

Cycle time of the non-pipelined processor (round off to two decimal places) is \_\_\_\_\_

$$\text{Cycle time} = \frac{1}{24\text{MHz}} = 0.5\text{ nsec}$$

# Stalls/instruction =

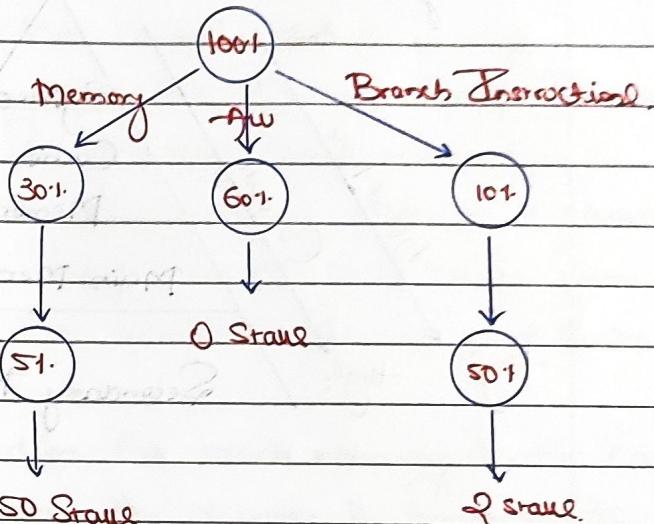
$$= 30 \times 0.05 \times 50$$

$$+ 60 \times 0$$

$$+ 10 \times 50 \times 2$$

$$\Rightarrow 0.85$$

$$\text{Avg Post. ET} = (1 + 0.85) \times 0.5 \\ = 1.85 \times 0.5 \\ = \underline{\underline{0.925 \text{ nsec}}}$$



Nonpipeline: 5 clock cycles

2.5 GHz

$$\text{Cycle time} = \frac{1}{2.5} \text{ sec} = \underline{\underline{0.4 \text{ nsec}}}$$

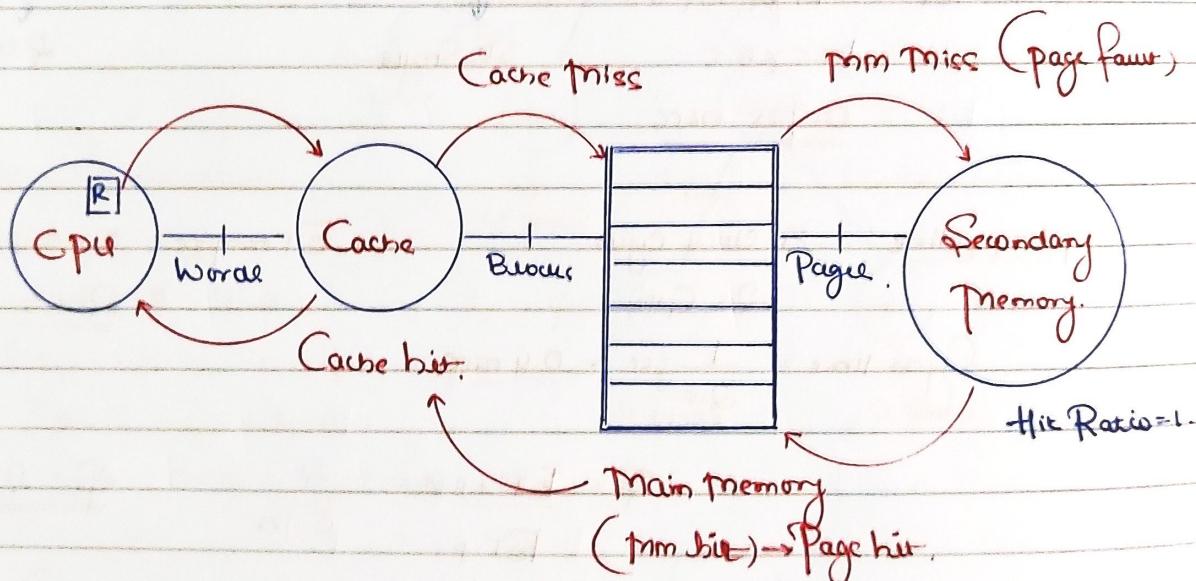
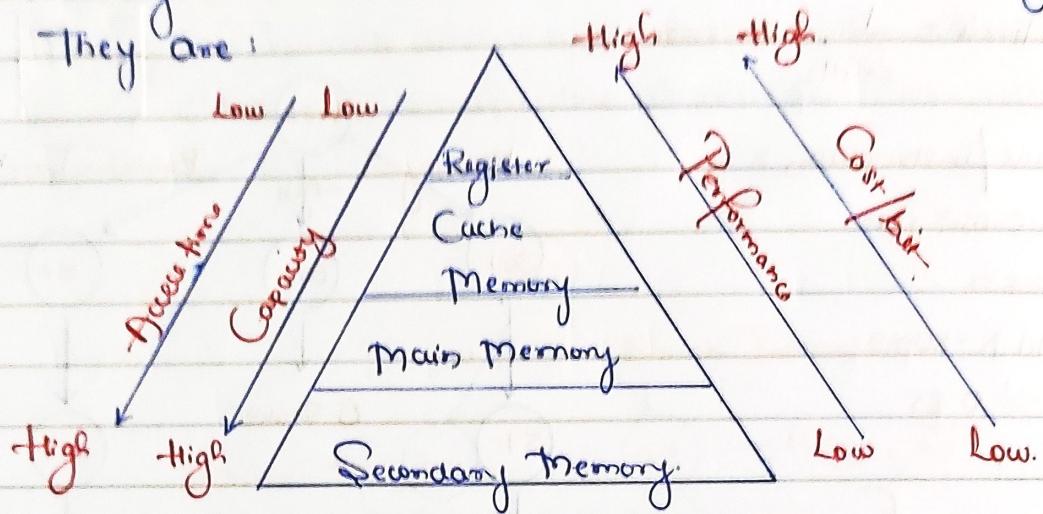
$$\text{ET}_{\text{nonpipe}} = 5 \times 0.4 \text{ nsec} \\ = \underline{\underline{2 \text{ nsec}}}$$

$$S = \frac{\text{ET}_{\text{NP}}}{\text{ET}_{\text{P}}} = \frac{2}{0.925} \therefore \underline{\underline{S = 2.16 \text{ (Ans)}}}$$

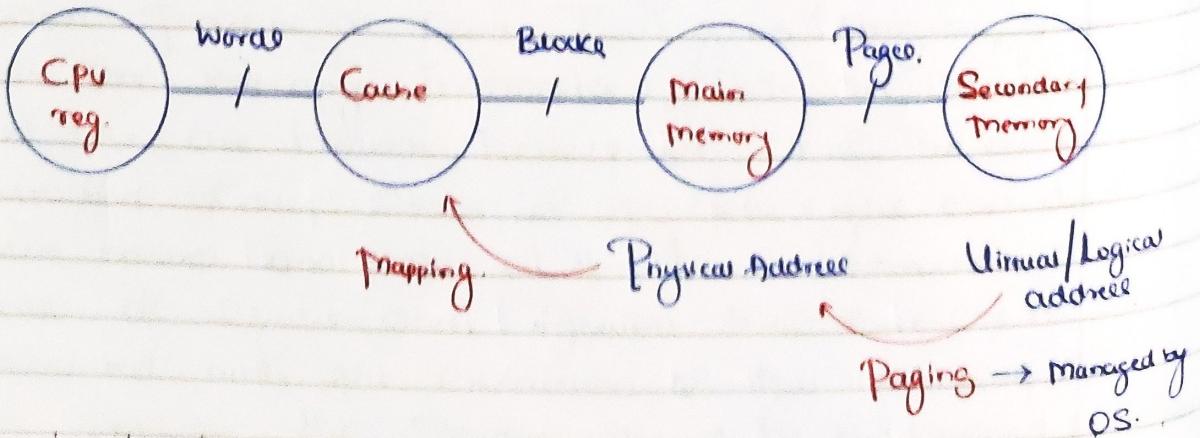
# CACHE MEMORY

Memory Hierarchy: Hierarchy design organizes the system supporting memory into 4 levels to minimize the accessing time.

They are:



$$\text{Hit Ratio} = \frac{\text{Number of hits}}{\text{Total number of access}}$$



M	T	W	T	F	S	S
Page No.:						YOUVA
Date:						

## Memory:

1. Cpu generate request initially refer to the Cache.
2. If the reference (respective data) found in the Cache that is called Cache hit (operation is carried here) then respective data give from Cache to Cpu in the form of words.
3. If the reference is not found in the Cache, that is known as Cache miss, the reference is forwarded to main memory.
4. If the reference found in the main memory then it is called Mm hit or page hit then respective data is given from main memory to Cache in the form of blocks and Cache to Cpu in the form of words.
5. If the reference is not found in the main memory that is called Mm miss or page fault then the response is forwarded to Secondary memory. So respective data is transferred from Secondary memory to main memory in the form of pages, main memory to Cache memory in the form of blocks, then Cache memory to Cpu in the form of words.

\* Note: Secondary memory is the last level of memory in which hit ratio is always "1".

\* The process of transfer the data from main memory to Cache memory is called "mapping".

## Average Memory Access Time [Tavg]

$$T_{avg} = h \times \text{Time taken by } [H] \text{ memory when there is a hit.} + (1-h) \left\{ \begin{array}{l} \text{Time taken by } \\ \text{memory when} \\ \text{there is a miss} \end{array} \right\}$$

- Q1. Consider a Cpu Generate 100 Request. Out of 100 requests got time hit and 10 times miss. When there is a hit then time taken is 20nsec and when there is a miss then time taken is 150nsec. Calculate the Tavg?

Ans.

Total CPU Request = 100

Hit = 90 times

Miss = 10 times

Hit will take 20 nsec

Miss will take 150 nsec.

$$\text{Total time} = 90 \times 20 + 10 \times 150$$

$$= 1800 + 1500$$

$$= 3300 \text{ nsec.}$$

$$T_{avg} = \frac{3300}{100} = 33 \text{ nsec (Ans)}$$

Total CPU Request = 100

Hit = 90 times, Miss = 10 times

$$\text{Hit Ratio} = \frac{90}{100} = 0.9$$

Time =  $H \times \text{Time taken when there is hit} + (1-H) (\text{Time taken when there is miss})$

$$= 0.9 \times 20 + 0.1 \times 150$$

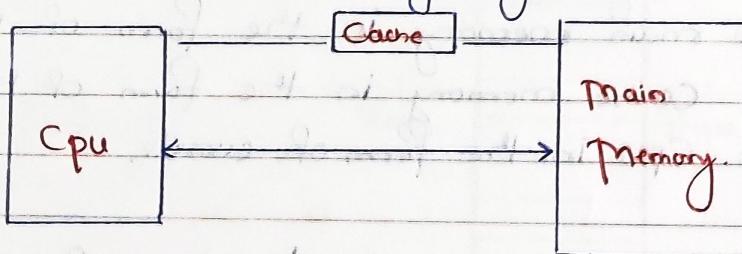
$$= 18 + 15$$

$$T_{avg} = 33 \text{ nsec (Ans)}$$

## Type of Memory Organisation:

1.

### Simultaneous Access Memory Organisation:

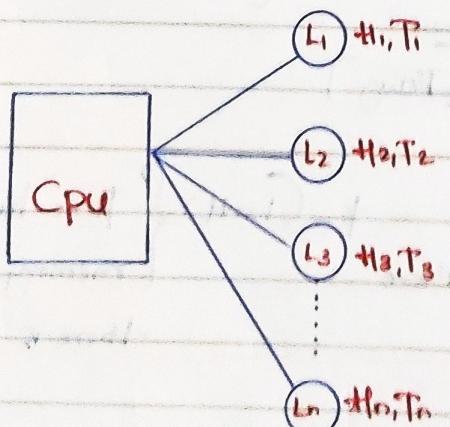


$$T_{avg} = b \cdot t_c + (1-b) \cdot t_m$$

→  $b$ : hit ratio

$t_c$ : Cache access time

$t_m$ : main memory access time



In the simultaneous access organisation all the level of memory is directly connected to CPU. But follow in sequence (Access a sequence) when there is a miss in level 1 memory then reference forward to level 2 memory.

When there is a hit in level 2 memory, then directly data is transferred from level 2 to CPU without copying into level 1 memory.

When there is a miss in level 1 and level 2 memory and hit in level 3 memory then directly data is transferred from level 3 memory to CPU without copying into level 1 and level 2 memory.

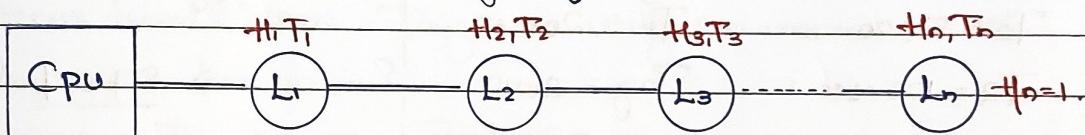
Here  $h_1, h_2, h_3, \dots$  are the hit ratios and  $T_1, T_2, T_3, \dots$  are the access of respective level memory.

The time required to access 1 word from memory.

$$T_{avg} = h_1 T_1 + (1-h_1) h_2 T_2 + (1-h_1)(1-h_2) h_3 T_3 + \dots (1-h_1)(1-h_2)\dots h_n T_n$$

\* In last level hit ratio = 1.

## ② Hierarchical Access Memory Organisations



When there is a miss in level 1 and level 2 memory and hit in level 3 memory then firstly data is transferred from level 3 [L3] memory to level 2 [L2] memory then level 2 [L2] memory to level 1 [L1] memory then from level 1 [L1] memory to CPU.

$$T_{avg} = h_1 T_1 + (1-h_1) h_2 (T_2 + T_1) + (1-h_1)(1-h_2) h_3 (T_3 + T_2 + T_1) + \dots (1-h_1)(1-h_2)\dots h_n [T_n + T_{n-1} + \dots + T_2 + T_1]$$

Note: In an question if they have mentioned the keyword 'hierarchical access' or 'level of access' or hierarchical then it means Using "Hierarchical Access".

Q2. Calculate the average access time with the cache access time 1 nsec, And main memory access time 100 nsec, Hit ratio 90.1? Using Hierarchical access?

Ans.  $t_c = 1 \text{ nsec}$        $H = 90.1 = 0.9$

$t_m = 100 \text{ nsec}$

$$T_{avg} = h * t_c + (1-h) (t_m + t_c)$$

$$= 0.9 * 1 + (1-0.9) (100+1)$$

$$= 0.9 + (0.1)(0.1)$$

$\Rightarrow$  11 nsec

Q. In a 2 level memory, level 1 memory is 5 times faster than level 2. And its access time is 10ns < Average access time. If level 1 access time is 20ns, what is the hit ratio? Using Signatures - one Access org?

