f₀, f₁, fib, loop

set f_0 to 0

set f_1 to 1

display f₀, f₁

for loop $\leftarrow 1$ to n

 $fib \leftarrow f_0 + f_1$

 $f_0 \leftarrow f_1$

 $f_1 \leftarrow fib$

display fib

end for

END

INTRODUCTION TO COMPUTER & IT (107)

UNIT 1: Basics of Computer and its Evolution

Introduction and Evolution of Computer

A computer is a programmable machine designed to sequentially and automatically carry out a sequence of arithmetic or logical operations. The particular sequence of operations can be changed readily, allowing the computer to solve more than one kind of problem. An important class of computer operations on some computing platforms is the accepting of input from human operators and the output of results formatted for human consumption. The interface between the computer and the human operator is known as the user interface. Conventionally a computer consists of some form of memory, at least one element that carries out arithmetic and logic operations, and a sequencing and control unit that can change the order of operations based on the information that is stored. Peripheral devices allow information to be entered from an external source, and allow the results of operations to be sent out. A computer's processing unit executes series of instructions that make it read, manipulate and then store data. Conditional instructions change the sequence of instructions as a function of the current state of the machine or its environment. The first electronic digital computers were developed in the mid-20th century (1940–1945). Originally, they were the size of a large room, consuming as much power as several hundred modern personal computers (PCs). In this era mechanical analog computers were used for military applications.

Modern computers based on integrated circuits are millions to billions of times more capable than the early machines, and occupy a fraction of the space. Simple computers are small enough to fit into mobile devices, and mobile computers can be powered by small batteries. Personal computers in their various forms are icons of the Information Age and are what most people think of as "computers". However, the embedded computers found in many devices from mp3 players to fighter aircraft and from toys to industrial robots are the most numerous.

General Functions of Computer

Computer is an advanced electronic device that takes raw data as input from the user and processes these data under the control of set of instructions (called program) and gives the result (output) and saves output for the future use. It can process both numerical and non-numerical (arithmetic and logical) calculations.

A computer has four functions:		
a. accepts data	Input	
b. processes data	Processing	
c. produces output	Output	
d. stores results	Storage	

Input (Data):

Input is the raw information entered into a computer from the input devices. It is the collection of letters, numbers, images etc.

Process:

Process is the operation of data as per given instruction. It is totally internal process of the computer system.

Output:

Output is the processed data given by computer after data processing. Output is also called as Result. We can save these results in the storage devices for the future use.

Storage:

Computer data storage, often called storage or memory, refers to computer components and recording media that retain digital data. Data storage is a core function and fundamental component of computers.

Computer System

All of the components of a computer system can be summarized with the simple equations.

COMPUTER SYSTEM = HARDWARE + SOFTWARE + USER

• <u>Hardware</u> = Internal Devices + Peripheral Devices

The hardware are the parts of the computer itself including the Central Processing Unit (CPU) and related microchips and micro-circuitry, keyboards, monitors, case and drives (hard, CD, DVD, floppy, optical, tape, etc...). Other extra parts called peripheral components or devices include mouse, printers, modems, scanners, digital

cameras and cards (sound, colour, video) etc... Together they are often referred to as a personal computer.

All physical parts of the computer (or everything that we can touch) are known as Hardware.

• <u>Software</u> = Programs

The software is the information that the computer uses to get the job done. Software needs to be accessed before it can be used. There are many terms used for the process of accessing software including running, executing, starting up, opening, and others.

Computer programs allow users to complete tasks. A program can also be referred to as an application and the two words are used interchangeably.

Software gives "intelligence" to the computer.

• **<u>USER</u>** = Person, who operates computer.

A user is an agent, either a human agent (end-user) or software agent, who uses a computer or network service. Users are also widely characterized as the class of people that use a system without complete technical expertise required to understand the system fully. Such users are also divided into users and power users. Both are terms of degradation but the latter connotes a "know-it-all" attitude. In projects in which the actor of the system is another system or a software agent, it is quite possible that there is no end-user for the system. In this case, the end-users for the system would be indirect end-users.

CHARACTERISTICS OF COMPUTER

Speed, accuracy, diligence, storage capability and versatility are some of the key characteristics of a computer. A brief overview of these characteristics is—

Speed The computer can process data very fast, at the rate of millions of instructions per second.

Some calculations that would have taken hours and days to complete otherwise, can be completed in a few seconds using the computer. For example, calculation and generation of salary slips of thousands of employees of an organization, weather forecasting that requires analysis of a large amount of data related to temperature, pressure and humidity of various places, etc.

<u>Accuracy</u> Computer provides a high degree of accuracy. For example, the computer can accurately give the result of division of any two numbers up to 10 decimal places.

Diligence When used for a longer period of time, the computer does not get tired or fatigued. It can perform long and complex calculations with the same speed and accuracy from the start till the end.

Storage Capability Large volumes of data and information can be stored in the computer and also retrieved whenever required. A limited amount of data can be stored, temporarily, in the primary memory. Secondary storage devices like fl oppy disk and compact disk can store a large amount of data permanently.

<u>Versatility</u> Computer is versatile in nature. It can perform different types of tasks with the same ease. At one moment you can use the computer to prepare a letter document and in the next moment you may play music or print a document. Computers have several limitations too. Computer can only perform tasks that it has been programmed to do. Computer cannot do any work without instructions from the user. It executes instructions as specified by the user and does not take its own decisions.

Application of Computers

Computers have proliferated into various areas of our lives. For a user, computer is a tool that provides the desired information, whenever needed. You may use computer to get information about the reservation of tickets (railways, airplanes and cinema halls), books in a library, medical history of a person, a place in a map, or the dictionary meaning of a word. The information may be presented to you in the form of text, images, video clips, etc. Some of the application areas of the computer are listed below:

Education Computers are extensively used, as a tool and as an aid, for imparting education. Educators use computers to prepare notes and presentations of their lectures. Computers are used to develop computer-based training packages, to provide distance education using the elearning software, and to conduct online examinations. Researchers use computers to get easy access to conference and journal details and to get global access to the research material.

Entertainment Computers have had a major impact on the entertainment industry. The user can download and view movies, play games, chat, book tickets for cinema halls, use multimedia for making movies, incorporate visual and sound effects using computers, etc. The users can also listen to music, download and share music, create music using computers, etc.

Sports A computer can be used to watch a game, view the scores, improve the game, play games (like chess, etc.) and create games. They are also used for the purposes of training players.

<u>Advertising</u> Computer is a powerful advertising media. Advertisement can be displayed on different websites, electronic-mails can be sent and reviews of a product by different customers can be posted. Computers are also used to create an advertisement using the visual

and the sound effects. For the advertisers, computer is a medium via which the advertisements can be viewed globally. Web advertising has become a significant factor in the marketing plans of almost all companies. In fact, the business model of Google is mainly dependent on web advertising for generating revenues.

<u>Medicine</u> Medical researchers and practitioners use computers to access information about the advances in medical research or to take opinion of doctors globally. The medical history of patients is stored in the computers. Computers are also an integral part of various kinds of sophisticated medical equipments like ultrasound machine, CAT scan machine, MRI scan machine, etc. Computers also provide assistance to the medical surgeons during critical surgery operations like laparoscopic operations, etc.

<u>Science and Engineering</u> Scientists and engineers use computers for performing complex scientific calculations, for designing and making drawings (CAD/CAM applications) and also for simulating and testing the designs. Computers are used for storing the complex data, performing complex calculations and for visualizing 3-dimensional objects. Complex scientific capplications like the launch of the rockets, space exploration, etc., are not possible without the computers.

Government The government uses computers to manage its own operations and also for egovernance. The websites of the different government departments provide information to the users. Computers are used for the fi ling of income tax return, paying taxes, online submission of water and electricity bills, for the access of land record details, etc.

Home Computers have now become an integral part of home equipment. At home, people use

computers to play games, to maintain the home accounts, for communicating with friends and relatives via Internet, for paying bills, for education and learning, etc. Microprocessors are embedded in house hold utilities like, washing machines, TVs, food processors, home theatres, security devices, etc. The list of applications of computers is so long that it is not possible to discuss all of them here. In addition to the applications of the computers discussed above, computers have also proliferated into areas like banks, investments, stock trading, accounting, ticket reservation, military operations, meteorological predictions, social networking, business organizations, police department, video conferencing, telepresence, book publishing, web newspapers, and information sharing.

History of computer science

The history of computer science began long before the modern discipline of computer science that emerged in the twentieth century, and hinted at in the centuries prior. The

progression, from mechanical inventions and mathematical theories towards the modern concepts and machines, formed a major academic field and the basis of a massive worldwide industry.

Mechanical computers:

A mechanical computer is built from mechanical components such as levers and gears, rather than electronic components. The most common examples are adding machines and mechanical counters, which use the turning of gears to increment output displays. More complex examples can carry out multiplication and division, and even differential analysis.

Abacus

The abacus, also called a counting frame, is a calculating tool used primarily in parts of Asia for performing arithmetic processes. Today, abaci are often constructed as a bamboo frame with beads sliding on wires, but originally they were beans or stones moved in grooves in sand or on tablets of wood, stone, or metal. The abacus was in use centuries before the adoption of the written modern numeral system and is still widely used by merchants, traders and clerks in Asia, Africa, and elsewhere. The user of an abacus is called an abacist.

Napier's Bones

Napier's bones is an abacus created by John Napier for calculation of products and quotients of numbers that was based on Arab mathematics and lattice multiplication. The abacus consists of a board with a rim; the user places Napier's rods in the rim to conduct multiplication or division. The board's left edge is divided into 9 squares, holding the numbers 1 to 9. The Napier's rods consist of strips of wood, metal or heavy cardboard. Napier's bones are three dimensional, square in cross section, with four different rods engraved on each one. A set of such bones might be enclosed in a convenient carrying case. A rod's surface comprises 9 squares, and each square, except for the top one, comprises two halves divided by a diagonal line. The first square of each rod holds a single digit, and the other squares hold this number's double, triple, quadruple, quintuple, and so on until the last square contains nine times the number in the top square. The digits of each product are written one to each side of the diagonal; numbers less than 10 occupy the lower triangle, with a zero in the top half. A set consists of 10 rods corresponding to digits 0 to 9. The rod 0, although it may look unnecessary, is obviously still needed for multipliers or multiplicands having 0 in them.

Slide rule

The slide rule is a mechanical computer. The slide rule is used primarily for multiplication and division, and also for functions such as roots, logarithms and trigonometry, but is not normally used for addition or subtraction. Slide rules come in a diverse range of styles and generally appear in a linear or circular form with a standardized set of markings (scales) essential to performing mathematical computations. Slide rules manufactured for specialized fields such as aviation or finance typically feature additional scales that aid in calculations common to that field. William Oughtred and others developed the slide rule in the 17th century based on the emerging work on logarithms by John Napier. Before the advent of the pocket calculator, it was the most commonly used calculation tool in science and engineering. The use of slide rules continued to grow through the 1950s and 1960s even as digital computing devices were being gradually introduced; but around 1974 the electronic scientific calculator made it largely obsolete and most suppliers left the business.

Pascal's calculator

Blaise Pascal invented the mechanical calculator in 1642. He conceived the idea while trying to help his father who had been assigned the task of reorganizing the tax revenues of the French province of Haute-Normandie; first called Arithmetic Machine, Pascal's Calculator and later Pascaline, it could add and subtract directly and multiply and divide by repetition. Pascal went through 50 prototypes before presenting his first machine to the public in 1645. He dedicated it to Pierre Séguier, the chancellor of France at the time. He built around twenty more machines during the next decade, often improving on his original design. Nine machines have survived the centuries, most of them being on display in European museums. In 1649 a royal privilege, signed by Louis XIV of France, gave him the exclusivity of the

design and manufacturing of calculating machines in France.

Stepped Reckoner

The Step Reckoner (or Stepped Reckoner) was a digital mechanical calculator invented by German mathematician Gottfried Wilhelm Leibniz around 1672 and completed in 1694. The name comes from the translation of the German term for its operating mechanism; staffelwalze meaning 'stepped drum'. It was the first calculator that could perform all four arithmetic operations: addition, subtraction, multiplication and division. Its intricate precision gearwork, however, was somewhat beyond the fabrication technology of the time; mechanical problems, in addition to a design flaw in the carry mechanism, prevented the machines from working reliably.

Jacquard loom

The Jacquard loom is a mechanical loom, invented by Joseph Marie Jacquard in 1801, that simplifies the process of manufacturing textiles with complex patterns such as brocade, damask and matelasse. The loom is controlled by punched cards with punched holes, each row of which corresponds to one row of the design. Multiple rows of holes are punched on each card and the many cards that compose the design of the textile are strung together in order. It is based on earlier inventions by the Frenchmen Basile Bouchon (1725), Jean Baptiste Falcon (1728) and Jacques Vaucanson (1740).

Charles Babbage's Difference engine

A difference engine is an automatic, mechanical calculator designed to tabulate polynomial functions. The name derives from the method of divided differences, a way to interpolate or tabulate functions by using a small set of polynomial coefficients. Both logarithmic and trigonometric functions, functions commonly used by both navigators and scientists, can be approximated by polynomials, so a difference engine can compute many useful sets of numbers. The historical difficulty in producing error free tables by teams of mathematicians and human "computers" spurred Charles Babbage's desire to build a mechanism to automate the process.

Analytical Engine

The Analytical Engine was a proposed mechanical general-purpose computer designed by English mathematician Charles Babbage. It was first described in 1837 as the successor to Babbage's difference engine, a design for a mechanical calculator. The Analytical Engine incorporated an arithmetical unit, control flow in the form of conditional branching and loops, and integrated memory, making it the first Turing-complete design for a general-purpose computer.

Charles Babbage (1791-1871) the Father of Computers

Charles Babbage is recognized today as the Father of Computers because his impressive designs for the Difference Engine and Analytical Engine foreshadowed the invention of the modern electronic digital computer. He led a fascinating life, as did all the folks involved in the history of computers. He also invented the cowcatcher, dynamometer, standard railroad gauge, uniform postal rates, occulting lights for lighthouses, Greenwich time signals, heliograph opthalmoscope.

Lady Augusta Ada Countess of Lovelace (First Computer Programmer)

Babbage owes a great debt to Lady Augusta Ada, Countess of Lovelace. Daughter of the famous romantic poet, Lord Byron, she was a brilliant mathematician who helped Babbage in his work. Above all, she documented his work, which Babbage never could bother to do. As a result we know about Babbage at all. Lady Augusta Ada also wrote programs to be run on Babbage's machines. For this, she is recognized as the first computer programmer.

Electro-Mechanical Computer:

Census Tabulating Machine

Herman Hollerith Develop The tabulating machine. The tabulating machine was an electrical device designed to assist in summarizing information and, later, accounting. Invented by Herman Hollerith, the machine was developed to help process data for the 1890 U.S. Census. It spawned a larger class of devices known as unit record equipment and the data processing industry. Herman Hollerith worked as a statistician for the U.S. Census Bureau in the 1880s and 1890s. The U.S. Constitution requires a census count every ten years so that the membership of the House of Representatives will be proportional to the population of each state. This is always a moving target, hence the ten year review of the current state of demographic affairs. The 1880 census took seven years to process. The end of the 19th/beginning of the 20th centuries was the period of highest rate of immigration to the United States. Hollerith deduced, and it didn't take a rocket scientist to conclude, that the next census would take longer than ten years, the results not available before the whole census counting thing had to start again.

So, as the saying goes, "necessity became the mother of invention" and Hollerith designed and built the Census Counting Machine illustrated here and in the next slide. Punched cards (a la Jacquard looms) were used to collect the census data (the origin of the IBM punched cards) and the cards were fed into a sorting machine before being read by the census counting machine which recorded and tabulated the results. Each card was laid on an open grid. A matrix of wires was lowered onto the card and wherever there was a hole in the card, a wire fell through, making an electrical connection which triggered a count on the appropriate dial(s) in the face of the machine. The 1890 census took just three months to process even though quite a bit more data was collected than ever before.

Hollerith was the first American associated with the history of computers. As you might expect, he was also the first to make a bunch of money at it. His company, the Tabulating Machine Company, became the Computer Tabulating Recording Company in 1913 after struggling in the market and merging with another company that produced a similar product. The company hired a gentleman named Thomas J. Watson in 1918 who was primarily instrumental in turning the company around. In 1924, the company was renamed International Business machines (IBM) Corporation.

Harvard Mark I

The IBM Automatic Sequence Controlled Calculator (ASCC), called the Mark I by Harvard University, was an electro-mechanical computer.

The electromechanical ASCC was devised by Howard H. Aiken, built at IBM and shipped to Harvard in February 1944. It began computations for the U.S. Navy Bureau of Ships in May and was officially presented to the university on August 7, 1944. The ASCC was built from switches, relays, rotating shafts, and clutches. It used 765,000 components and hundreds of miles of wire, comprising a volume of 51 feet (16 m) in length, eight feet (2.4 m) in height, and two feet (~61 cm) deep. It had a weight of about 10,000 pounds (4500 kg). The basic calculating units had to be synchronized mechanically, so they were run by a 50-foot (~15.5 m) shaft driven by a five-horsepower (4 kW) electric motor. From the IBM Archives: The Automatic Sequence Controlled Calculator (Harvard Mark I) was the first operating machine that could execute long computations automatically. A project conceived by Harvard University's Dr. Howard Aiken, the Mark I was built by IBM engineers in Endicott, N.Y.

The first computer bug

The lady is U.S. Rear Admiral Dr. Grace Murray Hopper, who worked with Howard Aiken from 1944 and used his machine for gunnery and ballistics calculation for the US Bureau of Ordnance's Computation project. One day, the program she was running gave incorrect results and, upon examination, a moth was found blocking one of the relays. The bug was removed and the program performed to perfection. Since then, a program error in a computer has been called a bug.

Electronic digital computers

The Turing Machine

The "Turing" machine was described by Alan Turing in 1936, who called it an "automaticmachine". The Turing machine is not intended as a practical computing technology, but rather as a hypothetical device representing a computing machine. Turing machines help computer scientists understand the limits of mechanical computation. A Turing machine is a device that manipulates symbols on a strip of tape according to a table of rules. Despite its simplicity, a Turing machine can be adapted to simulate the logic of any computer algorithm, and is particularly useful in explaining the functions of a CPU inside a computer.

Atanasoff–Berry Computer

The ABC was built by Dr. Atanasoff and graduate student Clifford Berry in the basement of the physics building at Iowa State College during 1939–42. The Atanasoff–Berry Computer (ABC) was the first electronic digital computing device. Conceived in 1937, the machine was not programmable, being designed only to solve systems of linear equations. It was

successfully tested in 1942. However, its intermediate result storage mechanism, a paper card writer/reader, was unreliable, and when inventor John Vincent Atanasoff left Iowa State College for World War II assignments, work on the machine was discontinued. The ABC pioneered important elements of modern computing, including binary arithmetic and electronic switching elements, but its special-purpose nature and lack of a changeable, stored program distinguish it from modern computers.

Colossus computer

Colossus was the world's first electronic, digital, programmable computer. Colossus and its successors were used by British codebreakers to help read encrypted German messages during World War II. They used thermionic valves (vacuum tubes) to perform the calculations.

Colossus was designed by engineer Tommy Flowers with input from Sidney Broadhurst, William Chandler, Allen Coombs and Harry Fensom. at the Post Office Research Station, Dollis Hill to solve a problem posed by mathematician Max Newman at Bletchley Park. The prototype, Colossus Mark 1, was shown to be working in December 1943 and was operational at Bletchley Park by February 1944. An improved Colossus Mark 2 first worked on 1 June 1944, just in time for the Normandy Landings. Ten Colossus computers were in use by the end of the war.

The Colossus computers were used to help decipher teleprinter messages which had been encrypted using the Lorenz SZ40/42 machine—British codebreakers referred to encrypted German teleprinter traffic as "Fish" and called the SZ40/42 machine and its traffic "Tunny". Colossus compared two data streams, counting each match based on a programmable Boolean function. The encrypted message was read at high speed from a paper tape. The other stream was generated internally, and was an electronic simulation of the Lorenz machine at various trial settings. If the match count for a setting was above a certain threshold, it would be sent as output to an electric typewriter.

ENIAC

ENIAC (Electronic Numerical Integrator And Computer) was conceived and designed by John Mauchly and J. Presper Eckert of the University of Pennsylvania. The team of design engineers assisting the development included Robert F. Shaw (function tables), Jeffrey Chuan Chu (divider/square-rooter), Thomas Kite Sharpless (master programmer), Arthur Burks (multiplier), Harry Huskey (reader/printer) and Jack Davis (accumulators).

ENIAC was the first general-purpose electronic computer. It was a Turing-complete digital computer capable of being reprogrammed to solve a full range of computing problems.

ENIAC was designed to calculate artillery firing tables for the United States Army's Ballistic Research Laboratory. When ENIAC was announced in 1946 it was heralded in the press as a "Giant Brain". It boasted speeds one thousand times faster than electro-mechanical machines, a leap in computing power that no single machine has since matched. This mathematical power, coupled with general-purpose programmability, excited scientists and industrialists. The inventors promoted the spread of these new ideas by conducting a series of lectures on computer architecture.

Generations of Computers

The history of computer development is often referred to in reference to the different generations of computing devices. A generation refers to the state of improvement in the product development process. This term is also used in the different advancements of new computer technology. With each new generation, the circuitry has gotten smaller and more advanced than the previous generation before it. As a result of the miniaturization, speed, power, and computer memory has proportionally increased. New discoveries are constantly being developed that affect the way we live, work and play. Each generation of computers is characterized by major technological development that fundamentally changed the way computers operate, resulting in increasingly smaller, cheaper, more powerful and more efficient and reliable devices. Read about each generation and the developments that led to the current devices that we use today.

First Generation - 1940-1956: Vacuum Tubes

The first computers used vacuum tubes for circuitry and magnetic drums for memory, and were often enormous, taking up entire rooms. They were very expensive to operate and in addition to using a great deal of electricity, generated a lot of heat, which was often the cause of malfunctions. First generation computers relied on machine language, the lowest-level programming language understood by computers, to perform operations, and they could only solve one problem at a time. Input was based on punched cards and paper tape, and output was displayed on printouts. The UNIVAC and ENIAC computers are examples of first-generation computing devices. The UNIVAC was the first commercial computer delivered to a business client, the U.S. Census Bureau in 1951.

Second Generation (1956-1963) Transistors

Transistors replaced vacuum tubes and ushered in the second generation of computers. The transistor was invented in 1947 but did not see widespread use in computers until the late 1950s. The transistor was far superior to the vacuum tube, allowing computers to become smaller, faster, cheaper, more energy-efficient and more reliable than their first-generation
predecessors. Though the transistor still generated a great deal of heat that subjected the computer to damage, it was a vast improvement over the vacuum tube. Second-generation computers still relied on punched cards for input and printouts for output. Second-generation computers moved from cryptic binary machine language to symbolic, or assembly, languages, which allowed programmers to specify instructions in words. High-level programming languages were also being developed at this time, such as early versions of COBOL and FORTRAN. These were also the first computers that stored their instructions in their memory, which moved from a magnetic drum to magnetic core technology.

Third Generation (1964-1971) Integrated Circuits

The development of the integrated circuit was the hallmark of the third generation of computers. Transistors were miniaturized and placed on silicon chips, called semiconductors, which drastically increased the speed and efficiency of computers. Instead of punched cards and printouts, users interacted with third generation computers through keyboards and monitors and interfaced with an operating system, which allowed the device to run many different applications at one time with a central program that monitored the memory. Computers for the first time became accessible to a mass audience because they were smaller and cheaper than their predecessors.

Fourth Generation (1971-Present) Microprocessors

The microprocessor brought the fourth generation of computers, as thousands of integrated circuits were built onto a single silicon chip. What in the first generation filled an entire room could now fit in the palm of the hand. The Intel 4004 chip, developed in 1971, located all the components of the computer—from the central processing unit and memory to input/output controls—on a single chip.

In 1981 IBM introduced its first computer for the home user, and in 1984 Apple introduced the Macintosh. Microprocessors also moved out of the realm of desktop computers and into many areas of life as more and more everyday products began to use microprocessors.

As these small computers became more powerful, they could be linked together to form networks, which eventually led to the development of the Internet. Fourth generation computers also saw the development of GUIs, the mouse and handheld devices.

Fifth Generation (Present and Beyond) Artificial Intelligence

Fifth generation computing devices, based on artificial intelligence, are still in development, though there are some applications, such as voice recognition, that are being used today. The use of parallel processing and superconductors is helping to make artificial intelligence a reality. Quantum computation and molecular and nanotechnology will radically change the

face of computers in years to come. The goal of fifth-generation computing is to develop devices that respond to natural language input and are capable of learning and selforganization.

Computer speed and Measurement Unit

Space measurement units

The size of device in computers does not reflect the space available to store data in it. There are larger devices that can store only a few data where as many tiny devices that stores unbelievable amount of data. Because how long, how thick etc cannot determine how much we can store inside, we need to find some other way to measure space.

Almost all of the computers use binary numbering systems (though there are some exceptions). Binary numbering system consists of only two digits -0 and 1 to represent any quantity. 10 in binary is equal to the 2 and 100 to 5. We will be learning this numbering system after some days.

Everything in computers is represented in strings of binary numbers. For example capital A is interpreted by computer as 0100 0001 and B is 0100 0010. All characters, numbers, symbols, images, sounds, animations, videos and everything, yes everything is converted into suitable binary code to store in computer or process by computer. So if there is any device that can store one binary digit (whether 0 or 1), its storage capacity is 1 bit. We've already learned that bit is the abbreviation of binary digit. Any device that has storage space to accommodate 5 binary digits has 5 bits space. You require thousands and millions of bits for a file and expressing the space available in bits only is really inconvenient because it will be an extremely large number. So, we have larger units that represent a group of lower units. A group of 4 binary digits is called a nibble (4 bits = 1 Nibble). Similarly a group of 8 bits is called a byte (1 byte = 8 bits). As you have seen the example above, each character requires 8 bits which is 1 byte. So 1 character requires 1 byte space. Now, if you have a text file whose size is 32 bytes, it means there are 32 x 8 binary digits (0s and 1s) stored in it. In the metric system 1000 lifts up the unit to the higher such as 1000 meter is 1 kilometer, 1000 liter is 1 kiloliter etc. In binary numbering system it is 1024 (=2 10) that converts to higher unit. Following table lists the different units and their values:

Space Measurement Units

Units	Equivalent
0 or 1	1 Bit
4 bits	1 Nibble

8 bits	1 Byte
1024 bytes	1 Kilobytes (KB)
1024 Kilobytes	1 Megabyte (MB)
1024 Megabytes	1 Gigabytes (GB)
1024 Gigabytes	1 Terabytes (TB)
1024 Terabytes	1 Petabyte (PB)
1024 Petabytes	1 Exabyte (EB)

Speed Measurement Units

Speed is related to time. Computer can perform millions of tasks in one second. So to compare the speed of computer operation (execution of programs and instructions) we require some units that can represent a very small fraction of time. Following are the units used to indicate fraction of seconds:

Speed Measurement Units

Units	Equivalent
1000 th of a second	1 Milliseconds (MS)
1000 th of a milliseconds	1 Microseconds (µs)
1000 th of a microseconds	1 Nanosecond (ns)
1000 th of a nanoseconds	1 Picoseconds (ps)
1000 th of a picoseconds	1 Femtoseconds (fs)

CLASSIFICATION OF COMPUTER

On the basis of size

The digital computers that are available nowadays vary in their sizes and types. The computers are broadly classified into four categories (Figure 1.8) based on their size and type

- (1) Microcomputers,
- (2) Mini-computers,
- (3) Mainframe computers, and
- (4) Supercomputer.

Figure: Classification of computers based on size and type

Microcomputers

Microcomputers are small, low-cost and single-user digital computer. They consist of CPU, input unit, output unit, storage unit and the software. Although microcomputers are standalone machines, they can be connected together to create a network of computers that can serve more than one user. IBM PC based on Pentium microprocessor and Apple Macintosh are some examples of microcomputers. Microcomputers include desktop computers, notebook computers or laptop, tablet computer, handheld computer, smart phones and notebook.

Desktop Computer or Personal Computer (PC) is the most common type of microcomputer. It is a stand-alone machine that can be placed on the desk. Externally, it consists of three units—key-board, monitor, and a system unit containing the CPU, memory, hard disk drive, etc. It is not very expensive and is suited to the needs of a single user at home, small business units, and organizations. Apple, Microsoft, HP, Dell and Lenovo are some of the PC manufacturers.

Notebook Computers or Laptop resemble a notebook. They are portable and have all the features of a desktop computer. The advantage of the laptop is that it is small in size (can be put inside a briefcase), can be carried anywhere, has a battery backup and has all the functionality of the desk-top. Laptops can be placed on the lap while working (hence the name). Laptops are costlier than the desktop machines.

<u>Netbook</u> These are smaller notebooks optimized for low weight and low cost, and are designed for accessing web-based applications. Starting with the earliest netbook in late 2007, they have gained significant popularity now. Netbooks deliver the performance needed to enjoy popular activities like streaming videos or music, emailing, Web surfing or instant messaging. The word netbook was created as a blend of Internet and notebook.

<u>**Tablet Computer**</u> has features of the notebook computer but it can accept input from a stylus or a pen instead of the keyboard or mouse. It is a portable computer. Tablet computer are the new kind of PCs.

<u>Handheld Computer or Personal Digital Assistant (PDA)</u> is a small computer that can be held on the top of the palm. It is small in size. Instead of the keyboard, PDA uses a pen or a stylus for input. PDAs do not have a disk drive. They have a limited memory and are less powerful. PDAs can be connected to the Internet via a wireless connection. Casio and Apple are some of the manufacturers of PDA. Over the last few years, PDAs have merged into mobile phones to create smart phones.

<u>Smart Phones</u> are cellular phones that function both as a phone and as a small PC. They may use a stylus or a pen, or may have a small keyboard. They can be connected to the Internet

wirelessly. They are used to access the electronic-mail, download music, play games, etc. Blackberry, Apple, HTC, Nokia and LG are some of the manufacturers of smart phones.

Minicomputers

Minicomputers are digital computers, generally used in multi-user systems. They have high processing speed and high storage capacity than the microcomputers. Minicomputers can support 4–200 users simultaneously. The users can access the minicomputer through their PCs or terminal. They are used for real-time applications in industries, research centers, etc. PDP 11, IBM (8000 series) are some of the widely used minicomputers.

Mainframe Computers

Mainframe computers are multi-user, multi-programming and high performance computers. They operate at a very high speed, have very large storage capacity and can handle the workload of many users. Mainframe computers are large and powerful systems generally used in centralized databases. The user accesses the mainframe computer via a terminal that may be a dumb terminal, an intelligent terminal or a PC. A dumb terminal cannot store data or do processing of its own. It has the input and output device only. An intelligent terminal has the input and output device, can do processing, but, cannot store data of its own. The dumb and the intelligent terminal use the processing power and the storage facility of the mainframe computer. Mainframe computers are used in organizations like banks or companies, where many people require frequent access to the same data. Some examples of mainframes are CDC 6600 and IBM ES000 series.

Supercomputers

Supercomputers are the fastest and the most expensive machines. They have high processing speed compared to other computers. The speed of a supercomputer is generally measured in FLOPS (Floating point Operations Per Second). Some of the faster supercomputers can perform trillions of calculations per second. Supercomputers are built by interconnecting thousands of processors that can work in parallel. Supercomputers are used for highly calculation-intensive tasks, such as, weather forecasting, climate research (global warming), molecular research, biological research, nuclear research and aircraft design. They are also used in major universities, military agencies and scientific research laboratories. Some examples of supercomputers are IBM Roadrunner, IBM Blue gene and Intel ASCI red. PARAM is a series of supercomputer assembled in India by C-DAC (Center for

Development of Advanced Computing), in Pune. PARAM Padma is the latest machine in this series. The peak computing power of PARAM Padma is 1 Tera FLOP (TFLOP).

On the basis of working principle

Digital computers:

Digital computers operates on inputs which are on-off type (being digit 1 and 0) and its outputs is also in form of on-off signals. Digital computers are based on counting operation. Any data to be manipulated by a digital computer must first be converted to a discrete(1, 0) representation. There is a practical limit to the accuracy of the readings of analog devices, usually to the nearest tenth of the unit of measure. Thus if water in a beaker was being heated and its temperature rose from 50 C to 51 C ,someone observing the thermometer might only be able to distinguish the temperature 50.0, 50.1, 50.2... 50.9, 51.0.

Analog computers:

The analog computer operate by measuring rather than counting. It measure continually, usually of a physical nature data such as lengths, voltages, or currents. It does not produce number but produces its results in the form of graph. It is more efficient in continues calculations. Analog machines are usually special purpose devices, dedicated to a single task.

Hybrid computers:

A hybrid computer is combination of both analog and digital computer i.e. a part of processing is done on analog computer and a part on digital computer. A hybrid computer combines the best characteristics of both analog and digital computer. It can accept input data in both analog and digital form. It is used for simulation application.

Computer name	Input	Output	Based on	Examples
Digital	On/Off, 1/0	On/Off, 1/0	Counting	General purpose PCs
Analog	Measure	Graphs pictures	Continuous	Weather
	elements		Measurements	forecasting speedometer
				FCC 1
				ECG machines, etc.
Hybrid	Both 0/1&	Both On/Off &	Counting &	Controlling &
	Measure	Graphs	Measurements	monitoring plants.
	elements			Petrol pumps,
				Modem, Simulation
				etc.

On the basis of brand

On the basis of brand , the computer can be classified as IBM PC, IBM compatibles and Apple/Macintosh computer.

IBM PC:

IBM PC is the largest computer manufacturing company establishing USA. The computer manufactures by IBM PC or branded computer. Personal Computer (PC) is the most important type of micro computer system. The micro computer manufacture by IBM company are called IBM PC. These computers are reliable, durable and have better quality but they are costly.

IBM compatibles:

The computers that have some functional characteristics and principles of IBM computer are called IBM compatibles. In other word, all the computer are manufactured by the another companies rather than IBM company are Known as IBM compatibles. All the software and hardware of IBM compatibles. These are cheaper and Their Parts are easily available in Market. they are also duplicate or assemble computer.

Apple/Macintosh Computer

All the computers manufacture by apple cooperation, a leading computers manufacturing computer of USA are known as apple/macintosh computers. These computer use their own software and hardware. They are totally different than that of IBM computers, In terms of both hardware and software. For e.g software developed for apple computer can't run or IBM computers and vice-versa. Similarly, floppy disk formatting in IBM computer can't be recognized by apple macintosh computer and vice-versa. It is popularly used in desktop publishing (DTP) houses as they provide better quality of graphic output.

Mobile Computing

Mobile computing' is a form of human–computer interaction by which a computer is expected to be transported during normal usage. Mobile computing has three aspects: mobile communication, mobile hardware, and mobile software. The first aspect addresses communication issues in ad-hoc and infrastructure networks as well as communication properties, protocols, data formats and concrete technologies. The second aspect is on the hardware, e.g., mobile devices or device components. The third aspect deals with the characteristics and requirements of mobile applications.

DEFINITIONS

Mobile computing is "taking a computer and all necessary files and software out into the field."

"Mobile computing: being able to use a computing device even when being mobile and therefore changing location. Portability is one aspect of mobile computing." "Mobile computing is the ability to use computing capability without a pre-defined location and/or connection to a network to publish and/or subscribe to information."

Mobile Computing is a variety of wireless devices that has the mobility to allow people to connect to the internet, providing wireless transmission to access data and information from where ever location they may be.

T DEPARTMENT PART

Computer - Overview

Today's world is an information-rich world and it has become a necessity for everyone to know about computers. A computer is an electronic data processing device, which accepts and stores data input, processes the data input, and generates the output in a required format. The purpose of this tutorial is to introduce you to Computers and its fundamentals.

Functionalities of a Computer

If we look at it in a very broad sense, any digital computer carries out the following five functions –

Step 1 – Takes data as input.

Step 2 – Stores the data/instructions in its memory and uses them as required.

- Step 3 Processes the data and converts it into useful information.
- **Step 4** Generates the output.
- **Step 5** Controls all the above four steps.



Advantages of Computers

Following are certain advantages of computers.

High Speed

- Computer is a very fast device.
- It is capable of performing calculation of very large amount of data.

- The computer has units of speed in microsecond, nanosecond, and even the picosecond.
- It can perform millions of calculations in a few seconds as compared to man who will spend many months to perform the same task.

Accuracy

- In addition to being very fast, computers are very accurate.
- The calculations are 100% error free.
- Computers perform all jobs with 100% accuracy provided that the input is correct.

Storage Capability

- Memory is a very important characteristic of computers.
- A computer has much more storage capacity than human beings.
- It can store large amount of data.
- It can store any type of data such as images, videos, text, audio, etc.

Diligence

- Unlike human beings, a computer is free from monotony, tiredness, and lack of concentration.
- It can work continuously without any error and boredom.
- It can perform repeated tasks with the same speed and accuracy.

Versatility

- A computer is a very versatile machine.
- A computer is very flexible in performing the jobs to be done.
- This machine can be used to solve the problems related to various fields.
- At one instance, it may be solving a complex scientific problem and the very next moment it may be playing a card game.

Reliability

- A computer is a reliable machine.
- Modern electronic components have long lives.
- Computers are designed to make maintenance easy.

Automation

- Computer is an automatic machine.
- Automation is the ability to perform a given task automatically. Once the computer receives a program i.e., the program is stored in the computer memory, then the

program and instruction can control the program execution without human interaction.

Reduction in Paper Work and Cost

- The use of computers for data processing in an organization leads to reduction in paper work and results in speeding up the process.
- As data in electronic files can be retrieved as and when required, the problem of maintenance of large number of paper files gets reduced.
- Though the initial investment for installing a computer is high, it substantially reduces the cost of each of its transaction.

Disadvantages of Computers

Following are certain disadvantages of computers.

No I.Q.

- A computer is a machine that has no intelligence to perform any task.
- Each instruction has to be given to the computer.
- A computer cannot take any decision on its own.

Dependency

• It functions as per the user's instruction, thus it is fully dependent on humans.

Environment

• The operating environment of the computer should be dust free and suitable.

No Feeling

- Computers have no feelings or emotions.
- It cannot make judgment based on feeling, taste, experience, and knowledge unlike humans.

ald ISO 9001:2015 & 14001:

Computer - Applications

In this chapter, we will discuss the application of computers in various fields. Business



A computer has high speed of calculation, diligence, accuracy, reliability, or versatility which has made it an integrated part in all business organizations.

Computer is used in business organizations for -

- Payroll calculations
- Budgeting
- Sales analysis
- Financial forecasting
- Managing employee database
- Maintenance of stocks, etc.

Banking



Today, banking is almost totally dependent on computers. Banks provide the following facilities –

- Online accounting facility, which includes checking current balance, making deposits and overdrafts, checking interest charges, shares, and trustee records.
- ATM machines which are completely automated are making it even easier for customers to deal with banks.

Insurance

+ =		Backspace	Print Screen	Home
	3	↓	Delete	End
",	Insu	urance	Page Up	Page Down
	S	hift	Insert	1
11	Alt	Ctrl	+	+

Insurance companies are keeping all records up-to-date with the help of computers. Insurance companies, finance houses, and stock broking firms are widely using computers for their concerns.

Insurance companies are maintaining a database of all clients with information showing -

- Procedure to continue with policies
- Starting date of the policies
- Next due installment of a policy
- Maturity date
- Interests due
- Survival benefits
- Bonus

Education



The computer helps in providing a lot of facilities in the education system.

- The computer provides a tool in the education system known as CBE (Computer Based Education).
- CBE involves control, delivery, and evaluation of learning.
- Computer education is rapidly increasing the graph of number of computer students.
- There are a number of methods in which educational institutions can use a computer to educate the students.
- It is used to prepare a database about performance of a student and analysis is carried out on this basis.

Marketing

In marketing, uses of the computer are following -



- Advertising With computers, advertising professionals create art and graphics, write and revise copy, and print and disseminate ads with the goal of selling more products.
- Home Shopping Home shopping has been made possible through the use of computerized catalogues that provide access to product information and permit direct entry of orders to be filled by the customers.

Healthcare

Computers have become an important part in hospitals, labs, and dispensaries. They are being used in hospitals to keep the record of patients and medicines. It is also used in scanning and diagnosing different diseases. ECG, EEG, ultrasounds and CT scans, etc. are also done by computerized machines.

Following are some major fields of health care in which computers are used.



- **Diagnostic System** Computers are used to collect data and identify the cause of illness.
- Lab-diagnostic System All tests can be done and the reports are prepared by computer.
- **Patient Monitoring System** These are used to check the patient's signs for abnormality such as in Cardiac Arrest, ECG, etc.

- **Pharma Information System** Computer is used to check drug labels, expiry dates, harmful side effects, etc.
- Surgery Nowadays, computers are also used in performing surgery.

Engineering Design

Computers are widely used for Engineering purpose.

One of the major areas is CAD (Computer Aided Design) that provides creation and modification of images. Some of the fields are –



- Structural Engineering Requires stress and strain analysis for design of ships, buildings, budgets, airplanes, etc.
- **Industrial Engineering** Computers deal with design, implementation, and improvement of integrated systems of people, materials, and equipment.
- Architectural Engineering Computers help in planning towns, designing buildings, determining a range of buildings on a site using both 2D and 3D drawings.

Military



Computers are largely used in defence. Modern tanks, missiles, weapons, etc. Military also employs computerized control systems. Some military areas where a computer has been used are –

- Missile Control
- Military Communication
- Military Operation and Planning
- Smart Weapons

Communication

Communication is a way to convey a message, an idea, a picture, or speech that is received and understood clearly and correctly by the person for whom it is meant. Some main areas in this category are -



- E-mail
- Chatting
- Usenet
- FTP
- Telnet
- Video-conferencing

Government

Computers play an important role in government services. Some major fields in this category are -

MANAG



- Budgets
- Sales tax department
- Income tax department
- Computation of male/female ratio
- Computerization of voters lists
- Computerization of PAN card
- Weather forecasting

Computer - Generations

Generation in computer terminology is a change in technology a computer is/was being used. Initially, the generation term was used to distinguish between varying hardware

technologies. Nowadays, generation includes both hardware and software, which together make up an entire computer system.

There are five computer generations known till date. Each generation has been discussed in detail along with their time period and characteristics. In the following table, approximate dates against each generation has been mentioned, which are normally accepted.

Following are the main five generations of computers.

S.No	Generation & Description
1	First Generation The period of first generation: 1946-1959. Vacuum tube based.
2	Second Generation The period of second generation: 1959-1965. Transistor based.
3	Third Generation The period of third generation: 1965-1971. Integrated Circuit based.
4	Fourth Generation The period of fourth generation: 1971-1980. VLSI microprocessor based.
5	Fifth Generation The period of fifth generation: 1980-onwards. ULSI microprocessor based.



Computer - Types

Computers can be broadly classified by their speed and computing power.

S.No.	Туре	Specifications
1	PC (Personal Computer)	It is a single user computer system having moderately powerful microprocessor
2	Workstation	It is also a single user computer system, similar to personal computer however has a more powerful microprocessor.
3	Mini Computer	It is a multi-user computer system, capable of supporting hundreds of users simultaneously.
4	Main Frame	It is a multi-user computer system, capable of supporting hundreds of users simultaneously. Software technology is different from minicomputer.
5	Supercomputer	It is an extremely fast computer, which can execute hundreds of millions of instructions per second.

PC (Personal Computer)



A PC can be defined as a small, relatively inexpensive computer designed for an individual user. PCs are based on the microprocessor technology that enables manufacturers to put an entire CPU on one chip. Businesses use personal computers for word processing, accounting, desktop publishing, and for running spreadsheet and database management applications. At home, the most popular use for personal computers is playing games and surfing the Internet.

Although personal computers are designed as single-user systems, these systems are normally linked together to form a network. In terms of power, now-a-days high-end models of the Macintosh and PC offer the same computing power and graphics capability as lowend workstations by Sun Microsystems, Hewlett-Packard, and Dell.

Workstation



Workstation is a computer used for engineering applications (CAD/CAM), desktop publishing, software development, and other such types of applications which require a moderate amount of computing power and relatively high quality graphics capabilities. Workstations generally come with a large, high-resolution graphics screen, large amount of RAM, inbuilt network support, and a graphical user interface. Most workstations also have mass storage device such as a disk drive, but a special type of workstation, called diskless workstation, comes without a disk drive.

Common operating systems for workstations are UNIX and Windows NT. Like PC, workstations are also single-user computers like PC but are typically linked together to form a local-area network, although they can also be used as stand-alone systems.

Minicomputer

It is a midsize multi-processing system capable of supporting up to 250 users simultaneously.



Mainframe

Mainframe is very large in size and is an expensive computer capable of supporting hundreds or even thousands of users simultaneously. Mainframe executes many programs concurrently and supports many simultaneous execution of programs.



Supercomputer

Supercomputers are one of the fastest computers currently available. Supercomputers are very expensive and are employed for specialized applications that require immense amount of mathematical calculations (number crunching).



For example, weather forecasting, scientific simulations, (animated) graphics, fluid dynamic calculations, nuclear energy research, electronic design, and analysis of geological data (e.g. in petrochemical prospecting).

Computer - Components

All types of computers follow the same basic logical structure and perform the following five basic operations for converting raw input data into information useful to their users.

S.No.	Operation	Description
1	Take Input	The process of entering data and instructions into the computer system.
2	Store Data	Saving data and instructions so that they are available for processing as and when required.
3	Processing Data	Performing arithmetic, and logical operations on data in order to convert them into useful information.
4	Output Information	The process of producing useful information or results for the user, such as a printed report or visual display.
5	Control the	Directs the manner and sequence in which all of



Input Unit

This unit contains devices with the help of which we enter data into the computer. This unit creates a link between the user and the computer. The input devices translate the information into a form understandable by the computer.

CPU (Central Processing Unit)

CPU is considered as the brain of the computer. CPU performs all types of data processing operations. It stores data, intermediate results, and instructions (program). It controls the operation of all parts of the computer.

CPU itself has the following three components -

- ALU (Arithmetic Logic Unit)
- Memory Unit
- Control Unit

Output Unit

The output unit consists of devices with the help of which we get the information from the computer. This unit is a link between the computer and the users. Output devices translate the computer's output into a form understandable by the users.

Computer - CPU(Central Processing Unit)

Central Processing Unit (CPU) consists of the following features

- CPU is considered as the brain of the computer.
- CPU performs all types of data processing operations.
- It stores data, intermediate results, and instructions (program).
- It controls the operation of all parts of the computer.



CPU itself has following three components.

- Memory or Storage Unit
- Control Unit
- ALU(Arithmetic Logic Unit)



REDIT

Memory or Storage Unit

This unit can store instructions, data, and intermediate results. This unit supplies information to other units of the computer when needed. It is also known as internal storage unit or the main memory or the primary storage or Random Access Memory (RAM).

Its size affects speed, power, and capability. Primary memory and secondary memory are two types of memories in the computer. Functions of the memory unit are –

- It stores all the data and the instructions required for processing.
- It stores intermediate results of processing.
- It stores the final results of processing before these results are released to an output device.
- All inputs and outputs are transmitted through the main memory.

Control Unit

This unit controls the operations of all parts of the computer but does not carry out any actual data processing operations.

Functions of this unit are -

• It is responsible for controlling the transfer of data and instructions among other units of a computer.

• It manages and coordinates all the units of the computer.

ANAG

- It obtains the instructions from the memory, interprets them, and directs the operation of the computer.
- It communicates with Input/Output devices for transfer of data or results from storage.
- It does not process or store data.

ALU (Arithmetic Logic Unit)

This unit consists of two subsections namely,

- Arithmetic Section
- Logic Section

Arithmetic Section

Function of arithmetic section is to perform arithmetic operations like addition, subtraction, multiplication, and division. All complex operations are done by making repetitive use of the above operations.

Logic Section

Function of logic section is to perform logic operations such as comparing, selecting, matching, and merging of data.

Computer - Input Devices

Following are some of the important input devices which are used in a computer -

- Keyboard
- Mouse
- Joy Stick
- Light pen
- Track Ball
- Scanner
- Graphic Tablet
- Microphone
- Magnetic Ink Card Reader(MICR)
- Optical Character Reader(OCR)
- Bar Code Reader
- Optical Mark Reader(OMR)

Keyboard

COPYRIGHT FIMT 2020

Keyboard is the most common and very popular input device which helps to input data to the computer. The layout of the keyboard is like that of traditional typewriter, although there are some additional keys provided for performing additional functions.



Keyboards are of two sizes 84 keys or 101/102 keys, but now keyboards with 104 keys or 108 keys are also available for Windows and Internet.

The keys on the keyboard are as follows -

S.No	Keys & Description
1	Typing Keys These keys include the letter keys (A-Z) and digit keys (09) which generally give the same layout as that of typewriters.
2	Numeric Keypad It is used to enter the numeric data or cursor movement. Generally, it consists of a set of 17 keys that are laid out in the same configuration used by most adding machines and calculators.
3	Function Keys The twelve function keys are present on the keyboard which are arranged in a row at the top of the keyboard. Each function key has a unique meaning and is used for some specific purpose.
4	Control keys These keys provide cursor and screen control. It includes four directional arrow keys. Control keys also include Home, End, Insert, Delete, Page Up, Page Down, Control(Ctrl), Alternate(Alt), Escape(Esc).
5	Special Purpose Keys Keyboard also contains some special purpose keys such as Enter, Shift, Caps Lock, Num Lock, Space bar, Tab, and Print Screen.

Mouse

Mouse is the most popular pointing device. It is a very famous cursor-control device having a small palm size box with a round ball at its base, which senses the movement of the mouse and sends corresponding signals to the CPU when the mouse buttons are pressed.

Generally, it has two buttons called the left and the right button and a wheel is present between the buttons. A mouse can be used to control the position of the cursor on the screen, but it cannot be used to enter text into the computer.



- Easy to use
- Not very expensive
- Moves the cursor faster than the arrow keys of the keyboard.

Joystick

Joystick is also a pointing device, which is used to move the cursor position on a monitor screen. It is a stick having a spherical ball at its both lower and upper ends. The lower spherical ball moves in a socket. The joystick can be moved in all four directions.



The function of the joystick is similar to that of a mouse. It is mainly used in Computer Aided Designing (CAD) and playing computer games.

Light Pen

Light pen is a pointing device similar to a pen. It is used to select a displayed menu item or draw pictures on the monitor screen. It consists of a photocell and an optical system placed in a small tube.



When the tip of a light pen is moved over the monitor screen and the pen button is pressed, its photocell sensing element detects the screen location and sends the corresponding signal to the CPU.

Track Ball

Track ball is an input device that is mostly used in notebook or laptop computer, instead of a mouse. This is a ball which is half inserted and by moving fingers on the ball, the pointer can be moved.



Since the whole device is not moved, a track ball requires less space than a mouse. A track ball comes in various shapes like a ball, a button, or a square.

Scanner

Scanner is an input device, which works more like a photocopy machine. It is used when some information is available on paper and it is to be transferred to the hard disk of the computer for further manipulation.



Scanner captures images from the source which are then converted into a digital form that can be stored on the disk. These images can be edited before they are printed.

Digitizer

Digitizer is an input device which converts analog information into digital form. Digitizer can convert a signal from the television or camera into a series of numbers that could be stored in a computer. They can be used by the computer to create a picture of whatever the camera had been pointed at.



Digitizer is also known as Tablet or Graphics Tablet as it converts graphics and pictorial data into binary inputs. A graphic tablet as digitizer is used for fine works of drawing and image manipulation applications.

Microphone

Microphone is an input device to input sound that is then stored in a digital form.



The microphone is used for various applications such as adding sound to a multimedia presentation or for mixing music.

Magnetic Ink Card Reader (MICR)

MICR input device is generally used in banks as there are large number of cheques to be processed every day. The bank's code number and cheque number are printed on the cheques with a special type of ink that contains particles of magnetic material that are machine readable.



This reading process is called Magnetic Ink Character Recognition (MICR). The main advantages of MICR is that it is fast and less error prone.

Optical Character Reader (OCR)

COPYRIGHT FIMT 2020

OCR is an input device used to read a printed text.



OCR scans the text optically, character by character, converts them into a machine readable code, and stores the text on the system memory.

Bar Code Readers

Bar Code Reader is a device used for reading bar coded data (data in the form of light and dark lines). Bar coded data is generally used in labelling goods, numbering the books, etc. It may be a handheld scanner or may be embedded in a stationary scanner.



Bar Code Reader scans a bar code image, converts it into an alphanumeric value, which is then fed to the computer that the bar code reader is connected to.

Optical Mark Reader (OMR)

OMR is a special type of optical scanner used to recognize the type of mark made by pen or pencil. It is used where one out of a few alternatives is to be selected and marked.



It is specially used for checking the answer sheets of examinations having multiple choice questions.

Computer - Output Devices

Following are some of the important output devices used in a computer.

- Monitors
- Graphic Plotter

• Printer

Monitors

Monitors, commonly called as **Visual Display Unit** (VDU), are the main output device of a computer. It forms images from tiny dots, called pixels that are arranged in a rectangular form. The sharpness of the image depends upon the number of pixels.

There are two kinds of viewing screen used for monitors.

- Cathode-Ray Tube (CRT)
- Flat-Panel Display

Cathode-Ray Tube (CRT) Monitor

The CRT display is made up of small picture elements called pixels. The smaller the pixels, the better the image clarity or resolution. It takes more than one illuminated pixel to form a whole character, such as the letter 'e' in the word help.



A finite number of characters can be displayed on a screen at once. The screen can be divided into a series of character boxes - fixed location on the screen where a standard character can be placed. Most screens are capable of displaying 80 characters of data horizontally and 25 lines vertically.

There are some disadvantages of CRT -

- Large in Size
- High power consumption

Flat-Panel Display Monitor

The flat-panel display refers to a class of video devices that have reduced volume, weight and power requirement in comparison to the CRT. You can hang them on walls or wear them on your wrists. Current uses of flat-panel displays include calculators, video games, monitors, laptop computer, and graphics display.



The flat-panel display is divided into two categories -

- Emissive Displays Emissive displays are devices that convert electrical energy into light. For example, plasma panel and LED (Light-Emitting Diodes).
- Non-Emissive Displays Non-emissive displays use optical effects to convert sunlight or light from some other source into graphics patterns. For example, LCD (Liquid-Crystal Device).

Printers

Printer is an output device, which is used to print information on paper. There are two types of printers –

- Impact Printers
- Non-Impact Printers

Impact Printers

Impact printers print the characters by striking them on the ribbon, which is then pressed on the paper.

Characteristics of Impact Printers are the following -

- Very low consumable costs
- Very noisy
- Useful for bulk printing due to low cost
- There is physical contact with the paper to produce an image

These printers are of two types -

- Character printers
- Line printers

Character Printers

Character printers are the printers which print one character at a time.

These are further divided into two types:

- Dot Matrix Printer(DMP)
- Daisy Wheel

Dot Matrix Printer

In the market, one of the most popular printers is Dot Matrix Printer. These printers are popular because of their ease of printing and economical price. Each character printed is in the form of pattern of dots and head consists of a Matrix of Pins of size (5*7, 7*9, 9*7 or 9*9) which come out to form a character which is why it is called Dot Matrix Printer.



- Inexpensive
- Widely Used
- Other language characters can be printed

Disadvantages

- Slow Speed
- Poor Quality

Daisy Wheel

Head is lying on a wheel and pins corresponding to characters are like petals of Daisy (flower) which is why it is called Daisy Wheel Printer. These printers are generally used for word-processing in offices that require a few letters to be sent here and there with very nice quality.



Advantages

- More reliable than DMP
- Better quality
- Fonts of character can be easily changed

Disadvantages

- Slower than DMP
- Noisy
- More expensive than DMP

Line Printers

Line printers are the printers which print one line at a time.



AG

These are of two types -

- Drum Printer
- Chain Printer

Drum Printer

This printer is like a drum in shape hence it is called drum printer. The surface of the drum is divided into a number of tracks. Total tracks are equal to the size of the paper, i.e. for a paper width of 132 characters, drum will have 132 tracks. A character set is embossed on the track. Different character sets available in the market are 48 character set, 64 and 96 characters set. One rotation of drum prints one line. Drum printers are fast in speed and can print 300 to 2000 lines per minute.

Advantages

• Very high speed

Disadvantages

• Very expensive

Characters fonts cannot be changed •

Chain Printer

In this printer, a chain of character sets is used, hence it is called Chain Printer. A standard character set may have 48, 64, or 96 characters.

Advantages

Character fonts can easily be changed.

S.

Different languages can be used with the same printer. MANAG

Disadvantages

Noisy •

Non-impact Printers

Non-impact printers print the characters without using the ribbon. These printers print a complete page at a time, thus they are also called as Page Printers.

These printers are of two types -

- Laser Printers
- **Inkjet Printers**

Characteristics of Non-impact Printers

- Faster than impact printers •
- They are not noisy •
- High quality
- Supports many fonts and different character size •

Laser Printers

These are non-impact page printers. They use laser lights to produce the dots needed to form the characters to be printed on a page. 015 & 14001:2015



Advantages

- Very high speed
- Very high quality output
- Good graphics quality
- Supports many fonts and different character size

MANA

Disadvantages

- Expensive
- Cannot be used to produce multiple copies of a document in a single printing

Inkjet Printers

Inkjet printers are non-impact character printers based on a relatively new technology. They print characters by spraying small drops of ink onto paper. Inkjet printers produce high quality output with presentable features.



They make less noise because no hammering is done and these have many styles of printing modes available. Color printing is also possible. Some models of Inkjet printers can produce multiple copies of printing also.

Advantages

- High quality printing
- More reliable

Disadvantages

- Expensive as the cost per page is high
- Slow as compared to laser printer

Computer - Memory

A memory is just like a human brain. It is used to store data and instructions. Computer memory is the storage space in the computer, where data is to be processed and instructions required for processing are stored. The memory is divided into large number of small parts called cells. Each location or cell has a unique address, which varies from zero to memory size minus one. For example, if the computer has 64k words, then this memory unit has $64 \times 1024 = 65536$ memory locations. The address of these locations varies from 0 to 65535. Memory is primarily of three types –

AAC ACCREDITED

- Cache Memory
- Primary Memory/Main Memory
- Secondary Memory

Cache Memory

Cache memory is a very high speed semiconductor memory which can speed up the CPU. It acts as a buffer between the CPU and the main memory. It is used to hold those parts of data and program which are most frequently used by the CPU. The parts of data and programs are transferred from the disk to cache memory by the operating system, from where the CPU can access them.



