# COMPUTER SCIENCE

Computer Organization and Architecture

## Cache Memory

Lecture_03

Vijay Agarwal sir

TOPICS TO BE COVERED

**o1** Memory Access

**o2** Cache Memory

Memory Hierarchemy

Cache Memory

$\quad\quad\rightarrow$ Simultaneous Access

$\quad\quad\rightarrow$ Hierarchical Access

[L.O.R]

Locality of Reference

(i) Temporaral LOR

(ii) Spatial LOR

Type of Cache

Numerical.

# LOR [Locality of Reference]

❏ Access the higher level of memory Data from level 1 Memory is called L.O.R,.

(1) Temporal LOR

(2) Spatial LOR.

(1) Temporal LOR: means the same word in the same block is reference by the CPU in near future (Frequently)[Eg: LRU]

Or

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

## LOR [Locality of Reference]

(2) Spatial LOR means adjacent word in the same block is referenced by the CPU in a sequence.

# L.O.R

Medical Store

$\Downarrow$

10pc Medicine

Only for 1 time U go Medical
Store
9pc Medicine is Near by you
at ur Home.

42
29
———

Till Now: 71 GATE QUESTION

30 → Cache
10 → Disk & DMA
————————
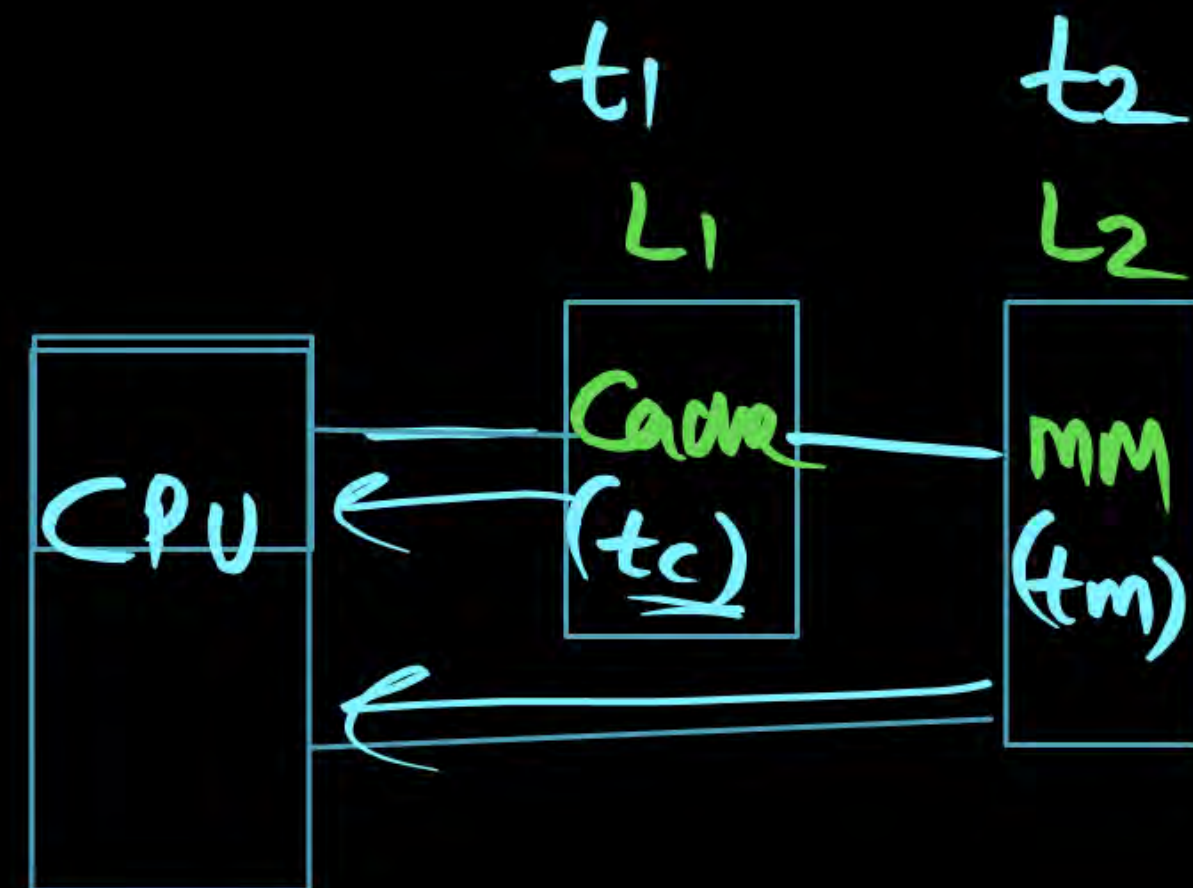111+ GATE PYQ in class.

