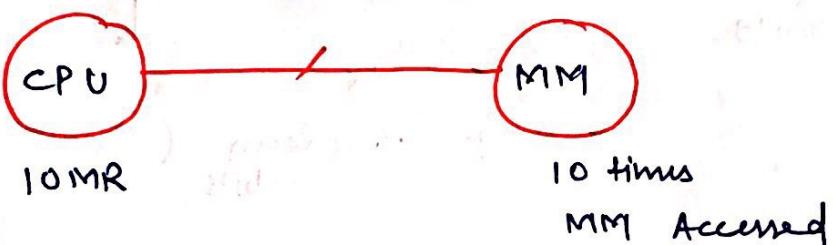


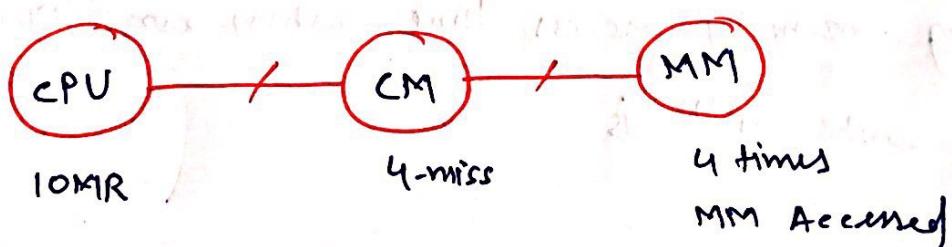
## Multi-level caches :-

- Multi-level cache structure is used in the structure design to reduce the miss-penalty.
- Miss-penalty means time required to transfer the data from higher level to lower level when there is a miss in lower-level memory.

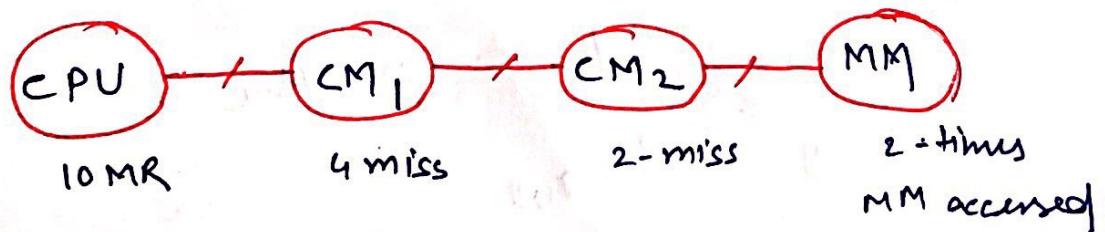
(a) Sys design without cache -



(b) sys design with cache -



(c) sys design with multi-level caches -



- In multi-level cache design -

$$L_1 \text{ CM size} < L_2 \text{ CM size}$$

$$L_1 \text{ CM Access time} < L_2 \text{ CM Access time}$$

$L_1 \text{ CM Association} > L_2 \text{ CM Association to CPU}$

- Two kinds of miss rates are calculated in the multi-level cache design.

(i) GMR (Global miss rate) =  $\frac{\text{no. of misses in cache}}{\text{Total no. of CPU generated references}}$

(ii) LMR (Local miss rate) =  $\frac{\text{no. of misses in cache}}{\text{Total no. of access in the cache}}$

$$GMR_{L_1} = \frac{4}{10} = 0.4$$

$$GMR_{L_2} = \frac{2}{10} = 0.2$$

$$LMR_{L_1} = \frac{4}{10} = 0.4$$

$$LMR_{L_2} = \frac{2}{4} = 0.5$$

- Avg. memory access time is calculated in terms of hit time, miss rate (local) and miss penalty is :

$$T_{avg} = \text{Hit Time}_{L_1} + (\text{miss-rate}_{L_1} * \text{miss-penalty}_{L_1})$$

$$\text{Miss-penalty}_{L_1} = \text{Hit Time}_{L_2} + (\text{miss-rate}_{L_2} * \text{miss-penalty}_{L_2})$$

$$\text{Miss-penalty}_{L_2} = \text{memory (MM) access time.}$$

- Avg. memory stalls created per inst<sup>n</sup> is calculated as:

$$= \left( \frac{\text{no. of misses in } L_1}{\text{Inst}^n} * \text{Hit-time } L_2 \right) + \\ \left( \frac{\text{no. of misses in } L_2}{\text{Inst}^n} * \underbrace{\text{MM access time}}_{\text{Miss-penalty } L_2} \right)$$

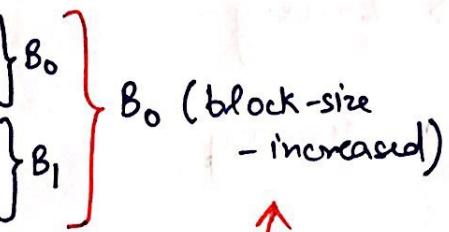
### Types of Cache-misses :-

(i) compulsory-miss :-  
 (Cold-start miss  
 or  
 first reference miss)

This miss will occur when the very first-time reference block is not present in the cache memory during the program execution. This misses can be minimized by increasing the block-size.

eg:-

0000	B
0001	B
0010	B
0011	B
⋮	⋮
1111	B



Prog:-

I<sub>1</sub>: MOV r<sub>0</sub>, [0000]  
 $\frac{B_0}{B_1}$

I<sub>2</sub>: MOV r<sub>1</sub>, [0010]  
 $\frac{B_1}{B_0}$

New Exec:-

I<sub>1</sub>: B<sub>0</sub> (miss)

I<sub>2</sub>: B<sub>0</sub> (Hit)

To minimize  
the misses, block  
size ( $2B \rightarrow 4B$ )

compulsory  
miss

Exec:-

I<sub>1</sub>: B<sub>0</sub> (miss)

I<sub>2</sub>: B<sub>1</sub> (miss)

\* Now [0000] and [0010] both lie in B<sub>0</sub>, block.

(ii) Capacity - miss :- This miss will occur when the cache is full. These misses can be minimized by increasing the cache capacity.

(iii)

Conflict - miss :-

(collision - miss  
or  
interference miss)

This miss will occur when too many blocks are mapped into a same cache-line or cache-set. These misses can be minimized by double the associativity of the cache i.e.

1-way  $\rightarrow$  2-way

2-way  $\rightarrow$  4-way

:

$1 < p < N$

Q. Consider cache-memory with a access time of 50 ns uses write-through protocol. MM access time is 200 ns. CPU generates 60% READ requests and remaining are used for WRITE op<sup>n</sup>. ~~H<sub>r</sub>~~ = 80%. What is avg. memory access time when considering both READ & WRITE op<sup>n</sup>.