$L_1 = T_1$	$S = \frac{P_{L1}}{P_{L2}} = \frac{1/T_1}{1/T_2}$	$T_{avg} = b * tc + (1-b) * tm$
$L_2 = T_2$		$30 = H * 20 + (1-H) * 100$
$S = \frac{T_2}{T_1}$		$= 20H + 100 - 100H$
$T_2 = S^* T_1$	$T_1 = T_{avg} - 10$	$80H = 70$
$T_2 = S * 20$	$T_{avg} = T_1 + 10$	$H = \frac{70}{80} = 0.875$
$T_2 = 100 \text{ nsec}$	$T_1 = 20 \text{ nsec}$	
	$T_{avg} = 20 + 10$ $= 30 \text{ nsec}$	$\Rightarrow \underline{87.5\%} \text{ (Ans)}$

Q4 Assume that for a certain processor, a read request takes 50 nsec on a Cache miss and 5 nsec on a Cache hit. Suppose while running a program, it was observed that 80% of the processor's read requests result in Cache hit. The average read access time in nsec is

Ans | Hit Ratio =  $80 = 0.8$  When there is hit, time = 5

Miss ratio =  $(1-0.8) = 0.2$  When there is miss, time = 50

$$\begin{aligned}
 T_{avg} &= 0.8 \times 5 + 0.2 \times 50 \\
 &= 4 + 10 \\
 &= \underline{\underline{14 \text{ msec}}} \quad (\text{Ans})
 \end{aligned}$$

## Hierarchical Access

$$T_{\text{aug}} = b^* + c + (1-b)[+m + tc]$$

$$\Rightarrow h_{TC} + h_{TM} = h_{TMC} + tc - h_{RC}$$

$$\Rightarrow fct \cdot t_m - b t_m$$

$$T_{avg} = tc + (1-t)c_m$$

eg. Hierarchical access bit ratio 80:1,  $t_c = 20\text{ns}$ ,  $t_m = 100\text{ nsec}$

$$\begin{aligned} T_{avg} &= b \cdot t_c + (1-b) [t_m + t_c] \\ &= 0.8 \times 20 + (1-0.8) (100+20) \\ &= 16 + 0.2 (120) \end{aligned}$$

$$T_{avg.} = \underline{\underline{40\text{nsec}}}$$

$$\begin{aligned} T_{avg} &= t_c + (1-b) \cdot t_m \\ &= 20 + (1-0.8) 100 \\ &= 20 + (0.2) 100 \end{aligned}$$

$$T_{avg.} = \underline{\underline{40\text{nsec}}}$$

Qs. Consider a System with 2 levels. Level 1 Access time is 20ns level 2 Access time = 150ns  $T_{avg} = 30$  using Simultaneous access.

(i) What is the bit ratio?

Ans.  $T_1 = 20\text{ns}, T_2 = 150\text{ nsec} \quad T_{avg} = 30\text{ nsec}$

$$T_{avg} = b \cdot t_1 + (1-b) \cdot t_2$$

$$30 = b \cdot 20 + (1-b) 150$$

$$30 = 20b + 150 - 150b$$

$$120 = 130b$$

$$\frac{b}{1-b} = \frac{120}{130} = 0.923 \Rightarrow \underline{\underline{92.3:1}}$$

(ii) If the bit ratio is made to 100:1 then what is the access time of L<sub>1</sub> and L<sub>2</sub> memory?

Ans.  $T_1 = 20\text{ns}, T_2 = 150\text{ nsec}$

L<sub>1</sub> and L<sub>2</sub> access remain same,  $b=100:1$

If bit ratio = 100:1

$$\begin{aligned} T_{avg} &= b \cdot t_1 + (1-b) \cdot t_2 \\ &= 1 \times 20 + (1-1) 150 \end{aligned}$$

$$T_{avg} = \underline{\underline{20\text{nsec}}}$$

Only  $T_{avg}$  changes

Q6. In the above question If  $T_{avg}$  is increased by 10:1 then what is the % change in the bit ratio?

Ans.  $T_1 = 20, T_2 = 150\text{nsec}, T_{avg} = 30\text{nsec} + 10:1 = 33\text{nsec}$

$$T_{avg} = h \cdot t_1 + (1-h) \cdot t_2$$

$$38 = h \cdot 20 + (1-h) \cdot 150$$

$$38 = 20h + 150 - 150h$$

$$130h = 112$$

$$h = \frac{112}{130} = \frac{t_1}{t_1+t_2} = 0.9$$

= 90%

Previous hit ratio = 92.3%.

→ Hit ratio decreased by 2.3%.

Note: \* 1 Word Access Time =  $T_{avg}$ .

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

$T_{avg}$

Data Transfer Rate.

\* Hierarchical Access Memory Org: 3 level.  $h_1$ : hit ratio in level 1  
 $(1-h_1)$ : miss ratio in level 1

$$T_{avg} = h_1 \cdot t_1 + (1-h_1) \cdot h_2 (t_2 + t_1) + (1-h_2) \cdot h_3 (t_3 + t_2 + t_1)$$

$(m_1)$

$(1-h_2)$ : miss ratio in level 2

$(m_2)$

$$T_{avg} = h_1 \cdot t_1 + m_1 \cdot h_2 (t_2 + t_1) + m_1 \cdot m_2 \cdot h_3 (t_3 + t_2 + t_1)$$

$m_1$ : miss ratio in level 1.

Q7. Assume that for a certain processor, a read request takes 5 ns on a Cache miss and 5 ns on a Cache hit. Suppose while running a program, it was observed that 80% of the processor's read request result in a Cache hit. The average read access time in nanoseconds is \_\_\_\_\_.

(i) What is the data transfer rate (performance) of this memory system (in words/sec)?

Ans.

$$\text{Data Transfer rate} = \frac{1}{T_{avg}} \text{ words/sec}$$

$$\frac{1000 \times 10^9}{14}$$

$$= \frac{1}{14 \times 10^9} \text{ words/sec}$$

$$= 71.42 \times 10^6$$

$$= \frac{1}{14} \times 10^9 \text{ words/sec}$$

$$= 7.142 \text{ million words/sec}$$

(Ans)

(ii) What is the bandwidth required of the memory system if word size is 8bit?

Anc **Bandwidth:** Word size = 8bit

$$72 \text{ million words/sec} = 72 \times 10^6 \text{ words/sec}$$

$$= 72 \times 10^6 \times 8 \text{ bit per sec}$$

$$= 72 \times 8 \text{ M bit/sec}$$

$$= 576 \text{ M bit/sec or } 72 \text{ Mbytes/sec (Ans)}$$

### Locality of Reference

Accessing the higher level of memory data from the lower memory (Cache memory) is called locality of reference (Faster Memory).

1. Temporal LOR

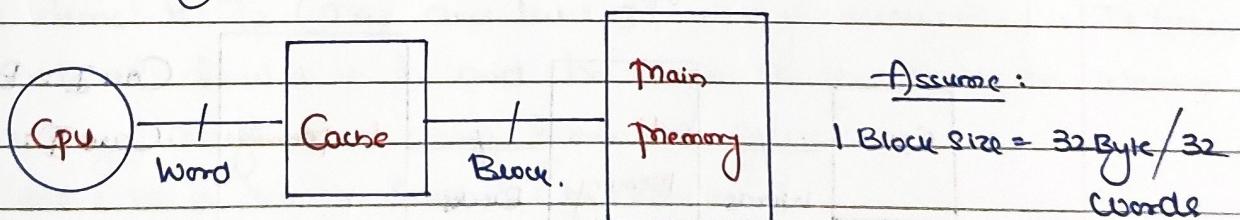
2. Spatial LOR

① **Temporal LOR:** This means the same word in the same block is referenced by the CPU in near future (Frequency).

Or

Same data which access again and again then that type of data stored in Temporal LOR.

② **Spatial LOR** means adjacent word is in the same block is referenced by the CPU in a Sequence.



If miss in Cache then 1 complete block is transferred from main memory to Cache, the respective word (which is requested by the CPU) given from Cache memory to CPU.

(Next reference there is a Cache hit)

Q.1 If block size = 32 word / 32 Byte & CPU require only 1 word?

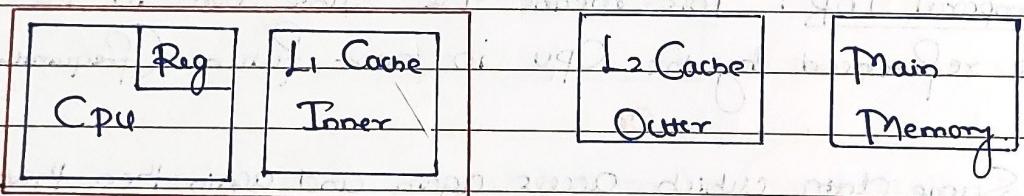
Anc In this process, 1 complete block of 32 word / 32 Byte transferred

from main memory to Cache memory then requested word is transferred from Cache memory to CPU.

Note: So in the next instructions time when CPU request some word or adjacent word then that request available in Cache (Cache hit).

### Types of Cache

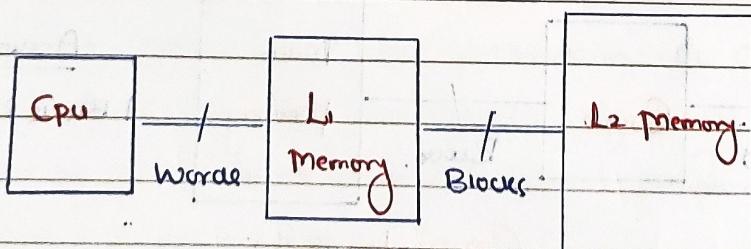
- ① Unified Cache: Instruction and data both are placed in same Cache.
- ② Split Cache: This Cache logically divided into two parts
  - (i) Instruction Cache (I-cache)
  - (ii) Data Cache (D-cache)
- ③ Multilevel Cache:



Size  $L_1 < L_2$

Speed  $L_1 > L_2$

### 2 Level of Memory (If locality of reference included)



Case I: Block Size = 1 word

Case II: Block size = n words

CPU always access the data from the fast memory (Level 1/Cache memory). If there is a miss in Level 1 (Cache) memory then one complete block is transferred from Level 2 (Slow) memory to Level 1 (Cache) memory and addressed word (which is requested) that respective word given from Cache memory to CPU.

TB: Block transfer time from L2 memory to L1 memory.

→ Case 1: Block size = 1 word

$$TB = T_2$$

→ Case 2: Block size = n words

$$TB = n \cdot T_2$$

If 2 level memory:

$$h_2 = 1$$

$$T_{avg} = h_1 t_1 + (1-h_1) (t_2 + t_1)$$

or

$$T_{avg} = t_1 + (1-h_1) t_2$$

If locality of reference included

$$(TB = n \cdot T_2)$$

$$T_{avg} = h_1 t_1 + (1-h_1) [TB + t_1]$$

or

$$T_{avg} = t_1 + (1-h_1) TB$$

If 3 level memory:

$$h_3 = 1$$

$$T_{avg} = h_1 t_1 + (1-h_1) h_2 (t_2 + t_1) + (1-h_1)(1-h_2) (t_3 + t_2 + t_1)$$

If locality of reference included:

$$T_{avg} = h_1 t_1 + (1-h_1) h_2 (TB_1 + t_1) + (1-h_1)(1-h_2) [TB_2 + T_1 + t_1]$$

Q1. In a 3 level memory, level 1 memory access time is  $T_1$ , level 2 memory access time is  $T_2$  ( $T_{B1}$ ) and level 3 memory access time is  $T_3$  ( $T_{B2}$ ). Hit ratio of level 1 is  $h_1$  and hit ratio of level 2 is  $h_2$ . What is the average access time using hierarchical access?

(i) If there is a hit in level 1 ( $h_1 = 100\%$ )?

$$h_1 = 100 \Rightarrow 1$$

$$(1-h_1) = 0$$

$$\underline{T_{avg} = T_1}$$

$$h_1 = 100 \Rightarrow h_1 = 1$$

$$(1-h_1) = 0 \Rightarrow 0 \cdot x = 0$$

$$\underline{T_{avg} = T_1}$$

(ii) If there is a miss in level 1 and hit in level 2 ( $h_2 = 100\%$ )

$$h_1 = 0 \& h_2 = 100 \Rightarrow h_2 = 1$$

$$1-h_2 = 0$$

$$\underline{T_{avg} = T_2 + T_1}$$

$$h_1 = 0, h_2 = 1 \quad (1-h_2) = 0$$

$$\underline{T_{avg} = T_2 + T_1}$$

- (iii) If there is a miss in level 1, and level 2 and hit in level 3  
 Ans.  $h_1=0, h_2=0, h_3=1$        $h_1=0, h_2=0, h_3=1$   
 $T_{avg} = T_B + T_{L2} + T_L$        $T_{avg} = T_{B2} + T_{L1} + T_L$

Q2 In a 2 level memory, level 1 memory access time is 30ns and level 2 memory access time is 250ns/word. Hit ratio of level 1 is 90%. If there is a miss in level 1 then 4 word block must be transferred (bword) from level 2 into level 1 and then addressed words is given to the CPU. What is the average access time?

Ans.  $h_1 = 90\% = 0.9$        $t_1 = 30\text{ nsec}$

$$t_2 = 250 \text{ nsec/word}$$

$$T_{B1} = 4 \text{ word} \times 250 \text{ ns/word} = 1000 \text{ nsec.}$$

$$T_{avg} = h_1 t_1 + (1-h_1) (T_{B1} + t_1)$$

$$= 0.9 [30] + (1-0.9) [1000+30]$$

$$= 27 + 0.3 \Rightarrow \underline{130 \text{ nsec}} \quad (\text{Ans})$$

Q3. A cache memory that has a hit rate of 0.8 has an access latency 10 ns and miss penalty 100 ns. An optimization is given done on the cache to reduce the miss rate. However, the optimization results in an increase of cache access latency to 15 ns, whereas the miss penalty is not affected. The minimum hit rate (rounded off to two decimal place) needed to after the optimization such that it should not increase the average memory access time is \_\_\_\_\_.

Ans  $h = 0.8, t_c = 10 \text{ nsec}, \text{miss penalty } (t_m) = 100 \text{ nsec}$

$$T_{avg} = h t_c + (1-h) (t_m + t_c)$$

$$= 0.8 \times 10 + (1-0.8) (100+10)$$

$$= 0.8 \times 8 + 0.2 (110) \Rightarrow \underline{30 \text{ nsec}}$$

New Optimization :  $t_{c\text{ new}} = 15 \text{ nsec}$

$$T_{avg\text{ (new)}} = h + t_c + (1-h) (t_m + t_c)$$

$$30 = h \times 15 + (1-h) (100+15)$$

$$30 = 15h + 115 - 115h + h$$

$$\frac{85}{100} \times 100 + \frac{1}{100} \times \frac{85}{100} = 0.85 \text{ (Ans)}$$

- Q4. A direct mapped cache memory of 1MB has a block size of 256 Bytes. The Cache has an access time of 3ns, and a hit rate of 94%. During a cache miss, it takes 20ns to bring the first word of a block from the main memory, while each subsequent word takes 5ns. The word size is 64 bits. The average memory access time in ns is \_\_\_\_\_.