# Type of Memory Org

Last level
Hit Ratio = 1
$H_2 = 1$

1. Simultaneous Access Memory Org.

$$H_1 T_1 + (1-H_1)(H_2)(T_2)$$

$$\boxed{T_{avg} = H t_1 + (1-H) t_2}$$

if in term of
cache & Main
Memory
then $T_{avg}$

$$\boxed{T_{avg} = h t_c + (1-h)[t_m]}$$

$t_1$     $t_2$

$L_1$     $L_2$

CPU ← Cache ($t_c$) — MM ($t_m$)

# Type of Memory Org

1. Simultaneous Access Memory Org.

$$1 \text{ word Access time} = T_{avg}$$

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

(Data transfer Rate)

OR

Performance

# Type of Memory Org

2. Hierarchical Access Memory Org.

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

↓ OR

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

$\rightarrow h t_c + t_m + t_c - h t_m - h t_c$

$$\boxed{t_c + (1-h) t_m}$$

$h t_c + (1-h) t_m + t_c - h t_c$

$t_c + (1-h) t_m$

$t_c$ : Cache Mem Access time

$t_m$ : Main Memory Access time

$h$ : Cache Hit Ratio

$(1-h)$ : Cache Miss Ratio

# Type of Memory Org

$$\text{Hit Ratio } t_1 = h_1$$
$$\text{Miss Ratio}(m_1) = (1-h_1)$$

## 2. Hierarchical Access Memory Org.

### 3 level

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

Last Level Hit Ratio always $= 1$

OR

miss Rate

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

$m_i$: miss Ratio of level $i$.

CPU — $L_1$ — $L_2$ — $L_3$

$t_1$    $t_2$    $t_3$

$h_1$    $h_2$    $h_3$

$$h_3 = 1$$

Last Level Hit Ratio $= 1$

$m_1$: Miss Ratio of level 1 $(1-h_1)$

$m_2$: Miss Ratio of level 2 $(1-h_2)$

## 2. Hierarchical Access Memory Org.

$1 \text{ Word Access time} = T_{avg}$

$$\text{Data transfer Rate} \atop (\text{Performance})$$
efficiency
$$\#\text{Words}/\text{sec} = \frac{1}{T_{avg}}$$

Remember

$1 \text{ Inst}^n \text{ ET} = 5.51 \text{ nsec} \quad \frac{5.51 \times 10^{-9}}{\text{sec}}$

In $1 \text{ sec} \longrightarrow \#\text{Inst}^n$

$5.51 \times 10^{-9} \text{ sec} \longrightarrow 1 \text{ Inst}^n$

$1 \text{ sec} \longrightarrow \frac{1}{5.51} \times 10^9 \text{ Inst}^n/\text{sec}$