Given :-

WT Protocol

$T_c = 50 \text{ ns}$

$T_m = 200 \text{ ns}$

$f_r = 60\%$

$f_w = 40\%$

$H_r = 80\%$

$H_w$  = not given

(Assume 1)

so  $H_w = 1$

this means

simultaneous  
access num.

$$\begin{aligned} T_{avg_r} &= H_r T_c + (1-H_r)(T_m + T_c) \\ &= 0.8(50) + 0.2(200) \\ &= 40 + 40 \\ &= 80 \text{ ns} \end{aligned}$$

$$T_{avg_w} = T_m = 200 \text{ ns}$$

$$\begin{aligned} T_{avg_{WT}} &= f_r \cdot T_{avg_r} + f_w \cdot T_{avg_w} \\ &= 0.6 * 80 + 0.4 * 200 \\ &= \underline{\underline{128 \text{ ns}}} \end{aligned}$$

Q. Consider cache memory having the hit ratio for read & write op<sup>n</sup> is 70% & 80% resp. CM access time is 40ns and MM access time is 400ns. When there is a miss in ~~cache~~ CM then 2 word block is copied from the MM to CM. CPU generates 40% READ req and remaining are used for WRITE op<sup>n</sup>. What is the efficiency of a memory in terms of ~~MWPS~~ (million words per sec) using -

(a) Write-through protocol

(b) Write-back protocol.

Given:-

$$H_r = 0.7$$

$$H_w = 0.8$$

$$T_c = 40\text{ns}$$

$$T_m = 400\text{ns}$$

$$\text{Block-size} = 2W$$

$$f_r = 0.4 \quad (\because \text{of clean-bits})$$

$$f_w = 0.6 \quad (\because \text{of dirty-bits})$$

since  $H_w \neq 1$



so, memory is

hierarchical  
access

$$T_w \text{ (simultaneous - WRITE)} = \max \left( \begin{array}{l} \text{word update in CM} \\ \text{word update in MM} \end{array} \right)$$

$$= \max ( \frac{T_c}{40\text{ns}}, \frac{200\text{ns}}{400\text{ns}} )$$

$$= 200\text{ns}$$

since 2W Block is transferred in  $T_m$  so for  $1W \rightarrow T_m/2$

$$\begin{aligned}
 (a) T_{avg_r} &= H_r T_c + (1-H_r) (T_m + T_c) \\
 &= 0.7(40) + 0.3 (40 + 400) \\
 &= 28 + 132 \\
 &= \underline{\underline{160 \text{ ns}}}
 \end{aligned}$$

$$\begin{aligned}
 T_{avg_w} &= H_w T_w + (1-H_w) (T_m + T_w) \\
 &= 0.8(200) + 0.2 (400 + 200) \\
 &= 160 + 120 \\
 &= \underline{\underline{280 \text{ ns}}}
 \end{aligned}$$

$$\begin{aligned}
 T_{avg_{WT}} &= (f_r * T_{avg_r}) + (f_w * T_{avg_w}) \\
 &= (0.4 * 160) + (0.6 * 280) \\
 &= 232 \text{ ns}
 \end{aligned}$$

$$\left( \eta_{WT} \right) = \frac{1}{T_{avg_{WT}}} \text{ words/sec} = \boxed{4.31 \text{ MWPS}}$$

$$\begin{aligned}
 (b) \frac{T_{avg_r}}{H_r T_c} &= H_r T_c + (1-H_r) \left[ \cdot \% \text{ dirty bits} (T_m + T_m + T_c) \right. \\
 &\quad \left. + \cdot \% \text{ clean bits} (T_m + T_c) \right] \\
 &= (0.7 * 40) + (1-0.7) \left[ 0.6 (400 + 400 + 40) + 0.4 (400 + 40) \right] = \underline{\underline{232 \text{ ns}}}
 \end{aligned}$$

$$T_{avg_w} = HW T_C + (1-HW) \left[ \begin{array}{l} \gamma_{\text{dirty}} \text{-b'ys} (T_m + T_m + T_C) + \\ \gamma_{\text{clean}} \text{-b'ys} (T_m + T_C) \end{array} \right]$$

$$= 0.8(40) + (1-0.8) \left[ 0.6(400+400+40) + 0.4(400+40) \right]$$

$$= \underline{168 \text{ ns}}$$

$$T_{avg_{WB}} = (f_r * T_{avg_{WB,r}}) + (f_w * T_{avg_w})$$

$$= (0.4 * 232) + (0.6 * 168)$$

$$= \underline{193.6 \text{ ns}}$$

$$\eta_{WB} = \frac{1}{T_{avg_{WB}}} = \boxed{5.16 \text{ MWPS}}$$

Q. Consider the following memory specification :

Mem	Access Time	Hit-Ratio
(Inst) I-cache	20ns	0.6
(Data) D-cache	40ns	0.8
MM	200ns	0.9
SM	600ns	1

CPU generates 60% req to access the inst<sup>n</sup> and 40% req to access the data. What is the total avg-memory access time when considering inst<sup>n</sup> & data access.

$$\begin{aligned}
 T_{avg} &= H_{IC} T_{IC} + (1-H_{IC}) H_m (T_m + T_{IC}) + (1-H_{IC}) \\
 &\quad (1-H_m) H_s (T_s + T_m + T_{IC}) \\
 &= (0.6 \times 20) + (0.4 \times 0.9 \times 220) + (0.4 \times 0.1 \\
 &\quad \times 820) \\
 &= \underline{124 \text{ ns}}
 \end{aligned}$$

$$\begin{aligned}
 T_{avg} &= H_{DC} T_{DC} + (1-H_{DC}) H_m (T_m + T_{DC}) + (1-H_{DC}) \\
 &\quad (1-H_m) H_s (T_s + T_m + T_{DC}) \\
 &= (0.8 \times 40) + (0.2 \times 0.9 \times 240) + (0.2 \times 0.1 \\
 &\quad \times 840) \\
 &= \underline{92 \text{ ns}}
 \end{aligned}$$

$$\begin{aligned}
 T_{avg} &= \left( f_{Inst^n} * T_{avg}_{IR} \right) + \left( f_{Data} * T_{avg}_{DR} \right) \\
 &= \boxed{111.2 \text{ ns}}
 \end{aligned}$$

Q. Consider computer system in which CPU generates 200 memory references (MR's), among them 60 miss op<sup>n</sup> present in L1 cache & and 30 - miss op<sup>n</sup> present in L2 cache. Hit-time of L1 memory is 4-cycles & Hit-time of L2 - cache is 20-cycles. Miss penalty of L2 memory is 100-cycles. What is the avg. memory access time.

$$LMR_{L_1} = \frac{60}{200}, \quad LMR_{L_2} = \frac{30}{60}$$

$$T_{avg} = \text{Hit Time}_{L_1} + (\text{miss-rate}_{L_1} * \text{miss-penalty}_{L_1})$$

$$\begin{aligned} \text{miss-penalty}_{L_1} &= \text{Hit-Time}_{L_2} + (\text{miss-rate}_{L_2} * \text{miss-penalty}_{L_2}) \\ &= 20 \text{ cycles} + \left( \frac{30}{60} * 100 \right) \\ &= \underline{\underline{70 \text{ cycles}}} \end{aligned}$$