Advantages

The advantages of cache memory are as follows -

- Cache memory is faster than main memory.
- It consumes less access time as compared to main memory.
- It stores the program that can be executed within a short period of time.
- It stores data for temporary use.

Disadvantages

The disadvantages of cache memory are as follows –

1. 1. 62

- Cache memory has limited capacity.
- It is very expensive.

Primary Memory (Main Memory)

Primary memory holds only those data and instructions on which the computer is currently working. It has a limited capacity and data is lost when power is switched off. It is generally made up of semiconductor device. These memories are not as fast as registers. The data and instruction required to be processed resides in the main memory. It is divided into two subcategories RAM and ROM.

Characteristics of Main Memory

- These are semiconductor memories.
- It is known as the main memory.
- Usually volatile memory.
- Data is lost in case power is switched off.
- It is the working memory of the computer.
- Faster than secondary memories.
- A computer cannot run without the primary memory.

Secondary Memory

This type of memory is also known as external memory or non-volatile. It is slower than the main memory. These are used for storing data/information permanently. CPU directly does not access these memories, instead they are accessed via input-output routines. The contents of secondary memories are first transferred to the main memory, and then the CPU can access it. For example, disk, CD-ROM, DVD, etc.



Characteristics of Secondary Memory

- These are magnetic and optical memories.
- It is known as the backup memory.
- It is a non-volatile memory.
- Data is permanently stored even if power is switched off.
- It is used for storage of data in a computer.
- Computer may run without the secondary memory.
- Slower than primary memories.

Random Access Memory

RAM (Random Access Memory) is the internal memory of the CPU for storing data, program, and program result. It is a read/write memory which stores data until the machine is working. As soon as the machine is switched off, data is erased.



Access time in RAM is independent of the address, that is, each storage location inside the memory is as easy to reach as other locations and takes the same amount of time. Data in the RAM can be accessed randomly but it is very expensive.

RAM is volatile, i.e. data stored in it is lost when we switch off the computer or if there is a power failure. Hence, a backup Uninterruptible Power System (UPS) is often used with
computers. RAM is small, both in terms of its physical size and in the amount of data it can hold.

RAM is of two types -

- Static RAM (SRAM)
- Dynamic RAM (DRAM)

Static RAM (SRAM)

The word **static** indicates that the memory retains its contents as long as power is being supplied. However, data is lost when the power gets down due to volatile nature. SRAM chips use a matrix of 6-transistors and no capacitors. Transistors do not require power to prevent leakage, so SRAM need not be refreshed on a regular basis.

There is extra space in the matrix, hence SRAM uses more chips than DRAM for the same amount of storage space, making the manufacturing costs higher. SRAM is thus used as cache memory and has very fast access.

Characteristic of Static RAM

- Long life
- No need to refresh
- Faster
- Used as cache memory
- Large size
- Expensive
- High power consumption

Dynamic RAM (DRAM)

DRAM, unlike SRAM, must be continually **refreshed** in order to maintain the data. This is done by placing the memory on a refresh circuit that rewrites the data several hundred times per second. DRAM is used for most system memory as it is cheap and small. All DRAMs are made up of memory cells, which are composed of one capacitor and one transistor.

Characteristics of Dynamic RAM

- Short data lifetime
- Needs to be refreshed continuously
- Slower as compared to SRAM
- Used as RAM

- Smaller in size
- Less expensive
- Less power consumption

Computer - Read Only Memory

ROM stands for **Read Only Memory**. The memory from which we can only read but cannot write on it. This type of memory is non-volatile. The information is stored permanently in such memories during manufacture. A ROM stores such instructions that are required to start a computer. This operation is referred to as **bootstrap**. ROM chips are not only used in the computer but also in other electronic items like washing machine and microwave oven.



Let us now discuss the various types of ROMs and their characteristics.

MROM (Masked ROM)

The very first ROMs were hard-wired devices that contained a pre-programmed set of data or instructions. These kind of ROMs are known as masked ROMs, which are inexpensive.

PROM (Programmable Read Only Memory)

PROM is read-only memory that can be modified only once by a user. The user buys a blank PROM and enters the desired contents using a PROM program. Inside the PROM chip, there are small fuses which are burnt open during programming. It can be programmed only once and is not erasable.

EPROM (Erasable and Programmable Read Only Memory)

EPROM can be erased by exposing it to ultra-violet light for a duration of up to 40 minutes. Usually, an EPROM eraser achieves this function. During programming, an electrical charge is trapped in an insulated gate region. The charge is retained for more than 10 years because the charge has no leakage path. For erasing this charge, ultra-violet light is passed through a quartz crystal window (lid). This exposure to ultra-violet light dissipates the charge. During normal use, the quartz lid is sealed with a sticker.

EEPROM (Electrically Erasable and Programmable Read Only Memory)

EEPROM is programmed and erased electrically. It can be erased and reprogrammed about ten thousand times. Both erasing and programming take about 4 to 10 ms (millisecond). In EEPROM, any location can be selectively erased and programmed. EEPROMs can be erased one byte at a time, rather than erasing the entire chip. Hence, the process of reprogramming is flexible but slow.

Advantages of ROM

The advantages of ROM are as follows –

- Non-volatile in nature
- Cannot be accidentally changed
- Cheaper than RAMs
- Easy to test
- More reliable than RAMs
- Static and do not require refreshing
- Contents are always known and can be verified

Computer - Motherboard

The motherboard serves as a single platform to connect all of the parts of a computer together. It connects the CPU, memory, hard drives, optical drives, video card, sound card, and other ports and expansion cards directly or via cables. It can be considered as the backbone of a computer.



Features of Motherboard

A motherboard comes with following features -

- Motherboard varies greatly in supporting various types of components.
- Motherboard supports a single type of CPU and few types of memories.
- Video cards, hard disks, sound cards have to be compatible with the motherboard to function properly.

• Motherboards, cases, and power supplies must be compatible to work properly together.

NAAC ACCREDITED

Popular Manufacturers

Following are the popular manufacturers of the motherboard.

INAMAG

- Intel
- ASUS
- AOpen
- ABIT
- Biostar
- Gigabyte
- MSI

Description of Motherboard

The motherboard is mounted inside the case and is securely attached via small screws through pre-drilled holes. Motherboard contains ports to connect all of the internal components. It provides a single socket for CPU, whereas for memory, normally one or more slots are available. Motherboards provide ports to attach the floppy drive, hard drive, and optical drives via ribbon cables. Motherboard carries fans and a special port designed for power supply.

There is a peripheral card slot in front of the motherboard using which video cards, sound cards, and other expansion cards can be connected to the motherboard.

On the left side, motherboards carry a number of ports to connect the monitor, printer, mouse, keyboard, speaker, and network cables. Motherboards also provide USB ports, which allow compatible devices to be connected in plug-in/plug-out fashion. For example, pen drive, digital cameras, etc.

Computer - Memory Units

Memory unit is the amount of data that can be stored in the storage unit. This storage capacity is expressed in terms of Bytes.

The following table explains the main memory storage units -

S.No.	Unit & Description
1	Bit (Binary Digit) A binary digit is logical 0 and 1 representing a passive or an active state of a component in an electric circuit.

2	Nibble A group of 4 bits is called nibble.
3	Byte A group of 8 bits is called byte. A byte is the smallest unit, which can represent a data item or a character.
4	Word A computer word, like a byte, is a group of fixed number of bits processed as a unit, which varies from computer to computer but is fixed for each computer. The length of a computer word is called word-size or word length. It may be as small as 8 bits or may be as long as 96 bits. A computer stores the information in the form of computer words.

The following table lists some higher storage units -

S.No.	Unit & Description		
1	Kilobyte (KB)1 KB = 1024 Bytes		
2	Megabyte (MB) 1 MB = 1024 KB		
3	GigaByte (GB) 1 GB = 1024 MB		
4	TeraByte (TB) 1 TB = 1024 GB		
5	PetaByte (PB) 1 PB = 1024 TB		

Computer - Ports

A port is a physical docking point using which an external device can be connected to the computer. It can also be programmatic docking point through which information flows from a program to the computer or over the Internet.

Characteristics of Ports

A port has the following characteristics -

- External devices are connected to a computer using cables and ports.
- Ports are slots on the motherboard into which a cable of external device is plugged in.
- Examples of external devices attached via ports are the mouse, keyboard, monitor, microphone, speakers, etc.



Let us now discuss a few important types of ports -

Serial Port

- Used for external modems and older computer mouse
- Two versions: 9 pin, 25 pin model
- Data travels at 115 kilobits per second

Parallel Port

- Used for scanners and printers
- Also called printer port
- 25 pin model
- IEEE 1284-compliant Centronics port

PS/2 Port

- Used for old computer keyboard and mouse
- Also called mouse port
- Most of the old computers provide two PS/2 port, each for the mouse and keyboard
- IEEE 1284-compliant Centronics port

Universal Serial Bus (or USB) Port

- It can connect all kinds of external USB devices such as external hard disk, printer, scanner, mouse, keyboard, etc.
- It was introduced in 1997.
- Most of the computers provide two USB ports as minimum.
- Data travels at 12 megabits per seconds.
- USB compliant devices can get power from a USB port.

VGA Port

- Connects monitor to a computer's video card.
- It has 15 holes.
- Similar to the serial port connector. However, serial port connector has pins, VGA port has holes.

Power Connector

- Three-pronged plug.
- Connects to the computer's power cable that plugs into a power bar or wall socket.

Firewire Port

- Transfers large amount of data at very fast speed.
- Connects camcorders and video equipment to the computer.
- Data travels at 400 to 800 megabits per seconds.
- Invented by Apple.
- It has three variants: 4-Pin FireWire 400 connector, 6-Pin FireWire 400 connector, and 9-Pin FireWire 800 connector.

Modem Port

• Connects a PC's modem to the telephone network.

Ethernet Port

- Connects to a network and high speed Internet.
- Connects the network cable to a computer.
- This port resides on an Ethernet Card.
- Data travels at 10 megabits to 1000 megabits per seconds depending upon the network bandwidth.

Game Port

- Connect a joystick to a PC
- Now replaced by USB

Digital Video Interface, DVI port

- Connects Flat panel LCD monitor to the computer's high-end video graphic cards.
- Very popular among video card manufacturers.

Sockets

• Sockets connect the microphone and speakers to the sound card of the computer.

Computer - Hardware

Hardware represents the physical and tangible components of a computer, i.e. the components that can be seen and touched.

Examples of Hardware are the following -

- Input devices keyboard, mouse, etc.
- Output devices printer, monitor, etc.
- Secondary storage devices Hard disk, CD, DVD, etc.
- Internal components CPU, motherboard, RAM, etc.



Relationship between Hardware and Software

- Hardware and software are mutually dependent on each other. Both of them must work together to make a computer produce a useful output.
- Software cannot be utilized without supporting hardware.
- Hardware without a set of programs to operate upon cannot be utilized and is useless.
- To get a particular job done on the computer, relevant software should be loaded into the hardware.
- Hardware is a one-time expense.
- Software development is very expensive and is a continuing expense.
- Different software applications can be loaded on a hardware to run different jobs.
- A software acts as an interface between the user and the hardware.
- If the hardware is the 'heart' of a computer system, then the software is its 'soul'. Both are complementary to each other.

14001:2015

Computer - Software

Software is a set of programs, which is designed to perform a well-defined function. A program is a sequence of instructions written to solve a particular problem.

There are two types of software -

- System Software
- Application Software

System Software

The system software is a collection of programs designed to operate, control, and extend the processing capabilities of the computer itself. System software is generally prepared by the computer manufacturers. These software products comprise of programs written in low-level languages, which interact with the hardware at a very basic level. System software serves as the interface between the hardware and the end users.

Some examples of system software are Operating System, Compilers, Interpreter, Assemblers, etc.



Here is a list of some of the most prominent features of a system software -

- Close to the system
- Fast in speed
- Difficult to design
- Difficult to understand
- Less interactive
- Smaller in size
- Difficult to manipulate
- Generally written in low-level language

Application Software

Application software products are designed to satisfy a particular need of a particular environment. All software applications prepared in the computer lab can come under the category of Application software.

Application software may consist of a single program, such as Microsoft's notepad for writing and editing a simple text. It may also consist of a collection of programs, often called a software package, which work together to accomplish a task, such as a spreadsheet package.

Examples of Application software are the following -

- Payroll Software
- Student Record Software

- Inventory Management Software
- Income Tax Software
- Railways Reservation Software
- Microsoft Office Suite Software
- Microsoft Word
- Microsoft Excel
- Microsoft PowerPoint



Features of application software are as follows -

- Close to the user
- Easy to design
- More interactive
- Slow in speed
- Generally written in high-level language
- Easy to understand
- Easy to manipulate and use
- Bigger in size and requires large storage space

Computer - Number System

When we type some letters or words, the computer translates them in numbers as computers can understand only numbers. A computer can understand the positional number system where there are only a few symbols called digits and these symbols represent different values depending on the position they occupy in the number.

Microsoft

The value of each digit in a number can be determined using -

- The digit
- The position of the digit in the number
- The base of the number system (where the base is defined as the total number of digits available in the number system)

Decimal Number System

The number system that we use in our day-to-day life is the decimal number system. Decimal number system has base 10 as it uses 10 digits from 0 to 9. In decimal number

system, the successive positions to the left of the decimal point represent units, tens, hundreds, thousands, and so on.

Each position represents a specific power of the base (10). For example, the decimal number 1234 consists of the digit 4 in the units position, 3 in the tens position, 2 in the hundreds position, and 1 in the thousands position. Its value can be written as

 $(1 \ge 1000) + (2 \ge 100) + (3 \ge 10) + (4 \ge 1)$ $(1 \ge 10^3) + (2 \ge 10^2) + (3 \ge 10^1) + (4 \ge 10^0)$ 1000 + 200 + 30 + 41234

As a computer programmer or an IT professional, you should understand the following number systems which are frequently used in computers.

S.No.	Number System and Description	
1	Binary Number System Base 2. Digits used : 0, 1	
2	Octal Number System Base 8. Digits used : 0 to 7	
3	Hexa Decimal Number System Base 16. Digits used: 0 to 9, Letters used : A- F	

Binary Number System

Characteristics of the binary number system are as follows -

- Uses two digits, 0 and 1
- Also called as base 2 number system
- Each position in a binary number represents a **0** power of the base (2). Example 2^0

RFILL

Last position in a binary number represents a x power of the base (2). Example 2^x where x represents the last position - 1.

01:2015 & 14001:2015

Example

Binary Number: 10101₂

Calculating Decimal Equivalent -

Step	Binary Number	Decimal Number
Step 1	101012	$((1 x 2^4) + (0 x 2^3) + (1 x 2^2) + (0 x 2^1) + (1 x 2^0))_{10}$

Step 2	101012	$(16+0+4+0+1)_{10}$
Step 3	101012	2110

Note -10101_2 is normally written as 10101.

Octal Number System

Characteristics of the octal number system are as follows -

- Uses eight digits, 0,1,2,3,4,5,6,7
- Also called as base 8 number system
- Each position in an octal number represents a $\mathbf{0}$ power of the base (8). Example 8^0
- Last position in an octal number represents a x power of the base (8). Example
 8^x where x represents the last position 1

Example

Octal Number: 125708

Calculating Decimal Equivalent -

Step	Octal Number	Decimal Number
Step 1	12570 ₈	$((1 x 8^4) + (2 x 8^3) + (5 x 8^2) + (7 x 8^1) + (0 x 8^0))_{10}$
Step 2	125708	$(4096 + 1024 + 320 + 56 + 0)_{10}$
Step 3	125708	549610

Note -12570_8 is normally written as 12570.

Hexadecimal Number System

Characteristics of hexadecimal number system are as follows -

- Uses 10 digits and 6 letters, 0, 1, 2, 3, 4, 5, 6, 7, 8, 9, A, B, C, D, E, F
- Letters represent the numbers starting from 10. A = 10. B = 11, C = 12, D = 13, E = 14, F = 15
- Also called as base 16 number system
- Each position in a hexadecimal number represents a **0** power of the base (16). Example, 16⁰
- Last position in a hexadecimal number represents a x power of the base (16).
 Example 16^x where x represents the last position 1

Example

Hexadecimal Number: 19FDE₁₆

Calculating Decimal Equivalent -

Step	Binary Number	Decimal Number
Step 1	19FDE ₁₆	$((1 x 16^4) + (9 x 16^3) + (F x 16^2) + (D x 16^1) + (E x 16^0))_{10}$
Step 2	19FDE ₁₆	$((1 x 16^4) + (9 x 16^3) + (15 x 16^2) + (13 x 16^1) + (14 x 16^0))_{10}$
Step 3	19FDE ₁₆	$(65536 + 36864 + 3840 + 208 + 14)_{10}$
Step 4	19FDE ₁₆	10646210

Note – 19FDE₁₆ is normally written as 19FDE.

Computer - Number Conversion

There are many methods or techniques which can be used to convert numbers from one base to another. In this chapter, we'll demonstrate the following –

- Decimal to Other Base System
- Other Base System to Decimal
- Other Base System to Non-Decimal
- Shortcut method Binary to Octal
- Shortcut method Octal to Binary
- Shortcut method Binary to Hexadecimal
- Shortcut method Hexadecimal to Binary

Decimal to Other Base System

Step 1 – Divide the decimal number to be converted by the value of the new base.

Step 2 – Get the remainder from Step 1 as the rightmost digit (least significant digit) of the new base number.

Step 3 – Divide the quotient of the previous divide by the new base.

Step 4 – Record the remainder from Step 3 as the next digit (to the left) of the new base number.

Repeat Steps 3 and 4, getting remainders from right to left, until the quotient becomes zero in Step 3.

The last remainder thus obtained will be the Most Significant Digit (MSD) of the new base number.

Example

Decimal Number: 29₁₀

Calculating Binary Equivalent -

Step	Operation	Result	Remainder
Step 1	29 / 2	14	1
Step 2	14 / 2	7	0
Step 3	7 / 2	3	S 21
Step 4	3 / 2	1	1
Step 5	1 / 2	0	1

NAAC ACCREDITED

As mentioned in Steps 2 and 4, the remainders have to be arranged in the reverse order so that the first remainder becomes the Least Significant Digit (LSD) and the last remainder becomes the Most Significant Digit (MSD).

Decimal Number : 29_{10} = Binary Number : 11101_2 .

Other Base System to Decimal System

Step 1 – Determine the column (positional) value of each digit (this depends on the position of the digit and the base of the number system).

Step 2 - Multiply the obtained column values (in Step 1) by the digits in the corresponding columns.

Step 3 – Sum the products calculated in Step 2. The total is the equivalent value in decimal.

Example

,9001:2015 & 14001:2015 Binary Number: 11101₂

Calculating Decimal Equivalent -

Step	Binary Number	Decimal Number
Step	111012	$((1 x 2^4) + (1 x 2^3) + (1 x 2^2) + (0 x 2^1) + (1 x 2^2))$

COPYRIGHT FIMT 2020

446 | Page

1		$(2^0))_{10}$
Step 2	111012	$(16 + 8 + 4 + 0 + 1)_{10}$
Step 3	111012	2910

Binary Number : 11101_2 = Decimal Number : 29_{10}

Other Base System to Non-Decimal System

Step 1 – Convert the original number to a decimal number (base 10).

Step 2 – Convert the decimal number so obtained to the new base number.

Example

Octal Number: 258

Calculating Binary Equivalent -

Step 1 - Convert to Decimal

Step	Octal Number	Decimal Number
Step 1	258	$((2 \times 8^1) + (5 \times 8^0))_{10}$
Step 2	258	$(16+5)_{10}$
Step 3	258	2110

10

Octal Number : 25_8 = Decimal Number : 21_{10}

Step 2 - Convert Decimal to Binary

Step	Operation	Result	Remainder
Step 1	21 / 2	10	1
Step 2	10 / 2	5	0
Step 3	5 / 2	2	
Step 4	2 / 2	8 140	0
Step 5	1 / 2	0	1

Decimal Number : 21_{10} = Binary Number : 10101_2

Octal Number : $25_8 = Binary Number : 10101_2$

Shortcut Method - Binary to Octal

Step 1 – Divide the binary digits into groups of three (starting from the right).

Step 2 – Convert each group of three binary digits to one octal digit.

Example

Binary Number : 10101₂

Calculating Octal Equivalent -

Step	Binary Number	Octal Number
Step 1	101012	010 101
Step 2	101012	2 ₈ 5 ₈
Step 3	101012	258

Binary Number : $10101_2 = \text{Octal Number} : 25_8$

Shortcut Method - Octal to Binary

Step 1 – Convert each octal digit to a 3-digit binary number (the octal digits may be treated as decimal for this conversion).

ALGEMEAN

Step 2 – Combine all the resulting binary groups (of 3 digits each) into a single binary number.

Example

Octal Number: 258

Calculating Binary Equivalent -

Step	Octal Number	Binary Number
Step 1	258	210 510
Step 2	258	0102 1012
Step 3	25 ₈	0101012

Octal Number : $25_8 = Binary Number : 10101_2$

Shortcut Method - Binary to Hexadecimal

Step 1 – Divide the binary digits into groups of four (starting from the right).

Step 2 – Convert each group of four binary digits to one hexadecimal symbol.

Example

Binary Number : 101012

Calculating hexadecimal Equivalent -

Step	Binary Number	Hexadecimal Number
Step 1	101012	0001 0101
Step 2	101012	$1_{10} 5_{10}$

Step 3	101012	1516
--------	--------	------

Binary Number : 10101_2 = Hexadecimal Number : 15_{16}

Shortcut Method - Hexadecimal to Binary

Step 1 – Convert each hexadecimal digit to a 4-digit binary number (the hexadecimal digits may be treated as decimal for this conversion).

Step 2 – Combine all the resulting binary groups (of 4 digits each) into a single binary number.

Example

Calculating Binary Equivalent –

Step	Hexadecimal Number	Binary Number
Step 1	1516	$1_{10} 5_{10}$
Step 2	1516	00012 01012
Step 3	1516	000101012

Hexadecimal Number : 15_{16} = Binary Number : 10101_2

Computer - Data and Information

Data can be defined as a representation of facts, concepts, or instructions in a formalized manner, which should be suitable for communication, interpretation, or processing by human or electronic machine.

Data is represented with the help of characters such as alphabets (A-Z, a-z), digits (0-9) or special characters (+,-,/,*,<,>,= etc.)

What is Information?

Information is organized or classified data, which has some meaningful values for the receiver. Information is the processed data on which decisions and actions are based.

For the decision to be meaningful, the processed data must qualify for the following characteristics -

- **Timely** Information should be available when required.
- Accuracy Information should be accurate.
- **Completeness** Information should be complete.



Data Processing Cycle

Data processing is the re-structuring or re-ordering of data by people or machine to increase their usefulness and add values for a particular purpose. Data processing consists of the following basic steps - input, processing, and output. These three steps constitute the data processing cycle.



- Input In this step, the input data is prepared in some convenient form for processing. The form will depend on the processing machine. For example, when electronic computers are used, the input data can be recorded on any one of the several types of input medium, such as magnetic disks, tapes, and so on.
- **Processing** In this step, the input data is changed to produce data in a more useful form. For example, pay-checks can be calculated from the time cards, or a summary of sales for the month can be calculated from the sales orders.
- **Output** At this stage, the result of the proceeding processing step is collected. The particular form of the output data depends on the use of the data. For example, output data may be pay-checks for employees.

150 9001:2015 & 14001:2015

Computer - Networking

A **computer network** is a system in which multiple computers are connected to each other to share information and resources.



Characteristics of a Computer Network

- Share resources from one computer to another.
- Create files and store them in one computer, access those files from the other computer(s) connected over the network.
- Connect a printer, scanner, or a fax machine to one computer within the network and let other computers of the network use the machines available over the network.

Following is the list of hardware's required to set up a computer network.

- Network Cables
- Distributors
- Routers
- Internal Network Cards
- External Network Cards

Network Cables

Network cables are used to connect computers. The most commonly used cable is Category

5 cable RJ-45.



Distributors

A computer can be connected to another one via a serial port but if we need to connect many computers to produce a network, this serial connection will not work.



The solution is to use a central body to which other computers, printers, scanners, etc. can be connected and then this body will manage or distribute network traffic.

Router

A router is a type of device which acts as the central point among computers and other devices that are a part of the network. It is equipped with holes called ports. Computers and other devices are connected to a router using network cables. Now-a-days router comes in wireless modes using which computers can be connected without any physical cable.



Network Card

Network card is a necessary component of a computer without which a computer cannot be connected over a network. It is also known as the network adapter or Network Interface Card (NIC). Most branded computers have network card pre-installed. Network cards are of two types: Internal and External Network Cards.

Internal Network Cards

Motherboard has a slot for internal network card where it is to be inserted. Internal network cards are of two types in which the first type uses Peripheral Component Interconnect (PCI) connection, while the second type uses Industry Standard Architecture (ISA). Network cables are required to provide network access.



External Network Cards

COPYRIGHT FIMT 2020

External network cards are of two types: Wireless and USB based. Wireless network card needs to be inserted into the motherboard, however no network cable is required to connect to the network.



Universal Serial Bus (USB)

USB card is easy to use and connects via USB port. Computers automatically detect USB card and can install the drivers required to support the USB network card automatically.





Computer - Operating System

The Operating System is a program with the following features -

- An operating system is a program that acts as an interface between the software and the computer hardware.
- It is an integrated set of specialized programs used to manage overall resources and operations of the computer.
- It is a specialized software that controls and monitors the execution of all other programs that reside in the computer, including application programs and other system software.



Objectives of Operating System

The objectives of the operating system are -

- To make the computer system convenient to use in an efficient manner.
- To hide the details of the hardware resources from the users.

- To provide users a convenient interface to use the computer system.
- To act as an intermediary between the hardware and its users, making it easier for the users to access and use other resources.
- To manage the resources of a computer system.
- To keep track of who is using which resource, granting resource requests, and mediating conflicting requests from different programs and users.
- To provide efficient and fair sharing of resources among users and programs.

Characteristics of Operating System

Here is a list of some of the most prominent characteristic features of Operating Systems -

- Memory Management Keeps track of the primary memory, i.e. what part of it is in use by whom, what part is not in use, etc. and allocates the memory when a process or program requests it.
- **Processor Management** Allocates the processor (CPU) to a process and deallocates the processor when it is no longer required.
- **Device Management** Keeps track of all the devices. This is also called I/O controller that decides which process gets the device, when, and for how much time.
- File Management Allocates and de-allocates the resources and decides who gets the resources.
- Security Prevents unauthorized access to programs and data by means of passwords and other similar techniques.
- Job Accounting Keeps track of time and resources used by various jobs and/or users.
- Control Over System Performance Records delays between the request for a service and from the system.
- Interaction with the Operators Interaction may take place via the console of the computer in the form of instructions. The Operating System acknowledges the same, does the corresponding action, and informs the operation by a display screen.
- Error-detecting Aids Production of dumps, traces, error messages, and other debugging and error-detecting methods.
- Coordination Between Other Software and Users Coordination and assignment of compilers, interpreters, assemblers, and other software to the various users of the computer systems.

Computer - Internet and Intranet

In this chapter, we will see what is Internet and Intranet, as well as discuss the similarities and differences between the two.

Internet

It is a worldwide/global system of interconnected computer networks. It uses the standard Internet Protocol (TCP/IP). Every computer in Internet is identified by a unique IP address. IP Address is a unique set of numbers (such as 110.22.33.114) which identifies a computer's location.

A special computer DNS (Domain Name Server) is used to provide a name to the IP Address so that the user can locate a computer by a name. For example, a DNS server will resolve a name https://www.tutorialspoint.com to a particular IP address to uniquely identify the computer on which this website is hosted.



Internet is accessible to every user all over the world.

Intranet

Intranet is the system in which multiple PCs are connected to each other. PCs in intranet are not available to the world outside the intranet. Usually each organization has its own Intranet network and members/employees of that organization can access the computers in their intranet.



Each computer in Intranet is also identified by an IP Address which is unique among the computers in that Intranet.

Similarities between Internet and Intranet

- Intranet uses the internet protocols such as TCP/IP and FTP.
- Intranet sites are accessible via the web browser in a similar way as websites in the internet. However, only members of Intranet network can access intranet hosted sites.
- In Intranet, own instant messengers can be used as similar to yahoo messenger/gtalk over the internet.

Differences between Internet and Intranet

- Internet is general to PCs all over the world whereas Intranet is specific to few PCs.
- Internet provides a wider and better access to websites to a large population, whereas Intranet is restricted.
- Internet is not as safe as Intranet. Intranet can be safely privatized as per the need.

How to Buy a Computer?

In this chapter, we will supply relevant information to help you buy a desktop on component by component basis. As desktops are highly customizable, so it is better to learn about the main parts and then visit the manufacturer or the retailer shop or site, instead of just looking at some specific model directly.

Popular desktop brands are Dell, Lenovo, HP and Apple. Always compare the desktops based on their specifications and base price.

Monitor



• Size – It is the diagonal size of the LCD screen. Larger the area, bigger the picture screen. A bigger picture is preferable for movie watching and gaming. It will increase the productivity as well.

- Resolution This is the number of pixels on the screen. For example, 24inch display is 1920x1200 (width by length) and 22-inch display is 1680x1050. High resolution provides better picture quality and a nice gaming experience.
- **Inputs** Now-a-days monitors can accept inputs from cable as well apart from the computer. They can also have USB ports.
- Stand Some monitors come with adjustable stands while some may not.
- **Recommended** 24 Inch LCD.

Operating System

- Operating System is the main software of the computer as everything will run on it in one form or other.
- There are primarily three choices: Windows, Linux, Apple OS X.
- Linux is free, however people generally do not use it for home purpose.
- Apple OS X works only on Apple desktops.
- Windows 7 is very popular among desktop users.
- Most of the computers come pre-equipped with Windows 7 Starter edition.
- Windows 8 is recently introduced and is available in the market.
- Windows 7 and Windows 8 come in multiple versions from starter, home basic, home premium, professional, ultimate, and enterprise editions.
- As the edition version increases, their features list and price increases.
- Recommended Windows 7 Home Premium.

Optical Drive (CD/DVD/Blu-ray)



- Optical drive is the drive on a computer, which is responsible for using CD, DVD, and Blu-ray discs.
- Now-a-days, DVD burners are industry standards.
- DVD Burner can burn CD, DVD and play them.
- DVD Burner is cheaper than Blu-ray drives.
- Blu-ray drives can play HD movies but are costlier component.
- **Recommended** DVD Burner.

Memory



- RAM is considered as Computer Memory as the performance of a computer is directly proportional to its memory and processor.
- Today's software and operating system require high memory.
- Today commonly used RAM is DDR3, which operates at 1066Mhz.
- As per Windows 7, 1 GB is the minimum RAM required to function properly.
- **Recommended** 4 GB.

Hard Drive



- Hard disk is used for storage purpose. Higher the capacity, more data you can save in it.
- Now-a-days computers are equipped with 500GB hard drive, which can be extended to 2TB.
- Most hard drives in desktop operate at the standard performance speed of 7200RPM.
- Recommended 500GB

CPU



• Frequency (GHz) – This determines the speed of the processor. More the speed, better the CPU.

- **Cores** Now-a-days CPUs come with more than one core, which is like having more than one CPU in the computer. Programs which can take advantage of multi-core environment will run faster on such machines.
- Brand Intel or AMD. Both are equivalent. Intel is in lead.
- Cache Higher the L1, L2 cache, better the CPU performance.
- Recommended Intel Core i3-3225 3.30 GHz Processor.

OperatingSystem-Overview

An Operating System (OS) is an interface between a computer user and computer hardware. An operating system is a software which performs all the basic tasks like file management, memory management, process management, handling input and output, and controlling peripheral devices such as disk drives and printers.

Some popular Operating Systems include Linux, Windows, OS X, VMS, OS/400, AIX, z/OS, etc.

Definition

An operating system is a program that acts as an interface between the user and the computer hardware and controls the execution of all kinds of programs.



015 & 1400

Following are some of important functions of an operating System.

- Memory Management
- Processor Management
- Device Management
- File Management
- Security
- Control over system performance
- Job accounting
- Error detecting aids

• Coordination between other software and users

Memory Management

Memory management refers to management of Primary Memory or Main Memory. Main memory is a large array of words or bytes where each word or byte has its own address.

Main memory provides a fast storage that can be accessed directly by the CPU. For a program to be executed, it must in the main memory. An Operating System does the following activities for memory management –

- Keeps tracks of primary memory, i.e., what part of it are in use by whom, what part are not in use.
- In multiprogramming, the OS decides which process will get memory when and how much.
- Allocates the memory when a process requests it to do so.
- De-allocates the memory when a process no longer needs it or has been terminated.

Processor Management

In multiprogramming environment, the OS decides which process gets the processor when and for how much time. This function is called **process scheduling**. An Operating System does the following activities for processor management –

- Keeps tracks of processor and status of process. The program responsible for this task is known as **traffic controller**.
- Allocates the processor (CPU) to a process.
- De-allocates processor when a process is no longer required.

Device Management

An Operating System manages device communication via their respective drivers. It does the following activities for device management –

- Keeps tracks of all devices. Program responsible for this task is known as the I/O controller.
- Decides which process gets the device when and for how much time.
- Allocates the device in the efficient way.
- De-allocates devices.

File Management

A file system is normally organized into directories for easy navigation and usage. These directories may contain files and other directions.

An Operating System does the following activities for file management -

- Keeps track of information, location, uses, status etc. The collective facilities are often known as **file system**.
- Decides who gets the resources.
- Allocates the resources.
- De-allocates the resources.

Other Important Activities

Following are some of the important activities that an Operating System performs -

- Security By means of password and similar other techniques, it prevents unauthorized access to programs and data.
- Control over system performance Recording delays between request for a service and response from the system.
- Job accounting Keeping track of time and resources used by various jobs and users.
- Error detecting aids Production of dumps, traces, error messages, and other debugging and error detecting aids.
- Coordination between other softwares and users Coordination and assignment of compilers, interpreters, assemblers and other software to the various users of the computer systems.

Types of Operating System

Operating systems are there from the very first computer generation and they keep evolving with time. In this chapter, we will discuss some of the important types of operating systems which are most commonly used.

Batch operating system

The users of a batch operating system do not interact with the computer directly. Each user prepares his job on an off-line device like punch cards and submits it to the computer operator. To speed up processing, jobs with similar needs are batched together and run as a group. The programmers leave their programs with the operator and the operator then sorts the programs with similar requirements into batches.

The problems with Batch Systems are as follows -

- Lack of interaction between the user and the job.
- CPU is often idle, because the speed of the mechanical I/O devices is slower than the CPU.

• Difficult to provide the desired priority.

Time-sharing operating systems

Time-sharing is a technique which enables many people, located at various terminals, to use a particular computer system at the same time. Time-sharing or multitasking is a logical extension of multiprogramming. Processor's time which is shared among multiple users simultaneously is termed as time-sharing.

The main difference between Multiprogrammed Batch Systems and Time-Sharing Systems is that in case of Multiprogrammed batch systems, the objective is to maximize processor use, whereas in Time-Sharing Systems, the objective is to minimize response time.

Multiple jobs are executed by the CPU by switching between them, but the switches occur so frequently. Thus, the user can receive an immediate response. For example, in a transaction processing, the processor executes each user program in a short burst or quantum of computation. That is, if **n** users are present, then each user can get a time quantum. When the user submits the command, the response time is in few seconds at most.

The operating system uses CPU scheduling and multiprogramming to provide each user with a small portion of a time. Computer systems that were designed primarily as batch systems have been modified to time-sharing systems.

Advantages of Timesharing operating systems are as follows -

- Provides the advantage of quick response.
- Avoids duplication of software.
- Reduces CPU idle time.

Disadvantages of Time-sharing operating systems are as follows -

- Problem of reliability.
- Question of security and integrity of user programs and data.
- Problem of data communication.

Distributed operating System

Distributed systems use multiple central processors to serve multiple real-time applications and multiple users. Data processing jobs are distributed among the processors accordingly. The processors communicate with one another through various communication lines (such as high-speed buses or telephone lines). These are referred as **loosely coupled systems** or distributed systems. Processors in a distributed system may vary in size and function. These processors are referred as sites, nodes, computers, and so on.

14001-201

The advantages of distributed systems are as follows -

- With resource sharing facility, a user at one site may be able to use the resources available at another.
- Speedup the exchange of data with one another via electronic mail.
- If one site fails in a distributed system, the remaining sites can potentially continue operating.
- Better service to the customers.
- Reduction of the load on the host computer.
- Reduction of delays in data processing.

Network operating System

A Network Operating System runs on a server and provides the server the capability to manage data, users, groups, security, applications, and other networking functions. The primary purpose of the network operating system is to allow shared file and printer access among multiple computers in a network, typically a local area network (LAN), a private network or to other networks.

Examples of network operating systems include Microsoft Windows Server 2003, Microsoft Windows Server 2008, UNIX, Linux, Mac OS X, Novell NetWare, and BSD.

The advantages of network operating systems are as follows -

- Centralized servers are highly stable.
- Security is server managed.
- Upgrades to new technologies and hardware can be easily integrated into the system.
- Remote access to servers is possible from different locations and types of systems.

The disadvantages of network operating systems are as follows

- High cost of buying and running a server.
- Dependency on a central location for most operations.
- Regular maintenance and updates are required.

Real Time operating System

A real-time system is defined as a data processing system in which the time interval required to process and respond to inputs is so small that it controls the environment. The time taken by the system to respond to an input and display of required updated information is termed as the **response time**. So in this method, the response time is very less as compared to online processing.

Real-time systems are used when there are rigid time requirements on the operation of a processor or the flow of data and real-time systems can be used as a control device in a dedicated application. A real-time operating system must have well-defined, fixed time constraints, otherwise the system will fail. For example, Scientific experiments, medical imaging systems, industrial control systems, weapon systems, robots, air traffic control systems, etc.

There are two types of real-time operating systems. ANAGE

Hard real-time systems

Hard real-time systems guarantee that critical tasks complete on time. In hard real-time systems, secondary storage is limited or missing and the data is stored in ROM. In these systems, virtual memory is almost never found.

Soft real-time systems

Soft real-time systems are less restrictive. A critical real-time task gets priority over other tasks and retains the priority until it completes. Soft real-time systems have limited utility than hard real-time systems. For example, multimedia, virtual reality, Advanced Scientific Projects like undersea exploration and planetary rovers, etc.

Operating System - Services

An Operating System provides services to both the users and to the programs.

- It provides programs an environment to execute.
- It provides users the services to execute the programs in a convenient manner. Following are a few common services provided by an operating system

1:2015 & 14001:2015

- Program execution
- I/O operations
- File System manipulation
- Communication
- Error Detection
- **Resource Allocation**
- Protection

Program execution

Operating systems handle many kinds of activities from user programs to system programs like printer spooler, name servers, file server, etc. Each of these activities is encapsulated as a process.

A process includes the complete execution context (code to execute, data to manipulate, registers, OS resources in use). Following are the major activities of an operating system with respect to program management –

- Loads a program into memory.
- Executes the program.
- Handles program's execution.
- Provides a mechanism for process synchronization.
- Provides a mechanism for process communication.
- Provides a mechanism for deadlock handling.

I/O Operation

An I/O subsystem comprises of I/O devices and their corresponding driver software. Drivers hide the peculiarities of specific hardware devices from the users.

An Operating System manages the communication between user and device drivers.

- I/O operation means read or write operation with any file or any specific I/O device.
- Operating system provides the access to the required I/O device when required.

File system manipulation

A file represents a collection of related information. Computers can store files on the disk (secondary storage), for long-term storage purpose. Examples of storage media include magnetic tape, magnetic disk and optical disk drives like CD, DVD. Each of these media has its own properties like speed, capacity, data transfer rate and data access methods.

A file system is normally organized into directories for easy navigation and usage. These directories may contain files and other directions. Following are the major activities of an operating system with respect to file management –

- Program needs to read a file or write a file.
- The operating system gives the permission to the program for operation on file.
- Permission varies from read-only, read-write, denied and so on.
- Operating System provides an interface to the user to create/delete files.

- Operating System provides an interface to the user to create/delete directories.
- Operating System provides an interface to create the backup of file system.

Communication

In case of distributed systems which are a collection of processors that do not share memory, peripheral devices, or a clock, the operating system manages communications between all the processes. Multiple processes communicate with one another through communication lines in the network.

The OS handles routing and connection strategies, and the problems of contention and security. Following are the major activities of an operating system with respect to communication -

- Two processes often require data to be transferred between them
- Both the processes can be on one computer or on different computers, but are connected through a computer network.
- Communication may be implemented by two methods, either by Shared Memory or by Message Passing.

Error handling

Errors can occur anytime and anywhere. An error may occur in CPU, in I/O devices or in the memory hardware. Following are the major activities of an operating system with respect to error handling –

- The OS constantly checks for possible errors.
- The OS takes an appropriate action to ensure correct and consistent computing.

Resource Management

In case of multi-user or multi-tasking environment, resources such as main memory, CPU cycles and files storage are to be allocated to each user or job. Following are the major activities of an operating system with respect to resource management –

- The OS manages all kinds of resources using schedulers.
- CPU scheduling algorithms are used for better utilization of CPU.

Protection

Considering a computer system having multiple users and concurrent execution of multiple processes, the various processes must be protected from each other's activities.

Protection refers to a mechanism or a way to control the access of programs, processes, or users to the resources defined by a computer system. Following are the major activities of an operating system with respect to protection -

- The OS ensures that all access to system resources is controlled.
- The OS ensures that external I/O devices are protected from invalid access attempts.
- The OS provides authentication features for each user by means of passwords.

Operating System - Properties

Batch processing

Batch processing is a technique in which an Operating System collects the programs and data together in a batch before processing starts. An operating system does the following activities related to batch processing –

- The OS defines a job which has predefined sequence of commands, programs and data as a single unit.
- The OS keeps a number a jobs in memory and executes them without any manual information.
- Jobs are processed in the order of submission, i.e., first come first served fashion.
- When a job completes its execution, its memory is released and the output for the job gets copied into an output spool for later printing or processing.



Advantages

- Batch processing takes much of the work of the operator to the computer.
- Increased performance as a new job get started as soon as the previous job is finished, without any manual intervention.

Disadvantages

• Difficult to debug program.

- A job could enter an infinite loop.
- Due to lack of protection scheme, one batch job can affect pending jobs.

Multitasking

Multitasking is when multiple jobs are executed by the CPU simultaneously by switching between them. Switches occur so frequently that the users may interact with each program while it is running. An OS does the following activities related to multitasking –

- The user gives instructions to the operating system or to a program directly, and receives an immediate response.
- The OS handles multitasking in the way that it can handle multiple operations/executes multiple programs at a time.
- Multitasking Operating Systems are also known as Time-sharing systems.
- These Operating Systems were developed to provide interactive use of a computer system at a reasonable cost.
- A time-shared operating system uses the concept of CPU scheduling and multiprogramming to provide each user with a small portion of a time-shared CPU.
- Each user has at least one separate program in memory.



- A program that is loaded into memory and is executing is commonly referred to as a **process**.
- When a process executes, it typically executes for only a very short time before it either finishes or needs to perform I/O.
- Since interactive I/O typically runs at slower speeds, it may take a long time to complete. During this time, a CPU can be utilized by another process.
- The operating system allows the users to share the computer simultaneously. Since each action or command in a time-shared system tends to be short, only a little CPU time is needed for each user.
- As the system switches CPU rapidly from one user/program to the next, each user is given the impression that he/she has his/her own CPU, whereas actually one CPU is being shared among many users.
Multiprogramming

Sharing the processor, when two or more programs reside in memory at the same time, is referred as **multiprogramming**. Multiprogramming assumes a single shared processor. Multiprogramming increases CPU utilization by organizing jobs so that the CPU always has one to execute.

The following figure shows the memory layout for a multiprogramming system.



An OS does the following activities related to multiprogramming.

- The operating system keeps several jobs in memory at a time.
- This set of jobs is a subset of the jobs kept in the job pool.
- The operating system picks and begins to execute one of the jobs in the memory.
- Multiprogramming operating systems monitor the state of all active programs and system resources using memory management programs to ensures that the CPU is never idle, unless there are no jobs to process.

Advantages

- High and efficient CPU utilization.
- User feels that many programs are allotted CPU almost simultaneously.

Disadvantages

- CPU scheduling is required.
- To accommodate many jobs in memory, memory management is required.

Interactivity

Interactivity refers to the ability of users to interact with a computer system. An Operating system does the following activities related to interactivity –

• Provides the user an interface to interact with the system.

- Manages input devices to take inputs from the user. For example, keyboard.
- Manages output devices to show outputs to the user. For example, Monitor.

The response time of the OS needs to be short, since the user submits and waits for the result.

Real Time System

Real-time systems are usually dedicated, embedded systems. An operating system does the following activities related to real-time system activity.

- In such systems, Operating Systems typically read from and react to sensor data.
- The Operating system must guarantee response to events within fixed periods of time to ensure correct performance.

Distributed Environment

A distributed environment refers to multiple independent CPUs or processors in a computer system. An operating system does the following activities related to distributed environment

—

- The OS distributes computation logics among several physical processors.
- The processors do not share memory or a clock. Instead, each processor has its own local memory.
- The OS manages the communications between the processors. They communicate with each other through various communication lines.

Spooling

Spooling is an acronym for simultaneous peripheral operations on line. Spooling refers to putting data of various I/O jobs in a buffer. This buffer is a special area in memory or hard disk which is accessible to I/O devices.

An operating system does the following activities related to distributed environment -

- Handles I/O device data spooling as devices have different data access rates.
- Maintains the spooling buffer which provides a waiting station where data can rest while the slower device catches up.
- Maintains parallel computation because of spooling process as a computer can perform I/O in parallel fashion. It becomes possible to have the computer read data from a tape, write data to disk and to write out to a tape printer while it is doing its computing task.



Advantages

- The spooling operation uses a disk as a very large buffer.
- Spooling is capable of overlapping I/O operation for one job with processor operations for another job

NAAC ACCREDI

Operating System - Processes

Process

A process is basically a program in execution. The execution of a process must progress in a sequential fashion. A process is defined as an entity which represents the basic unit of work to be implemented in the system. To put it in simple terms, we write our computer programs in a text file and when we execute this program, it becomes a process which performs all the tasks mentioned in the program. When a program is loaded into the memory and it becomes a process, it can be divided into four sections — stack, heap, text and data. The following image shows a simplified layout of a process inside main memory –

	Stack	
	ſ	
तेजा	1	H-F-C
11.011	Неар	3
150.9	Data	2015
	Text	NEW CONTRACT

S.N.	Component & Description
1	Stack The process Stack contains the temporary data such as method/function parameters, return address and local variables.

2	Heap This is dynamically allocated memory to a process during its run time.
3	Text This includes the current activity represented by the value of Program Counter and the contents of the processor's registers.
4	Data This section contains the global and static variables.

Program

A program is a piece of code which may be a single line or millions of lines. A computer program is usually written by a computer programmer in a programming language. For example, here is a simple program written in C programming language –

```
#include <stdio.h>
int main() {
    printf("Hello, World! \n");
    return 0;
}
```

A computer program is a collection of instructions that performs a specific task when executed by a computer. When we compare a program with a process, we can conclude that a process is a dynamic instance of a computer program. A part of a computer program that performs a well-defined task is known as an **algorithm**. A collection of computer programs, libraries and related data are referred to as a **software**.

Process Life Cycle

When a process executes, it passes through different states. These stages may differ in different operating systems, and the names of these states are also not standardized. In general, a process can have one of the following five states at a time.

S.N.	State & Description		
1	Start This is the initial state when a process is first started/created.		
2	Ready The process is waiting to be assigned to a processor. Ready processes are waiting to have the processor allocated to them by the operating system so that they can run. Process may come into this state after Start state or while running it by but interrupted by		

	the scheduler to assign CPU to some other process.
3	Running Once the process has been assigned to a processor by the OS scheduler, the process state is set to running and the processor executes its instructions.
4	Waiting Process moves into the waiting state if it needs to wait for a resource, such as waiting for user input, or waiting for a file to become available.
5	Terminated or Exit Once the process finishes its execution, or it is terminated by the operating system, it is moved to the terminated state where it waits to be removed from main memory.

Process Control Block (PCB)

A Process Control Block is a data structure maintained by the Operating System for every process. The PCB is identified by an integer process ID (PID). A PCB keeps all the information needed to keep track of a process as listed below in the table –

14.3

×C.

Wait

S.N.	Information & Description	
1 d	Process State The current state of the process i.e., whether it is ready, running, waiting, or whatever.	
2	Process privileges This is required to allow/disallow access to system resources.	
3	Process ID Unique identification for each of the process in the operating system.	
4	Pointer A pointer to parent process.	
5	Program Counter Program Counter is a pointer to the address of the next instruction	

	to be executed for this process.
6	CPU registers Various CPU registers where process need to be stored for execution for running state.
7	CPU Scheduling Information Process priority and other scheduling information which is required to schedule the process.
8	Memory management information This includes the information of page table, memory limits, Segment table depending on memory used by the operating system.
9	Accounting information This includes the amount of CPU used for process execution, time limits, execution ID etc.
10	IO status information This includes a list of I/O devices allocated to the process.

The architecture of a PCB is completely dependent on Operating System and may contain different information in different operating systems. Here is a simplified diagram of a PCB –



The PCB is maintained for a process throughout its lifetime, and is deleted once the process terminates.

Operating System - Process Scheduling

Definition

The process scheduling is the activity of the process manager that handles the removal of the running process from the CPU and the selection of another process on the basis of a particular strategy. Process scheduling is an essential part of a Multiprogramming operating

systems. Such operating systems allow more than one process to be loaded into the executable memory at a time and the loaded process shares the CPU using time multiplexing.

Process Scheduling Queues

The OS maintains all PCBs in Process Scheduling Queues. The OS maintains a separate queue for each of the process states and PCBs of all processes in the same execution state are placed in the same queue. When the state of a process is changed, its PCB is unlinked from its current queue and moved to its new state queue.

The Operating System maintains the following important process scheduling queues -

- Job queue This queue keeps all the processes in the system.
- **Ready queue** This queue keeps a set of all processes residing in main memory, ready and waiting to execute. A new process is always put in this queue.
- **Device queues** The processes which are blocked due to unavailability of an I/O device constitute this queue.



The OS can use different policies to manage each queue (FIFO, Round Robin, Priority, etc.). The OS scheduler determines how to move processes between the ready and run queues which can only have one entry per processor core on the system; in the above diagram, it has been merged with the CPU.

Two-State Process Model

Two-state process model refers to running and non-running states which are described below -

S.N.	State & Description
1	Running When a new process is created, it enters into the system as in the running state.
2	Not Running Processes that are not running are kept in queue, waiting for their turn to execute.

Each entry in the queue is a pointer to a particular process. Queue is implemented by using linked list. Use of dispatcher is as follows. When a process is interrupted, that process is transferred in the waiting queue. If the process has completed or aborted, the process is discarded. In either case, the dispatcher then selects a process from the queue to execute.

Schedulers

Schedulers are special system software which handle process scheduling in various ways. Their main task is to select the jobs to be submitted into the system and to decide which process to run. Schedulers are of three types –

- Long-Term Scheduler
- Short-Term Scheduler
- Medium-Term Scheduler

Long Term Scheduler

It is also called a **job scheduler**. A long-term scheduler determines which programs are admitted to the system for processing. It selects processes from the queue and loads them into memory for execution. Process loads into the memory for CPU scheduling.

The primary objective of the job scheduler is to provide a balanced mix of jobs, such as I/O bound and processor bound. It also controls the degree of multiprogramming. If the degree of multiprogramming is stable, then the average rate of process creation must be equal to the average departure rate of processes leaving the system.

On some systems, the long-term scheduler may not be available or minimal. Time-sharing operating systems have no long term scheduler. When a process changes the state from new to ready, then there is use of long-term scheduler.

Short Term Scheduler

It is also called as **CPU scheduler**. Its main objective is to increase system performance in accordance with the chosen set of criteria. It is the change of ready state to running state of the process. CPU scheduler selects a process among the processes that are ready to execute and allocates CPU to one of them.

Short-term schedulers, also known as dispatchers, make the decision of which process to execute next. Short-term schedulers are faster than long-term schedulers.

Medium Term Scheduler

Medium-term scheduling is a part of **swapping**. It removes the processes from the memory. It reduces the degree of multiprogramming. The medium-term scheduler is in-charge of handling the swapped out-processes.

A running process may become suspended if it makes an I/O request. A suspended processes cannot make any progress towards completion. In this condition, to remove the process from memory and make space for other processes, the suspended process is moved to the secondary storage. This process is called **swapping**, and the process is said to be swapped out or rolled out. Swapping may be necessary to improve the process mix.

ALGEMENT -

Comparison among Scheduler

S.N.	Long-Term Scheduler	Short-Term Scheduler	Medium-Term Scheduler
1	It is a job scheduler	It is a CPU scheduler	It is a process swapping scheduler.
2	Speed is lesser than short term scheduler	Speed is fastest among other two	Speed is in between both short and long term scheduler.
3	It controls the degree of multiprogramming	It provides lesser control over degree of multiprogramming	It reduces the degree of multiprogramming.
4	It is almost absent or minimal in time sharing system	It is also minimal in time sharing system	It is a part of Time sharing systems.
5	It selects processes from pool and loads them into memory for execution	It selects those processes which are ready to execute	It can re-introduce the process into memory and execution can be continued.

Context Switch

A context switch is the mechanism to store and restore the state or context of a CPU in Process Control block so that a process execution can be resumed from the same point at a later time. Using this technique, a context switcher enables multiple processes to share a single CPU. Context switching is an essential part of a multitasking operating system features.

When the scheduler switches the CPU from executing one process to execute another, the state from the current running process is stored into the process control block. After this, the

state for the process to run next is loaded from its own PCB and used to set the PC, registers, etc. At that point, the second process executing. can start CPU Process P1 Saves P1 State Process P2 Restore P2 State Process P2 Save P2 State Process P1

Context switches are computationally intensive since register and memory state must be saved and restored. To avoid the amount of context switching time, some hardware systems employ two or more sets of processor registers. When the process is switched, the following information is stored for later use.

Restore P1 State

- Program Counter
- Scheduling information
- Base and limit register value
- Currently used register
- Changed State
- I/O State information
- Accounting information

Operating System Scheduling algorithms

A Process Scheduler schedules different processes to be assigned to the CPU based on particular scheduling algorithms. There are six popular process scheduling algorithms which we are going to discuss in this chapter –

- First-Come, First-Served (FCFS) Scheduling
- Shortest-Job-Next (SJN) Scheduling
- Priority Scheduling
- Shortest Remaining Time
- Round Robin(RR) Scheduling

• Multiple-Level Queues Scheduling

These algorithms are either **non-preemptive or preemptive**. Non-preemptive algorithms are designed so that once a process enters the running state, it cannot be preempted until it completes its allotted time, whereas the preemptive scheduling is based on priority where a scheduler may preempt a low priority running process anytime when a high priority process enters into a ready state.

First Come First Serve (FCFS)

- Jobs are executed on first come, first serve basis.
- It is a non-preemptive, pre-emptive scheduling algorithm.
- Easy to understand and implement.
- Its implementation is based on FIFO queue.
- Poor in performance as average wait time is high.

Process	Arrival Time	Execute Time	Service Time
PO	0	5	0
P1	1	3	5
P2	2	8	8
P3	3	6	16



Wait time of each process is as follows -

Process	Wait Time : Service Time - Arrival Time	
PO	0 - 0 = 0	
P1	5 - 1 = 4	
P2	8 - 2 = 6	
Р3	16 - 3 = 13	

Average Wait Time: (0+4+6+13) / 4 = 5.75

Shortest Job Next (SJN)

- This is also known as **shortest job first**, or SJF
- This is a non-preemptive, pre-emptive scheduling algorithm.
- Best approach to minimize waiting time.
- Easy to implement in Batch systems where required CPU time is known in advance.

- Impossible to implement in interactive systems where required CPU time is not known.
- The processer should know in advance how much time process will take.

Process	Arrival Time	Execute Time	Service Time
PO	0	5	3
P1	1	3	0
P2	2	8	16
P3	3	6	8

р	1	PO	P3	P:	2
0	3	8	1	16	22

Wait time of each process is as follows -

Process	Wait Time : Service Time - Arrival Time
PO	3 - 0 = 3
P1	0 - 0 = 0
P2	16 - 2 = 14
Р3	8 - 3 = 5

Average Wait Time: (3+0+14+5) / 4 = 5.50

Priority Based Scheduling

- Priority scheduling is a non-preemptive algorithm and one of the most common scheduling algorithms in batch systems.
- Each process is assigned a priority. Process with highest priority is to be executed first and so on.
- Processes with same priority are executed on first come first served basis.
- Priority can be decided based on memory requirements, time requirements or any other resource requirement.

	Amvarime	Execute Time	Priority	Service Time
PO	0	5	1	9
P1	1	3	2	6
P2	2	8	1	14
P3	3	6	3	0
			1	

Wait time of each process is as follows –

Process	Wait Time : Service Time - Arrival Time
PO	9 - 0 = 9
P1	6 - 1 = 5
P2	14 - 2 = 12
Р3	0 - 0 = 0

Average Wait Time: (9+5+12+0) / 4 = 6.5

Shortest Remaining Time

- Shortest remaining time (SRT) is the preemptive version of the SJN algorithm.
- The processor is allocated to the job closest to completion but it can be preempted by a newer ready job with shorter time to completion.
- Impossible to implement in interactive systems where required CPU time is not known.
- It is often used in batch environments where short jobs need to give preference.

Round Robin Scheduling

- Round Robin is the preemptive process scheduling algorithm.
- Each process is provided a fix time to execute, it is called a quantum.
- Once a process is executed for a given time period, it is preempted and other process executes for a given time period.
- Context switching is used to save states of preempted processes.