$\frac{1000}{5.51} \times 10^6$

$= 181.14 \text{ MIPS}$

# If Locality of Reference is Considered.



① Block Size 1 Word

② Block Size More than 1 Word (Assume n Words)
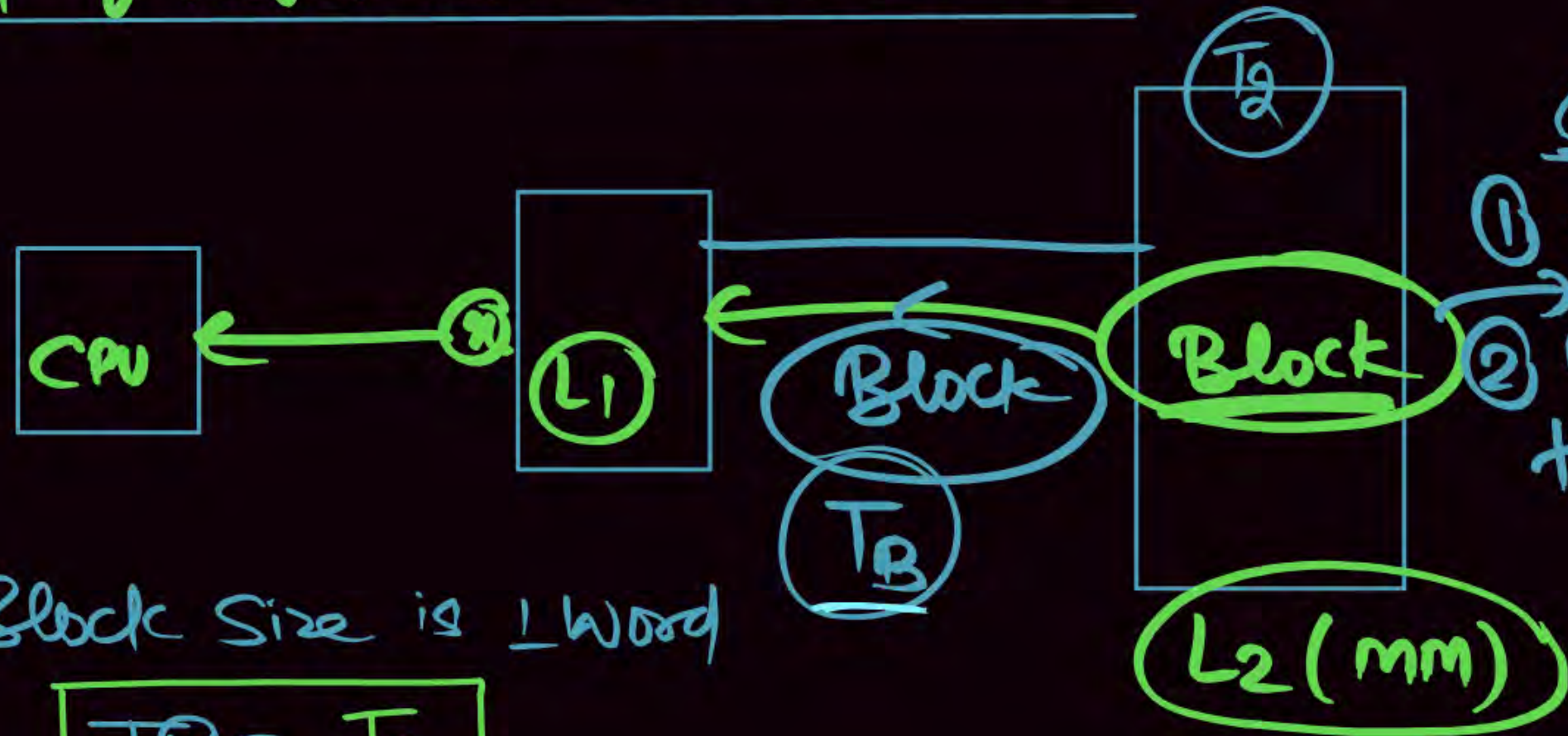
CPU Always Access the Data from the Cache (Faster/Level 1) Memory. If there is a Miss in Level 1 Memory & Hit in Level 2 (MM @ Slower Compare to L1) then One Complete Block is transfer from L2 Memory to L1 Memory & addressed Word (which Request/Demand) by the CPU given From Faster (Level 1/Cache) Memory.

