

→ Set Theory started from
Back of Note book

Cof GATE

<u>Topic</u>	<u>Level</u>	<u>No. of revisions</u>
① Caching	Moderate	
② Memory Interfacing	Easy and Tricky	
③ Secondary Memory	Easy	
④ Machine Instruction & Addressing modes	Moderate	
⑤ ALU, Data Path & Control Unit	Moderate	
⑥ I/O Interface	Moderate	
⑦ Pipelining.	Moderate	

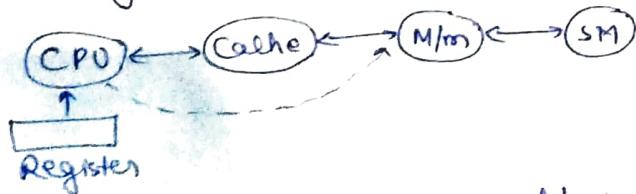
→ Do more and more questions on each topic from book.

Caching

(100% question in Gate)

Introduction to Cache Memory

- CPU is faster than main memory. So, if CPU tries to access data from main memory its speed gets reduced because main m/m can't reply that fast.
- So, we need some fast m/m. Register is the solution but its size is small and it's also too costly.
- Next solution is cache, it is slower than register but better than main m/m. Cache doesn't match the level with CPU in speed but it's better than main m/m.
- So, we have various levels of memories.



Main m/m contains a part of pages available in secondary m/m.
Cache contains some part of page.
Register contains some part of Cache.
So, CPU directly starts access from register and go forward to main m/m if required.

→ But we are considering registers here bcoz it's take negligible time.

→ If require element present in cache called 'Cache hit'.

And time required to fetch data from cache is 'hit latency'.

→ If data doesn't find in cache called as 'Cache miss' and time required to fetch data from main m/m called 'miss latency'.

→ If data doesn't present in main m/m called as page fault otherwise page hit.

→ Tag directory: It says the required page is present in tag or not.

→ In case of page fault we take page fault service means loading the page from Secondary memory.

→ So, there is need to thought that which part of main m/m needs to load in cache. For this we have principle of locality. It is of two types -

① Spatial Locality

→ It says that if you possess a cache miss then its chances to find the element in proximity is very high. It means you can find the element around the previous miss.

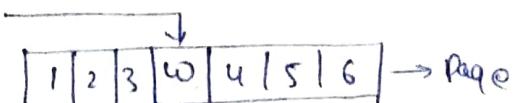
② Temporal Locality

→ It says that least recently used element is going to use again. So, we use LRU generally in page replacement algorithm.

→ Explanation of spatial locality.

Cont-

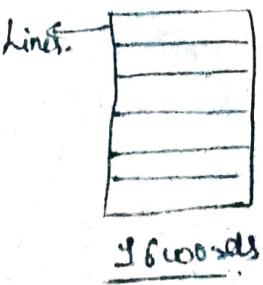
Say, we face a page fault because we don't find a word in main m/m. so we load the complete page from the secondary m/m



Say, W is the word that causes page fault. We load the complete page instead of loading only a single word bcoz of spatial locality means it's highly probable that other 6 words can be demanded soon.

Introduction to Direct mapping

- We are talking about main m/m and Cache mapping.
- Paging is just a concept to manage the m/m more efficiently.
- In OS we divide main m/m into → frame but here we divide it into blocks and Cache into lines i.e. block size is same as line size.
- we assume 1 word = 1 Byte.



0	0	1	2	3	blocks
1	4	5	6	7	
2	8	9	10	11	
3	12	13	14	15	
4					
5					
6					
7					
8					
9					
10					
11					
12					
13					
14					
15	60	61	62	63	64 words

Say block size = 4 words

$$\begin{aligned} \text{No. of blocks} &= \frac{64}{4} \\ &= 16 \text{ blocks.} \end{aligned}$$

→ 16 blocks with 4 words in each block.

$$\text{No. of lines} = \frac{16}{4} = 4 \text{ lines}$$

4 lines with 4 words in each line.

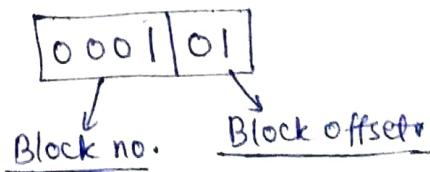
→ for 64 words we need 6 bit to represent.

→ So, CPU will generate 6-bit address.

Say CPU generates address as

$$[\underline{\underline{0}} \underline{\underline{0}} \underline{\underline{1}} \underline{\underline{0}} \underline{\underline{1}}]_2 \rightarrow (5)_{10}$$

→ Since block size = 4 words
So we divide the least significant (2^2) 2 bits as block offset.



So, it is task that which block from main m/m should be loaded in cache. For this we have various kind of mapping techniques.

① Direct mapping:

Cache	→ So main m/m to cache is a many-one kind of function.
0 0 4 8 12	
1 1 5 9 13	
2 2 6 10 14	
3 3 7 11 15	

Cache contains block numbers.

We have 36 blocks so we need 4 bits to represent each block.

	offset			Block no
0 -	0	0	0	0
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

for line '0' in cache we choose those block that have least two digits as 00.

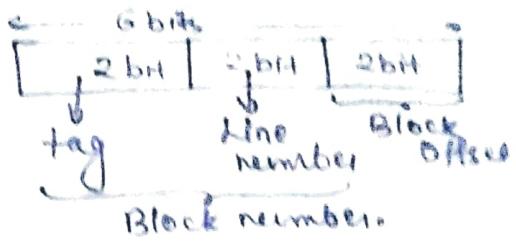
Similarly for various lines the last 2 least significant digits are

Line 1 \rightarrow 01

Line 2 \rightarrow 10

Line 3 \rightarrow 11

\rightarrow we take two bits only to take line no.



Ex:

Tags	Line number	Block offset	0, 1, 2, 3 are line numbers and each line contain 4 words.
100	0	0	
101	1	0	
110	2	0	
111	3	0	

In block offset we are free to put any of the four combinations.

①

0x4 = 0 to 1111 + 3 (0, 1, 2, 3)		
2 bits	2 bits	2 bits
00	00	00
00	00	01
00	00	10
00	00	11
01	01	00
01	01	01
01	01	10
01	01	11

②

5x4 = 20 to 1111 + 3 (20, 21, 22, 23)		
2 bits	2 bits	2 bits
01	01	00
01	01	01
01	01	10
01	01	11
01	01	11

③

10x4 = 40 to 1111 + 3 + 4 (41, 42, 43, 44)		
2 bits	2 bits	2 bits
10	10	00
10	10	01
10	10	10
10	10	11
10	10	11

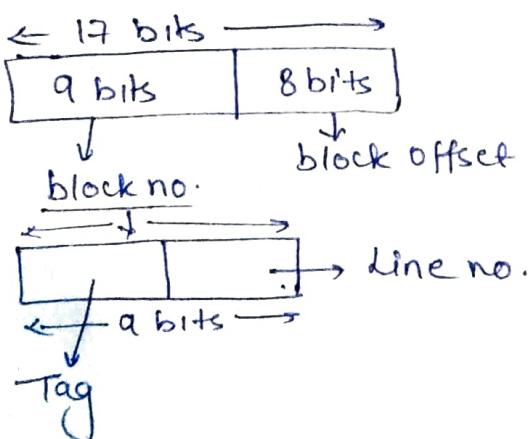
\rightarrow we give all four different values to tag but we can give same value to tag if we have sufficient no. of blocks to represent.