Ans. Block Size = 256 Byte

$$1 \text{ word size} = 64 \text{ bit} = 8 \text{ Byte} \quad \text{Total} = 94 \times 3 + (1-0.94)(2+20)$$

$$\text{No. of words} = \frac{256 \text{ Byte}}{8 \text{ Byte}} = 32 \text{ words}$$

$$= 2.82 + (0.06)(3+20+155)$$

1 word takes 20ns

$$= 13.5 \text{ nsec (Ans)}$$

Remaining 31 words take 5ns

## Basic Elements of Cache Design:

1. Memory Organisation
2. Mapping technique
  - i. Direct mapping
  - ii. Set associative mapping
  - iii. Fully associative mapping
3. Replacement algorithm
4. Updating technique
5. Multi Level Cache

- \* In the cache design, data is transferred from main memory to Cache memory in the form of blocks.
- \* So both memory (main memory and Cache memory) must be organised based on block size.
- \* Main memory and Cache memory are divided into pairs (equal size pairs) based on block size called Main Block and Cache Block (number of lines) respectively.

$$\# \text{ Main Block} = \frac{\text{Mem Size}}{\text{Block Size}}$$

$$\# \text{ Cache Block} = \frac{\text{Cm Size}}{\text{Block Size}}$$

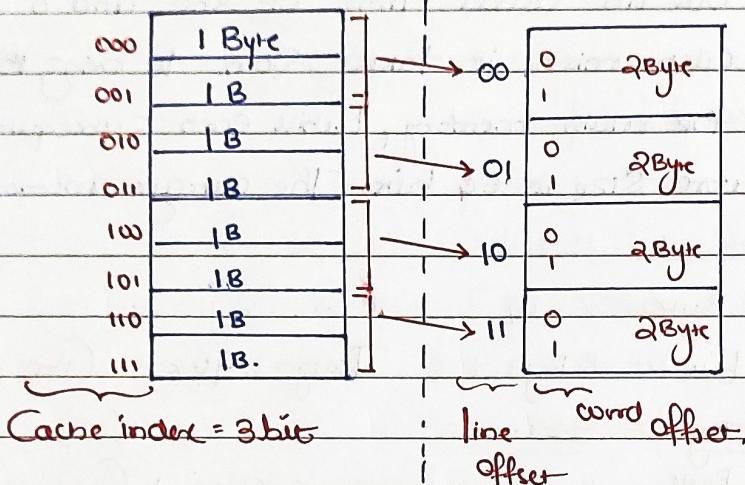
## Cache Memory:

## Before Organisation:

$$C_m = 8 \text{ Byte}$$

\* Consider a 8 Byte Cache with 2 Byte Block Size then  
Cm. organisation is as follows:

| After Organisation



~~#Cm Linee = .Cm Size~~

## Block size

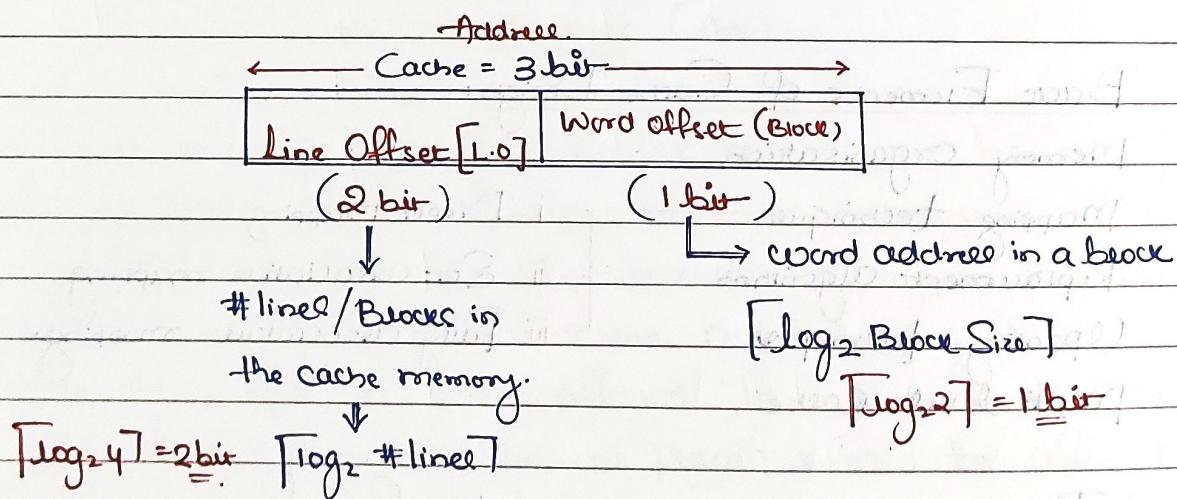
$$= \frac{8B}{2B} = \underline{\underline{4}}$$

$$\text{Cm Size} = 4 \times 2 \text{ Byte}$$

$$= \underline{\underline{8 \text{ Byte}}} \text{ (After)}$$

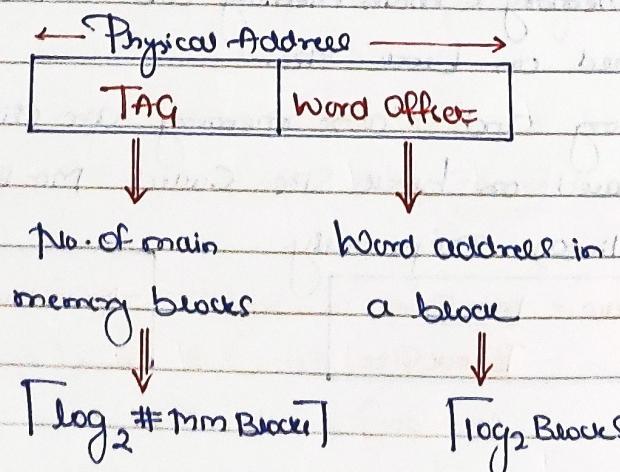
## Organisation

Before and after the organisation Cache memory Capacity is same but internal Structure is different.



Physical address is interpreted as

mm Size > Cm Size



Q1. Consider a Cache memory which is indexed using 16 bit address, organised into 32 byte block size and physical address is 32 bits.

- #Cm Block
- #Mm Block
- Cm Address Format
- Mm Address Format

Ans

$$\text{Cache Index} = 16 \text{ bits} \Rightarrow \text{Cm Size} = 2^{16} \text{ Bytes} = 64 \text{ KB}$$

$$\text{Block Size} = 32 \text{ Bytes}$$

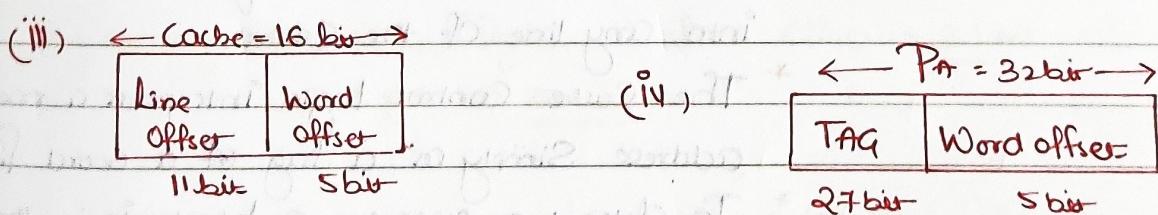
$$\text{Block/word offset} = \lceil \log_2 32 \rceil \\ = 5 \text{ bits}$$

$$\text{Physical address} = 32 \text{ bits}$$

$$\text{Mm Size} = 2^{32} \text{ Bytes} = 2^{20} \text{ Gbytes} \\ = 4 \text{ GBytes}$$

$$(i) \#Cm Block = \frac{\text{Cm Size}}{\text{Block Size}} = \frac{64 \text{ KB}}{32 \text{ B}} \\ = \frac{2^{16} \text{ B}}{2^5 \text{ B}} = 2^5 = 32 \text{ Lines}$$

$$(ii) \#Mm Block = \frac{\text{Mm Size}}{\text{Block Size}} = \frac{4 \text{ GB}}{32 \text{ B}} \\ = \frac{2^{32} \text{ B}}{2^5 \text{ B}} = 2^{27} \\ = 128 \text{ Mm Blocks}$$



Q2. Consider a System Mm = 256 MB, Cache = 128 KB, and Block Size = 512 Byte then

- #Lines
- L.O and W.O Format
- #Mm Block
- Tag and W.O Format

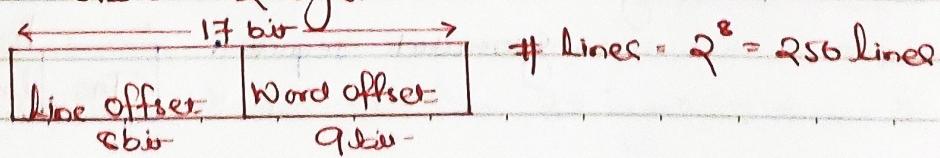
Ans.

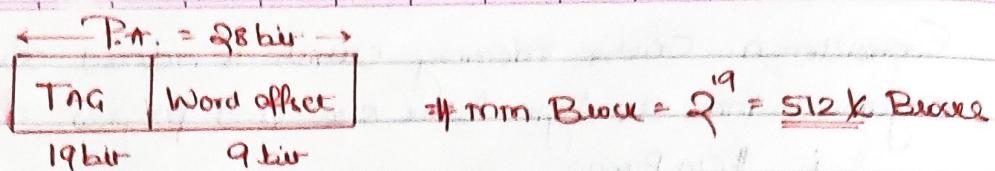
$$\text{Mm Size} = 256 \text{ MB} = 2^{28} \text{ Bytes}$$

$$\text{Cache Size} = 128 \text{ KB} = 2^{17} \text{ Bytes}$$

$$\text{Block Size} = 512 \text{ B} = 2^9 \text{ Bytes}$$

$$W.O = 9 \text{ bits}$$





Mapping: The process of transferring the data from main memory to Cache memory is called mapping. There are 3 types of mapping:

1. Direct mapping
2. Set Associative mapping
3. Fully associative mapping.

Mapping function: \* Because there are fewer Cache lines than main memory blocks, an algorithm is needed for mapping main memory blocks into Cache lines.

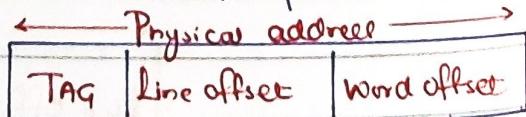
\* Three techniques can be used:

- ① Direct: \* The simplest technique  
\* Maps each block of main memory into only one possible Cache line.

- ② Associative: \* Permits each main memory block to be loaded into any line of the Cache.  
\* The Cache Control logic interprets a memory address simply as a tag of a word field.  
\* To determine whether a block is in the Cache, the Cache Control logic must simultaneously examine every line's Tag for a match.

- ③ Set Associative: \* A compromise that exhibits the strengths of both the direct and associative approaches while reducing their disadvantages.

Direct Mapping: In this Direct Cache Controller interprets the CPU generated request:



# Line = Cm Size  
Block Size

$$\begin{aligned}\text{Word offset} &= \log_2 (\text{BlockSize}) \\ \text{Line offset} &= \log_2 (\text{Line})\end{aligned}$$

- TAG = Physical Address - (Line Offset + Word Offset)
  - TAG Memory Size = # Lines x Tag bits (depend on the mapping technique)

Q3. Consider a direct mapping if the size of the Cache memory is 512 KB and main memory 512 MB and Cache line size (Block) is 64 KB calculate the bit required for



$$\text{Ans} \quad \text{Cache Size} = 512 \text{ KB} \Rightarrow 2^{19} \text{ Bytes}$$

$$\text{Mem Size} = 512 \text{ MB} \Rightarrow 2^{29} \text{ Byte} \quad PA = 29 \text{ bit}$$

$$\text{Block Size} = 64 \text{ KB} = 2^{16} \text{ Byte} \quad w_0 = 16 \text{ bit}$$

TAG	L.O	W.O
10 bits	3 bits.	16 bits.

$$\# \text{ lines} = \frac{\text{Cm Size}}{\text{Block Size}} = \frac{512 \text{ KB}}{64 \text{ KB}} = \frac{2^{19} \text{ Bits}}{2^{16} \text{ Bits}} = 8 \text{ lines}$$

$$\text{Page memory size} = \# \text{ lines} \times \text{Page bit} = 8 \times 10 \Rightarrow \underline{\underline{80 \text{ bits}}}$$

Q4 Consider a direct mapping, Cache Size = 64 Byte, Line Size = 8 Byte, mm = 256 Byte then # bits for P.A., TAG, L.O., W.O., Tag memory size?

The diagram illustrates the structure of a 3-bit address (PA). It consists of three adjacent boxes labeled TAG, L.O, and W.O. Above the boxes, a double-headed arrow spans the width of the three boxes, indicating that the total width of the address is 3 bits.

$$\# \text{ Lines} = \frac{\text{Cm Size}}{\text{Block Size}} = \frac{64B}{8B} = 8 \text{ Lines}$$

$$1.0 = 3 \text{ bits}$$

$$\begin{aligned} \text{Tag Memory Size} &= \# \text{Lines} \times \text{Tag bits} \\ &= 8 \times 2 \Rightarrow 16 \text{ bits} \end{aligned}$$

Q5. Consider a machine with a byte addressable main memory of  $2^{32}$  bytes, divided into blocks of size 32 bytes. Assume that a direct

mapped Cache having 512 Cache lines is used with this machine.  
The size of the tag field is bits is 18.

Ans



$$18 \text{ bit} \quad 9 \text{ bit} \quad 5 \text{ bit} \quad \Rightarrow \underline{18 \text{ bits}} \text{ (Ans)}$$

$$\text{TAG} = 32 - (9+5)$$

Q6. Consider a direct mapping, Cache Size = 128KB, Line Size = 64B.

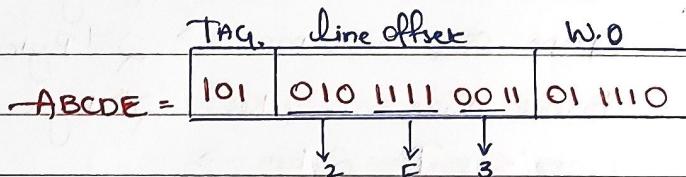
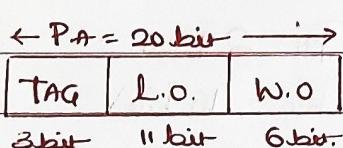
Main memory is 1MB then what is the line number of physical address  $(ABCDE)_{16}$ ?

Ans

$$1\text{MB} = 2^{20} \text{B} \rightarrow \text{P.A} = 20 \text{ bit}$$

$$\text{Line Size} = 64B = 2^6 \text{B} \Rightarrow \text{W.O} = 6 \text{ bit}$$

$$\# \text{Lines} = \frac{\text{Cache Size}}{\text{Block Size}} = \frac{128 \text{ KB}}{64 \text{ B}} = \frac{2^{19} \text{ B}}{2^6 \text{ B}} = 2^{\underline{13}} \text{ lines}$$



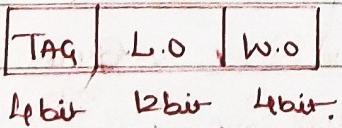
$$\begin{aligned} \text{TAG} &= 20 - (11+6) \\ &= \underline{3 \text{ bit}} \end{aligned}$$

$$= (2F3)_{16} \text{ (Ans)}$$

Q7. Consider a machine using a byte addressable main memory of  $2^{20}$  bytes, block size of 16 bytes and a direct mapped Cache having  $2^{12}$  Cache lines. Let the addresses of two consecutive bytes in main memory be  $(E201F)_{16}$  and  $(E2020)_{16}$ . What are the tag and Cache line address (in hex) for main memory address  $(E201F)_{16}$ ?