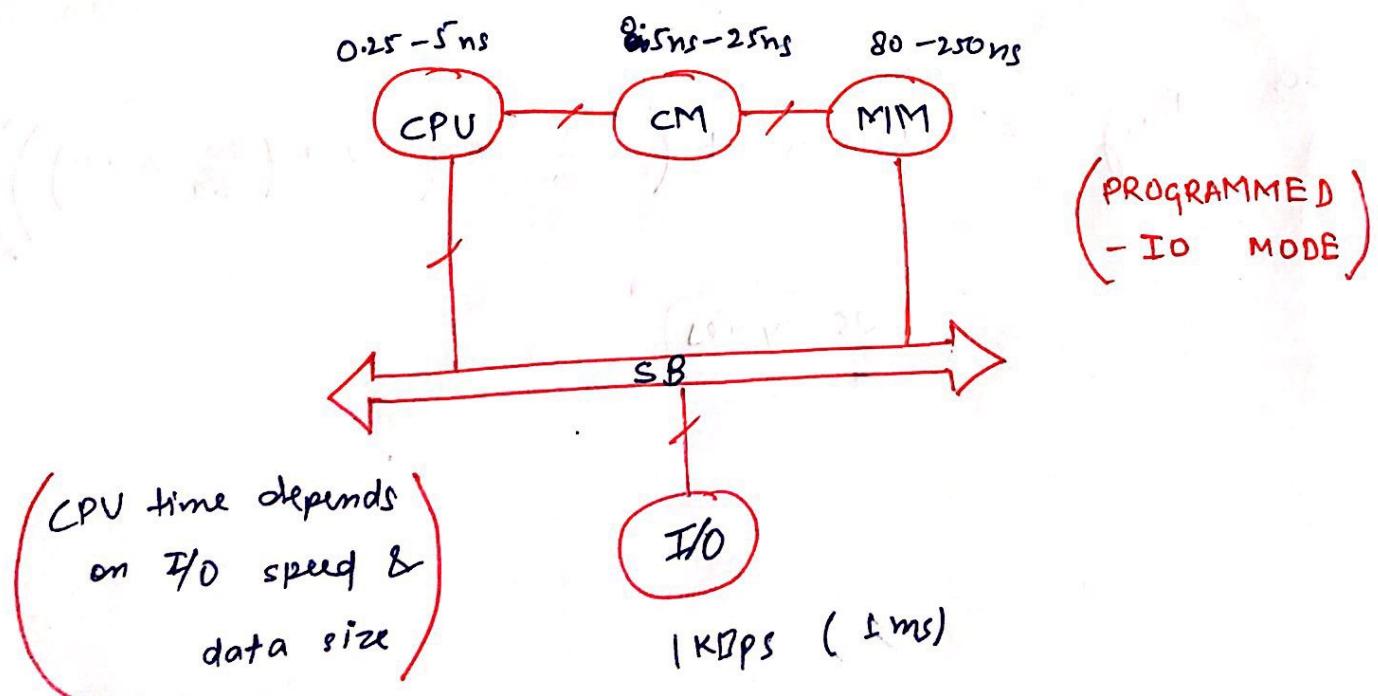
Memory access time

$$\begin{aligned} \text{so, } T_{avg} &= 4 \text{ cycles} + \left( \frac{60}{200} * \left( 20 + \left( \frac{30}{60} * 100 \right) \right) \right) \\ &= \boxed{25 \text{-cycles}} \end{aligned}$$

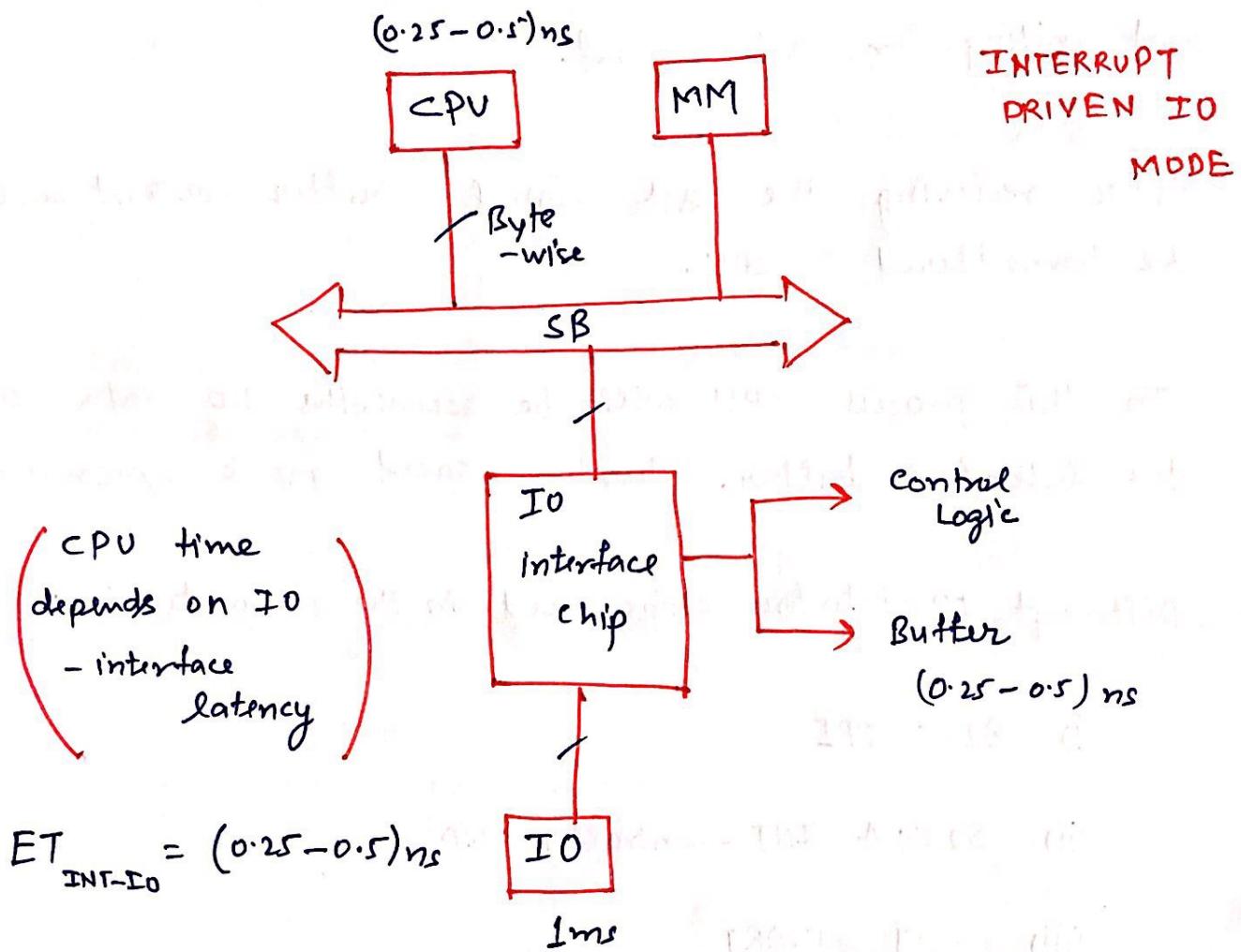
## Input - Output (I/O) organization :-

- I/O devices are electro-magnetic component & CPU is an electronic component, so there is a difference existed in terms of operating modes, data-transistor ratio and word formats.
- To synchronize the I/O properties with a CPU, high-speed interface chip is used called as I/O interface or I/O module.
- I/O interface chip is responsible for I/O op's. ∴ in the computer design all the I/O devices are interfaced to CPU & via I/O interface chip.