Direct mapping Problems

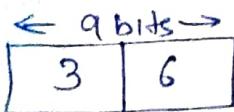
Problem 1 (Byte Addressable)

MM Size	Cache Size	Block size	Tag Bits	Tag Directory size
① 128 kB	16 kB	256 B	3	<u>3x2⁶</u>
② 32 GB	32 kB	1 kB	20	<u>20x2⁵</u>
③ 64 MB	512 kB	1 kB	7	<u>7x2⁹</u>
④ 16 GB	16 MB	4 kB	10	<u>10x2¹²</u>
⑤ 64 MB	<u>2¹⁶ B</u>	Can't guess	10	Can't guess
⑥ 2 ²⁶ B	512 kB	Can't guess	7	Can't guess

Solⁿ ① MM size = 2^{17} B
 Cache — = 2^{14} B
 Block — = 2^8 B
 No. of blocks = $\frac{2^{17}}{2^8} = 2^9$ blocks



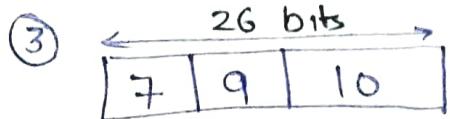
No. of lines = $\frac{2^{14}}{2^8} = 2^6$



2⁶ blocks of main memory mapped with cache.

~~Tag directory size~~
 $= (\text{Tag bits}) \times \text{no. of lines}$

② Solution is same as ①

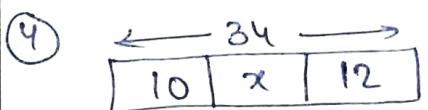


Block size = 2^{10} B

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

No. of lines = 2^9 (9 bits)

MM size = 2^{26} B
 $= 64 \text{ MB}$



No. of blocks = $\frac{2^4 \times 2^{30}}{2^{12}}$
 $= 2^{34-12} = 2^{22}$

16 GB = 2^{34} B

So, $10 + x + 12 = 34$

$x = 12$

Cache size 2^{24} B = 16 MB

No. of lines = $\frac{2^{24}}{2^{12}} = 2^{12}$

⑤ we can't find block size because we can guess many combination.

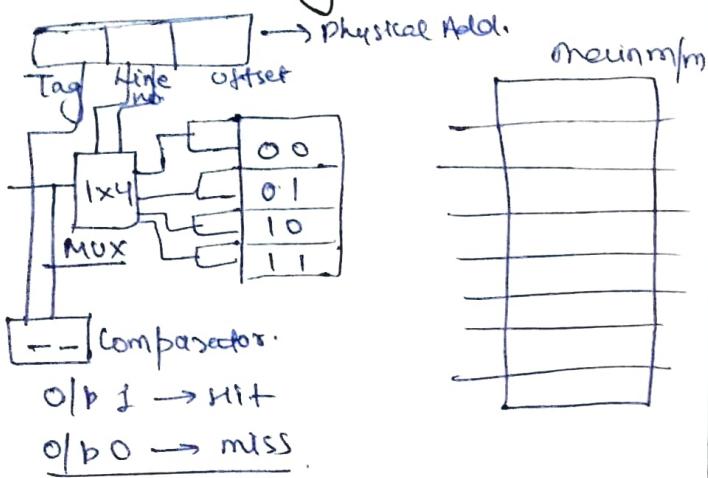
⑥ Cache size = 2^{19} B

$(\text{line no.}) \times (\text{line offset}) = 2^{19}$ B.

Ratio between main memory and Cache is tag size.

→ Tag directory size is mostly asked question.

Direct mapping H/W Implementation



We choose the line no and gets its tag bit as o/p through MUX. If MUX's o/p and tag bit in address is same then we get '1' else '0'.

Time taken in above procedure is hit latency.

As the tag bits increases no. of MUX gets increases.

All the MUX will get the same line no.

For k -bit tag we need k Multiplexers but only one Comparator (k -bits).

Since all the MUX are working in parallel so at a time only one MUX and only one comparator should work.

~~Hit latency = Comparator delay + MUX delay~~

Generally, $MUX \text{ delay} \ll \text{Comparator delay}$

So we can neglect it.

Ex: Main m/m size = 1GB
Cache — = 1MB
Comparator delay = 10 ns
Hit latency = ?

$$\text{Soln} \quad \text{Tags} = 2^{10} \quad (\frac{10 \text{ bits}}{\text{Line}})$$

$$\text{Delay} = 10 \times 10 = 100 \text{ ns}$$

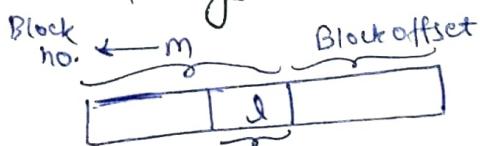
Hit latency = Delay (here).

Note :-

$$k = m \% n$$

m th block of main m/m present in k th line of Cache, where n is no. of lines in cache.

Disadvantage of direct mapping



$$\text{No. of blocks} = 2^m$$

$$\text{No. of lines} = 2^l$$

Say x is any m bits no. then

$$x \% 2^l = \text{line no.}$$

→ It gives line no. in which x is present

Say a cache have 4 lines and CPU generate block request as 5, 6, 4, 8, 9, 12, 15, 20. then find the line no. in which block get placed and also find the miss rate.

$12 \% / 4 = 0$	$5 \% / 4 = 1$
$15 \% / 4 = 0$	$6 \% / 4 = 2$
$20 \% / 4 = 0$	$4 \% / 4 = 0$
0 1 2 3	8 9 6 15

$$\begin{aligned} 5 \% / 4 &= 1 \\ 6 \% / 4 &= 2 \\ 4 \% / 4 &= 0 \\ 8 \% / 4 &= 0 \\ 9 \% / 4 &= 1 \end{aligned}$$

After filling the block no. 4 in line no. 0, we get the request of block no. 8 and then we replace the block no. 4 with block no. 8 even the line 3 was empty, but we can't use it. It is the drawback of direct mapping.

In worst case input should be like

4, 8, 12, 16, 20, 24 - -

0	4	8	12	16	20	24
1						
2						
3						

(Conflict miss)

Even we have plenty of space available but we can't use because of many-one mapping.

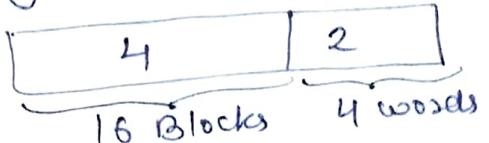
We are missing the block due to conflict on same line with another block.

→ Particular block can go in particular line.

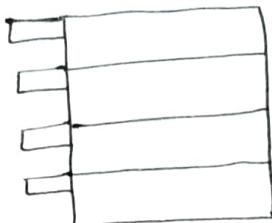
Associative mapping

- Solution to direct mapping.
- A block can get inside any line under some conditions.

Say we have 6 bit PA -

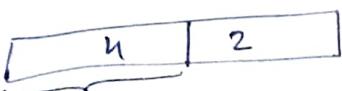


All 16 blocks can map to any line of cache, so it's many-many function.



→ Tag says that how many blocks can reside in a single line.

→ Here we can allow all of 16 blocks to use any of the line. So tag should be of 4 bits for 16 blocks.

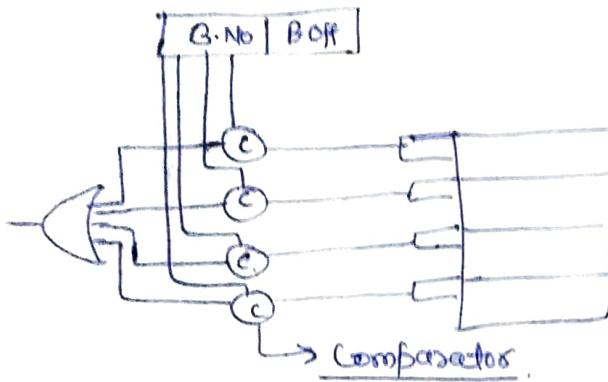


Block no = tag (here)

→ There is 'no conflict miss'. Conflict miss. block 4 already line में था but 8 को आने से उसे replace किया तो अब आपको लाने की block 4 की demand होगी तो यह miss कर देगा। परंतु यहाँ पर्सनल में पर्सनल की ओर से एक block की ओर से संकेत है।

→ Where should be search for a block?