# If Locality of Reference is Considered.



Case
① Block Size 1 word
② Block Size More than 1 word (Assume n words)

Case I : If Block Size is 1 word

$$TB = T_2$$

Case II : If Block Size is n words

$$TB = n * T_2$$

TB : Block Transfer time from $L_2$ Memory to $L_1$ Memory.

## For 2 Level

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

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

OR

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

## Locality of Reference.

$$T_{avg} = h t_1 + (1 - h) (\boxed{TB} + t_1)$$
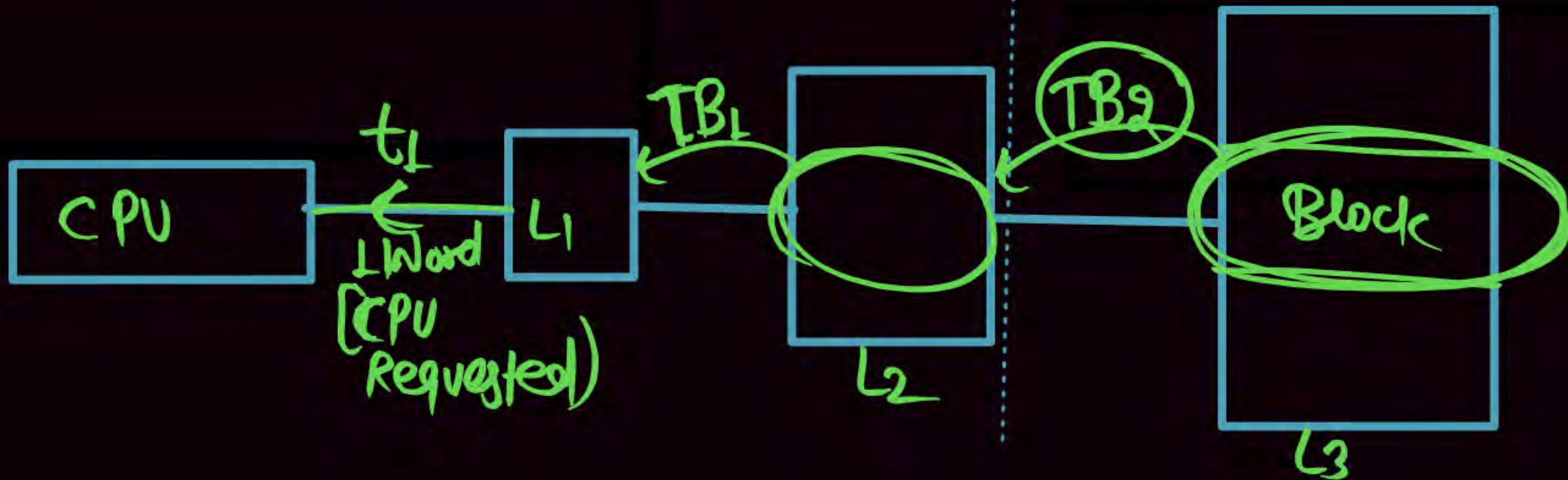
OR

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

# $T_{avg}$ for 3 Level

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

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

## Locality of Reference.

$$T_{avg} = h_1 t_1 + (1-h_1) h_2 (TB_1 + t_2) + (1-h_1)(1-h_2)(TB_2 + TB_1 + t_1)$$

CPU

$t_1$

1 word
(CPU
Requested)

$L_1$

$TB_1$

$L_2$

$TB_2$

Block

$L_3$

**Q.** In a 3 level memory, level 1 memory Access time is T1, level 2 memory Access time is T2(TB1) and level 3 memory Access time is T3(TB2). Hit ratio of level 1 is h1 and Hit ratio of level 2 is h2. What is the average Access time Using Hierarchical Access?

(i) If there is a hit in level1(h1=100%). $h_1 = 1$

(ii) If there is a miss in level1 & hit in level2(h2=100%) $h_2 = 1$ & $h_1 = 0$

(iii) If there is a miss in level1 and Level 2 & hit in level3.