### (a) Sys design without I/O Interface logic -



## (b) Sys design with I/O - interface chip logic



### Access Sequence :-

- CPU initializes the IO interface chip along with IO command later busy with other useful task.
- IO-interface control logic interprets the IO command and enables the IO op<sup>n</sup>.
- Based on the speed of a IO-device consume the time to prepare the data, later transfer data to Interface-buffer.

- When the data is available in buffer then IO-interface chip enables the interrupt signal to CPU and waiting for ack signal.
- After receiving the ack signal - buffer content will be transferred to CPU.
- In this process CPU will be accessing the IO-data from the interface buffer. Therefore speed-gap is synchronized.
- Different IO-interface chips used in the computer is :
  - (i) 8255 PPI
  - (ii) 8259A INT-controller chip
  - (iii) 8251 USART
  - (iv) 8257 DMA  
8237

### IO-modes :-

3-different IO-transfer modes are used in the system design to transfer data -

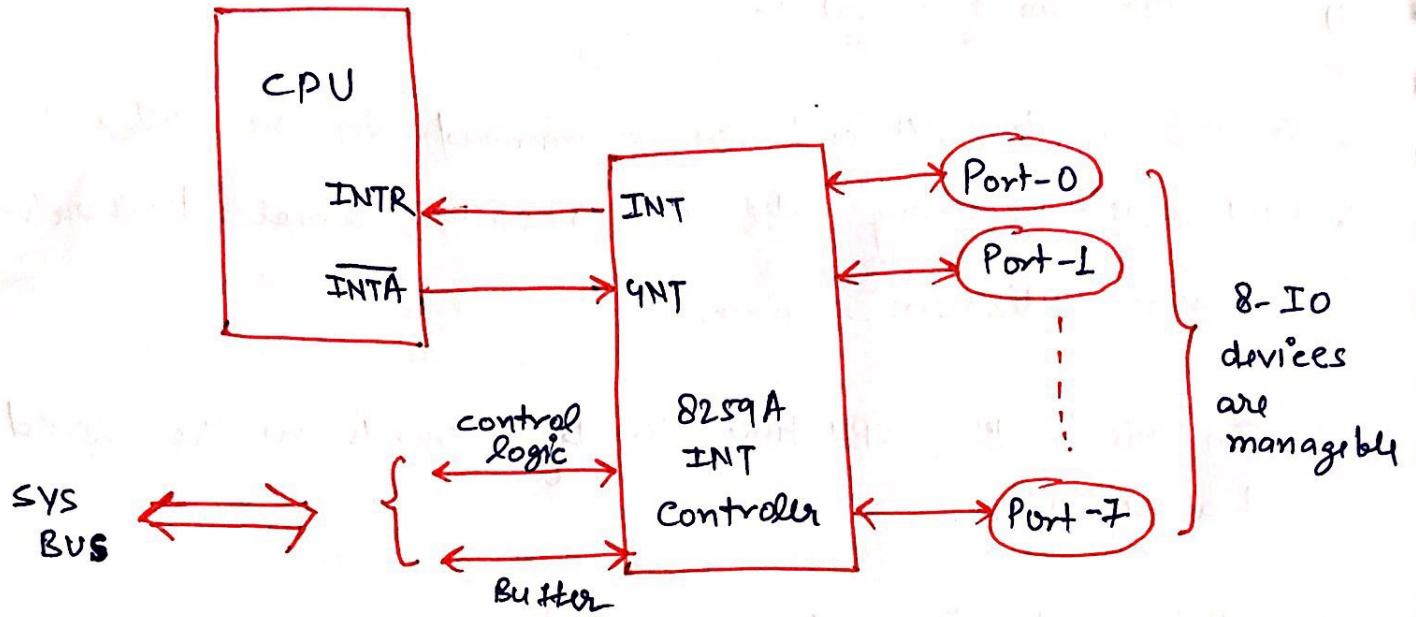
- (i) Programmed - IO
- (ii) Interrupt - driven IO
- (iii) DMA

### (i) Programmed - IO :-

- In this mode IO op<sup>n</sup> are programmed in the CPU, so CPU will be waiting until the IO op<sup>n</sup> is completed therefore processor utilization is poor.
- In this mode, CPU time directly depends on the speed of a IO device.
- This mode is used in the system design where application goals are important than the CPU time.

### (ii) Interrupt - driven IO :-

- In this mode IO op<sup>n</sup> are controlled by based on interrupt signals.
- In this mode, IO interface chip is responsible for IO op<sup>n</sup> so CPU time is utilized for other useful task. Therefore processor utilization is efficient.
- In this mode CPU time depends on the latency of a IO interface rather than the speed of a IO device.