→ So here we have many comparators connected to each tag.



So, we get the solⁿ of conflict miss but it increases the hardware cost bcoz of many comparators.

Hit latency → Time taken in comparator + Time taken by (OR) Gate

→ We can replace OR Gate with any other combinational Ckt if required.

Ex → MM = 32 GB, Block = 32 kB

Tag = ?

Propagation delay of Comparator = 10 ns

Prop. delay of OR gate = 1 ns

Hit latency = ?

(k = no. of tag bits)

Solⁿ Tag = $\frac{2^{35}}{2^{15}} = 2^{20}$ = 20 bits
k = 20

$$\text{Hit latency} = (10 \times 20 + 10) \text{ ns} \\ = 210 \text{ ns}$$

→ Tag = no. of blocks.

→ All the comparators are cost in parallel. So we consider it as only one comparator time.

→ Assume that all the lines are full and we want to place a new block then which block from which line will get replaced. For this decision we have block replacement algorithms.

Ex → (Byte Addressable)

	MM Size	Cache Size	Block Size	Tag size	Tag directory size
①	128 kB	16 kB	256 B	9	<u>$2^6 \times 9$</u>
②	32 GB	32 kB	1 kB	25	<u>$2^5 \times 25$</u>
③	<u>2^{17} B</u>	512 kB	1 kB	7	<u>$2^7 \times 2^9$</u>
④	16 GB	Count guess	4 kB	22	<u>$10 \times \text{Count guess}$</u>
⑤	64 MB	Count guess	Count guess	10	<u>$10 \times \text{Count guess}$</u>
⑥	Count guess	512 kB	Count guess	7	<u>$7 \times \text{Count guess}$</u>

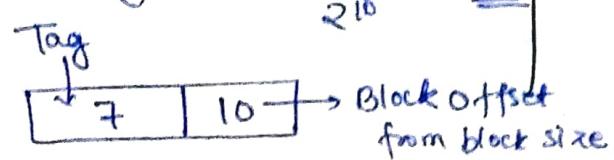
$$\text{① Tag size} = \frac{2^{17}}{2^8} = 9 \text{ bits}$$

$$\text{T. Directory size} = \frac{2^{14}}{2^8} \times 9 = 2^6 \times 9$$

$$\text{② Tag size} = \frac{2^{35}}{2^{10}} = 25 \text{ bits}$$

$$\text{No. of lines} = \frac{2^{15}}{2^{10}} = 2^5$$

$$\text{③ No. of lines} = \frac{2^{19}}{2^{10}} = 2^9$$



Similarly we can do (4), (5) and

(6) question option.

~~gmp~~
→ No. of comparators required is equals to no. of lines in cache.

→ In (3) option size of main m/m is less than cache size. It doesn't make sense. So the question is logically not correctly framed.

✓ To reduce no. of comparators we have set associative method because associative method was taking huge no. of comparators.

We divide cache into sets and we allow a block to sit in particular set but in a set it can sit anywhere. It reduces the comparison time.

Set associative mapping

Ex → MM size = 64 B

Cache $\frac{1}{1}$ = 32 B

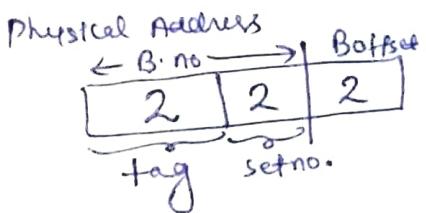
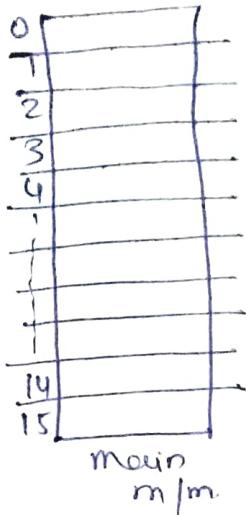
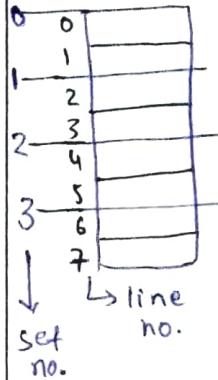
Block $\frac{1}{1}$ = 4 B

Set $\frac{1}{1}$ = 2 blocks/lines

(2-way set associative bcoz 2 blocks are here).

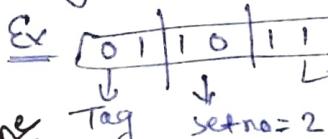
$$\text{No. of lines} = \frac{32B}{4B} = 8 \text{ lines}$$

$$\text{No. of sets} = \frac{8}{2} = 4 \text{ sets.}$$



Set no. = 2 bits bcoz of 4 sets

Tag = 2 bits.

Ex  offset = 3
set no. = 2

Tag = 2 bits bcoz $(6-4) = 2$
and PA = $2^6 = 64$ block in m/m.

On a K-way set associative we need K comparators of size of tag.

→ Here we require two comparators of 2 bit size.

Open
for revision

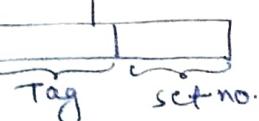
Problems on set associative mapping

MM Size	Cache Size	Block Size	Tag bits	Set Associative (Given)
128KB	16KB	256B	4	2-way
324B	32KB	1KB	2	4-way
8MB	512KB	1KB	7	8-way
16GB	64MB	4KB	10	4-way
64MB	2 ¹⁸ B	count guess	10	4-way
8MB	512KB	count guess	7	8-way

Soln

① Block offset = 8 bits

$$\text{Main m/m} = 2^{17} \text{ B}$$

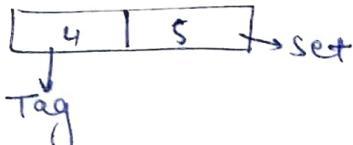


$$\text{No. of lines} = \frac{2^{14}}{2^8} = 2^6$$

In two way set associative

1 set contains 2 lines

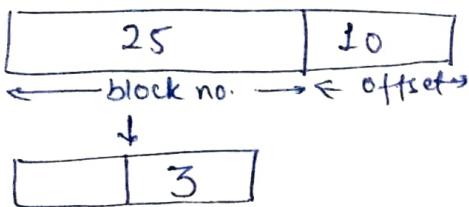
$$\text{No. of sets} = \frac{2^6}{2} = 2^5 \quad (5\text{-bit})$$



No. of comparators required = 2.

② Block offset = 10 bits (2^{10} B)

$$32 \text{ GB} = 2^{35} \text{ B} \quad \text{PA}$$



$$\text{No. of lines} = \frac{2^{15}}{2^{10}} = 2^5$$

In k-way set associative

$$\text{No. of set bits} = \frac{25}{2^2} = 2^3 \quad (3\text{ bits})$$

$$\text{Tag bits} = 25 - 3 = 22$$

No. of comparators = 4.

