

⇒ What is CO?

- CO is the consult with the way the hardware components operate and the way they are connected together to form comp. system.
- The various components are assumed to ^{be in} the place and task is to investigate the computer part operates with intended

⇒ Computer Architecture -

- It is consult with the structure and behaviour of the computer as seen by user.
- It includes the information formats, the instruction set and techniques for addressing the memory.
- The Architectural design of computer system is consult with specification of various functional modules such as processors, memories and structuring them together in comp. system.

→ Exa. of computer Architecture -

(1) Von- Neumann Architecture

(2) Harvard "

⇒ Computer Design -

→ Computer design is consult with hardware design of computer.

→ It covers the determination of what hardware should be used and how path is connected.

→ It is sometimes referred ^{to} as comp. Implementation

⇒ Von- Neumann Architecture -

→ It describes a general framework or structure that a computer's hardware, programming and data should follow.

→ Von-neumann envisioned the structure of computer system with following components.

(1) Central Arithmetic Unit.

(2) Memory

(3) control unit.

(4) Man-Machine Interface

- The most note-worthy concept of Von-Neuman Arch. was store programme principle.
- ⇒ store programme principle -
- It holds the data as well as instructions used to manipulate the data.
- The data and instructions should be stored together in the same memory area of the comp. and instructions are carried out sequential manner, One instruction at a time.
- ⇒ Harvard Architecture -
- It uses physically separate storage and signal pathways for their instructions and data. The term originated from harvard mark-I relay base comp.
- It stores instructions on punch tape (with 24-bits) and data on relay latches. (24 digits wide)

- In such architectures CPU can read both data and instruction from memory at same time, leading to increase in memory band-width.
- Ex. of hardend architecture - Micro-processor based comp.
- ⇒ Digital signal processing based comp. system -
- The data representation consists of four key elements.
 - (1) Data type
 - (2) Representation of data
 - (3) conversion ^{into} info.
- The term data refers to factual info. used for analysis and reasoning. It has no meaning in itself. On the other hand on information is a collection of data that needs to be communicated.
- Ex. Binary information is stored in memory on processor registers.

⇒ Registers -

- It contains data or control information. Control info. is a bit or a group of bits used to specify the sequence of command or signals, needed for manipulation of data in other registers.

⇒ Data -

- It can be defined as the no. and binary coded information, that are operated onto achieve required computation results.

⇒ Assignment -

- Diff b/w von-neumann and Harvard architecture.

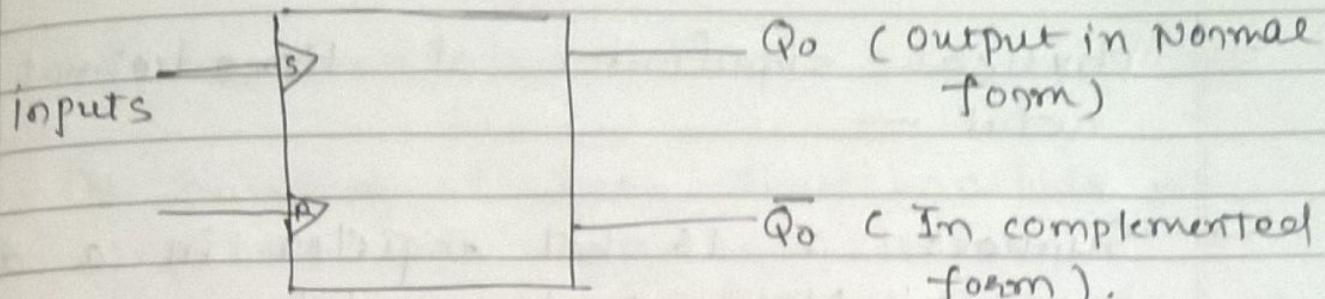
⇒ Data types -

- The data types found in the registers of digital computers mainly classified as
 - (1) Numbers used in numeric computation
 - (2) Letters of the alphabet used in data processing
 - (3) Other derivative symbols used for specific purposes
- All types data accepts binary nos. are represented in computer registers in binary coded form. This is because registers are made up of flip-flops (flip-flops are two state devices that can store only 0's and 1's).
The storage element is employed in clock sequential circuits are called flip-flop. A flip-flop stores 1-bit of information.)

⇒ Clock Sequential Circuits -

- In such circuits, the output is determined by periodic clock pulses. A flip-flop has two outputs,

One for
⇒ Normal value and other for the complemented value of bits stored



⇒ Data representation format -

⇒ Number System -

→ A Number system of base or radix represented by $(N)_x$ is a system that uses distinct symbols for x digits. Numbers are represented by a string of digit symbols.

- (1) Decimal Number system - It employs the radix 10 system. The symbols are from 0 to 9. To determine the quantity a no. represent, it is necessary to determine all digit by an integer power of radix and then form of weighted digits.
- (2) Binary Number system - It uses radix 2 (0 & 1)
- (3) Octal - It has radix 8 (0-7)
- (4) Hexadecimal - It has radix 16 (0 to 9)
A to F

Physical significance of Hexadecimal and octal -

Suppose, a 16-bit register in a digital comp. Now, each register may be composed of 16 binary storage cells. With each cell capable of holding a 0 and 1.

Thus, A 16-bit register can hold or store any binary no. from 0 to $2^{16}-1$

Binary coded Octal : e.g. $(10)_8 \Rightarrow 001\ 000$.	Binary coded Hexadic : e.g. $(14)_{16} \Rightarrow 0001\ 010$
--	--

Binary code -

It is a group of n bits that assume upto 2^n of distinct combinations 0's and 1's. with each combination representing 1 element of the set that is being coded.

$(\text{Number})^{\text{bits}} = \text{No. of elements}$

⇒ Alpha Numeric -

- for handing of data that consists not only no. but also of letters of alphabet and certain special characters
- It consists of a set of elements. That includes 10 decimal 26 alphabet letters and different special characters.

⇒ Assignment -

(1) What are char and their ASCII values?

(2) What is the difference among binary, octal, deci, hexa, BCD, BCDH.

⇒ complements - complements are used in Digital comp. to simply find subtraction ^{operation} and some logical Manipulation.

→ There are 2 types of complements, for any radix (R)

(1) R's comp-

(2) (R-1)'s comp.

$\Rightarrow (2-1)$'s comp. - Given a no. N and radix
n There has n digits

$$(2-1)'s \ comp. = (2^n - 1) - N$$

$$\text{e.g. } 2=10 \Rightarrow (10^n - 1) - N.$$

e.g. $n=4 \Rightarrow (2-1)$'s comp.

$$(2-1)'s \ comp. = (10^4 - 1) - N$$

e.g. 9's comp. - 546700 OP

(1) Difference between Von-nuemann and Harvard Architecture -

Von Nuemann

Harvard.

- The data and program are stored in the same memory. → The data and program memories are Separate.
- The code is executed serially and takes more clock cycles. → The code is executed in parallel.
- The programs can be optimized in lesser size. → The program tend to grow big in size.
- Used in conventional processors found in PCs and Servers. → Used in embedded systems and Mobile communication Systems.

(2) ASCII values of characters -

000	NUL
001	SOH
002	STX
003	ETX
004	EOT
005	ENQ
006	ACK

007	BEL
008	BS
009	HT
010	LF
011	VT
012	FF
013	CR

014	SO	045
015	SI	046
016	DLE	047
017	DC1	048
018	DC2	049
019	DC3	050
020	DC4	051
021	NAK	052
022	SYN	053
023	ETB	054
024	CAN	055
025	EM	056
026	SUB	057
027	ESC	058
028	FS	059
029	GS	060
030	RS	061
031	US	062
032	blank	063
033	!	064
034	"	065
035	#	066
036	\$	067
037	%	068
038	&	069
039	r	070
040	(071
041)	072
042	*	073
043	+	074
044	,	075

046	L	101	e
047	M	102	f
048	N	103	g
049	O	104	h
050	P	105	i
051	Q	106	j
052	R	107	k
053	S	108	l
054	T	109	m
055	U	110	n
056	V	111	o
057	W	112	p
058	X	113	q
059	Y	114	r
060	Z	115	s
061	[116	t
062	\	117	u
063]	118	v
064	↑	119	w
065	-	120	x
066	←	121	y
067	a	122	z
068	b	123	{
069	c	124	l
100	d	125	}
		126	~
		127	DEL

\Rightarrow 1's comp. —

\rightarrow $n=2$ then $(2-1)$'s comp. = $(2^n - 1) - N$

e.g. $n=4$

$$\downarrow \\ 2^4 = (10000)_2$$

$$2^4 - 1 = (1111)_2$$

\Rightarrow Note:- $(2-1)$'s comp. of octal and hexadecimal
are obtain by subtracting each
digit from 7 and F respectively

\Rightarrow 2's complement —

\rightarrow 1's comp. of $N = 2^n - N, N \neq 0$
and 2's comp. of $N = 0, N = 0$.

