# INDEX

Name : Rakesh Nama.  Subject : C.O.A

Class : _____  Sec : _____  Roll No. : _____

Cache memory

- Introduction to cache memory :



$10^{-12}$ sec

$10^{-9}$ sec

→ cache memory faster than main memory.

- cache hit : If required element present in cache, then it is called 'cache hit'.
  -that

- hit latency : Time taken to find out whether an element present on the cache or not, that is called "hit latency".

- cache miss : If required element not present in the cache, that is called 'cache miss'.

- Miss latency : Time taken to get something from main memory and then place it into the cache and then read that's called "miss latency".

- page hit : If required element present in main memory.

- page fault : If required element not present in main memory.

- Tag directory : Tag directory says that required element present in tag or not.

→ Before accessing main memory we first access page table, Before accessing cache memory we first access Tags.

→ Purpose of cache memory reduce to cost of the system.

locate locality of reference :
(i) special locality.
(ii) temporal locality.

• Locality of reference :- is a term for the phenomenon in which the same values or related storage locations are frequently accessed, depending on the memory access pattern.

( Locality of reference )

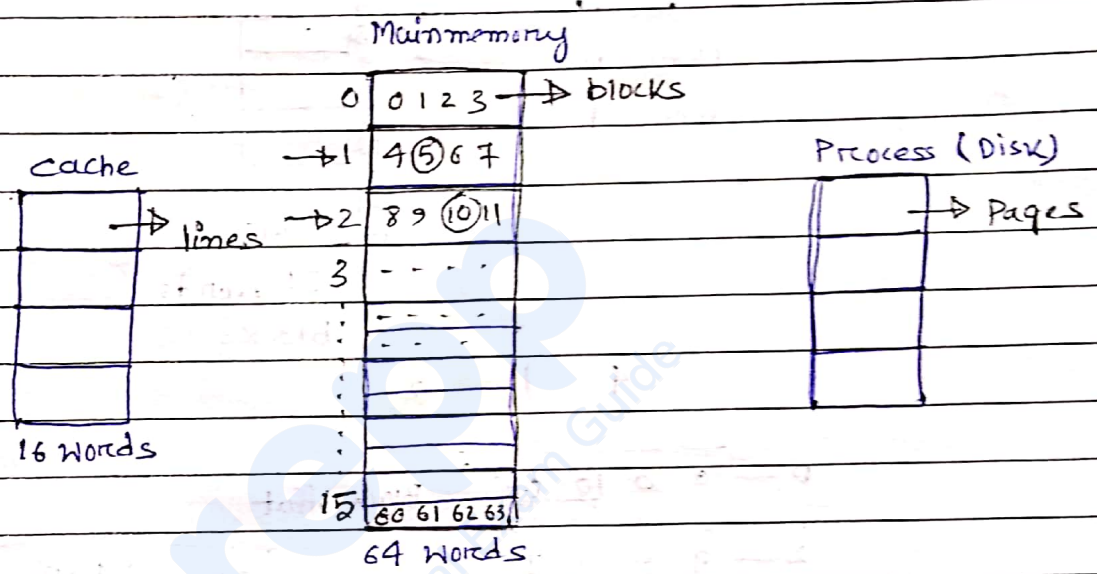| (Temporal locality) | (Spacial locality) |
|---|---|
| Recently refferenced items are likely to be referenced in the near future. | items with nearby addresses tend to be referenced close together in time. |

• Introduction to Direct mapping :

- **Introduction to Direct mapping:**

→ talking about Disk and main memory He talking about pageing.
→ talking about cache and " " . " " Blocks.

→ (Block size = lines size.)

Main memory

0 | 0 1 2 3 → blocks
→1 | 4 ⑤ 6 7
cache
→2 | 8 9 ⑩ 11          Process (Disk)
→ lines →2 | 8 9 ⑩ 11    → Pages
      3 | - - - -
        | - - - -
        | - - - -
16 words
      15 | 60 61 62 63

64 words

→ smallest addressable unit ^called word.
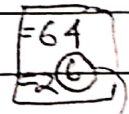                    in the memory
let's

$1W = 1B$ (means oute system is byte addressable)

Block size = 4 words

No. of Block in main memory = $\frac{64}{4}$ = (16 Block.)

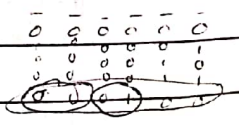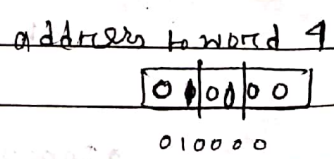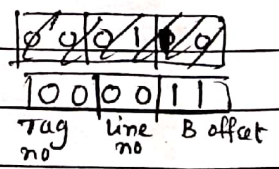No. of lines in cache = $\frac{16}{4}$ = (4 lines)
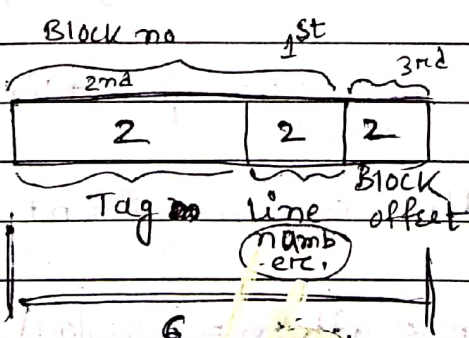
Physical address contain = ₀ 6 bits          $\boxed{64 \atop 2^6}$

**•** processor generate address = | 0 0 0 1 | 0 1 |  (to 50 word 5)
                                      B1K    B1K
                                      num    offset

= | 0 0 1 0 | 1 0 |  (to word 10)
   B.no    B off

• **Direct Mapping :**

cache         Main memory

$Tag_0 \leftarrow$ [00]   0 | 0 4 8 12      0 | 0 1 2 3
$Tag_1 \leftarrow$ [01]   1 | 1 5 9 13      1 | 4 5 6 7
$Tag_2 \leftarrow$ [10]   2 | 2 6 10 14     2 | 8 9 10 11
$Tag_3 \leftarrow$ [11]   3 | 3 7 11 15     3 | 12 13 14 15
               4 | 16 17 18 19
16 words           5
lines = 4          6
                7
                8
                9
               10
               11
               12
               13
               14
             15 | 60 61 62 63

       64 words
       Block = 16

    4        2

0 — 0 0 [0 0]    Block offset
1 — 0 0 0 1
2 — 0 0 1 0
3 — 0 0 1 1
4 — 0 1 [0 0] ✓
5 — 0 1 0 1
6 — 0 1 1 0
7 — 0 1 1 1
8 — 1 0 [0 0] ✓
9 — 1 0 0 1
10 — 1 0 1 0
11 — 1 0 1 1
12 — 1 1 [0 0] ✓
13 — 1 1 0 1
14 — 1 1 1 0
15 — 1 1 1 1

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

          0
          1
          2
          3
          4

32 16 8 4 2 1
0 0 0 1 0 0
* [0 0] [1 0] [0 1] — 9
   line   Tag   Blk
   offt   L.N   off

address to word 3

Block no     1st        [0 0 0 0 1 1]
     2nd           3rd
| 2 | 2 | 2 |      [0 0 | 0 0 | 1 1]
            Block    Tag   Line   B offset
Tag    line   offset      no     no
     number

         6        address to word 4

                [0 0 0 0 0 0]
                0 1 0 0 0 0

1 Bit ज्यादा   Byte = 8 bit.

KB = $2^{10}$ B

MB = $2^{10}$ × KB

GB = $10^{10}$ MB

atlantis **prepp**
Your Personal Exams Guide

Date ___
Page 5

## Direct Mapping   Problem – ①

[Problem – 1]

| | MM size | Cache size | Block size | Tag Bits | Tag Directory size |
|---|---|---|---|---|---|
| Q1 | 128 KB | 16 KB | 256 B | ✓3 bit | ✓$(3*2^6)$ bit |
| Q2 | 32 GB | 32 KB | 1KB | ✓20 bit | ✓$(20*2^5)$ " |
| Q3 | $2^{26}$ B | 512 KB | 1KB | 7 | $(7*2^9)$ " |
| Q4 | 16 GB | $2^{19}$ B | 1KB | 10 | ✓10 $\times 2^{14}$ |
| Q5 | 64 MB | $2^{16}$ B | can't guess | 10 | can't guess |
| Q6 | $2^{26}$ B | 512 KB | can't guess | 7 | can't guess |

Assuming that memory is Byte addressable.