Quantum = 3

	PO	P1	P2	P3	PO	P2	P3	P2
 0	3	6	9	12	14	17	 7 20	 0 22

Wait time of each process is as follows –

Process	Wait Time : Service Time - Arrival Time
PO	(0 - 0) + (12 - 3) = 9
P1	(3 - 1) = 2
P2	(6 - 2) + (14 - 9) + (20 - 17) = 12
Р3	(9 - 3) + (17 - 12) = 11

Average Wait Time: (9+2+12+11) / 4 = 8.5

Multiple-Level Queues Scheduling

Multiple-level queues are not an independent scheduling algorithm. They make use of other existing algorithms to group and schedule jobs with common characteristics.

- Multiple queues are maintained for processes with common characteristics.
- Each queue can have its own scheduling algorithms.
- Priorities are assigned to each queue.

For example, CPU-bound jobs can be scheduled in one queue and all I/O-bound jobs in another queue. The Process Scheduler then alternately selects jobs from each queue and assigns them to the CPU based on the algorithm assigned to the queue.

Operating System - Multi-Threading

What is Thread?

A thread is a flow of execution through the process code, with its own program counter that keeps track of which instruction to execute next, system registers which hold its current working variables, and a stack which contains the execution history.

A thread shares with its peer threads few information like code segment, data segment and open files. When one thread alters a code segment memory item, all other threads see that.

A thread is also called a **lightweight process**. Threads provide a way to improve application performance through parallelism. Threads represent a software approach to improving performance of operating system by reducing the overhead thread is equivalent to a classical process.

Each thread belongs to exactly one process and no thread can exist outside a process. Each thread represents a separate flow of control. Threads have been successfully used in implementing network servers and web server. They also provide a suitable foundation for parallel execution of applications on shared memory multiprocessors. The following figure shows the working of a single-threaded and a multithreaded process.





Single Process P with three threads

Difference between Process and Thread

S.N.	Process	Thread
1	Process is heavy weight or resource intensive.	Thread is light weight, taking lesser resources than a process.
2	Process switching needs interaction with operating system.	Thread switching does not need to interact with operating system.
3	In multiple processing environments, each process executes the same code but has its own memory and file resources.	All threads can share same set of open files, child processes.
4	If one process is blocked, then no other process can execute until the first process is unblocked.	While one thread is blocked and waiting, a second thread in the same task can run.
5	Multiple processes without using threads use more	Multiple threaded

	resources.	processes use fewer resources.
6	In multiple processes each process operates independently of the others.	One thread can read, write or change another thread's data.

- Advantages of Thread
 - Threads minimize the context switching time. •
 - Use of threads provides concurrency within a process.
 - Efficient communication.
 - It is more economical to create and context switch threads.
 - Threads allow utilization of multiprocessor architectures to a greater scale and efficiency.

Types of Thread

Threads are implemented in following two ways -

- User Level Threads User managed threads.
- Kernel Level Threads Operating System managed threads acting on kernel, an operating system core.

User Level Threads

In this case, the thread management kernel is not aware of the existence of threads. The thread library contains code for creating and destroying threads, for passing message and data between threads, for scheduling thread execution and for saving and restoring thread The contexts. application with single thread. starts a



Advantages

- Thread switching does not require Kernel mode privileges.
- User level thread can run on any operating system.
- Scheduling can be application specific in the user level thread.
- User level threads are fast to create and manage.

Disadvantages

- In a typical operating system, most system calls are blocking.
- Multithreaded application cannot take advantage of multiprocessing.

Kernel Level Threads

In this case, thread management is done by the Kernel. There is no thread management code in the application area. Kernel threads are supported directly by the operating system. Any application can be programmed to be multithreaded. All of the threads within an application are supported within a single process.

NAAC ACCREDITED

The Kernel maintains context information for the process as a whole and for individuals threads within the process. Scheduling by the Kernel is done on a thread basis. The Kernel performs thread creation, scheduling and management in Kernel space. Kernel threads are generally slower to create and manage than the user threads.

Advantages

- Kernel can simultaneously schedule multiple threads from the same process on multiple processes.
- If one thread in a process is blocked, the Kernel can schedule another thread of the same process.
- Kernel routines themselves can be multithreaded.

Disadvantages

- Kernel threads are generally slower to create and manage than the user threads.
- Transfer of control from one thread to another within the same process requires a mode switch to the Kernel.

Multithreading Models

Some operating system provide a combined user level thread and Kernel level thread facility. Solaris is a good example of this combined approach. In a combined system, multiple threads within the same application can run in parallel on multiple processors and a blocking system call need not block the entire process. Multithreading models are three types

- Many to many relationship.
- Many to one relationship.
- One to one relationship.

Many to Many Model

The many-to-many model multiplexes any number of user threads onto an equal or smaller number of kernel threads.

1462

The following diagram shows the many-to-many threading model where 6 user level threads are multiplexing with 6 kernel level threads. In this model, developers can create as many user threads as necessary and the corresponding Kernel threads can run in parallel on a multiprocessor machine. This model provides the best accuracy on concurrency and when a thread performs a blocking system call, the kernel can schedule another thread for execution.



Many to One Model

Many-to-one model maps many user level threads to one Kernel-level thread. Thread management is done in user space by the thread library. When thread makes a blocking system call, the entire process will be blocked. Only one thread can access the Kernel at a time, so multiple threads are unable to run in parallel on multiprocessors.

If the user-level thread libraries are implemented in the operating system in such a way that the system does not support them, then the Kernel threads use the many-to-one relationship modes.



One to One Model

There is one-to-one relationship of user-level thread to the kernel-level thread. This model provides more concurrency than the many-to-one model. It also allows another thread to run when a thread makes a blocking system call. It supports multiple threads to execute in parallel on microprocessors. Disadvantage of this model is that creating user thread requires the corresponding Kernel thread. OS/2, windows NT and windows 2000 use one to one relationship model.



COPYRIGHT FIMT 2020

Difference between User-Level & Kernel-Level Thread

S.N.	User-Level Threads	Kernel-Level Thread
1	User-level threads are faster to create and manage.	Kernel-level threads are slower to create and manage.
2	Implementation is by a thread library at the user level.	Operating system supports creation of Kernel threads.
3	User-level thread is generic and can run on any operating system.	Kernel-level thread is specific to the operating system.
4	Multi-threaded applications cannot take advantage of multiprocessing.	Kernel routines themselves can be multithreaded.

Operating System - Memory Management

Memory management is the functionality of an operating system which handles or manages primary memory and moves processes back and forth between main memory and disk during execution. Memory management keeps track of each and every memory location, regardless of either it is allocated to some process or it is free. It checks how much memory is to be allocated to processes. It decides which process will get memory at what time. It tracks whenever some memory gets freed or unallocated and correspondingly it updates the status.

This tutorial will teach you basic concepts related to Memory Management.

Process Address Space

The process address space is the set of logical addresses that a process references in its code. For example, when 32-bit addressing is in use, addresses can range from 0 to 0x7fffffff; that is, 2^31 possible numbers, for a total theoretical size of 2 gigabytes.

The operating system takes care of mapping the logical addresses to physical addresses at the time of memory allocation to the program. There are three types of addresses used in a program before and after memory is allocated -

S.N.	Memory Addresses & Description
1	Symbolic addresses The addresses used in a source code. The variable names, constants, and instruction labels are the basic elements of the symbolic address space.

2	Relative addresses At the time of compilation, a compiler converts symbolic addresses into relative addresses.
3	Physical addresses The loader generates these addresses at the time when a program is loaded into main memory.
X 7:	and a hand and a difference and the second in second 1. Charles and 1. I difference difference him dimension

Virtual and physical addresses are the same in compile-time and load-time address-binding schemes. Virtual and physical addresses differ in execution-time address-binding scheme.

The set of all logical addresses generated by a program is referred to as a **logical address space**. The set of all physical addresses corresponding to these logical addresses is referred to as a **physical address space**. The runtime mapping from virtual to physical address is done by the memory management unit (MMU) which is a hardware device. MMU uses following mechanism to convert virtual address to physical address.

- The value in the base register is added to every address generated by a user process, which is treated as offset at the time it is sent to memory. For example, if the base register value is 10000, then an attempt by the user to use address location 100 will be dynamically reallocated to location 10100.
- The user program deals with virtual addresses; it never sees the real physical addresses.

Static vs Dynamic Loading

The choice between Static or Dynamic Loading is to be made at the time of computer program being developed. If you have to load your program statically, then at the time of compilation, the complete programs will be compiled and linked without leaving any external program or module dependency. The linker combines the object program with other necessary object modules into an absolute program, which also includes logical addresses.

If you are writing a Dynamically loaded program, then your compiler will compile the program and for all the modules which you want to include dynamically, only references will be provided and rest of the work will be done at the time of execution.

At the time of loading, with **static loading**, the absolute program (and data) is loaded into memory in order for execution to start.

If you are using **dynamic loading**, dynamic routines of the library are stored on a disk in relocatable form and are loaded into memory only when they are needed by the program.

Static vs Dynamic Linking

As explained above, when static linking is used, the linker combines all other modules needed by a program into a single executable program to avoid any runtime dependency. When dynamic linking is used, it is not required to link the actual module or library with the program, rather a reference to the dynamic module is provided at the time of compilation and linking. Dynamic Link Libraries (DLL) in Windows and Shared Objects in Unix are good examples of dynamic libraries.

Swapping

Swapping is a mechanism in which a process can be swapped temporarily out of main memory (or move) to secondary storage (disk) and make that memory available to other processes. At some later time, the system swaps back the process from the secondary storage to main memory. Though performance is usually affected by swapping process but it helps in running multiple and big processes in parallel and that's the reason **Swapping is also known as a technique for memory compaction**.

Main Memory		Secondary Memo
	Process P1	
	Process P2	Process P3
	P1 goes for I/O wait	Process P4
	Process P1 Swap out	
	Swap in Process I	P3 Process Pn
	P1 comes back after I/O	
	Process P3 Swap out	
	Swap in Process I	21
	L	

The total time taken by swapping process includes the time it takes to move the entire process to a secondary disk and then to copy the process back to memory, as well as the time the process takes to regain main memory.

Let us assume that the user process is of size 2048KB and on a standard hard disk where swapping will take place has a data transfer rate around 1 MB per second. The actual transfer of the 1000K process to or from memory will take

```
2048KB / 1024KB per second
```

= 2 seconds

COPYRIGHT FIMT 2020

= 2000 milliseconds

Now considering in and out time, it will take complete 4000 milliseconds plus other overhead where the process competes to regain main memory.

Memory Allocation

Main memory usually has two partitions -

- Low Memory Operating system resides in this memory.
- High Memory User processes are held in high memory.

Operating system uses the following memory allocation mechanism.

S.N.	Memory Allocation & Description
1	Single-partition allocation In this type of allocation, relocation-register scheme is used to protect user processes from each other, and from changing operating-system code and data. Relocation register contains value of smallest physical address whereas limit register contains range of logical addresses. Each logical address must be less than the limit register.
2	Multiple-partition allocation In this type of allocation, main memory is divided into a number of fixed-sized partitions where each partition should contain only one process. When a partition is free, a process is selected from the input queue and is loaded into the free partition. When the process terminates, the partition becomes available for another process.
Fragme	entation
As pr little j consid Fragn Fragn	rocesses are loaded and removed from memory, the free memory space is broken into pieces. It happens after sometimes that processes cannot be allocated to memory blocks dering their small size and memory blocks remains unused. This problem is known as nentation.
S.N.	Fragmentation & Description

External fragmentation

Total memory space is enough to satisfy a request or to reside a process in it, but it is

1

 not contiguous, so it cannot be used.

 2
 Internal fragmentation

 Memory block assigned to process is bigger. Some portion of memory is left unused, as it cannot be used by another process.

The following diagram shows how fragmentation can cause waste of memory and a compaction technique can be used to create more free memory out of fragmented memory –

Fragmented memory before compact	ion

Memory after compac	tion		

External fragmentation can be reduced by compaction or shuffle memory contents to place all free memory together in one large block. To make compaction feasible, relocation should be dynamic.

The internal fragmentation can be reduced by effectively assigning the smallest partition but large enough for the process.

Paging

A computer can address more memory than the amount physically installed on the system. This extra memory is actually called virtual memory and it is a section of a hard that's set up to emulate the computer's RAM. Paging technique plays an important role in implementing virtual memory.

Paging is a memory management technique in which process address space is broken into blocks of the same size called **pages** (size is power of 2, between 512 bytes and 8192 bytes). The size of the process is measured in the number of pages.

Similarly, main memory is divided into small fixed-sized blocks of (physical) memory called **frames** and the size of a frame is kept the same as that of a page to have optimum utilization of the main memory and to avoid external fragmentation.

		Main Memory		Secondary Memor
		Operating System	`	
Process P		Process P – Page 4	FO	
First 100 bytes	Page 0	л.	F1	
Second 100 bytes	Page 1	Process P - Page 0	F2	
Third 100 bytes	Page 2	Process P – Page 2	F3	
Fourth 100 bytes	Page 3	Process P – Page 1	F4	
Fifth 100 bytes	Page 4	Brocott P - Page 7	ES	
Sixth 100 bytes	Page 5	Process P - Page /		
Seventh 100 bytes	Page 6	Process P - Page N		
Eight 100 bytes	Page 7			
and so on		Descer for other processes		
	Page N	Pages for other processes		
		Pages for other processes		
		Pages for other processes	FN	

Address Translation

Page address is called logical address and represented by page number and the offset.

Logical Address = Page number + page offset

Frame address is called **physical address** and represented by a **frame number** and the **offset**.

Physical Address = Frame number + page offset

A data structure called **page map table** is used to keep track of the relation between a page of a process to a frame in physical memory.



When the system allocates a frame to any page, it translates this logical address into a physical address and create entry into the page table to be used throughout execution of the program.

When a process is to be executed, its corresponding pages are loaded into any available memory frames. Suppose you have a program of 8Kb but your memory can accommodate only 5Kb at a given point in time, then the paging concept will come into picture. When a computer runs out of RAM, the operating system (OS) will move idle or unwanted pages of memory to secondary memory to free up RAM for other processes and brings them back when needed by the program.

This process continues during the whole execution of the program where the OS keeps removing idle pages from the main memory and write them onto the secondary memory and bring them back when required by the program.

Advantages and Disadvantages of Paging

Here is a list of advantages and disadvantages of paging -

- Paging reduces external fragmentation, but still suffer from internal fragmentation.
- Paging is simple to implement and assumed as an efficient memory management technique.
- Due to equal size of the pages and frames, swapping becomes very easy.
- Page table requires extra memory space, so may not be good for a system having small RAM.

Segmentation

Segmentation is a memory management technique in which each job is divided into several segments of different sizes, one for each module that contains pieces that perform related functions. Each segment is actually a different logical address space of the program.

When a process is to be executed, its corresponding segmentation are loaded into noncontiguous memory though every segment is loaded into a contiguous block of available memory.

Segmentation memory management works very similar to paging but here segments are of variable-length where as in paging pages are of fixed size.

A program segment contains the program's main function, utility functions, data structures, and so on. The operating system maintains a **segment map table** for every process and a list of free memory blocks along with segment numbers, their size and corresponding memory locations in main memory. For each segment, the table stores the starting address of the segment and the length of the segment. A reference to a memory location includes a value that identifies a segment and an offset.

segment and an onset.

	SN	Size	Memory Address	Operating System
Segment 1	1	400	100	
	2	200	500	100
	з	100	800	200
	N	×	NM	300
Segment 2	L	1		400
				500
				600
				700
Segment 3				800
Segment N				

Operating System - Virtual Memory

A computer can address more memory than the amount physically installed on the system. This extra memory is actually called **virtual memory** and it is a section of a hard disk that's set up to emulate the computer's RAM.

AGE 67

The main visible advantage of this scheme is that programs can be larger than physical memory. Virtual memory serves two purposes. First, it allows us to extend the use of physical memory by using disk. Second, it allows us to have memory protection, because each virtual address is translated to a physical address.

Following are the situations, when entire program is not required to be loaded fully in main memory.

- User written error handling routines are used only when an error occurred in the data or computation.
- Certain options and features of a program may be used rarely.
- Many tables are assigned a fixed amount of address space even though only a small amount of the table is actually used.
- The ability to execute a program that is only partially in memory would counter many benefits.
- Less number of I/O would be needed to load or swap each user program into memory.
- A program would no longer be constrained by the amount of physical memory that is available.
- Each user program could take less physical memory, more programs could be run the same time, with a corresponding increase in CPU utilization and throughput.

Modern microprocessors intended for general-purpose use, a memory management unit, or MMU, is built into the hardware. The MMU's job is to translate virtual addresses into physical addresses. A basic example is given below –



Virtual memory is commonly implemented by demand paging. It can also be implemented in a segmentation system. Demand segmentation can also be used to provide virtual memory.

-

Demand Paging

A demand paging system is quite similar to a paging system with swapping where processes reside in secondary memory and pages are loaded only on demand, not in advance. When a context switch occurs, the operating system does not copy any of the old program's pages out to the disk or any of the new program's pages into the main memory Instead, it just begins executing the new program after loading the first page and fetches that program's pages as they are referenced.



While executing a program, if the program references a page which is not available in the main memory because it was swapped out a little ago, the processor treats this invalid memory reference as a **page fault** and transfers control from the program to the operating system to demand the page back into the memory.

Advantages

Following are the advantages of Demand Paging -

- Large virtual memory.
- More efficient use of memory.
- There is no limit on degree of multiprogramming.

Disadvantages

• Number of tables and the amount of processor overhead for handling page interrupts are greater than in the case of the simple paged management techniques.

Page Replacement Algorithm

Page replacement algorithms are the techniques using which an Operating System decides which memory pages to swap out, write to disk when a page of memory needs to be allocated. Paging happens whenever a page fault occurs and a free page cannot be used for allocation purpose accounting to reason that pages are not available or the number of free pages is lower than required pages.

When the page that was selected for replacement and was paged out, is referenced again, it has to read in from disk, and this requires for I/O completion. This process determines the quality of the page replacement algorithm: the lesser the time waiting for page-ins, the better is the algorithm.

A page replacement algorithm looks at the limited information about accessing the pages provided by hardware, and tries to select which pages should be replaced to minimize the total number of page misses, while balancing it with the costs of primary storage and processor time of the algorithm itself. There are many different page replacement algorithms. We evaluate an algorithm by running it on a particular string of memory reference and computing the number of page faults,

Reference String

The string of memory references is called reference string. Reference strings are generated artificially or by tracing a given system and recording the address of each memory reference. The latter choice produces a large number of data, where we note two things.

- For a given page size, we need to consider only the page number, not the entire address.
- If we have a reference to a page **p**, then any immediately following references to page **p** will never cause a page fault. Page p will be in memory after the first reference; the immediately following references will not fault.
- For example, consider the following sequence of addresses 123,215,600,1234,76,96
- If page size is 100, then the reference string is 1,2,6,12,0,0

First In First Out (FIFO) algorithm

Missos

- Oldest page in main memory is the one which will be selected for replacement.
- Easy to implement, keep a list, replace pages from the tail and add new pages at the head.



Optimal Page algorithm

- An optimal page-replacement algorithm has the lowest page-fault rate of all algorithms. An optimal page-replacement algorithm exists, and has been called OPT or MIN.
- Replace the page that will not be used for the longest period of time. Use the time when a page is to be used.



Fault Rate = 6 / 12 = 0.50

Least Recently Used (LRU) algorithm

- Page which has not been used for the longest time in main memory is the one which will be selected for replacement.
- Easy to implement, keep a list, replace pages by looking back into time.





Page Buffering algorithm

- To get a process start quickly, keep a pool of free frames.
- On page fault, select a page to be replaced.
- Write the new page in the frame of free pool, mark the page table and restart the process.
- Now write the dirty page out of disk and place the frame holding replaced page in free pool.

Least frequently Used(LFU) algorithm

- The page with the smallest count is the one which will be selected for replacement.
- This algorithm suffers from the situation in which a page is used heavily during the initial phase of a process, but then is never used again.

Most frequently Used(MFU) algorithm

This algorithm is based on the argument that the page with the smallest count was • probably just brought in and has yet to be used. 4001-2015

6.....

Operating System - I/O Hardware

One of the important jobs of an Operating System is to manage various I/O devices including mouse, keyboards, touch pad, disk drives, display adapters, USB devices, Bitmapped screen, LED, Analog-to-digital converter, On/off switch, network connections, audio I/O, printers etc.

An I/O system is required to take an application I/O request and send it to the physical device, then take whatever response comes back from the device and send it to the application. I/O devices can be divided into two categories –

- Block devices A block device is one with which the driver communicates by sending entire blocks of data. For example, Hard disks, USB cameras, Disk-On-Key etc.
- Character devices A character device is one with which the driver communicates by sending and receiving single characters (bytes, octets). For example, serial ports, parallel ports, sounds cards etc

ALGEMEN,

Device Controllers

Device drivers are software modules that can be plugged into an OS to handle a particular device. Operating System takes help from device drivers to handle all I/O devices. The Device Controller works like an interface between a device and a device driver. I/O units (Keyboard, mouse, printer, etc.) typically consist of a mechanical component and an electronic component where electronic component is called the device controller. There is always a device controller and a device driver for each device to communicate with the Operating Systems. A device controller may be able to handle multiple devices. As an interface its main task is to convert serial bit stream to block of bytes, perform error correction as necessary.

Any device connected to the computer is connected by a plug and socket, and the socket is connected to a device controller. Following is a model for connecting the CPU, memory, controllers, and I/O devices where CPU and device controllers all use a common bus for communication.



Synchronous vs asynchronous I/O

- Synchronous I/O In this scheme CPU execution waits while I/O proceeds
- Asynchronous I/O I/O proceeds concurrently with CPU execution

Communication to I/O Devices

The CPU must have a way to pass information to and from an I/O device. There are three approaches available to communicate with the CPU and Device.

- Special Instruction I/O
- Memory-mapped I/O
- Direct memory access (DMA)

Special Instruction I/O

This uses CPU instructions that are specifically made for controlling I/O devices. These instructions typically allow data to be sent to an I/O device or read from an I/O device.

Memory-mapped I/O

When using memory-mapped I/O, the same address space is shared by memory and I/O devices. The device is connected directly to certain main memory locations so that I/O device can transfer block of data to/from memory without going through CPU.



While using memory mapped IO, OS allocates buffer in memory and informs I/O device to use that buffer to send data to the CPU. I/O device operates asynchronously with CPU, interrupts CPU when finished.

The advantage to this method is that every instruction which can access memory can be used to manipulate an I/O device. Memory mapped IO is used for most high-speed I/O devices like disks, communication interfaces.

Direct Memory Access (DMA)

Slow devices like keyboards will generate an interrupt to the main CPU after each byte is transferred. If a fast device such as a disk generated an interrupt for each byte, the operating system would spend most of its time handling these interrupts. So a typical computer uses direct memory access (DMA) hardware to reduce this overhead.

Direct Memory Access (DMA) means CPU grants I/O module authority to read from or write to memory without involvement. DMA module itself controls exchange of data between main memory and the I/O device. CPU is only involved at the beginning and end of the transfer and interrupted only after entire block has been transferred.

Direct Memory Access needs a special hardware called DMA controller (DMAC) that manages the data transfers and arbitrates access to the system bus. The controllers are programmed with source and destination pointers (where to read/write the data), counters to track the number of transferred bytes, and settings, which includes I/O and memory types, interrupts and states for the CPU cycles.



The operating system uses the DMA hardware as follows -

Step	Description
1	Device driver is instructed to transfer disk data to a buffer address X.
2	Device driver then instruct disk controller to transfer data to buffer.
3	Disk controller starts DMA transfer.
4	Disk controller sends each byte to DMA controller.
5	DMA controller transfers bytes to buffer, increases the memory address, decreases the counter C until C becomes zero.
6	When C becomes zero, DMA interrupts CPU to signal transfer completion.

Polling vs Interrupts I/O

A computer must have a way of detecting the arrival of any type of input. There are two ways that this can happen, known as **polling** and **interrupts**. Both of these techniques allow the processor to deal with events that can happen at any time and that are not related to the process it is currently running.

Polling I/O

Polling is the simplest way for an I/O device to communicate with the processor. The process of periodically checking status of the device to see if it is time for the next I/O operation, is called polling. The I/O device simply puts the information in a Status register, and the processor must come and get the information.

Most of the time, devices will not require attention and when one does it will have to wait until it is next interrogated by the polling program. This is an inefficient method and much of the processors time is wasted on unnecessary polls.

Compare this method to a teacher continually asking every student in a class, one after another, if they need help. Obviously the more efficient method would be for a student to inform the teacher whenever they require assistance.

Interrupts I/O

An alternative scheme for dealing with I/O is the interrupt-driven method. An interrupt is a signal to the microprocessor from a device that requires attention.

A device controller puts an interrupt signal on the bus when it needs CPU's attention when CPU receives an interrupt, It saves its current state and invokes the appropriate interrupt handler using the interrupt vector (addresses of OS routines to handle various events). When the interrupting device has been dealt with, the CPU continues with its original task as if it had never been interrupted.

Operating System - I/O Software's

I/O software is often organized in the following layers -

- User Level Libraries This provides simple interface to the user program to perform input and output. For example, stdio is a library provided by C and C++ programming languages.
- **Kernel Level Modules** This provides device driver to interact with the device controller and device independent I/O modules used by the device drivers.
- **Hardware** This layer includes actual hardware and hardware controller which interact with the device drivers and makes hardware alive.

A key concept in the design of I/O software is that it should be device independent where it should be possible to write programs that can access any I/O device without having to specify the device in advance. For example, a program that reads a file as input should be able to read a file on a floppy disk, on a hard disk, or on a CD-ROM, without having to modify the program for each different device.

User -		User I/O Libraries	
Kernel		Device Independent I/O	
L	Device Driver	Device Driver	Device Driver
lardware -	Device Controller	Device Controller	Device Controller
	USB Drive	Disk	Printer

ARGEMEN

Device Drivers

Device drivers are software modules that can be plugged into an OS to handle a particular device. Operating System takes help from device drivers to handle all I/O devices. Device drivers encapsulate device-dependent code and implement a standard interface in such a way that code contains device-specific register reads/writes. Device driver, is generally written by the device's manufacturer and delivered along with the device on a CD-ROM.

A device driver performs the following jobs –

- To accept request from the device independent software above to it.
- Interact with the device controller to take and give I/O and perform required error handling
- Making sure that the request is executed successfully

How a device driver handles a request is as follows: Suppose a request comes to read a block N. If the driver is idle at the time a request arrives, it starts carrying out the request immediately. Otherwise, if the driver is already busy with some other request, it places the new request in the queue of pending requests.

Interrupt handlers

An interrupt handler, also known as an interrupt service routine or ISR, is a piece of software or more specifically a callback function in an operating system or more specifically in a device driver, whose execution is triggered by the reception of an interrupt.

When the interrupt happens, the interrupt procedure does whatever it has to in order to handle the interrupt, updates data structures and wakes up process that was waiting for an interrupt to happen.

The interrupt mechanism accepts an address - a number that selects a specific interrupt handling routine/function from a small set. In most architectures, this address is an offset
stored in a table called the interrupt vector table. This vector contains the memory addresses of specialized interrupt handlers.

Device-Independent I/O Software

The basic function of the device-independent software is to perform the I/O functions that are common to all devices and to provide a uniform interface to the user-level software. Though it is difficult to write completely device independent software but we can write some modules which are common among all the devices. Following is a list of functions of device-independent I/O Software –

- Uniform interfacing for device drivers
- Device naming Mnemonic names mapped to Major and Minor device numbers
- Device protection
- Providing a device-independent block size
- Buffering because data coming off a device cannot be stored in final destination.
- Storage allocation on block devices
- Allocation and releasing dedicated devices
- Error Reporting

User-Space I/O Software

These are the libraries which provide richer and simplified interface to access the functionality of the kernel or ultimately interactive with the device drivers. Most of the user-level I/O software consists of library procedures with some exception like spooling system which is a way of dealing with dedicated I/O devices in a multiprogramming system.

I/O Libraries (e.g., stdio) are in user-space to provide an interface to the OS resident deviceindependent I/O SW. For example putchar(), getchar(), printf() and scanf() are example of user level I/O library stdio available in C programming.

Kernel I/O Subsystem

Kernel I/O Subsystem is responsible to provide many services related to I/O. Following are some of the services provided.

• Scheduling – Kernel schedules a set of I/O requests to determine a good order in which to execute them. When an application issues a blocking I/O system call, the request is placed on the queue for that device. The Kernel I/O scheduler rearranges the order of the queue to improve the overall system efficiency and the average response time experienced by the applications.

- **Buffering** Kernel I/O Subsystem maintains a memory area known as **buffer** that stores data while they are transferred between two devices or between a device with an application operation. Buffering is done to cope with a speed mismatch between the producer and consumer of a data stream or to adapt between devices that have different data transfer sizes.
- Caching Kernel maintains cache memory which is region of fast memory that holds copies of data. Access to the cached copy is more efficient than access to the original.
- Spooling and Device Reservation A spool is a buffer that holds output for a device, such as a printer, that cannot accept interleaved data streams. The spooling system copies the queued spool files to the printer one at a time. In some operating systems, spooling is managed by a system daemon process. In other operating systems, it is handled by an in kernel thread.
- Error Handling An operating system that uses protected memory can guard against many kinds of hardware and application errors.

Operating System - File System

File

A file is a named collection of related information that is recorded on secondary storage such as magnetic disks, magnetic tapes and optical disks. In general, a file is a sequence of bits, bytes, lines or records whose meaning is defined by the files creator and user.

File Structure

A File Structure should be according to a required format that the operating system can understand.

- A file has a certain defined structure according to its type.
- A text file is a sequence of characters organized into lines.
- A source file is a sequence of procedures and functions.
- An object file is a sequence of bytes organized into blocks that are understandable by the machine.
- When operating system defines different file structures, it also contains the code to support these file structure. Unix, MS-DOS support minimum number of file structure.

File Type

File type refers to the ability of the operating system to distinguish different types of file such as text files source files and binary files etc. Many operating systems support many types of files. Operating system like MS-DOS and UNIX have the following types of files –

Ordinary files

- These are the files that contain user information.
- These may have text, databases or executable program.
- The user can apply various operations on such files like add, modify, delete or even remove the entire file.

Directory files

• These files contain list of file names and other information related to these files.

Special files

- These files are also known as device files.
- These files represent physical device like disks, terminals, printers, networks, tape drive etc.

These files are of two types –

- Character special files data is handled character by character as in case of terminals or printers.
- Block special files data is handled in blocks as in the case of disks and tapes.

File Access Mechanisms

File access mechanism refers to the manner in which the records of a file may be accessed. There are several ways to access files –

- Sequential access
- Direct/Random access
- Indexed sequential access

Sequential access

A sequential access is that in which the records are accessed in some sequence, i.e., the information in the file is processed in order, one record after the other. This access method is the most primitive one. Example: Compilers usually access files in this fashion.

Direct/Random access

- Random access file organization provides, accessing the records directly.
- Each record has its own address on the file with by the help of which it can be directly accessed for reading or writing.
- The records need not be in any sequence within the file and they need not be in adjacent locations on the storage medium.

Indexed sequential access

- This mechanism is built up on base of sequential access.
- An index is created for each file which contains pointers to various blocks.
- Index is searched sequentially and its pointer is used to access the file directly.

Space Allocation

Files are allocated disk spaces by operating system. Operating systems deploy following three main ways to allocate disk space to files.

- Contiguous Allocation
- Linked Allocation
- Indexed Allocation

Contiguous Allocation

- Each file occupies a contiguous address space on disk.
- Assigned disk address is in linear order.
- Easy to implement.
- External fragmentation is a major issue with this type of allocation technique.

Linked Allocation

• Each file carries a list of links to disk blocks.

- Directory contains link / pointer to first block of a file.
- No external fragmentation
- Effectively used in sequential access file.
- Inefficient in case of direct access file.

Indexed Allocation

- Provides solutions to problems of contiguous and linked allocation.
- A index block is created having all pointers to files.
- Each file has its own index block which stores the addresses of disk space occupied by the file.
- Directory contains the addresses of index blocks of files.

Operating System - Security

Security refers to providing a protection system to computer system resources such as CPU, memory, disk, software programs and most importantly data/information stored in the computer system. If a computer program is run by an unauthorized user, then he/she may cause severe damage to computer or data stored in it. So a computer system must be protected against unauthorized access, malicious access to system memory, viruses, worms etc. We're going to discuss following topics in this chapter.

- Authentication
- One Time passwords
- Program Threats
- System Threats
- Computer Security Classifications

Authentication

Authentication refers to identifying each user of the system and associating the executing programs with those users. It is the responsibility of the Operating System to create a protection system which ensures that a user who is running a particular program is authentic. Operating Systems generally identifies/authenticates users using following three ways –

• Username / Password – User need to enter a registered username and password with Operating system to login into the system.

- User card/key User need to punch card in card slot, or enter key generated by key generator in option provided by operating system to login into the system.
- User attribute fingerprint/ eye retina pattern/ signature User need to pass his/her attribute via designated input device used by operating system to login into the system.

One Time passwords

One-time passwords provide additional security along with normal authentication. In One-Time Password system, a unique password is required every time user tries to login into the system. Once a one-time password is used, then it cannot be used again. One-time password are implemented in various ways.

- Random numbers Users are provided cards having numbers printed along with corresponding alphabets. System asks for numbers corresponding to few alphabets randomly chosen.
- Secret key User are provided a hardware device which can create a secret id mapped with user id. System asks for such secret id which is to be generated every time prior to login.
- Network password Some commercial applications send one-time passwords to user on registered mobile/ email which is required to be entered prior to login.

Program Threats

Operating system's processes and kernel do the designated task as instructed. If a user program made these process do malicious tasks, then it is known as **Program Threats**. One of the common example of program threat is a program installed in a computer which can store and send user credentials via network to some hacker. Following is the list of some well-known program threats.

- **Trojan Horse** Such program traps user login credentials and stores them to send to malicious user who can later on login to computer and can access system resources.
- **Trap Door** If a program which is designed to work as required, have a security hole in its code and perform illegal action without knowledge of user then it is called to have a trap door.
- Logic Bomb Logic bomb is a situation when a program misbehaves only when certain conditions met otherwise it works as a genuine program. It is harder to detect.

• Virus – Virus as name suggest can replicate themselves on computer system. They are highly dangerous and can modify/delete user files, crash systems. A virus is generatly a small code embedded in a program. As user accesses the program, the virus starts getting embedded in other files/ programs and can make system unusable for user

System Threats

System threats refers to misuse of system services and network connections to put user in trouble. System threats can be used to launch program threats on a complete network called as program attack. System threats creates such an environment that operating system resources/ user files are misused. Following is the list of some well-known system threats.

- Worm Worm is a process which can choked down a system performance by using system resources to extreme levels. A Worm process generates its multiple copies where each copy uses system resources, prevents all other processes to get required resources. Worms processes can even shut down an entire network.
- **Port Scanning** Port scanning is a mechanism or means by which a hacker can detects system vulnerabilities to make an attack on the system.
- **Denial of Service** Denial of service attacks normally prevents user to make legitimate use of the system. For example, a user may not be able to use internet if denial of service attacks browser's content settings.

Computer Security Classifications

As per the U.S. Department of Defense Trusted Computer System's Evaluation Criteria there are four security classifications in computer systems: A, B, C, and D. This is widely used specifications to determine and model the security of systems and of security solutions. Following is the brief description of each classification.

S.N.	Classification Type & Description
1	Type A Highest Level. Uses formal design specifications and verification techniques. Grants a high degree of assurance of process security.
2	Type B Provides mandatory protection system. Have all the properties of a class C2 system. Attaches a sensitivity label to each object. It is of three types.

	 B1 – Maintains the security label of each object in the system. Label is used for making decisions to access control. B2 – Extends the sensitivity labels to each system resource, such as storage objects, supports covert channels and auditing of events. B3 – Allows creating lists or user groups for access-control to grant access or revoke access to a given named object.
3	 Type C Provides protection and user accountability using audit capabilities. It is of two types. C1 – Incorporates controls so that users can protect their private information and keep other users from accidentally reading / deleting their data. UNIX versions are mostly Cl class. C2 – Adds an individual-level access control to the capabilities of a Cl level system.
4	Type D Lowest level. Minimum protection. MS-DOS, Window 3.1 fall in this category.

Operating System - Linux

Linux is one of popular version of UNIX operating System. It is open source as its source code is freely available. It is free to use. Linux was designed considering UNIX compatibility. Its functionality list is quite similar to that of UNIX.

Components of Linux System

Linux Operating System has primarily three components

- Kernel Kernel is the core part of Linux. It is responsible for all major activities of this operating system. It consists of various modules and it interacts directly with the underlying hardware. Kernel provides the required abstraction to hide low level hardware details to system or application programs.
- System Library System libraries are special functions or programs using which application programs or system utilities accesses Kernel's features. These libraries implement most of the functionalities of the operating system and do not requires kernel module's code access rights.

System Utility – System Utility programs are responsible to do specialized, individual level tasks.



Kernel Mode vs User Mode

A GELOTA Kernel component code executes in a special privileged mode called kernel mode with full access to all resources of the computer. This code represents a single process, executes in single address space and do not require any context switch and hence is very efficient and fast. Kernel runs each processes and provides system services to processes, provides protected access to hardware to processes.

Support code which is not required to run in kernel mode is in System Library. User programs and other system programs works in User Mode which has no access to system hardware and kernel code. User programs/ utilities use System libraries to access Kernel functions to get system's low level tasks.

Basic Features

Following are some of the important features of Linux Operating System.

- Portable Portability means software can works on different types of hardware in same way. Linux kernel and application programs supports their installation on any kind of hardware platform.
- Open Source Linux source code is freely available and it is community based • development project. Multiple teams work in collaboration to enhance the capability of Linux operating system and it is continuously evolving.
- Multi-User Linux is a multiuser system means multiple users can access system resources like memory/ ram/ application programs at same time.
- Multiprogramming Linux is a multiprogramming system means multiple applications can run at same time.
- **Hierarchical File System** Linux provides a standard file structure in which system files/ user files are arranged.

- Shell Linux provides a special interpreter program which can be used to execute commands of the operating system. It can be used to do various types of operations, call application programs. etc.
- Security Linux provides user security using authentication features like password protection/ controlled access to specific files/ encryption of data.

Architecture

The following illustration shows the architecture of a Linux system -



The architecture of a Linux System consists of the following layers -

- Hardware layer Hardware consists of all peripheral devices (RAM/ HDD/ CPU etc).
- **Kernel** It is the core component of Operating System, interacts directly with hardware, provides low level services to upper layer components.
- Shell An interface to kernel, hiding complexity of kernel's functions from users.
 The shell takes commands from the user and executes kernel's functions.
- Utilities Utility programs that provide the user most of the functionalities of an operating systems.

DataCommunication & Computer Network

DCN-Overview

A system of interconnected computers and computerized peripherals such as printers is called computer network. This interconnection among computers facilitates information sharing among them. Computers may connect to each other by either wired or wireless media.

Classification of Computer Networks

Computer networks are classified based on various factors. They includes:

- Geographical span
- Inter-connectivity
- Administration
- Architecture

Geographical Span

Geographically a network can be seen in one of the following categories:

- It may be spanned across your table, among Bluetooth enabled devices,. Ranging not more than few meters.
- It may be spanned across a whole building, including intermediate devices to connect all floors.
- It may be spanned across a whole city.
- It may be spanned across multiple cities or provinces.
- It may be one network covering whole world.

Inter-Connectivity

Components of a network can be connected to each other differently in some fashion. By connectedness we mean either logically, physically, or both ways.

- Every single device can be connected to every other device on network, making the network mesh.
- All devices can be connected to a single medium but geographically disconnected, created bus like structure.
- Each device is connected to its left and right peers only, creating linear structure.
- All devices connected together with a single device, creating star like structure.
- All devices connected arbitrarily using all previous ways to connect each other, resulting in a hybrid structure.

Administration

From an administrator's point of view, a network can be private network which belongs a single autonomous system and cannot be accessed outside its physical or logical domain.A network can be public which is accessed by all.

Network Architecture

Computer networks can be discriminated into various types such as Client-Server, peer-to-peer or hybrid, depending upon its architecture.

- There can be one or more systems acting as Server. Other being Client, requests the Server to serve requests. Server takes and processes request on behalf of Clients.
- Two systems can be connected Point-to-Point, or in back-to-back fashion. They both reside at the same level and called peers.
- There can be hybrid network which involves network architecture of both the above types.

Network Applications

Computer systems and peripherals are connected to form a network. They provide numerous advantages:

- Resource sharing such as printers and storage devices
- Exchange of information by means of e-Mails and FTP
- Information sharing by using Web or Internet
- Interaction with other users using dynamic web pages
- IP phones
- Video conferences
- Parallel computing
- Instant messaging

DCN-ComputerNetworkTypes

Generally, networks are distinguished based on their geographical span. A network can be as small as distance between your mobile phone and its Bluetooth headphone and as large as the internet itself, covering the whole geographical world,

Personal Area Network

A Personal Area Network (PAN) is smallest network which is very personal to a user. This may include Bluetooth enabled devices or infra-red enabled devices. PAN has connectivity range up to 10 meters. PAN may include wireless computer keyboard and mouse, Bluetooth enabled headphones, wireless printers and TV remotes.



For example, Piconet is Bluetooth-enabled Personal Area Network which may contain up to 8 devices connected together in a master-slave fashion.

Local Area Network

A computer network spanned inside a building and operated under single administrative system is generally termed as Local Area Network (LAN). Usually,LAN covers an organization' offices, schools, colleges or universities. Number of systems connected in LAN may vary from as least as two to as much as 16 million.

LAN provides a useful way of sharing the resources between end users. The resources such as printers, file servers, scanners, and internet are easily sharable among computers.



LANs are composed of inexpensive networking and routing equipment. It may contains local servers serving file storage and other locally shared applications. It mostly operates on private IP addresses and does not involve heavy routing. LAN works under its own local domain and controlled centrally.

LAN uses either Ethernet or Token-ring technology. Ethernet is most widely employed LAN technology and uses Star topology, while Token-ring is rarely seen.

LAN can be wired, wireless, or in both forms at once.

Metropolitan Area Network

The Metropolitan Area Network (MAN) generally expands throughout a city such as cable TV network. It can be in the form of Ethernet, Token-ring, ATM, or Fiber Distributed Data Interface (FDDI).

Metro Ethernet is a service which is provided by ISPs. This service enables its users to expand their Local Area Networks. For example, MAN can help an organization to connect all of its offices in a city.



Backbone of MAN is high-capacity and high-speed fiber optics. MAN works in between Local Area Network and Wide Area Network. MAN provides uplink for LANs to WANs or internet.

Wide AreaNetwork

As the name suggests, the Wide Area Network (WAN) covers a wide area which may span across provinces and even a whole country. Generally, telecommunication networks are Wide Area Network. These networks provide connectivity to MANs and LANs. Since they are equipped with very high speed backbone, WANs use very expensive network equipment.



WAN may use advanced technologies such as Asynchronous Transfer Mode (ATM), Frame Relay, and Synchronous Optical Network (SONET). WAN may be managed by multiple administration.

Internetwork

A network of networks is called an internetwork, or simply the internet. It is the largest network in existence on this planet. The internet hugely connects all WANs and it can have connection to LANs and Home networks. Internet uses TCP/IP protocol suite and uses IP as

its addressing protocol. Present day, Internet is widely implemented using IPv4. Because of shortage of address spaces, it is gradually migrating from IPv4 to IPv6.

Internet enables its users to share and access enormous amount of information worldwide. It uses WWW, FTP, email services, audio and video streaming etc. At huge level, internet works on Client-Server model.

Internet uses very high speed backbone of fiber optics. To inter-connect various continents, fibers are laid under sea known to us as submarine communication cable.

Internet is widely deployed on World Wide Web services using HTML linked pages and is accessible by client software known as Web Browsers. When a user requests a page using some web browser located on some Web Server anywhere in the world, the Web Server responds with the proper HTML page. The communication delay is very low.

Internet is serving many proposes and is involved in many aspects of life. Some of them are:

- Web sites
- E-mail
- Instant Messaging
- Blogging
- Social Media
- Marketing
- Networking
- Resource Sharing
- Audio and Video Streaming

DCN-NetworkLANTechnologies

Let us go through various LAN technologies in brief:

Ethemet

Ethernet is a widely deployed LAN technology. This technology was invented by Bob Metcalfe and D.R. Boggs in the year 1970. It was standardized in IEEE 802.3 in 1980. Ethernet shares media. Network which uses shared media has high probability of data collision. Ethernet uses Carrier Sense Multi Access/Collision Detection (CSMA/CD) technology to detect collisions. On the occurrence of collision in Ethernet, all its hosts roll back, wait for some random amount of time, and then re-transmit the data.

Ethernet connector is, network interface card equipped with 48-bits MAC address. This helps other Ethernet devices to identify and communicate with remote devices in Ethernet.

Traditional Ethernet uses 10BASE-T specifications. The number 10 depicts 10MBPS speed, BASE stands for baseband, and T stands for Thick Ethernet. 10BASE-T Ethernet provides transmission speed up to 10MBPS and uses coaxial cable or Cat-5 twisted pair cable with RJ-45 connector. Ethernet follows star topology with segment length up to 100 meters. All devices are connected to a hub/switch in a star fashion.

Fast-Ethemet

To encompass need of fast emerging software and hardware technologies, Ethernet extends itself as Fast-Ethernet. It can run on UTP, Optical Fiber, and wirelessly too. It can provide speed up to 100 MBPS. This standard is named as 100BASE-T in IEEE 803.2 using Cat-5 twisted pair cable. It uses CSMA/CD technique for wired media sharing among the Ethernet hosts and CSMA/CA (CA stands for Collision Avoidance) technique for wireless Ethernet LAN.

Fast Ethernet on fiber is defined under 100BASE-FX standard which provides speed up to 100 MBPS on fiber. Ethernet over fiber can be extended up to 100 meters in half-duplex mode and can reach maximum of 2000 meters in full-duplex over multimode fibers.

Giga-Ethemet

After being introduced in 1995, Fast-Ethernet could enjoy its high speed status only for 3 years till Giga-Ethernet introduced. Giga-Ethernet provides speed up to 1000 mbits/seconds. IEEE802.3ab standardize Giga-Ethernet over UTP using Cat-5, Cat-5e and Cat-6 cables. IEEE802.3ah defines Giga-Ethernet over Fiber.

VirtualLAN

LAN uses Ethernet which in turn works on shared media. Shared media in Ethernet create one single Broadcast domain and one single Collision domain. Introduction of switches to Ethernet has removed single collision domain issue and each device connected to switch works in its separate collision domain. But even Switches cannot divide a network into separate Broadcast domains.

Virtual LAN is a solution to divide a single Broadcast domain into multiple Broadcast domains. Host in one VLAN cannot speak to a host in another. By default, all hosts are placed into the same VLAN.



In this diagram, different VLANs are depicted in different color codes. Hosts in one VLAN, even if connected on the same Switch cannot see or speak to other hosts in different VLANs. VLAN is Layer-2 technology which works closely on Ethernet. To route packets between two different VLANs a Layer-3 device such as Router is required.

DCN-Computer Network Topologies

A Network Topology is the arrangement with which computer systems or network devices are connected to each other. Topologies may define both physical and logical aspect of the network. Both logical and physical topologies could be same or different in a same network.

Point-to-Point

Point-to-point networks contains exactly two hosts such as computer, switches or routers, servers connected back to back using a single piece of cable. Often, the receiving end of one host is connected to sending end of the other and vice-versa.



If the hosts are connected point-to-point logically, then may have multiple intermediate devices. But the end hosts are unaware of underlying network and see each other as if they are connected directly.

BusTopology

In case of Bus topology, all devices share single communication line or cable. Bus topology may have problem while multiple hosts sending data at the same time. Therefore, Bus topology either uses CSMA/CD technology or recognizes one host as Bus Master to solve the issue. It is one of the simple forms of networking where a failure of a device does not

affect the other devices. But failure of the shared communication line can make all other devices stop functioning.



Both ends of the shared channel have line terminator. The data is sent in only one direction and as soon as it reaches the extreme end, the terminator removes the data from the line. ARGEME

StarTopobgy

All hosts in Star topology are connected to a central device, known as hub device, using a point-to-point connection. That is, there exists a point to point connection between hosts and hub. The hub device can be any of the following:

- Layer-1 device such as hub or repeater
- Layer-2 device such as switch or bridge
- Layer-3 device such as router or gateway



As in Bus topology, hub acts as single point of failure. If hub fails, connectivity of all hosts to all other hosts fails. Every communication between hosts, takes place through only the hub.Star topology is not expensive as to connect one more host, only one cable is required and configuration is simple. 015 & 14001:2

Ring Topology

In ring topology, each host machine connects to exactly two other machines, creating a circular network structure. When one host tries to communicate or send message to a host which is not adjacent to it, the data travels through all intermediate hosts. To connect one more host in the existing structure, the administrator may need only one more extra cable.



Failure of any host results in failure of the whole ring. Thus, every connection in the ring is a point of failure. There are methods which employ one more backup ring.

MeshTopology

ABGEMEA In this type of topology, a host is connected to one or multiple hosts. This topology has hosts in point-to-point connection with every other host or may also have hosts which are in pointto-point connection to few hosts only.



Hosts in Mesh topology also work as relay for other hosts which do not have direct point-topoint links. Mesh technology comes into two types:

- Full Mesh: All hosts have a point-to-point connection to every other host in the network. Thus for every new host n(n-1)/2 connections are required. It provides the most reliable network structure among all network topologies.
- **Partially Mesh**: Not all hosts have point-to-point connection to every other host. • Hosts connect to each other in some arbitrarily fashion. This topology exists where we need to provide reliability to some hosts out of all.

TreeTopobgy

Also known as Hierarchical Topology, this is the most common form of network topology in use presently. This topology imitates as extended Star topology and inherits properties of bus topology.

This topology divides the network in to multiple levels/layers of network. Mainly in LANs, a network is bifurcated into three types of network devices. The lowermost is access-layer where computers are attached. The middle layer is known as distribution layer, which works as mediator between upper layer and lower layer. The highest layer is known as core layer, and is central point of the network, i.e. root of the tree from which all nodes fork.



All neighboring hosts have point-to-point connection between them.Similar to the Bus topology, if the root goes down, then the entire network suffers even.though it is not the single point of failure. Every connection serves as point of failure, failing of which divides the network into unreachable segment.

DaisyChain

This topology connects all the hosts in a linear fashion. Similar to Ring topology, all hosts are connected to two hosts only, except the end hosts.Means, if the end hosts in daisy chain are connected then it represents Ring topology.



Each link in daisy chain topology represents single point of failure. Every link failure splits the network into two segments. Every intermediate host works as relay for its immediate hosts.

Hybrid Topology

A network structure whose design contains more than one topology is said to be hybrid topology. Hybrid topology inherits merits and demerits of all the incorporating topologies.



The above picture represents an arbitrarily hybrid topology. The combining topologies may contain attributes of Star, Ring, Bus, and Daisy-chain topologies. Most WANs are connected by means of Dual-Ring topology and networks connected to them are mostly Star topology networks. Internet is the best example of largest Hybrid topology

DCN-ComputerNetwork Models

Networking engineering is a complicated task, which involves software, firmware, chip level engineering, hardware, and electric pulses. To ease network engineering, the whole networking concept is divided into multiple layers. Each layer is involved in some particular task and is independent of all other layers. But as a whole, almost all networking tasks depend on all of these layers. Layers share data between them and they depend on each other only to take input and send output.

Layered Tasks

In layered architecture of Network Model, one whole network process is divided into small tasks. Each small task is then assigned to a particular layer which works dedicatedly to process the task only. Every layer does only specific work.

In layered communication system, one layer of a host deals with the task done by or to be done by its peer layer at the same level on the remote host. The task is either initiated by layer at the lowest level or at the top most level. If the task is initiated by the-top most layer, it is passed on to the layer below it for further processing. The lower layer does the same thing, it processes the task and passes on to lower layer. If the task is initiated by lower most layer, then the reverse path is taken.



COPYRIGHT FIMT 2020

Every layer clubs together all procedures, protocols, and methods which it requires to execute its piece of task. All layers identify their counterparts by means of encapsulation header and tail.

OSIModel

Open System Interconnect is an open standard for all communication systems. OSI model is established by International Standard Organization (ISO). This model has seven layers:



- Application Layer: This layer is responsible for providing interface to the application user. This layer encompasses protocols which directly interact with the user.
- **Presentation Layer**: This layer defines how data in the native format of remote host should be presented in the native format of host.
- Session Layer: This layer maintains sessions between remote hosts. For example, once user/password authentication is done, the remote host maintains this session for a while and does not ask for authentication again in that time span.
- Transport Layer: This layer is responsible for end-to-end delivery between hosts.
- Network Layer: This layer is responsible for address assignment and uniquely addressing hosts in a network.
- **Data Link Layer**: This layer is responsible for reading and writing data from and onto the line. Link errors are detected at this layer.
- **Physical Layer**: This layer defines the hardware, cabling wiring, power output, pulse rate etc.

Internet Model

Internet uses TCP/IP protocol suite, also known as Internet suite. This defines Internet Model which contains four layered architecture. OSI Model is general communication model but Internet Model is what the internet uses for all its communication. The internet is independent of its underlying network architecture so is its Model. This model has the following layers:



- **Application Layer**: This layer defines the protocol which enables user to interact with the network.For example, FTP, HTTP etc.
- **Transport Layer**: This layer defines how data should flow between hosts. Major protocol at this layer is Transmission Control Protocol (TCP). This layer ensures data delivered between hosts is in-order and is responsible for end-to-end delivery.
- **Internet Layer**: Internet Protocol (IP) works on this layer. This layer facilitates host addressing and recognition. This layer defines routing.
- Link Layer: This layer provides mechanism of sending and receiving actual data.Unlike its OSI Model counterpart, this layer is independent of underlying network architecture and hardware.

DCN-ComputerNetworkSecurity

During initial days of internet, its use was limited to military and universities for research and development purpose. Later when all networks merged together and formed internet, the data useds to travel through public transit network.Common people may send the data that can be highly sensitive such as their bank credentials, username and passwords, personal documents, online shopping details, or confidential documents.

All security threats are intentional i.e. they occur only if intentionally triggered. Security threats can be divided into the following categories:

• Interruption

Interruption is a security threat in which availability of resources is attacked. For example, a user is unable to access its web-server or the web-server is hijacked.

• Privacy-Breach

In this threat, the privacy of a user is compromised. Someone, who is not the authorized person is accessing or intercepting data sent or received by the original authenticated user.

• Integrity

This type of threat includes any alteration or modification in the original context of communication. The attacker intercepts and receives the data sent by the sender and the attacker then either modifies or generates false data and sends to the receiver. The receiver receives the data assuming that it is being sent by the original Sender.

• Authenticity

This threat occurs when an attacker or a security violator, poses as a genuine person and accesses the resources or communicates with other genuine users.

No technique in the present world can provide 100% security. But steps can be taken to secure data while it travels in unsecured network or internet. The most widely used technique is Cryptography.



Cryptography is a technique to encrypt the plain-text data which makes it difficult to understand and interpret. There are several cryptographic algorithms available present day as described below:

- Secret Key
- Public Key
- Message Digest

Searct KeyEncryption

Both sender and receiver have one secret key. This secret key is used to encrypt the data at sender's end. After the data is encrypted, it is sent on the public domain to the receiver. Because the receiver knows and has the Secret Key, the encrypted data packets can easily be decrypted.

Example of secret key encryption is Data Encryption Standard (DES). In Secret Key encryption, it is required to have a separate key for each host on the network making it difficult to manage.

Public Key Encryption

In this encryption system, every user has its own Secret Key and it is not in the shared domain. The secret key is never revealed on public domain. Along with secret key, every

user has its own but public key. Public key is always made public and is used by Senders to encrypt the data. When the user receives the encrypted data, he can easily decrypt it by using its own Secret Key.

Example of public key encryption is Rivest-Shamir-Adleman (RSA).

Message Digest

In this method, actual data is not sent, instead a hash value is calculated and sent. The other end user, computes its own hash value and compares with the one just received. If both hash values are matched, then it is accepted otherwise rejected.

Example of Message Digest is MD5 hashing. It is mostly used in authentication where user password is cross checked with the one saved on the server.

DCN-Physical Layer Introduction

Physical layer in the OSI model plays the role of interacting with actual hardware and signaling mechanism. Physical layer is the only layer of OSI network model which actually deals with the physical connectivity of two different stations. This layer defines the hardware equipment, cabling, wiring, frequencies, pulses used to represent binary signals etc.

Physical layer provides its services to Data-link layer. Data-link layer hands over frames to physical layer. Physical layer converts them to electrical pulses, which represent binary data. The binary data is then sent over the wired or wireless media.

Signals

When data is sent over physical medium, it needs to be first converted into electromagnetic signals. Data itself can be analog such as human voice, or digital such as file on the disk.Both analog and digital data can be represented in digital or analog signals.

• Digital Signals

Digital signals are discrete in nature and represent sequence of voltage pulses. Digital signals are used within the circuitry of a computer system.

Analog Signals

Analog signals are in continuous wave form in nature and represented by continuous electromagnetic waves.

Transmission Impairment

When signals travel through the medium they tend to deteriorate. This may have many reasons as given:

• Attenuation

For the receiver to interpret the data accurately, the signal must be sufficiently strong. When the signal passes through the medium, it tends to get weaker. As it covers distance, it loses strength.

• Dispersion

As signal travels through the media, it tends to spread and overlaps. The amount of dispersion depends upon the frequency used.

Delay distortion

Signals are sent over media with pre-defined speed and frequency. If the signal speed and frequency do not match, there are possibilities that signal reaches destination in arbitrary fashion. In digital media, this is very critical that some bits reach earlier than the previously sent ones.

• Noise

Random disturbance or fluctuation in analog or digital signal is said to be Noise in signal, which may distort the actual information being carried. Noise can be characterized in one of the following class:

Thermal Noise

Heat agitates the electronic conductors of a medium which may introduce noise in the media. Up to a certain level, thermal noise is unavoidable.

• Intermodulation

When multiple frequencies share a medium, their interference can cause noise in the medium. Intermodulation noise occurs if two different frequencies are sharing a medium and one of them has excessive strength or the component itself is not functioning properly, then the resultant frequency may not be delivered as expected.

Crosstalk

This sort of noise happens when a foreign signal enters into the media. This is because signal in one medium affects the signal of second medium.