\Rightarrow 10's complement —

- (i) 2 (ii) 3 (iii) 8 (iv) 9.

10
01 (1's comp.)

10 (2's comp.)

11
00
1
01 (2's co.)

99
00
—
9
2
7

\Rightarrow 10's comp. -

2389

9999

$\begin{array}{r} 2389 \\ + 9999 \\ \hline 7611 \end{array}$

$\begin{array}{r} 1 \\ + 7611 \\ \hline 7611 \end{array}$

7611 (10's comp.)

$$= 2^n - N$$

$$= 10^4 - 2389$$

$$= 10000 - 2389$$

$$\begin{array}{r} 2389 \\ - 10000 \\ \hline 07611 \end{array}$$

$$= 7611$$

\Rightarrow 2's comp. - (101100)

010011 (1's comp.)

$\begin{array}{r} 1 \\ + 010011 \\ \hline 010100 \end{array}$ (2's comp.).

$$= 2^n - N$$

$$= 2^6 - 101100$$

$$= 1000000 - 101100$$

$$= 010100.$$

\rightarrow for 10^n it will be 1 followed by n 0's.

\rightarrow for $10^n - N$ it will be 10's comp. of (N) .

$\rightarrow 2^n - N$ is 2's comp. of N .

\Rightarrow 9's comp. (246700) = 353299.

10's comp. = 353399.

\Rightarrow find the 1's comp. (1101100).

~~= 0010011~~

~~0010100~~

~~= 0010011~~

→ Note: The comp. of the comp. of a number is the number itself.

⇒ R's comp. of R's comp. of N .

$$\text{R's comp. of } \text{N} = (\cancel{2^n - 1}) N \cdot 2^n - N$$

$$\begin{aligned}\text{R's comp. of R's comp.} &= R^n - (R^n - N) \\ &= N.\end{aligned}$$

⇒ Binary Point — It is needed to represent fractions, integers and Mixed integer fractions.

→ In registers, the complication are by the fact that binary points are characterised by position in registers.

→ There are two ways of specifying the position of ^{binary point in} register

- (1) fixed point representation —
- (2) floating point "

(1) It assumes that binary point is always fixed point. Most widely used positions are (i) Binary point in the extreme left ^{right} of the register to store no. as integer.

(ii) Binary point in the extreme left at the register to store no. as ~~float~~ ^{decimal} _{fractional}

(1) fixed Integer representation -

(A) signed numbers -

→ In ~~the~~ signed numbers,

(i) Positive numbers - Sign = 0

magnitude = Binary numbers

→

(ii) Negative numbers - Sign = 1

and rest is represented in
One of three possible ways

(1) Signed magnitude

(2) Signed 1's comp.

(3) Signed 2's comp.

$$\rightarrow +14 = 0000\ 1110.$$

$$\begin{array}{r} -14 = \\ \text{SM} \rightarrow 10001110, \\ \text{s 1'sc} \rightarrow 11110001 \\ \text{s 2'sc} \rightarrow 11110010 \end{array}$$

⇒ Signed Magnitude - It is used in ordinary
Architecture but not employed in computer
Arithmetic

⇒ Difficulty with 1's comp. -

→ In case of 1's comp., there are 2 0's
(+0, -0)

(1) Arithmetic Addition -

⇒ Subtraction - 2 n digits Unsigned nos.

$$\cancel{N \neq 0} \quad M - N, \quad \cancel{N \neq 0}$$

(i) Add Minuend M. to 8 comp.

$$\cancel{\text{of}} \quad \text{Subtrahend} = M + (8^n - N)$$

(ii) If $M \geq N$, sum will produce,
and end carry $^{(8^n)}$ which is
discarded and left out $M - N$.

(iii) If $M < N$, sum does not produce,
end carry, and carry
 $8^n - (N - M)$.

→ exa. $72532 - 13250$

10's comp. of subtrahend.

$ \begin{array}{r} 99999 \\ 13250 \\ \hline 86749 \\ \hline 86750 \end{array} $	$ \begin{array}{r} 72532 \\ 86750 \\ \hline 15928 \\ \hline 100000 \\ \hline 059282 \end{array} $
--	--

Ans. = 59282

$$\begin{array}{r} 13250 \\ - 72532 \\ \hline 59252 \end{array}$$

classmate

Date _____

Page _____

→ exa : $13250 - 72532 \quad (\text{MCN})$

$$= 2^n - (N - M)$$

$$= 10^5 - (72532 - 13250)$$

$$= 100000 - \cancel{02532} \quad 59252$$

=

→ exa. $x = (1010100)_2$

$$y = (1000011)_2$$

$$x - y$$

$$= 1010100 - 1000011$$

8's comp. of $y = 2^n - N$

$$= 2^7 - 1000011$$

$$= 10000000 - 1000011$$

$$= \begin{smallmatrix} 0 & 1 & 0 & 0 & 0 & 0 & 0 & 0 \\ + & 1 & 0 & 0 & 0 & 0 & 0 & 0 \\ \hline 1 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \end{smallmatrix}$$

$$- 1000011$$

$$\underline{\underline{00101101}}$$

⇒ Addition Algorithm -

→ signs of A & B are identical. \Rightarrow
Add two magnitude and attach
sign of A to the result

(1) Signs of A and B are different

(A) compare the magnitudes and
subtract the smaller no. from
larger no. choose the sign of
result to be same as A.
if $A > B$ or complement of
Sign of A if $A < B$.

(B) If Magnitudes are equal

Subtract B from A and make
the sign of result positive.

→ flip-flops are used to store "Sign bits".

⇒ Hardware components -

⇒ Parallel Adder - $A + B$ (Micro-operation)

⇒ Comparator - $A \geq B$, $A \leq B$, $A = B$.

⇒ Parallel Subtractor - $(A - B)$, $(B - A)$

⇒ Sign flip-flops (store sign bits).

Subtraction is achieved by additions and complement.

⇒ E-Register - It is used to check to determine the relative magnitudes of two numbers.

⇒ AVF (Add overflow flip-flop).

→ IF $M=0$ (Mode control) then output = B and Input = 0 and Adder output = A + B

→ IF $M=1$ then output will be 1's comp. of B and input = 1 and Adder output = $A + \bar{B} + 1$.

$$\Rightarrow X = 1010100$$

$$Y = 1000011$$

$$X+Y = \begin{array}{r} 1010100 \\ + 1000011 \\ \hline \textcircled{1} \quad 0010111 \end{array}$$

$$X-Y = 1010100 \\ + (2^n - Y)$$

$$= 1010100 + (2^7 - 1000011) \\ = 1010100 + (10000000 \\ - 1000011)$$

$$= 1010100 + 0111101 \\ = 10101001$$

⇒ 2's comp Arithmetic -

- (1) Addition: Add two nos including their sign bits and discard any carry out of the sign bit position

e.g. +6 00000110
 +13 00001101
 - 00010011

→ If no carry is generated, no is negative.

+6 00000110 13 = 00001101
 -13 11110011 -
 -7 11111001 1110010
 ↗ ↘
 2's comp. Sign bit representation
 00000110 by magnitude.
 - 00000111

→ -6 11111010
 -13 11110011
 -19 11101101

Note: Negative numbers must initially be in 2's comp and that if sum obtain after addition is negative, It is in 2's comp form.

(2) Subtraction: Take the 2's comp. of subtrahend (including the sign bit) and add to the minuend (including the sign bit). A carry out of the sign bit position is discarded.

(3) Binary nos in sign 1's comp. system are added and subtracted using same basic addition and subtraction rules as unsigned nos. Therefore computer need only one common hardware circuits to handle both types of arithmetic operation.

e.g.
$$\begin{array}{r} (-6) \\ - (-13) \\ \hline +7 \end{array} \quad \begin{array}{r} 11111010 \\ + 00001101 \\ \hline 00000111 \end{array}$$

⇒ Overflow -

When two nos of n digits are added and sum occupies of $n+1$ digits is called as "Overflow."

→ An overflow can not occur after an addition if one no. is positive and other is negative.

→ An overflow may occur if two nos added are both positive or both negative.

carry out.

carry in

e.g.
$$\begin{array}{r} +70 \\ +80 \\ \hline +150 \end{array}$$

$$\begin{array}{r} 01000110 \\ 01010000 \\ \hline 10010110 \end{array}$$

$$\begin{array}{r} -70 \\ -80 \\ \hline -150 \end{array}$$

$$\begin{array}{r} 10111010 \\ 10110000 \\ \hline 01101010 \end{array}$$

\Rightarrow Overflow Detection —

\rightarrow ~~Observing~~ Observing the carry into sign bit position and carry out of the sign bit position. If these two carries are not equal then there is a Overflow

$$C_{in} \oplus C_{out} = \begin{cases} 1, & \text{Overflow} \\ 0, & \text{no Overflow.} \end{cases}$$

\Rightarrow Algo for signed 2's comp

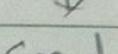
Subtract



Minuend in Ac
Subtrahend in Bn

$$Ac = Ac + Bn^1$$

$V \leftarrow$ overflow



End.