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

Locality of Reference OR

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

$$\boxed{\text{Hit} + \text{Miss} = 1}$$

80% Hit that means 20% Miss.

(i) $H_1 = 100\%$    $H_1 = L$    $(1-H_1) = \underline{0}$

$T_{avg} = \dfrac{H_1 T_1}{0} + \underset{0}{\underbrace{(1-H_1)}} H_2 (T_2 + T_1) + (1-H_1)(1-H_2)$
$\hspace{7cm} T_3 + T_2 + T_1$

(i) $\boxed{H_1 = 1}$ $\hspace{3cm}$ $(1-H_1) = 0$

$\boxed{T_{avg} = T_L}$ $\;$ Ans

(OR)
$T_{avg} = H_1 T_1 + (1-H_1) H_2 (TB_1 + T_1) + (1-H_1)(1-H_2)(TB_2 + TB_1 + T_1)$

$\boxed{T_{avg} = T_L}$ Ans

_____

(ii) $\boxed{H_2 = 1}$ $\boxed{H_1 = 0}$

$\boxed{T_{avg} = T_2 + T_1}$ Ans

$\boxed{T_{avg} = TB_1 + T_L}$ Ans

first Block transferred from $L_2$ to $L_1$ (TB1) then Respective (Requested) Word given $L_1$ to CPU ($T_L$)

_____

(iii) $H_1 = 0$ $\;$ $H_2 = 0$ $\;$ $\underline{H_3 = 1}$

$\boxed{T_{avg} = T_3 + T_2 + T_1}$ Ans

$\boxed{T_{avg} = TB_2 + TB_1 + T_L}$ Ans

**Q.** $h_2 = 1$

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 level1 then 4word block must be transferred(moved) from level 2 into level1 and then addressed word is given to CPU. What is the average Access time ?

(P W)

**Soln**

$h = 90\%$

$\boxed{h = 0.9}$

$T_1 = 30 \, nsec$

$T_2 = 250 \, ns/word$

Block Size = 4 word

$TB_L = 4 \times 250 = 1000 \, nsec$

$$\boxed{T_{avg} = h\, t_1 + (1-h)\, (\overset{(TB_L)}{t_2} + t_L)} \qquad \left(\text{Complete } \underset{T_2}{} \right)$$

$\Rightarrow 0.9 \times 30 + (1 - 0.9)(1000 + 30)$

$27 + 103 = \boxed{130 \, nsec}$ Ans

**OR**

$$\boxed{T_{avg} = t_1 + (1-h)\,\overset{t_{B_L}}{\cancel{t_2}}} \Rightarrow$$

**OR**

$30 + (1 - 0.9)\,1000 \Rightarrow 30 + 100 = \boxed{130 \, nsec}$ Ans

OR

**Q.** Assume that for a certain processor, a read request takes 50 nanoseconds on a cache miss and 5 nanoseconds on a cache hit. Suppose while running a program, it was observed that 80% of the processor's read requests result in a cache hit. The average read access time in nanoseconds is _____.  [GATE – 2015]

$$T_{avg} = 0.80 \times 5 + (1 - 0.80)\, 50$$
$$= 4 + (0.20)\, 50$$
$$= 4 + 10$$
$$\boxed{T_{avg} = 14\, nsec}\ \underline{Ans}$$

# NAT ①

Assume that for a certain processor, a read request takes 50 nanoseconds on a cache miss and 5 nanoseconds on a cache hit. Suppose while running a program, it was observed that 80% of the processor's read requests result in a cache hit. The average read access time in nanoseconds is __14 nsec__.

(i) What is Data Transfer rate (performance) of this memory system

(in words/sec)?