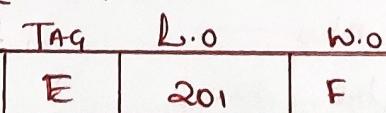
Ans

$$\text{P.A} = 20 \text{ bit}$$

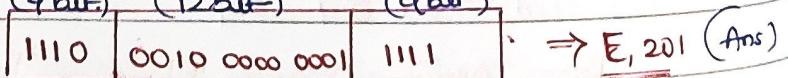


4bit 12bit 4bit.

$$(E201F)_{16}$$



(4bit) (12bit) (4bit)

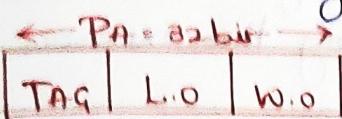


$\rightarrow E, 201 \text{ (Ans)}$

Q8. Consider a machine using a byte addressable main memory of  $2^{32}$  bytes divided into blocks of size 32 bytes. Assume that a direct mapped Cache having 512 Cache lines is used with this machine.

The size of the tag field in bits is \_\_\_\_\_

Ans



18bit 9bit 5bit.

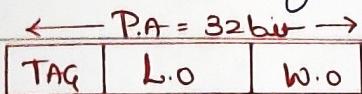
$$\text{TAG} = 32 - (9+5)$$

$\Rightarrow 18 \text{ bit}$  (Ans)

Q9. Consider a Computer System with a byte-addressable primary memory of size  $2^{32}$  bytes. Assume the computer system has a direct mapped Cache of size 32 KB ( $1 \text{ KB} = 2^{10}$  bytes) and each cache block of size 64 bytes.

The size of tag field is 17 bits.

Ans



17bit 9bit 6bit.

$$\begin{aligned} \text{TAG} &= 32 - (9+6) \\ &= 17 \text{ bit} \quad (\text{Ans}) \end{aligned}$$

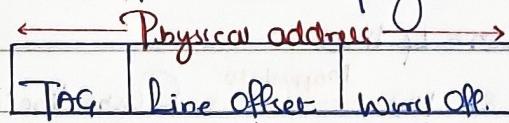
$$\# \text{Lines} = \frac{\text{Cm Size}}{\text{Block Size}} = \frac{32 \text{ KB}}{64 \text{ B}} = \frac{2^{13} \text{ B}}{2^6 \text{ B}}$$

$= 2^9 \text{ Lines}$

$$L.O = 9 \text{ bit}$$

### Direct Mapping (Continued)

In direct mapped Cache physical address is interpreted as



No. of mm

blocks are  
possible in one  
cache line. But no  
at same time

$$\begin{aligned} \text{TAG} &= \text{mm Size} \\ &\quad \text{Cm Size} \end{aligned}$$

Number of lines

$$\lceil \log_2 \# \text{Lines} \rceil$$

No. of address bits in

a block

$$\lceil \log_2 \text{Block Size} \rceil$$

eg:  $\text{TAG} = 2 \text{ bit} \Rightarrow 2^2 = 4 \text{ mm Block are mapped to one Cache line}$

In direct mapping an mapping function is used to transfer (copying) the data from main memory to Cache memory.

Mapping Function:  $K \bmod N = i$

K: mm Block no./mm request

N: Number of Cm Line

i: Cache line number (0 to n-1)

$k \bmod N = i$ ] means 'k' main memory blocks are transferred into  $k \bmod N$  Cache memory line numbers.

Assume we have 4 Cache lines in Cm and 16m blocks no 11 (B<sub>11</sub>)

(i) Then mm block no 11 (B<sub>11</sub>) are mapped into  $11 \bmod 4 \rightarrow$  Cm line no '3'.

(ii) If mm block = 7 (B<sub>7</sub>) then they are placed into  $7 \bmod 4 \rightarrow$  Cm line no '3'.

i.e here block no B<sub>11</sub> and Block no 7 are mapped into Cache line number 3. But any one block present in Cache line no 3.

Eg: Pmm = 50 Blocks, Cm = 10 lines

$$TAG = \frac{Pmm}{Cm} = \frac{50}{10} \rightarrow 5:1$$

$\rightarrow$  5 from Blocks are Fighting / mapped to one Cache system line. But out of 5 only Only one mm block present in cache at a time.

Eg: Pmm = 16 Block, Cm = 4 lines

Block No (0, 4, 8, 12)  $\xrightarrow{\text{Mapped to}}$  Cache line '0'

Block No (1, 5, 9, 13)  $\xrightarrow{\text{}}$  Cache line '1'

Block No (2, 6, 10, 14)  $\xrightarrow{\text{}}$  Cache line '2'

Block No (3, 7, 11, 15)  $\xrightarrow{\text{}}$  Cache line '3'

Q10 Why before mapping we organise b/w the memory?

Eg: Main Memory size = 16 B, Cache Memory = 8 Byte, Block Size = 2 Byte.

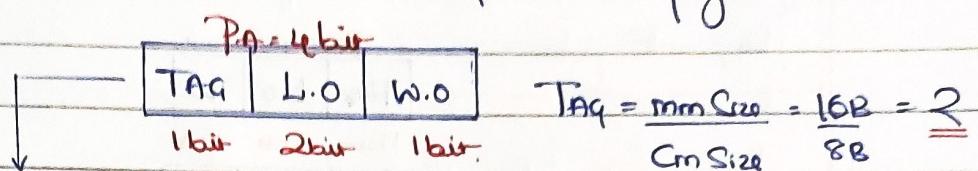
Ans

Before Mapping:  $\xleftarrow{\text{Cache = 3 bit}} \begin{array}{|c|c|} \hline L_0 & W_0 \\ \hline 2bit & 1bit \\ \hline \end{array}$

Main Memory:  $\xleftarrow{\text{PA = 4 bit}} \begin{array}{|c|c|} \hline TAG & W_0 \\ \hline 3bit & 1bit \\ \hline \end{array}$

Here TAG, bits tell the number of main memory blocks.

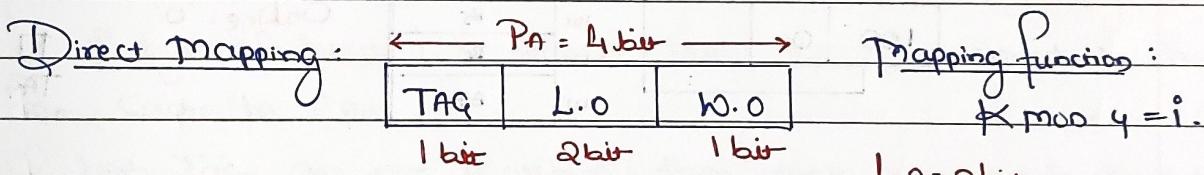
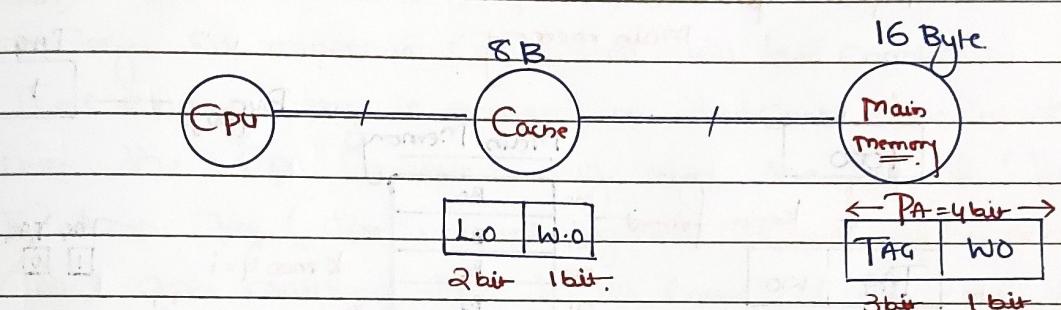
After Mapping: Cache controller interprets the physical address as:



→ Here TAG bit is 1 → 2 → 2 mm blocks

are mapped for one Cache line but only one mm block can be present in any one Cache line {2:1}

→ This TAG bit is used to check (Identify) whether there is Cache hit or Cache miss.



Consider the following program:

T1: MOV R0 [0000]

T2: MOV R1 [1000]

T3: ADD R0 R1

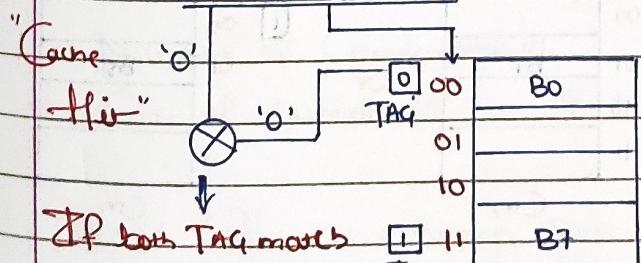
Cpu Generated Request Refer to Cache

J1: D000

$\xleftarrow{\text{PA = 4 bit}} \text{ Direct mapped Cache}$

$\xleftarrow{\text{PA = 4 bit}} \text{ (Known format)}$

1bit 2bit 1bit

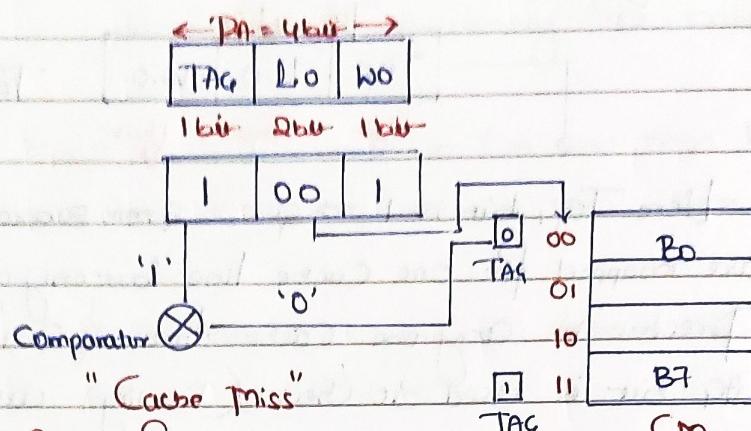


JP both TAG match

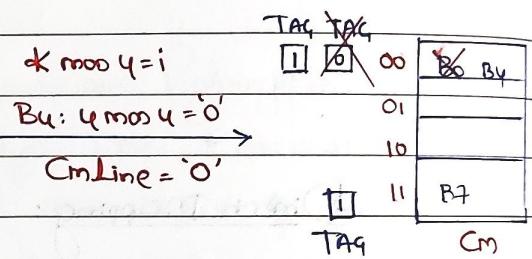
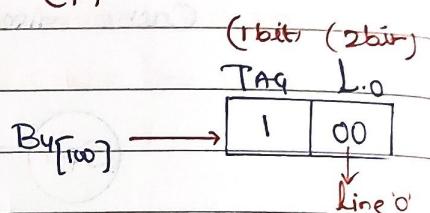
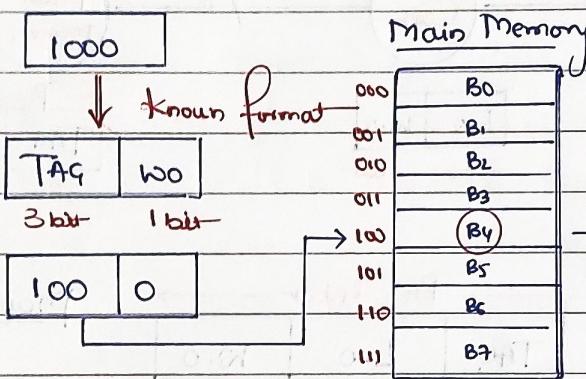
i.e Cache hit

then data is given from Cache to CPU directly.

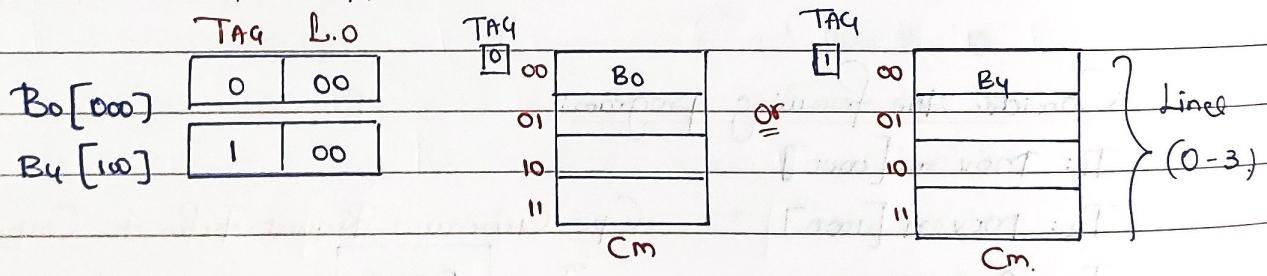
Q: 1000



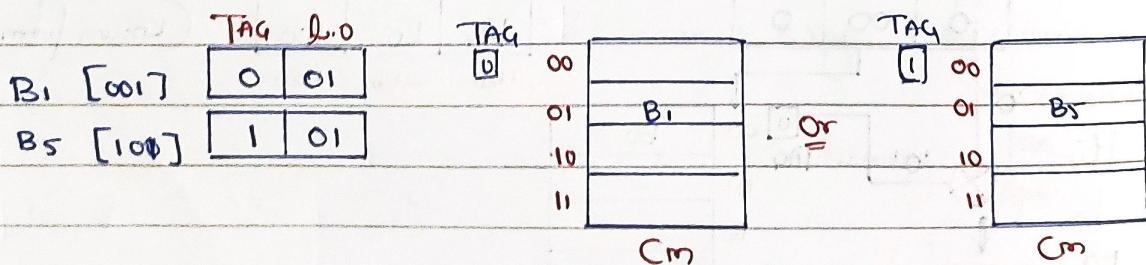
Reference forwarded to Main memory.



(i) B0 and B4 mapped to CM line '0'



(ii) B1 and B5 mapped to CM line '1'



(iii) B3 and B7 mapped to CM line '3'

(iv) B2 and B6 mapped to CM line '2'

$$\text{Tag Memory} = \frac{\text{Size}}{\text{# Lines}} \times \text{Tag bits} \quad \left\{ \begin{array}{l} \text{Depends upon the} \\ \text{mapping technique} \end{array} \right\}$$

### Complete Working:

1. CPU Generated request firstly refer to the Cache memory.
2. (Directly Mapped Cache) Cache controller interpreted the physical address as:
 

Physical Address		
TAG	Line Offset	Word Offset
3. Line offset field is directly connected to the address logic of Cache memory so respective Cache line will be enabled.
4. The existing tag is enabled in the Cache line is compared with the CPU Generated tag with the help of Tag comparison.
5. If both Tag (CPU Generated tag and ~~existing tag on enabled Cache line~~ Tag) are matching then it is said Cache hit (i.e. data is present in the Cache). So based on word offset (~~0<sup>th</sup> Byte of block or 1<sup>st</sup> byte of block~~) respective data, respective word is transferred from Cache to CPU.
6. If both Tag are not matching then that operation is called Cache miss, then reference is forwarded to main memory.
7. According to MM format, respective MM block is enabled and transfer from MM to CM with the help of direct mapping function [ $k \bmod N$ ] into particular Cache line the Cache to CPU.
8. Here Cache Controller maintains some extra bits for each Cache line and that stored in Tag memory / memory.

$$\text{Tag memory} = \text{No. of CM lines} \times \text{Tag bits}$$

→ Number of Comparators = 1.