Addition



Minuend in Ac
sub. in Bn

$$Ac = Ac + Bn$$

$V \leftarrow$ overflow



End.

→ floating point representation -

(1) Mantissa - It is a signed fixed point number it may be a fraction or an integer.

(2) Exponent - It designate the position of decimal or binary point

e.g. 1001.11 convert into fraction and exponent.

$$\underline{0.100111} \times \underline{2^4}$$

sign bit

0 1001110

Mantissa

000100

exponent.

→ convert in Fixed point represent.

$$(-0.100111)_2 \times 2^4$$

Note: A floating point no. is said to be normalised if MSB of mantissa is non-zero

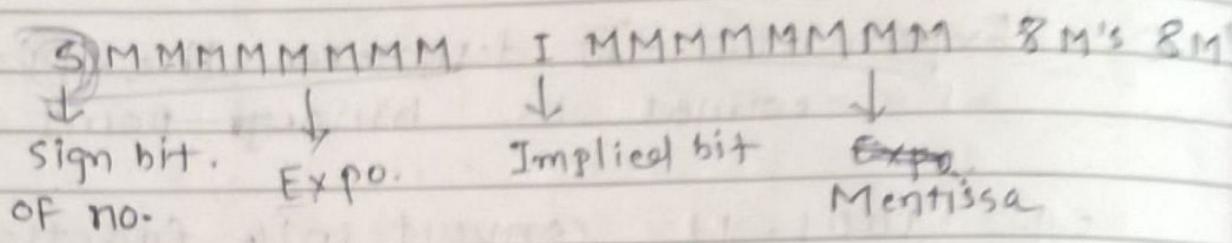
e.g. $(0001101000)_2$ ← Normalize it.

→ zero 1101000 can't be normalized.

$$(-.1101000)_2 \times 2^{-3}$$

⇒ Two main standard forms of floating point numbers are ANSI and IEEE

e.g. +13
Ansii Floating point representation -



e.g. $(1101)_2 = (0.1101) * 2^4$

↓ Ansii

0 0000100 11010000 00000000
0000000

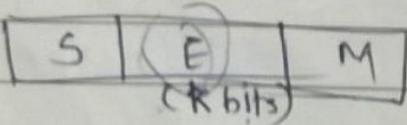
e.g. $(-17)_{10} = (-1.0001)_2$

$= (-0.10001)_2 * 2^5$

↓ ANSI

10000101 10001000 00000000
0000000

Note: Mantissa is represented in normalised sign magnitude fraction. The exponent is in biased form. If bias field has k bits then $\text{BIAS} = 2^{k-1}$.



The value of a no. is given by,

$$v = (-1)^s (1.M)_B * B^{E-BIAS} \quad (B = \text{Base})$$

Implicit normalized form increases the accuracy and value of no.

$$v = (-1)^s (1.M)_B * B^{E-BIAS}$$

⇒ IEEE - it provides the standard form
floating point no. is stored in
either in single precision and or
double precision (32 bits)
(64 bits)

→ The floating point no. can be denoted
with (a) Implicit normalization
(b) Fractional form.

32-bit Base: 2.

S(1)	E(8)	M(23)	Value
0/1	00...0	000...0	± 0
0/1	111...1	00...0	$\pm \infty$

By Default

→ If $E \neq 0, E \neq 255 \Rightarrow v = (-1)^s (1.M)_2 * 2^{E-127}$

→ If $E=0, M \neq 0 \Rightarrow v = (-1)^s (1.M)_2 * 2^{-126}$

By default

→ Other combination \Rightarrow NAN (Not a Number)

⇒ Floating point Arithmetic —

$$x = x_3 \text{ (. } x_M \text{)}_2 \times 2^{E_x - BIAS}$$

$$y = y_3 \text{ (. } y_M \text{)}_2 \times 2^{E_y - BIAS}$$

(i) $z = x + y = z_3 \text{ (. } z_M \text{)}_2 \times 2^{E_z - BIAS}$

where $z_3 = x_3 \oplus y_3$.

$$z_M = x_M + y_M$$

$$E_z = E_x + E_y$$

(ii) $z = x \times y = z_3 \text{ (. } z_M \text{)}_2 \times 2^{E_z - BIAS}$

where, $z_3 = x_3 \oplus y_3$.

$$z_M = x_M / y_M$$

$$E_z = E_x - E_y$$

(iii) Addition and Subtraction

→

(1) Check for zeros.

(2) Align the Mantissa.

(3) Add or subtract the Mantissa.

(4) Normalize the results.

Multiplication —

(1) Check for zeros.

(2) Add the exponents.

(3) Multiply the Mantissas.

(4) Normalize the products.

Division

- (1) Check for zeros.
- (2) Initialize Registers and
- (3) Align the dividend.
- (4) Subtract the exponents.
- (5) Divide the Mantissas.

= Booth Algorithm - It gives a procedure for multiplying binary integers in signed 2's comp. representation.

→ If there is a string of 0's, there is no addition but only shifting operations. If there is a string of 1s, from weight 2^k 2^m can be treated as $2^{k+1} - 2^m$

$$\text{e.g. } (+14)_{10} = \underline{0} \ 1110.$$

↑ 3 2 1 0
 $j=3 \leftarrow \rightarrow k=1$

k will be 1st "1" occur in binary from LSB to MSB

$$2^{j+1} - 2^k$$

$$= 2^4 - 2^1 = 14$$

Rule (1) Multiplicand is subtracted from Partial product upon encountering the first least significant 1 in string of 1s in the Multiplier

(2) The Multiplicand is added to the partial product upon encountering the first zero (provided there was a previous one) in a string of os in the Multiplier.

(3) The partial product doesn't change when the Multiplier bit is identical to previous Multiplier bit.

e.g. $(11) \times (110)$ $S_C = \text{No. of bits of BR}$ $BR = \text{Multiplicand}$
 $QR = \text{Multiplier.}$

Q_{n+1} BR = 10111 AC QR Q_{n+1} S_C
 $\overline{BR} + 1 = 01001$, on
Initial 0 10011 0 101

$$\begin{array}{r}
(1) \quad (0) \\
\text{Sub } b_2 \\
\text{Add } \overline{BR} + 1 \\
\hline
01001
\end{array}
\qquad
\begin{array}{r}
00000 \\
+ 01001 \\
\hline
11101
\end{array}
\qquad
\begin{array}{r}
10011 \\
\hline
10011
\end{array}$$

ashr (shifting) 00100 1100 1 100

$$\begin{array}{r}
1 \quad 1 \quad \text{ashr} \\
\qquad \qquad \qquad 00010 \quad 01100 \quad 1 \quad 011
\end{array}$$

$$\begin{array}{r}
0 \quad 1 \quad \text{Add } B_1 \\
\text{Sub } b_1 \\
\hline
11001
\end{array}$$

ashr 11100 10110 0 010

$$\begin{array}{r}
0 \quad 0 \quad \text{ashr} \\
\qquad \qquad \qquad 11110 \quad 01011 \quad 0 \quad 001
\end{array}$$

$$\begin{array}{r}
0 \quad 0 \quad \text{Sub } b_0 \\
\text{Sub } b_0 \\
\hline
00111
\end{array}$$

ashu 00011 10101 1 000

$$\text{so, } \text{MULT} = \text{AC} \cdot \text{BR}$$

$$= 00011 10101$$

1 0 1 1 1 1 1 0

→ e.g. 0 1 1 1 1 1 0

$$\text{conventional: } 2^1 + 2^2 + 2^3 + 2^4 + 2^5 + 2^6$$

6 Add operation, 21 shift

$$\text{Booth: } +2^{j+1} - 2^k$$

$$= +2^7 - 2^1$$

$$= +A - B = A + (-B)$$

1 Add, 1 sub, 8 shift

→ Sign Extension -

- Placing a smaller value in larger register is called sign extension.
- The sign extension doesn't change the value and sign of the number.
- In 1s, 2s comp. notation the sign bit is copied in all the extended places.
- In sign magnitude notation, the sign bit is moved to most MSB of extended register and remaining bits are filled with zero.

e.g. $(-5) = 10101$ (5-bit register)

1	0	1	0	1
---	---	---	---	---

$= 11101$ (copied Answer)

QUESTION:- What is sign extension?

ANSWER:-

$$(3)_10 \Rightarrow (1001)_2, 10001 (-) \Rightarrow (10101)_2$$

Memory - Organization.