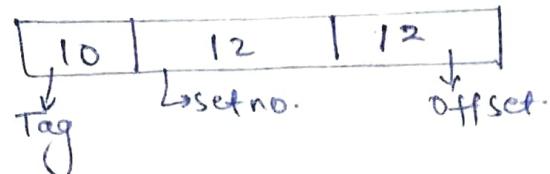
(3) same as above ② and ①

Ans Tag bits = 7

$$\text{Main m/m} = 8 \text{ MB} (2^{23} \text{ B})$$

~~2¹⁴~~

④ PA = 34 bits



$$\text{No. of sets} = 2^{12} = \frac{\text{No. of lines}}{2^2}$$

$$\text{No. of lines} = 2^{14}$$

$$\begin{aligned} \text{Cache size} &= \text{No. of lines} \times \text{line size} \\ &= 2^{14} \times 2^{12} \text{ B} \\ &= 2^{26} \text{ B. (64 MB)} \end{aligned}$$

~~PA = 2²⁶ B, Tag bits = 10~~

~~wrong~~



Cache size depends on line size
that is not given.

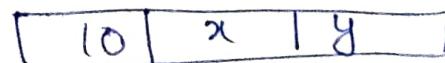
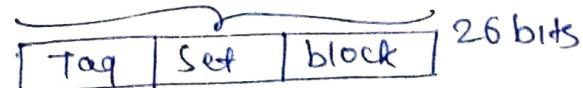
→ So, we can't guess bcoz there
can be multiple combinations
possible

~~PA~~



→ So we can't guess anything.

⑤



$$2^x = \text{no. of sets}, 2^y = \text{block size}$$

$$x+y = 26 - 10 = 16$$

$$\begin{aligned}
 \text{Cache size} &= \text{No. of sets} * \frac{\text{No. of lines}}{\text{per set}} + \text{line size} \\
 &= 2^x + 2^2 * 2^y \\
 &\quad \downarrow \\
 &\quad (\text{4-way}) \\
 &= 2^{x+y+2} = 4 \cdot 2^{(x+y)} \\
 &= 4 \cdot 2^{16} = \underline{2^{18} \text{ B}}
 \end{aligned}$$

→ But we can't find block size.

⑥ Similar approach as que. ⑤
we get

$$\underline{\text{MM size}} = \underline{8 \text{ MB}}$$

Gate 2012

Given, Cache size = 256 kB

4-way set associative used

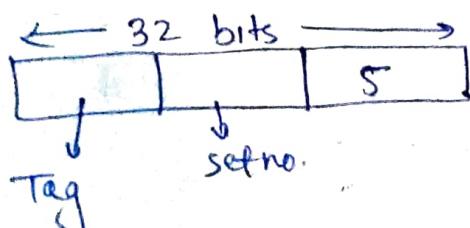
Block size = 32 B, PA = 32 bits

-) Each cache tag directory entry contains in addition to address tag, 2 valid bits, 1 modified bit and 1 replacement bit

To find:

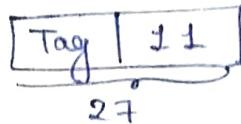
- ① No. of tag bits.
- ② Size of cache tag directory.

Soln

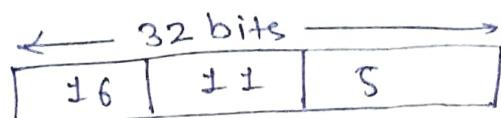


$$\text{No. of lines} = \frac{2^{18}}{2^5} = 2^{13}$$

$$\text{No. of sets} = \frac{2^{13}}{4} = 2^{11}$$



Tag = 16 bits.



Size of tag directory

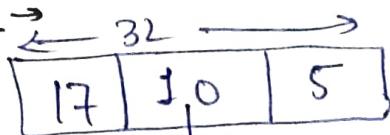
$$\begin{aligned}
 &= (16 + 2 + 1) * \text{no. of lines} \\
 &= 20 * 2^{13} \\
 &= \underline{160 \text{ kB}}
 \end{aligned}$$

Gate-2005 (Direct mapping)

Given, Cache size = 32 kB
Block — = 32 B
PA = 32 bits.

Find no. of bits of cache indexing and tag.

Soln →



cache / lineIndex or line no.

$$\text{No. of lines} = \frac{32 \text{ kB}}{32} = 2^{10}$$

$$\text{tag bits} = 32 - 15$$

= 17. Ans!!

Comparing all the mappings

$$\text{Ex} \rightarrow \text{MM} = 4\text{GB} = 2^{32}\text{B}$$

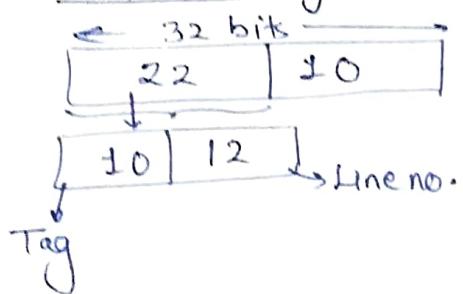
$$\text{Cache} = 1\text{MB} = 2^{22}\text{B}$$

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

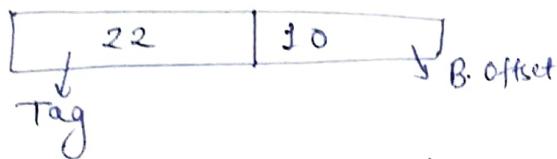
$$PA = 32 \text{ bits}$$

$$\text{No. of lines} = \frac{2^{22}}{2^{10}} = 2^{12}$$

① Direct mapping

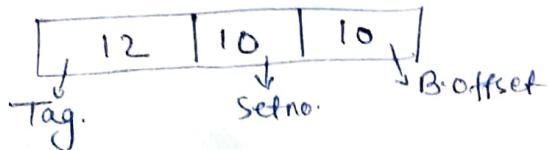


② Associative mapping



③ 4-way set associative

$$\text{No. of sets} = \frac{2^{12}}{4} = 2^{10}$$



~~✓~~ Direct mapping is k-way set associative with

$$k=1$$

~~✓~~ K-way set associative is fully associative with $k=N$ where, N is no. of lines.

→ k-way is a balance in between other two mapping.

	TAG	Comparators	Size of tagfield
Direct	10	1	10×2^1
Associate	22	2^{12}	22×2^1
k-way	12	4	12×2^1

Gate Questions

Q9) Main m/m = 2 CM Blocks

Cache = 2C Blocks

2-way set associativity.

$$\text{Soln} \quad \text{No. of sets} = \frac{2C}{2} = C \text{ sets}$$

k-th block will map to $k \bmod C$.

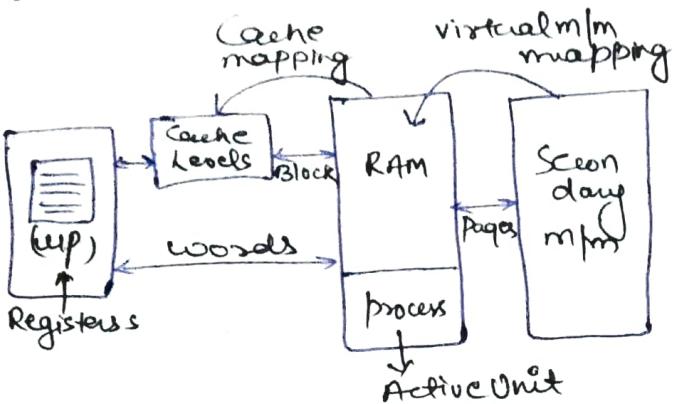
~~V.9m~~
~~Q.10~~
Note :- There are some videos with previous yr. gate. Do watch all the videos during solving practice questions bcoz we will definitely get a better solution.

→ Must watch all the videos.

→ Do some questions from standard books (10-15 Ques)

Memory Interfacing

Introduction



The purpose of main m/m is to bridge the speed mismatch between fastest processor to slow memory at reasonable costs.

→ SRAM faster than DRAM.