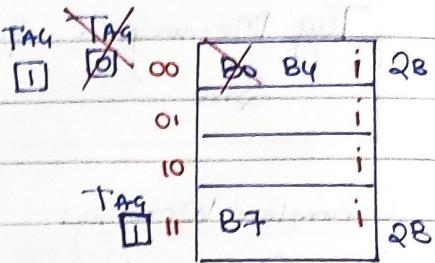
→ Size of Comparator = Tag bits.

Note: \* In a mapping process complete 1mm block is transferred from main memory to Cache memory.

\* But respective word/byte of that Particular block which is demanded by the CPU is given tag from Cache to CPU.

Consider the following Programs:

- I<sub>1</sub> : Mov [0000] 0<sup>th</sup> Byte of B<sub>0</sub>
- I<sub>2</sub> : Mov [0001] 1<sup>st</sup> Byte of B<sub>0</sub>
- I<sub>3</sub> : Mov [1000] 0<sup>th</sup> Byte of B<sub>4</sub>
- I<sub>4</sub> : Mov [1001] 1<sup>st</sup> Byte of B<sub>4</sub>.



I<sub>1</sub>: [0000] : Block B<sub>0</sub> (0<sup>th</sup> Byte)  $\rightarrow$  present in Cache : Cache hit.

I<sub>2</sub>: [0001] : Block B<sub>0</sub> (Cache Byte)  $\rightarrow$  present in Cache. Due to locality of reference.

I<sub>3</sub>: [1000] : B<sub>4</sub>  $\Rightarrow$  Cache miss so bring B<sub>4</sub> from main memory to Cache Memory.

$\hookrightarrow$  then 0<sup>th</sup> Byte of given to CPU (MM to Cache [B4] then Cache to CPU 0<sup>th</sup> Byte of B<sub>4</sub>)

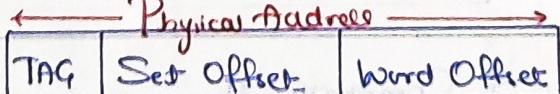
I<sub>4</sub>: [1001] : B<sub>4</sub> B<sub>4</sub> (1<sup>st</sup> byte of B<sub>4</sub>): Cache hit. Due to locality of reference.

### Disadvantages of Direct Mapping:

- \* In direct mapping each Cache line is able to hold only 1 main memory block at a time. So number of conflict miss increase.
- \* If CPU is referring the multiple blocks which are mapped into the same Cache line number (eg: B<sub>0</sub>, B<sub>4</sub>, B<sub>0</sub>, B<sub>4</sub>  $\Rightarrow$  Cm line no '0') then the number of conflict miss will be very high. Even other Cache line are empty. {Thrashing}.

### Set Associative Mapping:

Set associative Cache controller, interpreter the CPU generated request as follows:



$\rightarrow$  Word Offset =  $\log_2$  Block Size.

$$\# \text{Lines} = \frac{\text{Cm Size}}{\text{Block Size}}$$

$\rightarrow$  No. of Sets =  $\# \text{Lines}$   
N-way.

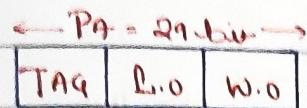
$\rightarrow$  Set offset =  $\log_2$  Sets

$\rightarrow$  TAG = Physical address - (S.O + W.O)

Q1. Consider a Direct Mapped Cache if the size of the Cache memory is 512KB and main memory 512MB and Cache line is 64KB then calculate the number of bits required for:

- (i) PA      (iii) Tag    (iv) L.O    (ii) W.O.

Ans.



10 bit    3 bit    16 bit.

$$P_A = 512\text{MB} = 2^{29}B \rightarrow P_A = 29 \text{ bit}$$

$$\text{Line Size} = 64\text{KB} = 2^{16}B \rightarrow W.O. = 16 \text{ bit}$$

$$\# \text{Lines} = \frac{\text{CnSize}}{\text{Block Size}} = \frac{512\text{KB}}{64\text{KB}} = \frac{2^9}{2^{16}} \text{B}$$

$$\text{Tag memory Size} = \# \text{Lines} \times \text{Tag bit} \quad \# \text{Lines} = 8$$

$$= 8 \times 10 \quad L.O. = 3 \text{ bit}$$

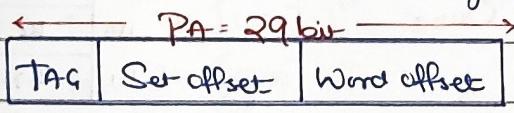
$$= \underline{\underline{80 \text{ bit}}}$$

$$\text{Tag} = P_A - (L.O + W.O)$$

$$= 29 - (3 + 16)$$

$$= \underline{\underline{10 \text{ bit}}}$$

Above Question uses 2-way Set associative mapping:



11 bit    2 bit    16 bit.

$$\# \text{Lines} = \frac{\text{CnSize}}{\text{Block Size}} = \frac{512\text{KB}}{64\text{KB}} = \frac{2^9}{2^{16}} \text{B}$$

$$\begin{aligned} \text{Tag memory Size} &= \# \text{Lines} \times \text{Tag bit} \\ &= 8 \times 11 \\ &= \underline{\underline{88 \text{ bits}}} \end{aligned}$$

$$\text{Set offset} = 2 \text{ bit}$$

$$\text{Tag} = P_A - (\text{S.o} + \text{W.o})$$

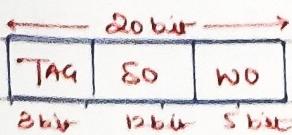
$$= 29 - (2 + 16)$$

$$= \underline{\underline{11 \text{ bit}}}$$

$$\begin{aligned} \text{Tag memory Size} &= \# \text{set} \times \text{Block per set} \\ &\times \text{Tag bit} \\ &= 4 \times 2 \times 11 \text{ bits} \\ &= \underline{\underline{88 \text{ bits}}} \end{aligned}$$

Q2. Consider a 2-way Set associative Cache Size = 256KB, Line size = 32B, MM = 1MB, then what is the Set number of physical address (ABCDE)<sub>16</sub>?

Ans.

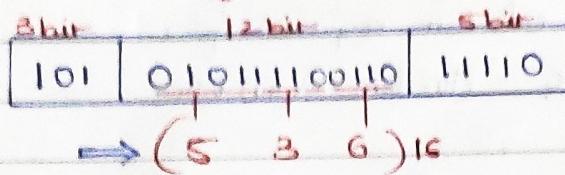


8bit    12bit    5bit

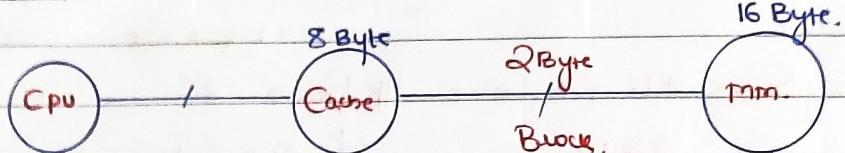
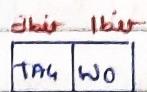
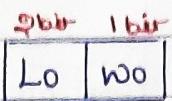
$$\# \text{Lines} = \frac{256\text{KB}}{32\text{B}} = \frac{2^{18}}{2^5} = 2^3 \text{ lines}$$

$$\# \text{Set} = \# \text{Lines} = \frac{2^8}{2^1} = 2^7$$

N-way  $\frac{2^8}{2^1}$  : So = 12 bits



Direct-Mapping Example.



PA = 4 bits			TAG		B0		B1		B2		B3		B4		B5		B6		B7							
TAG	LO	WO	00	00	i	2B	000	000	B0	001	001	B1	010	010	B2	011	B3	100	100	B4	101	B5	110	B6	111	B7
1 bit	2 bits	1 bit	01	01	i	2B	001	001	B1	010	010	B2	011	011	B3	100	B4	101	101	B5	110	B6	111	B7		
TAG	10				i	2B	011	011	B3	100	100	B4	101	101	B5	110	B6	111	111	B7	111	111	111	2B		
	1	11			B7	i	2B	100	B4	101	101	B5	110	110	B6	111	B7	111	111	B7	111	111	111	2B		
Cache																										
MM																										

Q3. MM = 16 B, CM = 8 B, Block Size = 2B, 2 way set associative mapping  
 then calculate bits for PA, TAG, S.O, W.O and Tag memory size.

Ans

$\leftarrow PA = 4 \text{ bits} \rightarrow$

TAG	S.O	W.O
2 bits	1 bit	1 bit

$$\# \text{Line} = \frac{\text{CM Size}}{\text{Block Size}} = \frac{8B}{2B} = 4 \text{ lines}$$

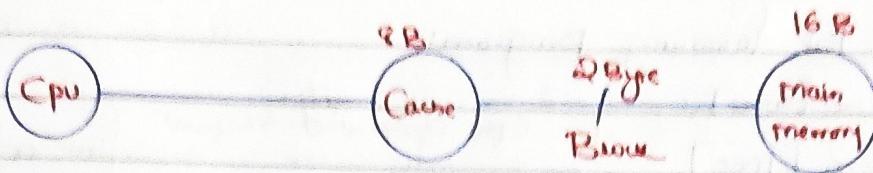
$$\# \text{Set} = \frac{\# \text{Line}}{\text{N-way}} = \frac{4}{2} = 2 \text{ set.}$$

Tag memory = 4x2

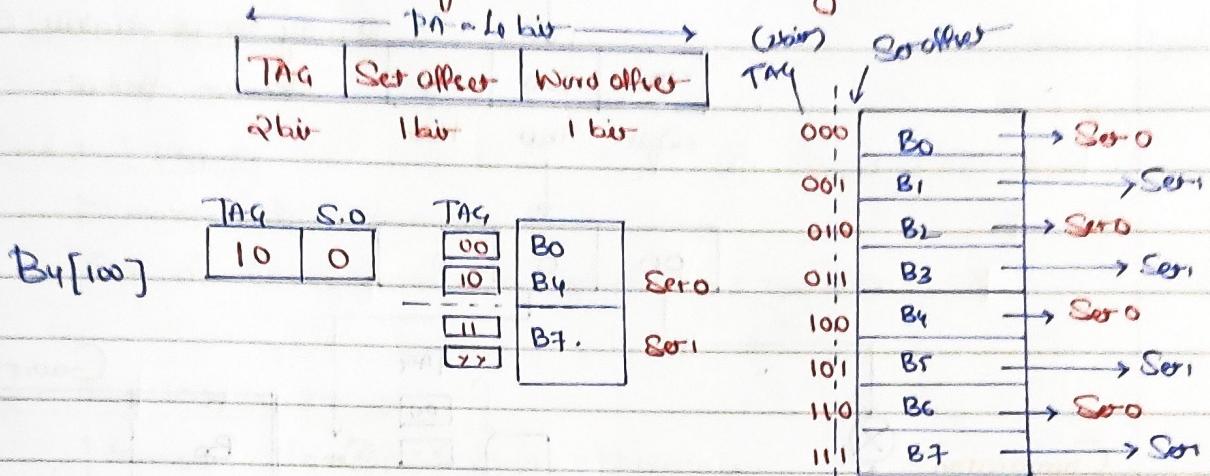
Size = 8 bits. (Ans.)

Set offset = 1 bit.

\* Here Tag = 2 bits i.e.  $2^2 \Rightarrow$  4 mm blocks are mapped to any one cache set.



2 way Set association mapping.



Advantages of Association Mapping

I<sub>1</sub>: Mov r0 [0000] 0<sup>th</sup> Byte of Block B0 I<sub>1</sub>: Hit (Cache hit)

I<sub>2</sub>: Mov r1 [1000] 0<sup>th</sup> " " B4 I<sub>2</sub>: miss (Cache miss)

I<sub>3</sub>: Mov r2 [0001] 1<sup>st</sup> " " B0 I<sub>3</sub>: hit

I<sub>4</sub>: Mov r3 [1001] 1<sup>st</sup> " " B4 I<sub>4</sub>: hit

I<sub>5</sub>: Mov r4 [0000] 0<sup>th</sup> " " B0 I<sub>5</sub>: hit

I<sub>6</sub>: Mov r5 [1000] 0<sup>th</sup> " " B4 I<sub>6</sub>: hit

6 Access: Total: 1 miss

5 hits

Note:

\* 2 way Set association: In one Cache memory Set we can store a mm block in 2 ways [2mm Blocks].

\* N-way Set association: In one Cm Set, we can store mm Block by N ways.

Set association mapping functions

Cache set = mm Request mod #Sets in Cache  
address

or

K mod S : i

K: mm Block number

i: Cache Set number.

S: # Cache Set

Consider the following program:

T1 : mov r0 [wo]

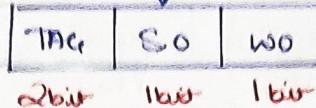
T2 : mov r1 [1000]

T3 : add r0, r1

In CPU generated request first refer to cache.

T1: 0000

2 way set associative Cache known format.



Comparator.



Tag mismatch!



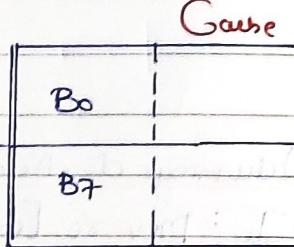
TAG

00

xx

Sero

Set-1

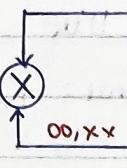


Requesting data Cache hit

given from cache to CPU.

T2: TAG S0 W0

Comparator.

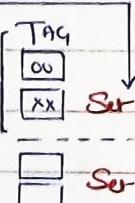


MUX

No tag match so

Cache misses, S0 reference forwarded to main memory.

Cause



TAG

00

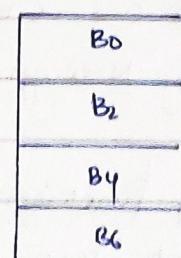
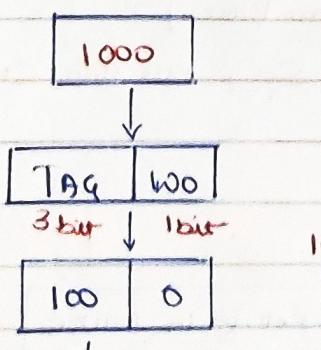
xx

Sero

Set-1

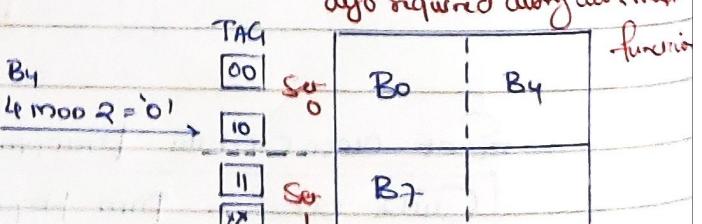
B0

B7



Main Memory.