• Impulse

This noise is introduced because of irregular disturbances such as lightening, electricity, short-circuit, or faulty components. Digital data is mostly affected by this sort of noise.

Transmission Media

The media over which the information between two computer systems is sent, called transmission media. Transmission media comes in two forms.

Guided Media

All communication wires/cables are guided media, such as UTP, coaxial cables, and fiber Optics. In this media, the sender and receiver are directly connected and the information is send (guided) through it.

• Unguided Media

Wireless or open air space is said to be unguided media, because there is no connectivity between the sender and receiver. Information is spread over the air, and anyone including the actual recipient may collect the information.

Channel Capacity

The speed of transmission of information is said to be the channel capacity. We count it as data rate in digital world. It depends on numerous factors such as:

- Bandwidth: The physical limitation of underlying media.
- Error-rate: Incorrect reception of information because of noise.
- Encoding: The number of levels used for signaling.

Multiplexing

Multiplexing is a technique to mix and send multiple data streams over a single medium. This technique requires system hardware called multiplexer (MUX) for multiplexing the streams and sending them on a medium, and de-multiplexer (DMUX) which takes information from the medium and distributes to different destinations.

Switching

Switching is a mechanism by which data/information sent from source towards destination which are not directly connected. Networks have interconnecting devices, which receives data from directly connected sources, stores data, analyze it and then forwards to the next interconnecting device closest to the destination.



DCN-Digital Transmission

Data or information can be stored in two ways, analog and digital. For a computer to use the data, it must be in discrete digital form.Similar to data, signals can also be in analog and digital form. To transmit data digitally, it needs to be first converted to digital form.

Digital-to-Digital Conversion

This section explains how to convert digital data into digital signals. It can be done in two ways, line coding and block coding. For all communications, line coding is necessary whereas block coding is optional.

LineCoding

The process for converting digital data into digital signal is said to be Line Coding. Digital data is found in binary format. It is represented (stored) internally as series of 1s and 0s.



Digital signal is denoted by discreet signal, which represents digital data. There are three types of line coding schemes available:



Uni-polar Encoding

Unipolar encoding schemes use single voltage level to represent data. In this case, to represent binary 1, high voltage is transmitted and to represent 0, no voltage is transmitted. It is also called Unipolar-Non-return-to-zero, because there is no rest condition i.e. it either represents 1 or 0.



COPYRIGHT FIMT 2020

Polar Encoding

Polar encoding scheme uses multiple voltage levels to represent binary values. Polar encodings is available in four types:

- Polar Non-Return to Zero (Polar NRZ)
 - It uses two different voltage levels to represent binary values. Generally, positive voltage represents 1 and negative value represents 0. It is also NRZ because there is no rest condition.

NRZ scheme has two variants: NRZ-L and NRZ-I.



NRZ-L changes voltage level at when a different bit is encountered whereas NRZ-I changes voltage when a 1 is encountered.

• Return to Zero (RZ)

Problem with NRZ is that the receiver cannot conclude when a bit ended and when the next bit is started, in case when sender and receiver's clock are not synchronized.



RZ uses three voltage levels, positive voltage to represent 1, negative voltage to represent 0 and zero voltage for none. Signals change during bits not between bits.

• Manchester

This encoding scheme is a combination of RZ and NRZ-L. Bit time is divided into two halves. It transits in the middle of the bit and changes phase when a different bit is encountered.

Differential Manchester

This encoding scheme is a combination of RZ and NRZ-I. It also transit at the middle of the bit but changes phase only when 1 is encountered.

Bipolar Encoding

Bipolar encoding uses three voltage levels, positive, negative and zero. Zero voltage represents binary 0 and bit 1 is represented by altering positive and negative voltages.



BlockCoding

To ensure accuracy of the received data frame redundant bits are used. For example, in even-parity, one parity bit is added to make the count of 1s in the frame even. This way the original number of bits is increased. It is called Block Coding.

Block coding is represented by slash notation, mB/nB.Means, m-bit block is substituted with n-bit block where n > m. Block coding involves three steps:

- Division,
- Substitution
- Combination.

After block coding is done, it is line coded for transmission.

Analog-to-Digital Conversion

Microphones create analog voice and camera creates analog videos, which are treated is analog data. To transmit this analog data over digital signals, we need analog to digital conversion.

Analog data is a continuous stream of data in the wave form whereas digital data is discrete.

To convert analog wave into digital data, we use Pulse Code Modulation (PCM).

PCM is one of the most commonly used method to convert analog data into digital form. It involves three steps:

- Sampling
- Quantization
- Encoding.

Sampling



The analog signal is sampled every T interval. Most important factor in sampling is the rate at which analog signal is sampled. According to Nyquist Theorem, the sampling rate must be at least two times of the highest frequency of the signal.

Quantization



Sampling yields discrete form of continuous analog signal. Every discrete pattern shows the amplitude of the analog signal at that instance. The quantization is done between the maximum amplitude value and the minimum amplitude value. Quantization is approximation of the instantaneous analog value.

Encoding



In encoding, each approximated value is then converted into binary format.

Transmission Modes

The transmission mode decides how data is transmitted between two computers. The binary data in the form of 1s and 0s can be sent in two different modes: Parallel and Serial. Parallel Transmission

	•		
		1	
0		0	(
		1	
		o_ /	

The binary bits are organized in-to groups of fixed length. Both sender and receiver are connected in parallel with the equal number of data lines. Both computers distinguish between high order and low order data lines. The sender sends all the bits at once on all lines.Because the data lines are equal to the number of bits in a group or data frame, a complete group of bits (data frame) is sent in one go. Advantage of Parallel transmission is high speed and disadvantage is the cost of wires, as it is equal to the number of bits sent in parallel.

Serial Transmission

In serial transmission, bits are sent one after another in a queue manner. Serial transmission requires only one communication channel.



Serial transmission can be either asynchronous or synchronous.

Asynchronous Serial Transmission

It is named so because there'is no importance of timing. Data-bits have specific pattern and they help receiver recognize the start and end data bits.For example, a 0 is prefixed on every data byte and one or more 1s are added at the end.

Two continuous data-frames (bytes) may have a gap between them.

Synchronous Serial Transmission

Timing in synchronous transmission has importance as there is no mechanism followed to recognize start and end data bits. There is no pattern or prefix/suffix method. Data bits are sent in burst mode without maintaining gap between bytes (8-bits). Single burst of data bits may contain a number of bytes. Therefore, timing becomes very important.

It is up to the receiver to recognize and separate bits into bytes. The advantage of synchronous transmission is high speed, and it has no overhead of extra header and footer bits as in asynchronous transmission.

DCN-Analog Transmission

To send the digital data over an analog media, it needs to be converted into analog signal. There can be two cases according to data formatting.

Bandpass:The filters are used to filter and pass frequencies of interest. A bandpass is a band of frequencies which can pass the filter.

Low-pass: Low-pass is a filter that passes low frequencies signals.

When digital data is converted into a bandpass analog signal, it is called digital-to-analog conversion. When low-pass analog signal is converted into bandpass analog signal, it is called analog-to-analog conversion.

Digital-to-Analog Conversion

When data from one computer is sent to another via some analog carrier, it is first converted into analog signals. Analog signals are modified to reflect digital data.

An analog signal is characterized by its amplitude, frequency, and phase. There are three kinds of digital-to-analog conversions:

Amplitude Shift Keying

In this conversion technique, the amplitude of analog carrier signal is modified to reflect binary data.



When binary data represents digit 1, the amplitude is held; otherwise it is set to 0. Both frequency and phase remain same as in the original carrier signal.

• Frequency Shift Keying

In this conversion technique, the frequency of the analog carrier signal is modified to reflect binary data.



This technique uses two frequencies, f1 and f2. One of them, for example f1, is chosen to represent binary digit 1 and the other one is used to represent binary digit 0. Both amplitude and phase of the carrier wave are kept intact.

Phase Shift Keying

In this conversion scheme, the phase of the original carrier signal is altered to reflect the binary data.



AAGEMEN

When a new binary symbol is encountered, the phase of the signal is altered. Amplitude and frequency of the original carrier signal is kept intact.

Quadrature Phase Shift Keying

QPSK alters the phase to reflect two binary digits at once. This is done in two different phases. The main stream of binary data is divided equally into two substreams. The serial data is converted in to parallel in both sub-streams and then each stream is converted to digital signal using NRZ technique. Later, both the digital signals are merged together.

Analog-to-Analog Conversion

Analog signals are modified to represent analog data. This conversion is also known as Analog Modulation. Analog modulation is required when bandpass is used. Analog to analog conversion can be done in three ways:



Amplitude Modulation

In this modulation, the amplitude of the carrier signal is modified to reflect the analog data.



Amplitude modulation is implemented by means of a multiplier. The amplitude of modulating signal (analog data) is multiplied by the amplitude of carrier frequency, which then reflects analog data.

The frequency and phase of carrier signal remain unchanged.

• Frequency Modulation

In this modulation technique, the frequency of the carrier signal is modified to reflect the change in the voltage levels of the modulating signal (analog data).



The amplitude and phase of the carrier signal are not altered.

Phase Modulation

In the modulation technique, the phase of carrier signal is modulated in order to reflect the change in voltage (amplitude) of analog data signal.



Phase modulation is practically similar to Frequency Modulation, but in Phase modulation frequency of the carrier signal is not increased. Frequency of carrier is signal is changed (made dense and sparse) to reflect voltage change in the amplitude of modulating signal.

DCN-Transmission Media

The transmission media is nothing but the physical media over which communication takes place in computer networks.

Magnetic Media

One of the most convenient way to transfer data from one computer to another, even before the birth of networking, was to save it on some storage media and transfer physical from one station to another. Though it may seem old-fashion way in today's world of high speed internet, but when the size of data is huge, the magnetic media comes into play.

For example, a bank has to handle and transfer huge data of its customer, which stores a backup of it at some geographically far-away place for security reasons and to keep it from uncertain calamities. If the bank needs to store its huge backup data then its,transfer through internet is not feasible. The WAN links may not support such high speed. Even if they do; the cost too high to afford.

In these cases, data backup is stored onto magnetic tapes or magnetic discs, and then shifted physically at remote places.

TwistedPairCable

A twisted pair cable is made of two plastic insulated copper wires twisted together to form a single media. Out of these two wires, only one carries actual signal and another is used for ground reference. The twists between wires are helpful in reducing noise (electro-magnetic interference) and crosstalk.
There are two types of twisted pair cables:

- Shielded Twisted Pair (STP) Cable
- Unshielded Twisted Pair (UTP) Cable

STP cables comes with twisted wire pair covered in metal foil. This makes it more indifferent to noise and crosstalk.

UTP has seven categories, each suitable for specific use. In computer networks, Cat-5, Cat-5e, and Cat-6 cables are mostly used. UTP cables are connected by RJ45 connectors.

Coaxial Cable

Coaxial cable has two wires of copper. The core wire lies in the center and it is made of solid conductor. The core is enclosed in an insulating sheath. The second wire is wrapped around over the sheath and that too in turn encased by insulator sheath. This all is covered by plastic cover.



Because of its structure, the coax cable is capable of carrying high frequency signals than that of twisted pair cable. The wrapped structure provides it a good shield against noise and cross talk. Coaxial cables provide high bandwidth rates of up to 450 mbps.

There are three categories of coax cables namely, RG-59 (Cable TV), RG-58 (Thin Ethernet), and RG-11 (Thick Ethernet). RG stands for Radio Government.

Cables are connected using BNC connector and BNC-T. BNC terminator is used to terminate the wire at the far ends.

PowerLines

Power Line communication (PLC) is Layer-1 (Physical Layer) technology which uses power cables to transmit data signals. In PLC, modulated data is sent over the cables. The receiver on the other end de-modulates and interprets the data.

Because power lines are widely deployed, PLC can make all powered devices controlled and monitored. PLC works in half-duplex.

There are two types of PLC:

- Narrow band PLC AC ACCREDITED
- Broad band PLC

Narrow band PLC provides lower data rates up to 100s of kbps, as they work at lower frequencies (3-5000 kHz). They can be spread over several kilometers.

Broadband PLC provides higher data rates up to 100s of Mbps and works at higher frequencies (1.8 – 250 MHz). They cannot be as much extended as Narrowband PLC.

Fiber Optics

Fiber Optic works on the properties of light. When light ray hits at critical angle it tends to refracts at 90 degree. This property has been used in fiber optic. The core of fiber optic cable is made of high quality glass or plastic. From one end of it light is emitted, it travels through it and at the other end light detector detects light stream and converts it to electric data.

Fiber Optic provides the highest mode of speed. It comes in two modes, one is single mode fiber and second is multimode fiber. Single mode fiber can carry a single ray of light whereas multimode is capable of carrying multiple beams of light.



Fiber Optic also comes in unidirectional and bidirectional capabilities. To connect and access fiber optic special type of connectors are used. These can be Subscriber Channel (SC), Straight Tip (ST), or MT-RJ.

DCN-Wireless Transmission

Wireless transmission is a form of unguided media. Wireless communication involves no physical link established between two or more devices, communicating wirelessly. Wireless signals are spread over in the air and are received and interpreted by appropriate antennas.

When an antenna is attached to electrical circuit of a computer or wireless device, it converts the digital data into wireless signals and spread all over within its frequency range. The receptor on the other end receives these signals and converts them back to digital data. A little part of electromagnetic spectrum can be used for wireless transmission.



Radio Transmission

Radio frequency is easier to generate and because of its large wavelength it can penetrate through walls and structures alike. Radio waves can have wavelength from 1 mm - 100,000 km and have frequency ranging from 3 Hz (Extremely Low Frequency) to 300 GHz (Extremely High Frequency). Radio frequencies are sub-divided into six bands.

Radio waves at lower frequencies can travel through walls whereas higher RF can travel in straight line and bounce back. The power of low frequency waves decreases sharply as they cover long distance. High frequency radio waves have more power.

Lower frequencies such as VLF, LF, MF bands can travel on the ground up to 1000 kilometers, over the earth's surface.



Radio waves of high frequencies are prone to be absorbed by rain and other obstacles. They use Ionosphere of earth atmosphere. High frequency radio waves such as HF and VHF bands are spread upwards. When they reach Ionosphere, they are refracted back to the earth.



Microwave Transmission

Electromagnetic waves above 100 MHz tend to travel in a straight line and signals over them can be sent by beaming those waves towards one particular station. Because Microwaves travels in straight lines, both sender and receiver must be aligned to be strictly in line-of-sight.

Microwaves can have wavelength ranging from 1 mm - 1 meter and frequency ranging from 300 MHz to 300 GHz.



Microwave antennas concentrate the waves making a beam of it. As shown in picture above, multiple antennas can be aligned to reach farther. Microwaves have higher frequencies and do not penetrate wall like obstacles.

Microwave transmission depends highly upon the weather conditions and the frequency it is using.

Infrared Transmission

Infrared wave lies in between visible light spectrum and microwaves. It has wavelength of 700-nm to 1-mm and frequency ranges from 300-GHz to 430-THz.

Infrared wave is used for very short range communication purposes such as television and it's remote. Infrared travels in a straight line hence it is directional by nature. Because of high frequency range, Infrared cannot cross wall-like obstacles.

Light Transmission

Highest most electromagnetic spectrum which can be used for data transmission is light or optical signaling. This is achieved by means of LASER.

Because of frequency light uses, it tends to travel strictly in straight line.Hence the sender and receiver must be in the line-of-sight. Because laser transmission is unidirectional, at both ends of communication the laser and the photo-detector needs to be installed. Laser beam is generally 1mm wide hence it is a work of precision to align two far receptors each pointing to lasers source.



Laser works as Tx (transmitter) and photo-detectors works as Rx (receiver).

Lasers cannot penetrate obstacles such as walls, rain, and thick fog. Additionally, laser beam is distorted by wind, atmosphere temperature, or variation in temperature in the path.

Laser is safe for data transmission as it is very difficult to tap 1mm wide laser without interrupting the communication channel.

DCN-Multiplexing

Multiplexing is a technique by which different analog and digital streams of transmission can be simultaneously processed over a shared link. Multiplexing divides the high capacity medium into low capacity logical medium which is then shared by different streams.

Communication is possible over the air (radio frequency), using a physical media (cable), and light (optical fiber). All mediums are capable of multiplexing.

When multiple senders try to send over a single medium, a device called Multiplexer divides the physical channel and allocates one to each. On the other end of communication, a Demultiplexer receives data from a single medium, identifies each, and sends to different receivers.

Frequency Division Multiplexing

When the carrier is frequency, FDM is used. FDM is an analog technology. FDM divides the spectrum or carrier bandwidth in logical channels and allocates one user to each channel. Each user can use the channel frequency independently and has exclusive access of it. All channels are divided in such a way that they do not overlap with each other. Channels are separated by guard bands. Guard band is a frequency which is not used by either channel.



Time Division Multiplexing

TDM is applied primarily on digital signals but can be applied on analog signals as well. In TDM the shared channel is divided among its user by means of time slot. Each user can transmit data within the provided time slot only. Digital signals are divided in frames, equivalent to time slot i.e. frame of an optimal size which can be transmitted in given time slot.

TDM works in synchronized mode. Both ends, i.e. Multiplexer and De-multiplexer are timely synchronized and both switch to next channel simultaneously.



When channel A transmits its frame at one end, the De-multiplexer provides media to channel A on the other end. As soon as the channel A's time slot expires, this side switches to channel B. On the other end, the De-multiplexer works in a synchronized manner and provides media to channel B. Signals from different channels travel the path in interleaved manner.

Wavelength Division Multiplexing

Light has different wavelength (colors). In fiber optic mode, multiple optical carrier signals are multiplexed into an optical fiber by using different wavelengths. This is an analog multiplexing technique and is done conceptually in the same manner as FDM but uses light as signals.



Further, on each wavelength time division multiplexing can be incorporated to accommodate more data signals.

Code Division Multiplexing

Multiple data signals can be transmitted over a single frequency by using Code Division Multiplexing. FDM divides the frequency in smaller channels but CDM allows its users to full bandwidth and transmit signals all the time using a unique code. CDM uses orthogonal codes to spread signals.

Each station is assigned with a unique code, called chip. Signals travel with these codes independently, inside the whole bandwidth. The receiver knows in advance the chip code signal it has to receive.

DCN-NetworkSwitching

Switching is process to forward packets coming in from one port to a port leading towards the destination. When data comes on a port it is called ingress, and when data leaves a port or goes out it is called egress. A communication system may include number of switches and nodes. At broad level, switching can be divided into two major categories:

- **Connectionless:** The data is forwarded on behalf of forwarding tables. No previous handshaking is required and acknowledgements are optional.
- **Connection Oriented:** Before switching data to be forwarded to destination, there is a need to pre-establish circuit along the path between both endpoints. Data is then forwarded on that circuit. After the transfer is completed, circuits can be kept for future use or can be turned down immediately.

Circuit Switching

When two nodes communicate with each other over a dedicated communication path, it is called circuit switching. There 'is a need of pre-specified route from which data will travels and no other data is permitted. In circuit switching, to transfer the data, circuit must be established so that the data transfer can take place.

Circuits can be permanent or temporary. Applications which use circuit switching may have to go through three phases:

- Establish a circuit
- Transfer the data



Circuit switching was designed for voice applications. Telephone is the best suitable example of circuit switching. Before a user can make a call, a virtual path between caller and callee is established over the network.

Message Switching

This technique was somewhere in middle of circuit switching and packet switching. In message switching, the whole message is treated as a data unit and is switching / transferred in its entirety.

A switch working on message switching, first receives the whole message and buffers it until there are resources available to transfer it to the next hop. If the next hop is not having enough resource to accommodate large size message, the message is stored and switch waits.



This technique was considered substitute to circuit switching. As in circuit switching the whole path is blocked for two entities only. Message switching is replaced by packet switching. Message switching has the following drawbacks:

- Every switch in transit path needs enough storage to accommodate entire message.
- Because of store-and-forward technique and waits included until resources are available, message switching is very slow.
- Message switching was not a solution for streaming media and real-time applications.

Packet Switching

Shortcomings of message switching gave birth to an idea of packet switching. The entire message is broken down into smaller chunks called packets. The switching information is added in the header of each packet and transmitted independently.

It is easier for intermediate networking devices to store small size packets and they do not take much resources either on carrier path or in the internal memory of switches.



Packet switching enhances line efficiency as packets from multiple applications can be multiplexed over the carrier. The internet uses packet switching technique. Packet switching enables the user to differentiate data streams based on priorities. Packets are stored and forwarded according to their priority to provide quality of service.

DCN-Data-link Layer Introduction

Data Link Layer is second layer of OSI Layered Model. This layer is one of the most complicated layers and has complex functionalities and liabilities. Data link layer hides the details of underlying hardware and represents itself to upper layer as the medium to communicate.

Data link layer works between two hosts which are directly connected in some sense. This direct connection could be point to point or broadcast. Systems on broadcast network are said to be on same link. The work of data link layer tends to get more complex when it is dealing with multiple hosts on single collision domain.

Data link layer is responsible for converting data stream to signals bit by bit and to send that over the underlying hardware. At the receiving end, Data link layer picks up data from hardware which are in the form of electrical signals, assembles them in a recognizable frame format, and hands over to upper layer.

Data link layer has two sub-layers:

- Logical Link Control: It deals with protocols, flow-control, and error control
- Media Access Control: It deals with actual control of media

Functionality of Data-link Layer

Data link layer does many tasks on behalf of upper layer. These are:

• Framing

Data-link layer takes packets from Network Layer and encapsulates them into Frames.Then, it sends each frame bit-by-bit on the hardware. At receiver' end, data link layer picks up signals from hardware and assembles them into frames.

Addressing

Data-link layer provides layer-2 hardware addressing mechanism. Hardware address is assumed to be unique on the link. It is encoded into hardware at the time of manufacturing.

• Synchronization

When data frames are sent on the link, both machines must be synchronized in order to transfer to take place.

Error Control

Sometimes signals may have encountered problem in transition and the bits are flipped.These errors are detected and attempted to recover actual data bits. It also provides error reporting mechanism to the sender.

• Flow Control

Stations on same link may have different speed or capacity. Data-link layer ensures flow control that enables both machine to exchange data on same speed.

• Multi-Access

When host on the shared link tries to transfer the data, it has a high probability of collision. Data-link layer provides mechanism such as CSMA/CD to equip capability of accessing a shared media among multiple Systems.

DCN-EnorDetection & Conection

There are many reasons such as noise, cross-talk etc., which may help data to get corrupted during transmission. The upper layers work on some generalized view of network architecture and are not aware of actual hardware data processing.Hence, the upper layers expect error-free transmission between the systems. Most of the applications would not function expectedly if they receive erroneous data. Applications such as voice and video may not be that affected and with some errors they may still function well.

Data-link layer uses some error control mechanism to ensure that frames (data bit streams) are transmitted with certain level of accuracy. But to understand how errors is controlled, it is essential to know what types of errors may occur.

TypesofEnors

There may be three types of errors:

• Single bit error



• Multiple bits error



Frame is received with more than one bits in corrupted state.

• Burst error



AC ACCREDITED

Frame contains more than1 consecutive bits corrupted.

Error control mechanism may involve two possible ways:

- Error detection
- Error correction

Enor Detection

Errors in the received frames are detected by means of Parity Check and Cyclic Redundancy Check (CRC). In both cases, few extra bits are sent along with actual data to confirm that bits received at other end are same as they were sent. If the counter-check at receiver' end fails, the bits are considered corrupted.

Parity Check

One extra bit is sent along with the original bits to make number of 1s either even in case of even parity, or odd in case of odd parity.

The sender while creating a frame counts the number of 1s in it. For example, if even parity is used and number of 1s is even then one bit with value 0 is added. This way number of 1s remains even. If the number of 1s is odd, to make it even a bit with value 1 is added.



The receiver simply counts the number of 1s in a frame. If the count of 1s is even and even parity is used, the frame is considered to be not-corrupted and is accepted. If the count of 1s is odd and odd parity is used, the frame is still not corrupted.

If a single bit flips in transit, the receiver can detect it by counting the number of 1s. But when more than one bits are erro neous, then it is very hard for the receiver to detect the error.

Cyclic Redundancy Check (CRC)

CRC is a different approach to detect if the received frame contains valid data. This technique involves binary division of the data bits being sent. The divisor is generated using polynomials. The sender performs a division operation on the bits being sent and calculates the remainder. Before sending the actual bits, the sender adds the remainder at the end of the actual bits. Actual data bits plus the remainder is called a codeword. The sender transmits data bits as codewords.



At the other end, the receiver performs division operation on codewords using the same CRC divisor. If the remainder contains all zeros the data bits are accepted, otherwise it is considered as there some data corruption occurred in transit.

EnorConection

In the digital world, error correction can be done in two ways:

- **Backward Error Correction** When the receiver detects an error in the data received, it requests back the sender to retransmit the data unit.
- Forward Error Correction When the receiver detects some error in the data received, it executes error-correcting code, which helps it to auto-recover and to correct some kinds of errors.

The first one, Backward Error Correction, is simple and can only be efficiently used where retransmitting is not expensive. For example, fiber optics. But in case of wireless transmission retransmitting may cost too much. In the latter case, Forward Error Correction is used.

To correct the error in data frame, the receiver must know exactly which bit in the frame is corrupted. To locate the bit in error, redundant bits are used as parity bits for error detection.For example, we take ASCII words (7 bits data), then there could be 8 kind of information we need: first seven bits to tell us which bit is error and one more bit to tell that there is no error.

For m data bits, r redundant bits are used. r bits can provide 2r combinations of information. In m+r bit codeword, there is possibility that the r bits themselves may get corrupted. So the number of r bits used must inform about m+r bit locations plus no-error information, i.e. m+r+1.

$2^{r} > = m + r + 1$

DCN-Data-link Control & Protocols

Data-link layer is responsible for implementation of point-to-point flow and error control mechanism.

FlowControl

When a data frame (Layer-2 data) is sent from one host to another over a single medium, it is required that the sender and receiver should work at the same speed. That is, sender sends at a speed on which the receiver can process and accept the data. What if the speed (hardware/software) of the sender or receiver differs? If sender is sending too fast the receiver may be overloaded, (swamped) and data may be lost.

Two types of mechanisms can be deployed to control the flow: AMAR

• Stop and Wait

This flow control mechanism forces the sender after transmitting a data frame to stop and wait until the acknowledgement of the data-frame sent is received.



Sliding Window

In this flow control mechanism, both sender and receiver agree on the number of data-frames after which the acknowledgement should be sent. As we learnt, stop and wait flow control mechanism wastes resources, this protocol tries to make use of underlying resources as much as possible.

EnorControl

When data-frame is transmitted, there is a probability that data-frame may be lost in the transit or it is received corrupted. In both cases, the receiver does not receive the correct data-frame and sender does not know anything about any loss. In such case, both sender and receiver are equipped with some protocols which helps them to detect transit errors such as loss of data-frame. Hence, either the sender retransmits the data-frame or the receiver may request to resend the previous data-frame.

Requirements for error control mechanism:

- Error detection The sender and receiver, either both or any, must ascertain that there is some error in the transit.
- **Positive ACK** When the receiver receives a correct frame, it should acknowledge it.
- **Negative ACK** When the receiver receives a damaged frame or a duplicate frame, it sends a NACK back to the sender and the sender must retransmit the correct frame.
- **Retransmission:** The sender maintains a clock and sets a timeout period. If an acknowledgement of a data-frame previously transmitted does not arrive before the timeout the sender retransmits the frame, thinking that the frame or it's acknowledgement is lost in transit.

There are three types of techniques available which Data-link layer may deploy to control the errors by Automatic Repeat Requests (ARQ):

• Stop-and-wait ARQ



The following transition may occur in Stop-and-Wait ARQ:

- The sender maintains a timeout counter.
- When a frame is sent, the sender starts the timeout counter.
 - If acknowledgement of frame comes in time, the sender transmits the next frame in queue.
 - If acknowledgement does not come in time, the sender assumes that either the frame or its acknowledgement is lost in transit. Sender retransmits the frame and starts the timeout counter.
- o If a negative acknowledgement is received, the sender retransmits the frame.
- Go-Back-N ARQ

Stop and wait ARQ mechanism does not utilize the resources at their best. When the acknowledgement is received, the sender sits idle and does nothing. In Go-Back-N ARQ method, both sender and receiver maintain a window.



The sending-window size enables the sender to send multiple frames without receiving the acknowledgement of the previous ones. The receiving-window enables the receiver to receive multiple frames and acknowledge them. The receiver keeps track of incoming frame's sequence number.

When the sender sends all the frames in window, it checks up to what sequence number it has received positive acknowledgement. If all frames are positively acknowledged, the sender sends next set of frames. If sender finds that it has received NACK or has not receive any ACK for a particular frame, it retransmits all the frames after which it does not receive any positive ACK.

• Selective Repeat ARQ

In Go-back-N ARQ, it is assumed that the receiver does not have any buffer space for its window size and has to process each frame as it comes. This enforces the sender to retransmit all the frames which are not acknowledged.

150 9001:2015 & 14001:2015



In Selective-Repeat ARQ, the receiver while keeping track of sequence numbers, buffers the frames in memory and sends NACK for only frame which is missing or damaged.

The sender in this case, sends only packet for which NACK is received.DCN-Network LayerIntroduction

Layer-3 in the OSI model is called Network layer. Network layer manages options pertaining to host and network addressing, managing sub-networks, and internetworking.

Network layer takes the responsibility for routing packets from source to destination within or outside a subnet. Two different subnet may have different addressing schemes or noncompatible addressing types. Same with protocols, two different subnet may be operating on different protocols which are not compatible with each other. Network layer has the responsibility to route the packets from source to destination, mapping different addressing schemes and protocols.

Layer-3Functionalities

Devices which work on Network Layer mainly focus on routing. Routing may include various tasks aimed to achieve a single goal. These can be:

- Addressing devices and networks.
- Populating routing tables or static routes.
- Queuing incoming and outgoing data and then forwarding them according to quality of service constraints set for those packets.
- Internetworking between two different subnets.
- Delivering packets to destination with best efforts.
- Provides connection oriented and connection less mechanism.

Network Layer Features

With its standard functionalities, Layer 3 can provide various features as:

- Quality of service management
- Load balancing and link management
- Security
- Interrelation of different protocols and subnets with different schema.
- Different logical network design over the physical network design.
- L3 VPN and tunnels can be used to provide end to end dedicated connectivity.

Internet protocol is widely respected and deployed Network Layer protocol which helps to communicate end to end devices over the internet. It comes in two flavors. IPv4 which has ruled the world for decades but now is running out of address space. IPv6 is created to replace IPv4 and hopefully mitigates limitations of IPv4 too.

DCN-Network Addressing

Layer 3 network addressing is one of the major tasks of Network Layer. Network Addresses are always logical i.e. these are software based addresses which can be changed by appropriate configurations.

A network address always points to host / node / server or it can represent a whole network. Network address is always configured on network interface card and is generally mapped by system with the MAC address (hardware address or layer-2 address) of the machine for Layer-2 communication.

There are different kinds of network addresses in existence:

- IP
- IPX
- AppleTalk

We are discussing IP here as it is the only one we use in practice these days.



IP addressing provides mechanism to differentiate between hosts and network. Because IP addresses are assigned in hierarchical manner, a host always resides under a specific

network. The host which needs to communicate outside its subnet, needs to know destination network address, where the packet/data is to be sent.

Hosts in different subnet need a mechanism to locate each other. This task can be done by DNS. DNS is a server which provides Layer-3 address of remote host mapped with its domain name or FQDN. When a host acquires the Layer-3 Address (IP Address) of the remote host, it forwards all its packet to its gateway. A gateway is a router equipped with all the information which leads to route packets to the destination host.

Routers take help of routing tables, which has the following information:

• Method to reach the network

Routers upon receiving a forwarding request, forwards packet to its next hop (adjacent router) towards the destination.

The next router on the path follows the same thing and eventually the data packet reaches its destination.

Network address can be of one of the following:

- Unicast (destined to one host)
- Multicast (destined to group)
- Broadcast (destined to all)
- Anycast (destined to nearest one)

A router never forwards broadcast traffic by default. Multicast traffic uses special treatment as it is most a video stream or audio with highest priority. Anycast is just similar to unicast, except that the packets are delivered to the nearest destination when multiple destinations are available.

DCN-Network Layer Routing

When a device has multiple paths to reach a destination, it always selects one path by preferring it over others. This selection process is termed as Routing. Routing is done by special network devices called routers or it can be done by means of software processes. The software based routers have limited functionality and limited scope.

A router is always configured with some default route. A default route tells the router where to forward a packet if there is no route found for specific destination. In case there are multiple path existing to reach the same destination, router can make decision based on the following information:

- Hop Count
- Bandwidth

- Metric
- Prefix-length
- Delay

Routes can be statically configured or dynamically learnt. One route can be configured to be preferred over others.

Unicastrouting

Most of the traffic on the internet and intranets known as unicast data or unicast traffic is sent with specified destination. Routing unicast data over the internet is called unicast routing. It is the simplest form of routing because the destination is already known. Hence the router just has to look up the routing table and forward the packet to next hop.



Broadcastrouting

By default, the broadcast packets are not routed and forwarded by the routers on any network. Routers create broadcast domains. But it can be configured to forward broadcasts in some special cases. A broadcast message is destined to all network devices.

Broadcast routing can be done in two ways (algorithm):

• A router creates a data packet and then sends it to each host one by one. In this case, the router creates multiple copies of single data packet with different destination addresses. All packets are sent as unicast but because they are sent to all, it simulates as if router is broadcasting.

This method consumes lots of bandwidth and router must destination address of each node.

• Secondly, when router receives a packet that is to be broadcasted, it simply floods those packets out of all interfaces. All routers are configured in the same way.



This method is easy on router's CPU but may cause the problem of duplicate packets received from peer routers.

Reverse path forwarding is a technique, in which router knows in advance about its predecessor from where it should receive broadcast. This technique is used to detect and discard duplicates.

Multicast Routing

Multicast routing is special case of broadcast routing with significance difference and challenges. In broadcast routing, packets are sent to all nodes even if they do not want it. But in Multicast routing, the data is sent to only nodes which wants to receive the packets.



The router must know that there are nodes, which wish to receive multicast packets (or stream) then only it should forward. Multicast routing works spanning tree protocol to avoid looping.

Multicast routing also uses reverse path Forwarding technique, to detect and discard duplicates and loops.

Anycast Routing

Anycast packet forwarding is a mechanism where multiple hosts can have same logical address. When a packet destined to this logical address is received, it is sent to the host which is nearest in routing topology.



Anycast routing is done with help of DNS server. Whenever an Anycast packet is received it is enquired with DNS to where to send it. DNS provides the IP address which is the nearest IP configured on it. AAGEMEA

Unicast Routing Protocols

There are two kinds of routing protocols available to route unicast packets:

Distance Vector Routing Protocol

Distance Vector is simple routing protocol which takes routing decision on the number of hops between source and destination. A route with less number of hops is considered as the best route. Every router advertises its set best routes to other routers. Ultimately, all routers build up their network topology based on the advertisements of their peer routers,

For example Routing Information Protocol (RIP).

Link State Routing Protocol

Link State protocol is slightly complicated protocol than Distance Vector. It takes into account the states of links of all the routers in a network. This technique helps routes build a common graph of the entire network. All routers then calculate their best path for routing purposes.for example, Open Shortest Path First (OSPF) and Intermediate System to Intermediate System (ISIS).

Multicast Routing Protocols

Unicast routing protocols use graphs while Multicast routing protocols use trees, i.e. spanning tree to avoid loops. The optimal tree is called shortest path spanning tree.

- **DVMRP** Distance Vector Multicast Routing Protocol
- MOSPF Multicast Open Shortest Path First
- **CBT** Core Based Tree
- **PIM** Protocol independent Multicast

Protocol Independent Multicast is commonly used now. It has two flavors:

PIM Dense Mode

This mode uses source-based trees. It is used in dense environment such as LAN.

• PIM Sparse Mode

This mode uses shared trees. It is used in sparse environment such as WAN.

Routing Algorithms

The routing algorithms are as follows:

Flooding

Flooding is simplest method packet forwarding. When a packet is received, the routers send it to all the interfaces except the one on which it was received. This creates too much burden on the network and lots of duplicate packets wandering in the network.

Time to Live (TTL) can be used to avoid infinite looping of packets. There exists another approach for flooding, which is called Selective Flooding to reduce the overhead on the network. In this method, the router does not flood out on all the interfaces, but selective ones.

Shortest Path

Routing decision in networks, are mostly taken on the basis of cost between source and destination. Hop count plays major role here. Shortest path is a technique which uses various algorithms to decide a path with minimum number of hops.

Common shortest path algorithms are:

- Dijkstra's algorithm
- Bellman Ford algorithm
- Floyd Warshall algorithm

DCN-Internetworking

In real world scenario, networks under same administration are generally scattered geographically. There may exist requirement of connecting two different networks of same kind as well as of different kinds. Routing between two networks is called internetworking. Networks can be considered different based on various parameters such as, Protocol, topology, Layer-2 network and addressing scheme.

In internetworking, routers have knowledge of each other's address and addresses beyond them. They can be statically configured go on different network or they can learn by using internetworking routing protocol.



Routing protocols which are used within an organization or administration are called Interior Gateway Protocols or IGP. RIP, OSPF are examples of IGP. Routing between different organizations or administrations may have Exterior Gateway Protocol, and there is only one EGP i.e. Border Gateway Protocol.

Tunneling

If they are two geographically separate networks, which want to communicate with each other, they may deploy a dedicated line between or they have to pass their data through intermediate networks.

Tunneling is a mechanism by which two or more same networks communicate with each other, by passing intermediate networking complexities. Tunneling is configured at both ends.



When the data enters from one end of Tunnel, it is tagged. This tagged data is then routed inside the intermediate or transit network to reach the other end of Tunnel. When data exists the Tunnel its tag is removed and delivered to the other part of the network.

Both ends seem as if they are directly connected and tagging makes data travel through transit network without any modifications.

Packet Fragmentation

Most Ethernet segments have their maximum transmission unit (MTU) fixed to 1500 bytes. A data packet can have more or less packet length depending upon the application. Devices in the transit path also have their hardware and software capabilities which tell what amount of data that device can handle and what size of packet it can process. If the data packet size is less than or equal to the size of packet the transit network can handle, it is processed neutrally. If the packet is larger, it is broken into smaller pieces and then forwarded. This is called packet fragmentation. Each fragment contains the same destination and source address and routed through transit path easily. At the receiving end it is assembled again.

If a packet with DF (don't fragment) bit set to 1 comes to a router which can not handle the packet because of its length, the packet is dropped.

When a packet is received by a router has its MF (more fragments) bit set to 1, the router then knows that it is a fragmented packet and parts of the original packet is on the way. If packet is fragmented too small, the overhead is increases. If the packet is fragmented too

large, intermediate router may not be able to process it and it might get dropped.

DCN-Network Layer Protocols

Every computer in a network has an IP address by which it can be uniquely identified and addressed. An IP address is Layer-3 (Network Layer) logical address. This address may change every time a computer restarts. A computer can have one IP at one instance of time and another IP at some different time.

Address Resolution Protocol (ARP)

While communicating, a host needs Layer-2 (MAC) address of the destination machine which belongs to the same broadcast domain or network. A MAC address is physically burnt into the Network Interface Card (NIC) of a machine and it never changes.

On the other hand, IP address on the public domain is rarely changed. If the NIC is changed in case of some fault, the MAC address also changes. This way, for Layer-2 communication to take place, a mapping between the two is required.



To know the MAC address of remote host on a broadcast domain, a computer wishing to initiate communication sends out an ARP broadcast message asking, "Who has this IP address?" Because it is a broadcast, all hosts on the network segment (broadcast domain) receive this packet and process it. ARP packet contains the IP address of destination host, the sending host wishes to talk to. When a host receives an ARP packet destined to it, it replies back with its own MAC address.

Once the host gets destination MAC address, it can communicate with remote host using Layer-2 link protocol. This MAC to IP mapping is saved into ARP cache of both sending and receiving hosts. Next time, if they require to communicate, they can directly refer to their respective ARP cache.

Reverse ARP is a mechanism where host knows the MAC address of remote host but requires to know IP address to communicate.

Internet Control Message Protocol (ICMP)

ICMP is network diagnostic and error reporting protocol. ICMP belongs to IP protocol suite and uses IP as carrier protocol. After constructing ICMP packet, it is encapsulated in IP packet. Because IP itself is a best-effort non-reliable protocol, so is ICMP.

Any feedback about network is sent back to the originating host. If some error in the network occurs, it is reported by means of ICMP. ICMP contains dozens of diagnostic and error reporting messages.

ICMP-echo and ICMP-echo-reply are the most commonly used ICMP messages to check the reachability of end-to-end hosts. When a host receives an ICMP-echo request, it is bound to send back an ICMP-echo-reply. If there is any problem in the transit network, the ICMP will report that problem.

Internet Protocol Version 4 (IPv4)

IPv4 is 32-bit addressing scheme used as TCP/IP host addressing mechanism. IP addressing enables every host on the TCP/IP network to be uniquely identifiable.

IPv4 provides hierarchical addressing scheme which enables it to divide the network into sub-networks, each with well-defined number of hosts. IP addresses are divided into many categories:

- Class A it uses first octet for network addresses and last three octets for host addressing
- Class B it uses first two octets for network addresses and last two for host addressing

- Class C it uses first three octets for network addresses and last one for host addressing
- **Class D** it provides flat IP addressing scheme in contrast to hierarchical structure for above three.
- Class E It is used as experimental.

IPv4 also has well-defined address spaces to be used as private addresses (not routable on internet), and public addresses (provided by ISPs and are routable on internet).

Though IP is not reliable one; it provides 'Best-Effort-Delivery' mechanism.

Internet Protocol Version 6(IPv6)

Exhaustion of IPv4 addresses gave birth to a next generation Internet Protocol version 6. IPv6 addresses its nodes with 128-bit wide address providing plenty of address space for future to be used on entire planet or beyond.

IPv6 has introduced Anycast addressing but has removed the concept of broadcasting. IPv6 enables devices to self-acquire an IPv6 address and communicate within that subnet. This auto-configuration removes the dependability of Dynamic Host Configuration Protocol (DHCP) servers. This way, even if the DHCP server on that subnet is down, the hosts can communicate with each other.

IPv6 provides new feature of IPv6 mobility. Mobile IPv6 equipped machines can roam around without the need of changing their IP addresses.

IPv6 is still in transition phase and is expected to replace IPv4 completely in coming years. At present, there are few networks which are running on IPv6. There are some transition mechanisms available for IPv6 enabled networks to speak and roam around different networks easily on IPv4. These are:

- Dual stack implementation
- Tunneling
- NAT-PT

DCN-Transport Layer Introduction

Next Layer in OSI Model is recognized as Transport Layer (Layer-4). All modules and procedures pertaining to transportation of data or data stream are categorized into this layer. As all other layers, this layer communicates with its peer Transport layer of the remote host. Transport layer offers peer-to-peer and end-to-end connection between two processes on remote hosts. Transport layer takes data from upper layer (i.e. Application layer) and then

:2015 & 14001:20

breaks it into smaller size segments, numbers each byte, and hands over to lower layer (Network Layer) for delivery.

Functions

- This Layer is the first one which breaks the information data, supplied by Application layer in to smaller units called segments. It numbers every byte in the segment and maintains their accounting.
- This layer ensures that data must be received in the same sequence in which it was sent.
- This layer provides end-to-end delivery of data between hosts which may or may not belong to the same subnet.
- All server processes intend to communicate over the network are equipped with wellknown Transport Service Access Points (TSAPs) also known as port numbers.

End-to-EndCommunication

A process on one host identifies its peer host on remote network by means of TSAPs, also known as Port numbers. TSAPs are very well defined and a process which is trying to communicate with its peer knows this in advance.



For example, when a DHCP client wants to communicate with remote DHCP server, it always requests on port number 67. When a DNS client wants to communicate with remote DNS server, it always requests on port number 53 (UDP).

The two main Transport layer protocols are:

Transmission Control Protocol

It provides reliable communication between two hosts.

User Datagram Protocol

It provides unreliable communication between two hosts.

DCN-Transmission Control Protocol

The transmission Control Protocol (TCP) is one of the most important protocols of Internet Protocols suite. It is most widely used protocol for data transmission in communication network such as internet.

Features

- TCP is reliable protocol. That is, the receiver always sends either positive or negative acknowledgement about the data packet to the sender, so that the sender always has bright clue about whether the data packet is reached the destination or it needs to resend it.
- TCP ensures that the data reaches intended destination in the same order it was sent.
- TCP is connection oriented. TCP requires that connection between two remote points be established before sending actual data.
- TCP provides error-checking and recovery mechanism.
- TCP provides end-to-end communication.
- TCP provides flow control and quality of service.
- TCP operates in Client/Server point-to-point mode.
- TCP provides full duplex server, i.e. it can perform roles of both receiver and sender.

Header

The length of TCP header is minimum 20 bytes long and maximum 60 bytes.

0 1 2 3	4 5 6	7 0 1 2 3 4 5 6	7	0 1 2 3 4 5 6 7 0 1 2 3 4 5 6 7	
Source Port				Destination Port	
Sequence Number					
Acknowledgement Number					
Data Offset	Reserv	Flags		Window Size	
	Checksum			Urgent	
		e)p t	tions	

- Source Port (16-bits) It identifies source port of the application process on the sending device.
- **Destination Port (16-bits)** It identifies destination port of the application process on the receiving device.
- Sequence Number (32-bits) Sequence number of data bytes of a segment in a session.
- Acknowledgement Number (32-bits) When ACK flag is set, this number contains the next sequence number of the data byte expected and works as acknowledgement of the previous data received.
- **Data Offset (4-bits)** This field implies both, the size of TCP header (32-bit words) and the offset of data in current packet in the whole TCP segment.
- **Reserved (3-bits)** Reserved for future use and all are set zero by default.

- Flags (1-bit each)
 - NS Nonce Sum bit is used by Explicit Congestion Notification signaling process.
 - **CWR** When a host receives packet with ECE bit set, it sets Congestion Windows Reduced to acknowledge that ECE received.
 - ECE -It has two meanings:
 - If SYN bit is clear to 0, then ECE means that the IP packet has its CE (congestion experience) bit set.
 - If SYN bit is set to 1, ECE means that the device is ECT capable.
 - **URG** It indicates that Urgent Pointer field has significant data and should be processed.
 - **ACK** It indicates that Acknowledgement field has significance. If ACK is cleared to 0, it indicates that packet does not contain any acknowledgement.
 - **PSH** When set, it is a request to the receiving station to PUSH data (as soon as it comes) to the receiving application without buffering it.
 - **RST** Reset flag has the following features:
 - It is used to refuse an incoming connection.
 - It is used to reject a segment.
 - It is used to restart a connection.
 - SYN This flag is used to set up a connection between hosts.
 - FIN This flag is used to release a connection and no more data is exchanged thereafter. Because packets with SYN and FIN flags have sequence numbers, they are processed in correct order.
- Windows Size This field is used for flow control between two stations and indicates the amount of buffer (in bytes) the receiver has allocated for a segment, i.e. how much data is the receiver expecting.
- Checksum This field contains the checksum of Header, Data and Pseudo Headers.
- Urgent Pointer It points to the urgent data byte if URG flag is set to 1.
- **Options** It facilitates additional options which are not covered by the regular header. Option field is always described in 32-bit words. If this field contains data less than 32-bit, padding is used to cover the remaining bits to reach 32-bit boundary.

Addressing

TCP communication between two remote hosts is done by means of port numbers (TSAPs). Ports numbers can range from 0 - 65535 which are divided as:

- System Ports (0 1023)
- User Ports (1024 49151)
- Private/Dynamic Ports (49152 65535)

Connection Management

TCP communication works in Server/Client model. The client initiates the connection and the server either accepts or rejects it. Three-way handshaking is used for connection management.



Establishment

Client initiates the connection and sends the segment with a Sequence number. Server acknowledges it back with its own Sequence number and ACK of client's segment which is one more than client's Sequence number. Client after receiving ACK of its segment sends an acknowledgement of Server's response.

Release

Either of server and client can send TCP segment with FIN flag set to 1. When the receiving end responds it back by ACKnowledging FIN, that direction of TCP communication is closed and connection is released.

Bandwidth Management

TCP uses the concept of window size to accommodate the need of Bandwidth management. Window size tells the sender at the remote end, the number of data byte segments the receiver at this end can receive. TCP uses slow start phase by using window size 1 and increases the window size exponentially after each successful communication.

For example, the client uses windows size 2 and sends 2 bytes of data. When the acknowledgement of this segment received the windows size is doubled to 4 and next sent the segment sent will be 4 data bytes long. When the acknowledgement of 4-byte data segment is received, the client sets windows size to 8 and so on.

If an acknowledgement is missed, i.e. data lost in transit network or it received NACK, then the window size is reduced to half and slow start phase starts again.

EnorControl&andFlowControl

TCP uses port numbers to know what application process it needs to handover the data segment. Along with that, it uses sequence numbers to synchronize itself with the remote host. All data segments are sent and received with sequence numbers. The Sender knows which last data segment was received by the Receiver when it gets ACK. The Receiver knows about the last segment sent by the Sender by referring to the sequence number of recently received packet.

If the sequence number of a segment recently received does not match with the sequence number the receiver was expecting, then it is discarded and NACK is sent back. If two segments arrive with the same sequence number, the TCP timestamp value is compared to make a decision.

Multiplexing

The technique to combine two or more data streams in one session is called Multiplexing. When a TCP client initializes a connection with Server, it always refers to a well-defined port number which indicates the application process. The client itself uses a randomly generated port number from private port number pools.

Using TCP Multiplexing, a client can communicate with a number of different application process in a single session. For example, a client requests a web page which in turn contains different types of data (HTTP, SMTP, FTP etc.) the TCP session timeout is increased and the session is kept open for longer time so that the three-way handshake overhead can be avoided.

This enables the client system to receive multiple connection over single virtual connection. These virtual connections are not good for Servers if the timeout is too long.

150 9001:2015 & 14001:2015

CongestionControl

When large amount of data is fed to system which is not capable of handling it, congestion occurs. TCP controls congestion by means of Window mechanism. TCP sets a window size telling the other end how much data segment to send. TCP may use three algorithms for congestion control:

• Additive increase, Multiplicative Decrease

- Slow Start
- Timeout React

Timer Management

TCP uses different types of timer to control and management various tasks:

Keep-alive timer:

- This timer is used to check the integrity and validity of a connection.
- When keep-alive time expires, the host sends a probe to check if the connection still exists.

Retransmission timer:

- This timer maintains stateful session of data sent.
- If the acknowledgement of sent data does not receive within the Retransmission time, the data segment is sent again.

Persist timer:

- TCP session can be paused by either host by sending Window Size 0.
- To resume the session a host needs to send Window Size with some larger value.
- If this segment never reaches the other end, both ends may wait for each other for infinite time.
- When the Persist timer expires, the host re-sends its window size to let the other end know.
- Persist Timer helps avoid deadlocks in communication.

Timed-Wait:

- After releasing a connection, either of the hosts waits for a Timed-Wait time to terminate the connection completely.
- This is in order to make sure that the other end has received the acknowledgement of its connection termination request.
- Timed-out can be a maximum of 240 seconds (4 minutes).

CrashRecovery

TCP is very reliable protocol. It provides sequence number to each of byte sent in segment. It provides the feedback mechanism i.e. when a host receives a packet, it is bound to ACK that packet having the next sequence number expected (if it is not the last segment).

When a TCP Server crashes mid-way communication and re-starts its process it sends TPDU broadcast to all its hosts. The hosts can then send the last data segment which was never unacknowledged and carry onwards.

DCN-User Datagram Protocol

The User Datagram Protocol (UDP) is simplest Transport Layer communication protocol available of the TCP/IP protocol suite. It involves minimum amount of communication mechanism. UDP is said to be an unreliable transport protocol but it uses IP services which provides best effort delivery mechanism.

In UDP, the receiver does not generate an acknowledgement of packet received and in turn, the sender does not wait for any acknowledgement of packet sent. This shortcoming makes this protocol unreliable as well as easier on processing.

RequirementofUDP

A question may arise, why do we need an unreliable protocol to transport the data? We deploy UDP where the acknowledgement packets share significant amount of bandwidth along with the actual data. For example, in case of video streaming, thousands of packets are forwarded towards its users. Acknowledging all the packets is troublesome and may contain huge amount of bandwidth wastage. The best delivery mechanism of underlying IP protocol ensures best efforts to deliver its packets, but even if some packets in video streaming get lost, the impact is not calamitous and can be ignored easily. Loss of few packets in video and voice traffic sometimes goes unnoticed.

Features

- UDP is used when acknowledgement of data does not hold any significance.
- UDP is good protocol for data flowing in one direction.
- UDP is simple and suitable for query based communications.
- UDP is not connection oriented.
- UDP does not provide congestion control mechanism.
- UDP does not guarantee ordered delivery of data.
- UDP is stateless.
- UDP is suitable protocol for streaming applications such as VoIP, multimedia streaming.

UDPHeader

UDP header is as simple as its function.

0 15	5 16 31	
Source Port	Destination Port	
Length	Checksum	

UDP header contains four main parameters:

- **Source Port** This 16 bits information is used to identify the source port of the packet.
- **Destination Port** This 16 bits information, is used identify application level service on destination machine.
- Length Length field specifies the entire length of UDP packet (including header). It is 16-bits field and minimum value is 8-byte, i.e. the size of UDP header itself.
- **Checksum** This field stores the checksum value generated by the sender before sending. IPv4 has this field as optional so when checksum field does not contain any value it is made 0 and all its bits are set to zero.

UDPapplication

Here are few applications where UDP is used to transmit data:

- Domain Name Services
- Simple Network Management Protocol
- Trivial File Transfer Protocol
- Routing Information Protocol
- Kerberos

DCN-ApplicationLayerIntroduction

Application layer is the top most layer in OSI and TCP/IP layered model. This layer exists in both layered Models because of its significance, of interacting with user and user applications. This layer is for applications which are involved in communication system.

A user may or may not directly interacts with the applications. Application layer is where the actual communication is initiated and reflects. Because this layer is on the top of the layer stack, it does not serve any other layers. Application layer takes the help of Transport and all layers below it to communicate or transfer its data to the remote host.

When an application layer protocol wants to communicate with its peer application layer protocol on remote host, it hands over the data or information to the Transport layer. The transport layer does the rest with the help of all the layers below it.



There'is an ambiguity in understanding Application Layer and its protocol. Not every user application can be put into Application Layer. except those applications which interact with the communication system. For example, designing software or text-editor cannot be considered as application layer programs.

On the other hand, when we use a Web Browser, which is actually using Hyper Text Transfer Protocol (HTTP) to interact with the network. HTTP is Application Layer protocol. Another example is File Transfer Protocol, which helps a user to transfer text based or binary files across the network. A user can use this protocol in either GUI based software like FileZilla or CuteFTP and the same user can use FTP in Command Line mode.

Hence, irrespective of which software you use, it is the protocol which is considered at Application Layer used by that software. DNS is a protocol which helps user application protocols such as HTTP to accomplish its work.

DCN-ClientServerModel

Two remote application processes can communicate mainly in two different fashions:

- **Peer-to-peer:** Both remote processes are executing at same level and they exchange data using some shared resource.
- Client-Server: One remote process acts as a Client and requests some resource from another application process acting as Server.

In client-server model, any process can act as Server or Client. It is not the type of machine, size of the machine, or its computing power which makes it server; it is the ability of serving request that makes a machine a server.



A system can act as Server and Client simultaneously. That is, one process is acting as Server and another is acting as a client. This may also happen that both client and server processes reside on the same machine.

Communication

Two processes in client-server model can interact in various ways:

- Sockets
- Remote Procedure Calls (RPC)

Sockets

In this paradigm, the process acting as Server opens a socket using a well-known (or known by client) port and waits until some client request comes. The second process acting as a Client also opens a socket but instead of waiting for an incoming request, the client processes 'requests first'.



When the request is reached to server, it is served. It can either be an information sharing or resource request.

Remote Procedure Call

This is a mechanism where one process interacts with another by means of procedure calls. One process (client) calls the procedure lying on remote host. The process on remote host is
said to be Server. Both processes are allocated stubs. This communication happens in the following way:

- The client process calls the client stub. It passes all the parameters pertaining to program local to it.
- All parameters are then packed (marshalled) and a system call is made to send them to other side of the network.
- Kernel sends the data over the network and the other end receives it.
- The remote host passes data to the server stub where it is unmarshalled.
- The parameters are passed to the procedure and the procedure is then executed.
- The result is sent back to the client in the same manner.

DCN-ApplicationProtocols

There are several protocols which work for users in Application Layer. Application layer protocols can be broadly divided into two categories:

- Protocols which are used by users.For email for example, eMail.
- Protocols which help and support protocols used by users. For example DNS.

Few of Application layer protocols are described below:

Domain Name System

The Domain Name System (DNS) works on Client Server model. It uses UDP protocol for transport layer communication. DNS uses hierarchical domain based naming scheme. The DNS server is configured with Fully Qualified Domain Names (FQDN) and email addresses mapped with their respective Internet Protocol addresses.

A DNS server is requested with FQDN and it responds back with the IP address mapped with it. DNS uses UDP port 53.

Simple Mail Transfer Protocol

The Simple Mail Transfer Protocol (SMTP) is used to transfer electronic mail from one user to another. This task is done by means of email client software (User Agents) the user is using. User Agents help the user to type and format the email and store it until internet is available. When an email is submitted to send, the sending process is handled by Message Transfer Agent which is normally comes inbuilt in email client software.

Message Transfer Agent uses SMTP to forward the email to another Message Transfer Agent (Server side). While SMTP is used by end user to only send the emails, the Servers normally use SMTP to send as well as receive emails. SMTP uses TCP port number 25 and 587.

Client software uses Internet Message Access Protocol (IMAP) or POP protocols to receive emails.

File Transfer Protocol

The File Transfer Protocol (FTP) is the most widely used protocol for file transfer over the network. FTP uses TCP/IP for communication and it works on TCP port 21. FTP works on Client/Server Model where a client requests file from Server and server sends requested resource back to the client.

FTP uses out-of-band controlling i.e. FTP uses TCP port 20 for exchanging controlling information and the actual data is sent over TCP port 21.