2-Level Memory

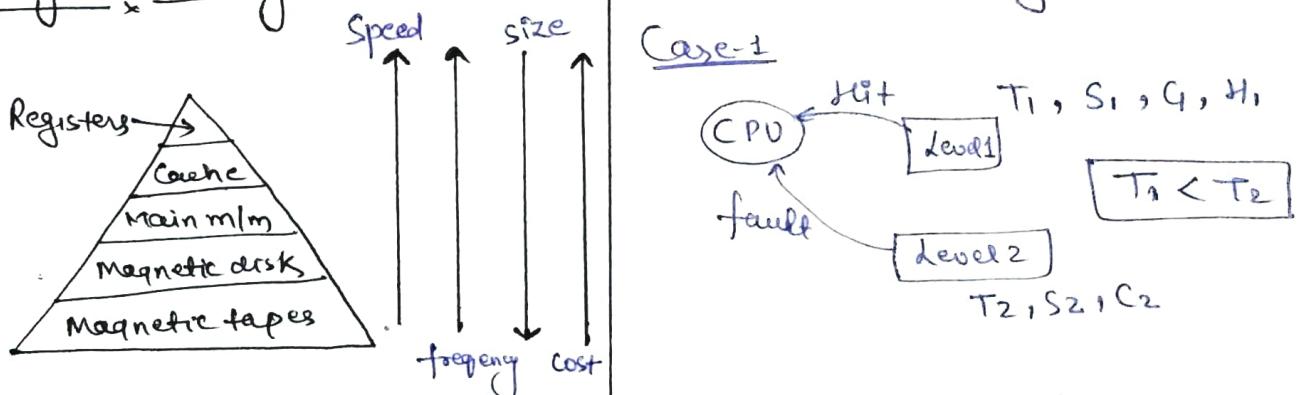
Information in

ith level C (i+1)th level

→ If processor refers to i^{th} level m/m is found then 'Hit' otherwise 'Miss' or 'Fault'.

→ Far There are two ways in which processor is connected to various levels of memory.

Case-1

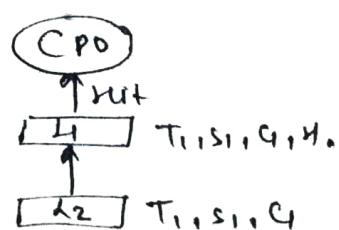


$T \rightarrow$ access time, $S \rightarrow$ size
 $C \rightarrow$ cost per hit $\propto \rightarrow$ hit rate

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

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

Case-2



Registers implemented with flip-flops
Cache implemented with static RAM. Main memory implemented with Dynamic RAM.

Cache
Main memory } Random Access
Magnetic Disk } Semi Random Access
Magnetic Tapes } Sequential access

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

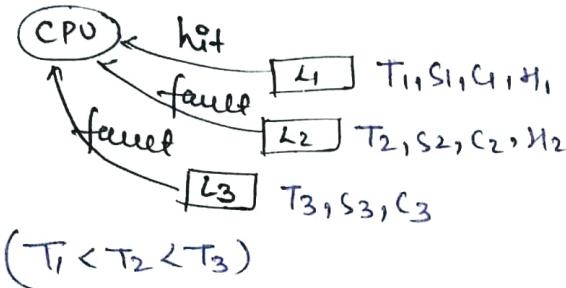
$$C_{avg} = \frac{C_1 S_1 + C_2 S_2}{S_1 + S_2}$$

→ In case of fault we have to go through both the levels.

→ Cost (C) measures in R_s/B , (Rupees per Byte) and size (S) in Bytes.

3-Level Memory

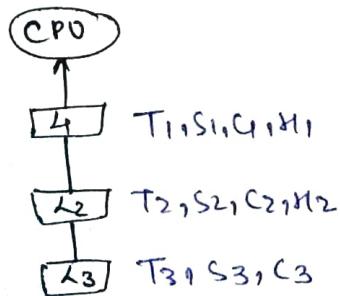
Case-1



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

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

Case-2



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

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

In the case 1:

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

→ In both the cases H_3 for L_3 is always 1, in 3-level memory otherwise instruction will never be found.

→ In case 2:

$$T_1 \leq T_{avg} \leq (T_1+T_2+T_3)$$

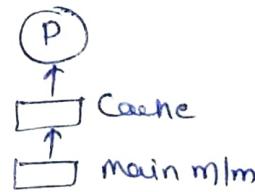
Ex: Cache and main m/m access time = 5 ns 800 ns respectively
Cache hit rate = 0.8
Find avg. m/m access time?

Soln Case-1



$$\begin{aligned} T_{avg} &= 0.8 \times 5 + 0.2 \times 800 \\ &= 4 + 20 \\ &= 24 \text{ ns} \end{aligned}$$

Case-2



$$\begin{aligned} T_{avg} &= 0.8 \times 5 + 0.2 (10^5) \\ &= (4 + 21) \text{ ns} \\ &= 25 \text{ ns} \end{aligned}$$

Cache Replacement Algo.

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

→ Replacement policy are aimed to minimize miss penalty references.

Replacement Policy

- Random: No specific criteria to replace the block.
- FIFO: The block which entered first is the candidate for replacement.
- LRU: The block which is not referred for longest time will get replaced.
- LFU: The block with fewer references is replaced.

Ex: Direct mapped cache with 8 blocks (0-7). Blocks request in order,

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 block will not be in cache at the end of sequence.

- a) 3 b) 18 c) 20 d) 30

Solⁿ

	0	1	2	3	4	5	6	7
0	Ø Ø 18(24)							
1	Ø Ø 17	17						
2	2 18 2	82						
3	3 2							
4	20							
5	5							
6	6 30							
7	(63)							

block no. % 8
= line no.

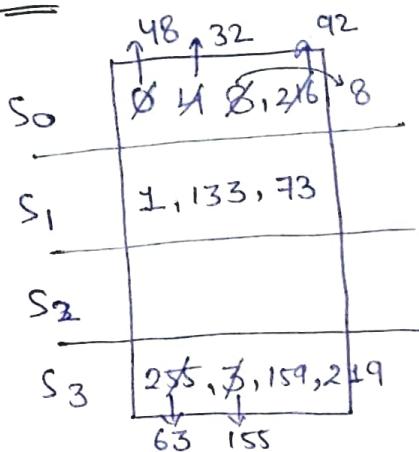
→ 48 is not there in cache at the end of sequence

Ex → 4-way set associative mapping with 16 cache block. Blocks are requested in order (0, 255, 1, 4, 3, 8, 133, 159, 216, 219, 63, 8, 48, 32, 73, 92, 155).

Which one of the main block will not be in cache if LRU used?

- a) 3 b) 8 c) 129 d) 216

Solⁿ



$$Sno = (\text{block no}) \bmod (\text{no. of sets})$$

↳ Set no.

$$\text{Set no} = \text{block no. mod } 4$$

Ans → a), c) and d)

→ when 8 comes 2nd time in S0, we write it again bcoz we were following LRU. So, 2nd time 8 can't get replaced by 92.

Ex → Consider a small 2-way set associative mapping with total 4 blocks for choosing the block to be replaced w.r.t LRU. The number of cache misses for the sequence of block address 8, 12, 0, 12, 8 is _____?

$$8, 12, 0, 12, 8 \text{ is } \underline{\hspace{2cm}} ?$$

Solⁿ

	8
S ₀	8, 12
S ₁	

Set no = $k \% 2$
→ block no.

8	12	0	12	8
↓	↓	↓	↓	↓
miss	Miss	miss	hit	miss

There are '4' misses.

→ In the beginning there were nothing in cache. So loading of first 8 and 12 was the miss bcz we load it from main m/m.

Q. Consider a 2-way set associative mapping consistency of 2^c memory blocks and 2^c cache blocks then, cache location for the memory block k is.

Solⁿ Cache size = 2^c blocks
MM — = 2^c blocks
Set size = 2 blocks.