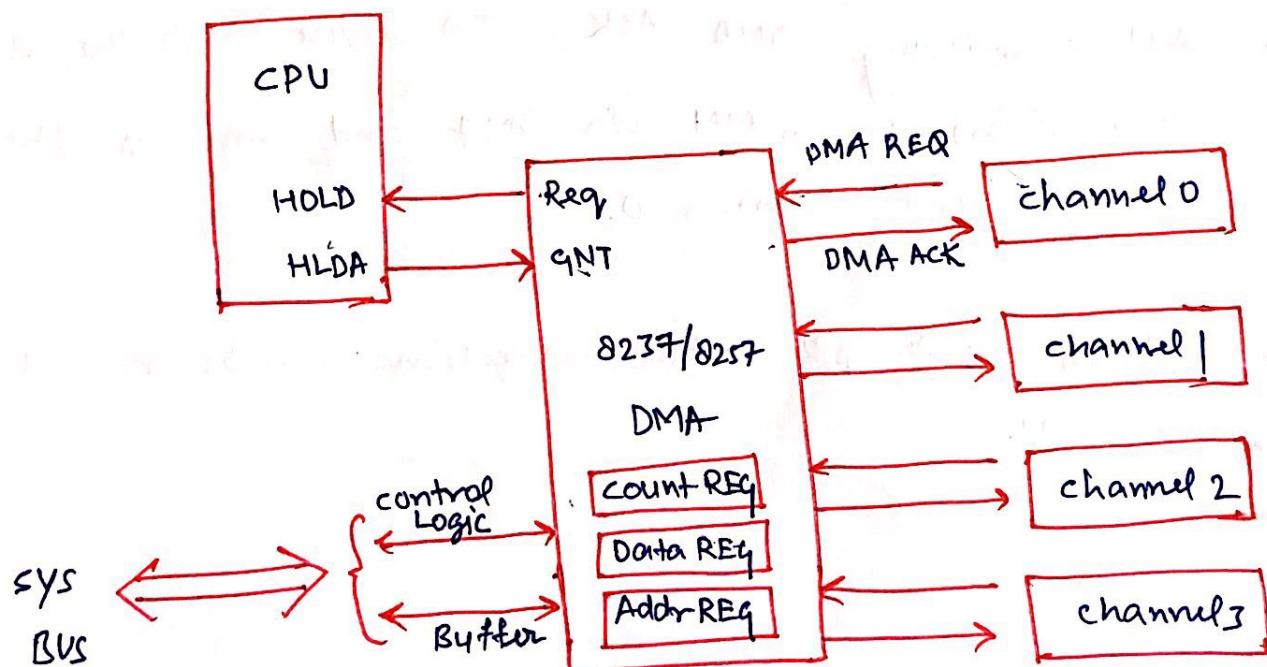
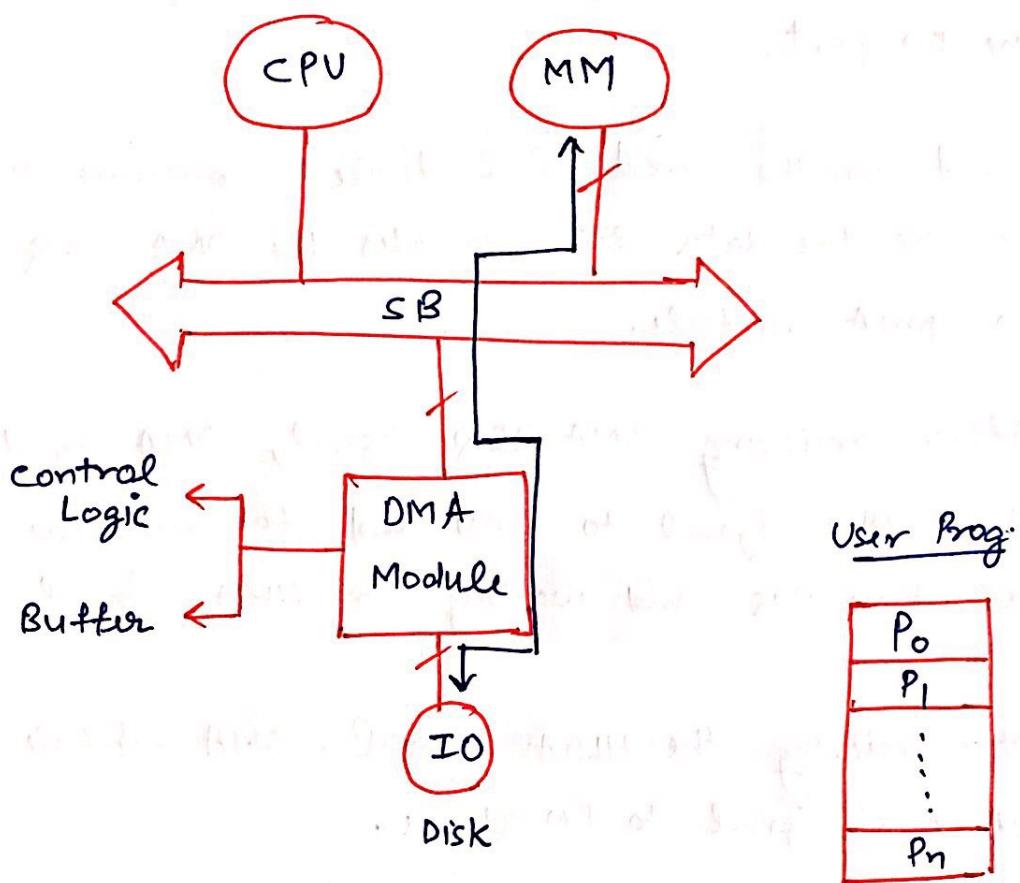


NOTE:- 8259 A chip is programmed into a cascading mode to manage more than 8 IO devices (max is 64 devices)

### (iii) Direct - Memory Access DMA :-

- In this mode bulk amount of data will be transferred from IO to main memory without involvement of CPU.
- when the program size is greater than MM size then virtual memory concept is used to increase the address span. This concept states that use the secondary storage space to store the user program.
- Secondary storage is an electromagnetic storage so it is classified under the ZO - category.

- In the computer design, secondary storage components are interface to system bus via DMA module, therefore during the execution of a program, data will be transferred from IO to main memory via DMA without involvement of a CPU.



- CPU initializes the DMA module along with a IO command later busy with other task.
- IO command contains port address, mem addr, control signals and count value.
- DMA control logic interprets the command & enables the IO port.
- Based on the speed of a device , consume the time to prepare the data later enables the DMA REQ signal to DMA module.
- After receiving DMA REQ signal, DMA module enables the HOLD signal to CPU and to gain the control of a system BUS and waiting for HLDA signal.
- After receiving the HLDA signal, DMA module enables the DMA ACK signal to IO device.
- After receiving DMA ACK, IO device starts the data transmission to MM via DMA and continue the op<sup>n</sup> until count becomes 0.
- After DMA op<sup>n</sup> , Bus connection will be re-established to CPU.

**NOTE:-** In the DMA op<sup>n</sup> CPU presents in two states -

- (i) Busy state
- (ii) Blocked or Hold state

- CPU is in Busy state until the IO device prepares the data.
- CPU is in Blocked state until the data is transferred to MM.
- Let x is a preparation time and y is transfer time then -

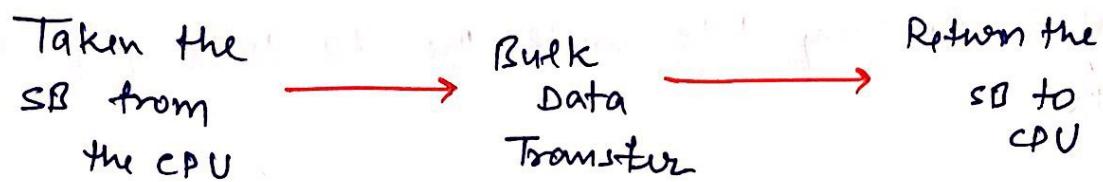
$$\% \text{ time}_{\text{CPU busy}} = \left( \frac{x}{x+y} \right) \times 100$$

$$\% \text{ time}_{\text{CPU blocked}} = \left( \frac{y}{x+y} \right) \times 100$$

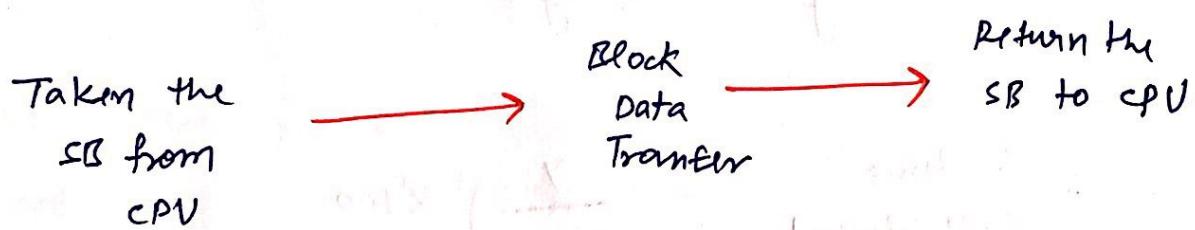
- DMA module is operating under 3 modes -