(ii) What is Bandwidth required of this memory system if word size is 8bit?

$$T_{avg} = 14\,nsec$$

(i) Data transfer Rate $= \dfrac{1}{T_{avg}}$ words/sec.

$$\Rightarrow \dfrac{1}{14 \times 10^{-9}} \text{ words/sec}$$

$$\Rightarrow \dfrac{1}{14} \times 10^9 \text{ words/sec}$$

$$\Rightarrow \dfrac{1000}{14} \times 10^6 \text{ words/sec}$$

$$\Rightarrow 71.42 \times 10^6$$

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

$$\Rightarrow 72 \text{ Millions Words/sec}$$

(OR)

$$71.42 \text{ Millions words/sec} \quad \text{Ans}$$

$$72 \text{ Millions Word}/sec.$$

$$\text{Word Size} = 8\,bit$$

(ii) Bandwidth = $72 \times 10^{6} \times \underline{8\,bit}/sec.$

$$= 576 \text{ Mbits}/sec$$

$$\Rightarrow 72 \times 10^{6} \text{ Byte}/sec$$

$$\Rightarrow 72 \text{ MBps} \quad @ \quad 71.43 \text{ MBps}.$$

**Q.2** 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 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 places) needed after the optimization such that it should not increase the average memory access time is _____.

$h = 0.8$

$(1-h) = 0.2$
$(1-0.8) \nearrow$

Cache Access time $(t_c) = 10\,nsec$

Miss Penalty (MM Access time) $t_m = 100\,nsec$ (MP)

$T_{avg} = t_c + (1-h)\,t_m^{(mp)}$

$= 10 + (1-0.8)\,100 = 10 + 0.2(100)$

$\boxed{T_{avg} = 30\,nsec}$

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

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

$= 8 + 0.2(110)$

$= 8 + 22$

$\boxed{T_{avg} = 30\,nsec}$

**Optimization:** $\Rightarrow$ Tavg & miss penalty Not Affected (Remain Same)

$$\boxed{t_{c_{new}} = 15\,nsec}$$

$h_{new} = ?$

(MP) miss penalty Not affected

ie $M.P/t_m = 100\,nsec.$

(Not affected) $Tavg_{new} = \underline{30\,nsec}$

$$\boxed{Tavg = t_{c_{new}} + (1-h_{new})(t_m)}$$

$30 = 15 + (1-h_{new})\,100$

$30 = 15 + 100 - 100\,h_{new}$

$85 = 100\,h_{new}$

$h_{new} = \dfrac{85}{100}$

$$\boxed{Tavg = h_{new} \times t_{c_{new}} + (1-h)(t_m + t_{c_{new}})}$$

$30 = h_{new} \times 15 + (1-h_{new})(100+15)$

$30 = 15h_{new} + 115 - 115h_{new}$

$85 = 100\,h_{new}$

$h_{new} = \dfrac{85}{100} = \boxed{0.85}\ \underline{Ans}$

$$\boxed{h_{new} = 0.85}\ \underline{Ans}$$

A direct mapped cache memory of 1 MB has a block size of 256 bytes. The cache has an access time of 3 ns and a hit rate of 94%. During a cache miss, it takes 20 ns to bring the first word of a block from the main memory, while each subsequent word takes 5 ns. The word size is 64 bits. The average memory access time in ns (round off to 1 decimal place) is 13.5 ns

**[GATE-2020]**

Block Size = 256 Byte
$t_C$ (Cache Access) = 3 nsec

1 word size = 64 bit ≈ 8 Byte

#Words = $\dfrac{256B}{8B}$ = **32 Words**

94%
h = 0.94

If Cache miss
first word = 20 nsec
mm to Cache

Each
Subsequent word = 5 nsec.
takes

$$T_{avg} = 0.94 \times 3 + (1-0.94)\left[3 + 20 + 31(5\,nsec)\right]$$
$$\Rightarrow 2.82 + 0.06\left[3 + 20 + 155\right]$$
$$= 2.82 + 0.06\left[178\right]$$
$$= \boxed{13.5\,nsec} \quad \text{Ans}$$

$$(\text{Max } \overset{\text{Data}}{\text{Transfer Rate}})$$

$$\text{Bandwidth} = 10 \text{ MByte}/\text{sec}$$

$$\text{In } 1 \text{ sec} = 10 \times 10^6 \text{ Byte } \underline{\text{Per Second}}$$

(Q) If the cycle time of the Memory is 500 nsec. then what is the Bandwidth? (The Maximum Rate which Memory Can Access _____ Byte/sec?)