The client requests the server for a file. When the server receives a request for a file, it opens a TCP connection for the client and transfers the file. After the transfer is complete, the server closes the connection. For a second file, client requests again and the server reopens a new TCP connection.

PostOfficeProtocol(POP)

The Post Office Protocol version 3 (POP 3) is a simple mail retrieval protocol used by User Agents (client email software) to retrieve mails from mail server.

When a client needs to retrieve mails from server, it opens a connection with the server on TCP port 110. User can then access his mails and download them to the local computer. POP3 works in two modes. The most common mode the delete mode, is to delete the emails from remote server after they are downloaded to local machines. The second mode, the keep mode, does not delete the email from mail server and gives the user an option to access mails later on mail server.

HyperTextTransferProtocol(HTTP)

The Hyper Text Transfer Protocol (HTTP) is the foundation of World Wide Web. Hypertext is well organized documentation system which uses hyperlinks to link the pages in the text documents. HTTP works on client server model. When a user wants to access any HTTP page on the internet, the client machine at user end initiates a TCP connection to server on port 80. When the server accepts the client request, the client is authorized to access web pages.

To access the web pages, a client normally uses web browsers, who are responsible for initiating, maintaining, and closing TCP connections. HTTP is a stateless protocol, which means the Server maintains no information about earlier requests by clients.

HTTP versions

- HTTP 1.0 uses non persistent HTTP. At most one object can be sent over a single TCP connection.
- HTTP 1.1 uses persistent HTTP. In this version, multiple objects can be sent over a single TCP connection.

DCN-Network Services

Computer systems and computerized systems help human beings to work efficiently and explore the unthinkable. When these devices are connected together to form a network, the capabilities are enhanced multiple-times. Some basic services computer network can offer are. ABGERTI

Directory Services

These services are mapping between name and its value, which can be variable value or fixed. This software system helps to store the information, organize it, and provides various means of accessing it.

• Accounting

In an organization, a number of users have their user names and passwords mapped to them. Directory Services provide means of storing this information in cryptic form and make available when requested.

Authentication & and Authorization

User credentials are checked to authenticate a user at the time of login and/or periodically. User accounts can be set into hierarchical structure and their access to resources can be controlled using authorization schemes.

Domain Name Services

DNS is widely used and one of the essential services on which internet works. This system maps IP addresses to domain names, which are easier to remember and recall than IP addresses. Because network operates with the help of IP addresses and humans tend to remember website names, the DNS provides website's IP address which is mapped to its name from the back-end on the request of a website name from the user.

File Services

File services include sharing and transferring files over the network.

• File Sharing

One of the reason which gave birth to networking was file sharing. File sharing enables its users to share their data with other users. User can upload the file to a specific server, which is accessible by all intended users. As an alternative, user can make its file shared on its own computer and provides access to intended users.

• File Transfer

This is an activity to copy or move file from one computer to another computer or to multiple computers, with help of underlying network. Network enables its user to locate other users in the network and transfers files.

C ACCREDITED

Communication Services

• Email

Electronic mail is a communication method and something a computer user cannot work without. This is the basis of today's internet features. Email system has one or more email servers. All its users are provided with unique IDs. When a user sends email to other user, it is actually transferred between users with help of email server.

Social Networking

Recent technologies have made technical life social. The computer savvy peoples, can find other known peoples or friends, can connect with them, and can share thoughts, pictures, and videos.

• Internet Chat

Internet chat provides instant text transfer services between two hosts. Two or more people can communicate with each other using text based Internet Relay Chat services. These days, voice chat and video chat are very common.

Discussion Boards

Discussion boards provide a mechanism to connect multiple peoples with same interests. It enables the users to put queries, questions, suggestions etc. which can be seen by all other users. Other may respond as well.

Remote Access

This service enables user to access the data residing on the remote computer. This feature is known as Remote desktop. This can be done via some remote device, e.g. mobile phone or home computer.

Application Services

These are nothing but providing network based services to the users such as web services, database managing, and resource sharing.

• Resource Sharing

To use resources efficiently and economically, network provides a mean to share them. This may include Servers, Printers, and Storage Media etc.

Databases

This application service is one of the most important services. It stores data and information, processes it, and enables the users to retrieve it efficiently by using queries. Databases help organizations to make decisions based on statistics.

Web Services

World Wide Web has become the synonym for internet. It is used to connect to the internet, and access files and information services provided by the internet servers. NNAGE

IT Act of India 2000

In May 2000, both the houses of the Indian Parliament passed the Information Technology Bill. The Bill received the assent of the President in August 2000 and came to be known as the Information Technology Act, 2000. Cyber laws are contained in the IT Act, 2000.

This Act aims to provide the legal infrastructure for e-commerce in India. And the cyber laws have a major impact for e-businesses and the new economy in India. So, it is important to understand what are the various perspectives of the IT Act, 2000 and what it offers.

The Information Technology Act, 2000 also aims to provide for the legal framework so that legal sanctity is accorded to all electronic records and other activities carried out by electronic means. The Act states that unless otherwise agreed, an acceptance of contract may be expressed by electronic means of communication and the same shall have legal validity and enforceability. Some highlights of the Act are listed below:

Chapter-II of the Act specifically stipulates that any subscriber may authenticate an electronic record by affixing his digital signature. It further states that any person can verify an electronic record by use of a public key of the subscriber.

Chapter-III of the Act details about Electronic Governance and provides inter alia amongst others that where any law provides that information or any other matter shall be in writing or in the typewritten or printed form, then, notwithstanding anything contained in such law, such requirement shall be deemed to have been satisfied if such information or matter is rendered or made available in an electronic form; and accessible so as to be usable for a subsequent reference. The said chapter also details the legal recognition of Digital Signatures

Chapter-IV of the said Act gives a scheme for Regulation of Certifying Authorities. The Act envisages a Controller of Certifying Authorities who shall perform the function of exercising supervision over the activities of the Certifying Authorities as also laying down standards and conditions governing the Certifying Authorities as also specifying the various forms and content of Digital Signature Certificates. The Act recognizes the need for recognizing foreign Certifying Authorities and it further details the various provisions for the issue of license to issue Digital Signature Certificates. Chapter-VII of the Act details about the scheme of things relating to Digital Signature Certificates. The duties of subscribers are also enshrined in the said Act. Chapter-IX of the said Act talks about penalties and adjudication for various offences. The penalties for damage to computer, computer systems etc. has been fixed as damages by way of compensation not exceeding Rs. 1,00,00,000 to affected persons. The Act talks of appointment of any officers not below the rank of a Director to the Government of India or an equivalent officer of state government as an Adjudicating Officer who shall adjudicate whether any person has made a contravention of any of the provisions of the said Act or rules framed there under. The said Adjudicating Officer has been given the powers of Civil a Court.

Chapter-X of the Act talks of the establishment of the Cyber Regulations Appellate Tribunal, which shall be an appellate body where appeals against the orders passed by the Adjudicating Officers, shall be preferred. Chapter-XI of the Act talks about various offences and the said offences shall be investigated only by a Police Officer not below the rank of the Deputy Superintendent of Police. These offences include tampering with computer source documents, publishing of information, which is obscene in electronic form, and hacking.

The Act also provides for the constitution of the Cyber Regulations Advisory Committee, which shall advice the government as regards any rules, or for any other purpose connected with the said act. The said Act also proposes to amend the Indian Penal Code, 1860, the Indian Evidence Act, 1872, The Bankers' Books Evidence Act, 1891, The Reserve Bank of India Act, 1934 to make them in tune with the provisions of the IT Act.

Advantages of Cyber Laws The IT Act 2000 attempts to change outdated laws and provides ways to deal with cyber crimes. We need such laws so that people can perform purchase transactions over the Net through credit cards without fear of misuse. The Act offers the much-needed legal framework so that information is not denied legal effect, validity or enforceability, solely on the ground that it is in the form of electronic records. In view of the

growth in transactions and communications carried out through electronic records, the Act seeks to empower government departments to accept filing, creating and retention of official documents in the digital format. The Act has also proposed a legal framework for the authentication and origin of electronic records / communications through digital signature. From the perspective of e-commerce in India, the IT Act 2000 and its provisions contain many positive aspects. Firstly, the implications of these provisions for the e-businesses would be that email would now be a valid and legal form of communication in our country that can be duly produced and approved in a court of law. Companies shall now be able to carry out electronic commerce using the legal infrastructure provided by the Act. Digital signatures have been given legal validity and sanction in the Act. The Act throws open the doors for the entry of corporate companies in the business of being Certifying Authorities for issuing Digital Signatures Certificates. The Act now allows Government to issue notification on the web heralding thus e-governance. The Act enables the companies to file any form, application or any other document with any office, authority, body or agency owned or controlled by the appropriate Government in electronic form by means of such electronic form as may be prescribed by the appropriate Government. The IT Act also addresses the important issues of security, which are so critical to the success of electronic transactions. The Act has given a legal definition to the concept of secure digital signatures that would be required to have been passed through a system of a security procedure, as stipulated by the Government at a later date.

Under the IT Act, 2000, it shall now be possible for corporates to have a statutory remedy in case if anyone breaks into their computer systems or network and causes damages or copies data. The remedy provided by the Act is in the form of monetary damages, not exceeding Rs. 1 crore.

Applications of IT

Every day, people use computers in new ways. Computers and other electronic devices are becoming increasingly affordable. They continue to be more powerful as informationprocessing tools as well as easier to use. Humans are continually becoming dependant on ITenabled devices for carrying out simple tasks like remembering a phone number to complex ones like flying a fighter plane. Information Technology has applications in almost all aspects of our life. Some of the important ones are:

Science and Engineering: Scientific progress in fields like biotechnology is almost entirely dependent on the use of computers and other microprocessor-controlled devices. Using

supercomputers, meteorologists predict future weather by using a combination of observations of weather conditions from many sources, a mathematical representation of the behavior of the atmosphere, and geographic data. Computer-aided design (CAD) and computer-aided manufacturing (CAM) programs have led to improved products in many fields, especially where designs tend to be very detailed. Computer programs make it possible for engineers to analyze designs of complex structures such as power plants and space stations.

Business & **Commerce**: One of the first and largest applications of computers is keeping and managing business and financial records. Most large companies keep the employment records of all their workers in large databases that are managed by computer programs. Similar programs and databases are used in business functions like billing customers; tracking payments received and payments to be made; and tracking supplies needed and items produced, stored, shipped, and sold. In fact, practically all the information companies need to do business involves the use of computers and Information Technology. Almost all the financial transactions in the world are done electronically. Newer technologies like m-commerce have enabled almost everybody to carry out routine financial transactions on the move.

On a smaller scale, many businesses have replaced cash registers with point-of-sale (POS) terminals. These POS terminals not only print a sales receipt for the customer but also send information to a computer database when each item is sold to maintain an inventory of items on hand and items to be ordered. Computers have also become very important in modern factories. Computer-controlled robots now do tasks that are hot, heavy, or hazardous. Robots are also used to do routine, repetitive tasks in which boredom or fatigue can lead to poor quality work.

With today's sophisticated hardware, software, and communications technologies, it is often difficult to classify a system as belonging uniquely to one specific application program. Organizations increasingly are consolidating their information needs into a single, integrated information system. Management Information System (MIS), with the Chief Information Officer (CIO) at its head, is a whole, new branch of enterprise management.

Education: The advent of Information Technology has changed the meaning of the term "literate", with computer literacy being almost as important as basic literacy in many cases. Computer education is an essential course at the primary level in most schools across the world. With more information getting digitized every day, and the internet making it accessible to anyone across the world, students are increasingly relying on electronic sources

of information rather than physical libraries for their needs. Instructional methodology has also undergone a sea change with use of images, animations, videos, presentations and elearning to complement traditional techniques.

Governance: The concept of e-governance is one of the most novel applications of Information Technology whereby it is changing the lives of millions across the globe. Computerization of Government activities makes it easier to supervise and audit, and makes the administration more responsive to the needs of society. It also bridges the divide between the Government and the people. Technologies like touch-screen kiosks help disseminate information on land records, photo identity cards, pending bills etc. and enable even illiterate people to take more informed decisions. India is leading the world in the effective use of IT for elections.

Medicine: Information Technology plays an important role in medicine. For example, a scanner takes a series of pictures of the body by means of computerized axial tomography (CAT) or magnetic resonance imaging (MRI). A computer then combines the pictures to produce detailed three-dimensional images of the body's organs. In addition, the MRI produces images that show changes in body chemistry and blood flow. Most critical life support equipment are programmed to respond to changes in the patient's status in split-seconds, thereby reducing the response time and risk of human error. Newer concepts like robotic surgery enable specialists to perform surgeries from remote locations. Genomic studies greatly depend on supercomputing power to develop technologies for the future.

Entertainment: IT has changed the lifestyle of most people. The convergence of various technologies has created various options for entertainment like games, streaming music and video, digital television broadcasts, satellite radio, animated movies etc. which can be accessed with the help of mobile phones, PDAs, notebook computers or on television either with a cable connection or wirelessly using newer-generation WiFi, CDMA or GPRS technologies. Information Technology plays a vital role in most of our daily activities. There is hardly anyone who has not been affected or influenced by IT. With each passing day, newer applications of IT are being developed which increase our interaction with and dependence on IT-enabled devices. Therefore, understanding this technology and using it creatively is imperative to human progress.

Computer Security Issues

• Hacking unauthorized access to or use of data, systems, server or networks, including any attempt to probe, scan or test the vulnerability of a system, server or network or to breach security or authentication measures without express authorization of the owner of the system,

server or network. Members of the University should not run computer programs that are associated with hacking without prior authorisation. Obtaining and using such programs is not typical of normal usage and may therefore otherwise be regarded as misuse.

- Use of University owned computer equipment, including the network, for illegal activities including copying Copyright material without permission. The vast majority of files shared on **P2P** (**peer-to-peer**) networks violate copyright law because they were posted without permission of the artist or label.
- Sending abusive e-mails or posting offensive Web pages.
- Creation or transmission of any offensive or indecent images.
- Giving unauthorised access to University computing resources e.g. allowing an account to be used by someone not authorised to use it.
- Deliberately creating or spreading computer viruses or worms.
- Unauthorised running of applications that involve committing the University to sharing its computing resources, e.g. network bandwidth, in an uncontrolled and unlimited way.





PHYSICS (109)

<u>Unit 1</u>

Newton's Laws of Motion

Sir Isaac Newton's three laws of motion describe the motion of massive bodies and how they interact. While Newton's laws may seem obvious to us today, more than three centuries ago they were considered revolutionary. Newton was one of the most influential scientists of all time. His ideas became the basis for modern physics. He built upon ideas put forth from the works of previous scientists including Galileo and Aristotle and was able to prove some ideas that had only been theories in the past. He studied optics, astronomy and math — he invented calculus. (German mathematician Gottfried Leibniz is also credited with developing it independently at about the same time.) Newton is perhaps best known for his work in studying gravity and the motion of planets. Urged on by astronomer Edmond Halley after admitting he had lost his proof of elliptical orbits a few years prior, Newton published his laws in 1687, in his seminal work "Philosophiæ Naturalis Principia Mathematica" (Mathematical Principles of Natural Philosophy) in which he formalized the description of how massive bodies move under the influence of external forces. Newton's laws pertain to the motion of massive bodies in an inertial reference frame, sometimes called a Newtonian reference frame, although Newton himself never described such a reference frame. An inertial reference frame can be described as a 3-dimensional coordinate system that is either stationary or in uniform linear motion. i.e., it is not accelerating or rotating. He found that motion within such an inertial reference frame could be described by three simple laws.

The First Law of Motion states, "A body at rest will remain at rest, and a body in motion will remain in motion unless it is acted upon by an external force." This simply means that things cannot start, stop, or change direction all by themselves. It takes some force acting on them from the outside to cause such a change. This property of massive bodies to resist changes in their state of motion is sometimes called *inertia*.

The Second Law of Motion describes what happens to a massive body when it is acted upon by an external force. It states, "The force acting on an object is equal to the mass of that object times its acceleration." This is written in mathematical form as F = ma, where F is force, *m* is mass, and a is acceleration. The bold letters indicate that force and acceleration are *vector* quantities, which means they have both magnitude and direction. The force can be a single force, or it can be the vector sum of more than one force, which is the net force after all the forces are combined. When a constant force acts on a massive body, it causes it to accelerate, i.e., to change its velocity, at a constant rate. In the simplest case, a force applied to an object at rest causes it to accelerate in the direction of the force. However, if the object is already in motion, or if this situation is viewed from a moving reference frame, that body might appear to speed up, slow down, or change direction depending on the direction of the force and the directions that the object and reference frame are moving relative to each other.

The Third Law of Motion states, "For every action, there is an equal and opposite reaction." This law describes what happens to a body when it exerts a force on another body. Forces always occur in pairs, so when one body pushes against another, the second body pushes back just as hard. For example, when you push a cart, the cart pushes back against you; when you pull on a rope, the rope pulls back against you; when gravity pulls you down against the ground, the ground pushes up against your feet; and when a rocket ignites its fuel behind it, the expanding exhaust gas pushes on the rocket causing it to accelerate.

If one object is much, much more massive than the other, particularly in the case of the first object being anchored to the Earth, virtually all of the acceleration is imparted to the second object, and the acceleration of the first object can be safely ignored. For instance, if you were to throw a baseball to the west, you would not have to consider that you actually caused the rotation of the Earth to speed up ever so slightly while the ball was in the air. However, if you were standing on roller skates, and you threw a bowling ball forward, you would start moving backward at a noticeable speed.

The three laws have been verified by countless experiments over the past three centuries, and they are still being widely used to this day to describe the kinds of objects and speeds that we encounter in everyday life. They form the foundation of what is now known as *classical mechanics*, which is the study of massive objects that are larger than the very small scales addressed by quantum mechanics and that are moving slower than the very high speeds addressed by relativistic mechanics.

Force and Inertia

Inertia.

Inertia is the resistance of any physical object to any change in its velocity. This includes changes to the object's speed, or direction of motion. An aspect of this property is the tendency of objects to keep moving in a straight line at a constant speed, when no forces act upon them.

In common usage, the term "inertia" may refer to an object's "amount of resistance to change in velocity" or for simpler terms, "resistance to a change in motion" (which is quantified by its mass), or sometimes to its momentum, depending on the context. The term "inertia" is more properly understood as shorthand for "the principle of inertia" as described by Newton in his first law of motion: an object not subject to any net external force moves at a constant velocity. Thus, an object will continue moving at its current velocity until some force causes its speed or direction to change. On the surface of the Earth, inertia is often masked by gravity and the effects of friction and air resistance, both of which tend to decrease the speed of moving objects (commonly to the point of rest). This misled the philosopher Aristotle to believe that objects would move only as long as force was applied to them. The principle of inertia is one of the fundamental principles in classical physics that are still used today to describe the motion of objects and how they are affected by the applied forces on them.

Force.

In physics, a **force** is any interaction that, when unopposed, will change the motion of an object. A force can cause an object with mass to change its velocity (which includes to begin moving from a state of rest), i.e., to accelerate. Force can also be described intuitively as a push or a pull. A force has both magnitude and direction, making it a vector quantity. It is measured in the SI unit of Newton's and represented by the symbol **F**. The original form of Newton's second law states that the net force acting upon an object is equal to the rate at which its momentum changes with time. If the mass of the object is constant, this law implies that the acceleration of an object is directly proportional to the net force acting on the object, is in the direction of the net force, and is inversely proportional to the mass of the object. Concepts related to force include: thrust, which increases the velocity of an object; drag, which decreases the velocity of an object; and torque, which produces changes in rotational speed of an object. In an extended body, each part usually applies forces on the adjacent parts; the distribution of such forces through the body is the internal mechanical stress. Such internal mechanical stresses cause no acceleration of that body as the forces balance one another. Pressure, the distribution of many small forces applied over an area of a body, is a simple type of stress that if unbalanced can cause the body to accelerate. Stress usually causes deformation of solid materials, or flow in fluids.

Newton's laws of motion

Newton's laws of motion are three physical laws that, together, laid the foundation for classical mechanics. They describe the relationship between a body and the forces acting upon it, and its motion in response to those forces. More precisely, the first law defines the force qualitatively, the second law offers a quantitative measure of the force, and the third asserts that a single isolated force doesn't exist. These three laws have been expressed in several ways, over nearly three centuries,[a] and can be summarized as follows:

First law: In an inertial frame of reference, an object either remains at rest or continues to move at a constant velocity, unless acted upon by a force.

Second law: In an inertial frame of reference, the vector sum of the forces **F** on an object is equal to the mass *m* of that object multiplied by the acceleration **a** of the object: **F** = *m***a**. (It is assumed here that the mass *m* is constant – see below.)

Third law: When one body exerts a force on a second body, the second body simultaneously exerts a force equal in magnitude and opposite in direction on the first body.

The three laws of motion were first compiled by Isaac Newton in his *Philosophiæ Naturalis Principia Mathematica*

(*Mathematical Principles of Natural Philosophy*), first published in 1687. Newton used them to explain and investigate the motion of many physical objects and systems. For example, in the third volume of the text, Newton showed that these laws of motion, combined with his law of universal gravitation, explained Kepler's laws of planetary motion.

Some also describe a **fourth law** which states that forces add up like vectors, that is, that forces obey the principle of superposition.

Overview

Newton's laws are applied to objects which are idealized as single point masses, in the sense that the size and shape of the object's body are neglected to focus on its motion more easily. This can be done when the object is small compared to the distances involved in its analysis, or the deformation and rotation of the body are of no importance. In this way, even a planet can be idealized as a particle for analysis of its orbital motion around a star. In their original form, Newton's laws of motion are not adequate to characterize the motion of rigid bodies and deformable bodies. Leonhard Euler in 1750 introduced a generalization of Newton's laws of motion for rigid bodies called Euler's laws of motion, later applied as well for deformable bodies assumed as a continuum. If a body is represented as an assemblage of discrete particles, each governed by Newton's laws of motion, then Euler's laws can be derived from Newton's laws. Euler's laws can, however, be taken as axioms describing the laws of motion for extended bodies, independently of any particle structure. Newton's laws hold only with respect to a certain set of frames of reference called Newtonian or inertial reference frames.

Some authors interpret the first law as defining what an inertial reference frame is; from this point of view, the second law holds only when the observation is made from an inertial reference frame, and therefore the first law cannot be proved as a special case of the second. Other authors do treat the first law as a corollary of the second. The explicit concept of an inertial frame of reference was not developed until long after Newton's death. In the given interpretation mass, acceleration, momentum, and (most importantly) force are assumed to be externally defined quantities. This is the most common, but not the only interpretation of the way one can consider the laws to be a definition of these quantities. Newtonian mechanics has been superseded by special relativity, but it is still useful as an approximation when the speeds involved are much slower than the speed of light.

Newton's laws read:

Law I: Everybody persists in its state of being at rest or of moving uniformly straight forward, except insofar as it is compelled to change its state by force impressed.

Law II: The alteration of motion is ever proportional to the motive force impressed; and is made in the direction of the right line in which that force is impressed.

Law III: To every action there is always opposed an equal reaction: or the mutual actions of two bodies upon each other are always equal, and directed to contrary parts.

A thAMAG

Newton's first law

The first law states that if the net force (the vector sum of all forces acting on an object) is zero, then the velocity of the object is constant. Velocity is a vector quantity which expresses both the object's speed and the direction of its motion; therefore, the statement that the object's velocity is constant is a statement that both its speed and the direction of its motion are constant.

The first law can be stated mathematically when the mass is a non-zero constant, as,

$$\sum \mathbf{F} = 0 \iff rac{\mathrm{d}\mathbf{v}}{\mathrm{d}t} = 0.$$

Newton's second law

The second law states that the rate of change of momentum of a body is directly proportional to the force applied, and this change in momentum takes place in the direction of the applied force.

$$\mathbf{F} = rac{\mathrm{d}\mathbf{p}}{\mathrm{d}t} = rac{\mathrm{d}(m\mathbf{v})}{\mathrm{d}t}.$$

The second law can also be stated in terms of an object's acceleration. Since Newton's second law is valid only for constant-mass systems, [20][21][22] *m* can be taken outside the differentiation operator by the constant factor rule in differentiation.

$$\mathbf{F}=mrac{\mathrm{d}\mathbf{v}}{\mathrm{d}t}=m\mathbf{a},$$

Thus, where **F** is the net force applied, m is the mass of the body, and **a** is the body's acceleration. Thus, the net force applied to a body produces a proportional acceleration. In other words, if a body is accelerating, then there is a force on it. An application of this notation is the derivation of G Subscript C.

The above statements hint that the second law is merely a definition of , not a precious observation of nature. However, current physics restate the second law in measurable steps:

(1) defining the term 'one unit of mass' by a specified stone,

(2) defining the term 'one unit of force' by a specified spring with specified length,

(3) measuring by experiment or proving by theory (with a principle that every direction of space are equivalent), that force can be added as a mathematical vector,

(4) finally conclude that . These steps hint the second law is a precious feature of nature.

The second law also implies the conservation of momentum: when the net force on the body is zero, the momentum of the body is constant. Any net force is equal to the rate of change of the momentum.

Any mass that is gained or lost by the system will cause a change in momentum that is not the result of an external force. A different equation is necessary for variable-mass

Newton's second law is an approximation that is increasingly worse at high speeds because of relativistic effects.

The change of momentum of a body is proportional to the impulse impressed on the body, and happens along the straight line on which that impulse is impressed.

Impulse

An impulse **J** occurs when a force **F** acts over an interval of time Δt , and it is given by

$$\mathbf{J} = \int_{\Delta t} \mathbf{F} \, \mathrm{d}t.$$

Since force is the time derivative of momentum, it follows that

$$\mathbf{J} = \Delta \mathbf{p} = m \Delta \mathbf{v}.$$

CDEDI

This relation between impulse and momentum is closer to Newton's wording of the second law. Impulse is a concept frequently used in the analysis of collisions and impacts.

A DAVALLING DAV

Newton's third law



An illustration of Newton's third law in which two skaters push against each other. The first skater on the left exerts a normal force N_{12} on the second skater directed towards the right, and the second skater exerts a normal force N_{21} on the first skater directed towards the left. The magnitudes of both forces are equal, but they have opposite directions, as dictated by Newton's third law.

The third law states that all forces between two objects exist in equal magnitude and opposite direction: if one object *A* exerts a force \mathbf{F}_A on a second object *B*, then *B* simultaneously exerts a force \mathbf{F}_B on *A*, and the two forces are equal in magnitude and opposite in direction: $\mathbf{F}_A = -\mathbf{F}_B$. The third law means that all forces are *interactions* between different bodies, or different regions within one body, and thus that there is no such thing as a force that is not accompanied by an equal and opposite force. In some situations, the magnitude and direction of the forces are determined entirely by one of the two bodies, say Body *A*; the force exerted by Body *A* on Body *B* is called the "action", and the force exerted by Body *B* on Body *A* is called the "reaction". This law is sometimes referred to as the *action-reaction law*,

with \mathbf{F}_A called the "action" and \mathbf{F}_B the "reaction". In other situations the magnitude and directions of the forces are determined jointly by both bodies and it isn't necessary to identify one force as the "action" and the other as the "reaction". The action and the reaction are simultaneous, and it does not matter which is called the *action* and which is called *reaction*; both forces are part of a single interaction, and neither force exists without the other.

The two forces in Newton's third law are of the same type (e.g., if the road exerts a forward frictional force on an accelerating car's tires, then it is also a frictional force that Newton's third law predicts for the tires pushing backward on the road).

From a conceptual standpoint, Newton's third law is seen when a person walks: they push against the floor, and the floor pushes against the person. Similarly, the tires of a car push against the road while the road pushes back on the tires—the tires and road simultaneously push against each other. In swimming, a person interacts with the water, pushing the water backward, while the water simultaneously pushes the person forward—both the person and the water push against each other. The reaction forces account for the motion in these examples. These forces depend on friction; a person or car on ice, for example, may be unable to exert the action force to produce the needed reaction force.

Newton used the third law to derive the law of conservation of momentum; from a deeper perspective, however, conservation of momentum is the more fundamental idea (derived via Noether's theorem from Galilean invariance), and holds in cases where Newton's third law appears to fail, for instance when force fields as well as particles carry momentum, and in quantum mechanics.

Newton's third law



An illustration of Newton's third law in which two skaters push against each other. The first skater on the left exerts a normal force N_{12} on the second skater directed towards the right, and the second skater exerts a normal force N_{21} on the first skater directed towards the left.

The magnitudes of both forces are equal, but they have opposite directions, as dictated by Newton's third law.

The third law states that all forces between two objects exist in equal magnitude and opposite direction: if one object A exerts a force \mathbf{F}_A on a second object B, then B simultaneously exerts a force \mathbf{F}_B on A, and the two forces are equal in magnitude and opposite in direction: $\mathbf{F}_A =$ $-\mathbf{F}_{B}$. The third law means that all forces are *interactions* between different bodies, or different regions within one body, and thus that there is no such thing as a force that is not accompanied by an equal and opposite force. In some situations, the magnitude and direction of the forces are determined entirely by one of the two bodies, say Body A; the force exerted by Body A on Body B is called the "action", and the force exerted by Body B on Body A is called the "reaction". This law is sometimes referred to as the action-reaction law, with \mathbf{F}_A called the "action" and \mathbf{F}_B the "reaction". In other situations the magnitude and directions of the forces are determined jointly by both bodies and it isn't necessary to identify one force as the "action" and the other as the "reaction". The action and the reaction are simultaneous, and it does not matter which is called the *action* and which is called *reaction*; both forces are part of a single interaction, and neither force exists without the other. The two forces in Newton's third law are of the same type (e.g., if the road exerts a forward frictional force on an accelerating car's tires, then it is also a frictional force that Newton's third law predicts for the tires pushing backward on the road).

From a conceptual standpoint, Newton's third law is seen when a person walks: they push against the floor, and the floor pushes against the person. Similarly, the tires of a car push against the road while the road pushes back on the tires—the tires and road simultaneously push against each other. In swimming, a person interacts with the water, pushing the water backward, while the water simultaneously pushes the person forward—both the person and the water push against each other. The reaction forces account for the motion in these examples. These forces depend on friction; a person or car on ice, for example, may be unable to exert the action force to produce the needed reaction force. Newton used the third law to derive the law of conservation of momentum; from a deeper perspective, however, conservation of momentum is the more fundamental idea (derived via Noether's theorem from Galilean invariance), and holds in cases where Newton's third law appears to fail, for instance when force fields as well as particles carry momentum, and in quantum mechanics.

Fundamental Forces of Nature

The Four Fundamental Forces of Nature are Gravitational force, Weak Nuclear force, Electromagnetic force and Strong Nuclear force. The weak and strong forces are effective only over a very short range and dominate only at the level of subatomic particles. Gravity and Electromagnetic force have infinite range. Let's see each of them in detail.

The Four Fundamental Forces and their strengths

- 1. Gravitational Force Weakest force; but infinite range.
- 2. Weak Nuclear Force Next weakest; but short range.
- 3. Electromagnetic Force Stronger, with infinite range.
- 4. Strong Nuclear Force Strongest; but short range.

Gravitational Force

The gravitational force is weak, but very long ranged. Furthermore, it is always attractive. It acts between any two pieces of matter in the Universe since mass is its source.

Weak Nuclear Force

The weak force is responsible for radioactive decay and neutrino interactions. It has a very short range and. As its name indicates, it is very weak. The weak force causes Beta decay ie. the conversion of a neutron into a proton, an electron and an antineutrino.

Electromagnetic Force

The electromagnetic force causes electric and magnetic effects such as the repulsion between like electrical charges or the interaction of bar magnets. It is long-ranged, but much weaker than the strong force. It can be attractive or repulsive, and acts only between pieces of matter carrying electrical charge. Electricity, magnetism, and light are all produced by this force.

Strong Nuclear Force

The strong interaction is very strong, but very short-ranged. It is responsible for holding the nuclei of atoms together. It is basically attractive, but can be effectively repulsive in some circumstances. The strong force is 'carried' by particles called gluons; that is, when two

particles interact through the strong force, they do so by exchanging gluons. Thus, the quarks inside of the protons and neutrons are bound together by the exchange of the strong nuclear force.

Weight of body in lift

Apparent weight

When body is at rest with no acceleration, R = W. Reading on the weighing machine reflects the true weight, W, force of gravity acting on our body (mg).

When we are in equilibrium, the normal reaction is equal to the weight.

Case of lift with upward acceleration

- As acceleration of lift is upward, the resultant force is upward.
- R = W + ma

Since R is greater than W, the weighing machine shows a reading greater than the actual force due to gravity (W), the person feels heavier or its apparent weight is heavier.

The person would also feel the same way when slowing down during a descent.

Case of lift with downward acceleration

– As acceleration of lift is downward, the resultant force is downward. – R = W - ma

Since R is less than W, the weighing machine shows a reading which is a reading lesser than the actual force due to gravity – the person feels lighter or its apparent weight is lighter.

The person would also feel the same way when slowing down during an ascent.

<u>Weightlessness</u>

- If the acceleration a is equal to g, the lift is free-falling, then we have R = 0.
- The machine would register a zero reading and is not in contact with the body. Therefore, apparent weight is zero and the body experiences weightlessness.
- A body is said to be free-falling and experiencing apparent weightlessness if the only force acting on it is its true weight (mg) and its acceleration, a is equal to g.

A TALL NO. TALK & R. MAR, MAR

For a force-time graph, area under graph is the impulse (Change in momentum)

Average force

 $\langle F \rangle = \Delta p \Delta t \langle F \rangle = \Delta p \Delta t$

For N number of collisions: $\langle F \rangle = \Delta p \times N\Delta t \langle F \rangle = \Delta p \times N\Delta t$ Area under average force graph = area under actual force graph. $\langle F \rangle = \Delta \langle F \rangle = \Delta$ in momentum in one collision × collision frequency $\langle F \rangle = v \times dmdt \langle F \rangle = v \times dmdt$

Equilibrium of concurrent forces

What is Equilibrium?

Equilibrium means "no acceleration". Since a force is a "push" or "pull" exerted on a body, equilibrium means that the total of all forces acting on a body must be zero.

According to Newton's second law, $\mathbf{F} = \mathbf{m} * \mathbf{a}$ (Remember! In Newton's second law \mathbf{F} is the TOTAL force on the body) Since we are studying STATICS, from now on we assume every body is in equilibrium.

Equilibrium of Concurrent Forces

Concurrent means that the forces intersect through a single point. If forces are concurrent, we can add them together as vectors to get the resultant. If the body is not accelerating, it must be in equilibrium, so that means the resultant is zero. For concurrent forces, the body is a point. So for concurrent forces in equilibrium, the forces

should all add up to give zero.

If a body is not accelerating is in **equilibrium**, so resultant of all forces = 0.

AAC ACCREDI

A typical concurrent force situation is a lifting eye. The pulling forces in any cables must pass through the centre of the eye. If there is only one eyebolt (correctly positioned over the centre of gravity) and the load is suspended, the bolt force must pass through the same centre. Hence all forces pass through one point (the centre of the eye), so we have concurrent equilibrium. All forces on a suspended load are concurrent. (Assuming the load remains level when lifted). It is possible to maintain equilibrium even when the cables are at different angles. In the example below, Cable B must have less tension than Cable A;



Diagrams

1. The Space Diagram (SD)

The initial problem is usually sketched. This illustration or picture shows the layout and dimensions. If this diagram was drawn to scale, the units would be length (mm, m etc). It is nice to be accurate, but it does not *have* to be to scale.

2. The Free Body Diagram (FBD)

The Free Body Diagram is a strict diagram that isolates the body for study. See <u>Free Body</u> <u>Diagrams</u> for more information. The idea of the FBD is to focus on one particular part or group of parts (called the body) and replace every external member with the force they would

Isolate the body. (An outline is best because we are supposed to forget about the inside of the body)
 Locate border crossings. Identify the contact points where forces are crossing the boundary.

Gravity acts through centre.

3. Line of Action. Some types of connections have a known direction. E.g. Cables have force running through the centreline.

AL D. G.

4. "To the Body". Since Newton's 3rd law has every action with an opposite reaction, we must eliminate half the forces. Identify those forces that are applied "to the body", and eliminate those done "by the body".

If the FBD were drawn to scale, the body might be length (mm, m etc) and the forces might be another scale (N, kN etc).

Warning! Do not get Linear dimensions and Force dimensions mixed up. You cannot add metres

and Newtons together!

3. The Force Polygon (FP)

The force polygon must be drawn strictly to scale, and everything is a Force. The only information coming from the FBD is;

ICA 0001.001C 0 14001.0011

- Force magnitudes
- Force Angles

Warning! Do not attempt to bring any FBD Lengths into the FP. There are no metres in the Force Polygon.

COPYRIGHT FIMT 2020

In some cases the Free Body diagram does not even look like the original. This is most obvious for concurrent forces. Since all forces go through one point, we can treat the body as a DOT!



Cable connection in a structure, specially designed to make the centreline of every cables intersect at one point.

Example Diagrams. These cranes are not accelerating, so they are in equilibrium. Therefore all the forces on any body should add up to zero. The body is actually the connection point which is probably a lifting eye of a hook. The FBD shows as much as we know from the Space diagram - in this case angles are known but only one magnitude. The force polygon should form a closed loop (since resultant = 0), so this defines the lengths (and hence the F2. magnitudes) of F1 and helpful when working with force CAD programs are very polygons.



Space Diagram

45° F2 45° 30° 49.05 kN

Free Body Diagram



Force Polygon

Special Contact Points

When drawing the Free Body Diagram we must include all the forces that cross the boundary (outline) of the body. Some of these contact points have special clues about the

direction	of	the	force	and	location	of	the	force.
-----------	----	-----	-------	-----	----------	----	-----	--------

- 1. **Cable Joint**: The force must run through the centreline of the cable
- 2. Frictionless Joint: The force must be perpendicular to the surface.
- Wheels and Rollers: The force must run through the centre of the axle. Free running wheels are frictionless so force is perpendicular to the surface and all forces pass through the centre of the axle.
- 4. **Pulleys:** The tension in the cable is the same on each side of the pulley and all forces pass can be made to go through the centre of the axle.
- 5. **Friction**: The force can be in any direction
- 6. Pin Joint: The force can be in any direction

Contact points are also called Support Reactions; Here is a table (Ignore the last one at this stage).

Number of Forces acting on a Body

One Force

This is impossible for equilibrium. The forces are supposed to add up to zero (unless the body is accelerating. E.g. A falling rock).

Two Forces

If a body has only 2 forces, they must be co-linear. E.g. A linkage between 2 pivot pins must have the force running through the line of the pins. (This assumes gravity force is ignored, otherwise you have three forces)

2015 & 14001-2015

Three Forces

If a body has exactly 3 forces, they must be concurrent. This is called the Three Force Principle. This can be very handy in solving problems because many mechanisms have bodies with 2 or 3 forces.

Four or more Forces...

COPYRIGHT FIMT 2020

We cannot assume the forces will be concurrent, unless specially made that way. (Like the five-way cable connection below). When forces are not concurrent they can create rotations, which we deal with in a later chapter. (Non Concurrent forces)

Five deliberately concurrent forces

The Equilibrium Equations

Equilibrium simply says the resultant is zero. Mathematically, this can be stated that the Fx and Fy components are zero. So, for concurrent forces in 2 dimensions (planar), equilibrium means that...

A AAADDOITED



Very often we know the angle of the forces but not the magnitudes. When solving mathematically, this means we will need to use simultaneous equations.

Lemi's Theorem

In physics, **Lami's theorem** is an equation relating the magnitudes of three coplanar, concurrent and non-collinear forces, which keeps an object in static equilibrium, with the angles directly opposite to the corresponding forces.

According to the theorem, $rac{A}{\sinlpha}=rac{B}{\sineta}=rac{C}{\sin\gamma}$

where *A*, *B* and *C* are the magnitudes of the three coplanar, concurrent and non-collinear forces, F_A , F_B , F_C , which keep the object in static equilibrium, and α , β and γ are the angles directly opposite to the forces. Lami's theorem is applied in static analysis of mechanical and structural systems. The theorem is named after Bernard Lamy.



Friction

Friction is the resistance to motion of one object moving relative to another. It is not a fundamental force, like gravity or electromagnetism. Instead, scientists believe it is the result of the electromagnetic attraction between charged particles in two touching surfaces.

Causes of Friction

Friction is a force resisting motion of an object when in contact with another. This resistive force is caused by the surface roughness of the contact area of the materials, molecular attraction or adhesion between materials, and deformations in the materials. The cause of friction may be any or all of these items and this applies to sliding, rolling and fluid frictions.

Questions you may have about friction include:

- How does surface roughness cause friction?
- How do deformations cause friction?
- How does molecular attraction cause friction?

Surface roughness

Most friction results because the surfaces of materials being rubbed together are not completely smooth. If you looked at what seems to be a smooth surface under a highpowered microscope, you would see bumps, hills and valleys that could interfere with sliding motion. Of course, the rougher the surface, the more the friction.

Close-up view of surface roughness

Treads add to friction

Treads or grooves on one or both sliding surfaces can increase the friction, especially if the treads have sharp edges and are not parallel with the line of motion. The most common use if treads are seen in automobile and bicycle tires, as used in rolling friction. You also may see them on pads intended to keep surfaces from sliding.

Sharp edges of treads add to sliding friction

The number and types of grooves or treads is an added factor to the friction equation.

Molecular adhesion

Another factor in friction can be caused by molecular adhesion or attraction. Ultrasmooth materials and "sticky" materials fall in this category.

Ultra-smooth

If both surfaces are ultra-smooth and flat, the friction from surface roughness becomes negligible, but then friction from molecular attraction comes into play. This can often become greater than friction if the surfaces where relatively rough.

Sticky materials

Rubber is an example of a material that can have friction caused by molecular attraction. Discounting resistance due to deformations with rubber, it is its stickiness factor that causes it to grip so well and have so much friction.

Fluids

Fluids often exhibit molecular adhesion, increasing the friction. This adhesion force is often seen in the capillary effect. This is where water will be pulled up a glass tube by the forces of molecular adhesion. That same force can slow down fluid motion.

One example is how a coin will easily slide down a ramp. But if you wet the coin, it will stay in place. That is because of the molecular friction of the fluid on the hard surfaces.

The motion of two fluids or two sections of a fluid against each other is also slowed down by the molecular attraction factor. This type of fluid friction is usually not considered as friction and is studied under the complex field of fluid dynamics.

Deformations

Soft materials will deform when under pressure. This also increased the resistance to motion. For example, when you stand on a rug, you sink in slightly, which causes resistance when you try to drag your feet along the rug's surface. Another example is how rubber tires flatten out at the area on contact with the road.

When materials deform, you must "plow" through to move, thus creating a resistive force.



Pushing object on soft surface

When the deformation becomes large, such that one object sinks into the other, streamlining can affect the friction, similar to what happens in fluid friction.

ALGENTI

Four Types of Friction

Friction is the force that opposes motion between any surfaces that are in contact. There are four types of friction: static, sliding, rolling, and **fluid** friction. Static, sliding, and rolling friction occur between <u>solid</u> surfaces. Fluid friction occurs in liquids and gases. All four types of friction are described below.

Static Friction

Static friction acts on objects when they are resting on a surface. For example, if you are hiking in the woods, there is static friction between your shoes and the trail each time you put down your foot. Without this static friction, your feet would slip out from under you, making it difficult to walk. In fact, that's exactly what happens if you try to walk on ice. That's because ice is very slippery and offers very little friction.



Sliding Friction

Sliding friction is friction that acts on objects when they are sliding over a surface. Sliding friction is weaker than static friction. That's why it's easier to slide a piece of furniture over the floor after you start it moving than it is to get it moving in the first place. Sliding friction can be useful. For example, you use sliding friction when you write with a pencil. The pencil

"lead" slides easily over the paper, but there's just enough friction between the pencil and paper to leave a mark.

Rolling Friction

Rolling friction is friction that acts on objects when they are rolling over a surface. Rolling friction is much weaker than sliding friction or static friction. This explains why most forms of ground transportation use wheels, including bicycles, cars, 4-wheelers, roller skates, scooters, and skateboards. Ball bearings are another use of rolling friction. You can see what they look like in the **Figure** <u>below</u>. They let parts of a wheel or other <u>machine</u> roll rather than slide over on another.



The ball bearings in this wheel reduce friction between the inner and outer cylinders when they turn.

Fluid Friction

Fluid friction is friction that acts on objects that are moving through a fluid. A **fluid** is a substance that can flow and take the shape of its container. Fluids include liquids and gases. If you've ever tried to push your open hand through the water in a tub or pool, then you've experienced fluid friction. You can feel the resistance of the water against your hand. Look at the skydiver in the **Figure** below. He's falling toward Earth with a parachute. Resistance of the air against the parachute slows his descent. The faster or larger a moving object is, the greater is the fluid friction resisting its motion. That's why there is greater air resistance against the parachute than the skydiver's body.



THE FIVE LAWS OF FRICTION

 When an object is moving, the friction is proportional and perpendicular to the normal force (N)

ACCREDIT

- 2. Friction is independent of the area of contact so long as there is an area of contact.
- 3. The coefficient of static friction is slightly greater than the coefficient of kinetic friction.
- 4. Within rather large limits, kinetic friction is independent of velocity.
- 5. Friction depends upon the nature of the surfaces in contact.

Angle of friction and angle of repose



Angle of Friction beween any two surface in contact is defined as the angle which the resultant

of the force of limiting friction F_{lim} and normal reaction N makes with the direction of normal reaction N as shown in figure. It is marked in the figure as θ .

The value of angle of friction depends on the material and nature of surfaces in contact. It can be seen from figure, tan θ = AC/OA = F_{lim} / N = μ , where μ is coefficient of limiting friction

COPYRIGHT FIMT 2020

609 | Page



The angle of repose or angle of sliding α is defined as the minimum angle of inclination of a plane with the horizontal such that a body placed on the plane just begins to slide down.

Its value depends on the material and nature of the surface in contact.

from the figure it can be seen that,

f =mg sinα(1)

 $N = mg \cos \alpha$ (2)

hence we get, f / N = tan α = μ

Centrifugal Force vs. Centripetal Force

Centrifugal and Centripetal forces may both be called a force, but one of them is really not a force at all. There is a relationship between the two forces, however. Centrifugal results from inertia, the tendency of an object to resist any change in its state of motion or when it is at rest, though it is technically not a force. Centrifugal describes an object as it flies outward along a curved path, away from the center of the curve. Often times it is called an "apparent force", mainly because it feels like a force.

On the other hand, centripetal force is a true force that will offset the centrifugal "force" stopping the motion of the object from its flying outward, keeping it in motion instead at a consistent speed along the curved or circular path. Centripetal is the force that prevents the moon from floating out of the Earth's orbit.

Banking of Roads

The phenomenon of raising outer edge of the curved road above the inner edge is to provide necessary centripetal force to the vehicles to take a safer turn and the curved road is called Banking of Roads....

Definition

The phenomenon of raising outer edge of the curved road above the inner edge is to provide necessary centripetal force to the vehicles to take a safer turn and the curved road is MANAG called Banking of Roads.

Introduction

When a vehicle goes round a curved road, it requires some centripetal force. While rounding the curve, the wheels of the vehicle have a tendency to leave the curved path and regain the straight line path. Force of friction between wheels and the roads opposes this tendency of the wheels. This force of friction therefore, acts towards the centre of circular track and provides the necessary centripetal force.



Fig. (1) Vehicle moving on level road

In fig (1), it is shown that a vehicle of weight 'mg' (acts vertically downwards) is moving on a level curved road. R1 and R2 are the forces of normal reaction of the road on the wheels. These are vertically upward since road is leveled. Hence,

R1 + R2 = mg

Let F1 & F2 are forces of frictions between tyre and road directed towards centre of curved road.

 \therefore F1 = μ R1

COPYRIGHT FIMT 2020

And

$$F2 = \mu R2$$

where μ is coefficient of friction between tyres and road.

13.1

3.4

If 'v' is the velocity of the vehicle while rounding the curve, the centripetal force required is mv^2/r . As this force is provided by the force of friction therefore

1.4

$$\frac{mv^{2}}{r} \leq (F_{1} + F_{2})$$

$$\leq (\mu R_{1} + \mu R_{2})$$

$$\leq \mu (R_{1} + R_{2})$$

$$\frac{mv^{2}}{r} \leq \mu mg$$

$$v^{2} \leq \mu rg$$

$$v \leq \sqrt{\mu rg}$$

Hence the maximum velocity with which a vehicle can go round a level curve; without skidding is

$$U = \sqrt{\mu rg}$$

Banking of Roads

In the above discussion, we see that the maximum permissible velocity with which a vehicle can go round a level curved road depends on μ , the coefficient of friction between tyres and road. The value of μ decreases when road is wet or extra smooth or tyres of the vehicle are worn out. Thus force of friction is not a reliable source for providing the required centripetal force to the vehicle. Especially in hilly areas where the vehicle has to move constantly along the curved track, the maximum speed at which it can run will be very low. If any attempt is made to run it at a greater speed, the vehicle is likely to skid and go out of track. In order that the vehicle can go round the curved track at a reasonable speed without skidding, the sufficient centripetal force is managed for it by raising the outer edge of the track a little above the inner edge. It is called banking of the circular track or **Banking of Roads**.
Consider a vehicle of weight 'Mg' moving round a curved path of radius 'r' with speed 'V' on a road banked through angle θ . If OA is banked road and OX is horizontal line, then $\angle AOX = \theta$ is called angle of **banking of road**. Refer **Fig (2)**



Fig.(2) Vehicle moving on Banked Road

Following forces are involved:

- 1. The weight 'Mg' acting vertically downwards
- The reaction 'R' of the ground to the vehicle acting along normal to the banked road OA in upward direction
- 3. The vertical component R.Cos θ of R will balance the weight of the vehicle.
- 4. The horizontal component R.Sin θ of R will provide necessary centripetal force to the vehicle.

Thus,

R.cos θ = Mg

And

$$R\sin\theta = \frac{Mv^2}{r} \cdot \cdot$$

On dividing equation (1) & equation (2), we get

. (2)

....(1)

 $Mv^2 lr$ $R \sin \theta$ Rcos 0 $\tan \theta = \frac{v^2}{2}$(3) rg

Knowing 'v' and 'r1', we can calculate θ . If 'h' is the height AX of outer edge of the road then from fig.(3),

:2015 & 140



$$OX = \sqrt{OA^2 - AX^2} = \sqrt{b^2 - h^2}$$

$$\tan \theta = \frac{AX}{OX} = \frac{h}{\sqrt{b^2 - h^2}} \qquad \dots (4)$$

From equations (3) & (4) we get

 $\tan \theta = \frac{v^2}{rg} = \frac{h}{\sqrt{b^2 - h^2}}$

From above eqn. we can calculate h. usually h < < b. Therefore. h2 is negligible, hence

$$\tan\theta = \frac{v^2}{rg} = \frac{h}{b}$$

Roads are generally banked for the average speed of vehicles passing over them. However, if the speed of a vehicle is somewhat less or more than this, the self adjusting state friction will operate between tyre and road and vehicle will not skid.

Unit 2

Work, Energy & Power

Work.

When a force acts upon an object to cause a displacement of the object, it is said that **work** was done upon the object. There are three key *ingredients* to work - force, displacement, and cause. In order for a force to qualify as having done *work* on an object, there must be a displacement and the force must *cause* the displacement. There are several good examples of work that can be observed in everyday life - a horse pulling a plow through the field, a father pushing a grocery cart down the aisle of a grocery store, a freshman lifting a backpack full of books upon her shoulder, a weightlifter lifting a barbell above his head, an Olympian launching the shot-put, etc. In each case described here there is a force exerted upon an object to cause that object to be displaced.

Work Equation

Mathematically, work can be expressed by the following equation.

$W = F \bullet d \bullet \cos \Theta$

where **F** is the force, **d** is the displacement, and the angle (**theta**) is defined as the angle between the force and the displacement vector. Perhaps the most difficult aspect of the above equation is the angle "theta." The angle is not just *any 'ole angle*, but rather a very specific

angle. The angle measure is defined as the angle between the force and the displacement. To gather an idea of it's meaning, consider the following three scenarios.

 Scenario A: A force acts rightward upon an object as it is displaced rightward. In such an instance, the force vector and the displacement vector are in the same direction.
 Thus, the angle between F and d is 0 degrees.



- Scenario B: A force acts leftward upon an object that is displaced rightward. In such an instance, the force vector and the displacement vector are in the opposite direction. Thus, the angle between F and d is 180 degrees.
- Scenario C: A force acts upward on an object as it is displaced rightward. In such an instance, the force vector and the displacement vector are at right angles to each other.
 Thus, the angle between F and d is 90 degrees.



Whenever F and d are in the same direction,

COPYRIGHT FIMT 2020

The Meaning of Theta

When determining the measure of the angle in the work equation, it is important to recognize that the angle has a precise definition - it is the angle between the force and the displacement vector. Be sure to avoid mindlessly using *any 'ole angle* in the equation. A common physics lab involves applying a force to displace a cart up a ramp to the top of a chair or box. A *force* is applied to a cart to *displace* it up the incline at constant speed. Several incline angles are typically used; yet, the force is always applied parallel to the incline. The displacement of the cart is also parallel to the incline. Since F and d are in the same direction, the angle theta in the work equation is 0 degrees. Nevertheless, most students experienced the strong temptation to measure the angle of incline and use it in the equation. Don't forget: the angle in the equation is not just *any 'ole angle*. It is defined as the angle between the force and the displacement vector.

The Meaning of Negative Work

On occasion, a force acts upon a moving object to hinder a displacement. Examples might include a car skidding to a stop on a roadway surface or a baseball runner sliding to a stop on the infield dirt. In such instances, the force acts in the direction opposite the objects motion in order to slow it down. The force doesn't cause the displacement but rather *hinders* it. These situations involve what is commonly called *negative work*. The *negative* of negative work refers to the numerical value that results when values of F, d and theta are substituted into the work equation. Since the force vector is directly opposite the displacement vector, theta is 180 degrees. The cosine(180 degrees) is -1 and so a negative value results for the amount of work done upon the object.

Units of Work

Whenever a new quantity is introduced in physics, the standard metric units associated with that quantity are discussed. In the case of work (and also energy), the standard metric unit is the **Joule** (abbreviated **J**). One Joule is equivalent to one Newton of force causing a displacement of one meter. In other words,

The Joule is the unit of work. 1 Joule = 1 Newton * 1 meter

1 J = 1 N * m

In fact, any unit of force times any unit of displacement is equivalent to a unit of work. Some nonstandard units for work are shown below. Notice that when analyzed, each set of units is equivalent to a force unit times a displacement unit.

Conservative Forces

A force is said to be conservative if the work done by or against it in moving an object is independent of the object's path. The work done by a conservative force depends only on the initial and final positions. Gravity is one example, the spring force another. The work done by a nonconservative force depends on the path through which the force acts. The most common example is this kind of force is friction. You can study the action of a conservative force, gravity, and a non-conservative force, friction, in the <u>work applet</u> in this chapter. As long as you lift the box in this applet off the ground, you can make it take any path, and when you return the box to its initial position, the total work done will be 0. But if you slide it back and forth along the ground, the work will always increase.

Power:- Power is the rate at which work is done.

Or

Power is the rate at which energy is converted from one form to another.

11.20

The unit of power is the Watt (W).



Kinetic energy is energy possessed by an object in motion. The earth revolving around the sun, you walking down the street, and molecules moving in space all have kinetic energy. Kinetic energy is directly proportional to the mass of the object and to the square of its velocity: *K.E.* = $1/2 m v^2$. If the mass has units of kilograms and the velocity of meters per

second, the kinetic energy has units of kilograms-meters squared per second squared. Kinetic energy is usually measured in units of Joules (J); one Joule is equal to 1 kg m^2 / s^2 .

Potential energy is energy an object has because of its position relative to some other object. When you stand at the top of a stairwell you have more potential energy than when you are at the bottom, because the earth can pull you down through the force of gravity, doing work in the process. When you are holding two magnets apart they have more potential energy than when they are close together. If you let them go, they will move toward each other, doing work in the process. The formula for potential energy depends on the force acting on the two objects. For the gravitational force the formula is **P.E.** = mgh, where m is the mass in kilograms, g is the acceleration due to gravity (9.8 m / s² at the surface of the earth) and h is the height in meters. Notice that gravitational potential energy has the same units as kinetic energy, kg m² / s². In fact, *all* energy has the same units, kg m² / s², and is measured using the unit Joule (J).

Work-Energy Theorem

The work-energy theorem states that the work done by all forces acting on a particle equals the change in the particle's kinetic energy.

The work 'W' done by the net force on a particle is equal the change in the particle's kinetic energy (KE).

$$\mathrm{d} = rac{\mathrm{v_f^2} - \mathrm{v_i^2}}{2\mathrm{a}}$$

Check the detailed work-energy theorem derivation given below.

Let us consider a case where the resultant force 'F' is constant in both direction and magnitude and is parallel to the velocity of the particle. The particle is moving with constant acceleration along a straight line. The relationship between the acceleration and the net force is given by the equation "F = ma" (Newton's second law of motion), and the particle's displacement 'd', can be determined from the equation:

$$\mathbf{v_f^2} = \mathbf{v_i^2} + 2\mathbf{ad}$$

Obtaining,

COPYRIGHT FIMT 2020

$W=\Delta \mathrm{KE}=rac{1}{2}\mathrm{m}\mathrm{v}_\mathrm{f}^2-rac{1}{2}\mathrm{m}\mathrm{v}_\mathrm{i}^2$

The work of the net force is calculated as the product of its magnitude (F=ma) and the particle's displacement. Substituting the above equations yields:

 $W=Fd=ma\frac{v_f^2-v_i^2}{2a}=\frac{1}{2}mv_f^2-\frac{1}{2}mv_i^2=KE_f-KE_i=\Delta KE$

Conservation of gravitational P.E. into K.E

1. Work done against gravity - gravitational potential energy

Let's now consider the work done when we lift an object. In order to lift an object that has mass m, we have to apply an upward force mg to overcome the downward force of gravity. If this force raises the object through a height h, then the work done is:

 $W = Fd = mg \times h = mgh$



Figure 5 (a) Placing a suitcase on a luggage rack involves doing work against gravity. (b) The stored energy is released if the suitcase falls off the rack. So if an object of mass *m* is raised through a height *h*, the work done on the object is equal to *mgh*, and so this amount of energy is transferred to the object. (Notice that this equation is identical to the one describing an object falling under gravity, Equation 7.) Of course, this ties in very well with everyday observations. If you lift a heavy suitcase onto a luggage rack in a train, or a heavy bag of shopping onto a table, you are very aware that you are doing work against gravity. You will also be aware that more work is required to lift a more massive object, or the same object to a greater height, and these 'observations' are consistent with the work done being equal to *mgh*.