**Q1:**

Main memory size = 128 KB

$\qquad$ = $2^{10}$ × 128 B

$\qquad$ = $2^{17}$ B

Block size = 256 B

$\qquad$ = $2^8$ B

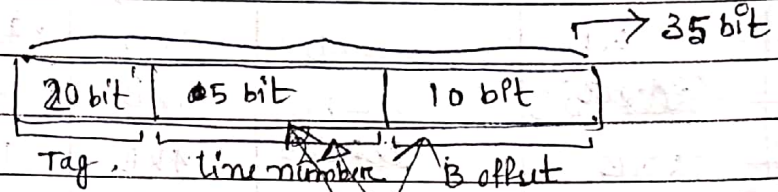no. of Block = $\dfrac{2^{17}}{2^8}$ = $2^9$ - bit

Cache size = 16 KB = $2^{14}$ B

$\dfrac{2^{14}}{2^8}$ = 6

→ line number

$\overset{17 B}{\longleftrightarrow}$

| 3 bit | 6 bit | 8 bit |
|---|---|---|
| Tag | Block number (9) | B offset |

no. of lines = $\dfrac{C\ size}{line\ size}$ = $\dfrac{2^{14}}{2^8}$ = $2^6$ - bit

$\boxed{\dfrac{Line}{size} = \dfrac{Block}{size}}$

Tag Directory size = (Tag size × no. of lines)

$\qquad$ = $(3 * 2^6)$ bit ✓

**Q2:**

$MM \text{ size} = 32 GB$

$= 2^{10} * 32 MB$

$= 2^{10} * 2^{10} * 32 KB$

$= 2^{10} * 2^{10} * 2^{10} * 32 * B$

$= 2^{30+5} B = 2^{35} B$

$\longrightarrow 35 \text{ bit}$

| 20 bit | 05 bit | 10 bit |
|--------|--------|--------|
| Tag | line number | B offset |

$Block \text{ size} = 1 KB$

$= 2^{10} B$  (10 bit for Block B offset)

✱ $Cache \text{ size} = 32 KB$

$= 2^{10} * 2^5 B$

$= 2^{15} B$

$no. \text{ of line} = \dfrac{C \text{ size}}{\text{line size}}$

$= \dfrac{2^{15}}{2^{10}}$

$= 2^5$

$Tag \text{ directory size} = Tag * no. \text{ of lines}$

$= (20 * 2^5)$

**Q3:**

| | | 26 B | |
|---|---|---|---|
| Tag | line no | B offset | |
| 7 | 9 | 10 | |

$Main \text{ memory size} = ? \ (64)$

$Cache \text{ size} = 512 KB = 2^{19} B$

$Block \text{ size} = 1 KB = 2^{10} B$  (10 bit B offset)

$Tag = 7$

✓ $Tag \text{ directory size} = 7 * no. \text{ of line}$

$= 7 * \dfrac{2^{19}}{2^{10}} = (7 * 2^9) B \text{ bit}$

$no. \text{ of line}$

$\dfrac{2^{19}}{2^{10}} = 2^9$

no. of block = $\dfrac{MM \text{ size}}{block \text{ size}}$

$address \text{ size} = 7 + 9 + 10 = 26$

$Main \text{ memory size} = 2^{26} B$

$= 2^{32} B = 2^6 MB = 64 MB$

$= 2^{30} * 2^2 B = 4 GB$

**Q4:**

$$MM \; size = 16 \, GB = 2^{30} * 2^4 \, B = 2^{34} B$$

| 10 | 12 | 12 |
|----|----|----|
| Tag | line num | Block offset = line offset |

34 bit

$$\begin{array}{r} 34 \\ -22 \\ \hline 12 \end{array}$$

cache size = ?

Block size = $4KB = 2^{10} * 2^2 \, B = 2^{12} B$

tag = 10

tag = Directory.

$$2^{12} * 2^{12} = 2^{24} B$$

cache size = $2^{24}$ ~~= 20 BLOCK 4KB~~

no of lines = $\dfrac{c \; size}{line \; size} = \dfrac{2^{24}}{2^{12}} = 2^{12}$

Tag directory = ~~become~~ $10 * 2^{12}$

= ~~100 bits~~ $(10 * 2^{12})$ bit.

**Q5:**

$$MM \; size = 64 \, MB = 2^{20} * 2^6 \, B = 2^{26} B \; (PA)$$

cache size = ?

Block size = ?
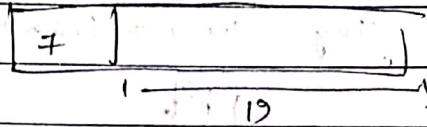
tag = 10

$$\begin{array}{r} 2^{26} \\ 26 \end{array}$$

26 B

| 10 | | |
|----|----|----|
| Tag | | 16 |

cache size = $2^{16} B$

**Q6:**

Cache size = 512 KB = $2^{19}$ B

Tag = 7 bit

| 7 | |
|---|---|

$\longleftarrow$ 19 $\longrightarrow$

$19 + 7 = \boxed{26}$

✓ MM size = $2^{26}$ B

- **Direct Mapping HW Implementation** —

CPU generated address $\rightarrow$

| Tag | LN | BD |
|-----|-----|-----|
| 0 | 0 | 01 |

*Comparator basic gate*

XNOR.

| a | b | y |
|---|---|---|
| 0 | 0 | 1 |
| 1 | 1 | 1 |

1×4 mux
1ns

| 00 | |
|----|---|
| 01 | ✓ |
| 10 | |
| 11 | |

Comparator
1ns
90ns

1 - hit
0 - miss

Total time taken = (latency of MUX + latency of comparator)
= 1ns + 1ns
= 2ns. (nanosecond)

→ no of MUX and comparator depends on no. of Tag bits.

**Direct mapping** → 

Tag size = K bit

no. of MUX require = K

no. of Comparator = 1 × (K bit compare)
↳ always (In Direct map)

**Q1:**

MM size = 1 GB

Cache size = 1 MB

Comparator size = 10 Kns.

Hit latency = ?

⟢ *Here latency MUXs negligible*

⟹ Hit latency a

MM = $2^{30}$ B

** Tag = $\dfrac{\text{Main M size}}{\text{Cache size}}$

$= \dfrac{1 GB}{1 MB}$ ⟹ 2(10)KB 1KB ⟹ $2^{10}$ B    $K = 10$

Tag: [ 10 | | ]

Hit Latency = K × 10 Kns. 10 Kns

$= 10 × 10$ ns

$= 100$ ns.

→ [ line number ≠ Block number / no. of lines ] ⟩ line number = $\dfrac{(B \text{ number})\%}{(\text{no. of lines})}$

$K = m \% n$

• Disadvantage of direct mapping =

→ Conflict miss problem.

ex:    Block req by CPU –

| | | |
|---|---|---|
| 0 | 4  8,12,20 | 5, 6, 4, 8, 9, 12, 15, 20 |
| 1 | 5  9 | |
| 2 | 6 | |
| 3 | 15 → | |

$LN = (B.N)\% (NL)$

→ place was empty for long time, but line no. 0 used heavily many time, this is conflict miss problem.

Cache size = no. of lines = 4

**ex:**

Block request by CPU

| 0 | 4 | 8 12 16 20 |
|---|---|---|
| 1 | | |
| 2 | | |
| 3 | | |

Cache.

4  8  12  16  20

$K = m \% n$

$n =$ num of lines in cache

$4 \% 4 = 0$

$8 \% 4 = 0$

line no, 1, 2, 3, empty but then line 0 are heavily used it is conflict miss problem.

- **Associative mapping = (Intro)**

can solve

→ By Associative mapping ∧ Conflict miss problem.

→Tag

Phical address →

| . 4 | 1 2 |
|---|---|
| B·NO | Block offset |

$2^4 = 16 B$

[BNO [BO]



Com parator

4 → (0-15)
→ (0-15)
→ (0-15)
→ (0-15)

Cache

1 - cache hit
0 - cache miss

→ 4 lines — 4 comparator need, so hardware cost increase along with freedom.

→ *no of comparator requires = no of cache lines.

→ We get freedom to place Block any where in cache, but requirement of comparator increase.

**Q1:**

Main memory size = 32 GB
Block size = 32 KB
Tag(K) = ?
Propagation delay (PD) of Comp = 10 Kns
PD of OR Gate = 10ns.
Hit latency = ?

$\Rightarrow$

| 35 | |
|---|---|
| 20 | 15 |

Tag(K)

K = 20

MM = 32 GB
$= 2^{30+5}$ B
$= 2^{35}$ B

BS = 32 KB
$= 2^{15}$ KB

hit latency = PD of comparator + PD of OR gate
$=$ 10 Kns + 10 ns
$= (10 * 20)$ ns + 10 ns
$= 210$ ns.