(Sol$^n$) 500 nsec (cycle time) Access _____ 1 Byte

In (1 sec) _____ $\dfrac{1}{500 \times 10^{-9}}$ Byte/sec

$\Rightarrow \dfrac{1}{500} \times 10^9$ Byte/sec

$\Rightarrow \dfrac{1000 \times 10^6}{500}$ Byte/sec

$= 2 \text{ MBps}$ Ans

or 16 Mbits/sec

④

A certain processor deploys a single-level cache. The cache block size is 8 words and the word size is 4 bytes. The memory system uses a 60-MHz clock. To service a cache miss, the memory controller first takes 1 cycle to accept the starting address of the block, it then takes 3 cycles to fetch all the eight words of the block, and finally transmits the words of the requested block at the rate of 1 word per cycle. The maximum bandwidth for the memory requested block at the rate of 1 word per cycle. The maximum bandwidth for the memory system when the program running on the processor issues a series of read operations is ___160___ $\times 10^6$ bytes/sec. [GATE-2019-CS: 2M]

Ans (160)

Cache Block Size = 8 Words

1 Word Size = 4 Byte

Clock Frequency = 60 MHz

$$\boxed{\text{Cycle time} = \frac{1}{60 \times 10^6} \text{ sec}}$$

Cache Block Size = 8 Words

$\Rightarrow 8 \times 4$ Byte

$$\boxed{\text{Cache Block Size} = 32 \text{ Byte}}$$

$\frac{\text{Transfer}}{1 \text{ word}/1 \text{ cycle}}$

So 8 Word ⇓ 8 Cycle

Total Time Taken to Transfer Block = (Accept) 1 Cycle + (Complete 8 Words) 3 Cycle + 8 Cycle

= 12 Cycle

$= 12 \times \frac{1}{60}$ usec $= \frac{1}{5} \times 10^{-6}$ sec

$$\text{Cycle time} = \frac{1}{60 \times 10^6} \text{ sec.}$$

$$\underline{12 \text{ Cycle}} \longrightarrow 32 \text{ Byte.}$$

$$12 \times \frac{1}{60 \times 10^6} \text{ sec} \longrightarrow 32 \text{ Byte}$$