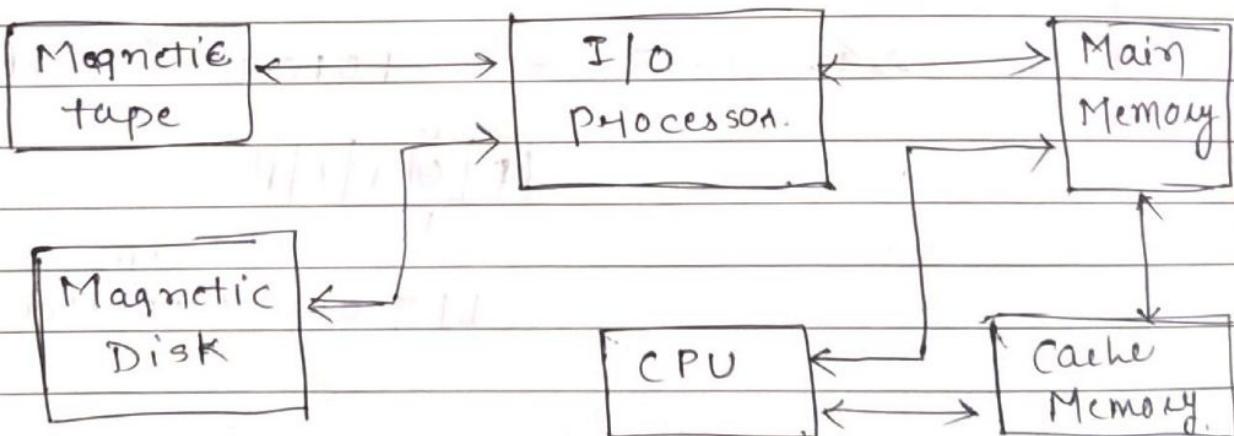
⇒ Memory hierarchy —

⇒ (1) Main Memory — The memory unit that communicates directly with CPU.

(2) Auxiliary Memory — Devices or memory which provides backup storage are called auxiliary memory.

→ They are used for storing system programmes, large data files and other backup information.

Auxiliary Memory (M/W)



→ Only programs and data currently needed by processor resides in main memory.

⇒ Cache Memory — It compensates for speed diff. between main memory and auxiliary memory. It makes

current programs and data available to CPU at a rapid rate.

Goal of Memory hierarchy - To increase the performance rate of computer and obtain highest possible average access speed with minimization of the cost of entire system.

Explanation about diagram -

- (1) The cache holds those parts of programs and data that heavily used while auxiliary memory holds those parts of programs and data not presently used by CPU.
- (2) CPU has direct access to both cache and main memory but not to auxiliary memory.
- (3) The transfer from auxiliary to main memory is usually done by means of direct memory access of large blocks of data.
- (4) The typical access time ratio b/w Cache and main memory is about 1 to 7.
- (5) Auxiliary memory access is usually thousands times that of main memory.

- The Memory hierarchy is based on four parameters (1) T (2) S (3) C (4) F.

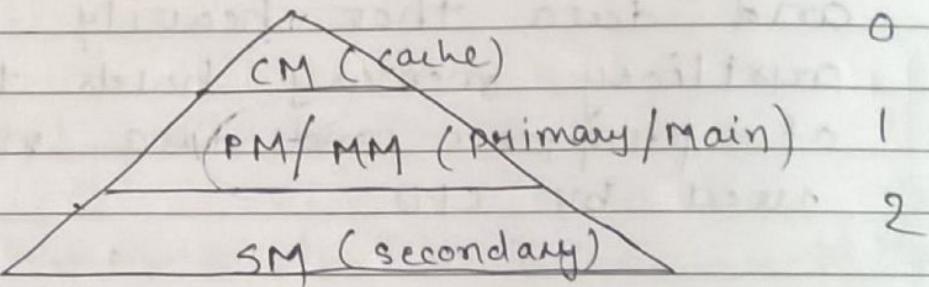
T - time

S - size

C - cost

F - frequency

- In the structure of Memory hierarchy, the i^{th} level is physically positioned ~~higher than~~ $(i+1)^{th}$ level memory



$$T_i < T_{i+1}$$

$$S_i < S_{i+1}$$

$$c_i > c_{i+1}$$

$$f_i > f_{i+1}$$

$$I_i \subset I_{i+1} \text{ (information)}$$

- The performance of Memory hierarchy is given by hit ratio.

- It is availability of required info at required level.

The hit ratio of last level is always one.

- The hit ratio is inversely proportional to average access time and directly proportional to size.
- Larger main memory and larger cache size is an ideal choice.
- The side effects of the Memory hierarchy is data inconsistency (same information is available differently at diff. levels)
- The proper updation technique resolves the problem of data inconsistency.
- ⇒ Mathematical expression of Memory Hierarchy.

Let T_1, S_1, C_1, H_1 are the specifications of level 1 memory in a two level system. Similarly the specification of Level 2 memory are T_2, S_2, C_2, H_2

- Total size of Memory $S_T = S_1 + S_2$.
- Average cost $C_{avg} = \frac{C_1 S_1 + C_2 S_2}{S_1 + S_2}$.
- Average access time,

$$T_{avg} = \frac{H_1 T_1 + H_2 T_2}{H_1 + H_2} = H_1 T_1 + (1 - H_1) T_2$$

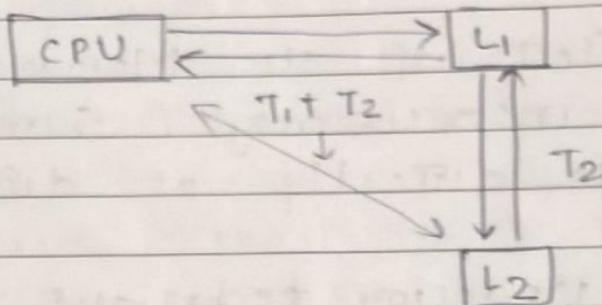
when $T_1 \ll T_2$.

$$T_{avg} = H_2 T_2 + (1 - H_2) T_1$$

when $T_1 \gg T_2$.

→ If the process is present in L_1 , then $H_1 = 1$. If the process is present in L_2 Then $H_2 = 1$ and $H_1 = 0$.

⇒ Strict hierarchy -



$$\text{Avg for CPU} = H_1 \cdot T_1 + (1-H_1) (T_1 + T_2).$$

→ Inspite of moving a single word, He can move a clock of words from L_2 to L_1 . So that next time when any thing asked, it will be present in L_1 itself where the time to move from L_2 to L_1 is given by $T_B = N * T_2$. (N = size of word-blocks).

so, Now the avg access time will be,

$$\checkmark = H_1 T_1 + (1-H_1) (T_B + T_1).$$

$$= H_1 T_1 + (1-H_1) (\frac{T_B}{N} + T_1)$$

$$= H_1 T_1 + (1-H_1) T_B + (1-H_1) T_1$$

$$= T_1 (H_1 + 1 - H_1) + (1-H_1) T_B$$

$$= T_1 + (1-H_1) T_B.$$

when the word is not available in L_1 , then T_B is the extra time to move from L_2 to L_1 and $N = \text{Size of blocks}$.

For Multiple levels, Average access time.

$$T_{avg} = T_1 + (1-H_1)T_{B_1} + (1-H_1)(1-H_2)T_{B_2} + \dots$$

If the word is available in L_1 , then T_1 is access time. Not in L_1 , then $H_1=0$ and $H_2=1$.

$$(1) T_{avg} = T_1, H_1=1$$

$$(2) T_{avg} = T_1 + T_{B_1}, H_1=0$$

$$(3) T_{avg} = T_1 + T_{B_1} + T_{B_2}, H_1=1, H_2=0.$$

Ques- consider a two level memory system. The cache access time is 100ns and it is 12 times faster than Main memory what will be the minimum hit ratio that ensures average access time should not exceed 120ns.

for two level memory system,

$$T_{avg} = H_1 T_1 + (1-H_1) T_2$$

here, $T_1 = T_c$ and $T_2 = T_m$, $T_m = 12 \times T_c$.
 $= 100 \text{ ns}$ $= 1200 \text{ ns}$.

$$\text{so, } T_{avg} = H_1 (100) + (1-H_1) (1200) \quad \text{--- (1)}$$

$$120 = H$$

$$H \propto \frac{1}{T_{avg}}$$

Average access time is max $\Rightarrow H$ will be min.

$$\therefore C_{120} = H_{\max} \cdot (100) + (1 - H_{\min}) \times 1200$$

$$\therefore H_{\min} = 120 = 100H + 1200 - 1200H$$

$$\therefore 120 = 1200 - 1100H$$

$$\therefore 1100H = 1200 - 120$$

$$\therefore H = 98\%$$

(b) If the heat ratio is 100% what will be avg. access time of cache.

$$H \propto \frac{1}{T_{\text{avg}}}$$

→ Heat ratio does not effect the individual access time.

so, T_c won't change. It will remain same = 100 ns

Ques-2 Consider a three level memory system with heat ratios 0.8, 0.9 and 1. Access times are 10ns, 20ns and 100ns. If the required word is not available in level 1 or ~~level 2~~ word block is transferred from level 2 to level 1 and handed over to the processor. If it is not present in level 2 then associated 8-word block is first transferred.

from level-3 to level-2 and ~~consent block~~^{cored} is transferred from level-2 to level 1. What will be the avg access time.

$$T_{avg} = T_1 + (1-H_1)T_{B1} + (1-H_1)(1-H_2)T_{B2}$$

here $T_{B1} = n * T_2$ (From Level 2 to Level 1)

$$\therefore T_{B1} = 2 T_2 = 40 \text{ ns.}$$

$$\begin{aligned} \text{and } T_{B2} &= 8 T_3 \\ &= 8 \times 1000 \\ &= 8000 \text{ ns} \end{aligned}$$

$$T_1 = 10 \text{ ns}, \quad H_1 = 0.8, \quad H_2 = 0.9.$$

$$\begin{aligned} T_{avg} &= 10 + (1-0.8)(40) + (1-0.8)(1-0.9)(8000) \\ &= 18 + 160 \end{aligned}$$

$$\therefore T_{avg} = 178 \text{ ns}$$

Ques consider a cache which is 5 times faster than memory with hit ratio 80%. The avg. access time of cache is increased by 20% from 80ns.