**Q2:**

| | MM size | cache size | Block Size | Tag size | Tag directory size |
|---|---|---|---|---|---|
| ① — | 128 KB | 16 KB | 256 B | 9 bits | $(9 \times 2^6)$ bits |
| ② — | 32 GB | 32 KB | 1 KB | 25 bits | $(25 \times 32)$ bits |
| ③ — | 128 MB | 512 KB | 1 KB | 17 | $(17 \times 2^9)$ bits |
| ④ — | 16 GB | ? (can't) | 4 KB | 22 bits | (can't guess) |
| ⑤ — | 64 MB | ? (can't) (any) | 64 KB | 10 | (can't guess) |
| ⑥ — | not possi | 512 KB | not possi | 7 | not possible. |

① MM size = 128 KB

   CS = Cache size = 16 KB          no of lines = $\frac{CS}{LS(BS)}$

   Block size (BS) = 256 B                     = $\frac{2^{14}}{2^8}$

   Tag size = ?                                = $2^6$

   Tag directory size = ?

⟹



          ← Tag = 9

   Tag directory size = Tag size * no of lines.
                      = $(9 * 2^6)$ bits

   (no of com req = $2^6$ = 64)

   BB = 128
      = $2^7 * 2^{10}$ B
      = $2^{17}$ B

   BS = 256 B = $2^8$ B

② MM = $2^{35}$ B

   CS = $2^{15}$ B          Tag = (35 - 10) = 25

   BS = $2^{10}$ B

   no. of Line = $\frac{2^{15}}{2^{10}}$          Tag directory = Tag * no of lines
                                                = (25 * 25)

              = $2^5$

③ MM = ?                    MM size = $2^{10+17}$

   CS = $2^{19}$ B                  = $2^{27}$ B = $2^7$ 128 MB

   BS = $2^{10}$ B                  = 128 MB

   Ts = 17                  no. of line = $\frac{2^{19}}{2^{10}}$ = $2^9$

   Tag directory = ?

                            Tag D.S = $(17 * 2^9)$ bits.

④ MM = $2^{34}$ B

cache CS = ? ✗

BS = $2^{12}$ B

✓ TS = ? 

T.DS = ? ✗

Tag size = $(34 - 12)$

= $(22)$

⑤ MM = $2^{26}$ B | BS = $2^{26-10}$

CS = ?

BS = ? $(2^{16})$ | = $2^{16}$ = 64 KB

Tag = 10

Tag D·S = ?

⑥ MM = ? ✗

CS = $2^{15}$ B

TS = 7

BS = ? ✗

T.DS = ? ✗   not possible.

3-Way ut

Sets



3 {

3 {

cache

memory

Scanned by CamScanner

• • set Associative mapping = (advi no. of comparator reduced)

EX: (how set associative work)

MM size = 64 B

cache s = 32 B

Block S = 4 B

✓ set size = 2 Blocks (lines)

✓ (2-way set associative) OR

$$lines = \frac{CS}{BS} = \frac{32}{4} = 8 \text{ lines}$$

$$sets = \frac{Lines}{sets} = \frac{8}{2} = 4 \text{ sets}$$

give address (P-A) = | 01 | 10 | 11 |
                              SNo



| set 0 | Lin-0 | B·NO- 0 |
|       | L-1   | BN - 1  |
| s1    | L-2   | 2       |
|       | L-3   | 3       |
| s2    | L-4   | 4       |
|       | L-5   |         |
| s3    | L-6   |         |
|       | L-7   |         |

cache M

B·NO      B·offset

P·A = | 2 | 2 | 2 |
       Tag  Set·No

                          61
                          62
              BN- 63

main M

**Q1:**

| MM size | cache size | Block size | Tag bits | Tag directory size | set associative |
|---|---|---|---|---|---|
| ① 128 KB | 16 KB | 256 B | 4 | $(4*2^6)$ bits | 2-way ut |
| ② 32 GB | 32 KB | 1KB | 22 | $(22*2^5)$ bits | 4 |
| ③ $2^3$ MB | 512 KB | 1 KB | 7 | $(10*2^9)$ bit | 8 |
| ④ 16 GB | 64 MB | 4 KB | 10 | $(10*2^{14})$ bit | 4 |
| ⑤ 64 MB | 256 KB | X | 10 | X | 4 |
| ⑥ 8 MB | 512 KB | X | 7 | X | 8 |

①



$$MM_S = 2^{17} B \qquad PA = 17\text{-bit}$$
$$CS = 2^{14} B$$
$$BS = 2^8 B$$

$$\text{no. of lines} = \frac{CS}{BS} = 2^6 \text{ lines.}$$

$$\text{set no. of sets} = \frac{\text{no. of lines}}{\text{set size}}$$

$$= \frac{2^6}{2} = 2^5 = (32 \text{ sets})$$

$$PA = \text{Tag} + \text{set no} + B.\text{offset}$$
$$\Rightarrow 17 = \text{Tag} + 5 + 8$$
$$\Rightarrow \text{Tag} = 4$$

$$\text{Tag directory} = \text{Tagsize} * \text{no. of lines.}$$
$$= (4 * 2^6) \text{ bits.}$$

② 

$$MM = 32GB = 2^{35}B$$

$$CS = 32KB = 2^{15}B$$

$$BS = 1KB = 2^{10}B$$

$$T = ? \quad (22)$$

$$T.D = ?$$

$$Set\text{-}as = 4\text{-way}$$

$$no.\,of\,line = \frac{CS}{BS}$$
$$\qquad = \frac{2^{15}}{2^{10}} = 2^5$$

$$no.\,of\,set = \frac{no.\,of\,line}{Set\,size} = \frac{2^5}{2^2} = 2^3$$

$$\rightarrow P.A = Tag + s.no + B.off$$

$$\rightarrow 35 = Tag + 3 + 10$$

$$\rightarrow Tag = 22$$

$$T.D = (22 * 2^5)$$

③ $$MM = ?$$

$$CS = 512KB = 2^{19}B$$

$$BS = 1KB = 2^{10}B$$

$$T = 7$$

$$T.D = ?$$

$$set\text{-}as = 8\text{-way}$$

$$\qquad \quad set 2^3$$

| 7 | 6 | 10 |
|---|---|---|
| Tag | Set no | B.0 |

$$= 23\,bit$$

&

$$no.\,of\,line = \frac{2^{19}}{2^{10}} = 2^9$$

$$no.\,of\,sets = \frac{no.\,of\,line}{size\,of\,sets}$$

$$MM\,size = 2^{23} = 2^3 MB$$

$$= \frac{2^9}{2^3} = 2^6$$

$$T.D = T * no.\,of\,lines.$$

$$= (7 * 2^9)\,bits.$$

④ $MM = 16 GB = 2^{34} B$

$CS = ?$

$BS = 4KB = 2^{12} B$

$T = 10$

$T \cdot D = ?$

set ass $= 4$-way set.

$$\begin{array}{|c|c|c|} \hline 10 & 12 & 12 \\ \hline \end{array}$$
Tag   set no   B. of

$no. of sets = 2^{12}$

$T \cdot D = 8 \, T * no. \, of \, lines$

$= 10 * 2^{14}$

$= \boxed{160 KB}$

$no. \, of \, sets = \dfrac{no. \, of \, lines}{9 \, set \, size}$

$no. \, of \, lines = set \, size * no. \, of \, sets$

$C \cdot S = no. \, of \, lines * line \, size$

$= 2^{14} * 2^{12}$

$= 2^{26} B$

$= 2^6 MB$

$= \cancel{256} \; \boxed{64 MB}$

$= 4 * 2^{12}$

$= 2^{14}$

⑤ $MM = 64 MB = 2^{26} B$

$CS = - \quad (2^{18} B)$

$BS = -$

$Tag = 10$

$tag \cdot D = -$

$set \cdot ass = 4 \, way$

$$\begin{array}{|c|c|c|} \hline 10 & ? & ? \\ \hline \end{array}$$
Tag   set NO   B-offset

$$\begin{array}{|c|c|c|} \hline 10 & x & y \\ \hline \end{array}$$
$\underbrace{\qquad\qquad}_{16}$

⑥ MM = Y

CS = 512 KB

BS = -

Tag = 9

T. D = -

set ass = 8 way.

$$\begin{array}{|c|c|c|c|c|c|c|} \hline & & & & & & \\ \hline \end{array}$$
Tag

Cache size $= no. \, set * lines \, per \, set$
$\qquad\qquad\qquad * line \, size (BS)$

$= 2^x * 2^2 * 2^y$

$= 4 . 2^{x+y}$

$= 4 . 2^{16}$

$= 2^{18} = 2^8 KB = 256 K$

**Q2:**

| MM size | Cache size | Block size | Tag size | Tag directory | Set-associative |
|---------|-----------|-----------|----------|---------------|----------------|
| 64MB | — | — | 10 | — | 4-way |
| — | 512KB | — | 7 | — | 8-way |

⑥  MM = ? — ($2^{23}$ B)