- (i) Burst mode
- (ii) Cycle-stealing mode
- (iii) Inter-leaving mode

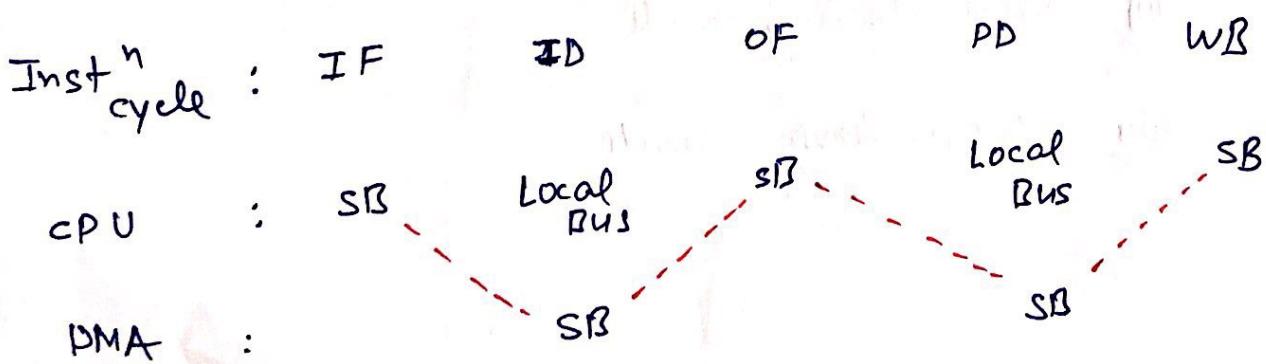
- In Burst mode of DMA, Bulk amount of data at a time will be transmitted to MM after receiving the sys BUS from the CPU. Therefore CPU waiting time is more.



- In cycle-stealing mode of DMA, data file is organized into a small units named as block, later transfer the data file to MM by req. the BUS control multiple times. Therefore CPU will be suspended in a small time interval.



- In Inter-leaving mode of DMA, data will be transferred to MM in form of words when the CPU busy with local op<sup>n</sup>.



Q. Consider 10KBps IO-device interface to 32-bit CPU in a programmed IO-mode. Data-transmission takes place b/w the CPU and I/O in a word-wise sequence. I/O interface overhead is  $\frac{40\mu s}{400\mu s}$ . How much performance is improved when the device is operating under interrupt mode under programmed IO mode.

### Prog - IO

10KB — 1sec

1W

$$(4B) \Rightarrow \frac{4B}{10KB} \text{ sec} = 400\mu s = ET_{Prog}$$

### INT - IO

$$\begin{aligned} \text{CPU time} &= \text{IO-interface latency} \\ &= 40\mu s. \end{aligned}$$

$$\text{i.e } ET_{INT} = 40\mu s$$

$$S = \frac{ET_{Prog}}{ET_{INT}} = \frac{400}{40} = 10$$

$$S = 10$$

Q. Consider 100 KBps IO device interfaced to 32-bit CPU in a cycle stealing mode of DMA whenever 8W data is available in the buffers then it is transferred to MM. Machine cycle time is 2MS. How much % of CPU time is consumed in the DMA opn.

$$\% \text{ time CPU is blocked} = \left( \frac{y}{x+y} \right) \times 100$$

$x$  is preparation time  
(Depends on IO-device speed)

$$100 \text{ KB} \rightarrow 1 \text{ ms}$$

$$1B \rightarrow \frac{1}{100K}$$

$$8W(32B) \rightarrow \frac{32}{100K} = \underline{\underline{320 \text{ MS}}}$$

$$x = 320 \text{ MS}$$

$y$  is transfer time (MM speed)

(1M/C cycle is used to transfer 1 word data)

so  $\rightarrow$  8 M/C cycles req., to transfer 8 words

$$y = 8 \times 2$$

$$y = 16 \text{ MS}$$

$$\% \text{ time CPU blocked} = \frac{16}{320+16} \times 100 \\ = \boxed{4.76 \%}$$

Q. Consider 1MBPS IO-device interfaced to 32-bit CPU using DMA interface. DMA contains 8-bit count register, 32-bit data reg & 24-bit addr reg. Data file size is 4KB. How many DMA cycles are required to transfer the data file to MM.

DMA op<sup>n</sup> is controlled by the count register-

count reg = 8-bit reg  
= Maintains 256 counts (max)

In every count, 1W data is transmitted

$$\frac{\text{Data size}}{\text{DMA cycle}} = \frac{256 \text{ W}}{1} = \frac{256 \times 4B}{1} \\ = 1024B \\ = 1KB$$