According to the law of conservation of energy, energy can't just disappear. When work is done on a toy train, the energy supplied is converted into kinetic energy (and some internal energy when friction is taken into account), yet a suitcase placed on a luggage rack is obviously stationary. In general, when an object is raised to a greater height, work is done on the object and the energy transferred is stored; the amount of energy stored is mgfch, where Ah is the change in height. This stored energy is given the name of gravitational potential energy. The term 'potential' signifies that this energy has the 'potential' for doing work when the object is lowered. However, as gravitational potential energy is a bit longdwinded we will usually omit the 'potential' and refer to this energy just as gravitational energy $E_{\rm g}$. You will meet other forms of potential energy later in this block. Now in all problems that involve gravitational energy it is the *changes* of gravitational energy that are important. A *change* of height leads to a *change* of gravitational energy; the absolute value of the gravitational energy is unimportant. So if we drop our 12 kg suitcase through 2 m in a room at the top of a tower block, it will have the same kinetic energy, and the same speed, just before it hits the floor as if it falls 2 m from the luggage rack to the floor of a train. And in neither case can we say that the suitcase has no gravitational energy after falling, since we could push it off a balcony, or out of the train, and its gravitational energy would decrease as it fell further.

We can emphasise the importance of these changes by writing an equation for gravitational energy in terms of changes. Thus:

change in gravitational energy = $mg \times$ change in height

or using the Δ (delta) notation:

$$\Delta E_{g} = mg\Delta h$$

Remember that ΔE_g means 'the change in E_g ', and Δh means 'the change in h' Clearly when an object is raised to a greater height, its gravitational energy increases, and when it falls to a lower height, then its gravitational energy decreases.

2. Gravitational energy and energy conservation

The concept of gravitational energy greatly simplifies calculations concerned with the effect of gravity on the motion of objects, particularly where no other forces are concerned. For example, if you take your book (which by now is getting rather dogeared) and throw it vertically upwards in the air, it will slow down as it travels higher, and eventually reach a point where it is momentarily stationary. It will then accelerate downwards to where you (hopefully) catch it before it hits the ground. By using the idea of the change in gravitational energy, coupled with the law of conservation of energy, we can calculate such quantities as the speed of the book at a given height and the maximum height it will reach. These calculations assume that no air resistance acts on the book after it is thrown, so, in accordance with the law of conservation of energy, the sum of its kinetic energy and its gravitational energy must be constant. To say the same thing in another way — any increase in gravitational energy will be accompanied by a decrease in kinetic energy of equal size, and vice versa.

This will become clearer if we take a specific example. A book of mass *m* is thrown vertically upwards with an initial speed *u*. After it has risen through a height Δh it has a speed *v*.

Since energy is conserved, the increase in gravitational energy must be equal to the decrease in kinetic energy, so we can write down the following equation:

$$mg\Delta h = \frac{1}{2}mu^2 - \frac{1}{2}mv^2$$

Of course this equation doesn't apply only to books; it also applies to the movement of any object when gravity is the only force acting on it.



Figure An unsuccessful attempt to rescue a damsel in distress.

Potential Energy of the spring

- To study this ,consider an electric spring of negligibly small mass .One end of the spring is attached to the rigid wall and another end of spring is attached to a block of mass m which can move on smooth frictionless horizontal surface
 - Consider the figure given below



Figure 6. PE of a spring

Here un-stretched or un-compressed position of the spring is taken at x=0

- We now take the block from its un-stretched position to a point P by stretching the spring
- At this point P restoring force is exerted by the spring on the block trying it bring it back to the equilibrium position.
- Similar restoring force developed in the spring when we try to compress it
- For an ideal spring ,this restoring force F is proportional to displacement x and direction of restoring force is opposite to that displacement
- Thus force and displacement are related as
 F α x
 or F= kx (16)

where K is called the spring constant and this equation (16) is known as Hook's law.negative sign indicates that force oppose the motion of the block along x

• To stretch a spring we need to apply the external force which should be equal in magnitude and opposite to the direction of the restoring force mentioned above i.e for stretching the spring $F_{ext}=Kx$ Similarly for compressing the spring $-F_{ext}=-Kx$

or F_{ext}=Kx (both F and x are being negative)

- Work done in both elongation and compression of spring is stored in the spring as its PE which can be easily calculated
- If the spring is stretched through a distance x from its equilibrium position x=0 then W=JF_{ext}dx

Since Doth Text and UX have Same direction Nov	Since	both	F _{ext} and	dx	have	same	direction	Now
--	-------	------	----------------------	----	------	------	-----------	-----

W=∫Kxdx

On integrating within the limits x=0 to x=x We have

 $W = Kx^2/2$ (17)

- This work done is positive as force is towards the right and spring also moves towards the right
- Same amount of external is done on the spring when it is compressed through a distance x
- Work done as calculated in equation (17) is stored as Potential Energy of the spring. Therefore

 $U = Kx^2/2$ (18)

Elastic potential energy is the potential energy stored by stretching or compressing an elastic object by an external force such as the stretching of a spring. It is equal to the work done to stretch the spring which depends on the spring constant k and the distance stretched.

According to Hooke's law, the force applied to stretch the spring is directly proportional to the amount of stretch.



In other words, Force required to stretch the spring is directly proportional to its displacement. It is given as

$$F = -kx$$

Wherein,

k = spring constant

x = displacement

The Elastic Potential Energy Formula of the spring stretched is given as

$$P.E = \frac{1}{2}kx^2$$

Where,

P.E = elastic potential energy and it's expressed in Joule.

Collision

Collision is short-duration interaction between two bodies or more than two bodies simultaneously causing change in motion of bodies involved due to internal forces acted between them during this. Collisions involve forces (there is a change in velocity). The magnitude of the velocity difference just before impact is called the **closing speed**. All collisions conserve momentum. What distinguishes different types of collisions is whether they also conserve kinetic energy. The Line of impact is the line that is colinear to the common normal of the surfaces that are closest or in contact during impact. This is the line along which internal force of collision acts during impact, and Newton's coefficient of restitution is defined only along this line. Collisions are of three types:

- 1. perfectly elastic collision
- 2. inelastic collision
- 3. perfectly inelastic collision.

Specifically, collisions can either be *elastic*, meaning they conserve both momentum and kinetic energy, or *inelastic*, meaning they conserve momentum but not kinetic energy. An inelastic collision is sometimes also called a *plastic collision*. A "perfectly inelastic" collision (also called a "perfectly plastic" collision) is a limiting case of inelastic collision in which the two bodies coalesce after impact. The degree to which a collision is elastic or inelastic is quantified by the coefficient of restitution, a value that generally ranges between zero and one. A perfectly elastic collision has a coefficient of restitution of one; a perfectly inelastic collision has a coefficient of zero.

Types of collisions

There are two types of collisions between two bodies - 1) Head-on collisions or onedimensional collisions - where the velocity of each body just before impact is along the line of impact, and 2) Non-head-on collisions, oblique collisions or two-dimensional collisions - where the velocity of each body just before impact is not along the line of impact.

According to the coefficient of restitution, there are two special cases of any collision as written below:

1. A perfectly elastic collision is defined as one in which there is no loss of kinetic energy in the collision. In reality, any macroscopic collision between objects will convert some kinetic energy to internal energy and other forms of energy, so no large-scale impacts are perfectly elastic. However, some problems are sufficiently close to perfectly elastic that they can be approximated as such. In this case, the coefficient of restitution equals one.

2. An inelastic collision is one in which part of the kinetic energy is changed to some other form of energy in the collision. Momentum is conserved in inelastic collisions (as it is for elastic collisions), but one cannot track the kinetic energy through the collision since some of it is converted to other forms of energy. In this case, coefficient of restitution does not equal one.

Elastic Collision in 1D

- Consider two particles whose masses are m₁ and m₂ respectively and they collide each other with velocity u₁ and u₂ and after collision their velocities become v₁ and v₂ respectively.
- Collision between these two particles is head on elastic collision. From law of conservation of momentum we have m₁u₁ + m₂u₂ = m₁v₁ + m₂v₂ (1) and from law of conservation of kinetic energy for elastic collision we have

$$\frac{1}{2}m_1\mathbf{u}_1^2 + \frac{1}{2}m_2\mathbf{u}_2^2 = \frac{1}{2}m_1\mathbf{v}_1^2 + \frac{1}{2}m_2\mathbf{v}_2^2$$
(2)

rearranging	equation	1	and	2	we	get
m1(u 1- v 1)=				m ₂	(3)	
and						
$m_1(\mathbf{u}_1^2-\mathbf{v}_1^2)=$	$= m_2(\mathbf{v}_2^2 - \mathbf{u}_2^2)$	(4))			
dividing	equation	4	by	3	we	get

 $u_1 + v_1 = u_2 + v_2$

u₂ - **u**₁ =

 $-(\mathbf{v}_2 - \mathbf{v}_1)$ (5)

where $(\mathbf{u}_2 - \mathbf{u}_1)$ is the relative velocity of second particle w.r.t. first particle before collision and $(\mathbf{v}_2 - \mathbf{v}_1)$ is the relative velocity of second particle w.r.t. first after collision.

From equation 5 we come to know taht in a perfectly elastic collision the magnitude of relative velocity remain unchanged but its direction is reversed. With the help of above equations we can find the values of v₂ and v₁, so from equation 5 v₁ = v₂ - u₁ + u₂ (6)

$$\mathbf{v}_2 = \mathbf{v}_1 + \mathbf{u}_1 - \mathbf{u}_2 \tag{7}$$

Now putting the value of \mathbf{v}_1 from equation 6 in equation 3 we get $m_1(\mathbf{u}_1 - \mathbf{v}_2 + \mathbf{u}_1 - \mathbf{u}_2) = m_2(\mathbf{v}_2 - \mathbf{u}_2)$ On solving the above equation we get value of \mathbf{v}_2 as

$$\mathbf{v}_{2} = \left(\frac{2m_{1}}{m_{1} + m_{2}}\right)\mathbf{u}_{1} + \left(\frac{m_{2} - m_{1}}{m_{1} + m_{2}}\right)\mathbf{u}_{2}$$
(8)

• Similarly putting the value of v_2 from equation 7 in equation 3 we get

$$\mathbf{v}_{1} = \left(\frac{2m_{2}}{m_{1} + m_{2}}\right)\mathbf{u}_{2} + \left(\frac{m_{1} - m_{2}}{m_{1} + m_{2}}\right)\mathbf{u}_{1}$$
(9)

Total kinetic energy of particles before collision is

$$KE_i = \frac{1}{2}m_1\mathbf{u}_1^2 + \frac{1}{2}m_2\mathbf{u}_2^2$$

and total K.E. of particles after collision is

$$KE_f = \frac{1}{2}m_1\mathbf{v}_1^2 + \frac{1}{2}m_2\mathbf{v}_2^2$$

Ratio of initial and final K.E. is

$$\frac{KE_i}{KE_f} = \frac{\frac{1}{2}m_1\mathbf{u}_1^2 + \frac{1}{2}m_2\mathbf{u}_2^2}{\frac{1}{2}m_1\mathbf{v}_1^2 + \frac{1}{2}m_2\mathbf{v}_2^2} = 1$$

• Special cases

Case I: When the mass of both the particles are equal i.e., $m_1 = m_2$ then from equation 8 and 9, $v_2=u_1$ and $v_1=u_2$. Thus if two bodies of equal masses suffer head

on elastic collision then the particles will exchange their velocities. Exchange of momentum between two particles suffering head on elastic collision is maximum when mass of both the particles is same.

Case II: when the target particle is at rest i,e $u_2=0$ From equation (8) and (9)

$$v_{2} = \left(\frac{2m_{1}}{m_{1} + m_{2}}\right)u_{1} \qquad ---(10)$$
$$v_{1} = \left(\frac{m_{1} - m_{2}}{m_{1} + m_{2}}\right)u_{1} \qquad ---(11)$$

Hence some part of the KE which is transformed into second particle would be

$$\frac{\frac{1}{2}m_2v_2^2}{\frac{1}{2}m_1u_1^2} = \frac{\frac{1}{2}m_2\left(\frac{2m_1u_1}{m_1+m_2}\right)^2}{\frac{1}{2}m_1u_1^2}$$
$$= \frac{4m_1m_2}{(m_1+m_2)^2} = \frac{4\frac{m_2}{m_1}}{(1+\frac{m_2}{m_1})^2} \qquad --(12)$$

when $m_1=m_2$, then in this condition $v_0=0$ and $v_2=u_1$ and part of the KE transferred would be

Therefore after collisom first particle moving with initial velocity u_1 would come to rest and the second particle which was at rest would start moving with the velocity of first particle. Hence in this case when $m_1=m_2$ transfer of energy is 100% if $m_1 > m_2$ or $m_1 < m_2$, then energy transformation is not 100%

Case III:

Case IV:

if $m_1 >>> m_2$ and $u_2=0$ then from equation (10) and (11) $v_1 \cong u_1$ and $v_2=2u_1$ (14)

Therefore when a heavy particle collide with a very light particle at rest ,then the heavy particle keeps on moving with the same velocity and the light particle come in motion with a velocity double that of heavy particle

Elastic collision in 2D

Deflection of an moving particle by a particle at rest during perfectly elastic collision in two dimension

• Let m_1 and m_2 be the two mass particle in a laboratory frame of refrence and m_1 collide with m_2 which is initially is at rest.Let the velocity of mass m_1 before collison be $\mathbf{u_1}$ and after the collison it moves with a velocity $\mathbf{v_1}$ and is delfected by the angle θ_1 withs its incident direction and m_2 after the collision moves with the velocity $\mathbf{v_2}$ and it is deflected by an angle θ_2 with its incident direction



Figure 6. Elastic collision in two dimension in laboratory frame of refrance

• From law of conservation of linear momentum , for components along x-axis $m_1u_1=m_1v_1cos\theta_1+m_2v_2cos\theta_2---(1)$

Forcomponentsalongy-axis $0=m_1v_1sin\theta_1 - m_2v_2sin\theta_2$ --(2)

COPYRIGHT FIMT 2020

And from law of conservation of energy $(1/2)m_1u_1^2=(1/2)m_1v_1^2+(1/2)m_2v_2^2$ ---(3)

 Analyzing above equations we come to know that we have to find values of four unknown quantities v₁,v₂,θ₁,θ₂ with the help of above three equations which is not possible. So we cannot predict the variable as they are four of them. However if we measure any one variable ,then other variable can be uniquely determined from the above equation

Inelastic Collision in 1 dimension

Inelastic collision of two particles which stick together

- Consider two particles whose masses are m₁ and m₂.Let u₁ and u₂ be the respective velocities before collision
- Let both the particles stick together after collision and moves with the same velocity v.Then from law of conservation of linear momentum m₁u₁ +m₂u₂=(m₁+m₂)v

or

$$\mathbf{v} = \frac{m_1 \mathbf{u}_1 + m_2 \mathbf{u}_2}{m_1 + m_2} \qquad --(1)$$

• If we consider second particle to be stationary or at rest then $\mathbf{u}_2=0$ then

 $m_1 u_1 = (m_1 + m_2)$

or

$$\mathbf{v} = \frac{m_1 \mathbf{u}_1}{m_1 + m_2} \qquad ---(2)$$

Hence $|v| < |u_1|$

• Kinetic energy before collision is $KE_1=(1/2)m_1u_1^2$

After collison KE of the system is $KE_2=(1/2)(m_1+m_2)v^2$

KE₂ =
$$\frac{1}{2}$$
 (m₁ + m₂) $\left(\frac{m_1 u_1}{m_1 + m_2}\right)^2 = \frac{1}{2} \frac{m_1^2 u_1^2}{(m_1 + m_2)}$ ---(3)

hence

$$\frac{K_2}{K_1} = \frac{\frac{1}{2} \frac{m_1^2 u_1^2}{(m_1 + m_2)}}{\frac{1}{2} m_1 u_1^2} = \frac{m_1}{m_1 + m_2} < 1 \qquad \qquad ---(4)$$

Hence from equation we come to know that $K_2 < K_1$ hence energy loss would be there after the collision of the particles.

Perfectly inelastic collision in one dimension there is a loss of kinetic energy

For collision, you should know two important concept

- 1. Law of conservation of momentum
- 2. work energy theorem,

Perfectly inelastic collision :- when two body move with different speed and after some time they collide and be a single body moves with some speed which is different from both Bodies . This type of collision known as perfectly inelastic collision.

Let Two bodies m_1 and m_2 moves with speed v_1 and v_2 , collide after some interval of time and new bodies of mass ($\underline{m_1} + m_2$) form and it moves with v velocity as shown



now, There is no external force act on Bodies so, linear momentum of bodies are conserved .

so, initial momentum = final momentum

 $m_1v_1 + m_2v_2 = (m_1 + m_2)v - - - - (1)$

Now, change in kinetic energy = final kinetic energy - intial kinetic energy

$$= 1/2(m_1 + m_2)v^2 - 1/2m_1v_1^2 - 1/2m_2v_2^2$$

Putting equation (1)

$$= 1/2(m_1 + m_2)[(m_1v_1 + m_2v_2)/(m_1 + m_2)]^2 - 1/2m_1v_1^2 - 1/2m_2v_2^2$$

= $1/2[m_1v_1 + m_2v_2]^2/(m_1 + m_2)$ - $1/2m_1v_1^2$ - $1/2m_2v_2^2$

= $1/2[(m_1^2v_1^2 + m_2^2v_2^2 + 2m_1m_2v_1v_2 - m_1^2v_1^2 - m_1m_2v_1^2 - m_2^2v_2^2 - m_1m_2v_2^2)]/(m_1+m_2)$

$$= \frac{1}{2} \left[\frac{2m_1m_2v_1v_2 - m_1m_2v_1^2 - m_1m_2v_2^2}{(m_1 + m_2)} \right]$$

$$= -1/2 \{m_1m_2/(m_1+m_2)\}(v_1^2+v_2^2-2v_1v_2)$$

$$= -1/2\mu (v_1 - v_2)^2$$

where $\mu = m_1 m_2 / (m_1 + m_2)$, here negative sign shows that kinetic energy is lost .

RELL

Hence, proved //

Unit 3

Electric charge

Electric charge is the physical property of matter that causes it to experience a force when There placed in an electromagnetic field. are of electric two types charge: positive and negative (commonly carried by protons and electrons respectively). Like charges repel each other and unlike charges attract each other. An object with an absence of net charge is referred to as neutral. Early knowledge of how charged substances interact is now called classical electrodynamics, and is still accurate for problems that do not require consideration of quantum effects.

Electric charge is a conserved property; the net charge of an isolated system, the amount of positive charge minus the amount of negative charge, cannot change. Electric charge is

carried by subatomic particles. In ordinary matter, negative charge is carried by electrons, and positive charge is carried by the protons in the nuclei of atoms. If there are more electrons than protons in a piece of matter, it will have a negative charge, if there are fewer it will have a positive charge, and if there are equal numbers it will be neutral. Charge is *quantized*; it comes in integer multiples of individual small units called the elementary charge, *e*, about 1.602×10^{-19} coulombs,^[11] which is the smallest charge which can exist freely (particles called quarks have smaller charges, multiples of 1/3e, but they are only found in combination, and always combine to form particles with integer charge). The proton has a charge of +e, and the electron has a charge of -e. An electric charge has an electric field, and if the charge is moving it also generates a magnetic field. The combination of the electric and magnetic field is called the electromagnetic field, and its interaction with charges is the source of the electromagnetic force, which is one of the four fundamental forces in physics. The study of photon-mediated interactions among charged particles is called quantum electrodynamics.

The SI derived unit of electric charge is the coulomb (C) named after French physicist Charles-Augustin de Coulomb. In electrical engineering, it is also common to use the ampere-hour (Ah); in physics and chemistry, it is common to use the elementary charge (e as a unit). Chemistry also uses the Faraday constant as the charge on a mole of electrons. The symbol Q often denotes charge.

Electron theory of electrification:

Electrification is the process that produces electric charges on an object. There are many different actions that can produce an electric charge. When this charge is produced on an object it is called an *electrostatic charge*. An electrostatic charge is a charge confined to an object; the charge itself is not moving. It is also referred to as *static electricity*, stationary electricity in the form of an electric charge at rest. (Static tends to make one think of the noise heard on a radio or the crackle of a discharge, but static just means *stationary*. So static electricity is charge that is NOT moving from one object to another. When you hear the radio noise or the crackle described earlier, the charge IS moving. At that point it is a flow of electric charge, not static.)

Charges are produced almost exclusively by the relative numbers of the two primary charged particles in atoms, protons and electrons. A neutral object has a net charge of zero, which means it contains equal numbers of positive and negative charges. *The Basic Law of Electrostatics* is that objects that are similarly charged repel each other; objects that are oppositely charged attract each other. (Same as magnetic poles.) The masses of the three basic subatomic particles are as follows:

AAC ACCREDITE

Electron mass: 9.1094 x 10^{-31} kg

Proton mass: 1.67

1.6726 x 10⁻²⁷ kg

Neutron mass: 1.6749 x 10⁻²⁷ kg

Protons and neutrons are similar in mass, but the mass of an electron is approximately 1/1800 of a proton. Protons and neutrons are called the nucleons and are found in the nucleus. They are bound together by the *strong nuclear force*. The strong nuclear force is a force that is only strongly attractive when the protons and neutrons are very close to one another. This force overcomes the repulsive force that would tend to push the protons apart. (That is why as atoms get larger and have more and more protons, the nuclei also have more and more neutrons. The neutrons add strong nuclear force without adding electrostatic forces.)

While protons and electrons are both charged, net charges are the result of the transfer of electrons only. Changing proton number is a *nuclear* change. This occurs, but it is a different topic than electrical changes. An object that has a positive charge has a net deficiency of electrons. Some electrons have *been removed from* the material if it was originally neutral. An object with a negative charge has a net excess of electrons. The object has *gained* electrons if it was originally neutral.

Frictional electricity:

Frictional electricity is another name for tribo electricity - electricity generated by friction. It's most commonly associated with static electricity. Charge separation often occurs through rolling or sliding contact or collision, usually between different materials. Common examples of frictional electricity and electrostatic charging include combing dry hair or sliding your shoes across a carpet on a dry winter day.

Electrostatic charge separation can even occur through collisions between particles of the same material that have different phases. For example, falling droplets of liquid water and upward-blowing ice crystals within a thunderstorm, create a build-up of positive charge in the upper portion of a thunderstorm cloud, and negative charge in the lower portion of the cloud.

Properties of electric charge

Basic Properties of Electric Charge

As charges are of two types, positive and negative, there are other certain basic properties they follow. If the size of charged bodies is so small, we consider them as point charges. The basic properties of electric charges are as follows:

- Charges are additive in nature
- Charge is a conserved quantity
- Quantization of charge

Charges are additive in nature



Image 3: Adding charges in a system

Charges are additive in nature means they're like scalars and can be added directly. For An **Example** consider a system which consists of two charges namely q_1 and q_2 . Now we wish to find the total charge of the system. The total charge of the system will be the algebraic sum of

 q_1 and q_2 i.e. $q_1 + q_2$. The same thing holds for a number of charges in a system. Let's say a system contains $q_1,q_2,q_3,q_4,\ldots,q_n$, then the net charge of the entire system will be

 $= q_1 + q_2 + q_3 + q_4 + \dots + q_n$

The charge is a scalar quantity as it has only magnitude and no direction. The charge is just as other fundamental properties of the system like mass. The only difference between mass and charge is that charge is both positive and negative, while mass is always positive.

Example:

The charges of a system are +3 C, +2 C, +5 C and -4 C respectively. What would be the net charge of the system?

We know that net charge of a system is algebraic sum of individual charges. Let the total charge of the system be "Q". Then

Q = 3 C + 2 C + 5C - 4C= 6 C

Charge is a conserved quantity



Image 4: Charge is conserved

The charge is a conserved quantity which means charge can neither be created nor be destroyed but can be transferred from one body to another by certain methods like conduction and induction. As charging involves rubbing two bodies, it is actually a transfer of electrons from one body to another. We can't create a charge in a body but eventually can transfer them to another body with some convenient methods.

In a system when charges are distributed accordingly, by the principle of conservation the net charge of the system remains constant. As an example if 5 C is the total charge of the system, then it can be redistributed as 1C, 2C and 2C or in any other possible permutation, but by conservation principle the net charge of system will always be 5 C. Although the charge carriers may be destroyed in a system but the net charge will remain conserved.

Example

Sometimes a neutrino decays to give one electron and one proton by default in nature. The net charge of the system will be zero as electrons and protons are of same magnitude and opposite signs. Then the net charge of the system before the creation of electron and proton (that is zero) equals to a net charge of the system after the creation of electron and proton (which is again zero). This proves the conservation principle.

Quantization of Charge

The charge of one electron is so small that it makes it impractical to use as a unit of electrical charge.

COULOMB is the practical unit adopted for measuring charges.

1 Coulomb = 6.28 x 10¹⁸ Electrons that is: 6,280,000,000,000,000,000

ELECTRONS

Image 5: Number of electrons in 1 Coulomb

Quantization of charge means that charge is a quantized quantity and is expressed as integral multiples of the basic unit of charge (e - charge on one electron). Suppose charge on a body is q, then it can be written as

$\mathbf{q} = \mathbf{n}\mathbf{e}$

where n is an integer and not fraction or irrational number, like 'n' can be any positive or negative integer like 1, 2, 3, -5 etc.

The basic unit of charge is the charge acquired by an electron or proton. By convention we take charge on the electron as negative and denote it as "-e" and charge on a proton is simply "e". The quantization of charge principle was first proposed by English experimentalist Faraday when he put forward his experimental laws of electrolysis. The principle was finally demonstrated and proved by Millikan in 1912.

1 A Coulomb of charge contains around 6×10^{18} electrons. Particles don't have a high magnitude of charge and we use micro coulombs or milli coulombs in order to express charge of a particle.

$1 \ \mu C = 10^{-6} \ C$ $1 \ mC = 10^{-3} \ C$

The principle of quantization can be used to calculate the total amount of charge present in a body and also to calculate a number of electrons or protons in a body.

Suppose a system has n_1 number of electrons and n_2 number of protons, then total amount of charge will be $n_2e - n_1e$.

Coulomb's Law

Coulomb's Law gives an idea about the force between two point <u>charges</u>. By the word point charge, we mean that in physics, the size of linear charged bodies is very small as against the distance between them. Therefore, we consider them as point charges as it becomes easy for us to calculate the force of attraction/ repulsion between them.



Charles-Augustin de Coulomb, a French physicist in 1784, measured the force between two point charges and he came up with the theory that the force is inversely proportional to the square of the distance between the charges. He also found that this force is directly proportional to the product of charges (magnitudes only). We can show it with the following explanation. Let's say that there are two charges q_1 and q_2 . The distance between the charges is 'r', and the force of attraction/repulsion between them is 'F'. Then

$F \propto q_1 q_2$

Or, F $\propto 1/r^2$

$$F = k q_1 q_2 / r^2$$

where k is proportionality constant and equals to $1/4 \pi \epsilon_0$. Here, ϵ_0 is the epsilon naught and it signifies permittivity of a vacuum. The value of k comes $9 \times 10^9 \text{ Nm}^2/\text{ C}^2$ when we take the S.I

unit of value of ε_0 is 8.854 × 10⁻¹² C² N⁻¹ m⁻². According to this theory, like charges repel each other and unlike charges attract each other. This means charges of same sign will push each other with repulsive forces while charges with opposite signs will pull each other with attractive force.

Vector Form of Coulomb's Law

The physical quantities are of two types namely <u>scalars</u> (with the only magnitude) and vectors (those quantities with magnitude and direction). Force is a vector quantity as it has both magnitude and direction. The Coulomb's law can be re-written in the form of vectors. Remember we denote the vector "F" as F, vector r as r and so on. Let there be two charges q_1 and q_2 , with position vectors r_1 and r_2 respectively. Now, since both the charges are of the same sign, there will be a repulsive force between them. Let the force on the q_1 charge due to q_2 be F_{12} and force on q_2 charge due to q_1 charge be F_{21} . The corresponding vector from q_1 to q_2 is r_{21} vector.

$r_{21} = r_2 - r_1$

To denote the direction of a vector from position vector r_1 to r_2 , and from r_2 to r_1 as:

$$\hat{r}_{21} = \frac{r_{21}}{r_{21}} \ . \ \hat{r}_{12} = \frac{r_{12}}{r_{12}} \ . \ \hat{r}_{21} = \ \hat{r}_{12}$$

Now, the force on charge q_2 due to q_1 , in vector form is:

$$F_{21} = \frac{1}{4\pi\epsilon_0} \frac{q_1 q_2}{r_{21}^2} \hat{r}_{21}$$

The above equation is the vector form of Coulomb's Law.

Remarks on Vector Form of Coulomb's Law

While applying Coulomb's Law to find out the force between two point charges, we have to be careful of the following remarks. The vector form of the equation is independent of signs of both the charges, as both the forces are opposite in nature.

The repulsive force F_{12} , that is the force on charge q_1 due to q_2 and another repulsive force F_{21} that is the force on charge q_2 due to q_1 are opposite in signs, due to change in position vector.

 $F_{12} = -F_{21}$

COPYRIGHT FIMT 2020

This is because the position vector in case of force F_{12} is r_{12} and position vector in case of force F_{21} is r_{21} , now

$r_{21} = r_2 - r_1$

$r_{12} = r_1 - r_2$

Since both r_{21} and r_{12} are opposite in signs, they make forces of opposite signs too. This proves that Coulomb's Law fits into Newton's Third Law i.e. every action has its equal and opposite reaction. Coulomb's Law provides the force between two charges when they're present in a vacuum. This is because charges are free in a vacuum and don't get interference from other matter or particles.

Limitations of Coulomb's Law

Coulomb's Law is derived under certain assumptions and can't be used freely like other general formulas. The law is limited to following points:

- We can use the formula if the charges are static (in rest position)
- The formula is easy to use while dealing with charges of regular and smooth shape, and it becomes too complex to deal with charges having irregular shapes
- The formula is only valid when the solvent molecules between the particle are sufficiently larger than both the charges

Superposition Principle:

An electric field is created by every charged particle in the universe in the space surrounding it. The originated field can be calculated with the help of Coulomb's law. The principle of superposition allows for the combination of two or more electric fields.

"The principle of superposition states that every charge in space creates an electric field at point independent of the presence of other charges in that medium. The resultant electric field is a vector sum of the electric field due to individual charges."

In the next section, let us discuss how the superposition principle is applied in electrostatics.

Principle of Superposition in Electrostatics

The superposition principle is helpful when there is large number of charges in a system. Let's consider the following case,



For our convenience let us consider one positive charge, and two negative charges exerting a force on it, from the superposition theorem we know that the resultant force is the vector sum of all the forces acting on the body, therefore the force F_r , the resultant force can be given as follows,

 $Fr \rightarrow = 14\pi \in [Qq1r212r^{12} + Qq2r213r^{13}]$ Where,

r^12 and r^13 are the unit vectors along the direction of q1 and q2.

 \in is the permittivity constant for the medium in which the charges are placed in.

Q, q_1 and q_2 are the magnitude of the charges respectively.

 r_{12} and r_{13} are the distances between the charges Q and q_1 & Q and q_2 respectively.

Continuous Charge Distribution:

We know that the smallest form of charge we can obtain would be +e or -e i.e. the charge of an electron or a proton, hence charges are quantized. <u>Continuous charge distribution</u> means that all charges are closely bound together having very less space between each other.

There are different ways in which charges can be distributed:

- 1. Linear charge distribution.
- 2. Surface charge distribution.
- 3. Volume charge distribution.

Linear charge distribution:

Linear charge distribution is when the charges get distributed uniformly along a length, like around the circumference of a circle or along a straight wire, linear charge distribution is denoted by the symbol λ .

 $\lambda = dq/dl$ and it is measured in Coulombs per meter.

Surface charge distribution:

When a charge is distributed over a specific area, like the surface of a disk, it is called as surface charge distribution, it is denoted by Greek letter σ .

Surface charge distribution is measured Coulombs per square meter or Cm⁻².

Volume charge distribution:

When a charge is distributed uniformly over a volume it is said to be volume charge distribution, like distribution of charge inside a sphere, or a cylinder. It is denoted by ρ .

Volume charge distribution is measured in coulombs per cubic meter or Cm⁻³

Electric field intensity

It was stated that the electric field concept arose in an effort to explain action-at-a-distance forces. All charged objects create an electric field that extends outward into the space that surrounds it. The charge alters that space, causing any other charged object that enters the space to be affected by this field. The strength of the electric field is dependent upon how charged the object creating the field is and upon the distance of separation from the charged object.

The Force per Charge Ratio

Electric field strength is a vector quantity; it has both magnitude and direction. The magnitude of the electric field strength is defined in terms of how it is measured. Let's

suppose that an electric charge can be denoted by the symbol \mathbf{Q} . This electric charge creates an electric field; since \mathbf{Q} is the source of the electric field, we will refer to it as the **source charge**. The strength of the source charge's



electric field could be measured by any other charge placed somewhere in its surroundings.

The charge that is used to measure the electric field strength is referred to as a **test charge** since it is used to *test* the field strength. The test charge has a quantity of charge denoted by the symbol \mathbf{q} . When placed within the electric field, the test charge will experience an electric force - either attractive or repulsive. As is usually the case, this force will be denoted by the symbol \mathbf{F} . The magnitude of the electric field is simply defined as the force per charge on the test charge.

Electric Field Strength = $\frac{Force}{Charge}$

If the electric field strength is denoted by the symbol \mathbf{E} , then the equation can be rewritten in symbolic form as

$\mathbf{E} = \frac{\mathbf{F}}{\mathbf{q}}$

The standard metric units on electric field strength arise from its definition. Since electric field is defined as a force per charge, its units would be force units divided by charge units. In this case, the standard metric units are Newton/Coulomb or N/C.

In the above discussion, you will note that two charges are mentioned - the source charge and the test charge. Two charges would always be necessary to encounter a force. In the electric world, it takes two to attract or repel. The equation for electric field strength (\mathbf{E}) has one of the two charge quantities listed in it. Since there are two charges involved, a student will have to be ultimately careful to use the correct charge quantity when computing the electric field strength. The symbol \mathbf{q} in the equation is the quantity of charge on the test charge (not the source charge). Recall that the electric field strength is defined in terms of how it is measured or tested; thus, the test charge finds its way into the equation. Electric field is the force per quantity of charge *on the test charge*.

The electric field strength is not dependent upon the quantity of charge on the test charge. If you think about that statement for a little while, you might be bothered by it. (Of course if you don't think at all - ever - nothing really bothers you. Ignorance is bliss.) After all, the quantity of charge on the test charge (\mathbf{q}) is in the equation for electric field. So how could electric field strength not be dependent upon \mathbf{q} if \mathbf{q} is in the equation? Good question. But if you think about it a little while longer, you will be able to answer your own question. (Ignorance might be bliss. But with a little extra thinking you might achieve insight, a state much better than bliss.) Increasing the quantity of charge on the test charge - say, by a factor of 2 - would increase the denominator of the equation by a factor of 2. But according to Coulomb's law, more charge also means more electric force (**F**). In fact, a twofold increase in **q** would be accompanied by a twofold increase in **F**. So as the denominator in the equation increases by a factor of two (or three or four), the numerator increases by the same factor. These two changes offset each other such that one can safely say that the electric field strength is not dependent upon the quantity of charge on the test charge. So regardless of what test charge is used, the electric field strength at any given location around the source charge **Q** will be measured to be the same.

Another Electric Field Strength Formula

The above discussion pertained to defining electric field strength in terms of how it is measured. Now we will investigate a new equation that defines electric field strength in terms of the variables that affect the electric field strength. To do so, we will have to revisit the Coulomb's law equation. Coulomb's law states that the electric force between two charges is directly proportional to the product of their charges and inversely proportional to the square of the distance between their centers. When applied to our two charges - the source charge (\mathbf{Q}) and the test charge (\mathbf{q}) - the formula for electric force can be written as

$$\mathbf{F} = \frac{\mathbf{k} \cdot \mathbf{q} \cdot \mathbf{Q}}{\mathbf{d}^2}$$

where $k = 9.0 \cdot 10^9 \text{ N} \cdot \text{m}^2/\text{C}^2$ d = separation distance between charges (meters)

If the expression for electric force as given by Coulomb's law is substituted for force in the above E = F/q equation, a new equation can be derived as shown below.

$$E = \frac{F}{q} = \frac{k \cdot q \cdot Q / d^2}{q} = \frac{k \cdot Q}{d^2}$$
$$E = \frac{k \cdot Q}{d^2}$$

Note that the derivation above shows that the test charge \mathbf{q} was canceled from both numerator and denominator of the equation. The new formula for electric field strength (shown inside the box) expresses the field strength in terms of the two variables that affect it. The electric field strength is dependent upon the quantity of charge on the source charge (\mathbf{Q}) and the distance of separation (\mathbf{d}) from the source charge.

Electric lines of force

We know that, when a unit charge or point charge is placed in the electric field of another charged particle, it will experience a force. The direction of this force can be represented by the imaginary lines. These imaginary lines are called electric lines of force. Electric lines of force are also called as electric field lines. The concept of electric lines of forces was introduced by Michael Faraday in 1837. The direction of electric lines of force for positive and negative charge is shown in the below figure. For positive charge, the electric lines of force move away from the centre of the charge. But in case of negative charge, the electric lines of force move towards the centre of the charge.



Opposite charges attract and like charges repel Opposite charges attract

If two opposite charges are placed close to each other, they get attracted because the force present between them is attractive.

Let us consider two apposite charges as shown in below figure. Below figure clearly shows that for positive charge the electric lines of force moves away from the centre of positive charge and for negative charge the electric lines of force moves towards the centre of the negative charge.



COPYRIGHT FIMT 2020

If these two opposite charges are placed close to each other, the positive charge moves in the direction of electric lines of force and enters into electric field of negative charge. Here, the positive charge gets pulled towards the negative charge because the electric lines of force for negative charge are also in the same direction. Therefore, the two opposite charges get attracted.

Like charges repel

If two positive charges are placed close to each other, they get repelled because the force present between them is repulsive. Let us consider two positive charges as shown in below figure. Below figure clearly shows that for both the positive charges electric lines of force moves away from the centre of the positive charges. If these two positive charges are placed close to each other, both the charges will try to move in the direction of electric lines of force. The positive charge at left side will try to move towards the positive charge at right side, but the electric lines of force of the right side positive charge oppose this movement. In the similar way, positive charge at right side will also experience a opposing force from left side positive charge. Hence, both the charges will experience a repulsive force from each other.



If two negative charges are placed close to each other, they get repelled because the force present between them is repulsive. Let us consider two negative charges as shown in below figure. Below figure clearly shows that for both the negative charges electric lines of force moves towards the centre of the negative charges. If these two negative charges are placed close to each other, both the charges will try to move in the direction of electric lines of force. The negative charge at left side will experience a pulling force from the right side negative charge, but the electric lines of force for the left negative charge is in opposite direction. Hence, it will moves away from the right side negative charge. In the similar way,

the right side negative charge will try to moves away from left side negative charge. Therefore, both the charges will move away from each other.



Properties of electric lines of force

1) The electric lines of force start from a positive charge and ends on a negative charge. 2) The electric lines of force always enter or leave the charged surface normally. lines of 3) Electric force can never intersect each other. electric lines of 4) The force conductor. cannot pass through а 5) When two opposite charges are placed close to each other, the electric lines of force between them will become shorten present in length. 6) When two like charges are placed closer to each other, the electric lines of force present between them will become enlarged in length.

Electrostatics

Electrostatics, as the name implies, is the study of stationary electric charges. A rod of plastic rubbed with fur or a rod of glass rubbed with silk will attract small pieces of paper and is said to be **electrically charged**. The charge on plastic rubbed with fur is defined as **negative**, and the charge on glass rubbed with silk is defined as **positive**.



Line integral of electric field

Area Integral



An area integral of a vector function E can be defined as the <u>integral</u> on a surface of the <u>scalar product</u> of E with area element dA. The direction of the area element is defined to be perpendicular to the area at that point on the surface.

The outward directed surface integral over an entire closed surface

is denoted

 $\oint \vec{E} \cdot \vec{dA}$



It is appropriate for such physical applications as Gauss' law.

Line Integral



Vector functions such as electric field and magnetic field occur in physical applications, and scalar products of these vector functions with another vector such as distance or path length appear with regularity. When such a product is summed over a path length where the magnitudes and directions change, that sum becomes an integral called a line integral.

$$\int_{A}^{B} \vec{E} \cdot \vec{ds} = \int_{A}^{B} E \cos \theta ds$$

A line integral is also used for the general definition of work in mechanics.

Applications of Line Integrals

The line integral of electric field around a closed loop is equal to the voltage generated in that loop (Faraday's law):

COPYRIGHT FIMT 2020

$$\oint \vec{E} \cdot \vec{ds} = -\frac{d\Phi_{\rm B}}{dt}$$

Such an integral is also used for the calculation of voltage difference since voltage is work per unit charge. Calculating the voltage difference near a point charge is a good example.

The line integral of a force over a path is equal to the work done by that force on the path.

MANAG

$$W_{ab} = \int_{a}^{b} \vec{F} \cdot \vec{ds}$$

Electrostatic potential

An electric potential (also called the *electric field potential*, potential drop or the electrostatic potential) is the amount of work needed to move a unit of charge from a reference point to a specific point inside the field without producing an acceleration. Typically, the reference point is the Earth or a point at infinity, although any point can be used. In classical electrostatics, the electrostatic field is a vector quantity which is expressed as the gradient of the electrostatic potential energy of any charged particle at any location (measured in joules) divided by the charge of that particle (measured in coulombs). By dividing out the charge on the particle a quotient is obtained that is a property of the electric field itself.

This value can be calculated in either a static (time-invariant) or a dynamic (varying with time) electric field at a specific time in units of joules per coulomb (J C^{-1}), or volts (V). The electric potential at infinity is assumed to be zero. In electrodynamics, when time-varying fields are present, the electric field cannot be expressed only in terms of a scalar potential. Instead, the electric field can be expressed in terms of both the scalar electric potential and the magnetic vector potential. The electric potential and the magnetic vector potential. The electric potential are mixed under Lorentz transformations. Practically, electric potential is always a continuous function in space; Otherwise, the special derivative of it will yield a field with infinite magnitude, which is practically impossible. Even an idealized Point Charge has potential, which is continuous everywhere except the origin; The Electric field across an idealized Surface charge is
not continuous, but it's not infinite at any point. Therefore, the electric potential is continuous across an idealized Surface charge. An idealized Linear Charge has potential, which is continuous everywhere except on the Linear Charge.

Electric potential due to a point charge

• Consider a positive test charge +q is placed at point O shown below in the figure.



- We have to find the electric potential at point P at a distance r from point O.
- If we move a positive test charge q' from infinity to point P then change in electric potential energy would be

$$U_P - U_{\infty} = \frac{qq'}{4\pi\varepsilon_0 r}$$

• Electric potential at point P is

$$V_p = \frac{U_p - U_a}{q'} = \frac{q}{4\pi\varepsilon_0 r}$$
(8)

• Potential V at any point due to arbitrary collection of point charges is given by

$$V = \frac{1}{4\pi\varepsilon_0} \sum_{i=1}^n \frac{q_i}{r_i} \tag{9}$$

- here we see that like electric field potential at any point independent of test charge used to define it.
- For continous charge distributions summation in above expressin will be replaced by the integration

$$V = \frac{1}{4\pi\varepsilon_0} \int \frac{dq}{r} \tag{10}$$

where dq is the differential element of charge distribution and r is its distance from the point at which V is to be calculated.

Potential due to an electric dipole

- We already know that electric dipole is an arrangement which consists of two equal and opposite charges +q and -q separated by a small distance 2a.
- Electric dipole moment is represented by a vector **p** of magnitude 2qa and this vector points in direction from -q to +q.
- To find electric potential due to a dipole consider charge -q is placed at point P and charge +q is placed at point Q as shown below in the figure.



• Since electric potential obeys superposition principle so potential due to electric dipole as a whole would be sum of potential due to both the charges +q and -q. Thus

$$V = \frac{1}{4\pi\varepsilon_0} \left(\frac{q}{r_1} - \frac{q}{r_2} \right) \tag{15}$$

where r_1 and r_2 respectively are distance of charge +q and -q from point R.

•	Now draw line I	PC perpandicular	to RO and line	QD perpandicu	lar to RO as
shown	in	figure.	From	triangle	POC
cosθ=OC/OP =					OC/a
therefo	ore	OC=acosθ	similar	ly	OD=acosθ
Now					,
r ₁ =	QR≅RD	=	OR-OD	=	r-acosθ
r ₂ =	PR≅RC	=	OR+OC	=	r+acosθ

$$V = \frac{q}{4\pi\varepsilon_0} \left(\frac{1}{r - a\cos\theta} - \frac{1}{r + a\cos\theta} \right) = \frac{q}{4\pi\varepsilon_0} \left(\frac{2a\cos\theta}{r^2 - a^2\cos^2\theta} \right)$$

since magnitude of dipole is
$$|\mathbf{p}| = 2qa$$

COPYRIGHT FIMT 2020

$$V = \frac{1}{4\pi\varepsilon_0} \left(\frac{p\cos\theta}{r^2 - a^2\cos^2\theta} \right)$$
(16)



$$V = \frac{p \cos \theta}{4\pi\varepsilon_0 r^2} \tag{17}$$

again since $p\cos\theta = \mathbf{p}\cdot\mathbf{r}^{*}$ where, \mathbf{r}^{*} is the unit vector along the vector OR then electric potential of dipole is

$$V = \frac{\boldsymbol{p} \cdot \hat{\boldsymbol{r}}}{4\pi\varepsilon_0 r^2} \tag{18}$$

for r>>a

- From above equation we can see that potential due to electric dipole is inversly proportional to r² not ad 1/r which is the case for potential due to single charge.
- Potential due to electric dipole does not only depends on r but also depends on angle between position vector r and dipole moment p.

Work done in rotating an electric dipole in an electric field

- Consider a dipole placed in a uniform electric field and it is in equilibrium position. If we rotate this dipole from its equilibrium position, work has to be done.
- Suppose electric dipole of moment p is rotated in uniform electric field E through an angle θ from its equilibrium position. Due to this rotation couple acting on dipole changes.
- If at any instant dipole makes an angle φ with uniform electric field then torque acting on dipole is

```
Г=pEsinф
```

again work done in rotating this dipole through an infitesimaly small angle dφ is dW=torque x angular displacement =pEsinφdφ

Total work done in rotating the dipole through an angle θfrom its equilibrium position

$$W = \int_{0}^{\theta} pE\sin\varphi d\varphi = pE[-\cos\varphi]_{0}^{\theta} = pE(1-\cos\theta)$$
(21)

This is the required formula for work done in rotating an electric dipole placed in uniform electric field through an angle θ from its equilibrium position.

Potential energy of dipole placed in uniform electric field

 Again consider equation 20 which gives the work done in rotating electric dipole through an infinetesimly small angle dφ is dW=pEsinφdφ

which is equal to the change in potential energy of the system $dW=dU=pEsin\phi d\phi$

• If angle $d\phi$ is changed from 90⁰ to θ then in potential energy would be

$$W = \int_{0}^{\theta} pE \sin \varphi d\varphi = pE(-\cos\varphi) \int_{0}^{\theta} = pE(1-\cos\theta)$$

• We have choosen the value of ϕ going from $\pi/2$ to θ because at $\pi/2$ we can take potential energy to be zero (axis of dipole is perpandicular to the field). Thus $U(90^0)=0$ and above equation becomes

$$U(\theta) - U(90^{\circ}) = \int_{90^{\circ}}^{\theta} pE \sin \varphi d\varphi = pE(-\cos \varphi) \Big|_{90^{\circ}}^{\theta} = -pE \cos \theta = -p \cdot E$$

Gauss's law



Gauss's law states that the net flux of an <u>electric field</u> in a closed surface is directly proportional to the enclosed electric charge. It is one of the four equations of Maxwell's laws of electromagnetism. It was initially formulated by Carl Friedrich Gauss in the year 1835 and relates the electric fields at the points on a closed surface and the net charge enclosed by that surface. The electric flux is defined as the electric field passing through a given area multiplied by the area of the surface in a plane perpendicular to the field. Yet another statement of Gauss's law states that the net flux of a given electric field through a given surface, divided by the enclosed charge should be equal to a constant. Usually, a positive electric charge is supposed to generate a positive electric field. The law was released in 1867 as part of a collection of work by the famous German mathematician, Carl Friedrich Gauss.

Gauss Law Equation

Let us now study Gauss's law through an integral equation. Gauss's law in integral form is given below:

$$\int \mathbf{E} \cdot d\mathbf{A} = \mathbf{Q}/\varepsilon_0 \qquad \dots \qquad (1)$$

Where,

- **E** is the electric field vector
- Q is the enclosed electric charge
- ε₀ is the electric permittivity of free space
- A is the outward pointing normal area vector

Flux is a measure of the strength of a field passing through a surface. Electric flux is defined as

 $\Phi = \int \mathbf{E} \cdot d\mathbf{A} \quad \dots (2)$

We can understand the electric field as flux density. Gauss's law implies that the net electric flux through any given closed surface is zero unless the volume bounded by that surface contains a net charge.

Gauss's law for electric fields is most easily understood by neglecting electric displacement (d). In matters, the dielectric permittivity may not be equal to the permittivity of free-space (i.e. $\epsilon \neq \epsilon_0$). In the matter, the density of electric charges can be separated into a "free" charge density (ρ_f) and a "bounded" charge density (ρ_b), such that:

 $P = \rho_f + \rho_b$

Statement of Gauss"s Theorem :

The net-outward normal electric flux through any closed surface of any shape is equal to $1/\epsilon_0$ times the total charge contained within that surface, i.e.,

 $\int_{S} \vec{E} \cdot d\vec{S} = \frac{1}{\varepsilon_{0}} \sum q$ Where \int_{S} indicates the surface integral

over the whole of the closed surface, q is the algebraic sum of all the charges (i.e., net charge in coulombs) enclosed by surface S.



Proof of Gauss"s Theorem:

Let a point charge +q coulomb be placed at O within the closed surface. Let E be the electric field strength at P. Let

OP= r and the permittivity of free space or vaccuum be ε_0 .



Consider a small area $d\vec{S}$ of the surface surrounding the point P. Then the electric flux through $d\vec{S}$ is given by $d\phi = \vec{E} \cdot d\vec{S}$

But the electric field strength at P, $E = \frac{1}{4\pi\epsilon_n} \cdot \frac{q}{r^2} \hat{r} = \frac{1}{4\pi\epsilon_n} \cdot \frac{q\dot{r}}{r^3}$

Since unit vector $\hat{\mathbf{f}} = \frac{\hat{\mathbf{f}}}{\hat{\mathbf{f}}}$ $d\phi = \vec{E} \cdot d\vec{S} = \frac{1}{4\pi\epsilon_0} q \frac{\vec{r} \cdot d\vec{S}}{r^3} = \frac{q}{4\pi\epsilon_0} d\Omega$ where $d\Omega = \frac{\vec{r} \cdot d\vec{S}}{r^3} = \frac{dS\cos\theta}{r^2}$ is the solid angle subtended by area dS at point

O. Here Θ is the angle between $d\vec{S}$ and \vec{E} Hence electric flux through whole of closed surface.

$$\phi = \oint \vec{E} \bullet d\vec{S} = \frac{q}{4\pi\epsilon_0} \times \oint d\Omega$$

But $\oint d\Omega$ is the solid angle due to the entire closed surface S at an internal point $\odot=4\pi$

$$\therefore \phi = \frac{q}{4\pi\epsilon_*} \cdot 4\pi = \frac{1}{\epsilon_*} q$$

If there are several charges, $+q_1$, $+q_2$, q_3 ,...... $-q_1$ ', $-q_2$ ', $-q_3$ '..... inside the closed surface, each will contribute to the total electric flux. For positive charges the flux will be outward and hence positive; for negative charges the flux will be inward and negative. Therefore, the total electric flux in such a case $= \frac{1}{\varepsilon_0}q_1 + \frac{1}{\varepsilon_0}q_2 + \frac{1}{\varepsilon_0}q_3.... - \frac{1}{\varepsilon_0}q_1' - \frac{1}{\varepsilon_0}q_2' - \frac{1}{\varepsilon_0}q_3'.....$ $= \frac{1}{\varepsilon_0}(q_1 + q_2 + q_3 + - q_1' - q_2' - q_3'....) = \frac{1}{\varepsilon_0}\sum q_1$ where $\sum q$ is the algebraic sum of the charges within the closed surface.

Hence total electric flux through any closed surface is equal to 1/2, times the total sherze (in equilarly) enclosed within the surface which is Gauss's law

to $1/\epsilon_0$ times the total charge (in coulomb) enclosed within the surface which is Gauss's law.

Capacitors

Capacitors are simple passive device that can store an electrical charge on their plates when connected to a voltage source. The capacitor is a component which has the ability or "capacity" to store energy in the form of an electrical charge producing a potential difference (*Static Voltage*) across its plates, much like a small rechargeable battery.

There are many different kinds of capacitors available from very small capacitor beads used in resonance circuits to large power factor correction capacitors, but they all do the same thing, they store charge. In its basic form, a capacitor consists of two or more parallel conductive (metal) plates which are not connected or touching each other, but are electrically separated either by air or by some form of a good insulating material such as waxed paper, mica, ceramic, plastic or some form of a liquid gel as used in electrolytic capacitors. The insulating layer between a capacitors plates is commonly called the **Dielectric**.



A Typical Capacitor

Due to this insulating layer, DC current cannot flow through the capacitor as it blocks it allowing instead a voltage to be present across the plates in the form of an electrical charge. The conductive metal plates of a capacitor can be either square, circular or rectangular, or they can be of a cylindrical or spherical shape with the general shape, size and construction of a parallel plate capacitor depending on its application and voltage rating. When used in a direct current or DC circuit, a capacitor charges up to its supply voltage but blocks the flow of current through it because the dielectric of a capacitor is non-conductive and basically an insulator. However, when a capacitor is connected to an alternating current or AC circuit, the flow of the current appears to pass straight through the capacitor with little or no resistance.

There are two types of electrical charge, a positive charge in the form of Protons and a negative charge in the form of Electrons. When a DC voltage is placed across a capacitor, the positive (+ve) charge quickly accumulates on one plate while a corresponding and opposite negative (-ve) charge accumulates on the other plate. For every particle of +ve charge that arrives at one plate a charge of the same sign will depart from the -ve plate. Then the plates remain charge neutral and a potential difference due to this charge is established between the two plates. Once the capacitor reaches its steady state condition an electrical current is unable to flow through the capacitor itself and around the circuit due to the insulating properties of the dielectric used to separate the plates.

The flow of electrons onto the plates is known as the capacitors **Charging Current** which continues to flow until the voltage across both plates (and hence the capacitor) is equal to the applied voltage Vc. At this point the capacitor is said to be "fully charged" with electrons. The strength or rate of this charging current is at its maximum value when the plates are fully discharged (initial condition) and slowly reduces in value to zero as the plates charge up to a potential difference across the capacitors plates equal to the source voltage. The amount of potential difference present across the capacitor depends upon how much charge was deposited onto the plates by the work being done by the source voltage and also by how much capacitance the capacitor has and this is illustrated below.



The parallel plate capacitor is the simplest form of capacitor. It can be constructed using two metal or metallised foil plates at a distance parallel to each other, with its capacitance value in Farads, being fixed by the surface area of the conductive plates and the distance of separation between them. Altering any two of these values alters the the value of its capacitance and this forms the basis of operation of the variable capacitors. Also, because capacitors store the energy of the electrons in the form of an electrical charge on the plates the larger the plates and/or smaller their separation the greater will be the charge that the capacitor holds for any given voltage across its plates. In other words, larger plates, smaller distance, more capacitance.

By applying a voltage to a capacitor and measuring the charge on the plates, the ratio of the charge Q to the voltage V will give the capacitance value of the capacitor and is therefore given as: C = Q/V this equation can also be re-arranged to give the familiar formula for the quantity of charge on the plates as: $Q = C \times V$. Although we have said that the charge is stored on the plates of a capacitor, it is more exact to say that the energy within the charge is stored in an "electrostatic field" between the two plates. When an electric current flows into the capacitor, it charges up, so the electrostatic field becomes much stronger as it stores more energy between the plates.

Likewise, as the current flowing out of the capacitor, discharging it, the potential difference between the two plates decreases and the electrostatic field decreases as the energy moves out of the plates. The property of a capacitor to store charge on its plates in the form of an electrostatic field is called the **Capacitance** of the capacitor. Not only that, but capacitance is also the property of a capacitor which resists the change of voltage across it.

The Capacitance of a Capacitor

Capacitance is the electrical property of a capacitor and is the measure of a capacitors ability to store an electrical charge onto its two plates with the unit of capacitance being the **Farad** (abbreviated to F) named after the British physicist Michael Faraday.

Capacitance is defined as being that a capacitor has the capacitance of **One Farad** when a charge of **One Coulomb** is stored on the plates by a voltage of **One volt**. Note that capacitance, C is always positive in value and has no negative units. However, the Farad is very large units of measurement to use on its own so sub-multiples of the Farad are generally used such as micro-farads, nano-farads and pico-farads, for example.

Standard Units of Capacitance

- Microfarad (μF) 1μF = 1/1,000,000 = 0.000001 = 10⁻⁶ F
- Nanofarad (nF) 1nF = 1/1,000,000,000 = 0.000000001 = 10⁻⁹ F
- Picofarad (pF) 1pF = 1/1,000,000,000,000 = 0.00000000001 = 10⁻¹² F

Then using the information above we can construct a simple table to help us convert between pico-Farad (pF), to nano-Farad (nF), to micro-Farad (μ F) and to Farads (F) as shown.