CS = 512 KB = $2^{19}$ B

BS = —

TS = 7

TDS = —

Set assn = 8-way

| Tag | Set | Block |
|-----|-----|-------|

| 7 | $x$ | $y$ |

cache size = no. of sets $*$ line's per set $*$ Block size

$2^{19} = 2^x * 2^3 * 2^y$

$\Rightarrow 2^{16} = 2^{x+y}$

$\Rightarrow x+y = 16$

P.A = $7 + (x+y)$

$= 7 + 16 = \boxed{23}$

MM = $2^{23}$ = 8 MB

- ## Comparing all the mappings :-

$MM \ size = 4GB \ = 2^{32} \ B$

$Cache = 4MB \ = 2^{22} \ B$

$Block \ size = 1KB = 2^{10} \ B$

$P.A = 32 \ Bits.$

$no. \ of \ lines = \dfrac{CS}{BS} = \dfrac{2^{22}}{2^{10}} = \boxed{2^{12} \ lines}$

$no. \ of \ set = \dfrac{no. \ line \ \#}{size \ of \ set} = \dfrac{2^{12}}{2^{2}}$

$= 2^{10}$

**Direct mapping :-**

| | B.NO | | B.O |
|---|---|---|---|
| 10 | 12 | | 10 |

Tag    line. no

← 32 bits →

**associative Mapping :-**

← 32 bits →

| 22 | 10 |
|---|---|

B.NO (Tag)    BO Offset

← 32 bits →

B.NO

**4-way set asso :-**

| 12 | 10 | 10 |
|---|---|---|

Tag    set NO    BO

Relationship between

| | TAG | Comparator replaces | Size of comparator (Tag directory) | Lines | size of comparator |
|---|---|---|---|---|---|
| Direct | 10 | 1 | $10 * 2^{12}$ | $2^{12}$ | 10 |
| Associative | 22 | $2^{12}$ | $22 * 2^{12}$ | $2^{12}$ | 22 |
| 4-way | 12 | 4 | $12 * 2^{12}$ | $2^{12}$ | 12 |

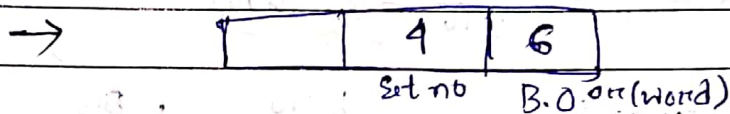**Questions** - gate

g-1995   (Q.1)   Cache size = 4K words  = $2^{12}$ words
Block size = 64 words. = $2^6$ words
set size = 4 blocks.
The number of bits in "set" and "word" field
of MM address are :-
$\underset{(4)}{}$   $\underset{(6)}{}$

→
| | | 4 | 6 |
|---|---|---|---|

Set no   B.O or (word)

$$no\ of\ sets = \frac{no\ of\ lines}{set\ size} \qquad no. of\ lines = \frac{2^{12}}{2^6} = 2^6$$
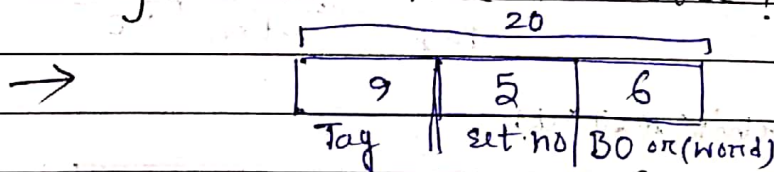
$$= \frac{2^6}{2^2} = 2^4 \text{ (sets)}$$

(Q.2)

4-way set associative.
cache lines = 128
lines size = 64 words
P.A = 20 bits
Tag, set and word fields are ?

$$\overbrace{\qquad\qquad}^{20}$$

→
| 9 | 5 | 6 |
|---|---|---|

Tag   set no   BO or (word)

$$no. of\ lines = \frac{cache\ size}{lines\ size}$$

$$\hookrightarrow cache\ size = 2^6 * 2^7 = 2^{13}\ words.$$

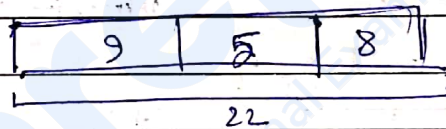$$sets = \frac{lines}{set\ size} \Rightarrow \frac{2^7}{2^2} \Rightarrow 2^5$$

Tag = 9,   set = 5,   words(BO) = 6.

**Q3)** Blocks in cache = 128
4-way set associative
MM contains $2^{14}$ blocks.
Block size is 256 eight bit words.

(i) How many bits are required for addressing MM ? (22)

(ii) How many bits are needed to represent the TAG, set and word fields.
 (9)    (5)    (8)

→ cache lines = $2^7$

 set size = 4

MM Blocks = $2^{14}$

 Block size = $2^8$ bit/word,

| 9 | 5 | 8 |
|---|---|---|

22

$sets = \dfrac{2^7}{2^2} = 2^5$

no. of lines

MM size = $2^{14} \times 2^8$ words

$= 2^{22}$ words

**Q4)**

Direct mapped cache
cache size = 32 KB = $2^{15}$ B
Block size = 32 B = $2^5$ B
PA = 32 bits

The number of bits needed for cache indexing and tag bits are respectively.

32

| 17 | 10 | 5 |
|----|----|---|
| tag | line no | B.O |

no. of lines = $\dfrac{2^{15}}{2^5} = 2^{10}$

cache indexing = 10 bit
tag bit = 17 .  Ans (10, 17)

(9 - 2006) (Q5)
1.20-1.21

Consider two cache organization

① first one                    ② 2nd one
set associative                 different
Cache size = 32 KB = $2^{15}$ B    Cache size same
2-way set associative           Direct mapped
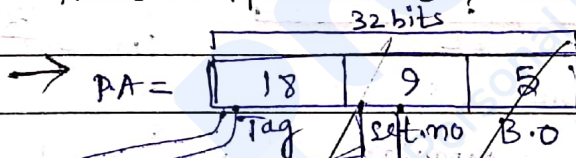Block size = 32 B = $2^5$ B
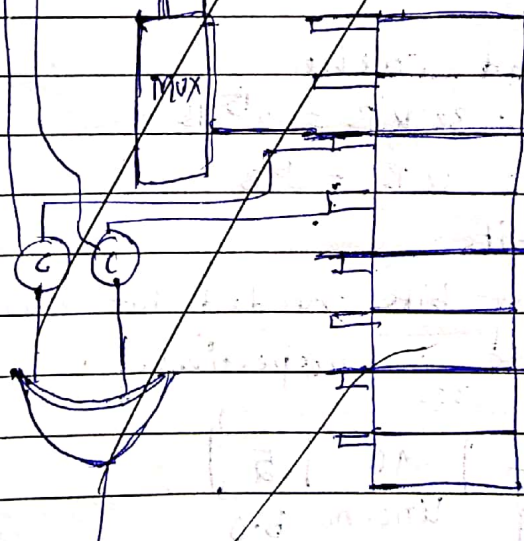
P.A = 32 bits

2-to-1 MUX latency = 0.6 ns.
k bit comparator has latency = $k/10$ ns.

The hit latency of set associative organization is $h_1$
and " " Direct mapped is $h_2$.
find $h_1$ and $h_2$?