Set no = cache location = $k \bmod c$

Sets = $\frac{2^c}{2} = c$

Ans → $k \bmod c$.

Q → 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($i=0$; $i<100$; $i+1$)
 $A[i] = A[i]+10$

Solⁿ Cache size = 32 words
block — = 8 words,
No. of blocks = 4 blocks

B ₀		
B ₁		
B ₂		
B ₃		

Cache have 4 blocks of size 8 words each. Each element take 4 words. so each block can hold 2 array elements.

In the statement

$$A[i] = A[i] + 2$$

↓ \rightarrow for reading
for writing(hit) (miss)

$A[i]$ referred two times first for reading and second for writing.

Since cache is empty so when element choose for reading its a miss and then when it refers for writing its a hit bcz after a miss it will get loaded to cache.

→ Not complete solution bcz of doubt.

→ watch the video in downloaded video as 'Array and cache ex-1'

So on calling $A[0]$ it will fill two elements in block as

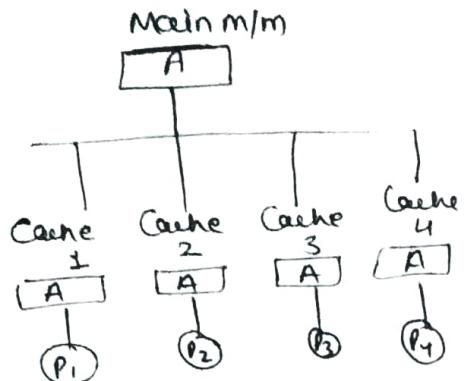
$A[0] | A[1]$ and since both elements occur two times so it will occur as

$A_0 A_0 A_1 A_1$ because as A_0

Miss Hit comes all the 3 other will get 9 hit. similarly 50 pairs formed by 100 elements.

Cache Coherence Problem

Multiple copies of same data can exist in different caches at the same time, and if processors are allowed to update their own copies freely an inconsistent view of m/m can result.



Each processor is running a job and all the jobs are accessing the common variable (A) at the same time. It may arise inconsistency while updating data.

Methods to avoid Cache Coherence Problems

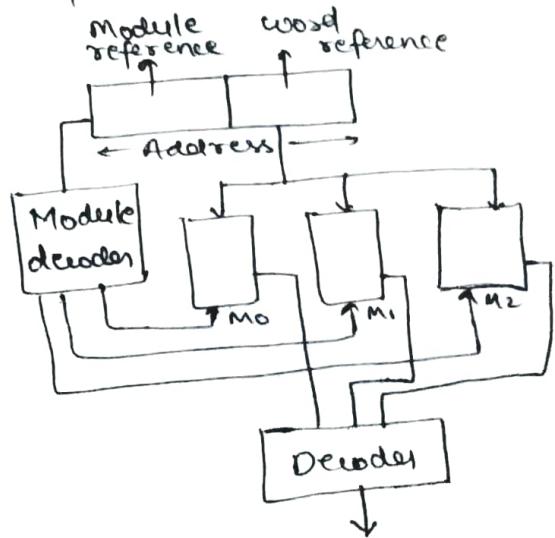
- ① write update - write through
- ② write update - write back
- ③ write invalidate - write through
- ④ write invalidate - write back

→ write through means update the value in main memory just after updating value in cache
→ write back means update the value in main memory at the end when you either replace the block or flush out the memory.

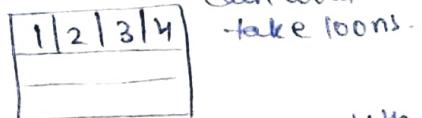
- write update means update the value of variable in all the cache at the same time
- write Invalidate means after update the value tells the other processors that the variable copy is invalidate

Memory Interleaving

→ Reduces avg access time and improves data transfer rate.



Consider the cache that have 4 words like, and searching each word take loops.

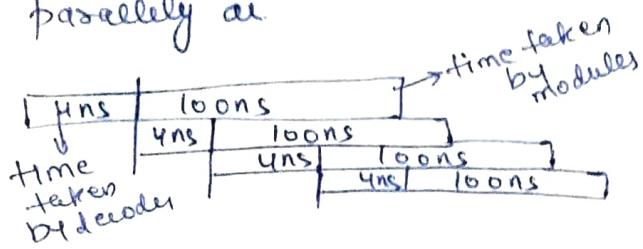


Assume you want to access 4th words then first you hit 1 & it take loops. similarly upto 4 you will hit 2, 3 then 4 it will take $4 \times 100\text{ ns} = 400\text{ ns}$.

But if all modules are different as:



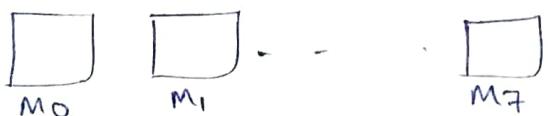
then, some time will take for modular decoder and some time will take by modules but everything will occur parallelly as



$$(4 \times 4) + 100\text{ ns} = 116\text{ ns only}$$

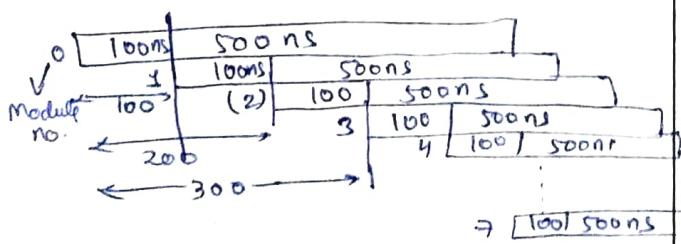
~~Ques~~ Gate 2014 COA (3.48 PYQ)

Soln → Total 8 modules



time taken by modules = 500ns

time taken by decoder = 100ns



We are generating request for next modules after 100ns so,
in 100ns → 1 request

$$1 \text{ s} = \frac{1}{100 \times 10^{-9}}$$

$$1 \text{ ms} = \frac{10^{-3}}{100 \times 10^{-9}} = \frac{10^{-3}}{10^{-7}} \\ = 10000 \text{ request}$$

→ Solution of Gate 2006 is given but you can't understand anything from it. So try to do it by self.

Gate 2006 : CPU has a cache of 64 bytes and main mem have k banks [See ques. in PYQ]

Soln



Each of size 2 bytes. All bytes are allotted sequentially. Total 48 bytes. All blocks accessed parallelly.

→ 8ns to read one block.

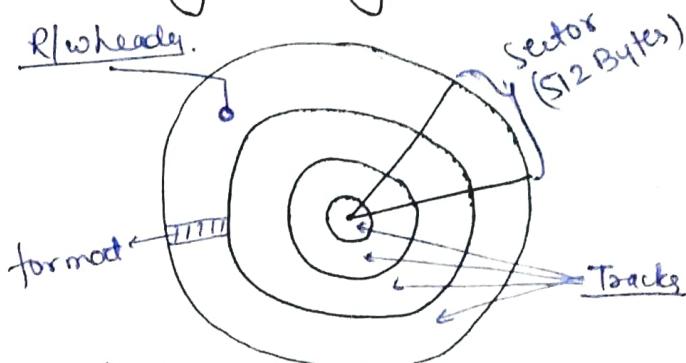
→ Decoder takes $\frac{k}{2}$ ns (12ns).

Total time = $12 + 60 = 92$
because decoder activates all the block at a time

$$\text{Ans } 2 \times 92 = 184$$

Secondary Memory

R/W head.



Standard size of sector is 512 Bytes.