(a) What is the cache access time.

(b) what is the main memory access time.

(c) What is the new hit ratio.

$$T_m = 5 \times T_c.$$

$$H_1 = 0.8.$$

$$(q) T_{avg} = H_1 T_c + (1-H_1) T_m$$

$$\therefore 50 = 0.8 (T_c) + (1-0.8) (5 T_c)$$

$$\therefore 50 = 0.8 T_c + T_c$$

$$\therefore T_c = 27.77 \text{ ns}, (b) T_m = 5 \times 27.77 = 138.85 \text{ ns}$$

→ Avg. access time is increased by 20%.

$$(T_{avg})_{new} = \frac{50 \times 20}{100} = \frac{1000}{100} = 10 \text{ ns} + 50$$

$$\therefore (T_{avg})_{new} = 60 \text{ ns}$$

~~$$(e) (H)_{new} = 0.64 \quad (\because H \propto \frac{1}{T_{avg}})$$~~

~~$$T_{avg} = H T_c + (1-H) T_m$$~~

~~$$\therefore 60 = (0.64) T_c + (1-0.64) (5 T_c)$$~~

~~$$\therefore 60 = 0.64 T_c + 1.8 T_c$$~~

~~$$\therefore (H)_{new} =$$~~

~~$$\therefore 60 = H_{new} (27.77) + (1-H_{new})(138.85)$$~~

~~$$\therefore 60 = 27.77 H_{new} + 138.85 - 138.85 H_{new}$$~~

~~$$\therefore 111.08 H_{new} = 78.85$$~~

(c) $\therefore H_{new} = 0.7098$

Cache Memory -

- It is the smallest and fastest memory component in hierarchy and maintains the locality of reference.
- The cache and main memory is divided into fixed sized partitions (blocks). The address mapping decides which block of the main memory is to be placed in which position.
- There are three types of address mapping.
 - (1) Direct mapping.
 - (2) Associative mapping.
 - (3) Set-associative mapping.
- Address mapping converts physical address to cache address.
- Locality of reference is of two types.
 - (1) ~~special~~ Spatial (space position)
 - (2) Temporal (In time).

(1) Direct Mapping -

- It is the simplest one and requires less no. of Tag bits.
- It is the fastest one but it requires more hardware.

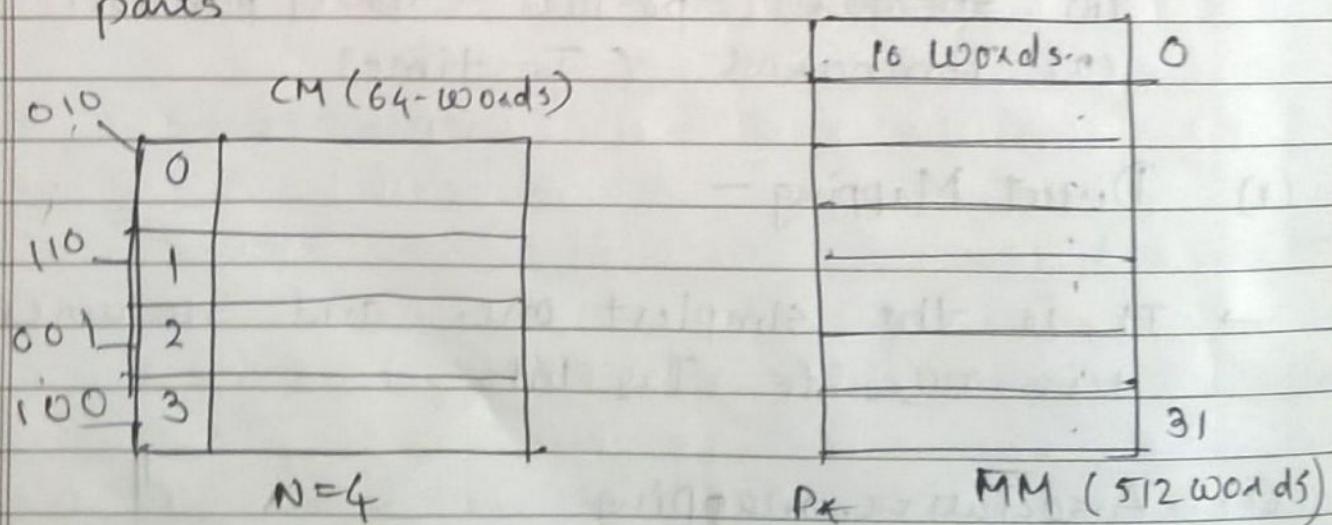
(3) Set-associative mapping -

- It gives the advantage of direct and associative mapping.

(1) Direct-mapping -

- k^{th} block of MM (main memory) has to be placed in $(k \bmod N)^{\text{th}}$ cache position. Many main memory block complete for same cache position. Only one among them is to be placed at a time in cache.

- Tag bits denote which main memory block is currently placed in cache. Therefore, every cache block is attached with tag field.
- Physical Address is divided into three parts



In the above exa. word offset = 4 bits
(position of word within the block)

Cache offset = 2 bits. (since 4 blocks are in cache)

Tag = 3 bits.

Total no. of bits = 9 bits (since 512 words in main memory). $= 2^N = 2^9 = 512$

⇒ MM Blocks -

[0, 4, 8, 12, 16, 20, 24, 28] → 0

[1, 5, 9, 13, 17, 21, 25, 29] → 1

[2, 6, 10, 14, 18, 22, 26, 30] → 2

[3, 7, 11, 15, 19, 23, 27, 31] → 3

} Cache

→ Tag bit will tell that which block of MM is placed in which block of CM.
It is of 3 bits.

Tag = {000, 001, 010, 011, 100, 101, 110, 111}.

ie check whether 200 words of MM is available or in cache or not.

[011 001000]

↑ ↑ ↑
Tag. Cache word offset.

200 words of MM is not available in cache because Tag bit is 011.

ie check whether 121 words of MM is available in cache or not.

[001 11 1001]

↑ ↑ ↑
Tag Cache word
offset.

121 words of MM is not available because 001 Tag is available but 11 (\Rightarrow) is not related to 001 2 is related to 001.

ie 92 words.

[001 01 1100]

"No".

Example 96 words

[001 10 0000],

"Yes".

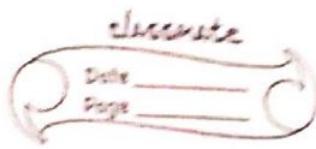
⇒ Limitations of Direct-Mapping -

- Suppose, a cache has four blocks with CPU MIF as given : 0 4 0 8 4 0 0 8
- For placing first block in cache it is placed in 0th place of cache. Now, Again references of blocks are placed in 0th block of cache by removing the previous one.

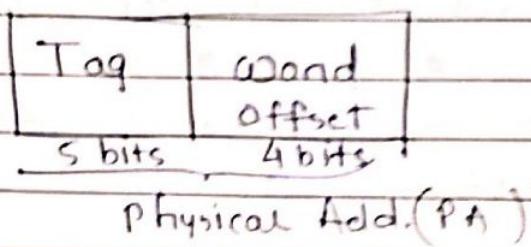
- (1) The No. of page replacements are more.
- (2) The time is slower.
- (3)

⇒ Associative Mapping —

- It is a solution for limitations of direct mapping.
- Any block of main memory can be placed anywhere in a cache.
- Since, It involves less block movement it is the fastest mapping technique.



- Associative mapping requires more hardware therefore costliest technique.
- The no. of Tag comparators is equal to No. of blocks in the cache.



- The no. of Tag bits are more in associative mapping

$$Tag_A = Tag_M + \text{Cache offset}$$

Actual

$$\frac{\text{Size of Cache}}{\text{Data M/W}} + \frac{\text{Tag M/W (Memory)}}{N * \text{Tag}}$$

- $N * \text{Tag} = \text{No. of blocks in } \frac{\text{no. of Tag bits}}{\text{cache}}$

$$R: 0\ 4\ 0\ 9\ 4\ 0\ 2\ 2$$

0	0
1	4
2	8
3	2

Limitations -

- It is a costliest technique.