32 bits

PA = | 18 | 9 | 5 |
        Tag  set no  B.O

no of comparator requires = $2^1$

$\frac{32}{19}$
 18

MUX

C  C

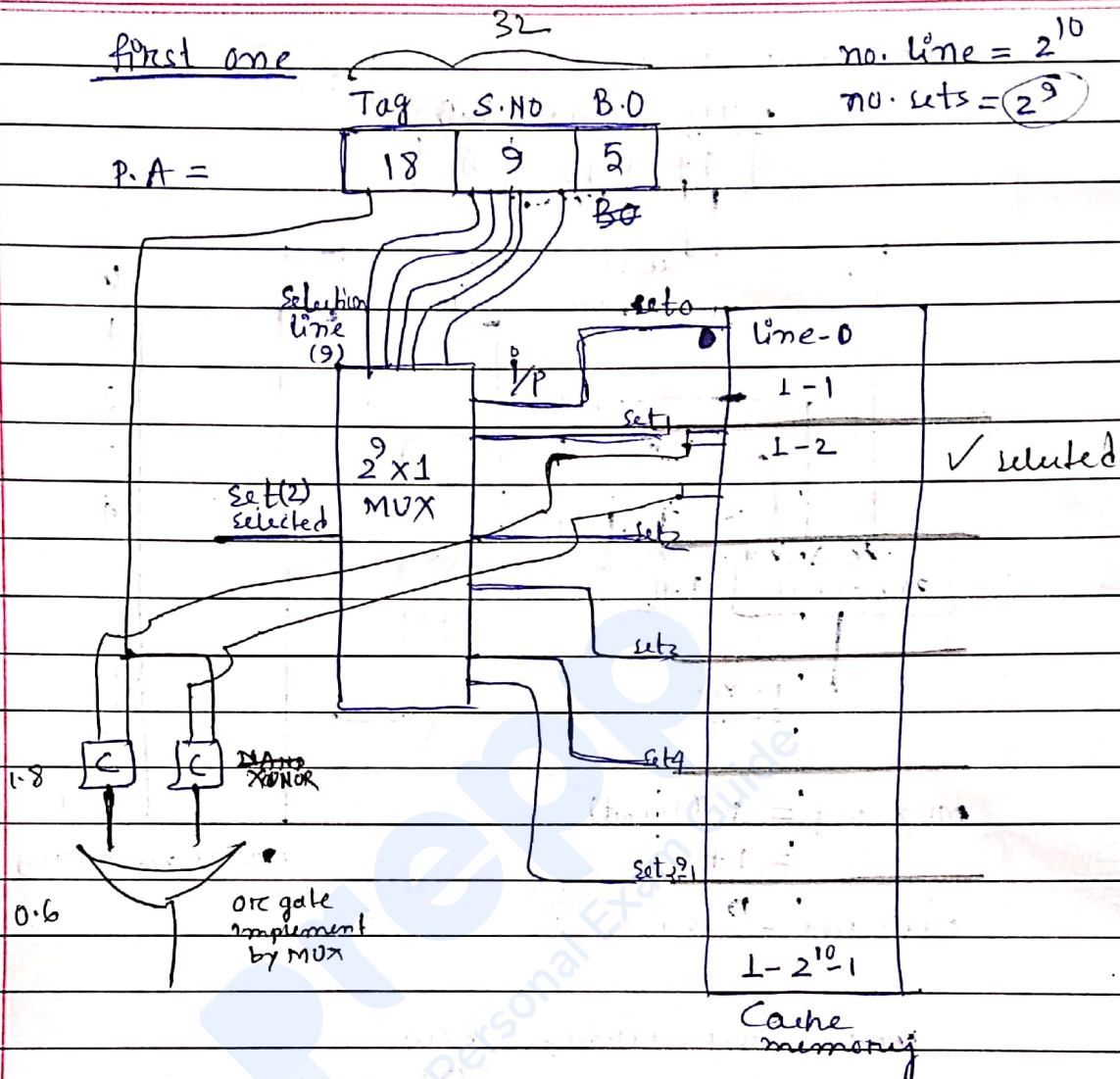line no = $\frac{2^{15}}{2^5} = 2^{10}$

no. of set = $\frac{2^{10}}{2^1} = 2^9$

here, MUX required
= 2 × Tag bits
= K × Tag bits

hit latency =

first one

$$\overbrace{\hspace{3cm}}^{32}$$

no. line $= 2^{10}$

no. sets $= \boxed{2^9}$

| Tag | S.NO | B.O |
|-----|------|-----|
| 18 | 9 | 5 |

P.A =

B.O

Selection line (9)

Set(2) selected

$2^9 \times 1$ MUX

i/p

Set(2) selected

NAND XNOR

1.8

C   C

OR gate implement by MUX

0.6

| | line-0 |
|---|--------|
| set₀ | |
| set₁ | L-1 |
| | L-2 ✓ selected |
| set₂ | |
| set₃ | |
| set₄ | |
| set ₂₉ ₋₁ | |
| | L-2¹⁰-1 |

Cache memory

MUX req = Tag hit * 2

K-way set associative = (MUX.req) = T * K.

size of each MUX = $2^5$ to 1

Comparator need = K

comparator latency = $K/10 = {}^{18}/10 = 1.8$ ns.

$h_1 = 1.8 + 0.6$
$\quad = 2.4$ ns

**Second one -**

$$\begin{array}{|c|c|c|}\hline \text{Tag} & \text{L.NO} & \text{B.O} \\ \hline 17 & 10 & 5 \\ \hline \end{array} \; -32\,bit$$



Selection line

$2^{10}$ to 1

i/p

Comparator 1·7

1 - Yes
0 - No

0
1
2
3
4

$2^{10}-1$

cache memory

$$MUX \; ref = K \, (\text{Tag bit})$$
$$= 17$$

$$each \; size = 2^{10} \; to \; 1$$

$$comparator \; latency = K/10$$
$$= 17/10 = 1·7$$

$$\boxed{h_2 = 1·7\,nS} \checkmark$$

g-2011

**(Q6)** Direct mapping
cache size = 8KB
BS = 32 byte.
PA = 32 bits.

the cache controller maintains tag information
for each cache block comprising of following.

1 valid bit, 1 modified bit and as many
bits as the minimum needed to identify block
mapped in the memory block mapped in the cache

What is the total size of memory needed at the cache controller of store meta-data (tags) for the cache?
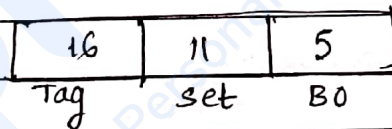
$$\text{lines} = \frac{CS}{BS} = \frac{2^{13}}{2^5} = 2^8$$

| 19 | 8 | 5 |
|----|---|---|

$(19+1+1) = 21$

$$\text{T.Directory} = (21 * 2^8) \text{ bits}$$
$$= 21 * 256$$
$$= 5376 \text{ bits}$$

P-229
1.40-1.41
g-2002

**(Q7)** cache size $= 256 KB = 2^{18} B$

set size $= 4$

$B.S = 32 B = 2^5 B$

$P.A = 32 \text{ bits}$

| 16 | 11 | 5 |
|----|----|---|
| Tag | set | BO |

$$\text{no of lines} = \frac{CS}{BS}$$
$$= 2^{13}$$

① The number of bits in the tag field of an address — (16).

$$\text{sets} = \frac{\text{lines}}{2^2}$$
$$= \frac{2^{13}}{2^2} = 2^{11}$$

② The size of the cache tag directory is — lines * tag

$$2^{13} * (16+2+11) \qquad = 2^{13} * 16$$
$$= 2^{13} * 20 \qquad = 2^{13} * 4$$
$$= 2^3 * 20 KB \qquad = 2^{17}$$
$$= 160 KB \qquad = 2^7 KB$$
$$= 128 KB$$

9-2019 (88)

4-way set associativity

Cache size = 16 KB $= 2^{14}$ B

Block size = 8 Words $= (8 * 32)$ bits $= \dfrac{2^8}{2^3}$ Bo $= 2^5$ B

Word size = 32 bit

DAS = 4 GB $= 2^{32}$ B

Tag bits = ?



$$
\begin{array}{|c|c|c|}
\hline
20 & 7 & 5 \\
\hline
\text{Tag} & \substack{\text{Set} \\ \text{no}} & \text{B·O} \\
\end{array}
$$

(Tag bits = 20)

$lines = \dfrac{CS}{BS}$

$= 2^9$

$Set = \dfrac{lines}{set size} = \dfrac{2^9}{2^2}$

$= 2^7$

## Computer Architecture:

It deals with Instructions, addressing modes, ALU, pipelining etc (Internal Design)
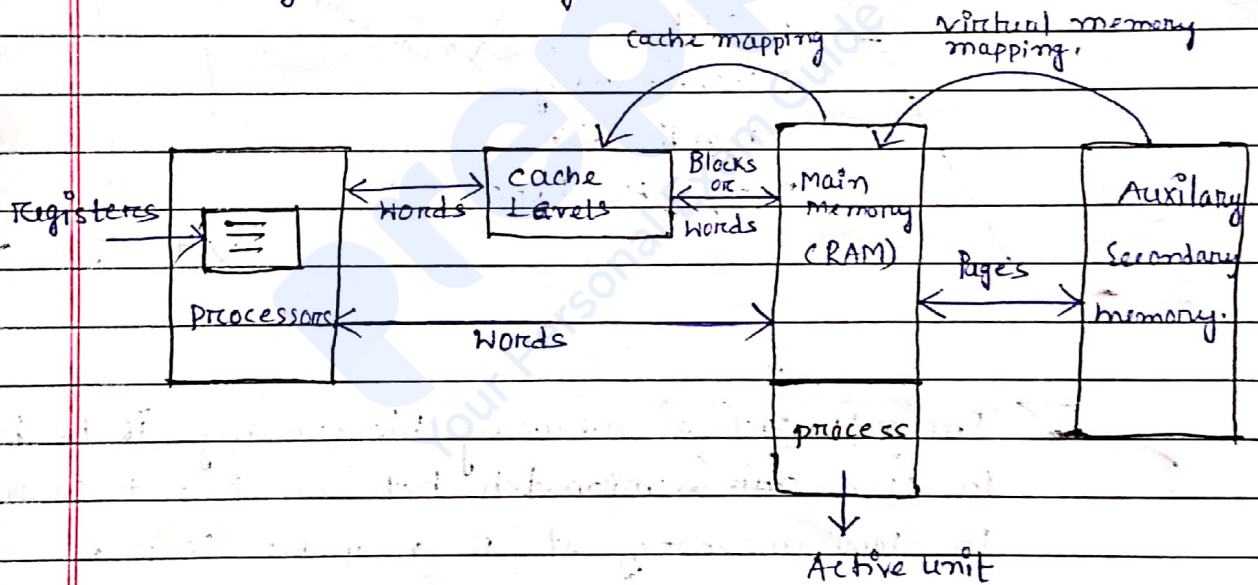
## Computer Organization:

It deals with how various memory and I/O interact with a system.