Capacitance of a Parallel Plate Capacitor

The capacitance of a parallel plate capacitor is proportional to the area, A in metres² of the smallest of the two plates and inversely proportional to the distance or separation, d (i.e. the dielectric thickness) given in metres between these two conductive plates. The generalised the capacitance of equation for a parallel plate capacitor is given as: $C = \varepsilon(A/d)$ where ε represents the absolute permittivity of the dielectric material being used. The permittivity of a vacuum, ε_0 also known as the "permittivity of free space" has the value of the constant 8.84 x 10^{-12} Farads per metre.

To make the maths a little easier, this dielectric constant of free space, ε_0 , which can be written as: $1/(4\pi \times 9 \times 10^9)$, may also have the units of picofarads (pF) per metre as the constant giving: 8.84 for the value of free space. Note though that the resulting capacitance value will be in picofarads and not in farads.

Generally, the conductive plates of a capacitor are separated by some kind of insulating material or gel rather than a perfect vacuum. When calculating the capacitance of a capacitor, we can consider the permittivity of air, and especially of dry air, as being the same value as a vacuum as they are very close.



Capacitance Example No1

A capacitor is constructed from two conductive metal plates 30cm x 50cm which are spaced 6mm apart from each other, and uses dry air as its only dielectric material. Calculate the capacitance of the capacitor.

Using:
$$C = \varepsilon_0 \frac{A}{d}$$

where: $\varepsilon_0 = 8.84 \times 10^{-12}$
 $A = 0.3 \times 0.5 m^2$ and $d = 6 \times 10^{-3} m$
 $C = \frac{8.84 \times 10^{-12} \times (0.3 \times 0.5)}{6 \times 10^{-3}} = 0.221 nF$

Then the value of the capacitor consisting of two plates separated by air is calculated as 221pF or 0.221nF

The Dielectric of a Capacitor

As well as the overall size of the conductive plates and their distance or spacing apart from each other, another factor which affects the overall capacitance of the device is the type of dielectric material being used. In other words the "Permittivity" (ϵ) of the dielectric.

The conductive plates of a capacitor are generally made of a metal foil or a metal film allowing for the flow of electrons and charge, but the dielectric material used is always an insulator. The various insulating materials used as the dielectric in a capacitor differ in their ability to block or pass an electrical charge. This dielectric material can be made from a number of insulating materials or combinations of these materials with the most common types used being: air, paper, polyester, polypropylene, Mylar, ceramic, glass, oil, or a variety of other materials. The factor by which the dielectric material, or insulator, increases the capacitance of the capacitor compared to air is known as the **Dielectric Constant**, **k** and a dielectric material with a high dielectric constant is a better insulator than a dielectric material with a lower dielectric constant. Dielectric constant is a dimensionless quantity since it is relative to free space.

The actual permittivity or "complex permittivity" of the dielectric material between the plates is then the product of the permittivity of free space (ε_0) and the relative permittivity (ε_r) of the material being used as the dielectric and is given as:

Complex Permittivity

$\boldsymbol{\varepsilon} = \boldsymbol{\varepsilon}_{o} \times \boldsymbol{\varepsilon}_{r}$

In other words, if we take the permittivity of free space, ε_0 as our base level and make it equal to one, when the vacuum of free space is replaced by some other type of insulating material, their permittivity of its dielectric is referenced to the base dielectric of free space giving a multiplication factor known as "relative permittivity", ε_r . So the value of the complex permittivity, ε will always be equal to the relative permittivity times one. Typical units of dielectric permittivity, ε or dielectric constant for common materials are: Pure Vacuum = 1.0000, Air = 1.0006, Paper = 2.5 to 3.5, Glass = 3 to 10, Mica = 5 to 7, Wood = 3 to 8 and Metal Oxide Powders = 6 to 20 etc. This then gives us a final equation for the capacitance of a capacitor as:

Capacitance, C = $\frac{\varepsilon_0 \varepsilon_r A}{d}$ Farads

One method used to increase the overall capacitance of a capacitor while keeping its size small is to "interleave" more plates together within a single capacitor body. Instead of just one set of parallel plates, a capacitor can have many individual plates connected together thereby increasing the surface area, A of the plates. For a standard parallel plate capacitor as shown above, the capacitor has two plates, labelled A and B. Therefore as the number of capacitor plates is two, we can say that n = 2, where "n" represents the number of plates.

Then our equation above for a single parallel plate capacitor should really be:

Capacitance, C = $\frac{\varepsilon_0 \varepsilon_r (n-1)A}{d}$ Farads

However, the capacitor may have two parallel plates but only one side of each plate is in contact with the dielectric in the middle as the other side of each plate forms the outside of the capacitor. If we take the two halves of the plates and join them together we effectively only have "one" whole plate in contact with the dielectric.

As for a single parallel plate capacitor, n - 1 = 2 - 1 which equals 1 as $C = (\varepsilon_0 * \varepsilon_r x \ 1 \ x \ A)/d$ is exactly the same as saying: $C = (\varepsilon_0 * \varepsilon_r * A)/d$ which is the standard equation above.

Now suppose we have a capacitor made up of 9 interleaved plates, then n = 9 as shown.

Multi-plate Capacitor



Now we have five plates connected to one lead (A) and four plates to the other lead (B). Then BOTH sides of the four plates connected to lead B are in contact with the dielectric, whereas only one side of each of the outer plates connected to A is in contact with the dielectric. Then as above, the useful surface area of each set of plates is only eight and its capacitance is therefore given as:

$$C = \frac{\varepsilon_{o}\varepsilon_{r}(n-1)A}{d} = \frac{\varepsilon_{o}\varepsilon_{r}(9-1)A}{d} = \frac{\varepsilon_{o}\varepsilon_{r}8A}{d}$$

Modern capacitors can be classified according to the characteristics and properties of their insulating dielectric:

- Low Loss, High Stability such as Mica, Low-K Ceramic, Polystyrene.
- Medium Loss, Medium Stability such as Paper, Plastic Film, High-K Ceramic.
- Polarized Capacitors such as Electrolytic's, Tantalum's.

Voltage Rating of a Capacitor

All capacitors have a maximum voltage rating and when selecting a capacitor consideration must be given to the amount of voltage to be applied across the capacitor. The maximum amount of voltage that can be applied to the capacitor without damage to its dielectric material is generally given in the data sheets as: WV, (working voltage) or as WV DC, (DC working voltage). If the voltage applied across the capacitor becomes too great, the dielectric will break down (known as electrical breakdown) and arcing will occur between the capacitor plates resulting in a short-circuit. The working voltage of the capacitor depends on the type of dielectric material being used and its thickness. The DC working voltage of a capacitor is just that, the maximum DC voltage and NOT the maximum AC voltage as a capacitor with a DC voltage rating of 100 volts DC cannot be safely subjected to an alternating voltage of 100

volts. Since an alternating voltage that has an RMS value of 100 volts will have a peak value of over 141 volts! ($\sqrt{2} \times 100$).

Then a capacitor which is required to operate at 100 volts AC should have a working voltage of at least 200 volts. In practice, a capacitor should be selected so that its working voltage either DC or AC should be at least 50 percent greater than the highest effective voltage to be applied to it. Another factor which affects the operation of a capacitor is **Dielectric Leakage**. Dielectric leakage occurs in a capacitor as the result of an unwanted leakage current which flows through the dielectric material. Generally, it is assumed that the resistance of the dielectric is extremely high and a good insulator blocking the flow of DC current through the capacitor (as in a perfect capacitor) from one plate to the other.

However, if the dielectric material becomes damaged due excessive voltage or over temperature, the leakage current through the dielectric will become extremely high resulting in a rapid loss of charge on the plates and an overheating of the capacitor eventually resulting in premature failure of the capacitor. Then never use a capacitor in a circuit with higher voltages than the capacitor is rated for otherwise it may become hot and explode.

Capacitors in series and parallel combinations

For practical applications, two or more capacitors are often used in combination and their total capacitance C must be known. To find total capacitance of the arrangement of capacitor we would use equation

Q=CV



• Figure below shows two capacitors connected in parallel between two points A and



 Right hand side plate of capacitors would be at same common potential V_A. Similarly left hand side plates of capacitors would also be at same common potential V_B. Thus in this case potential difference V_{AB}=V_A-V_B would be same for both the capacitors, and charges Q₁ and Q₂ on both the capacitors are not necessarily equal. So,

 $Q_1=C_1V$ and $Q_2=C_2V$

- Thus charge stored is divided amongst both the capacitors in direct proportion to their capacitance.
- Total charge on both the capacitors is, $\label{eq:Q2} Q=Q_1+Q_2=V(C_1+C_2)$

and

 $Q/V=C_1+C_2$ So system is equivalent to a single capacitor of capacitance C=Q/V

where,

- When capacitors are connected in parallel their resultant capacitance C is the sum of their individual capacitances.
- The value of equivalent capacitance of system is greater than the greatest individual one.
- If there are number of capacitors connected in parallel then their equivalent capacitance would be C=C₁+C₂+ C₃........... (10)
- (ii) Series combination of capacitors
- Figure 7 below shows two capacitors connected in series combination between points
 A
 and
 B.



- Both the points A and B are maintained at constant potential difference V_{AB}.
- In series combination of capacitors right hand plate of first capacitor is connected to left hand plate of next capacitor and combination may be extended for any number of capacitors.
- In series combination of capacitors all the capacitors would have same charge.

- Now potential difference across individual capacitors are given by $V_{AR}{=}Q/C_1$

and,

 $V_{RB}=Q/C_2$

• Sum of V_{AR} and V_{RB} would be equal to applied potential difference V so, $V=V_{AB}=V_{AR}+V_{RB}=Q(1/C_1 + 1/C_2)$

or,

$$\frac{V}{Q} = \frac{1}{C_1} + \frac{1}{C_2} = \frac{Q}{C}$$

where

$$\frac{1}{C} = \frac{1}{C_1} + \frac{1}{C_2}$$

i.e., resultant capacitance of series combination C=Q/V, is the ratio of charge to total potential difference across the two capacitors connected in series.

- So, from equation 12 we say that to find resultant capacitance of capacitors connected in series, we need to add reciprocals of their individual capacitances and C is always less then the smallest individual capacitance.
- Result in equation 12 can be summarized for any number of capacitors i.e.,

12 1- 1. 24

$$\frac{1}{C} = \frac{1}{C_1} + \frac{1}{C_2} + \frac{1}{C_3} + \dots$$

Energy stored in a capacitor

- Consider a capacitor of capacitance C, completely uncharged in the begging.
- Charming process of capacitor requires expenditure of energy because while charging a capacitor charge is transferred from plate at lower potential to plate at higher potential.
- Now if we start charging capacitor by transporting a charge dQ from negative plate ti the positive plate then work is done against the potential difference across the plate.
- If q is the amount of charge on the capacitor at any stage of charging process and φ is the potential difference across the plates of capacitor then magnitude of potential difference is φ=q/C.

- Now work dW required to transfer dq is dW=\phidq=qdq/C
- To charge the capacitor starting from the uncharged state to some final charge Q work required is Integrating from 0 to Q

W=(1/C)∫qdq =(Q²)/2C

(14a)

=(CV²)/2

=QV/2

Which is the energy stored in the capacitor and can also be written as $U=(CV^2)/2$ ---(15)

- From equation 14c, we see that the total work done is equal to the average potential V/2 during the charging process , multiplied by the total charge transferred
- If C is measured in Farads ,Q in coulumbs and V in volts the energy stored would in Joules
- A parallel plate capacitor of area A and separation d has capacitance $C=\epsilon_0 A/d$
- electric field in the space between the plates is E=V/d V=Ed or Putting above values of and С in equation 14b find V we $W=U=(1/2)(\epsilon_0 A/d)(Ed)^2$

 $=(1/2)\varepsilon_0 E^2(Ad)$

 $=(1/2)\varepsilon_0 E^2.V ---(16)$

- If u denotes the energy per unit volume or energy density then $u=(1/2)\epsilon_0 E^2 \, x \, \text{volume}$
- The result for above equation is generally valid even for electrostatic field that is not constant in space.

Effect of Dielectric on Capacitance

Dielectrics

Dielectrics are basically insulating and non-conducting substances. They are bad conductors of electric current. Dielectrics are capable of holding electrostatic charges while emitting minimal energy. This energy is usually in the form of heat. The common examples of dielectrics include mica, plastics, porcelain, metal oxides and glass etc. It is important for you to note that dry air is also a dielectric.

What is the Dielectric Constant?

When we put a dielectric slab in between two plates of a parallel plate capacitor, the ratio of the applied electric field strength to the strength of the reduced value of electric field capacitor is called the dielectric constant. It is given as

$K = E_0/E$

 E_0 is greater than or equal to E, where E_0 is the field with the slab and E is the field without it. The larger the dielectric constant, the more charge can be stored. Completely filling the space between capacitor plates with a dielectric, increases the capacitance by a factor of the dielectric constant:

$C = KC_{o}$,

Where C_0 is the capacitance with no slab between the plates. This is all about a quick recap. Now let us move ahead and see what effect dielectrics have on the capacitance.

Effect of Dielectric on Capacitance

We usually place dielectrics between the two plates of parallel plate capacitors. They can fully or partially occupy the region between the plates. When we place the dielectric between the two plates of a parallel plate capacitor, the electric field polarises it. The surface charge densities are σ_p and $-\sigma_p$. When we place the dielectric fully between the two plates of a capacitor, then it's



The electric field inside a capacitor is as follows:

$$E = \frac{\sigma - \sigma_E}{\varepsilon_0}$$

Hence we have:

$$V=\!\frac{\sigma d}{\epsilon_0 k}\!=\!\frac{Q d}{A \epsilon_0 k}$$

Therefore:

$$C=\frac{Q}{V}=\frac{A\epsilon_0k}{d}=\frac{A\epsilon}{d}$$

 \mathcal{E} is the permittivity of the substance. The potential difference between the plates is given by

$$V = Ed = \frac{\sigma - \sigma_p}{\epsilon_0} d$$

For linear dielectrics:

$$\sigma-\sigma_p=\frac{\sigma}{k}$$

Where k is a dielectric constant of the substance, K = 1.

$$K = \frac{\epsilon}{\epsilon_0}$$

How does the dielectric increase the capacitance of a capacitor?

The electric field between the plates of parallel plate capacitor is directly proportional to capacitance C of the capacitor. The strength of the electric field is reduced due to the presence of dielectric. If the total charge on the plates is kept constant, then the potential difference is reduced across the capacitor plates. In this way, dielectric increases the capacitance of the capacitor.



Electric Current:

An **electric current** is the rate of flow of electric charge past a point or region. An electric current is said to exist when there is a net flow of electric charge through a region. In electric circuits this charge is often carried by electrons moving through a wire. It can also be carried by ions in an electrolyte, or by both ions and electrons such as in an ionized gas (plasma). The SI unit of electric current is the ampere, which is the flow of electric charge across a surface at the rate of one coulomb per second. The ampere (symbol: A) is an SI base unit Electric current is measured using a device called an ammeter. Electric currents cause Joule heating, which creates light in incandescent light bulbs. They also create magnetic fields, which are used in motors, inductors and generators. The moving charged particles in an electric current are called charge carriers. In metals, one or more electrons from each atom are loosely bound to the atom, and can move freely about within the metal. These conduction electrons are the charge carriers in metal conductors.

Current

Current is the flow of electrical charge carriers like electrons. Current flows from negative to positive points. The SI unit for measuring electric current is the ampere (A). One ampere of current is defined as one coulomb of electrical charge moving past a unique point in a second. Electric current is widely used in household and industrial appliances. There are two types of electric current, namely alternating and direct current. In alternating current, the flow of current reverses its direction periodically. Alternating current in a circuit is represented by the sine wave. Direct current, unlike alternating current, flows in the same direction continuously. An example of direct current would be the current provided by a battery. In order to calculate the current flow through a conductor, Ohm's law is used. According to Ohm's law, the current through a conductor between two given points is also directly proportional to the potential difference between the points. The constant used in the proportionality is called resistance and the mathematical equation is I =V/R.

Electric current produces heating and magnetic effects. When current passes through a conductor, there is some heat generation due to ohmic loss in the conductor. This property is put to use for creating light in incandescent light bulbs. The stronger the current, the higher would be intensity of the magnetic field. Electric current is measured with the help of an ammeter.

Voltage, Current, and Resistance

An electric circuit is formed when a conductive path is created to allow electric charge to continuously move. This continuous movement of electric charge through the conductors of a circuit is called a *current*, and it is often referred to in terms of "flow," just like the flow of a liquid through a hollow pipe. The force motivating charge carriers to "flow" in a circuit is called *voltage*. Voltage is a specific measure of potential energy that is always relative between two points. When we speak of a certain amount of voltage being present in a circuit, we are referring to the measurement of how much *potential* energy exists to move charge carriers from one particular point in that circuit to another particular point. Without reference to *two* particular points, the term "voltage" has no meaning. Current tends to move through the conductors with some degree of friction, or opposition to motion. This opposition to motion is more properly called *resistance*. The amount of current in a circuit depends on the

amount of voltage and the amount of resistance in the circuit to oppose current flow. Just like voltage, resistance is a quantity relative between two points. For this reason, the quantities of voltage and resistance are often stated as being "between" or "across" two points in a circuit.

The Ohm's Law Equation

Ohm's principal discovery was that the amount of electric current through a metal conductor in a circuit is directly proportional to the voltage impressed across it, for any given temperature. Ohm expressed his discovery in the form of a simple equation, describing how voltage, current, and resistance interrelate:

$$E = 1 R$$

In this algebraic expression, voltage (E) is equal to current (I) multiplied by resistance (R). Using algebra techniques, we can manipulate this equation into two variations, solving for I and for R, respectively:

$$1 = \frac{E}{R} \qquad \qquad R = \frac{E}{1}$$

Combination of resistors in series connection

Consider three resistors R_1 , R_2 , R_3 which are connected in series. Here the charge first flows through R_1 and enters R_2 and finally reaches R_3 .



Combination of three resistors in series

By ohm's law, the potential difference across $R_1 = V_1 = I R_1$

The potential difference across $R_2 = V_2 = I R_2$.

The potential difference across $R_3 = V_3 = I R_3$.

Thus, the potential difference V across this series connection of resistors

$$\mathbf{V} = \mathbf{V}_1 + \mathbf{V}_2 + \mathbf{V}_3$$

 $= I R_1 + I R_2 + I R_3$

COPYRIGHT FIMT 2020

 $= I (R_1 + R_2 + R_3)$

Thus, in case of series connection, the equivalent resistance, $R_{eq} = V/I = (R_1 + R_2 + R_3)$.

For n number of resistors connected in series, the equivalent resistance $R_{eq} = R_1 + R_2 + R_3 \dots R_n$.

Equivalent resistance is the total resistance of the circuit. It is the single value of resistance which may replace the number of the resistors in the circuit without changing current and voltage in the network. Thus, in a series connection, the total resistance of a circuit is determined by adding the resistance of each individual resistors.

For Example, consider a series circuit which consists of three resistors with resistance 5 Ω , 10 Ω , 5 Ω respectively with a 15v battery.

So, the total resistance $R = R_{1+}R_{2+}R_{3} = 5 + 10 + 5 = 20\Omega$

We know that V = I R.

Current I = V / R = 15/20 = 0.75A. Though the current in the series network is the same, the voltage drop across each resistor is different. Each resistor with different resistance provide different voltage drop and we can find the total voltage by ohm's law, V = I R. Let's take the previous example. The voltage mentioned here is 15v. We can verify that by calculating in this way. V = I R₁ + I R₂ + I R₃ = 0.75 (5 +10 + 5) = 15v. This is the exact measurement of voltage we provided in this example. It is found that the resistor with more resistance has more voltage drop.

Combination of resistors in Parallel Connection

Consider three resistors R_1 , R_2 , R_3 which are connected in parallel. The charge splits into three and flows through R_1 , R_2 and R_3 .



Combination of three resistors in parallel

Current $I = I_1 + I_2 + I_3$.

The potential difference applied to $R1 = V = I_1 R_1$

The potential difference across $R_2 = V = I_2 R_2$

The potential difference across $R_3 = V = I_3 R_3$

Thus
$$I = I_1 + I_2 + I_3$$

$$= V/R_1 + V/R_2 + V/R_3$$

 $= V (1/R_1 + 1/R_2 + 1/R_3)$

If this parallel combination is replaced by an equivalent resistance, Req

Then $I = V/R_{eq}$

 $1/R_{eq} = 1/R_1 + 1/R_2 + 1/R_3$

So, in a parallel connection, the total resistance of a circuit is determined by adding the reciprocal of the resistance of each individual resistors.

For Example, consider a parallel circuit which consists of three resistors with resistance 5Ω , 10Ω , 5Ω respectively with a 15v battery.

Thus, the total resistance, $1/R = 1/R_1 + 1/R_2 + 1/R_3 = 1/5 + 1/10 + 1/5 = 5/10$. So $R = 2\Omega$.

In parallel connection, the total resistance or the equivalent resistance will always be less than the smallest resistor present in the circuit. The value of equivalent resistance will be between the smallest resistance in the circuit and the smallest resistance divided by the number of resistors present in the circuit. In this example, the smallest resistor has 5 Ω resistance and the value for total resistance is 2 Ω which clearly justifies the above mentioned fact.

The voltage across each resistor is 15v. Now to find the current across each branch, the formula is I = V / R.

 $I_1 = 15/5 = 3A$ $I_2 = 15/10 = 1.5A$ $I_3 = 15/5 = 3A$

Total current I = 3 + 1.5 + 3 = 7.5A

Combination of series and parallel resistors

Consider a circuit in which R₂ and R₃ are in parallel and R₁ is in series with R₂ and R₃.



Combination of series and parallel resistors

First consider R_2 and R_3 and thus $1/R_{23eq} = 1/R_2 + 1/R_3$ $R_{23eq} = R_2 R_3 / R_2 + R_3$ $R_{123eq} = R_{23eq} + R_1$ Thus, current $I = V/R_{123eq} = V/[R_1 + (R_2 R_3 / R_2 + R_3)]$ $= V (R_2 + R_3) / R_1R_2 + R_1 R_3 + R_2 R_3.$

NAAC ACCREDITED

Kirchhoffs Circuit Law

In 1845, a German physicist, **Gustav Kirchhoff** developed a pair or set of rules or laws which deal with the conservation of current and energy within electrical circuits. These two rules are commonly known as: *Kirchhoffs Circuit Laws* with one of Kirchhoffs laws dealing with the current flowing around a closed circuit, **Kirchhoffs Current Law**, (**KCL**) while the other law deals with the voltage sources present in a closed circuit, **Kirchhoffs Voltage Law**, (**KVL**).

Kirchhoffs First Law – The Current Law, (KCL)

Kirchhoffs Current Law or KCL, states that the "*total current or charge entering a junction* or node is exactly equal to the charge leaving the node as it has no other place to go except to leave, as no charge is lost within the node". In other words the algebraic sum of ALL the currents entering and leaving a node must be equal to zero, $I_{(exiting)} + I_{(entering)} = 0$. This idea by Kirchhoff is commonly known as the **Conservation of Charge**.

Kirchhoffs Current Law



Here, the three currents entering the node, I_1 , I_2 , I_3 are all positive in value and the two currents leaving the node, I_4 and I_5 are negative in value. Then this means we can also rewrite the equation as;

 $I_1 + I_2 + I_3 - I_4 - I_5 = 0$

COPYRIGHT FIMT 2020

The term **Node** in an electrical circuit generally refers to a connection or junction of two or more current carrying paths or elements such as cables and components. Also for current to flow either in or out of a node a closed circuit path must exist. We can use Kirchhoff's current law when analysing parallel circuits.

Kirchhoffs Second Law – The Voltage Law, (KVL)

Kirchhoffs Voltage Law or KVL, states that "*in any closed loop network, the total voltage around the loop is equal to the sum of all the voltage drops within the same loop*" which is also equal to zero. In other words the algebraic sum of all voltages within the loop must be equal to zero. This idea by Kirchhoff is known as the **Conservation of Energy**.

Kirchhoffs Voltage Law



Starting at any point in the loop continue in the **same direction** noting the direction of all the voltage drops, either positive or negative, and returning back to the same starting point. It is important to maintain the same direction either clockwise or anti-clockwise or the final voltage sum will not be equal to zero. We can use Kirchhoff's voltage law when analysing series circuits.

When analysing either DC circuits or AC circuits using **Kirchhoffs Circuit Laws** a number of definitions and terminologies are used to describe the parts of the circuit being analysed such as: node, paths, branches, loops and meshes. These terms are used frequently in circuit analysis so it is important to understand them.

Wheatstone Bridge

The Wheatstone bridge (or resistance bridge) circuit can be used in a number of applications and today, with modern operational amplifiers we can use the *Wheatstone Bridge Circuit* to interface various transducers and sensors to these amplifier circuits.

The Wheatstone Bridge circuit is nothing more than two simple series-parallel arrangements of resistances connected between a voltage supply terminal and ground producing zero voltage difference between the two parallel branches when balanced. A Wheatstone bridge circuit has two input terminals and two output terminals consisting of four resistors configured in a diamond-like arrangement as shown. This is typical of how the Wheatstone bridge is drawn.

The Wheatstone Bridge



When balanced, the Wheatstone bridge can be analysed simply as two series strings in parallel. In our tutorial about **Resistors in Series**, we saw that each resistor within the series chain produces an **IR** drop, or voltage drop across itself as a consequence of the current flowing through it as defined by Ohms Law. Consider the series circuit below.



As the two resistors are in series, the same current (i) flows through both of them. Therefore the current flowing through these two resistors in series is given as: V/R_T .

 $I = V \div R = 12V \div (10\Omega + 20\Omega) = 0.4A$

The voltage at point C, which is also the voltage drop across the lower resistor, R_2 is calculated as:

 $V_{R2} = I \times R_2 = 0.4A \times 20\Omega = 8$ volts

COPYRIGHT FIMT 2020

Then we can see that the source voltage V_S is divided among the two series resistors in direct proportion to their resistances as $V_{R1} = 4V$ and $V_{R2} = 8V$. This is the principle of voltage division, producing what is commonly called a potential divider circuit or voltage divider network. Now if we add another series resistor circuit using the same resistor values in parallel with the first we would have the following circuit.



As the second series circuit has the same resistive values of the first, the voltage at point D, which is also the voltage drop across resistor, R_4 will be the same at 8 volts, with respect to zero (battery negative), as the voltage is common and the two resistive networks are the same. But something else equally as important is that the voltage difference between point C and point D will be zero volts as both points are at the same value of 8 volts as: C = D = 8 volts, then the voltage difference is: 0 volts

When this happens, both sides of the parallel bridge network are said to be **balanced** because the voltage at point C is the same value as the voltage at point D with their difference being zero. Now let's consider what would happen if we reversed the position of the two resistors, R_3 and R_4 in the second parallel branch with respect to R_1 and R_2 .



With resistors, R_3 and R_4 reversed, the same current flows through the series combination and the voltage at point D, which is also the voltage drop across resistor, R_4 will be:

 $V_{R4} = 0.4A \times 10\Omega = 4$ volts

Now with V_{R4} having 4 volts dropped across it, the voltage difference between points C and D will be 4 volts as: C = 8 volts and D = 4 volts. Then the difference this time is: 8 - 4 = 4 volts

The result of swapping the two resistors is that both sides or "arms" of the parallel network are different as they produce different voltage drops. When this happens the parallel network is said to be **unbalanced** as the voltage at point C is at a different value to the voltage at point D.

Then we can see that the resistance ratio of these two parallel arms, ACB and ADB, results in a voltage difference between **0** volts (balanced) and the maximum supply voltage (unbalanced), and this is the basic principal of the Wheatstone Bridge Circuit.

So we can see that a Wheatstone bridge circuit can be used to compare an unknown resistance R_X with others of a known value, for example, R_1 and R_2 , have fixed values, and R_3 could be variable. If we connected a voltmeter, ammeter or classically a galvanometer between points C and D, and then varied resistor, R_3 until the meters read zero, would result in the two arms being balanced and the value of R_X , (substituting R_4) known as shown.

Wheatstone Bridge Circuit



By replacing R_4 above with a resistance of known or unknown value in the sensing arm of the Wheatstone bridge corresponding to R_X and adjusting the opposing resistor, R_3 to "balance" the bridge network, will result in a zero voltage output. Then we can see that balance occurs when:

$$\frac{R_1}{R_2} = \frac{R_3}{R_x} = 1 \text{ (Balanced)}$$

The Wheatstone Bridge equation required to give the value of the unknown resistance, R_X at balance is given as:

$$V_{\text{OUT}} = (V_{\text{C}} - V_{\text{D}}) = (V_{\text{R}2} - V_{\text{R}4}) = 0$$

$$R_{\text{C}} = \frac{R_2}{R_1 + R_2} \text{ and } R_{\text{D}} = \frac{R_4}{R_3 + R_4}$$
At Balance: $R_{\text{C}} = R_{\text{D}}$ So, $\frac{R_2}{R_1 + R_2} = \frac{R_4}{R_3 + R_4}$

$$\therefore R_2(R_3 + R_4) = R_4(R_1 + R_2)$$

$$R_2R_3 + R_2R_4 = R_1R_4 + R_2R_4$$

$$\therefore R_4 = \frac{R_2R_3}{R_1} = R_X$$

Where resistors, R_1 and R_2 are known or preset values.

Wheatstone Bridge Example No1

The following unbalanced Wheatstone Bridge is constructed. Calculate the output voltage across points C and D and the value of resistor R₄ required to balance the bridge circuit.

REDI



The value of resistor, R₄ required to balance the bridge is given as:

$$R_{4} = \frac{R_{2}R_{3}}{R_{1}} = \frac{120\Omega \times 480\Omega}{80\Omega} = 720\Omega$$

COPYRIGHT FIMT 2020

We have seen above that the **Wheatstone Bridge** has two input terminals (A-B) and two output terminals (C-D). When the bridge is balanced, the voltage across the output terminals is 0 volts. When the bridge is unbalanced, however, the output voltage may be either positive or negative depending upon the direction of unbalance.

Wheatstone Bridge Light Detector

Balanced bridge circuits find many useful electronics applications such as being used to measure changes in light intensity, pressure or strain. The types of resistive sensors that can be used within a wheatstone bridge circuit include: photoresistive sensors (LDR's), positional sensors (potentiometers), piezoresistive sensors (strain gauges) and temperature sensors (thermistor's), etc. There are many wheatstone bridge applications for sensing a whole range of mechanical and electrical quantities, but one very simple wheatstone bridge application is in the measurement of light by using a photoresistive device. One of the resistors within the bridge network is replaced by a light dependent resistor, or LDR.

An LDR, also known as a cadmium-sulphide (Cds) photocell, is a passive resistive sensor which converts changes in visible light levels into a change in resistance and hence a voltage. Light dependent resistors can be used for monitoring and measuring the level of light intensity, or whether a light source is ON or OFF.

A typical Cadmium Sulphide (CdS) cell such as the ORP12 light dependent resistor typically has a resistance of about one Megaohm (M Ω) in dark or dim light, about 900 Ω at a light intensity of 100 Lux (typical of a well lit room), down to about 30 Ω in bright sunlight. Then as the light intensity increases the resistance reduces. By connecting a light dependant resistor to the Wheatstone bridge circuit above, we can monitor and measure any changes in the light levels as shown.

Wheatstone Bridge Light Detector



The LDR photocell is connected into the Wheatstone Bridge circuit as shown to produce a light sensitive switch that activates when the light level being sensed goes above or below the pre-set value determined by V_{R1} . In this example V_{R1} either a 22k or 47k Ω potentiometer. The op-amp is connected as a voltage comparator with the reference voltage V_D applied to the non-inverting pin. In this example, as both R_3 and R_4 are of the same 10k Ω value, the reference voltage set at point D will therefore be equal to half of Vcc. That is Vcc/2. The potentiometer, V_{R1} sets the trip point voltage V_C , applied to the inverting input and is set to the required nominal light level. The relay turns "ON" when the voltage at point C is less than the voltage at point D. Adjusting V_{R1} sets the voltage at point C to balance the bridge circuit at the required light level or intensity. The LDR can be any cadmium sulphide device that has a high impedance at low light levels and a low impedance at high light levels. Note that the circuit can be used to act as a "light-activated" switching circuit or a "dark-activated" switching circuit simply by transposing the LDR and R_3 positions within the design.

The **Wheatstone Bridge** has many uses in electronic circuits other than comparing an unknown resistance with a known resistance. When used with Operational Amplifiers, the Wheatstone bridge circuit can be used to measure and amplify small changes in resistance, R_X due, for example, to changes in light intensity as we have seen above.

But the bridge circuit is also suitable for measuring the resistance change of other changing quantities, so by replacing the above photo-resistive LDR light sensor for a thermistor, pressure sensor, strain gauge, and other such transducers, as well as swapping the positions of the LDR and V_{R1} , we can use them in a variety of other Wheatstone bridge applications.

Also more than one resistive sensor can be used within the four arms (or branches) of the bridge formed by the resistors R_1 to R_4 to produce "full-bridge", "half-bridge" or "quarterbridge circuit arrangements providing thermal compensation or automatic balancing of the Wheatstone bridge.

150 9001-2015 & 14001-2015

Slide Wire Bridge

A meter bridge also called a slide wire bridge is an instrument that works on the principle of a Wheatstone bridge. A meter bridge is used in finding the unknown resistance of a conductor as that of in a Wheatstone bridge.

What Is Wheatstone Bridge?

A Wheatstone bridge is a kind of electrical circuit used in measuring an electrical resistance, which is unknown by balancing its two legs of the bridge circuit, where one of the legs includes an unknown component. Samuel Hunter Christie created this instrument in the year 1833 and was improved and also simplified by Sir Charles Wheatstone in the year 1843. The digital multimeters in today's world provide the simplest forms in measuring the resistance. The Wheatstone Bridge can still be used in measuring light values of resistances around the range of milli-Ohms.

How Is A Meter Bridge Used In Finding The Unknown Resistance?

A meter bridge is an apparatus utilized in finding the unknown resistance of a coil. The below figure 12 is the diagram of a useful meter bridge instrument.



In the above figure, R is called as the Resistance, P is the Resistance coming across AB, S is the Unknown Resistance, Q is the Resistance between the joints BD.

AC is the long wire measuring 1m in length and it is made of constantan or manganin having a uniform area of the cross-section Such that L1 + L2 = 100

Assuming that $L1 = L \Rightarrow L2 = 100 - L$

Relation obtains the unknown resistance 'X' of the given wire: X = RL2/L1 = R(100 - L)/L

And the specific resistance of the material for a given wire is obtained by the relation = (3.14) r2X/l where, r = the radius of the cable and also l = length of the wire.

The devices required in finding the unknown resistance of a conductor using a meter bridge are:

- Meter bridge
- Resistance box
- Galvanometer
- Unknown Resistance of a length 1 m
- Screw gauge
- Connecting Wires
- Jockey
- One way key

In the meter bridge, one of the lateral kinds of resistances is replaced by a wire having a length of the uniform cross section of about 1m. The other pair consists of one known and an unknown pair of resistances. The one part of the galvanometer is connected in between both resistances, whereas the other part of the wire is finding the null point where the galvanometer is not showing any deflection. At this point, the bridge is said to be balanced.

ARGEME

Procedure For Finding The Unknown Resistance Using Meter Bridge

- Collect the instruments and prepare connections as shown in the above figure.
- Take some suitable kind of resistance 'R' from the resistance box.
- Touch jockey at the point A; look that there exists a deflection in galvanometer on one of the sides, then contact the jockey on point C of wire, then the deflection in galvanometer has to be on another side.
- Find the position of the null point having deflection in the galvanometer that becomes zero. Note the length AB (I) BC = (100 – I).
- Continue the above method for some different values of the 'R'. Note at least some 5 readings.
- Consider the point where galvanometer shows a 0 deflection; this is called balance point.
- Now, Measure the length of given wire by the use of ordinary scale and radius of the wire by the utilization of a screw gauge, (Take at least five readings).
- Calculate Mean Resistance of Single Unknown Resistance = Total Sum of resistances of Unknown resistance from the above five readings)/5.

UNIT – IV

Structure of Atom:

An atom is the smallest unit of matter that retains all of the chemical properties of an element. Atoms combine to form molecules, which then interact to form solids, gases, or liquids. For example, water is composed of hydrogen and oxygen atoms that have combined to form water molecules. Many biological processes are devoted to breaking down molecules into their component atoms so they can be reassembled into a more useful molecule.

Atomic Particles

Atoms consist of three basic particles: protons, electrons, and neutrons. The nucleus (center) of the atom contains the protons (positively charged) and the neutrons (no charge). The outermost regions of the atom are called electron shells and contain the electrons (negatively charged). Atoms have different properties based on the arrangement and number of their basic particles. The hydrogen atom (H) contains only one proton, one electron, and no neutrons. This can be determined using the atomic number and the mass number of the element (see the concept on atomic numbers and mass numbers).



Structure of an atom: Elements, such as helium, depicted here, are made up of atoms. Atoms are made up of protons and neutrons located within the nucleus, with electrons in orbitals surrounding the nucleus.

Atomic Mass

Protons and neutrons have approximately the same mass, about 1.67×10^{-24} grams. Scientists define this amount of mass as one atomic mass unit (amu) or one Dalton. Although similar in mass, protons are positively charged, while neutrons have no charge. Therefore, the number of neutrons in an atom contributes significantly to its mass, but not to its charge. Electrons are much smaller in mass than protons, weighing only 9.11×10^{-28} grams, or about 1/1800 of an atomic mass unit. Therefore, they do not contribute much to an element's overall atomic

mass. When considering atomic mass, it is customary to ignore the mass of any electrons and calculate the atom's mass based on the number of protons and neutrons alone. Electrons contribute greatly to the atom's charge, as each electron has a negative charge equal to the positive charge of a proton. Scientists define these charges as "+1" and "-1. " In an uncharged, neutral atom, the number of electrons orbiting the nucleus is equal to the number of protons inside the nucleus. In these atoms, the positive and negative charges cancel each other out, leading to an atom with no net charge.

Thomson Model of an atom

The description of Thomson's atomic model is one of the many scientific models of the atom. It was proposed by J.J Thomson in the year 1904 just after the discovery of electrons. However, at that time the atomic nucleus was yet to be discovered. So, he proposed a model on the basis of known properties available at that time. The known properties are:

ITTING PACALLEVITED

- Atoms are neutrally charged
- Negatively charged particles called electrons are present in an atom. Learn about Charged particles in Matter in more detail here.

Thomson's Atomic Model- Postulates

- According to the postulates of Thomson's atomic model, an atom resembles a sphere of positive charge with electrons (negatively charged particles) present inside the sphere.
- The positive and negative charge is equal in magnitude and therefore an atom has no charge as a whole and is electrically neutral.
- Thomson's atomic model resembles a spherical plum pudding as well as a watermelon. It
 resembles a plum pudding because the electrons in the model look like the dry fruits
 embedded in a sphere of positive charge just like a spherical plum pudding. The model
 has also been compared to a watermelon because the red edible part of a watermelon
 was compared to the sphere having a positive charge and the black seeds filling the
 watermelon looked similar to the electrons inside the sphere.


Thomson's Atomic Model

Limitations of Thomson's Atomic Model

- Thomson's atomic model failed to explain how the positive charge holds on the electrons inside the atom. It also failed to explain an atom's stability.
- The theory did not mention anything about the nucleus of an atom.
- It was unable to explain the scattering experiment of Rutherford.

Alpha-Particle Scattering and Rutherford's Nuclear Model of Atom

In 1911, Rutherford, along with his assistants, H. Geiger and E. Marsden, performed the Alpha Particle scattering experiment, which led to the birth of the 'nuclear model of an atom' – a major step towards how we see the atom today.

J.J Thomson's Plum-pudding Model

In 1897-98, the first model of an atom was proposed by J.J. Thomson. Famously known as the Plum-pudding model or the watermelon model, he proposed that an atom is made up of a positively charged ball with electrons embedded in it. Further, the negative and positive charges were equal in number, making the atom electrically neutral. Figure 1 shows what Thomson's plum-pudding model of an atom looked like. Ernest Rutherford, a former research student working with J.J. Thomson, proposed an experiment of scattering of alpha particles by atoms to understand the structure of an atom. Rutherford, along with his assistants – H. Geiger and E. Marsden – started performing experiments to study the structure of an atom. In 1911, they performed the Alpha particle scattering experiment, which led to the birth of the 'nuclear model of an atom' – a major step towards how we see the atom today.



Thomson's 'plum-pudding' model of the atom

Figure 1.

in Di

GEMEAL

The Alpha Particle Scattering Experiment

They took a thin gold foil having a thickness of 2.1×10^{-7} m and placed it in the centre of a rotatable detector made of zinc sulfide and a microscope. Then, they directed a beam of 5.5MeV alpha particles emitted from a radioactive source at the foil. Lead bricks collimated these alpha particles as they passed through them. After hitting the foil, the scattering of these alpha particles could be studied by the brief flashes on the screen. Rutherford and his team expected to learn more about the structure of the atom from the results of this experiment.



Observations

Here is what they found:

- Most of the alpha particles passed through the foil without suffering any collisions
- Around 0.14% of the incident alpha particles scattered by more than 1°
- Around 1 in 8000 alpha particles deflected by more than 90°

These observations led to many arguments and conclusions which laid down the structure of the nuclear model on an atom.

Conclusions and arguments

The results of this experiment were not in sync with the plum-pudding model of the atom as suggested by Thomson. Rutherford concluded that since alpha particles are positively charged, for them to be deflected back, they needed a large repelling force. He further argued that for this to happen, the positive charge of the atom needs to be concentrated in the centre, unlike scattered in the earlier accepted model. Hence, when the incident alpha particle came very close to the positive mass in the centre of the atom, it would repel leading to a deflection. On the other hand, if it passes through at a fair distance from this mass, then there would be no deflection and it would simply pass through. He then suggested the 'nuclear model of an atom' wherein the entire positive charge and most of the mass of the atom is concentrated in the nucleus. Also, the electrons are moving in orbits around the nucleus akin to the planets and the sun. Further, Rutherford also concluded from his experiments that the size of the nucleus is between 10^{-15} and 10^{-14} m.

According to Kinetic theory, the size of an atom is around 10⁻¹⁰ m or around 10,000 to 100,000 times the size of the nucleus proposed by Rutherford. Hence, the distance of the electrons from the nucleus should be around 10,000 to 100,000 times the size of the nucleus. This eventually implies that most of the atom is empty space and explains why most alpha particles went right through the foil. And, these particles are deflected or scattered through a large angle on coming close to the nucleus. Also, the electrons having negligible mass, do not affect the trajectory of these incident alpha particles.

Alpha Particle Trajectory

The trajectory traced by an alpha particle depends on the impact parameter of the collision. The impact parameter is simply the perpendicular distance of each alpha particle from the centre of the nucleus. Since in a beam all alpha particles have the same kinetic energy, the scattering of these particles depends solely on the impact parameter. Hence, the particles with a small impact parameter or the particles closer to the nucleus, experience large angle of scattering. On the other hand, those with a large impact parameter suffer no deflection or scattering at all. Finally, those particles having ~zero impact parameter or a head-on collision with the nucleus rebound back.

Coming to the experiment, Rutherford and his team observed that a really small fraction of the incident alpha particles was rebounding back. Hence, only a small number of particles were colliding head-on with the nucleus. This, subsequently, led them to believe that the mass of the atom is concentrated in a very small volume.

Electron Orbits

In a nutshell, Rutherford's nuclear model of the atom describes it as:

- An electrically neutral sphere with
 - A small and positively charged nucleus at the centre
 - Surrounded by revolving electrons in their dynamically stable orbits

The centripetal force that keeps the electrons in their orbits is an outcome of:

- The electrostatic force of attraction between-
 - The positively charged nucleus and
 - The negatively charged revolving electrons.

Postulates of Bohr's atomic model

The physicist Niels Bohr said, "Anyone who is not shocked by quantum theory has not understood it." He also said, "We must be clear that when it comes to atoms, language can only be used as in poetry." So what exactly is this Bohr atomic model? Let us find out! Bohr atomic model and the models after that explain the properties of atomic electrons on the basis of certain allowed possible values. The model explained how an atom absorb or emit radiation when electrons on subatomic level jump between the allowed and stationary states. Germanborn physicists James Franck and Gustav Hertz obtained the experimental evidence of the presence of these states.

Bohr Atomic Model

A Danish physicist named Neil Bohr in 1913 proposed the Bohr atomic model. He modified the problems and limitations associated with Rutherford's model of an atom. Earlier in Rutherford Model, Rutherford explained in an atom a nucleus is positively charged and is surrounded by electrons (negatively charged particles). The electrons move around in a predictable path

called **orbits**. Bohr modified Rutherford's model where he explained that electrons move around in fixed orbital shells. Furthermore, he explained that each orbital shell has fixed energy levels. Therefore, Rutherford basically explained a nucleus of an atom whereas Bohr took the model one step ahead. He explained about electrons and the different energy levels associated with it.

According to Bohr Atomic model, a small positively charged nucleus is surrounded by revolving negatively charged electrons in fixed orbits. He concluded that electron will have more energy if it is located away from the nucleus whereas the electrons will have less energy if it located near



Bohr's Model of an Atom

Postulates of the Bohr Atomic Model

- Electrons revolve around the nucleus in a fixed circular path termed "orbits" or "shells" or "energy level."
- The orbits are termed as "stationary orbit."
- Every circular orbit will have a certain amount of fixed energy and these circular orbits were termed orbital shells. The electrons will not radiate energy as long as they continue to revolve around the nucleus in the fixed orbital shells.
- The different energy levels are denoted by integers such as n=1 or n=2 or n=3 and so on. These are called as quantum numbers. The range of quantum number may vary and begin from the lowest energy level (nucleus side n=1) to highest energy level. Learn the concept of an Atomic number <u>here</u>.
- The different energy levels or orbits are represented in two ways such as 1, 2, 3, 4... or K,
 L, M, N.... shells. The lowest energy level of the electron is called the ground state. Learn the concept of Valency here in detail <u>here</u>.

 The change in energy occurs when the electrons jump from one energy level to other. In an atom, the electrons move from lower to higher energy level by acquiring the required energy. However, when an electron loses energy it moves from higher to lower energy level.

Therefore,

- 1st orbit (energy level) is represented as K shell and it can hold up to 2 electrons.
- 2nd orbit (energy level) is represented as L shell and it can hold up to 8 electrons.
- 3rd orbit (energy level) is represented as M shell and it can contain up to 18 electrons.
- 4th orbit (energy level) is represented as N Shell and it can contain maximum 32 electrons.

The orbits continue to increase in a similar manner.

Distribution of Electrons in Orbits or Shells:

Electronic distribution of various orbits or energy levels can be calculated by the formula $2n^2$. Here, 'n' denotes the number of orbits.

- The number of electrons in K shell (1st orbit) can be calculated by 2n²= 2 x 1² = 2. Thus, maximum number of electrons in 1st orbit = 2
- Similarly, The number of electrons in L shell (2nd orbit)= 2 x 2² = 8. Thus, maximum number of electrons in 2nd orbit = 8

We can determine the maximum number of electrons in a similar way. Read about Thomson's Model of an Atom, the very first model of an Atom by J.J. Thomsons.

Limitations of Bohr's Model of an Atom:

Bohr atomic model had few limitations. They are:

- Failure to explain Zeeman Effect (how atomic spectra are affected by magnetic fields).
- It contradicts Heisenberg Uncertainty Principle.

• Unable to explain how to determine the spectra of larger atoms.

Semiconductors

Semiconductors are materials which have a conductivity between conductors (generally metals) and nonconductors or insulators (such as most ceramics). **Semiconductors** can be pure elements, such as silicon or germanium, or compounds such as gallium arsenide or cadmium selenide.

Semiconductor, any of a class of crystalline solids intermediate in electrical conductivity between a conductor and an insulator. Semiconductors are employed in the manufacture of various kinds of electronic devices, including diodes, transistors, and integrated circuits. Such devices have found wide application because of their compactness, reliability, power efficiency, and low cost. As discrete components, they have found use in power devices, optical sensors, and light emitters, including solid-state lasers. They have a wide range of current- and voltage-handling capabilities and, more important, lend themselves to integration into complex but readily manufacturable microelectronic circuits. They are, and will be in the foreseeable future, the key elements for the majority of electronic systems, serving communications, signal processing, computing, and control applications in both the consumer and industrial markets.

Semiconductor Materials

Solid-state materials are commonly grouped into three classes: insulators, semiconductors, and conductors. (At low temperatures some conductors, semiconductors, and insulators may become superconductors.) The figure shows the conductivities σ (and the corresponding resistivities $\rho = 1/\sigma$) that are associated with some important materials in each of the three classes. Insulators, such as fused quartz and glass, have very low conductivities, on the order of 10^{-18} to 10^{-10} siemens per centimetre; and conductors, such as aluminum, have high conductivities, typically from 10^4 to 10^6 siemens per centimetre. The conductivities of semiconductors are between these extremes and are generally sensitive to temperature, illumination, magnetic fields, and minute amounts of impurity atoms. For example, the addition of about 10 atoms of boron (known as a dopant) per million atoms of silicon can increase its electrical conductivity a thousandfold (partially accounting for the wide variability shown in the preceding figure).



Typical range of conductivities for insulators, semiconductors, and conductors. *Encyclopædia* Britannica, Inc.

The study of semiconductor materials began in the early 19th century. The elemental semiconductors are those composed of single species of atoms, such as silicon (Si), germanium (Ge), and tin (Sn) in column IV and selenium (Se) and tellurium (Te) in column VI of the periodic table. There are, however, numerous compound semiconductors, which are composed of two or more elements. Gallium arsenide (GaAs), for example, is a binary III-V compound, which is a combination of gallium (Ga) from column III and arsenic (As) from column V. Ternary compounds can be formed by elements from three different columns—for instance, mercury indium telluride (HgIn₂Te₄), a II-III-VI compound. They also can be formed by elements from two columns, such as aluminum gallium arsenide (AlxGa₁ – xAs), which is a ternary III-V compound, where both Al and Ga are from column III and the subscript *x* is related to the composition of the two elements from 100 percent Al (x = 1) to 100 percent Ga (x = 0). Pure silicon is the most important material for integrated circuit applications, and III-V binary and ternary compounds are most significant for light emission.

150 9001:2015 & 1400

Periodic table of the elements

Li Be 11 12 Na Mg 3 19 20 21 K Ca Sc 37 38 39 Rb Sr Y	4 22 Ti 40	5 23 V 41	6 24 Cr	7 25 Mn	8 26 Fe	9 27 Co	10 28 Ni	11 29	12 30	13 Al 31	14 Si 32	N 15 P 33	16 S 34	17 CI 35	18 18 Ar 36
19 20 21 K Ca Sc 37 38 39 Rb Sr Y	22 Ti 40	23 V 41	24 Cr	25 Mn	26 Fe	27 Co	28 Ni	29	30	31	32	33	34	35	36
37 38 39 Rb Sr Y	40	41	-				1.41	Cu	Zn	Ga	Ge	As	Se	Br	Kr
	Zr	Nb	42 Mo	43 Tc	44 Ru	45 Rh	46 Pd	47 Ag	48 Cd	49 In	50 Sn	51 Sb	52 Te	53 	54 Xe
55 56 57 Cs Ba La	72 Hf	73 Ta	74 W	75 Re	76 Os	77 Ir	78 Pt	79 Au	80 Hg	81 TI	82 Pb	83 Bi	84 Po	85 At	86 Rn
87 88 89 Fr Ra Ac	104 Rf	105 Db	106 Sg	107 Bh	108 Hs	109 Mt	110 Ds	111 Rg	112 Cn	113 Nh	114 Fl	115 Mc	116 Lv	117 Ts	118 Og
	58	59	60	61	62	63	64	65	66	67	68	69	70	71	1

*Numbering system adopted by the International Union of Pure and Applied Chemistry (IUPAC). © Encyclopædia Britannica, Inc.

periodic table Modern version of the periodic table of the elements.

Prior to the invention of the bipolar transistor in 1947, semiconductors were used only as two-terminal devices, such as rectifiers and photodiodes. During the early 1950s germanium was the major semiconductor material. However, it proved unsuitable for many applications, because devices made of the material exhibited high leakage currents at only moderately elevated temperatures. Since the early 1960s silicon has become by far the most widely used semiconductor, virtually supplanting germanium as a material for device fabrication. The main reasons for this are twofold: (1) silicon devices exhibit much lower leakage currents, and (2) silicon dioxide (SiO₂), which is a high-quality insulator, is easy to incorporate as part of a silicon-based device. Thus, silicon technology has become very advanced and pervasive, with silicon devices constituting more than 95 percent of all semiconductor products sold worldwide.

Many of the compound semiconductors have some specific electrical and optical properties that are superior to their counterparts in silicon. These semiconductors, especially gallium arsenide, are used mainly for optoelectronic and certain radio frequency (RF) applications.

Electronic Properties

The semiconductor materials described here are single crystals; i.e., the atoms are arranged in a three-dimensional periodic fashion. Part A of the figure shows a simplified twodimensional representation of an intrinsic (pure) silicon crystal that contains negligible impurities. Each silicon atom in the crystal is surrounded by four of its nearest neighbours. Each atom has four electrons in its outer orbit and shares these electrons with its four neighbours. Each shared electron pair constitutes a covalent bond. The force of attraction between the electrons and both nuclei holds the two atoms together. For isolated atoms (e.g., in a gas rather than a crystal), the electrons can have only discrete energy levels. However, when a large number of atoms are brought together to form a crystal, the interaction between the atoms causes the discrete energy levels to spread out into energy bands. When there is no thermal vibration (i.e., at low temperature), the electrons in an insulator or semiconductor crystal will completely fill a number of energy bands, leaving the rest of the energy bands empty. The highest filled band is called the valence band. The next band is the conduction band, which is separated from the valence band by an energy gap (much larger gaps in crystalline insulators than in semiconductors). This energy gap, also called a bandgap, is a region that designates energies that the electrons in the crystal cannot possess. Most of the important semiconductors have bandgaps in the range 0.25 to 2.5 electron volts (eV). The bandgap of silicon, for example, is 1.12 eV, and that of gallium arsenide is 1.42 eV. In contrast, the bandgap of diamond, a good crystalline insulator, is 5.5 eV.



Three bond pictures of a semiconductor. Encyclopædia Britannica, Inc.

At low temperatures the electrons in a semiconductor are bound in their respective bands in the crystal; consequently, they are not available for electrical conduction. At higher temperatures thermal vibration may break some of the covalent bonds to yield free electrons that can participate in <u>current</u> conduction. Once an electron moves away from a covalent bond, there is an electron vacancy associated with that bond. This vacancy may be filled by a neighbouring electron, which results in a shift of the vacancy location from one crystal site to another. This vacancy may be regarded as a fictitious particle, dubbed a "<u>hole</u>," that carries a positive charge and moves in a direction opposite to that of an electron. When an <u>electric field</u> is applied to the semiconductor, both the free electrons (now residing in the conduction band) and the holes (left behind in the valence band) move through the crystal, producing an

electric current. The electrical conductivity of a material depends on the number of free electrons and holes (charge carriers) per unit volume and on the rate at which these carriers move under the influence of an electric field. In an intrinsic semiconductor there exists an equal number of free electrons and holes. The electrons and holes, however, have different mobilities; that is, they move with different velocities in an electric field. For example, for intrinsic silicon at room temperature, the electron <u>mobility</u> is 1,500 square centimetres per volt-second ($cm^2/V \cdot s$)—i.e., an electron will move at a velocity of 1,500 centimetres per second under an electric field of one volt per centimetre—while the hole mobility is 500 $cm^2/V \cdot s$. The electron and hole mobilities in a particular semiconductor generally decrease with increasing temperature.



electron hole: movement Movement of an electron hole in a crystal lattice. *Encyclopædia* Britannica, Inc.

Electrical conduction in intrinsic semiconductors is quite poor at room temperature. To produce higher conduction, one can intentionally introduce impurities (typically to a concentration of one part per million host atoms). This is called doping, a process that increases conductivity despite some loss of mobility. For example, if a silicon atom is replaced by an atom with five outer electrons, such as arsenic (*see* part B of the figure), four of the electrons form covalent bonds with the four neighbouring silicon atoms. The fifth electron becomes a conductor because of the addition of the electron. The arsenic atom is the donor. Similarly, part C of the figure shows that, if an atom with three outer electrons, such as boron, is substituted for a silicon atom, an additional electron is accepted to form four covalent bonds around the boron atom, and a positively charged hole is created in the valence band. This creates a *p*-type semiconductor, with the boron constituting an acceptor.

The P-N Junction

If an abrupt change in impurity type from acceptors (p-type) to donors (n-type) occurs within a single crystal structure, a p-n junction is formed (*see* parts B and C of the figure). On the p side, the holes constitute the dominant carriers and so are called majority carriers. A few thermally generated electrons will also exist in the p side; these are termed minority carriers. On the n side, the electrons are the majority carriers, while the holes are the minority carriers. Near the junction is a region having no free charge carriers. This region, called the depletion layer, behaves as an insulator.



(A) Current-voltage characteristics of a typical silicon *p*-*n* junction. (B) Forward-bias and (C) reverse-bias conditions. (D) The symbol for a *p*-*n* junction.*Encyclopædia Britannica, Inc.*

The most important characteristic of *p*-*n* junctions is that they rectify. Part A of the figure shows the current-voltage characteristics of a typical silicon p-n junction. When a forward bias is applied to the *p*-*n* junction (i.e., a positive voltage applied to the *p*-side with respect to the *n*-side, as shown in part B of the figure), the majority charge carriers move across the junction so that a large current can flow. However, when a reverse bias is applied (as in part C of the figure), the charge carriers introduced by the impurities move in opposite directions away from the junction, and only a small leakage current flows. As the reverse bias is increased, the leakage current remains very small until a critical voltage is reached, at which point the current suddenly increases. This sudden increase in current is referred to as the nondestructive phenomenon if junction breakdown, usually а the resulting power dissipation is limited to a safe value. The applied forward voltage is typically less than one volt, but the reverse critical voltage, called the breakdown voltage, can vary from less than one volt to many thousands of volts, depending on the impurity concentration of the junction and other device parameters.

Although other junction types have been invented (including p-n-p and n-p-n), p-n junctions remain fundamental to semiconductor devices. For further details on applications of these basic semiconductor properties, *see* transistor and integrated circuit.