Set-associative Mapping -

In this mapping, k^{th} block of main memory to be placed in $k \pmod s$ set, where $s = \text{Total no. of sets in cache}$.

$$s = \frac{N}{P}, P = \text{Implemented as } p\text{-way set-association}$$

Within & within the set, it can be placed anywhere in set.

The no. of Tag comparators is equal to the size of set P .

The p -way set association divides the phy-add. into 3 components.

$\log_2 s$ bits.

Tag	Set offset	Word offset

PA

If $P=1$, it is direct mapping.

If $P=N$, associative mapping.

(Mb - Mega bits)
Byte

Consider 1 MB of Cache memory
and 1 GB of Main memory. Both are divided into 32-word blocks.

Each word contains 32 bits.

- (1) What will be the no. of bits required to address main memory.
- (2) What will be the no. of tag bits for following mapping techniques.
Direct, associative, Set-associative.
(8-way)
- (3) How many comparators and of what size are required for each case.
- (4) What is the actual size of cache memory required for above mapping.

Consider a cache memory which is applied with direct mapping.

The no. of tag bits is equal to the no. of block in cache.

Each blocks contain N words

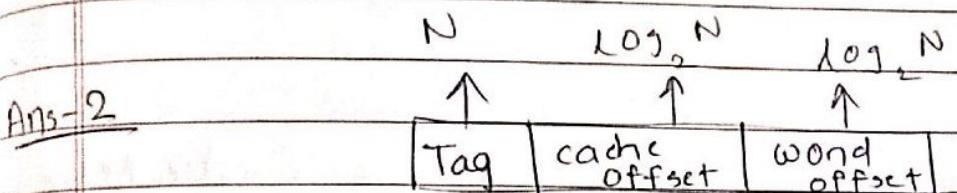
N = total cache blocks

How many words are there in main memory.

Ques The k^{th} block of main-memory has to be placed in which cache set if 2-way set association is used for the entire $2C$ cache blocks.

Ans - 3 $s = \frac{N}{P} = \frac{2C}{2} = c.$

so, k^{th} block is to be place in cache memory
 $= k \bmod c.$



Here Tag bits = The no. of blocks in cache,

\therefore Tag bits = N bits

If 4 blocks in blocks in cache \Rightarrow 2 bits

N blocks $\Rightarrow \log_2 N$.

$$PA = N + \log_2 N + \log_2 N$$

$$= N + \log_2 N^2.$$

$$\therefore 2^K = 2^{N + \log_2 N^2}$$

$$10^9 = 2^{30} \text{ bits}$$

Ans-1

$$512M = 1MB$$

$$5MM = 1GB$$

Block size = 32 words

No. of bits of in each word = 32 bits.

$$\rightarrow \text{No. of words} = 1000 \text{ KB}$$

$$= \frac{1000 \times 1000}{8} \text{ Bytes}$$

$$= \frac{10^8}{8} \text{ bytes}$$

$$= \frac{10^6}{4}$$

$$= \frac{2^{20}}{2^2}$$

= 2^{18} words in Cache Memory.

$$= 256 \text{ K words}$$

$$\rightarrow \text{No. of blocks} = \frac{2^{18}}{2^5}$$

$$= 2^{13} \text{ blocks}$$

$$= 2^{10}, 2^3$$

$$= 8 \text{ K blocks.}$$

\Rightarrow for Main Memory,

$$\text{No. of words} = \frac{10^9 \text{ bytes}}{2^2 \text{ bytes}}$$

$$= \frac{2^{30}}{2^2}$$

$$= 2^{28} \text{ bytes words}$$

$$= 2^{10} \cdot 2^{18} \text{ words}$$

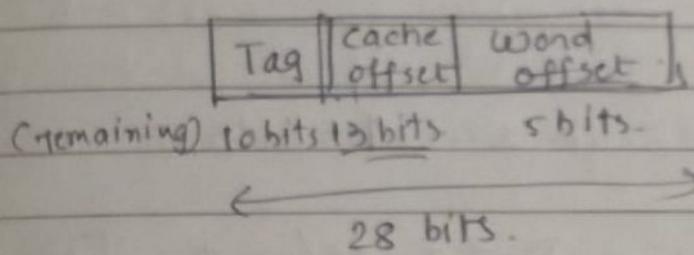
~~$$= 2^{25} \text{ words.} = 2^{18} \text{ K words.} = 256 \text{ M words.}$$~~

$$\begin{aligned}
 \text{No. of blocks} &= \frac{2^{28}}{2^5} \\
 &= 2^{23} \text{ blocks} \\
 &= 2^{10} \cdot 2^{13} \\
 &= 2^{13} \text{ K blocks} \\
 &= 8 \text{ M blocks.}
 \end{aligned}$$

(a) No. of words = 2^{28} words.

$$\log_2 2^{28} = 28 \text{ bits}$$

(b) (1) Direct Mapping :



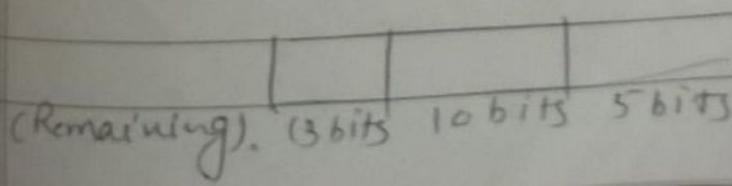
(2) Association : Tag bits = $10 + 13 = 23$ bits
 Word off = 5 bits

(3) Set-associative mapping :

$$\underline{P = 8}$$

$$s = N/P = \frac{2^{13}}{2^3} = 2^{10}.$$

Set offset = 10 bits



Note: Higher association. Higher will be the tag bits

(c) In Direct mapping:

1 Tag comparison having 10 bits. (because $\log_2 1024 = 10$)

In Associative mapping:

8 ~~K~~ Tag comp. = no. of blocks in cache.
~~812~~ has 23 bits

In SAM:

8 ~~R~~ Tag comp. have 23 bits

(d) Actual size = Data M/W + Tag M/W

= 1 MB + $(\frac{\text{No. of blocks} \times 10 \text{ bits}}{8})$ bytes

8 ~~R~~ - to convert bit in bytes

$$(i) D = 1 \text{ MB} + \frac{2^{13} \times 10}{8} \text{ bytes}$$

$$(ii) A = 1 \text{ MB} + \frac{2^{13} \times 23}{8} \text{ bytes}$$

$$(iii) S = 1 \text{ MB} + \frac{2^{13} \times 13}{8} \text{ bytes}$$

- Updation Techniques - They are used to deal with cache coherence problem
- The write-through Updation simultaneously update the Cache and Main
- The Time for Updation = Max (T_C, T_M)
- It gives better performance for less no. of Updation. In write-back Updation the main memory updation is done only when consecutive updated block chosen for replacement.
- If the block chosen for replacement is not modified, then Incoming block can be simply over-written in the case.
- The dirty bit is used to indicate the status of the block in Cache.
- The time for updation is given by,

$$T_{\text{Upd}} = T_B + T_C$$

↓ →
 New block Current Updation

This ~~block~~ formula is used to clean block.

for dirty block, $T_{\text{Upd}} = 2T_B + \alpha T_C$

↓
 T_B(old) + T_B(new)

- The current Updation is always done in Cache Memory.

Ques. Consider a hypothetical processor which issues 25% references for update. It uses 2 level memory hierarchy with
 $T_c = 10\text{ns}$

$$T_m = 100\text{ns}$$

$$T_B = 10\text{ns}$$

$$\text{HR} = 0.2 \text{ (Read heat)}$$

$$H_w = 0.9 \text{ (Write head)}$$

If the set word is not available then a four word block is to be moved from main memory to cache (either for read or write) what is the performance of this memory with write through strategy.

$$\text{Performance} = \frac{1}{T_{avg}} \text{ words/sec}$$

$$\rightarrow T_{avg} = 25\% \cdot T_{avg}(w) + 75\% \cdot T_{avg}(\text{Read}) \quad \text{--- (1)}$$

$$T_{avg}(w) = H_w \cdot \text{Topdation} + (1-H_w) (\text{Topdation} T_B)$$

$$= 0.9 \cdot \text{Max}(T_c, T_m) + (0.1) (\text{Max}(T_c, T_m) +$$

$$T_{avg}(w) = 14\text{ns}$$

$$T_{avg}(\text{R}) = H_R \cdot T_c + (1-H_R) (T_c + T_B)$$

$$= 9\text{ns}$$

$$BR = -13 = 11101$$

$$= 10011$$

$$\overline{BR+1} = 01101, QR = 0111$$

-13 ~~7~~

classmate
Put Multiplier in
"Magnitude form"
2 is Complement

Qn	Qn+1	BR	AC	QR	Qn+1	SC
----	------	----	----	----	------	----

putting values in eq. (1).