→ Now cache is full replaced also required along with main function



→ In direct mapping,

- Hit Latency = Latency of tag comparator.

→ In set associative mapping:

Hit Latency = Hit latency of tag comparator + latency of multiplexer.

Direct Mapping:

Number of Comparator = 1

Size of Comparator = # bits in Tag (Tag bits)

Set associative mapping:

Number of Comparator = N-way

Size of Comparator = # bits in Tag (Tag bits)

Important Points about Set-Associative:

- \* # Line = L & N way Set associative  $\Rightarrow$   $\text{Sets} [S] = \frac{L}{N}$
- \* If  $N=1$ : Direct Mapping
- \* If  $N=L$ : Association Mapping; ( $S=1$ ) i.e. Only 1 set
- \* Set associative mapping follows direct mapping as well as Set associative mapping. Using mapping functions as:

Within the Set MM Block can be placed anywhere.

Note: In the Set associative mapping replacement algorithm is required along with the mapping function when Cache Set is full.

Tag bits in N ways = Tag bits in direct mapping +  $\log_2 N$

e.g.: Consider a direct mapping if the size of Cache memory is 512KB and main memory 512MB and Cache line size is 2KB then calculate the number of bit required for Tag?

Ans.

PA = 29 bit		
Tag	L.O	W.O
10 bit	8bit	11bit

$$\# \text{Lines} = \frac{\text{Cm Size}}{\text{Line}} = \frac{512 \text{KB}}{2 \text{KB}} = 256$$

$$L.O = 8 \text{bit}$$



eg: Consider a 2-way Set associative mapping, if the size of the Cache memory is 512KB and main memory 512MB and Cache line size is 2KB then calculate the number of bits required for Tag?

Ans.

PA = 29 bit			# Lines = 2 <sup>8</sup>
Tag	S.O	W.O	# Set = $\frac{2^8}{2^1} = 2^7$
11bit	7bit	11bit	

$$29 = (7+11)$$

$$TAG = 11 \text{bit}$$

### 2 Way Set Association

$$\text{Tagbit} = \text{Tagbit in } [ \text{log}_2 \text{N-way}] \text{ direct map}$$

$$= 10 + [\log_2^2]$$

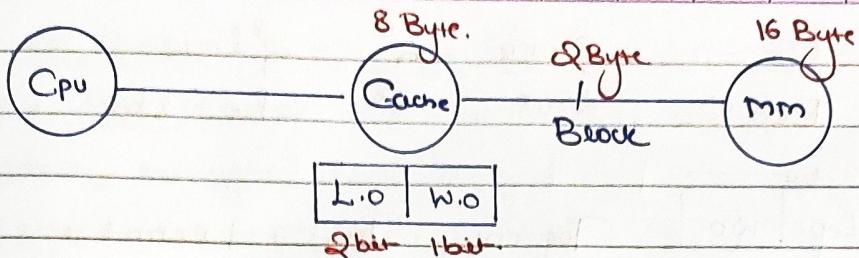
$$= 11 \text{bits}$$

### Associative Mapping:

- \* In this no mapping function is used to transfer the data from MM to CM.
- \* Any block of main memory can be placed anywhere in cache memory.
- \* No conflict miss.
- \* This cache is designed without address Called as Content Addressable memory.
- \* In associative mapping more tag bits are required along with more memory size.
- \* In associative mapping more hardware required, Expensive (N tag comparators, here N is number of lines, for each comparator is required.)
- \* In associative Cache design, Counter Sequence is used to map the data, means any main memory block can be mapped to any cache memory line in a sequence.
- \* In associative mapping physical address is interpreted as

Tag	W.O
-----	-----

eg:



### Direct Mapping

PA = 4 bit
Tag   L.O   W.O
1bit 2bit 1bit

PA = 4 bit
Tag   L.O   W.O
2bit 1bit 1bit

### Associative / Fully Associative

PA = 4 bit
Tag, W.O
3bit 1bit

$$K_{\text{moo}} N = 1$$

$$K_{\text{moo}} q = 1$$

$$\text{Tag memory} = \# \text{lines} \times \text{Tag bit} = 4 \times 2 = 8 \text{ bits}$$

$$= 4 \times 1 = 4 \text{ bits}$$

$$= 4 \times 1 = 4 \text{ bits}$$

$$K_{\text{moo}} S = i$$

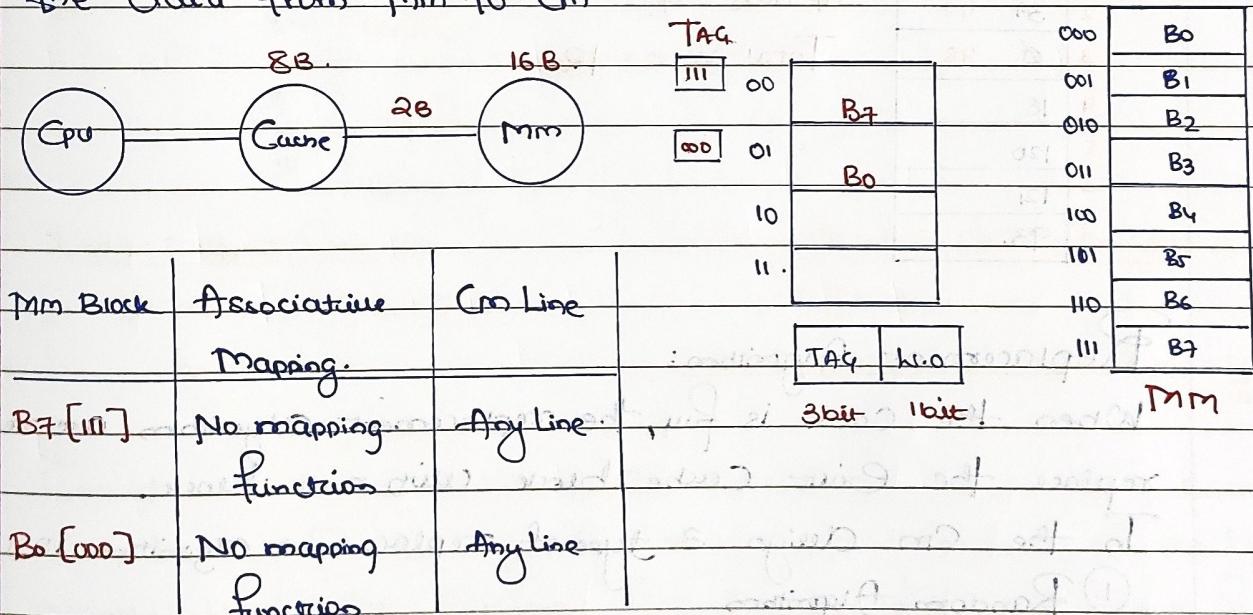
$$K_{\text{moo}} Q = 1$$

$$\text{Tag memory} = 4 \times 2 = 8 \text{ bits}$$

No mapping function

$$\text{Tag memory} = 4 \times 3 = 12 \text{ bits}$$

In associative mapping no mapping function is used to transfer the data from MM to CM



### Direct Mapped Cache

$$\text{Cache Size} = 64 \text{ KB}$$

$$\text{Block Size} = 64 \text{ Word}$$

$$= 64 \times 64 \text{ Byte}$$

$$= 256 \text{ Byte}$$

$$MM = 4 \text{ Byte} = 2^{32} \text{ Byte}$$

$$W.O = 8 \text{ Bit}$$

$$\# \text{Lines} = \frac{\text{Cm Size}}{\text{Block Size}} = \frac{64 \text{ KB}}{256 \text{ B}} = \frac{2^{10}}{2^8} = 2^2 \text{ Lines}$$

$$= 4.0 = 8 \text{ bits}$$

$\leftarrow P_A = 32 \text{ bits} \rightarrow$		
TAG	L.O	W.O
16 bits	8 bits	8 bits

Tag entry = Tag bits + extrabit  
 $= 16 + 1 + 1$

$$= 32 - (8+8)$$

$$= 16 \text{ bits}$$

$$\begin{aligned} \text{Tag memory Size} &= \# \text{Lines} \times \text{Tag bits} \\ &= 2^8 \times 16 \\ &= 256 \times 16 \\ &= 4096 \text{ bits} \end{aligned}$$

Q1. Consider fully associative Cache Consist 8 Block And 1mm excess  
 Contains 128 Block and request made by the Cpu : 119, 84, 37, 0, 16, 0, 84, 120, 121, 93, 37, 0, 43, 39, 47, 48.

Calculate the # of Compulsory And Capacity miss?

Ans.

0	119	43
1	84	39
2	37	47
3	0	48
4	16	
5	120	
6	121	
7	93	

stall : 0, 84, 37, 0

$$\# \text{miss} = 4$$

$$\text{Total miss} = 12.$$

### Replacement Algorithm:

When the Cache is full, the replacement algorithms are required to replace the Exit Cache block with new block.

In the Cm design 3 type of replacement algorithms is used.

- ① Random Algorithm
- ② FIFO Algorithm
- ③ LRU replacement.

In the random algorithm, Any cache block can be replaced based on the random selection.

## Types of Miss:

In the CM design 3 types of misses are present:

① Compulsory miss - (cold start miss / first reference miss)

This miss will occur when the very first reference to the cache itself is a miss.

② Capacity miss - This will occur only when the capacity is full.

③ Conflict miss - (Collision miss / reference miss)

This miss will occur when the too many blocks are placed into same cache line or same cache set.

Q2 Consider 4 block cache memory (initially empty) with the following MM blocks referenced: 7, 8, 10, 15, 7, 8, 16, 7, 8, 10 Identify hit ratio using:

Ans.

7  
8  
10 }  
15  
Compulsory miss

MM request = 10

Cache line = 4.

7 → conflict  
8 → conflict  
16 → compulsory miss  
7 {  
8 }  
10 }  
Conflict.

Direct mapped Cache:

$$\text{Hit ratio} = \frac{3}{10} = 0.3$$

2 Way Set Association with LRU

$$\# \text{ lines} = 4$$

$$\text{MM request} = 10 \quad \text{Hit Ratio} = \frac{4}{10} = 0.4$$

Q3. Consider a small two-way set associative Cache memory, consisting of 4 blocks. For choosing the block to be replaced, use the least recently used (LRU) scheme. The number of cache misses for the following sequence of block address is 8, 12, 10, 12, 8.

Ans

# CM Lines = 4      2 way Set association

$$\# \text{ sets} = 4/2 = 2 \quad k \bmod 2 \text{ in } 1 \quad 12 \bmod 2 = 0 \rightarrow \text{hit}$$

$$8 \bmod 2 = 0 \rightarrow M$$

$$12 \bmod 2 = 0 \rightarrow M$$

0, 1, 2, 0 → M (cause set is full, apply LRU replace algo.)

Set 0	8	8
Set 1	12	

∴ Cache miss = 4

How to reduce all the miss:

- ① Compulsory miss: We can reduce the compulsory miss by increasing the block size (large block size).
- ② Capacity miss: Capacity miss occurs when cache is full. If we increase the cache size then capacity miss can be reduced.
- ③ Conflict miss: Conflict miss can be reduced by increasing the associativity.

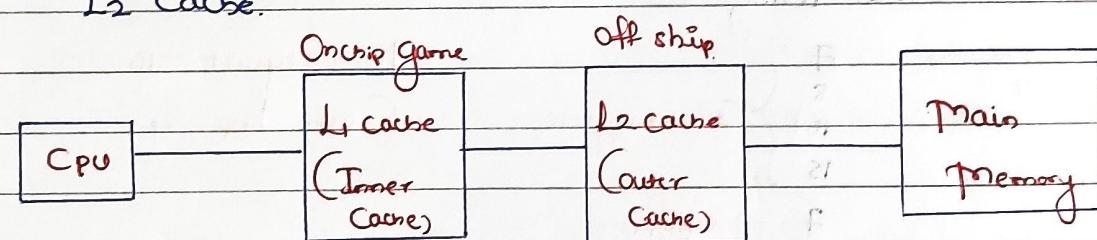
Multilevel Cache:

2 Level Cache:

- L1 Cache
- L2 Cache

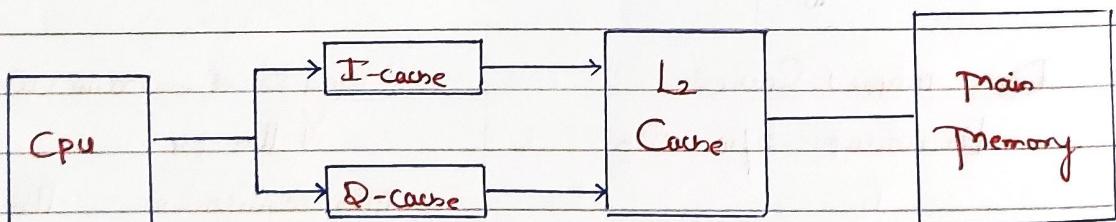
L1 Cache

D Cache



Size  $L_2 > L_1$

Speed  $L_1 > L_2$



I-Cache: Instruction Cache

D-Cache: Data Cache

Working Principle of multilevel Cache:

- ① Principal of Inclusion: Data present in L1 Cache must be present in L2 Cache.
- ② Principal of Exclusion: Data present in L1 Cache must be different from L2 Cache.

Why Multi-level Cache is used:

Level 1 Cache access = 10nsec  
Time

Main memory access time = 300ns.

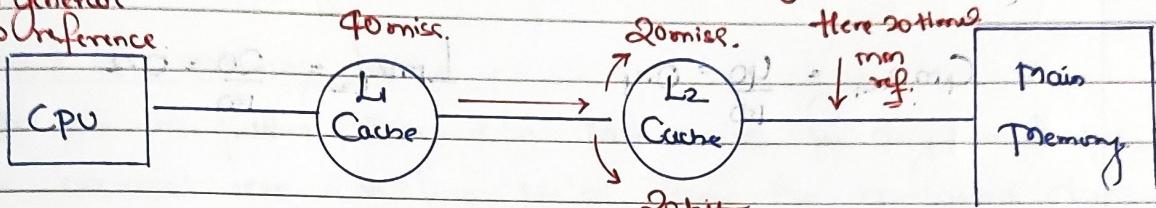
L<sub>1</sub> Cache access time = 10 nsec

L<sub>2</sub> cache access time = 15 nsec

From access time = 300ns

So Tag will reduce

CPU generates  
100 reference



and wait until t<sub>ref</sub> + t<sub>access</sub> 20 time used.

If only L<sub>1</sub> Cache then if there is miss in Level 1 Cache then directly refer to main memory

(Here higher time in main memory access so Tag increase).

But in 2 level Cache, if miss is Level 1 Cache then tag miss forwarded to Level 2 Cache (so some hit is Level 2) then after Level 2 there is a miss then we go to main memory. So Tag can be reduced (because L<sub>1</sub> and L<sub>2</sub> Cache access time very low).