Semiconductors types / classifications

There are two basic groups or classifications that can be used to define the different semiconductor types:

- *Intrinsic material:* An intrinsic type of semiconductor material made to be very pure chemically. As a result it possesses a very low conductivity level having very few number of charge carriers, namely holes and electrons, which it possesses in equal quantities.
- **Extrinsic material:** Extrinisc types of semiconductor are those where a small amount of impurity has been added to the basic intrinsic material. This 'doping' uses an element from a different periodic table group and in this way it will either have more or less electrons in the valence band than the semiconductor itself. This creates either an excess or shortage of electrons. In this way two types of semiconductor are available: Electrons are negatively charged carriers.
 - <u>N-type</u>: An N-type semiconductor material has an excess of electrons. In this way, free electrons are available within the lattices and their overall movement in one direction under the influence of a potential difference results in an electric current flow. This in an N-type semiconductor, the charge carriers are electrons.
 - <u>P-type:</u> In a P-type semiconductor material there is a shortage of electrons, i.e. there are 'holes' in the crystal lattice. Electrons may move from one empty position to another and in this case it can be considered that the holes are moving. This can happen under the influence of a potential difference and the holes can be seen to flow in one direction resulting in an electric current flow. It is actually harder for holes to move than for free electrons to move and therefore the mobility of holes is less than that of free electrons. Holes are positively charged carriers.

Energy band theory in solids

In a single isolated atom, the electrons in each orbit have definite energy associated with it. But in case of solids all the atoms are close to each other, so the energy levels of outermost orbit electrons are affected by the neighboring atoms. When two single or isolated atoms are bring close to each other then the outermost orbit electrons of two atoms are interact or shared with each other. i.e, the electrons in the outermost orbit of one atom experience a attractive force from the nearest or neighboring atomic nucleus. Due to this the energies of the electrons will not be in same level, the energy levels of electrons are changed to a value which is higher or lower than that of the original energy level of the electron. The electron in same orbit exhibits different energy levels. The grouping of this different energy levels is called energy band. However, the energy levels of inner orbit electrons are not much affected by the presence of neighboring atoms. ADGEMA

Important energy bands in solids

There are number of energy bands in solids but three of them are very important. These three energy bands are important to understand the behavior of solids. These energy bands are

- Valence band
- Conduction band
- Forbidden band or forbidden gap



001:2015 & 14001:2015 Valence band

The energy band which is formed by grouping the range of energy levels of the valence electrons or outermost orbit electrons is called as valence band. Valence band is present below the conduction band as shown in figure. Electrons in the valence band have lower energy than the electrons in conduction band.

The electrons present in the valence band are loosely bound to the nucleus of atom.

• Conduction band

The energy band which is formed by grouping the range of energy levels of the free electrons is called as conduction band. Generally, the conduction band is empty but when external energy is applied the electrons in the valence band jumps in to the conduction band and becomes free electrons. Electrons in the conduction band have higher energy than the electrons in valence band.

The conduction band electrons are not bound to the nucleus of atom.

• Forbidden gap

The energy gap which is present between the valence band and conduction band by separating these two energy bands is called as forbidden band or forbidden gap. In solids, electrons cannot stay in forbidden gap because there is no allowed energy state in this region. Forbidden gap is the major factor for determining the electrical conductivity of a solid. The classification of materials as insulators, conductors and semiconductors is mainly depends on forbidden gap.

RHAGEM

The energy associated with forbidden band is called energy gap and it is measured in unit electron volt (eV).

$1 \text{ eV} = 1.6 \times 10^{-19} \text{ J}$

The applied external energy in the form of heat or light must be equal to to the forbidden gap in order to push an electron from valence band to the conduction band.

Classification of materials based on forbidden gap

Forbidden gap plays a major role for determining the electrical conductivity of material. Based on the forbidden gap materials are classified in to three types, they are

- Insulators
- Conductors
- semiconductors

• Insulators

The materials which does not allow the flow of electric current through them are called as insulators. Insulators are also called as poor conductors of electricity.



Normally, in insulators the valence band is fully occupied with electrons due to sharing of outer most orbit electrons with the neighboring atoms. Whereas conduction band is empty, I.e, no electrons are present in conduction band. The forbidden gap between the valence band and conduction band is very large in insulators. The energy gap of insulator is approximately equal to 15 electron volts (eV). The electrons in valence band electrons in to conduction band large amount of external energy is applied which is equal to the forbidden gap. But in insulators, this is practically impossible to move the valence band electrons in to conduction band.

Rubber, wood, diamond, plastic are some examples of insulators. Insulators such as plastics are used for coating of electrical wires. These insulators prevent the flow of electricity to unwanted points and protect us from electric shocks.

• Conductors

The materials which easily allow the flow of electric current through them are called as conductors. Metals such as copper, silver, iron, aluminum etc. are good conductors of electricity.



Copyright@2013-2014, Physics and Radio-Electronics, All rights reserved

In a conductor, valence band and conduction band overlap each other as shown in figure. Therefore, there is no forbidden gap in a conductor. A small amount of applied external energy provides enough energy for the valence band electrons to move in to conduction band. Therefore, more number of valence band electrons can easily moves in to the conduction band. When valence band electrons moves to conduction band they becomes free electrons. The electrons present in the conduction band are not attached to the nucleus of a atom. In conductors, large number of electrons are present in conduction band at room temperature, I.e, conduction band is almost full with electrons. Where as valence band is partially occupied with electrons. The electrons present in the conduction band moves freely by carrying the electric current from one point to other.

Semiconductors

The material which has electrical conductivity between that of a conductor and an insulator is called as semiconductor. Silicon, germanium and graphite are some examples of semiconductors.



Copyright@2013-2014, Physics and Radio-Electronics, All rights reserved

In semiconductors, the forbidden gap between valence band and conduction band is very small. It has a forbidden gap of about 1 electron volt (eV). At low temperature, the valence band is completely occupied with electrons and conduction band is empty because the electrons in the valence band do not have enough energy to move in to conduction band. Therefore, semiconductor behaves as an insulator at low temperature. However, at room temperature some of the electrons in valence band gains enough energy in the form of heat and moves in to conduction band. When the temperature is goes on increasing, the number of valence band electrons moving in to conduction band is also increases. This shows that electrical conductivity of the semiconductor increases with increase in temperature. I.e. a semiconductor has negative temperature co-efficient of resistance.

The resistance of semiconductor decreases with increase in temperature.

Distinction between metals, semiconductors and insulators

Distinction between metals, semiconductors and insulators, According to band theory, the electrons in a solid can possess bands of energies called allowed bands of energies and these electrons may not possess some other bands of energies called forbidden bands of energies. The allowed bands of energies and forbidden bands of energies are present alternatively one after another for the electrons of a solid. The top-most band is called Conduction band and the next band below Conduction band is valance band. These two bands are separated by forbidden bands.

1. Insulator: The valence band of those materials remains full of electrons. The conduction band of those materials remains empty. The forbidden energy gap between the conduction band and the valence band is widest. The difference is more than 4 eV. Crossing the

forbidden energy gap from valence band to conduction band require large amount of energy. Mica, glass, eboniote etc are the examples of insulators.

2. Conductors: The valence band and the conduction band overlap each other. There is no forbidden energy gap here so Eg=0. At absolute zero temperature large number of electrons remains in the conduction band. The resistance of conductor is very low; large number charge carriers are available here. So, the electricity can pass easily through the conductors. Aluminum, Silver, etc are good conductors

3. Semiconductors: A semiconductor remains partially full valence band and partially full conduction band at the room temperature The energy gap is narrower. The conduction band remains full empty of a semiconductor where the valence band remains full of electrons at absolute zero temperature. The value of Eg =1.1eV for silicon crystal and Eg=0.7eV for germanium. It can easily overcome due to thermal agitation or light. So, silicon and germanium are insulators at absolute zero temperature. On the other hand with the increasing of temperature the electrical conductivity of semiconductors increase.



Majority & minority carriers What is charge carrier?

Generally, carrier refers to any object that carry another object from one place to another place. For example, in countries such as India, Singapore and Brazil: Tiffin box or Tiffin carriers are widely used for carrying food from one place to another place. Here, the Tiffin box acts as a carrier that carries the food from one place to another place.

Let us take another example; People use vehicles such as buses, trains, airplanes, etc. to travel from one place to another place. Here, the vehicles act as carriers that carry people from one place to another place. In the similar way, particles such as free electrons and holes carry the charge or electric current from one place to another place.

Negative charge carriers

The negative charge carriers such as free electrons are the charge carriers that carry negative charge with them while moving from one place to another place. Free electrons are the electrons that are detached from the parent atom and moves freely from one place to another place.



Positive charge carriers

The positive charge carriers such as holes are the charge carriers that carry positive charge with them while moving from one place to another place. Holes are the vacancies in valence band that moves from one place to another place within the valence band.

Majority and minority charge carriers definition

The charge carriers that are present in large quantity are called majority charge carriers. The majority charge carriers carry most of the electric charge or electric current in the semiconductor. Hence, majority charge carriers are mainly responsible for electric current flow in the semiconductor.

The charge carriers that are present in small quantity are called minority charge carriers. The minority charge carriers carry very small amount of electric charge or electric current in the semiconductor.

Charge carriers in intrinsic semiconductor

The semiconductors that are in pure form are called intrinsic semiconductors. In intrinsic semiconductor the total number of negative charge carriers (free electrons) is equal to the total number of positive charge carriers (holes or vacancy).

Total negative charge carriers = Total positive charge carriers

Majority and minority charge carriers in n-type semiconductor

When the pentavalent atoms such as Phosphorus or Arsenic are added to the intrinsic semiconductor, an n-type semiconductor is formed. In n-type semiconductor, large number of free electrons is present. Hence, free electrons are the majority charge carriers in the n-type semiconductor. The free electrons (majority charge carriers) carry most of the electric charge or electric current in the n-type semiconductor.

In n-type semiconductor, very small number of holes is present. Hence, holes are the minority charge carriers in the n-type semiconductor. The holes (minority charge carriers) carry only a small amount of electric charge or electric current in the n-type semiconductor.

The total number of negative charge carriers (free electrons) in n-type semiconductor is greater than the total number of positive charge carriers (holes) in the n-type semiconductor.

Total negative charge carriers > Total positive charge carriers



Copyright © 2013-2015, Physics and Radio-Electronics, All rights reserved

Majority and minority charge carriers in p-type semiconductor

When the trivalent atoms such as Boron or Gallium are added to the intrinsic semiconductor, a <u>p-type semiconductor</u> is formed. In p-type semiconductor, large number of holes is present. Hence, holes are the majority charge carriers in the p-type semiconductor. The holes (majority charge carriers) carry most of the electric charge or electric current in the p-type semiconductor. In p-type semiconductor, very small number of free electrons is present. Hence, free electrons are the minority charge carriers in the p-type semiconductor. The free electrons (minority charge carriers) carry only a small amount of electric current in the p-type semiconductor. The total number of negative charge carriers (free electrons) in p-type semiconductor is less than the total number of positive charge carriers (holes) in the p-type semiconductor.

Total negative charge carriers < Total positive charge carriers

Charge Carriers in Semiconductors

When an electric field is applied to a metal, negatively charged electrons are accelerated and carry the resulting current. In a semiconductor the charge is not carried exclusively by electrons. Positively charged **holes** also carry charge. These may be viewed either as vacancies in the otherwise filled valence band, or equivalently as positively charged particles. Since the Fermi-Dirac distribution is a step function at absolute zero, pure semiconductors will have all the states in the valence bands filled with electrons and will be insulators at absolute zero. This is depicted in the **E-k** diagram below; shaded circles represent filled momentum states and empty circles unfilled momentum states. In this diagram **k**, rather than *k*, has been used to denote that the wave vector is actually a vector, i.e., a tensor of the first rank, rather than a scalar.



COPYRIGHT FIMT 2020

If the band gap is sufficiently small and the temperature is increased from absolute zero, some electrons may be thermally excited into the conduction band, creating an electron-hole pair. This is as a result of the smearing out of the Fermi-Dirac distribution at finite temperature. An electron may also move into the conduction band from the valence band if it absorbs a photon that corresponds to the energy difference between a filled state and an unfilled state. Any such photon must have an energy that is greater than or equal to the band gap between the valence band and the conduction band, as in the diagram below.



Whether thermally or photonically induced, the result is an electron in the conduction band and a vacant state in the valence band.



If an electric field is now applied to the material, all of the electrons in the solid will feel a force from the electric field. However, because no two electrons can be in the exact same quantum state, an electron cannot gain any momentum from the electric field unless there is a vacant momentum state adjacent to the state being occupied by the electron. In the above schematic, the electron in the conduction band can gain momentum from the electric field, as can an electron adjacent to the vacant state left behind in the valence band. In the diagram below, both of these electrons are shown moving to the right.



The result of this is that the electrons have some net momentum, and so there is an overall movement of charge. This slight imbalance of positive and negative momentum can be seen in the diagram below, and it gives rise to an electric current.



The vacant site in the valence band which has moved to the left can be viewed as being a particle which carries positive electric charge of equal magnitude to the electron charge. This is therefore a **hole**. It should be appreciated that these schematics do not represent electrons 'hopping' from site to site in real space, because the electrons are not localised to specific sites in space. These schematics are in momentum space. As such, holes should not be thought of as moving through the semiconductor like dislocations when metals are plastically deformed – it suffices to view them simply as particles which carry positive charge. The opposite process to the creation of an electron-hole pair is called **recombination**. This occurs when an electron drops down in energy from the conduction band to the valence band. Just as the creation of an electron-hole pair may be induced by a photon, recombination can produce a photon. This is the principle behind semiconductor optical devices such as light-emitting diodes (LEDs), in which the photons are light of visible wavelength.

Intrinsic and Extrinsic Semiconductors

In most pure semiconductors at room temperature, the population of thermally excited charge carriers is very small. Often the concentration of charge carriers may be orders of magnitude

lower than for a metallic conductor. For example, the number of thermally excited electrons cm^{-3} in silicon (Si) at 298 K is 1.5×10^{10} . In gallium arsenide (GaAs) the population is only 1.1×10^{6} electrons cm^{-3} . This may be compared with the number density of free electrons in a typical metal, which is of the order of 10^{28} electrons cm^{-3} . Given these numbers of charge carriers, it is no surprise that, when they are extremely pure, silicon and other semiconductors have high electrical resistivities, and therefore low electrical conductivities. This problem can be overcome by doping a semiconducting material with impurity atoms. Even very small controlled additions of impurity atoms at the 0.0001% level can make very large differences to the conductivity of a semiconductor. It is easiest to begin with a specific example. Silicon is a group IV element, and has 4 valence electrons per atom. In pure silicon the valence band is completely filled at absolute zero. At finite temperatures the only charge carriers are the electrons in the conduction band and the holes in the valence band that arise as a result of the thermal excitation of electrons to the conduction band. These charge carriers are called *intrinsic* charge carriers, and necessarily there are equal numbers of electrons and holes. Pure silicon is therefore an example of an **intrinsic semiconductor**.

If a very small number of atoms of a group V element such as phosphorus (P) are added to the silicon as substitutional atoms in the lattice, additional valence electrons are introduced into the material because each phosphorus atom has 5 valence electrons. These additional electrons are bound only weakly to their parent impurity atoms (the bonding energies are of the order of hundredths of an eV), and even at very low temperatures these electrons can be promoted into the conduction band of the semiconductor. This is often represented schematically in band diagrams by the addition of 'donor levels' just below the bottom of the conduction band, as in the schematic below.



The presence of the dotted line in this schematic does not mean that there now exist allowed energy states within the band gap. The dotted line represents the existence of additional electrons which may be easily excited into the conduction band. Semiconductors that have been doped in this way will have a surplus of electrons, and are called *n*-type semiconductors. In such semiconductors, electrons are the majority carriers.

Conversely, if a group III element, such as aluminium (Al), is used to substitute for some of the atoms in silicon, there will be a deficit in the number of valence electrons in the material. This introduces electron-accepting levels just above the top of the valence band, and causes more holes to be introduced into the valence band. Hence, the majority charge carriers are positive holes in this case. Semiconductors doped in this way are termed *p*-type semiconductors.



Doped semiconductors (either *n*-type or *p*-type) are known as **extrinsic semiconductors**. The activation energy for electrons to be donated by or accepted to impurity states is usually so low that at room temperature the concentration of majority charge carriers is similar to the concentration of impurities. It should be remembered that in an extrinsic semiconductor there is an contribution to the total number of charge carriers from intrinsic electrons and holes, but at room temperature this contribution is often very small in comparison with the number of charge carriers introduced by the controlled impurity doping of the semiconductor.

Doping

Doping means the introduction of impurities into a semiconductor crystal to the defined modification of conductivity. Two of the most important materials silicon can be doped with, are boron (3 valence electrons = 3-valent) and phosphorus (5 valence electrons = 5-valent). Other materials are aluminum, indium (3-valent) and arsenic, antimony (5-valent). The dopant is integrated into the lattice structure of the semiconductor crystal, the number of outer electrons define the type of doping. Elements with 3 valence electrons are used for p-type doping, 5-valued elements for n-doping. The conductivity of a deliberately contaminated silicon crystal can be increased by a factor of 10^6 .

n-doping

The 5-valent dopant has an outer electron more than the silicon atoms. Four outer electrons combine with ever one silicon atom, while the fifth electron is free to move and serves as charge carrier. This free electron requires much less energy to be lifted from the valence band into the conduction band, than the electrons which cause the intrinsic conductivity of silicon. The dopant, which emits an electron, is known as an electron donor (donare, lat. = to give).

The dopants are positively charged by the loss of negative charge carriers and are built into the lattice, only the negative electrons can move. Doped semimetals whose conductivity is based on free (negative) electrons are n-type or n-doped. Due to the higher number of free electrons those are also named as majority charge carriers, while free mobile holes are named as the minority charge carriers.

n-doping with phosphorus



Arsenic is used as an alternative to phosphorus, because its diffusion coefficient is lower. This means that the dopant diffusion during subsequent processes is less than that of phosphorus and thus the arsenic remains at the position where it was introduced into the lattice originally.

p-doping

In contrast to the free electron due to doping with phosphorus, the 3-valent dopant effect is exactly the opposite. The 3-valent dopants can catch an additional outer electron, thus leaving a hole in the valence band of silicon atoms. Therefore the electrons in the valence band become mobile. The holes move in the opposite direction to the movement of the electrons. The necessary energy to lift an electron into the energy level of indium as a dopant, is only 1 % of the energy which is needed to raise a valence electron of silicon into the conduction band.

With the inclusion of an electron, the dopant is negatively charged, such dopants are called acceptors (acceptare, lat. = to add). Again, the dopant is fixed in the crystal lattice, only the

positive charges can move. Due to positive holes these semiconductors are called pconductive or p-doped. Analog to n-doped semiconductors, the holes are the majority charge carriers, free electrons are the minority charge carriers.

p-doping with boron



The free place on the boron atom is filled with an electron. Therefore a new hole ("defect electron") is generated. This holes move in the opposite direction to the electrons

Doped semiconductors are electrically neutral. The terms n- and p-type doped do only refer to the majority charge carriers. Each positive or negative charge carrier belongs to a fixed negative or positive charged dopant. N- and p-doped semiconductors behave approximately equal in relation to the current flow. With increasing amount of dopants, the number of charge carriers increases in the semiconductor crystal. Here it requires only a very small amount of dopants. Weakly doped silicon crystals contain only 1 impurity per 1,000,000,000 silicon atoms, high doped semiconductors for example contain 1 foreign atom per 1,000 silicon atoms.

Electronic band structure in doped semiconductors

By the introduction of a dopant with five outer electrons, in n-doped semiconductors there is an electron in the crystal which is not bound and therefore can be moved with relatively little energy into the conduction band. Thus in n-doped semiconductors the donator energy level is close to the conduction band edge, the band gap to overcome is very small. Analog, through introduction of a 3-valent dopant in a semiconductor, a hole is available, which may be already occupied at low-energy by an electron from the valence band of the silicon. For pdoped semiconductors the acceptor energy level is close the valence band.

Band model of doped semiconductors



1400

COPYRIGHT FIMT 2020



N-Type Semiconductor



The addition of pentavalent impurities such as antimony, arsenic or phosphorus contributes free electrons, greatly increasing the conductivity of the intrinsic semiconductor. Phosphorus may be added by diffusion of phosphine gas (PH3).



P-Type Semiconductor

The addition of trivalent impurities such as boron, aluminum or gallium to an intrinsic semiconductor creates deficiencies of valence electrons, called "holes". It is typical to use B₂H₆ diborane gas to diffuse boron into the silicon material.



COPYRIGHT FIMT 2020

713 | Page



NAAC ACCREDITED

Bands for Doped Semiconductors

The application of <u>band theory</u> to <u>n-type</u> and <u>p-type</u> semiconductors shows that extra levels have been added by the impurities. In n-type material there are electron energy levels near the top of the band gap so that they can be easily excited into the conduction band. In p-type material, extra holes in the band gap allow excitation of valence band electrons, leaving mobile holes in the valence band.



PN Junction Diode

A PN-junction diode is formed when a p-type semiconductor is fused to an n-type semiconductor creating a potential barrier voltage across the diode junction.

The effect described in the previous tutorial is achieved without any external voltage being applied to the actual PN junction resulting in the junction being in a state of equilibrium. However, if we were to make electrical connections at the ends of both the N-type and the P-type materials and then connect them to a battery source, an additional energy source now exists to overcome the potential barrier. The effect of adding this additional energy source

results in the free electrons being able to cross the depletion region from one side to the other. The behaviour of the PN junction with regards to the potential barrier's width produces an asymmetrical conducting two terminal device, better known as the **PN Junction Diode**.

A *PN Junction Diode* is one of the simplest semiconductor devices around, and which has the characteristic of passing current in only one direction only. However, unlike a resistor, a diode does not behave linearly with respect to the applied voltage as the diode has an exponential current-voltage (I-V) relationship and therefore we can not described its operation by simply using an equation such as Ohm's law. If a suitable positive voltage (forward bias) is applied between the two ends of the PN junction, it can supply free electrons and holes with the extra energy they require to cross the junction as the width of the depletion layer around the PN junction is decreased. By applying a negative voltage (reverse bias) results in the free charges being pulled away from the junction resulting in the depletion layer width being increased. This has the effect of increasing or decreasing the effective resistance of the junction itself allowing or blocking current flow through the diode.

Then the depletion layer widens with an increase in the application of a reverse voltage and narrows with an increase in the application of a forward voltage. This is due to the differences in the electrical properties on the two sides of the PN junction resulting in physical changes taking place. One of the results produces rectification as seen in the PN junction diodes static I-V (current-voltage) characteristics. Rectification is shown by an asymmetrical current flow when the polarity of bias voltage is altered as shown below.

Junction Diode Symbol and Static I-V Characteristics



But before we can use the PN junction as a practical device or as a rectifying device we need to firstly **bias** the junction, ie connect a voltage potential across it. On the voltage axis above, "Reverse Bias" refers to an external voltage potential which increases the potential barrier.

An external voltage which decreases the potential barrier is said to act in the "Forward Bias" direction.

There are two operating regions and three possible "biasing" conditions for the standard **Junction Diode** and these are:

- 1. Zero Bias No external voltage potential is applied to the PN junction diode.
- 2. Reverse Bias The voltage potential is connected negative, (-ve) to the P-type material and positive, (+ve) to the N-type material across the diode which has the effect of Increasing the PN junction diode's width.
- 3. Forward Bias The voltage potential is connected positive, (+ve) to the P-type material and negative, (-ve) to the N-type material across the diode which has the effect of **Decreasing** the PN junction diodes width.

Zero Biased Junction Diode

When a diode is connected in a **Zero Bias** condition, no external potential energy is applied to the PN junction. However if the diodes terminals are shorted together, a few holes (majority carriers) in the P-type material with enough energy to overcome the potential barrier will move across the junction against this barrier potential. This is known as the "**Forward Current**" and is referenced as I_F

Likewise, holes generated in the N-type material (minority carriers), find this situation favourable and move across the junction in the opposite direction. This is known as the "**Reverse Current**" and is referenced as I_R . This transfer of electrons and holes back and forth across the PN junction is known as diffusion, as shown below.



The potential barrier that now exists discourages the diffusion of any more majority carriers across the junction. However, the potential barrier helps minority carriers (few free electrons

in the P-region and few holes in the N-region) to drift across the junction. Then an "Equilibrium" or balance will be established when the majority carriers are equal and both moving in opposite directions, so that the net result is zero current flowing in the circuit. When this occurs the junction is said to be in a state of "**Dynamic Equilibrium**".

The minority carriers are constantly generated due to thermal energy so this state of equilibrium can be broken by raising the temperature of the PN junction causing an increase in the generation of minority carriers, thereby resulting in an increase in leakage current but an electric current cannot flow since no circuit has been connected to the PN junction.

Reverse Biased PN Junction Diode

When a diode is connected in a **Reverse Bias** condition, a positive voltage is applied to the N-type material and a negative voltage is applied to the P-type material. The positive voltage applied to the N-type material attracts electrons towards the positive electrode and away from the junction, while the holes in the P-type end are also attracted away from the junction towards the negative electrode. The net result is that the depletion layer grows wider due to a lack of electrons and holes and presents a high impedance path, almost an insulator. The result is that a high potential barrier is created thus preventing current from flowing through the semiconductor material.





This condition represents a high resistance value to the PN junction and practically zero current flows through the junction diode with an increase in bias voltage. However, a very small **leakage current** does flow through the junction which can be measured in micro-amperes, (μA). One final point, if the reverse bias voltage Vr applied to the diode is increased to a sufficiently high enough value, it will cause the diode's PN junction to overheat and fail due to the avalanche effect around the junction. This may cause the diode to

COPYRIGHT FIMT 2020

become shorted and will result in the flow of maximum circuit current, and this shown as a step downward slope in the reverse static characteristics curve below.



Reverse Characteristics Curve for a Junction Diode

Sometimes this avalanche effect has practical applications in voltage stabilising circuits where a series limiting resistor is used with the diode to limit this reverse breakdown current to a preset maximum value thereby producing a fixed voltage output across the diode. These types of diodes are commonly known as Zener Diodes and are discussed in a later tutorial.

Forward Biased PN Junction Diode

When a diode is connected in a **Forward Bias** condition, a negative voltage is applied to the N-type material and a positive voltage is applied to the P-type material. If this external voltage becomes greater than the value of the potential barrier, approx. 0.7 volts for silicon and 0.3 volts for germanium, the potential barriers opposition will be overcome and current will start to flow. This is because the negative voltage pushes or repels electrons towards the junction giving them the energy to cross over and combine with the holes being pushed in the opposite direction towards the junction by the positive voltage. This results in a characteristics curve of zero current flowing up to this voltage point, called the "knee" on the static curves and then a high current flow through the diode with little increase in the external voltage as shown below.

:2015 & 14001:20

Forward Characteristics Curve for a Junction Diode



The application of a forward biasing voltage on the junction diode results in the depletion layer becoming very thin and narrow which represents a low impedance path through the junction thereby allowing high currents to flow. The point at which this sudden increase in current takes place is represented on the static I-V characteristics curve above as the "knee" point.

Reduction in the Depletion Layer due to Forward Bias



This condition represents the low resistance path through the PN junction allowing very large currents to flow through the diode with only a small increase in bias voltage. The actual potential difference across the junction or diode is kept constant by the action of the depletion layer at approximately 0.3v for germanium and approximately 0.7v for silicon junction diodes. Since the diode can conduct "infinite" current above this knee point as it effectively becomes a short circuit, therefore resistors are used in series with the diode to limit its current flow. Exceeding its maximum forward current specification causes the device to dissipate more power in the form of heat than it was designed for resulting in a very quick failure of the device.

Junction Diode Summary

The PN junction region of a Junction Diode has the following important characteristics:

- Semiconductors contain two types of mobile charge carriers, "Holes" and "Electrons".
- The holes are positively charged while the electrons negatively charged.

- A semiconductor may be doped with donor impurities such as Antimony (N-type doping), so that it contains mobile charges which are primarily electrons.
- A semiconductor may be doped with acceptor impurities such as Boron (P-type doping), so that it contains mobile charges which are mainly holes.
- The junction region itself has no charge carriers and is known as the depletion region.
- The junction (depletion) region has a physical thickness that varies with the applied voltage.
- When a diode is Zero Biased no external energy source is applied and a natural Potential Barrier is developed across a depletion layer which is approximately 0.5 to 0.7v for silicon diodes and approximately 0.3 of a volt for germanium diodes.
- When a junction diode is **Forward Biased** the thickness of the depletion region reduces and the diode acts like a short circuit allowing full current to flow.
- When a junction diode is **Reverse Biased** the thickness of the depletion region increases and the diode acts like an open circuit blocking any current flow, (only a very small leakage current).

We have also seen above that the diode is two terminal non-linear device whose I-V characteristic are polarity dependent as depending upon the polarity of the applied voltage, V_D the diode is either *Forward Biased*, $V_D > 0$ or *Reverse Biased*, $V_D < 0$. Either way we can model these current-voltage characteristics for both an ideal diode and for a real silicon diode as shown:

Junction Diode Ideal and Real Characteristics


Light Emitting Diode (LED) What is light?

Before going into how LED works, let's first take a brief look at light self. Since ancient times man has obtained light from various sources like sunrays, candles and lamps.

In 1879, Thomas Edison invented the incandescent light bulb. In the light bulb, an electric current is passed through a filament inside the bulb.

When sufficient current is passed through the filament, it gets heated up and emits light. The light emitted by the filament is the result of electrical energy converted into heat energy which in turn changes into light energy.



Unlike the light bulb in which electrical energy first converts into heat energy, the electrical energy can also be directly converted into light energy.

In Light Emitting Diodes (LEDs), electrical energy flowing through it is directly converted into light energy.

Light is a type of <u>energy</u> that can be released by an <u>atom</u>. Light is made up of many small particles called photons. Photons have energy and momentum but no mass.

Atoms are the basic building blocks of matter. Every object in the universe is made up of atoms. Atoms are made up of small particles such as electrons, protons and neutrons.

Electrons are negatively charged, protons are positively charged, and neutrons have no charge.

The attractive force between the protons and neutrons makes them stick together to form nucleus. Neutrons have no charge. Hence, the overall charge of the nucleus is positive.



Physics and Radio-Electronics

The negatively charged electrons always revolve around the positively charged nucleus because of the electrostatic force of attraction between them. Electrons revolve around the nucleus in different orbits or shells. Each orbit has different energy level.

For example, the electrons orbiting very close to the nucleus have low energy whereas the electrons orbiting farther away from the nucleus have high energy.

The electrons in the lower energy level need some additional energy to jump into the higher energy level. This additional energy can be supplied by the outside source. When electrons orbiting the nucleus gains energy from outside source they jump into higher orbit or higher energy level.

The electrons in the higher energy level will not stay for long period. After a short period, the electrons fall back to lower energy level. The electrons which jump from higher energy level to lower energy level will releases energy in the form of a photon or light. In some materials, this energy lose is released mostly in the form of heat. The electron which loses greater energy will releases a greater energy photon.

What is Light Emitting Diode (LED)?

Light Emitting Diodes (LEDs) are the most widely used semiconductor diodes among all the different types of semiconductor diodes available today. Light emitting diodes emit

either visible light or invisible infrared light when forward biased. The LEDs which emit invisible infrared light are used for remote controls.

A light Emitting Diode (LED) is an optical semiconductor device that emits light when voltage is applied. In other words, LED is an optical semiconductor device that converts electrical energy into light energy.



When Light Emitting Diode (LED) is forward biased, free electrons in the conduction band recombines with the holes in the valence band and releases energy in the form of light. The process of emitting light in response to the strong electric field or flow of electric current is called electroluminescence. A normal p-n junction diode allows electric current only in one direction. It allows electric current when forward biased and does not allow electric current when reverse biased. Thus, normal p-n junction diode operates only in forward bias condition. Like the normal p-n junction diodes, LEDs also operates only in forward bias condition. To create an LED, the n-type material should be connected to the negative terminal of the battery and p-type material should be connected to the positive terminal of the battery. In other words, the n-type material should be negatively charged and the p-type material should be positively charged.

The construction of LED is similar to the normal p-n junction diode except that gallium, phosphorus and arsenic materials are used for construction instead of silicon or germanium materials. In normal p-n junction diodes, silicon is most widely used because it is less sensitive to the temperature. Also, it allows electric current efficiently without any damage. In some cases, germanium is used for constructing diodes. However, silicon or germanium diodes do not emit energy in the form of light. Instead, they emit energy in the form of heat. Thus, silicon or germanium is not used for constructing LEDs.

Layers of LED

A Light Emitting Diode (LED) consists of three layers: p-type semiconductor, n-type semiconductor and depletion layer. The p-type semiconductor and the n-type semiconductor are separated by a depletion region or depletion layer.

P-type semiconductor

When trivalent impurities are added to the intrinsic or pure semiconductor, a p-type semiconductor is formed In p-type semiconductor, holes are the majority charge carriers and free electrons are the minority charge carriers. Thus, holes carry most of the electric current in p-type semiconductor.

N-type semiconductor

When pentavalent impurities are added to the intrinsic semiconductor, an n-type semiconductor is formed. In n-type semiconductor, free electrons are the majority charge carriers and holes are the minority charge carriers. Thus, free electrons carry most of the electric current in n-type semiconductor.

Depletion layer or region

Depletion region is a region present between the p-type and n-type semiconductor where no mobile charge carriers (free electrons and holes) are present. This region acts as barrier to the electric current. It opposes flow of electrons from n-type semiconductor and flow of holes from p-type semiconductor.

To overcome the barrier of depletion layer, we need to apply voltage which is greater than the barrier potential of depletion layer. If the applied voltage is greater than the barrier potential of the depletion layer, the electric current starts flowing.

How Light Emitting Diode (LED) works?

Light Emitting Diode (LED) works only in forward bias condition. When Light Emitting Diode (LED) is forward biased, the free electrons from n-side and the holes from p-side are pushed towards the junction.

When free electrons reach the junction or depletion region, some of the free electrons recombine with the holes in the positive ions. We know that positive ions have less number of electrons than protons. Therefore, they are ready to accept electrons. Thus, free electrons recombine with holes in the depletion region. In the similar way, holes from p-side recombine with electrons in the depletion region.



Light Emitting Diode (LED)
Physics and Radio-Electronics

Because of the recombination of free electrons and holes in the depletion region, the <u>width</u> of depletion region decreases. As a result, more charge carriers will cross the <u>p-n junction</u>. Some of the charge carriers from p-side and n-side will cross the p-n junction before they recombine in the depletion region. For example, some free electrons from n-type semiconductor cross the p-n junction and recombines with holes in p-type semiconductor. In the similar way, holes from p-type semiconductor cross the p-n junction and recombines with free electrons in the n-type semiconductor. Thus, recombination takes place in depletion region as well as in p-type and n-type semiconductor. The free electrons in the conduction band releases energy in the form of light before they recombine with holes in the valence band. In silicon and germanium diodes, most of the energy is released in the form of heat and emitted light is too small. However, in materials like gallium arsenide and gallium phosphide the emitted photons have sufficient energy to produce intense visible light.

How LED emits light?

When external voltage is applied to the valence electrons, they gain sufficient energy and breaks the bonding with the parent atom. The valence electrons which breaks bonding with

the parent atom are called free electrons. When the valence electron left the parent atom, they leave an empty space in the valence shell at which valence electron left. This empty space in the valence shell is called a hole. The energy level of all the valence electrons is almost same. Grouping the range of energy levels of all the valence electrons is called valence band. In the similar way, energy level of all the free electrons is called conduction band. The energy level of free electrons in the conduction band is high compared to the energy level of valence electrons or holes in the valence band. Therefore, free electrons in the conduction band need to lose energy in order to recombine with the holes in the valence band. The free electrons in the conduction band do not stay for long period. After a short period, the free electrons lose energy in the form of light and recombine with the holes in the valence band. Each recombination of charge carrier will emit some light energy.

The energy lose of free electrons or the intensity of emitted light is depends on the forbidden gap or energy gap between conduction band and valence band. The semiconductor device with large forbidden gap emits high intensity light whereas the semiconductor device with small forbidden gap emits low intensity light. In other words, the brightness of the emitted light is depends on the material used for constructing LED and forward current flow through the LED. In normal silicon diodes, the energy gap between conduction band and valence band is less. Hence, the electrons fall only a short distance. As a result, low energy photons are released. These low energy photons have low frequency which is invisible to human eye. In LEDs, the energy gap between conduction band and valence, the free electrons fall to a large distance. As a result, high energy photons are released. These high energy photons have high frequency which is visible to human eye. The efficiency of generation of light in LED increases with increase in injected current and with a decrease in temperature.

In light emitting diodes, light is produced due to recombination process. Recombination of charge carriers takes place only under forward bias condition. Hence, LEDs operate only in forward bias condition. When light emitting diode is reverse biased, the free electrons

(majority carriers) from n-side and holes (majority carriers) from p-side moves away from the junction. As a result, the width of depletion region increases and no recombination of charge carriers occur. Thus, no light is produced. If the reverse bias voltage applied to the LED is highly increased, the device may also be damaged. All diodes emit photons or light but not all diodes emit visible light. The material in an LED is selected in such a way that the wavelength of the released photons falls within the visible portion of the light spectrum.

Light emitting diodes can be switched ON and OFF at a very fast speed of 1 ns.

Light emitting diode (LED) symbol

The symbol of LED is similar to the normal p-n junction diode except that it contains arrows pointing away from the diode indicating that light is being emitted by the diode.



LEDs are available in different colors. The most common colors of LEDs are orange, yellow, green and red.

The schematic symbol of LED does not represent the color of light. The schematic symbol is same for all colors of LEDs. Hence, it is not possible to identify the color of LED by seeing its symbol.

LED construction

One of the methods used to construct LED is to deposit three semiconductor layers on the substrate. The three semiconductor layers deposited on the substrate are n-type semiconductor, p-type semiconductor and active region. Active region is present in between the n-type and p-type semiconductor layers.



When LED is forward biased, free electrons from n-type semiconductor and holes from ptype semiconductor are pushed towards the active region.

When free electrons from n-side and holes from p-side recombine with the opposite charge carriers (free electrons with holes or holes with free electrons) in active region, an invisible or visible light is emitted. In LED, most of the charge carriers recombine at active region. Therefore, most of the light is emitted by the active region. The active region is also called as depletion region.

Biasing of LED

The safe forward voltage ratings of most LEDs is from 1V to 3 V and forward current ratings is from 200 mA to 100 mA.

If the voltage applied to LED is in between 1V to 3V, LED works perfectly because the current flow for the applied voltage is in the operating range. However, if the voltage applied to LED is increased to a value greater than 3 volts. The depletion region in the LED breaks down and the electric current suddenly rises. This sudden rise in current may destroy the device.

To avoid this we need to place a <u>resistor</u> (R_s) in series with the LED. The resistor (R_s) must be placed in between voltage source (Vs) and LED.



Physics and Radio-Electronics

The resistor placed between LED and voltage source is called current limiting resistor. This resistor restricts extra current which may destroy the LED. Thus, current limiting resistor protects LED from damage.

The current flowing through the LED is mathematically written as

$$I_F = \frac{V_s - V_D}{R_s}$$

Where,

I_F = Forward current

V_s = Source voltage or supply voltage

V_D = Voltage drop across LED

R_s = Resistor or current limiting resistor

Voltage drop is the amount of voltage wasted to overcome the depletion region barrier (which leads to electric current flow).

The voltage drop of LED is 2 to 3V whereas silicon or germanium diode is 0.3 or 0.7 V.

Therefore, to operate LED we need to apply greater voltage than silicon or germanium diodes.

COPYRIGHT FIMT 2020

Light emitting diodes consume more energy than silicon or germanium diodes to operate.

Output characteristics of LED

The amount of output light emitted by the LED is directly proportional to the amount of forward current flowing through the LED. More the forward current, the greater is the emitted output light. The graph of forward current vs output light is shown in the figure.



Visible LEDs and invisible LEDs

LEDs are mainly classified into two types: visible LEDs and invisible LEDs.

Visible LED is a type of LED that emits visible light. These LEDs are mainly used for display or illumination where LEDs are used individually without photosensors.

Invisible LED is a type of LED that emits invisible light (infrared light). These LEDs are mainly used with photosensors such as photodiodes.

What determines the color of an LED?

The material used for constructing LED determines its color. In other words, the wavelength or color of the emitted light depends on the forbidden gap or energy gap of the material.

Different materials emit different colors of light. Gallium arsenide LEDs emit red and infrared light. Gallium nitride LEDs emits bright blue light. Yttrium aluminium garnet LEDs emit white light. Gallium phosphide LEDs emit red, yellow and green light. Aluminium

gallium nitride LEDs emit ultraviolet light. Aluminum gallium phosphide LEDs emit green light.

Advantages of LED

- The brightness of light emitted by LED is depends on the current flowing through the LED. Hence, the brightness of LED can be easily controlled by varying the current. This makes possible to operate LED displays under different ambient lighting conditions.
- 2. Light emitting diodes consume low energy.
- 3. LEDs are very cheap and readily available.
- 4. LEDs are light in weight.
- 5. Smaller size.
- 6. LEDs have longer lifetime.
- 7. LEDs operates very fast. They can be turned on and off in very less time.
- 8. LEDs do not contain toxic material like mercury which is used in fluorescent lamps.
- 9. LEDs can emit different colors of light.

Disadvantages of LED

1. LEDs need more power to operate than normal p-n junction diodes.

9001:2015 & 140

2. Luminous efficiency of LEDs is low.

Applications of LED

The various applications of LEDs are as follows

- 1. Burglar alarms systems
- 2. Calculators
- 3. Picture phones
- 4. Traffic signals
- 5. Digital computers
- 6. Multimeters
- 7. Microprocessors
- 8. Digital watches

- 9. Automotive heat lamps
- 10. Camera flashes
- 11. Aviation lighting

Types of Diodes

The various types of diodes are as follows:

- 1. Zener diode
- 2. Avalanche diode
- 3. Photodiode
- 4. Light Emitting Diode
- 5. Laser diode
- 6. <u>Tunnel diode</u>
- 7. Schottky diode
- 8. Varactor diode
- 9. <u>P-N junction diode</u>

Transistors – Basics, Types & Biasing Modes Introduction to Transistor:

ACCREDITED

Earlier, the critical and important component of an electronic device was a vacuum tube; it is an electron tube used to control electric current. The vacuum tubes worked but they are bulky, require higher operating voltages, high power consumption, yield lower efficiency and cathode electron-emitting materials are used up in operation. So, that ended up as heat which shortened the life of the tube itself. To overcome these problems, John Bardeen, Walter Brattain and William Shockley were invented a transistor at Bell Labs in the year of 1947. This new device was a much more elegant solution to overcome many of the fundamental limitations of vacuum tubes.

Transistor is a semiconductor device that can both conduct and insulate. A transistor can act as a switch and an amplifier. It converts audio waves into electronic waves and resistor, controlling electronic current. Transistors have very long life, smaller in size, can operate on lower voltage supplies for greater safety and required no filament current. The first transistor was fabricated with germanium. A transistor performs the same function as a vacuum tube triode, but using semiconductor junctions instead of heated electrodes in a vacuum chamber. It is the fundamental building block of modern electronic devices and found everywhere in modern electronic systems.

Transistor Basics:

A transistor is a three terminal device. Namely,

- Base: This is responsible for activating the transistor.
- Collector: This is the positive lead.
- Emitter: This is the negative lead.

The basic idea behind a transistor is that it lets you control the flow of current through one channel by varying the intensity of a much smaller current that's flowing through a second channel.

OT DESIGN

Types of Transistors:

There are two types of transistors in present; they are bipolar junction transistor (BJT), field effect transistors (FET). A small current is flowing between the base and the emitter; base terminal can control a larger current flow between the collector and the emitter terminals. For a field-effect transistor, it also has the three terminals, they are gate, source, and drain, and a voltage at the gate can control a current between source and drain. The simple diagrams of BJT and FET are shown in figure below:



Field Effect Transistors(FET)

Bipolar Junction Transistor:

A Bipolar Junction Transistor (BJT) has three terminals connected to three doped semiconductor regions. It comes with two types, P-N-P and N-P-N. P-N-P transistor,

consisting of a layer of N-doped semiconductor between two layers of P-doped material. The base current entering in the collector is amplified at its output.

That is when PNP transistor is ON when its base is pulled low relative to the emitter. The arrows of PNP transistor symbol the direction of current flow when the device is in forward active mode. N-P-N transistor consisting a layer of P-doped semiconductor between two layers of N-doped material. By amplifying current the base we get the high collector and emitter current.

That is when NPN transistor is ON when its base is pulled low relative to the emitter. When the transistor is in ON state, current flow is in between the collector and emitter of the transistor. Based on minority carriers in P-type region the electrons moving from emitter to collector. It allows the greater current and faster operation; because of this reason most bipolar transistors used today are NPN.



• Field Effect Transistor (FET):

The field-effect transistor is a unipolar transistor, N-channel FET or P-channel FET are used for conduction. The three terminals of FET are source, gate and drain. The basic n-channel and p-channel FET's are shown above. For an n-channel FET, the device is constructed from n-type material. Between the source and drain then-type material acts as a resistor.

This transistor controls the positive and negative carriers with respect to holes or electrons. FET channel is formed by moving of positive and negative charge carriers. The channel of FET which is made by silicon. There are many types of FET's, MOSFET, JFET and etc. The applications of FET's are in low noise amplifier, buffer amplifier and analog switch.

Bipolar Junction Transistor Biasing



Transistors are the most important semiconductor active devices essential for almost all circuits. They are used as electronic switches, amplifiers etc in circuits. Transistors may be NPN, PNP, FET, JFET etc which have different functions in electronic circuits. For the proper working of the circuit, it is necessary to bias the transistor using resistor networks. Operating point is the point on the output characteristics that shows the Collector-Emitter voltage and the Collector current with no input signal. The Operating point is also known as the Bias point or Q-Point (Quiescent point).

Biasing is referred to provide resistors, capacitors or supply voltage etc to provide proper operating characteristics of the transistors. DC biasing is used to obtain DC collector current at a particular collector voltage. The value of this voltage and current are expressed in terms of the Q-Point. In a transistor amplifier configuration, the IC (max) is the maximum current that can flow through the transistor and VCE (max) is the maximum voltage applied across the device. To work the transistor as an amplifier, a load resistor RC must be connected to the collector. Biasing set the DC operating voltage and current to the correct level so that the AC input signal can be properly amplified by the transistor. The correct biasing point is somewhere between the fully ON or fully OFF states of the transistor. This central point is the Q-Point and if the transistor is properly biased, the Q-point will be the central operating point of the transistor. This helps the output current to increase and decrease as the input signal swings through the complete cycle.

For setting the correct Q-Point of the transistor, a collector resistor is used to set the collector current to a constant and steady value without any signal in its base. This steady DC

operating point is set by the value of the supply voltage and the value of the base biasing resistor. Base bias resistors are used in all the three transistor configurations like common base, common collector and Common emitter configurations.



Fig.1 Current Biasing

Fig.2 Feedback Biasing

Fig.3 Double Feedback Biasing



Modes of biasing:

Following are the different modes of transistor base biasing:

1. Current biasing:

As shown in the Fig.1, two resistors RC and RB are used to set the base bias. These resistors establish the initial operating region of the transistor with a fixed current bias. The transistor forward biases with a positive base bias voltage through RB. The forward base-Emitter voltage drop is 0.7 volts. Therefore the current through RB is $I_B = (V_{cc} - V_{BE}) / I_B$

2. Feedback biasing:

Fig.2 shows the transistor biasing by the use of a feedback resistor. The base bias is obtained from the collector voltage. The collector feedback ensures that the transistor is always biased in the active region. When the collector current increases, the voltage at the collector drops.

This reduces the base drive which in turn reduces the collector current. This feedback configuration is ideal for transistor amplifier designs.

3. Double Feedback Biasing:

Fig.3 shows how the biasing is achieved using double feedback resistors. By using two resistors RB1 and RB2 increases the stability with respect to the variations in Beta by increasing the current flow through the base bias resistors. In this configuration, the current in RB1 is equal to 10 % of the collector current.

4. Voltage Dividing Biasing:

Fig.4 shows the Voltage divider biasing in which two resistors RB1 and RB2 are connected to the base of the transistor forming a voltage divider network. The transistor gets biases by the voltage drop across RB2. This kind of biasing configuration is used widely in amplifier circuits.

5. Double Base Biasing:

Fig.5 shows a double feedback for stabilization. It uses both Emitter and Collector base feedback to improve the stabilization through controlling the collector current. Resistor values should be selected so as to set the voltage drop across the Emitter resistor 10% of the supply voltage and the current through RB1, 10% of the collector current.

Advantages of Transistor:

- 1. Smaller mechanical sensitivity.
- 2. Lower cost and smaller in size, especially in small-signal circuits.
- 3. Low operating voltages for greater safety, lower costs and tighter clearances.
- 4. Extremely long life.
- 5. No power consumption by a cathode heater.
- 6. Fast switching.

Difference Between NPN and PNP Transistor

The transistors PNP and NPN are BJTs and it is a basic electrical component, used in various <u>electrical and electronic circuits to build the projects</u>. The operation of the PNP and

NPN transistors mainly utilizes holes and electrons. These transistors can be used as amplifiers, switches and oscillators. In PNP transistor, the majority charge carriers are holes, where in NPN the majority charge carriers are electrons. Except, <u>FETs have only one sort of charge carrier</u>. The major difference between NPN and PNP transistor is, an NPN transistor gets the power when the flow of current runs through the base terminal of the transistor.

In NPN transistor, the flow of current runs from the collector terminal to the emitter terminal. A PNP transistor switches ON, when there is no flow of current at the base terminal of the transistor. In PNP transistor, the flow of current runs from the emitter terminal to the collector terminal. As a result, a PNP transistor switch ON by a low signal, where NPN transistor switches ON by a high signal.



Difference between PNP and NPN

Difference between NPN and PNP Transistor

The main difference between <u>NPN and PNP transistors</u> includes what are PNP and NPN transistors, construction, working and its applications.

What is a PNP Transistor?

The term 'PNP' stands for positive, negative, positive and also known as sourcing. The PNP transistor is a BJT; in this transistor the letter 'P' specifies the polarity of the voltage necessary for the emitter terminal. The second letter 'N' specifies the polarity of the base terminal. In this kind of transistor, the majority charge carriers are holes. Mainly, this transistor works as the same as the NPN transistor.



PNP Transistor

The required materials used to build the emitter (E), base (B) and collector(C) terminals in this transistor are diverse from those used in the NPN transistor. The BC terminals of this transistor are constantly reversed biased, then the –Ve voltage should be used for the collector terminal. Consequently, the base-terminal of the PNP transistor must be –Ve with respect to the emitter- terminal, and the collector terminal must be –Ve than the base terminal

NAAC ACCREDIT

PNP Transistor Construction

The PNP transistor construction is shown below. The main characteristics of both the transistors are similar except that the biasing of the current & voltage directions are inverted for any one of the achievable 3-configurations namely common base, common emitter and common collector.



PNP Transistor Construction

The voltage between the VBE (base and emitter terminal) is -Ve at the base terminal & +Ve at the emitter terminal. Since for this transistor, the base terminal constantly biased -Ve with respect to the emitter terminal. Also, the VBE is positive with respect to the collector VCE.

The voltage sources connected to this transistor is shown in the above figure. The emitter terminal is connected to the 'Vcc' with the load resistor 'RL'. This resistor stops the current

flow through the device, which is allied to the collector terminal. The base voltage 'VB' is connected to the 'RB' base resistor, which is biased negative with respect to the emitter terminal. To root the base current to flow through a PNP transistor, the base terminal of the transistor should be more negative than the base terminal by approximately 0.7volts (or) a Si device.

The **primary difference between PNP and NPN transistor** is the correct biasing of the transistor joints. The directions of current and the voltage polarities are constantly reverse to each other.

What is an NPN Transistor?

The term 'NPN' stands for negative, positive, negative and also known as sinking. **The NPN transistor is a BJT**, in this transistor, the initial letter 'N' specifies a negatively charged coating of the material. Where, 'P' specifies a completely charged layer. The two transistors have a positive layer, which are situated in the middle of two negative layers. Generally, NPN transistor is used in various electrical circuits for switching and strengthens the signals that exceed through them.



NPN Transistor

The NPN transistor includes three terminals like base, emitter and collector. These three terminals can be used to connect the transistor to the circuit board. When the current flows through this transistor, the base terminal of the transistor gets the electrical signal. The collector terminal creates a **stronger electric current**, and the emitter terminal exceeds this stronger current on to the circuit. In PNP transistor, the current runs through the collector to the emitter terminal.

Usually, NPN transistor is used because it is so simple to generate. For an NPN transistor to function properly, it requires to be created from a semiconductor object, which holds some

current. But not the max amount as extremely conductive materials such as metal. Silicon is one of the most normally used in semiconductors. These transistors are the simple transistors to build out of silicon. The NPN transistor is used on a computer circuit board to translate the information into binary code, and this procedure is proficient through a plethora of tiny switches flipping On & OFF on the boards. A powerful electric signal twists the switch on, while a lack of a signal makes the switch off.

Construction of NPN Transistor

The construction of this transistor is shown below. The voltage at the transistor's base is +Ve and –Ve at the transistors emitter terminal. The base terminal of the transistor is positive at all times with respect to the emitter, and also collector voltage supply is +Ve with respect to the transistor's emitter terminal. In this transistor, the collector terminal is linked to the VCC through the RL



NPN Transistor Construction

This resistor restricts the current flow through the highest base current. In NPN transistor, the electrons flow through the base represents transistor action. The main characteristic of this transistor action is the connection between the i/p and o/p circuits. Because, the amplifying properties of transistor come from the resultant control that the base utilizes upon the collector to emitter current.

The NPN transistor is a current activated device. When the transistor is turned ON, the huge current IC supplies between the collector & emitter terminals in the transistor. But, this only occurs when a tiny biasing current 'Ib' flows through the transistor's base terminal. It is a bipolar transistor; the current is the relation of two currents (Ic/Ib), named the DC current

gain of the device. It is specified with "hfe" or these days beta. The beta value can be huge up to 200 for typical transistors. When the NPN transistor is used in an active region, then base current 'Ib' offers the i/p and collector current 'IC' gives the o/p. The current gain of the NPN transistor from the C to the Eis called alpha (Ic/Ie), and it is a purpose of the transistor itself. As the Ie (emitter current) is the sum of a tiny base current and huge collector current. The worth of the alpha is very close to unity, and for a typical low power signal transistor the value ranges from about 0.950- 0.999.

Main Difference Between PNP and NPN

PNP and NPN transistors are three terminal device, which are made up of doped materials, frequently used in switching and amplifying applications. There are a combined of PN junction diodes in every bipolar junction transistor. When the couple of diodes connected, then it shapes a sandwich. That seat a kind of semiconductor in the middle of the similar two types.



Difference between NPN and PNP Transistor

So, there are only two kinds of bipolar sandwich, that are namely PNP & NPN. In semiconductor devices, the NPN transistor has typically high electron mobility evaluated to the mobility of a hole. Thus, it allows a huge amount of current & works very fast. And also, the construction of this transistor is simple from silicon.

 Both the transistors are collected of special materials and the flow of current in these transistors is also different.

- In an NPN transistor, the flow current runs from the collector terminal to the Emitter terminal, whereas in a PNP, the flow of current runs from the emitter terminal to the collector terminal.
- PNP transistor is made up of two P-type material layers with a layer of sandwiched of Ntype. The NPN transistor is made up of two N-type material layers with a layer of sandwiched of P-type.
- In an NPN-transistor, a +ve voltage is set to the collector terminal to generate a flow of current from the collector. For PNP transistor, a +ve voltage is set to the emitter terminal to generate flow of current from the emitter terminal to collector.
- The main working principle of an NPN transistor is, when the current is increased to the base terminal, then the transistor switches ON & it performs fully from the collector terminal to emitter terminal.
- When you reduce the current to the base, the transistor switches ON and the flow of current is so low. The transistor no longer works across the collector terminal to emitter terminal, and turns OFF.
- The main working principle of a PNP transistor is, when the current exists at the base of the PNP transistor, and then the transistor turns OFF. When there is no flow of current at the base of the transistor, then the transistor switches ON.

Benefits or advantages of Transistor

Following are the benefits or advantages of Transistor:

➡Input impedance is highest and output impedance is lowest for common collector BJT amplifier. Darlington pair is used where very high impedance is needed. Darlington pair offers

very high current gain.

- \Rightarrow It is used as current controlled current source.
- ➡It is used for fast switching applications.
- \blacksquare It is available at very low cost.
- \blacksquare It is very smaller in size.

➡It has longer life.

Ht uses low voltage for its operation. Hence it offers more safey.

There is no power consumption by cathode heater.

Drawbacks or disadvantages of Transistor

Following are the drawbacks or disadvantages of Transistor:
→Due to its small size, it is difficult to trace out faulty ones due to failure. Moreover it is very difficult to unsolder and replace new ones.

Manufacturing techniques are very complex and requires clean room environment.

➡Transistor has non zero ON resistance. Hence when it is ON, voltage across transistor is never zero. Moreover during OFF state also, there is flow of small leakage current. Hence it does not work as efficiently as mechanical switch or electrical switch or relay.

Integrated Circuit (IC)

An integrated circuit (IC) is a small semiconductor-based electronic device consisting of fabricated transistors, resistors and capacitors. Integrated circuits are the building blocks of most electronic devices and equipment. An integrated circuit is also known as a chip or microchip. A integrated circuit is built with the primary objective of embedding as many transistors as possible on a single semiconductor chip with numbers reaching in the billions as of 2012. According to their design assembly, integrated circuits have undergone several generations of

advancements and developments such as:

- Small Scale Integration (SSI): Ten to hundreds of transistors per chip
- Medium Scale Integration (MSI): Hundreds to thousands of transistors per chip
- Large Scale Integration (LSI): Thousands to several hundred thousand transistors per chip
- Very Large Scale Integration (VLSI): Up to 1 million transistors per chip
- Ultra Large Scale Integration (ULSI): This represents a modern IC with millions and billions of transistors per chip

An IC can be further classified as being digital, analog or a combination of both. The most common example of a modern IC is the computer processor, which consists of billions of fabricated transistors, logic gates and other digital circuitry.





COPYRIGHT FIMT 2020

745 | Page