Total capacity of hard disk equals to

$$(\text{No. of tracks}) \times (\text{No. of sectors}) \times \frac{\text{Bytes per track}}{\text{No. of sectors}} \\ = \text{Bytes}$$

Some space get used in formatting so, we can't use the complete space for storage purpose

$$\text{Usable storage} = (\text{Total size}) - (\text{format size})$$

Time taken by R/w head to move from one track to another track is called 'seek time'

Time taken by head to reach particular sector in the track called 'rotational delay'.

$$(\text{Total}) \text{ access} = T_{\text{seek}} + T_{\text{rd}} + T_{\text{trans}}$$

T_{trans} → time spent in target sector
if T_{seek} not given assume as '0',
for T_{rd} use ($\frac{1}{2}$ Rotation time)

Example

① A disk pack have 16 surfaces, 128 tracks / surface, 256 sectors per track and 512 Bytes / sector.

a) find Capacity of disk pack?

Each disk have two surfaces and we can store data on both surfaces.

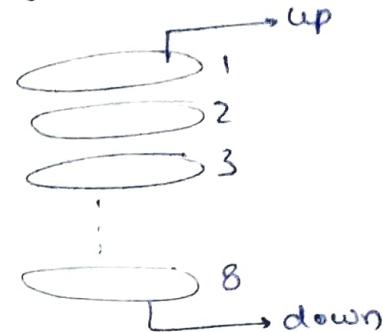
$$16 \times 128 \times \frac{1}{\text{sec}} \times \frac{256}{\text{sec}} \times \frac{512 \text{ B}}{\text{sec}} \\ \Rightarrow 16 \times 128 \times 256 \times 512 \\ = 2^{28} = \underline{256 \text{ MB}}$$

b) Bits required to address the sectors.

$$\rightarrow \text{No. of sectors} = 2^4 \times 2^7 \times 2^8 \\ = 2^{19}$$

So, 19 bits are required.

Note → In previous question we have 8 disk so, that we can get 16 surfaces as



In case if question says that we are not considering the boundary surfaces means (up of disk 1) and (down of disk 8) then we will add a new disk to introduce two more surfaces. In that case we get 16 surfaces from 9 disks.

c) If format overhead is 32 Bytes / sector. What is the formatted disk space

$$\text{Wasted space} = 2^{19} \times 2^5 \\ = 2^{24} \text{ Bytes} \\ = \underline{16 \text{ MB}}$$

formatted disk space

$$= (256 - 16) \text{ MB} \\ = \underline{240 \text{ MB}}$$

d) Det the diameter of innermost track is 21 cm. What is the max. recording density.

→ Each track have same no. of sectors and each sector have same size. So, outer track have low recording density bcoz

m/m divided in high area as
Compare to innermost track

Max recording density → inner most back

Min. recording density → outer most back

$$\begin{aligned}\text{Track capacity} &= 256 \times 512 \\ &= 2^8 \times 2^9 \\ &= \underline{2^{17} \text{ Bytes}}.\end{aligned}$$

$$\begin{aligned}\text{Circumference} &= 2\pi r = \pi d \\ &= \frac{22}{7} \times 21.3 \\ &= \underline{66 \text{ cm}}.\end{aligned}$$

$$\begin{aligned}\text{Max density} &= \frac{128 \text{ KB}}{66} \\ &= \underline{1.9 \text{ KB/cm}}.\end{aligned}$$

(e) The disk is rotating at 3600 Rpm. What is the data transfer rate?

for each surface we can read one track and for 16 parallel surfaces we can read 16 surfaces at a time.

$$\rightarrow 3600 \text{ Rotation} \rightarrow 3 \text{ m}$$

$$1 \text{ Rotation} = \frac{1}{60} \text{ sec}$$

and it means we read the complete track once in $\frac{1}{60}$ sec.

$$\text{So, } \frac{1}{60} \text{ sec} \rightarrow \text{read } 128 \text{ KB}$$

$$\rightarrow 1 \text{ sec} \rightarrow 128 \times 60 \text{ KB}$$

$$\rightarrow (\text{Data Rate})_{\text{total}} = 128 \times 60 \times 16 \\ = \underline{120 \text{ MBPS}}$$

→ But there was nothing is given. So we go with 16 R/w heads for 16 surfaces.

(f) Using one read/write head, the disk is rotating at 6000 Rpm. What is data transfer rate?

$$6000 \text{ Rotations} \rightarrow 60 \text{ sec}$$

$$1 \text{ Rotation} \rightarrow \frac{1}{60} \text{ sec}$$

$$\frac{1}{60} \text{ sec} \rightarrow 128 \text{ KB}$$

$$1 \text{ sec} \rightarrow 128 \times 100 \text{ KB}$$

$$\begin{aligned}(\text{Data Rate}) &\approx 128 \times 100 \text{ KB} \\ &\approx \underline{13 \text{ MBPS}}.\end{aligned}$$

(g) If rotation speed is 3000 rpm, what is average access time with seek time of 13.5 ms?

$$\begin{aligned}T_{\text{access}} &= T_{\text{seek}} + T_{\text{rotation delay}} \\ &\quad + T_{\text{trans}}^6 \text{ (Data transfer time)} \\ &= 13.5 + \frac{1}{2}(\text{RT}) \\ &= (13.5 + 10) \text{ ms} \\ &= \underline{23.5 \text{ ms}}.\end{aligned}$$

(h) Find the avg. access time for transferring 512 bytes of data with following specification.

$$\text{Avg seek time} = 5 \text{ msec}$$

$$\text{disk rotation} = 6000 \text{ rpm}$$

$$\text{Data rate} = 40 \text{ KB/sec}$$

$$\text{Controller overhead} = 0.1 \text{ msec}$$

Cont-

$$(T_{\text{avg}})_{\text{avg}} = T_{\text{seek}} + T_{\text{rot}} + T_{\text{trans}}$$

$$+ 0.1$$

$$\Rightarrow 5 + 0.1 + T_{\text{rot}} + T_{\text{trans}}$$

$$6000 \rightarrow 60 \text{ sec}$$

$$1R = \frac{1}{100} \text{ sec}$$

$$\frac{1}{2}R = \frac{1}{200} \text{ sec} = \frac{0.5 \times 10^{-2}}{T_{\text{rot}}} = \frac{5 \text{ ms}}{T_{\text{rot}}}$$

$$T_{\text{trans}} = \frac{512}{40 \times 1024 \times 1024} = 12.8 \text{ ms}$$

$$= 12.5 \text{ ms}$$

$$(T_{\text{avg}})_{\text{avg}} = (5 + 0.1 + 5 + 12.5) \text{ ms}$$

$$= 22.6 \text{ ms}$$

Q: A certain moving arm disk storage with one head has following specification.

$$\text{No of tracks / surface} = 200$$

$$\text{disk rotation speed} = 2400 \text{ rpm}$$

$$\text{track capacity} = 62,500 \text{ bits}$$

$$\text{avg. latency} = p \text{ msec}$$

$$\text{data transfer rate} = q \text{ bits/sec}$$

Find the value of p and q ?