$$T_{avg} = 25\% (140) + 75\% (90)$$

$$= 102.5 \text{ ns}$$

$$\text{Performance} = \frac{1}{T_{avg}} = 9.75 \times 10^3 \text{ words/sec}$$

$$= 9.75 \times 10^6 \text{ words/sec.}$$

Ques Perform the same previous que with write back policy.
At any point of time 20% cache blocks are modified and cache is full

~~T_{avg}~~ Performance = $\frac{1}{T_{avg}}$ words/sec.

$$T_{avg(R)} = HR * T_c + (1-HR) \left[\underbrace{(T_B + T_c) * 80\%}_{\text{clean}} + 20\% * \underbrace{(T_B(\text{old}) + T_B(\text{new}) + T_c)}_{\text{Dirty}} \right]$$

$$= 0.8 \times 10 + (0.2) \left[\frac{410 \times 80}{100} + \frac{20}{100} (400 + 900 + 10) \right]$$

$$= 106 \text{ ns./words}$$

$$T_{avg(w)} = Hw * T_{update} + (1-Hw) \left[(T_B + T_c) * 80\% + 20\% * \frac{(T_B(\text{old}) + T_B(\text{new}) + T_c)}{T_c} \right]$$

$$= 58 \text{ ns./words}$$

$$T_{avg} = 25\% T_{avg(w)} + 75\% T_{avg(R)} = 94 \text{ ns./words}$$

$$\text{Performance} = 10.6 \times 10^6 \text{ words/sec.}$$

⇒ Block replacement techniques -

- These techniques are aimed to choose such a cache block for replacement that may result in less or no penalty.
 - There are 3 techniques for block replacement.
 - (1) FIFO - The block that spent longer time in cache is chosen for replacement on the basis of arrival time. This strategy is ~~not~~ implemented using queue.
 - (2) LRU (Least Recently Used) - The block that is not recently used is chosen for replacement. It can be implemented with a queue with slight adjustment on hit.
- Note: (1) If there is no hit, FIFO and LRU are same.
- (2) Hit is the Recently Referenced References.

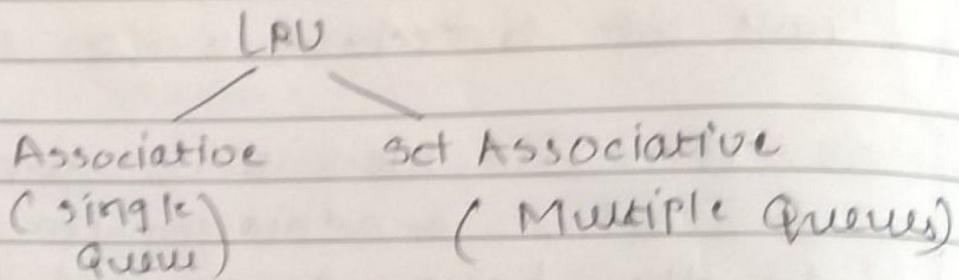
- (3) Direct Mapping - $(k \bmod n)^{\text{th}}$ block is chosen for replacement, where k = Recently referred Main memory block.

Q-1 No. of hits

Q-2 which blocks are present?

Q-3 " " " not present?

⇒ LRU is of two types.



Que Consider a four block cache. The following main memory blocks are referred. Which replacement strategy use best performance? (No. of miss)

Referred blocks: 4, 5, 7, 12, 4, 5, 13, 4, 5, 7, 12, 8

a) FIFO

					Hit
4*	5*	7*	12*	4	5
4	4	4	4	4	4
	5	5	5	5	5
		7	7	7	7
			12	12	12
13*	4*	5*	7*	12*	13*
13	13	13	13	12	12
5	4	4	4	4	13
7	7	5	5	5	5
12	12	12	7	7	7

* - Shows that first time it is coming

No. of miss = 10.

No. of Hit = 2.

Last Replacement = 4.

Blocks present =

12
13
5
7

(2) LRU -

4*	5*	7*	12*	4	5	13*	4*	5	7*
4	4	4	4	4	4	4	4	4	4
5	5	5	5	5	5	5	5	5	5
	7	7	7	7	7	13	13	13	13
		1	12	12	12	12	12	12	7

12未 13未

No. of Miss = 8

Hit = 4

4	13
5	5
12	12
7	7

(3) Direct Mapping - ($k \bmod 4$)

	4*	5*	7*	12*	4*	5	13*	4	5*	4	12*	13*
0	4	4	4	12	4	4	4	4	4	4	12	12
1		5	5	5	5	5	13	13	5	5	5	13
2												
3			7	7	7	7	7	7	7	7	7	7

No. of miss = 9.

Hit = 3.

Last Replacement = 5.

Performance: FIFO < DM < LRU.

$$S = \frac{P}{P} = \frac{16}{4} = 4$$

Ques Consider a 4 way set associative cache which was empty initially. Cache contains 16-blocks & main memory contains 175. Use LRU replacement policy with following References.

Ref: 0, 255, 1, 4, 3, 8, 133, 159, 216, 129, 63, 8, 48, 32, 73, 92, 155

formula $k \bmod s$, $s = \frac{N}{P} = \frac{16}{4 \text{ way}} = 4$
 $= k \bmod 4$

⇒ Memory hierarchy - The objective of Memory hierarchy is to prove the matching of the data transfer rate of faster processor to that of lowest level memory of reasonable cost.

⇒ Memory Unit: It is an essential unit in a digital computer since, it is needed for storing programs that are executed by CPU.

⇒ Main Memory: The Memory Unit that communicates directly with CPU is called Main memory. ex. RAM & ROM. It is also called "Primary Memory."

→ The principle tech. used for main memory is based on semiconductor JCS.

- The static RAM consists essentially of internal flip flops that stores binary info.
- The stored information remains valid as long as power is apply to the Unit.
- The dynamic RAM stores the binary info. in the form of electric charges that are apply to capacitor.
- The stored charge tend to and the capacitors must be periodically recharged by refreshing the dynamic memory.
- Refreshing is done by cycling through the words every few milliseconds to restore the decaying charge.
- The dynamic RAM offers reduced power consumption and larger storage capacity in a single chip.
- The static RAM is easier to use and has shorter read and write cycles.

- classmate
Date _____
Page _____
- Cache Memory - A special very high speed m/m called cache is used to increase the speed of processing by making current programs and data available to CPU at Rapid Rate.
 - The Cache Memory is employed in comp. system to compensate for the Speed diff. b/w main memory access time and processor logic.
 - The speed of Main memory is very low in comparison with speed of Modern processors.
 - For good performance, the processor can't spent much of its time waiting to access instructions and data in Main memory. An efficient solution is to used a fast cache m/m, which essentially makes the main memory appears to the processor to be faster than it really is.
 - The effect of the cache mechanism is based on a property of computer program called Locality of Program.
 - Many instru. in localised area of prog. are executed Repeatedly during some time period and the remainder of program is accessed relatively infrequently.

this is ref. as locality of program.

References.

→ It Manifests in two ways

- (1) Temporal (2) Spatial

(1) The temporal aspect of locality of ref suggests that whenever an information (Instr./data) ^{is first needed this item} is fetched into the cache where it will be hopefully remain until it is needed again.

(2) The special aspect suggests that instead of fetching one item from MM to the cache. It is useful to fetch several items that reside at adjacent addresses as well. We will use the term block to ref. to a set of contiguous address locations of some size. Another term that is often used to ref. to a cache block is Cache line.

→ The correspondence b/w the main memory blocks and those in cache is specified by mapping function. When the cache is full and memory word that is not in cache is ref. the cache control hardware must decide which block should be removed to create space for the new block that contains the ref. words. The collection of rule for

making this decision constitutes the Replacement algorithm.

- The processor doesn't need to know about the existence of MM. It simply issues Read and write locations in the memory. In a read operation MM is not involved. For write operation the system proceeds in two ways.
 - (1) Write-through Mechanism - The cache location and MM location are updated simultaneously.
 - The second tech. is to update only the cache location and is to mark it as updated with an associative flag bit often called dirty or Modified bit.
 - The MM location of the word is updated later when the block containing this marked word is to be removed from the cache to make room for a new block. This technique is known as write back / copy back.
- When the addressed word in a Read operation is not in a cache a Read miss occurs.
 - The block of words that contains the requested word is copied from the main memory into cache. After entire block is loaded into the cache

the particular word requested is forwarded to the processor.

- Alternatively the word is may be sent to the processor as soon as it is send from the MM. The later approach called load through/ early restart reduces the processors waiting period

⇒ Cache coherence Problem —

- (1) will have multiple copies of same data in diff caches.
- (2) The data in caches modified locally.
- (3) Cache become inconsistent b/w caches and b/w cache and MM. from which cache data was copied.
- The transformation of data from MM to CM is ref. as a mapping process

⇒ Address Mapping —

- When a tag add. is presented to the cache it must be quickly compare to the stored tags to determine whether the matching tag is currently assigned to the cache.
- The fastest tech. for implementing tag comparisons is associative or content addressing which permits the ip tag to be compared simultaneously to all tags in the cache tag RM.