### Multi Level Cache

- \* To reduce the miss penalty multi-level Caches are used in the System design.
- \* The number of cycles required to transfer the data from higher levels to L<sub>1</sub> due to miss operations is called as miss penalty.

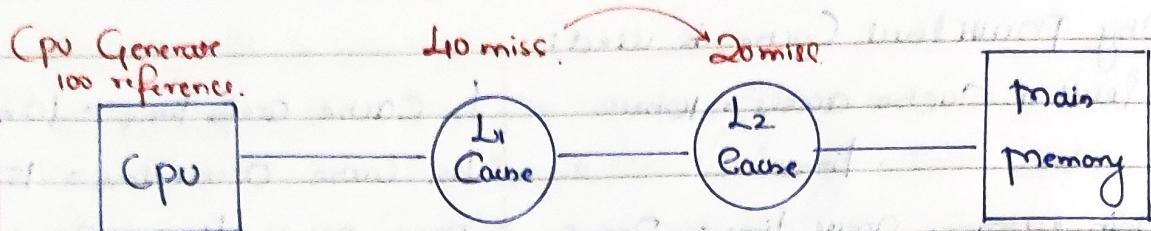
Miss Rate

Local Miss Rate (LmR)

$$LmR = \frac{\# \text{ misses in Cache}}{\# \text{ access to that cache}}$$

Global Miss Rate (GmR)

$$GmR = \frac{\# \text{ misses in the Cache}}{\text{Total } \# \text{ CPU generated references}}$$



$$LMR_{L1} = \frac{40}{100} = 0.4$$

$$LMR_{L2} = \frac{20}{40} = 0.5$$

$$GMR_{L1} = \frac{40}{100} = 0.4$$

$$LMR_{L2} = \frac{20}{100} = 0.2$$

### Access time in Multilevel Cache:

Level 1:

$$T_{avg} = h_1 t_1 + (1-h_1) (t_m + t_1)$$

$$T_{avg} = t_1 + (1-h_1) t_m$$

Level 2:

$$T_{avg} = h_1 t_1 + (1-h_1) h_2 (t_2 + t_1) + (1-h_1) (1-h_2) (t_m + t_2 + t_1)$$

Or

$$T_{avg} = t_1 + (1-h_1) \{ t_2 + (1-h_2) t_m \}$$

hit time in level 1, miss rate of level 1, hit time of level 2, miss rate of level 2, access time of level 2

### Updating Technique:

- \* CPU performs read or write operation on Cache memory.
- \* When the data (mm block) is available in Cache then it is known as read hit and write hit.
- \* If the data is not available in the Cache then it is called read miss and write miss then we have to perform Read allocate and write allocate.

Read Allocate: When there is a read miss in the Cache then transferring that data from main memory to Cache memory is

Called read allocate.

Write Allocate: When there is a write miss in the Cache then transferring that data (mm block) from main memory to Cache memory is called write allocate.

Need of updating technique:

Assume in future there is a Cache miss for any mm block, which are mapped into Cm line No'0. Then the updated data (of block B0) is replaced by new block. (All the updation goes far).

That means coherent work done by the CPU and load into the Cache or lost. Because, that block is replaced by the new main memory block.

Now in future if we bring Block B0 then we get the old value.

Updating Technique:

① Write Through: In the write through both Cache memory and main memory update simultaneously so no problem of Coherence (No inconsistency).

$$tw = \max \left\{ \begin{array}{l} \text{word update time in Cache, word update time in MM} \\ \text{time} \end{array} \right\}$$

$$\text{Tavg read} = hr \cdot tc + (1-hr) (tm + tc)$$

↓      ↓      ↓      ↓      ↓  
 Read hit    Read data    Read miss    Read allocate    Read data from CM.

$$\text{Tavg write} = bw \cdot tc + (1-bw) (tm + tw)$$

↓      ↓      ↓      ↓  
 Write hit    Word update    Write miss    Write allocate    Word update

$$\text{Tavg wr} = \frac{\% \text{ of read}}{\text{Operations}} \times \text{Tavg read} + \frac{\% \text{ of write}}{\text{Operations}} \times \text{Tavg write.}$$

$$\eta_{wr} = \frac{1}{T_{avg\ wr.}}$$

Disadvantage in WT: If there is a frequent update in cache, then WT give worst performance because every update, we have to update Pm and Cm block.

Simultaneous Access:

$$T_{avg\ read} = br + (1-br)tm$$

$$T_{avg\ write} = tw$$

- ② Write Block Update Technique: In the WB updating technique main memory update later, Only cache memory update. So in WB Cache memory update and Pm update later so coherence is present).

But this coherence does not create any problem

- \* Because in WB, each Cache line maintains some extra bit (i.e. valid/invalid, dirty bit (modified bit)) to maintain the status of Cm Block.
- \* If the Cm Block is updated then modified bit (dirty bit) is set to 1. If block is not updated then clean bit is set (dirty bit=0).
- \* Before replacing the Cm Block, we check the status of the block.
- \* If a block is an updated block (modified bit set) then firstly this updated clean cache block into main memory (higher memory) then replace the Cm block with any new (requested) block.
- \* If block is not updated block then directly replace with new main memory (requested block).

$$T_{avg\_read} = brtc + (1-br) \left\{ \begin{array}{l} \cdot \text{data} \quad (tr + tm + tc) + \cdot \text{clean} \\ \text{modified} \quad \text{bit} \\ (bm + tc) \end{array} \right\}$$

$$T_{avg\_write} = bwtc + (1-br) \left\{ \begin{array}{l} \cdot \text{data} \quad (tw + tm + tc) + \cdot \text{clean} \\ \text{modified} \quad \text{bit} \\ (tm + tc) \end{array} \right\}$$

$$T_{avg\_wb} = Fr \times T_{avg\_read} + Fw \times T_{avg\_write}$$

$$\eta_{wb} = \frac{1}{T_{avg\_wb}}$$

Fr = % of read operations

Fw = Frequency (%) of write operations

- Q1. Let WB and WT be two set associative Cache organizations that use LRU algorithms for cache block replacement. WB is a write back Cache and WT is a write through Cache. Which of the following statements is/are true or false?

- Ans.
- A. Each Cache block in WB and WT has a dirty bit.
  - B. Every write hit in WB leads to data transfer from Cache to main memory.
  - C. Eviction of a block from WT will not lead to data transfer from Cache to main memory.
  - D. A read miss in WB will never lead to eviction of a dirty block from WB.

### While Cache Accessing:

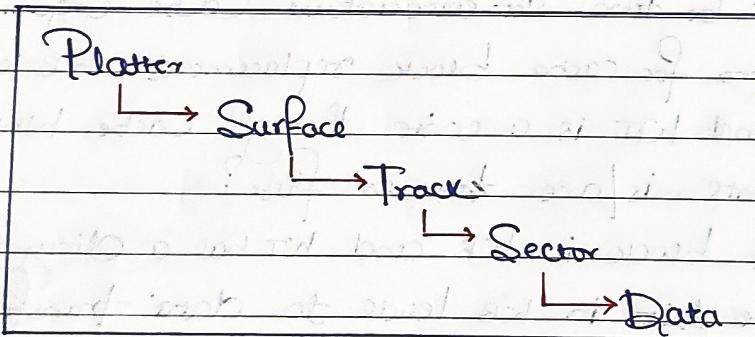
1. Physical Index Cache: MM address is used in address mapping.
2. Virtual Index Cache: VM address is used in address mapping.

(VIRT)

→ Virtual Index Virtual Tag.

## SECONDARY MEMORY & I/O INTERFACE

- \* Magnetic disks provide the bulk of Secondary Storage for modern Computer Systems. Conceptually, disks are relatively simple.
- \* Each disk platter has a flat circular shape, like a CD. Common platter diameters range from 1.85 to 3.5 inch.
- \* The two surfaces of a platter are covered with magnetic heads. Are attached to disk arms that move all the heads as a unit.
- \* The surface of a platter logically divided into circular tracks, which are subdivided into sectors.
- \* The set of tracks that are at one arm position makes up a cylinder. There may be thousands of concentric cylinders in a disk drive, and each track may contain hundreds of sectors. The storage capacity of common disk drives is measured in disk drives.



### Moving-Head Disk mechanism

- \* Each platter has 2 surfaces
- \* Read/write head with each surface to read/write data.
- \* These read/write head are connected with the arms.
- \* All these arms are connected to arm assembly and rotate at the same time.

e.g: If there are 20 surfaces then:

# Read/write head : 20

# Arm : 20

# Arm Assembly : 1

Note: - All read/write head (Arm) rotating/moving at the same time but only 1 R/W head performing read/write operation.

Q1. Consider a Disk which has 16 platters. Each platter has two surfaces. Every surface has 1K Track. Each track is further divided into 512 sectors. And every sector can store the 8KB data. Then calculate.

(i) What is the capacity of the disk?

Ans: 16 platters, each platter has 2 surfaces.

Total number of surfaces =  $32 [16 \times 2]$  surfaces

# Tracks per surface = 1K, total = 32K tracks

# Sectors/tracks = 512

Total # sectors =  $32K \times 512$  sectors

Each Sector Capacity = 8KB

Disk Capacity =  $32K \times 512 \times 8KB$

= 128GB

(ii) How many bits are required to identify any particular sector of disk?

Ans: # bits to identify =  $16 \times 2 \times 1K \times 512$

Particular Sector =  $2^4 \times 2^1 \times 2^{10} \times 2^9$

=  $2^{24}$

= 24 bit

Q2. Consider a Disk Pack with 16 surfaces, 128 tracks per surface and 256 sectors per track. 512 bytes of data are stored in a bit serial manner in a sector. The capacity of the disk pack and the number of bits required to specify a particular sector in the disk are respectively

Ans: Disk Capacity =  $16 \times 128 \times 256 \times 512$  Byte

$$= 2^4 \times 2^7 \times 2^8 \times 2^9 \text{ Byte}$$

$$= 2^{28} \text{ Byte}$$

$$= \underline{\underline{256MB}}$$

# bits required =  $16 \times 128 \times 256$

to identify particular =  $2^4 \times 2^7 \times 2^8$

sector =  $2^{19}$

= 19 bit

## Disk I/O Operations:

1. Seek time: The amount of time taken to move the read/write head from its current position to the desired track is known as Seek time.
2. Rotational Latency: The amount of time taken to rotate the disk when the read/write head comes out to exact position.  
Rotational Latency is considered as  $\frac{1}{2}$  Rotation Time
3. Transfer time: The amount of time taken to transfer the required data.
4. Transfer rate: The number of bytes read for unique line is called as transfer rate of disk.

- \* The read/write head can never be outside the track it can be pointing to any particular track of the surface.
- \* Read/write head can move in forward and backward direction and disk in one direction (either clockwise or anti-clockwise).

eg: 600 Rpm. (600 rotation per minute)

$$1 \text{ Rotation} = \frac{60}{60} = \frac{1}{10} \text{ sec} \times \frac{1000}{1000} = 100 \times 10^3 \text{ sec} \\ = 100 \text{ msec}$$

$$1 \text{ Rotation time} = \frac{1}{10} \text{ sec or } 100 \text{ msec.}$$

$$\text{Avg rotational latency} = \frac{1}{2} \times 100$$

$$= 50 \text{ msec}$$

eg: Assume 1' Rotation time  $1/6$  sec, 1 Track Capacity = 4KB then calculate data transfer time for x Byte and data transfer rate?

$$(i) \text{ Data Transfer time: } 4\text{KB} \longrightarrow \frac{1}{6} \text{ sec}$$

$$1 \text{ Byte} \longrightarrow \frac{1}{6 \times 4\text{KB}} \text{ sec}$$

$$x \text{ Byte} \longrightarrow x \text{ Byte} / 6 \times 4\text{KB} \text{ sec.}$$

(ii) Data Transfer Rate:  $4 \text{ KB} \rightarrow \frac{1}{6} \text{ sec}$

$\frac{1}{6} \text{ sec} \rightarrow 4 \text{ KB}$

$1 \text{ sec} \rightarrow 4 \text{ KB} \times 6 \text{ sec}$

= 24 KB/sec

$$\boxed{\text{Disk Access Time} = \text{S.T} + \text{R.L} + \text{D.T.T} + \text{Overhead} \\ (\text{D.A.T})}$$

- Q1. Consider a disk system, which has an average seek time of 30ms and rotational rate of the disk is 360 RPM. Each track of the disk has 512 sectors, each of the size 512 Byte, then calculate.
- (i) Average Seek time
  - (ii) Average rotational latency
  - (iii) Data transfer time for 4 sectors (Continue)?
  - (iv) Data transfer rate?

Ans

360 RPM

$$\text{Avg. R.L} = \frac{1}{2} \times \frac{1}{6} = \frac{1}{12} \text{ sec}$$

$$1 \text{ Rotation} = \frac{60}{360} = \frac{1}{6} \text{ sec.}$$

$$= 0.083 \text{ sec}$$

$$\text{Avg. Seek time} = \underline{\underline{30 \text{ msec}}}$$

$$1 \text{ Track Capacity} = \# \text{ sectors/track} \times \text{each sector}$$

Capacity

$$= 512 \times 512 \text{ Byte} = 2^9 \times 2^9 \text{ B} = 2^{18} \text{ B}$$

$$= \underline{\underline{256 \text{ KB}}}$$

$$\begin{aligned} \text{Data Transfer Amount} &= 4 \text{ sectors} \rightarrow 4 \times 512 \text{ B} \\ &= 2^2 \times 2^9 \text{ B} \\ &= \underline{\underline{2 \text{ KB}}} \end{aligned}$$

Data Transfer Time:

$$256 \text{ KB} \rightarrow \frac{1}{6} \text{ sec}$$

$$2 \text{ KB} \rightarrow \frac{2 \text{ KB}}{6 \times 256 \text{ KB}} \text{ sec}$$

$$2 \text{ KB} = \frac{2 \text{ KB}}{6 \times 256 \text{ KB}} \text{ sec} \Rightarrow \frac{1}{768} \text{ sec.}$$

Data Transfer Rate:

$$256 \text{ KB} \rightarrow \frac{1}{6} \text{ sec}$$

$$\frac{1}{6} \text{ sec} \rightarrow 256 \text{ KB}$$

$$\begin{aligned} \text{In 1 sec} &\rightarrow 6 \times 256 \text{ KBps} \\ &= \underline{\underline{1536 \text{ KBps}}} \end{aligned}$$

Fixed Sector Capacity: Each sector (innermost and outermost sector) has fixed storage. (Storage density is variable)

Variable Sector Capacity: Each sector has variable storage. Innermost and outermost sector capacity is variable. (Storage density is fixed).

\* In Today's time we generally use "fixed sector capacity".

Need of Cylinder:

Assume that we have a very large file size (in GB) and 1 track size is 32 MB.

If we store a file on a disk then we can store 32 MB on one track then we have to go next track of the same surface so on storing data. So we have to suffer unnecessary very high seek latency (time).

So we store the data cylinder wise as we know each surface r/w head point to same address at a time.

So first we store {track '0' of surface '0'  
track '0' of surface '1'  
"  
"  
track '0' of surface '2000' }  
} 1 Seek time (latency)

With the help of cylinder we can store 2000 track data in one seek time.