Solⁿ $p = \frac{1}{2}(\text{Rotation time})$

$$2400 \rightarrow 60 \text{ sec}$$

$$1R = \frac{1}{60} \text{ sec}$$

$$\frac{1}{2}R = \frac{1}{120} \text{ sec}$$

$$p = 12.5$$

$$(2 \times 12.5 \text{ ms}) \rightarrow 1R \rightarrow T_{\text{track}}$$

$$(2 \times 12.5) \rightarrow 62500 \text{ bits}$$

$$25 \text{ ms} \rightarrow 62500 \text{ bits}$$

$$1s = \frac{62500}{25 \times 10^3}$$

$$= 2.5 \times 10^6 \text{ bits/sec}$$

$$Q = 2500000$$

Q: Ques: A HD has 63 sectors per track, 10 platters

$$\text{i)} < 400, 16, 29 >$$

$$\text{ii)} < 2039 > .$$

Note: Do it by self first

Machine Instructions and

Addressing Modes

(12)

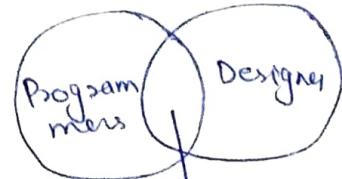
Introduction

Embedded System \leftrightarrow (uP)

Computer System

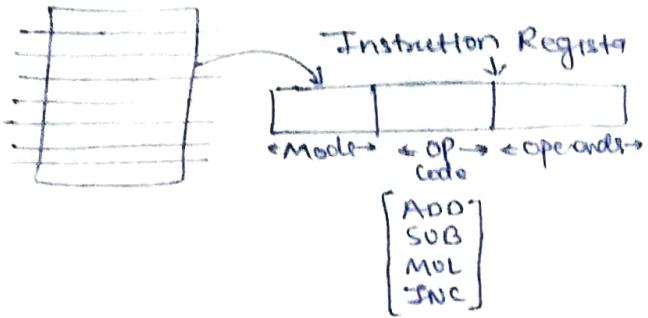
Main m/m, registers, ALU and control unit. We select uP to design a computer architecture according to cost selection might vary.

Generate control signal according to instruction called instruction decoding.



(Machine Instructions)

Machine instr. discussed by both before designing.



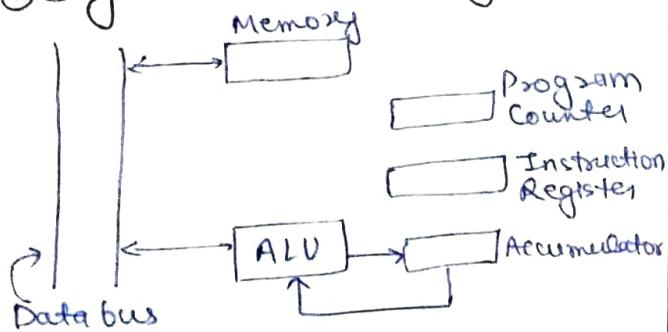
→ Mode defines about the operands.

Computer Organizations based on registers

Types of CPU organizations:-

- 1) Single accumulator organization.
- 2) General register organization
- 3) Stack Organization.

① Single accumulator organization



→ We use single address instruction, means we will have either 0 instruction or 1 instruction.

Ex → INC A : is a 0 address register bcz it has no need of mem address.

ADD X : → means add X to value of accumulator so, its a 1 address register.

$m[x]$ → operand present in m/m 'm' at address 'X'.

$$ADD X = m[X] + \text{Accumulator}$$

→ And finally value get stored in Accumulator.

② General register organization.

→ We have lot of registers for general purpose.

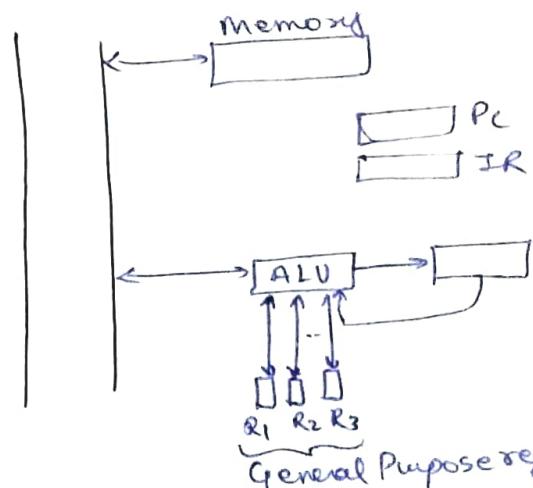
Ex - ADD R1, R2, R3 ;

$$R_3 \leftarrow R_1 + R_2$$

→ No. of instruction for executing program is very less, as compare to single accumulator.

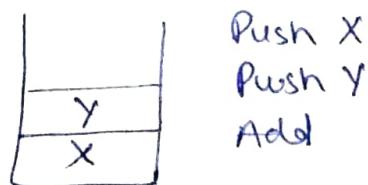
→ It is costly and faster.

→ We can have 2 address or 3 address instruction.



③ Stack Organization

→ Only zero instruction possible except Push or Pop.



→ We just give opcode Add without giving any operand.

→ No. of operands in instruction decide the size of instruction.

$$\text{Operands} \uparrow \leftrightarrow \text{Size} \uparrow$$

Addressing Mode

opcode	operand	mode
--------	---------	------

mode defines that how to treat operand.

operand can be variable and constant both but generally refers variables.

opcode is operation code to perform operation.

To reduce size of instruction we keep the address in register despite of m/m.

We can manage pointers, loops, counters etc.

Implied Add. Mode (Implicit Mode)

Used for 0 or 1 address instrn

operand specified implicitly in the instruction.

Ex → INC A, ↗ Add in stack organization
 CRC A, CL A, ADD organ ↗ complement
 ↳ Clear carry

Immediate Addressing Mode

No effective address required

operand provided as constants

We use constant.

Ex → Add R1, #3 ↗ data

We are not using variable here.

Constant size equals to length of operand.

Register Add. Mode

operand present in register.

Register no. written in instruction

Data present inside the registers.

Reduce the size of operand and instrn

→ High speed access.

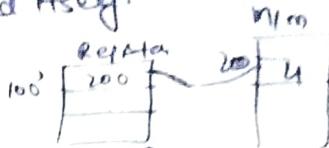
④ Register Indirect Mode

Register contains address of operand instead of operand itself.

Add R1

R100

Data = 4



Add (R1)

Re Act MERT

→ Effective address required.

→ Generally registers are less as compare to m/m, so give reference of m/m using registers to reduce the bits in operand.

⑤ Auto Increment or decrement

→ Special case of register indirect add. modes

→ for continuous sequential access.

→ Registers are collection of flip flops and we can use them as counters.

Ex ADD 500; ↗ Direct addressing mode

→ Absolute addressing mode.

→ Actual address given

→ Use to access variables

→ Direct m/m address given

→ Disadv → Large instruction size

② If address field is fixed size then variable value size is fixed.

⑦ Relative addressing mode

→ Use in program control instrn like, looping, branching etc

$$EA = \text{Offset} + \text{Program counter}$$

→ Beware of the value of Program counters

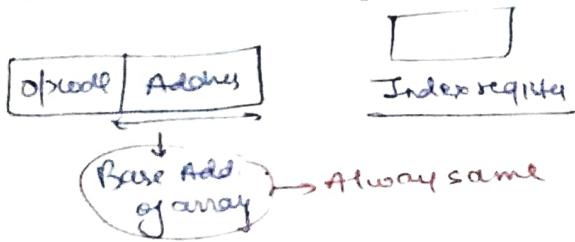
Base register addressing mode

- Used in 'program relocation'
- Useful in swap in/out of process.
- After swapping process add. might get changed.

$$EA = \frac{\text{Base Reg. Value} + \text{Displacement}}{\downarrow \text{Starting Add.}}$$

⑧ Indexed Add. Mode

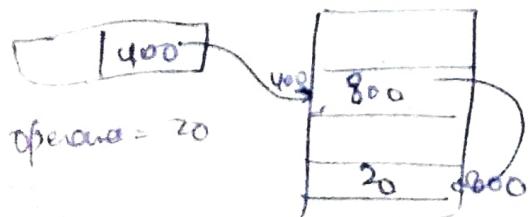
- Used to implement array efficiently.
- Multiple registers required to implement.
- Any element can be accessed without changing instruⁿo



$$EA = \text{Base Add} + \text{Index Register}$$

⑨ Indirect Add. Modes

- Use to implement pointers and passing parameters.
- 2 m/m access required.



$$EA = AC + M[M[x]]$$