## Computer design:

It deals with hardware design.

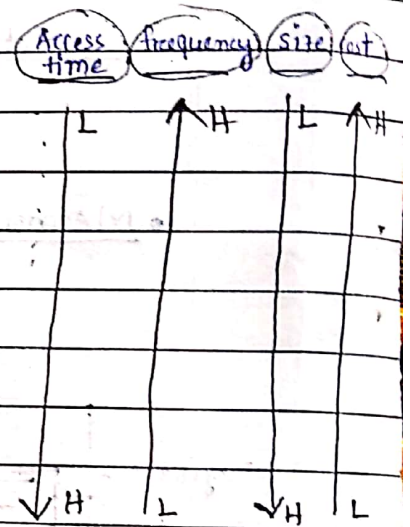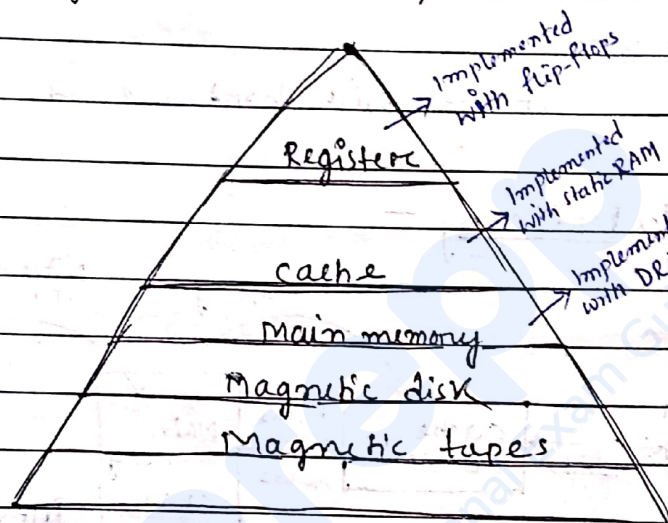• ## Memory Interfacing —

- <u>Memory Hierarchy :</u>

cache levels
Main memory } Random Access.

Magnetic disk → semi Random Access.

Magnetic tapes → sequential Access.

| | Access time | Frequency | Size | cost |
|---|---|---|---|---|
| | L | H | L | H |

Register — Implemented with flip-flops

cache — Implemented with static RAM

Main memory — Implemented with DRAM

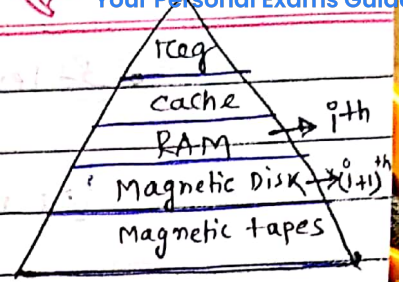Magnetic disk

Magnetic tapes

| | H | L | H | L |

→ The purpose of memory hierarchy is to bridge the speed mismatch between fastest processor to slow memory at reasonable cost.

→ The goal of memory hierarchy is to minimize average access time of entire memory system.

- <u>Level memory —</u>

• 2 level memory –

Information in

$i^{th}$ $(i+1)^{th}$ level

Tag
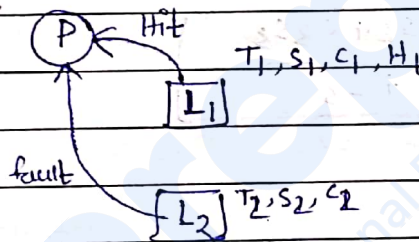cache → $i^{th}$
RAM →
Magnetic Disk → $(i+1)^{th}$
Magnetic tapes

→ If processor referes to ith level memory is found then "Hit" otherwise "Miss (or) fault".

→ There are two way in which the processor is connected to various levels of memory

Case 1: $T_1 < T_2$

Cases

P  Hit  $T_1, S_1, C_1, H_1$

$L_1$

fault

$L_2$ $T_2, S_2, C_2$

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

$$Cost_{avg} = \frac{C_1 S_1 + C_2 S_2}{S_1 + S_2}$$

Hit rate $= \dfrac{x}{100}$

Miss rate $= \left(1 - \dfrac{x}{100}\right)$

$T_1 \rightarrow$ Time to access
$S_1 \rightarrow$ size of level 1 memory.
$C_1 \rightarrow$ cost per bit.
$H_1 \rightarrow$ Hit rate.

$100 \rightarrow x$

$$T_{avg} = \frac{x T_1 + (100 - x) T_2}{100}$$
$$= H_1 T_1 + (1 - H_1) T_2$$
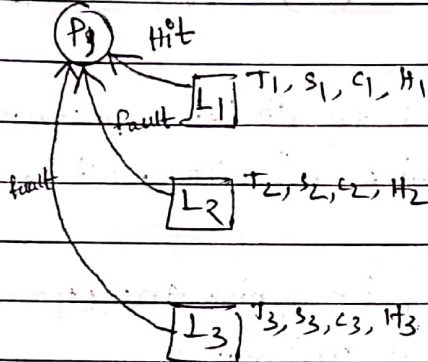
Case 2:

P

$L_1$ $T_1, S_1, C_1, H_1$

$L_2$ $T_2, S_2, C_2$ $\boxed{T_{avg} = H_1 T_1 + (1 - H_1)(T_1 + T_2)}$

$$Cost_{avg} = \frac{S_1 C_1 + S_2 C_2}{S_1 + S_2}$$

- ## 3-level memory –

case – 1

$$T_1 < T_2 < T_3$$



$P_g$ ↻ Hit

fault $\boxed{L_1}$ $T_1, S_1, C_1, H_1$

fault $\boxed{L_2}$ $T_2, S_2, C_2, H_2$
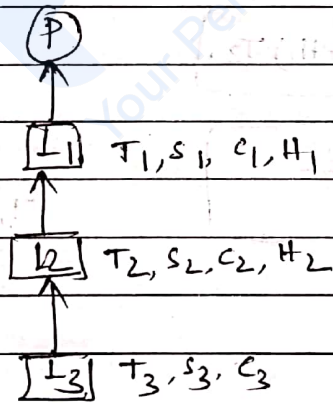
$\boxed{L_3}$ $T_3, S_3, C_3, H_3$

Best case time | Worst time taken

$$T_1 \leq T_{avg} \leq T_3$$

$$T_{avg} = H_1 T_1 + (1-H_1) H_2 T_2 + (1-H_1)(1-H_2) T_3$$

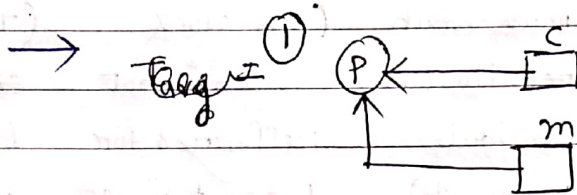$$C_{avg} = \frac{C_1 S_1 + C_2 S_2 + C_3 S_3}{S_1 + S_2 + S_3}$$

Case – 2

$P$

$\boxed{L_1}$ $T_1, S_1, C_1, H_1$

$\boxed{L_2}$ $T_2, S_2, C_2, H_2$

$\boxed{L_3}$ $T_3, S_3, C_3$

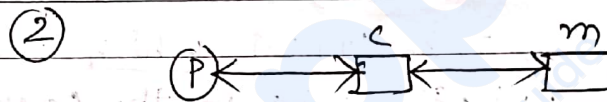Best time taken | Worst time taken

$$T_1 < T_{avg} < (T_1 + T_2 + T_3)$$

$$T_{avg} = H_1 T_1 + (1-H_2) H_2 (T_1 + T_2) + (1-H_1)(1-H_2)(T_1 + T_2 + T_3)$$

$$C_{avg} = \frac{S_1 C_1 + S_2 C_2 + S_3 C_3}{S_1 + S_2 + S_3}$$

**Q:** The average memory access time for a machine with a cache hit rate of 80% where the cache access time is 5ns and memory access time is 100ns is?

→ $T_{avg} = $ (1) $P \leftarrow C$

$m$

$$T_{avg} = H_1 T_c + (1-H_1) T_m$$
$$= (0.8)(5ns) + (1-0.8)(100ns)$$
$$= \boxed{24\ ns}$$

(2) $P \leftarrow C \leftarrow m$

$$T_{avg} = H_1 T_c + (1-H_1)(T_c + T_m)$$
$$= (0.8)(5ns) + (1-0.8)(5ns + 100ns)$$
$$= 4ns + (0.2)(105ns)$$
$$= 4ns + 21.0\ ns$$
$$= \boxed{25\ ns}$$

• **Cache replacement policy :—**

↳ Replacement policy is required for associative mapping and set associative mapping but not for direct mapping.