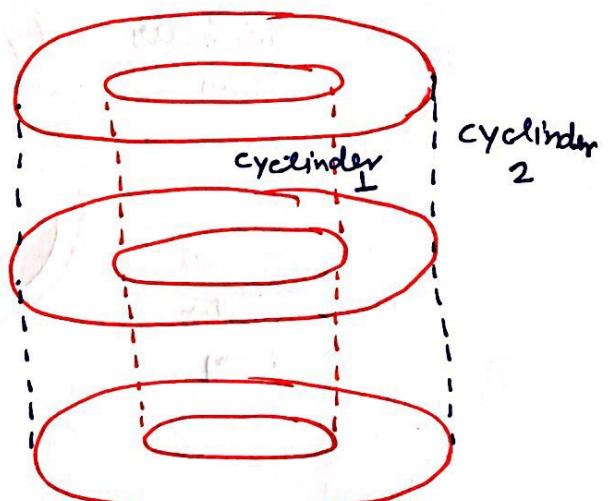
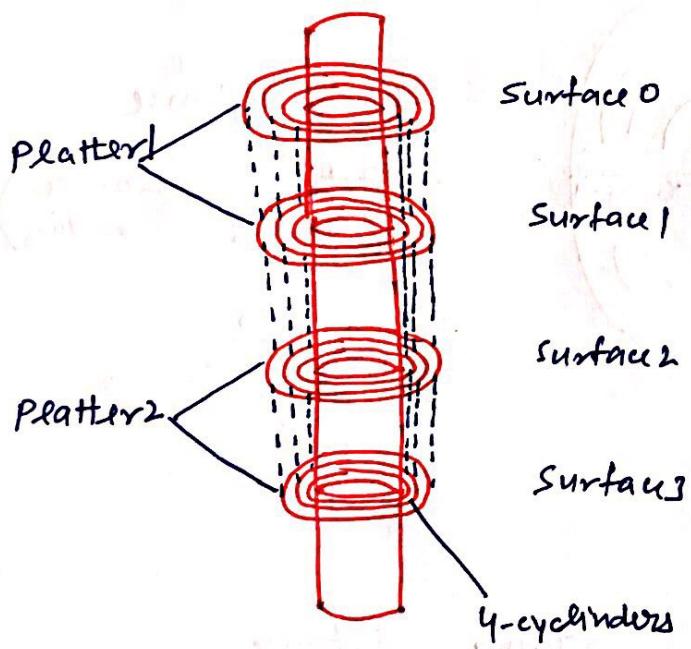
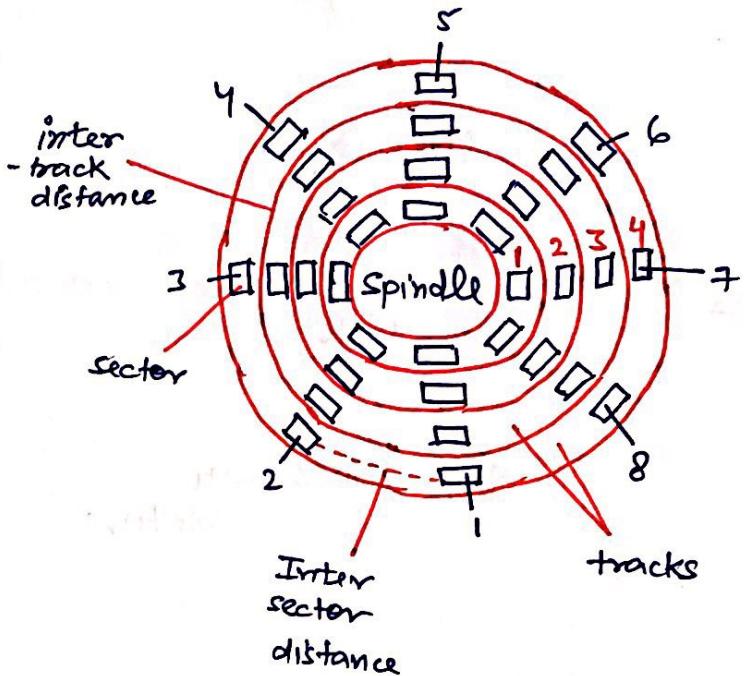
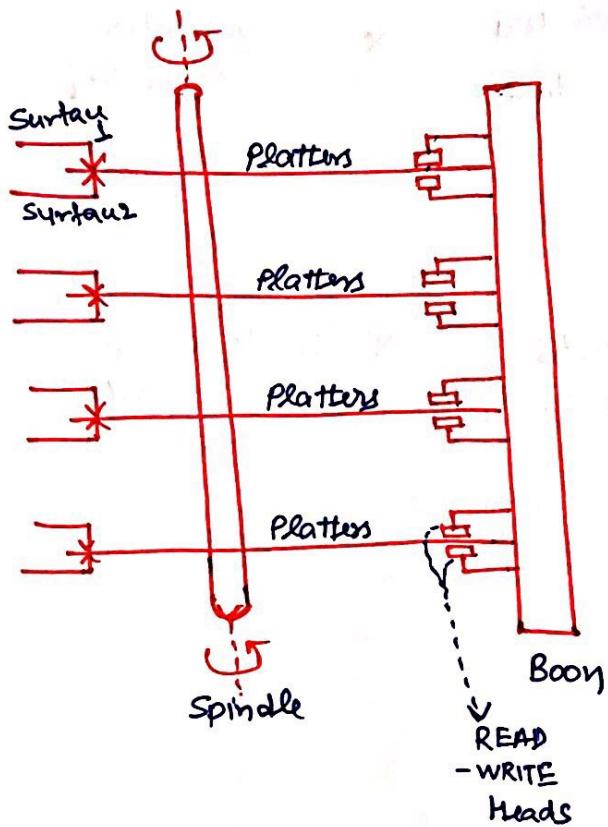
$$\text{file-size} = 4KB$$

$$\text{DMA cycles req for } 4KB = 4 \text{ cycles.}$$

## Disk Hard Structure :-

- It is a direct access electromagnetic storage component.
- It contains a bunch of magnetic coated platters to store the data.
- Each platters contains two surfaces. Each surface contains set of tracks.
- Each track contain set of sectors.
- Sector holds the data.
- In the hard disk each surface has its own READ - WRITE head , use to access the data.
- In the hard disk surfaces, tracks and sectors are the addressable units.
- In the hard disk cylinder is formed by connecting the same track no. in all the surfaces.

Platters — surfaces — tracks — sectors — bytes



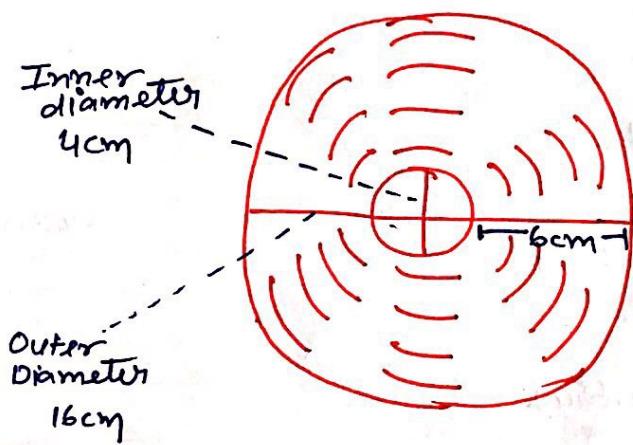
- no. cylinders in the disc = no. tracks in the surface
- Track capacity = no. sectors per track  $\times$  no. bytes per sector
- Cylinder capacity = no. surfaces in disk  $\times$  track capacity

$$\text{Disk capacity} = \text{no. of cylinders in the disc} \times \text{cylinder capacity}$$

### Generalization:-

$$= \text{no. Surfaces in the disk} \times \text{no. of tracks per surface} \times \text{sectors per track} \\ \times \text{bytes per sector}$$

### Surface Analysis:-



$$\begin{aligned}\text{Storage distance} &= (\text{Outer radius} - \text{Inner radius}) \\ &= 8 - 2 \\ &= 6\text{cm}\end{aligned}$$

$$\begin{aligned}\text{no. of tracks per surface} &= \frac{\text{Storage Distance}}{\text{Inter track Distance}} \\ &= \frac{6\text{cm}}{0.2\text{mm}} \\ &= 300\end{aligned}$$

- Different adjustments are required in the disk op^n.
- Seek time: Amount of time required to move the head point from the initial position to desired track in the surface.
  - Rotational Latency: Amount of time required to adjust the head point at the beginning of desired sector in the track. In best case 1 revolution is required to adjust the head points. In worst case half of the revolution is required to adjust the head point.
  - Transfer time: Amount of time req to access the data from hard disk
  - Overhead: Additional delays in the disk op^n.
- Total Time req to access the data from the hardisk
- $$(T_{avg} = \frac{\text{seek time} + \text{Avg Rotational latency}}{\text{Transfer time}} + \text{Overhead})$$

**NOTE :-** When the sector address is expressed in a 3D address format then sector no. in the hard disk is calculated as,

$$\text{sector no} = (i, j, k)$$

cylinder no.      |      Surface no.      Sector no.

$$\text{Sector no} = k + s(j + i * t)$$

count of  
a sector  
no. w.r.t.  
surface      |  
 count of  
respective  
surface      |  
 count of  
current  
cylinder

**Q:** Consider hard disk with a rotational rate of 9800 rpm. Avg seek time of a disk is 8.4 ms. Disk contains 64 platters, surface contains 128 tracks, track contains 256 sectors. Sector holds 512B data. 32KB file is stored in the hard-disk.