\* Same track number in all the surfaces will form a cylinder.



- Number of platters = 2
- Number of Surface = 4
- Number of tracks per Surface = 5
- Number of Cylinders = 8 (0-7)
- Number of Sectors per track = 8
- Number of Sectors per cylinder =  $2 \times 8 = 32$



8 Sectors  
per track

In this Example each track contains 8 Sector. This is for 1<sup>st</sup> Cylinder (Cylinder=0)

First we drawed (cross) all sector of track 0<sup>th</sup> track of surface 0

Then we drawe (cross) " " " " 0 of Surface 2<sup>nd</sup>

" " " " " " 0 " " " 2  
" " " " " " " " " 3

T Cylinder : Cylinder Number 0

Step 1 : Surface number 0 (I<sup>st</sup> surface) track no '0' (1<sup>st</sup> track) all 8 Sector are drawed (crossed).

Step 2 : Surface number 1 (II<sup>nd</sup> surface) track no 0 (1<sup>st</sup> track) all 8 Sector are drawed (crossed).

Step 3 : Surface number 2 (III<sup>rd</sup> surface) track no. 0 (1<sup>st</sup> track) all 8 Sector are drawed (crossed).

Step 4 : Surface number 3 (IV<sup>th</sup> surface) track no 0 (1<sup>st</sup> track) all 8 Sectors are drawed (crossed).

Now move to next Cylinder, and so on.

T Cylinder : Cylinder Number 0

$\langle L, h, S \rangle$

0<sup>th</sup> Surface {  $\begin{cases} \langle 0, 0, 0 \rangle \text{ 1<sup>st</sup> Sector} \\ \langle 0, 0, 1 \rangle \text{ 2<sup>nd</sup> Sector} \\ \vdots \\ \langle 0, 0, 7 \rangle \text{ 8<sup>th</sup> Sector} \end{cases}$

# Sector/track : 8

# Surface : 4

# track per Surface = 5

# Cylinder = 5 (0-4)

# Surface = 4 (0-3)

1<sup>st</sup> Surface {  $\langle 0, 1, 0 \rangle \text{ 9<sup>th</sup> Sector}$

And so on upto 32<sup>nd</sup> Sector.

- All 32 Sector stored in 4 Surface of cylinder no. 0.
- Now we move to the next cylinder 2<sup>nd</sup> cylinder (Cylinder number 1). And so on to all cylinders.
- In this manner we store total 160 Sector which stored in (32 Sector per cylinder x 5 cylinder).

eg: What is the sector no. of  $\langle 4, 3, 7 \rangle$

(i) To cross 4 cylinder =  $4 \times 4 \times 8 = 128$   
 $(0, 1, 2, 3)$

(ii) To Cross 3 Surface =  $3 \times 8 = 24$

(iii) To Cross 7<sup>th</sup> Surface =  $\frac{7}{159}$  (Ans)

Q. A hard disk has 63 sectors per track, 10 platters each with 2 recording surfaces and 1000 cylinders. The address of a sector is given as a triple  $\langle c, h, s \rangle$  where C is the cylinder number, h is the surface number and s is the sector number. Thus the 0<sup>th</sup> sector is addressed as  $\langle 0, 0, 0 \rangle$  the 1<sup>st</sup> sector as  $\langle 0, 0, 1 \rangle$  and so on.  
 # Surface =  $10 \times 2 = 20$  Surface  
 # Sector/track = 63  
 # Sector per cylinder =  $20 \times 63$   
 $= 1260$  Sector/cylinder.

(i) The address  $\langle 400, 16, 29 \rangle$  corresponds to Sector number:  $\langle 400, 16, 19 \rangle$

(i) To Cross 400 cylinder =  $400 \times 1260 = 504,000$

(ii) To Cross 16 Surface =  $16 \times 63 = 1008$

(iii) To Cross 29 Sector = 29

505087

Alternate Method

Sector Number =  $s + sT + h + sT + Tc + c$

ST: # Sector per track

TC: # track per cylinder

SC: # Sector per cylinder

$$\hookrightarrow ST \times TC$$

$$\begin{aligned} \text{Sector Number} &= 29 + 63 \times 16 + 63 \times 20 \times 400 \\ &= \underline{\underline{505037}} \quad (\text{Ans}) \end{aligned}$$

(ii) The address of 1039<sup>th</sup> Sector is :

$$\langle 0 + 15 \times 63 + 31 \rangle = 976$$

C = [ Given Sector No ]

$$\langle 0 + 16 \times 63 + 30 \rangle = 1038$$

Total # Sector/cylinder

$$\langle 0 + 16 \times 63 + 31 \rangle = 1039$$

$\frac{1039}{1260} = 0$  (1<sup>st</sup> cylinder)

$$\langle 0 + 17 \times 63 + 31 \rangle = 1102$$

$\therefore \underline{\underline{\langle 0, 16, 31 \rangle}}$  (Ans)

$$h = \left( \frac{1039 - 1260}{63} \right) / 63 = \left[ \frac{1039}{63} \right]$$

How much Sector remaining

$$= 1039 - 16 \times 63 = 16 \text{ Surface}$$

$$= 1039 - 1008 = 31$$

$$= \underline{\underline{31}}$$

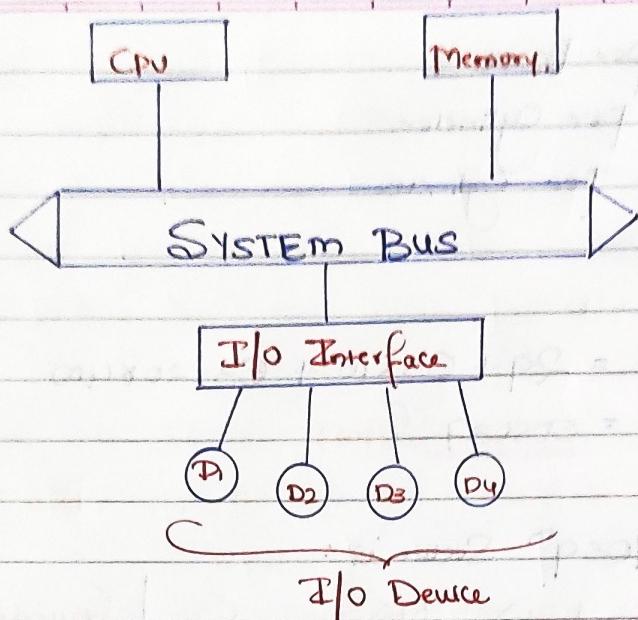
Note:

$\langle C, h, s \rangle$

$$\begin{array}{l} C, h, s = [ \begin{array}{c} \text{Given Sector} \\ \# \text{Sector/cylinder} \end{array} ] \\ (\text{Cylinder Number}) \end{array}$$

$$\begin{array}{l} h = [ \begin{array}{c} \text{Given Sector} / \# \text{sector} \\ \text{per cylinder} \end{array} ] / \text{Total} \\ (\text{Surface Number}) \end{array}$$

$$\text{Sector Number } y = x^* \# \text{Sector}$$



### Input - Output Interface

Input - output interface provides a method for transferring information between internal storage and external I/O devices.

Peripherals Connected to a Computer Needs Special Communications links for interfacing them with the Central processing Unit. The purpose of the communications link is to resolve the differences that exist between the Central Computer and each peripheral. The major differences are:

- ① Peripherals are electrochemical and electromagnetic devices and their manner of operation is different from the operation of the CPU and memory, which are electronic devices. Therefore, a conversion of signal value may be required.
- ② The data transfer rate of peripherals is usually slow than the transfer rate of the CPU and consequently, a synchronisation mechanism may be needed.
- ③ Data code and formats in peripheral differ from the word format in the CPU and memory.
- ④ The operating mode of peripherals are different from each other and each must be controlled so as not to disturb the operations of other peripherals connected to the CPU.

## I/O Organisation

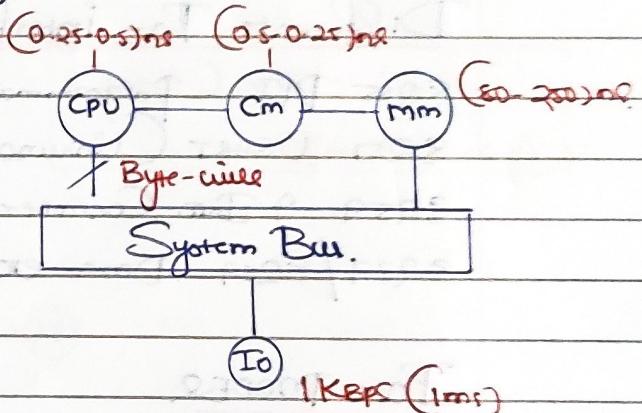
- ① I/O devices are electro-magnetic components and CPU is a electronic component. So, there is a difference that exist is the term of operating modes, data transfer rate and word formats.
- ② To synchronise the I/O speed with a CPU, high speed interface chip is used named as I/O interface or I/O module.
- ③ I/O interface chip is responsible for I/O operations so, the computer system design I/O devices are connected to system bus via I/O interface chip.

## System without I/O-interface (Programmed-I/O).

1KB - 1 sec

1B - ?

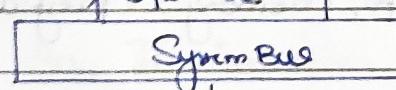
$$ETIO = \frac{1B}{1KB} \text{ sec} = \frac{10^3}{10^3} \text{ sec} = \underline{\underline{1 \text{ ms}}}$$



## System with I/O-interface (INT-Driver-I/O).

(0.25-0.5) ns  
Byte-wise

mm (80-200)ns



I/O  
Interface name

Control Logic

Buffer

(0.25-0.5)ns

IO 1Kbps (1ms)

## Access Sequence / Working Process

- ① CPU initializes the I/O interface chip along with a I/O command and then CPU will go and perform other useful task.
- ② I/O-interface, control logic, interprets the I/O-Command and

Accordingly I/O port will be enabled for the operation.

- (3) Based on the speed of a I/O device and amount of data to be transfer consume the time to prepare the data, then data is transferred from I/O device to a interface buffer.
- (4) When the data is available in the buffer I/O interface generates the interrupt signals and send to the CPU and waiting for ack signal.
- (5) After receiving the ack signal, buffer content will be transferred to CPU. To this process, CPU will be accessing the I/O data from interface buffer therefore Speed gap is synchronized & time is saved.

Different I/O-interface Chip used in the Computer Design is:

8255 PPI (Programmable Peripheral Interface)

8251 USART (Universal Sync. asynchronous receiver transmitter).

8259 A Intr. Controller

8237 / 8257 DMA etc.

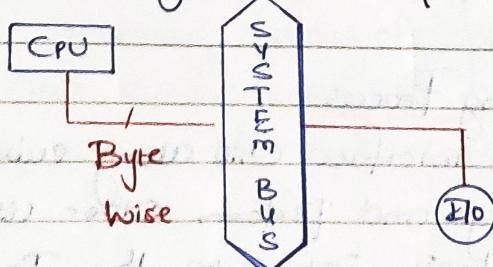
### I/O MODES.

Three types of I/O transfer modes are present in the Computer System, Used to transfer the data from the I/O to other Component of a Computer (CPU, memory). They are:

1. Programmed - I/O.
2. Interrupt Driven I/O
3. DMA (Direct Memory Access).

### Programmed - I/O

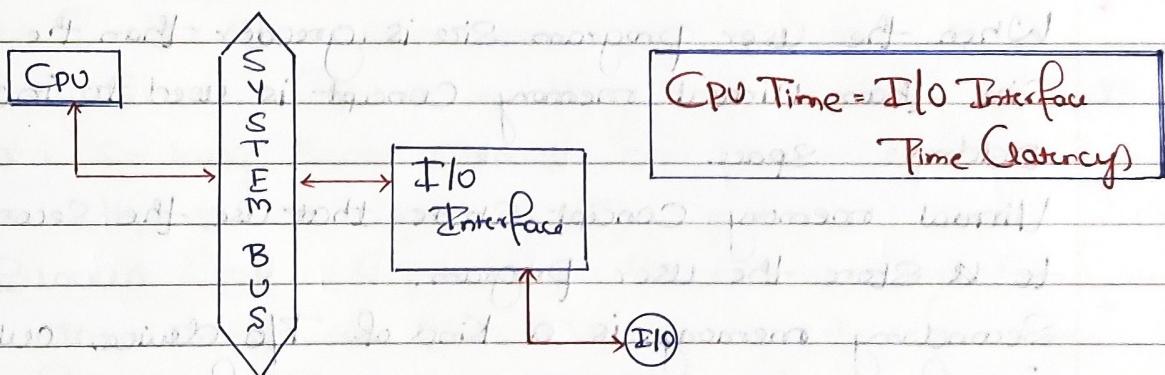
There is no high speed interface logic used



$$\text{CPU Time} = \frac{\text{Speed of I/O}}{\text{(Data transfer Time)}}$$

- \* In this mode, I/O devices are directly connected to CPU without I/O interface chip.
- \* In this mode, CPU takes the responsibility to complete the I/O operation, so CPU will be blocked (Waiting) until the I/O operation is completed.
- \* In this mode CPU utilization is inefficient.
- CPU time depends on the speed of a I/O device and the size of a data unit to be transferred.
- This mode is suitable in the System Centric Applications where the I/O time is important than CPU time.

## (2) Interrupt Driven I/O



- \* In this mode, I/O operations are controlled based on the interrupt signals.
- \* I/O devices are connected to a system bus via I/O-interface chip so, I/O-interface takes the responsibility of a I/O open.
- \* In this mode processor utilization is more efficient so CPU executing the other useful task during the I/O open.
- \* CPU depends on the latency of a interface chip rather than speed of an I/O device.

Drawback: More Complexity in interrupt implementation.

Solution: Use priority. Set highest priority interrupt.

Who will assign:

1. Daisy Chain Method (Static Approach): Fixed: Starvation Occur
2. Polling Method (Dynamic Approach): Changing: No Starvation

- ① Programmed I/O
- ② Interrupt Driven I/O

\* In these two we cannot access the memory directly (we have to access memory using the instructions of CPU).

\* These are used for small amounts of data transfer.

### ③ Direct Memory Access:

\* In this mode, bulk amount of data will be transferred from the I/O to main memory without involvement of a CPU. The Processor And DMA Controller Use the System bus in Transfer Slave mode.

\* When the User program size is greater than the main memory size than Virtual memory Concept is used to increase the address space.

\* Virtual memory Concept State that use the Secondary memory to store the User programs.

\* Secondary memory is a kind of I/O device, which is interfaced to a DMA module, therefore during the program execution data will be transferred from I/O to main memory via DMA controller without involvement of a CPU.

### Access Sequence (DMA)

- ① CPU initializes the DMA module along with a I/O command later busy with other useful task.
- ② I/O Command Contains the information about Port Address, memory address, Control Signal and Count Value.
- ③ DMA module Control logic interprets the command and enables the respective port for the operation.
- ④ Based on the speed of a device consumes the time to prepare the data later enable the "DMA REQ" signal.
- ⑤ After receiving this signal, DMA module enable the hold signal to CPU to gain the control of a system bus and