↳ Replacement policies are aimed to minimize miss penality for ~~references~~ future references.

## Replacement Policies

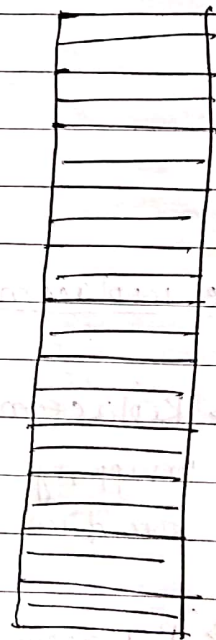| Random | FIFO | LRU | LFU |
|---|---|---|---|
| (No specific criteria to replace the block) | (The block which entered first is the candidate replacement) | (The block which is not refferenced for longest time which is replaced) | (the block which fewer references is replaced) |

Q:1 consider a direct mapped cache with 8 cache blocks (0-7), if the memory block requests are in the orders 3, 5, 2, 8, 0, 6, 3, 9, 16, 20, 17, 25, 18, 30, 24, 2, 63, 5, 82, 17, 24 which one of the memory blocks will be not be in the cache at the end of the sequence?

(a) 3     (b) 18 ✓     (c) 20     (d) 30.

⟹

| | |
|---|---|
| 0 | 8 0 16 24 |
| 1 | 9 17 25 17 |
| 2 | 2 18 2 82 |
| 3 | 3 ✓ |
| 4 | 20 ✓ |
| 5 | 5 |
| 6 | 6 30 ✓ |
| 7 | 63 |

no. of Block = 8
(Line)
cache

$$i = j \bmod 8$$

line
$i \rightarrow$ block no. of cache.

$j \rightarrow$ block no. of MM.

number of lines on cache (8).

main memory

**Q:2** Consider a 1-way set associative mapping with 16 cache blocks, the more the memory block request are in the order (0, 255, 1, 1, 3, 8, 133, 159, 216, 219, 48, 32, 73, 92, 155) which one of the following memory block, will ~~access~~ be ✓ in the cache if LRU is used.
present

**A) 3** ~~B) 8~~ **C) 129** **D) 216.**

⟹.

(LRU – least recently used)



$S_0$ | ~~0, 1, 8, 216~~, ~~48, 32, 92~~
8, 48, 32, 92

$S_1$ | 1, 133, 73,

$S_2$ |

$S_3$ | (255, 3, 159), 219, ~~216~~ 63, 155

Set number = (i) mod 4

i → Block number.

Cache memory (4-way set associative)
(means each set contain 4 lines)

**Q:3** Consider a small 2-way set associative mapping with a total of 4 blocks, for choosing the block to be replace use LRU scheme, the number of cache misses for the following sequence of block addresses 8, 12, 0, 12, 8 is _4_ ?
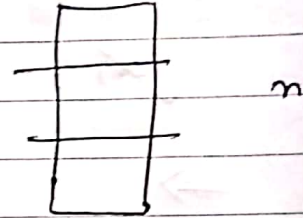M   M   M   H   M

⟹

$S_0$ | ~~8, 12~~, ~~0, 8~~ 12, 8

$S_1$ |

cache

s.no = (i) mod 2

✓ miss rate = $\left(\dfrac{4}{5} \times 100\right)$ = 80 %
✓ Hit rate = $\left(\dfrac{1}{5} \times 100\right)$ = 20 %

**Q:4** Consider a 2-way set associative mapping consisting of $2^s$ memory blocks and 2c cache blocks, the cache allocation for the memory block 'K' is ____

a) K mod 2c

b) K mod $2^c$

c) K mod c

d) K mod K

no. of line in cache = 2C

2-way set -associative.

no. of set = $\frac{2c}{2}$ = c

Set no = K mod c

**Q:5**

Consider the cache has 4 blocks, for the memory reference (5, 12, 13, 17, 4, 12, 13, 17, 2, 13, 19, 13, 43, 61, 19) what is the hit ratio for the following cache replacement algorithms —

(i) FIFO  (ii) LRU  (iii) Direct mapping.

(iv) 2-Way set Associate (LRU).

→ (i) FIFO.  (5, 12, 13, 17, 4, 12, 13, 17, 2, 13, 19, 13, 43, 61, 19)

| 5 | 43 |
|---|---|
| 12 | 61 |
| 13 | 19 |
| 17 | 13 |

hit ratio = $(\frac{5}{15} \times 100)$

miss ratio = $(\frac{10}{15} \times 100)$

→ (II) $\boxed{LRU}$. $(\overset{m}{5}, \overset{m}{12}, \overset{m}{13}, \overset{m}{17}, \overset{m}{4}, \overset{\checkmark}{12}, \overset{\checkmark}{13}, \overset{\checkmark}{17}, \overset{m}{2}, \overset{\checkmark}{13}, \overset{m}{19}, \overset{\checkmark}{13}, \overset{m}{43}, \overset{m}{61}, \overset{\checkmark}{19})$

$(5, 12, 13, 17, 4, 17, 13, 17, 4, 13, 19, 13, 43, 61, 19)$

hit ratio = $(\frac{6}{15} \times 100)$     miss ratio = $(\frac{9}{15} \times 100)$

(III)

→ $\boxed{\text{Direct mapping}}$ $(\overset{m}{5}, \overset{m}{12}, \overset{m}{13}, \overset{m}{17}, \overset{m}{4}, \overset{m}{12}, \overset{m}{13}, \overset{m}{17}, \overset{m}{2}, \overset{m}{13}, \overset{m}{19}, \overset{\checkmark}{13}, \overset{m}{43}, \overset{m}{61}, \overset{m}{19})$

| 0 | 12 | 4 12 2 |
| 1 | 5 | 13 17 13 17 13 61 |
| 2 | | |
| 3 | 19 | 43 19 |

$\frac{9(6)}{21} \Big| 15$
$\frac{\cdot 1}{\cdot 0}$

line no = (B.NO) mode 4

cache hit ratio = $(\frac{1}{15} \times 100)$

miss ratio = $(\frac{14}{15} \times 100)$

→ (IV) $\boxed{\text{2-Way set Associate (LRU)}}$ —

$(\overset{m}{5}, \overset{m}{12}, \overset{m}{13}, \overset{m}{17}, \overset{m}{4}, \overset{\checkmark}{12}, \overset{\checkmark}{13}, \overset{\checkmark}{17}, \overset{m}{2}, \overset{\checkmark}{13}, \overset{m}{19}, \overset{\checkmark}{13}, \overset{m}{43}, \overset{m}{61}, \overset{m}{19})$

| $S_0$ | 12, 4 | 12, 2 |
| $S_1$ | 5, 13 | 17, 13 17, 13, 19, 13, 43, 61, 19 |

line no = (i mod 2)

i → Block no.

hit ratio = $(\frac{5}{15} \times 100)$

miss ratio = $(\frac{10}{15} \times 100)$.

**Q. 6**

A hierarcal memory system has the following specification. 20MB main storage with access time of 300ns, 256B cache with access time of 50ns. Word size 4B, page size 8 words. What will be the hit ratio if the page address trace of a program has the pattern 0,1,2,3,0,1,3,0,12,4 following LRU page replacement technique.

→

Cache size = 256 B

Word size = 4B

page size = 8 words = (8×4)B = 32 B

$$No.\ of\ cache\ page = \frac{size\ of\ cache}{cache\ page\ size} = \frac{2^8}{2^5} = 2^3$$

$$(0,1,2,3,0,1,3,0,1,2,4)$$

$$hit\ ratio = \left(\frac{3}{8} \times 100\right)$$

**Q. 7**

consider an array A[100] and each element occupies 4 words, A 32-word cache is used and divided into 8-word blocks.

a) What is the hit ratio for the statement.
$$for\ (i=0,\ i<100,\ i++)$$
$$A[i] = A[i] + 10.$$

→ cache size = 32 W

Block size = 8 W

$$no.\ of\ Block\ in\ cache = \frac{32}{8} = 4\ Block.$$

| $B_0$ | $A_0$ | $A_1$ |
|---|---|---|
| $B_1$ | $A_2$ | $A_3$ |
| $B_2$ | $A_4$ | $A_5$ |
| $B_3$ | $A_6$ | $A_7$ |

cache.

$\dfrac{A_0}{m}, \dfrac{A_0}{H}, \dfrac{A_1}{H}, \dfrac{A_1}{H}$

$\dfrac{A_2}{m}, \dfrac{A_2}{H}, \dfrac{A_3}{H}, \dfrac{A_3}{H}$

hit ratio $= \left(\dfrac{3}{4} \times 100\right)$ ✓ $= 75\%$ hit rate.