Calculate the following:

- Disk Capacity.
- Time req to access file
- " " " 100 random sectors

(d) Disk - data transfer rate ( bandwidth )

$$\begin{aligned}\text{(a) Disk capacity} &= 64 \times 2 \times 128 \times 256 \times 512 \\ &= 2^{31} B \\ &= \boxed{29B}\end{aligned}$$

$$\text{(b) Seek-time} = \underline{8.4 \text{ ms}}$$

$$9800 \text{ revolution} = 1 \text{ min} = 60 \text{ sec}$$

$$1 \text{ revolution} = \frac{60}{9800} \text{ s}$$

$$= 6.12 \text{ ms}$$

$$\therefore \text{avg. Rotational-latency} = \frac{6.12}{2} = \underline{\underline{3.06 \text{ ms}}}$$

1-revolution time  $\rightarrow$  1-track data  
( 256 sectors )

$$1 \text{ sector} \Rightarrow 512 B$$

$$\text{for } 32KB \rightarrow \frac{32KB}{512} = 64 \text{ sectors}$$

so,  $6.12 \text{ ms} \rightarrow 256 \text{ sectors access}$

$$1 \text{ sector access} = \frac{6.12}{256} \text{ ms}$$

so, time req. for access 64 sectors data -

$$\text{Transfer time} = \frac{6.12 \text{ ms}}{256} \times 64 \\ = \underline{\underline{1.53 \text{ ms}}}$$

$$T_{\text{avg}} = (8.4 + 3.06 + 1.53 + 0) \text{ ms} \\ = \boxed{12.99 \text{ ms}}$$

(c) Random sectors accessing requires the adjustment to every sector.

$$\text{for 1 sector Transfer time} = \frac{6.12}{256} \\ \text{one-sector} \\ = 0.0239 \text{ ms}$$

$$T_{\text{avg}} = (8.4 + 3.06 + 0.0239 + 0) \text{ ms} \\ \text{for 1 sector}$$

$$T_{\text{avg}} = (8.4 + 3.06 + 0.0239 + 0) \text{ ms} \\ \text{for 100 random sectors} \\ = \boxed{148.39 \text{ ms}}$$

#### (d) Bandwidth

Surface data rate : In 1 revolution time — 1 track data accessed

$$6.12 \text{ ms} \quad \underline{\quad} \quad 256 \times 512 \text{ B data}$$

$$\text{so, Data-rate} = \frac{256 \times 512 \text{ B}}{6.12 \text{ ms}}$$

Since Total  $\approx 128$  surfaces are present

$$\text{so Total data rate} = \frac{256 \times 512 \text{ B} \times 128}{6.12 \text{ ms}}$$

### Spatial Locality in the Memory:-

- Spatial Locality means adjacent data words accessing in the same block, in a sequence.
- In the array implementation proves spatial locality is realized in the memory structure.
- Array is an ordered set of homogeneous data elements, stored in the memory in a row-wise sequence.

eg:- char a[4][4];

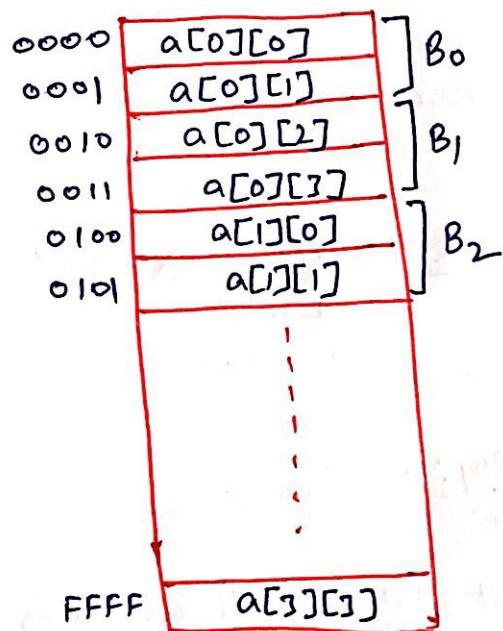
$a[0][0]$ ,  $a[0][1]$ ,  $a[0][2]$ ,  $a[0][3]$

$a[1][0]$ ,  $a[1][1]$ ,  $a[1][2]$ ,  $a[1][3]$

$a[2][0]$ ,  $a[2][1]$ ,  $a[2][2]$ ,  $a[2][3]$

$a[3][0]$ ,  $a[3][1]$ ,  $a[3][2]$ ,  $a[3][3]$

### Memory storage



Assume block size =  $2B$

### Array elements Accessing

(a) Row-major Indexing:- i.e.  $a[0][0]$ ,  $a[0][1]$ ,  $a[0][2]$

Assume CM is empty -

Access :-

$$a[0][0] : M ; \quad B_0 (a[0][0], a[0][1]) \rightarrow CM$$

$$a[0][1] : H ;$$

$$a[0][2] : M ; \quad B_1 (a[0][2], a[0][3]) \rightarrow CM$$

$$a[0][3] : H ;$$

Hit follows miss  $\rightarrow$  M = 50%

$$\text{Element size} = 1B \quad \text{no. of miss / row} = \underline{2}$$

$$\text{Block size} = 2B$$

$$\text{no. of elements / block} = 2B / 1B = 2$$

$$\text{no. of Block / row} = 2$$

$$\text{So, Total no. of misses / array} = \text{no. of rows} * \text{no. of misses / row}$$

$$= 4 * 2$$

$$= \boxed{8}$$

(b) Col-major Indexing :- i.e.  $a[0][0], a[1][0], a[2][0], \dots$

Assume CM is empty

$a[0][0] : M \quad B_0 \rightarrow CM$

$a[1][0] : M \quad B_1 \rightarrow CM$

$a[2][0] : M \quad B_2 \rightarrow CM$

$a[3][0] : M \quad B_3 \rightarrow CM$

\* Random Block references replace the existed Blocks in CM.

Every element in the column is miss.

Element size = 1B

Block size = 2B

no. of elements / Block =  $\frac{2B}{1B}$       no. of Element / col = 4

no. of Blocks / col = 4

so,    no. of miss / col = 4

Total misses / array = no. of cols \* no. of miss / col  
=  $4 * 4$   
= 16

## Conclusion

Storage sequence	Access sequence	Spatial Locality
Row-wise	Row-wise	YES
Row-wise	Col-wise	NO
Col-wise	Row-wise	NO
Col-wise	Col-wise	YES