(A) Associative Mapping: More flexible mapping Method in which a main Memory block can be replaced into any cache block position. The tag bits of an address received from the processor are compared to the tag bits of each block of the cache to see if the desired block is present or not.

- An associative cache employs a tag that is a block add. as the key. At the start of Memory access the incoming tag is compared simultaneously to all the tag stored in the cache tag memory.
- If a match occurs, a match indicating signal triggers the cache to service the requested M. access.
- A no match signal identifies a cache miss and the memory access requested to be forwarded to Main Memory for service. A cache block containing the target add is then sent to from M to C. and at the same time a data word is sent to the CPU or transferred from CPU to Cache. In response to the original access request.

(B) Direct Mapping -

- An Alternative and simpler address mapping tech for cache
- The lower order s bit of each block identifies the Unique Cache set that can store a block. the remaining high order bit constitute a tag and only these bits need to be stored in the cache tag memory. Then the lower order d bits form the displacement address of the word with its block.

Note: This is easy to implement but it is not very flexible.

- The main drawback of direct mapping is that the caches hit ratio drop sharply if two or more on to the same

(C) Set associative mapping - The combinations of direct and associative tech.

- The blocks of the cache are group into sets and the mapping allows a block of main memory to reside in any block of a specified set at the same time the hardware cost is reduced by decreasing the size of associative search

the tag field of the add. must then be associatively compared to the tags of the two blocks of the set to check if the desired block is present or not. This 2-way association search is simple to implement.

- The no. of blocks per set is a parameter that can be selected to suit the requirement of a particular comp.
- ⇒ Virtual Memory - It is a concept used in large computer systems that permit the users to construct his programs as though he had a large memory space to the totality of auxiliary memory. A virtual M. is used to give the programmer the illusion that he has very large memory at his disposal even though the computer actually has a relatively small main memory. An address used by the programmer is called a "virtual add" and set of such add. is called "add space".
- M. Mapping table is used to Transform virtual add. to Physical address.
- The virtual Memory is divided into pages and the main memory is divided into blocks. Generally the size of page

must be equal to the size of blocks. Random access memory page table tech. uses the no. of Memory location required to store memory mapping table equal to the no. of pages available in U.M.

- In associative M. Page table tech. the no of M. locations required to store memory mapping table equal to the blocks available in M.m.
- The wastage of Memory is minimum in the case of associative mem. Page table tech.
- The most commonly used page replacement algorithm in the VM are
 - (1) FIFO
 - (2) LRU.

I/O data transfer

- Interrupting
- ⇒ Interfacing — writing s/w inst. and designing hardware ckt to enable the CPU to communicate with peripheral devices is called Interfacing and the hardware ckt is called Interfacing ckt.
- ⇒ Need for Interfacing — Most of the peripherals are electromechanical dev. and their manner of op. is diff from the op. of CPU and Memory which are electronic devices.
- The data transfer rate of peripheral is ~~weak~~ must lower than the transfer rate in central comp.
- The op. of peripherals must be synchronised with op. of CPU and M. Unit.
- Data format in peripheral differ from the word format in central processor
- The op. of each peripheral must be controlled so as not to disturb the op. of the controlled comp. and other peripherals conn to the system

- I/O process: It is sometimes called "A Data-channel controller"
 - There are 2-ways to connect I/O devices to CPU
 - (1) Memory Mapped I/O
 - (2) Isolated I/O or I/O mapped I/O
 - (1) Programmed I/O — CPU has direct control over I/O for
 - sending status
 - Read/write commands
 - transferring data
 - CPU waits for I/O module to complete operation
 - limitation: wastage of time
- ⇒ steps - how program I/O work —

- (1) CPU request I/O op.
- (2) I/O module performs op
- (3) I/O module sets status bits
- (4) CPU checks status bits periodically
- (5) I/O module doesn't interrupt CPU
- (6) CPU may wait or come back later

→ ~~end~~

⇒ I/O mapping —

- (1) M.M. I/O — devices and memory share an address space
- I/O looks just like memory read/write op.

- No special commands are given for I/O.
- Large selection of Memory access commands are available.

(a) Isolated I/O -

- Separate add. space for device and M.
- Needs I/O or memory select lines.
- There are special commands for I/O.
- Limited set of Memory access commands are available.

⇒ Interrupt driven I/O -

- It overcomes CPU waiting.
- No repeated CPU checking of device.
- I/O module interrupts when ready.

⇒ Basic Op. of interrupt driven I/O.

(1) CPU read command

(2) I/O module gets data from peripheral while CPU does other work.

(3) I/O module interrupts CPU

(4) CPU requests data

(5) I/O module transfers data

→ In Memory mapped I/O technique the I/O -

under that assumption they will be given address. Same control lines are used to activate memory location at

i/o devices.

- In i/o on isolated i/o tech, the i/o devices are given separate add. space and control signals to activate M1 loadings and i/o devices.
 - Data transfer b/w CPU and peripheral is handled in one of three possible modes
 - (1) Data transfer under program controlled
 - (2) Interrupt initiated data transfer
 - (3) Direct Memory access transfer
 - Program controlled op. are the result of i/o instructions written in comp. programmes
 - each data item transfer is initiated by an instruction in the program.
- ⇒ Disadvantage of program controlled data transfer.
- (1) The processor stays in a program mode until the i/o unit indicates that it is ready.
 - This is a time consuming process since it keeps the processor busy needlessly.
 - In interrupt initiated data transfer, when the P.P is ready for data transfer it generates an interrupt request to the processor. When the processor stops momentarily the task it is doing, branches to a service routine to process

12.1, 12.2, 12.3, 12.5, 12.6

11.2, 11.3, 11.4, 11.6,

the data transfer and then returns to the task it was performing.

- In direct Memory access the interface transfers data into and out of Memory unit through the memory bus generally to transfer bulk amount of data from memory to peripheral or from peripheral to CPU. There are two formats of data transfer.
 - Parallel
 - Serial.
- In parallel mode, data bits (usually a byte) are transferred parallelly over the communication line (known as buses). Thus all the bits of a byte are transferred simultaneously within time frame allotted for the transmission.
- In serial data transfer each data bit is sent sequentially over a single data line. In order to implement serial data transfer, the sender and receiver must divide the time frame allotted for the transmission of a byte into some intervals during which each bit is send and receive.
- In serial transmission, the information is transferred in the form of frames. The frame consist of three parts
 - start bit
 - character code
 - stop bit

\Rightarrow RAID

Numericals - Write back, write through
classmate

Ch. - 3, 10, 12, 11, basics of
RAM, ROM

Date _____

Page _____

X Associative Memory, I/O processor.

- start bit is always at logic 0.
stop bit " " 1.
- The data transfer in both the parallel and serial mode of OP can be either synchronous or Asyn.
- In synchronous mode, both source and destination unit works in synchronous with same control signal.
- In A synchronous mode, the source unit and destination unit have their own independent control signals.
- Asynchronous parallel data transfer can be of 3 types
 - Strobe control
 - two wire hand shaking method.
 - three wire hand shaking method.
- Strobe control is of two types (1) Source initiated (2) Destination initiated.
- Similarly two wire hand shaking method is also of two types (1) Source " (2) Destination "
- Exchange of control signals b/w source and destination during the data transfer called as hand shaking.
- The disad. of hand shaking is that it is not possible to conn. ~~more~~ than one destination

Unit to a single source Unit transmission is used to interface ^{serial} I/O device to a bus structure.

- In DMA transfer the CPU initializes the DMA by sending a M. Address and no. of words to be transferred.
- DMA is preferred data transfer for high speed peripheral devices such as Magnetic disks.
- Three modes of DMA operation.
 - continuous DMA / ~~burst~~ mode
 - cycle stealing
 - interleaved
- In burst mode a sequence of arbitarily level data is transferred in a single continuous burst during which the DMA controller is the master of the system.
- In cycle stealing tech, the transfer is done by first checking if the Memory Unit is not used by CPU. and then the DMA steals one memory cycle to access a word in Memory.
- In interleaved DMA, the processor and DMA access of Memory can be interleaved on alternate half cycle of the clock. The M/W during first half is access by DMA

while CPU accesses during the 2nd half.

- ⇒ **Interrupts** - It is an exceptional event which causes the CPU to temporarily suspend the current program being executed. The control is subsequently transferred to some other program referred to as interrupt service routine which specifies the actions to be taken if the exceptional event occurs.

while CPU accesses during the 2nd half.

- ⇒ **Interrupts** - It is an exceptional event which causes the CPU to temporarily suspend the current program being executed. The control is subsequently transferred to some other program referred to as interrupt service routine which specifies the actions to be taken if the exceptional event occurs.