**Q:8**

Considers an array has 100 elements and each elements occupies 4 words. A 32 bit word cache is used and divided into a blocks of 8 words. What is the hit rate of.

```
# for ( i = 0 ; i < 10 ; i++)
    for ( j = 0 ; j < 10 ; j ++)

    A [i] [j] = A [i] [j] + 10.
```

→ Cache size = 32 Word
Block size = 8 W
no. of Block = 4.

$i = 0$ to $10$ . . . . . . . . .
$i = 1$ : . . . . . . .
$j = 2$ . . . .

array $= [A_0, A_1, A_2, \ldots \ldots A_{100}]$

| $B_0$ | $A_0$ | $A_4$ |
|---|---|---|
| $B_1$ | | |
| $B_2$ | | |
| $B_3$ | | |

cache.

Rmo
$\boxed{00}$ $\boxed{01}$ $\boxed{02}$ $\boxed{03}$ $\boxed{04}$ $\boxed{05}$
Cmo
$\boxed{00}$ $\boxed{10}$ $\boxed{20}$ $\boxed{30}$

```
00 01 02 ...09
10 11 12 ...19
20
30
:
90 91 92 ...99
```

**RMO :**

| 00 | 01 |
|----|----|
|    |    |
|    |    |
|    |    |

| m | H | H | H |
|---|---|---|---|
| 00 | 00 | 01 | 01 |
| R | W | R | W |

$TM = (50)1$

$TH = (50)43$

$T \, Reference = (50)1$

$Hit \; ratio = \left(\dfrac{3}{4} \times 100\right)$

$= 75 \%$

**(Column major order)**

**CMO :**

| m | H | m | H |
|---|---|---|---|
| 00 | 00 | 01 | 01 |

$TM = (50)2$

$TH = (50)2$

$TR = (50)4$

$Hit \; ratio = \left(\dfrac{2}{4} \times 100\right)$

$= 50 \%$

- **Cache Coherence problem :**

→ cache coherence problem : Multiple copies of same data can exist in different caches simultaneously. and if processors are allow to Update their own copies freely an insonsistent view of memory can result.

$C_1 \to$ cache

$P_1 \to$ processor

MM (main memory)

$A = 5$

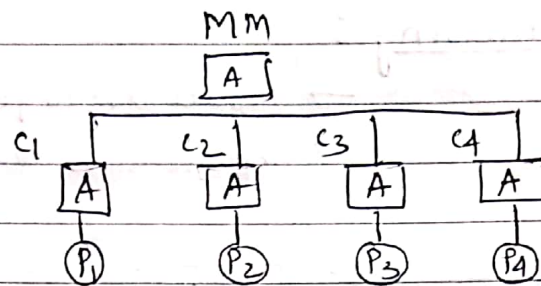| C₁ | C₂ | C₃ | C₄ |
|----|----|----|----|
| $A = A+1$ = 6 | $A = A-1$ = 4 | $A$ | $A$ |

$P_1$   $P_2$   $P_3$   $P_4$

- Methods to avoid cache coherence problem –

→
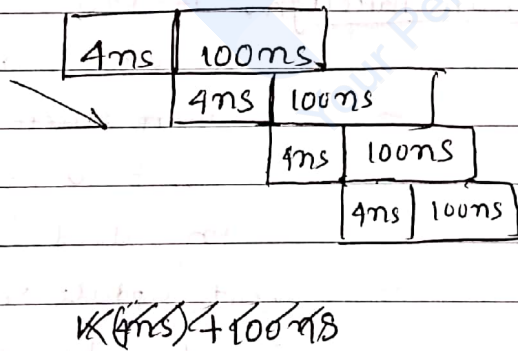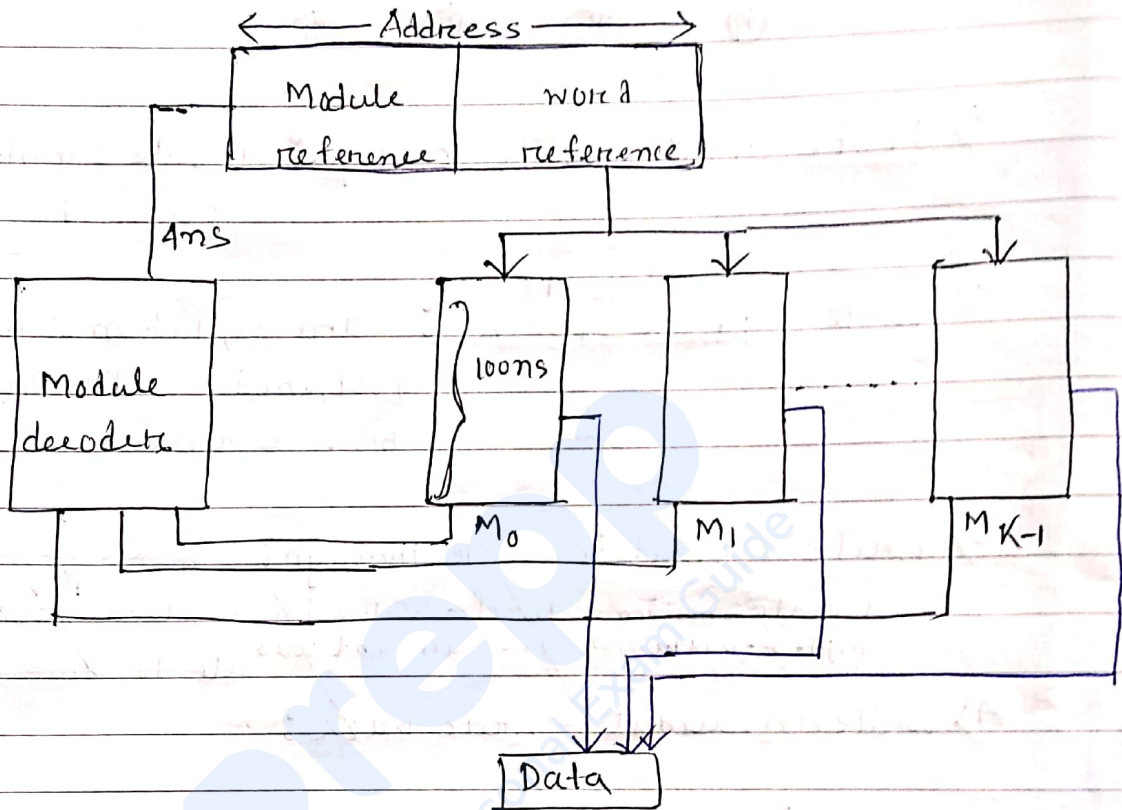
There are 4-method to avoid cache coherence problem

1) Write Update – write through.
2) Write Update – write back.
3) Write invalidate – write through.
4) Write invalidate – Write back.

MM



1) <u>Write update</u> - <u>write through</u>: update simultaniously of a word in cache & - M memory

2) <u>Write update</u> - <u>Write Back</u>: The updation of MM is postponed until the anouciated block is replaced.

3) <u>Write Invalidate</u> - <u>write through</u>: ~~changing a value in~~ update/simultaniously of a value in ~~one cache~~ this value will be invalidate ~~to other cache~~. MM. ~~and~~ and other cach can't used this value.

4) write Invalidate - write Back:

3) <u>write invalidate</u> - <u>write through</u>: if processor $P_1$ ~~was~~ used a value and change, then ~~other other~~ Processor $P_2$ ~~can't~~ the value will be invalied for other processor and update simultaniously of value in ~~mm~~ MM.

4) <u>Write Invalidate</u> - <u>write Back</u>: if $P_1$ ~~and~~ change a value then this value will be Invalid to other processor and updation of MM is postponed until the anouciated block is replaced.
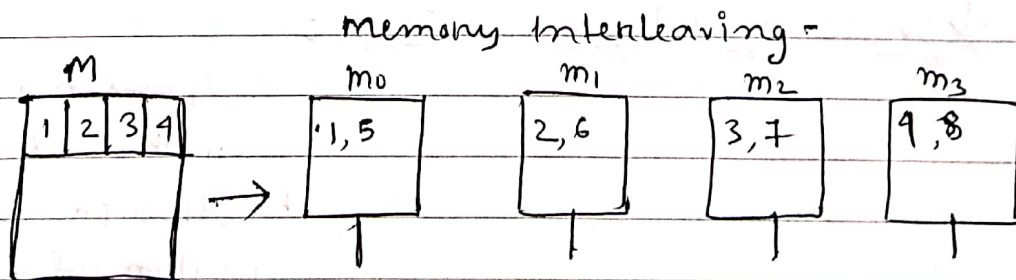
- **Memory interleaving :**
  → "Reduce average access time".
  → "Improve data transfer rate".



If read 'K' words then, total time taken,- $T = K(4ns) + 100ns$.

AND

without interleaving consept total time, $T = K(100ns)$.

$K(4ns) + 100ns$

**Memory interleaving -**

TO DOWNLOAD THE COMPLETE PDF

# CLICK ON THE LINK GIVEN BELOW

WWW.GATENOES.IN

## GATE CSE NOTES