$$\text{In } \underline{1 \text{ Sec}} \longrightarrow \frac{32 \text{ Byte}}{12 \times \frac{1}{60 \times 10^6}} \Rightarrow \frac{32B \times 60 \times 10^6}{12} \text{ Byte/sec}$$
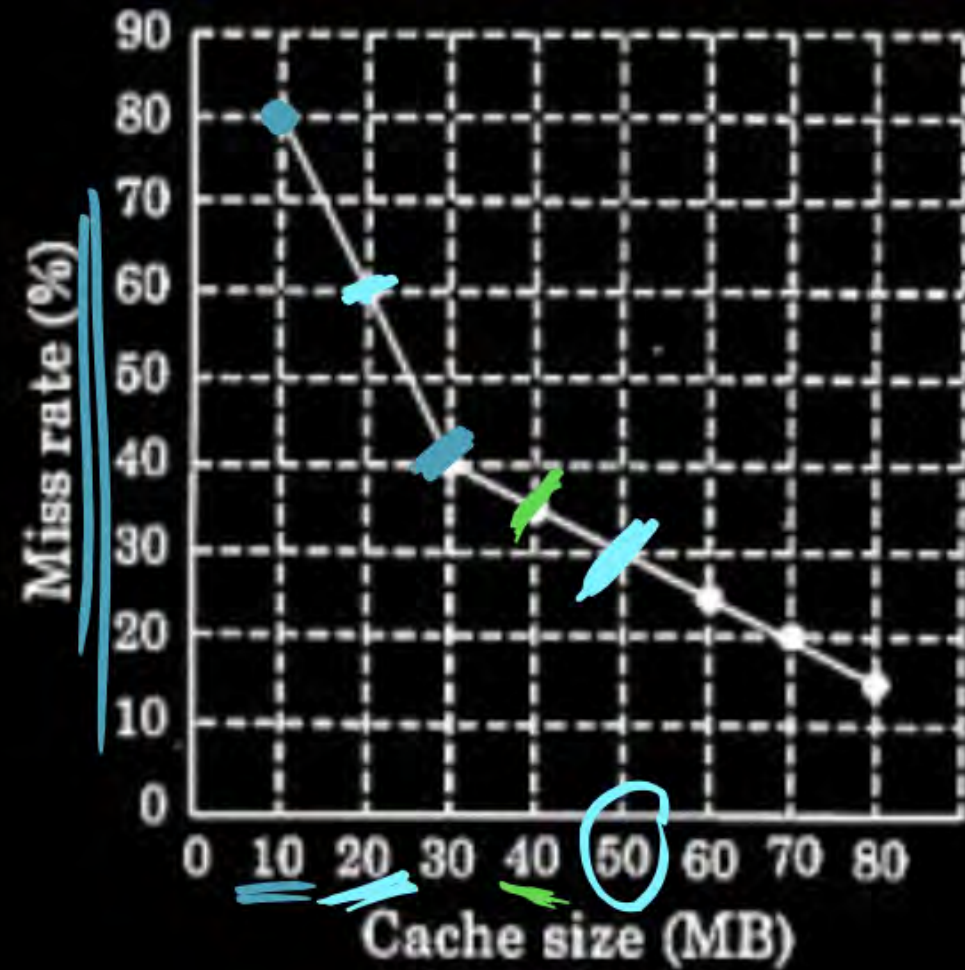
$$= 160 \times 10^6 \text{ Byte/sec} \underline{\text{ Ans}}$$

A file system uses an in-memory cache to cache disk blocks. The miss rate of the cache is shown in the figure. The latency to read a block from the cache is 1 ms and to read a block from the disk is 10 ms. Assume that the cost of checking whether a block exists in the cache is negligible. Available cache sizes are in multiples of 10 MB.



Miss rate (%) vs Cache size (MB)

The smallest cache size required to ensure an average read latency of less than 6 ms is

30 MB. Ans

[GATE-2016(Set2)-CS: 2M]

Cache Access time $[t_c]$ = 1 msec.

Main Memory Access time $(t_m)$ = 10 msec.    Hit Rate = 0.2.

When Cache Size 10MB then Miss Rate 80% [0.8]

$T_{avg} = h * t_c + (1-h)(t_m + t_c)$    ⇒ 0.2 + 0.8(11)

$= 0.2 \times 1 + 0.8(10+1)$    = 0.2 + 8.8 = 9 msec

<u>Case I</u> When Cache Size = 20mB $\qquad$ $\boxed{h = 0.4}$

then Miss Rate $(1-h) = 60\%$ $(0.6)$

$T_{avg} = h * t_c + (1-h)(t_m + t_c) = 0.4 \times 1 + 0.6[10+1] \Rightarrow 0.4 + 6.6$

$\boxed{T_{avg} = 7 \, msec}$ When Cache Size is 20mB.

<u>Case II</u> : CacheSize = 30mB then Miss = 40% $\underline{Hit = 0.6}$

$T_{avg} = h t_c + (1-h)(t_m + t_c) \Rightarrow 0.6 \times 1 + 0.4 \overset{(1+10)}{(11)} = 0.6 + 4.4 = 5 msec$

$\boxed{T_{avg} = 5 msec}$ When Cache Size is $\underline{30mB}$.
↳ smallest

Case IV :    Cache Size = 40MB  then  Miss Rate 35%.    Hit = 0.65

$T_{avg}$ = 0.65 × L + 0.35 ($\perp\perp$)

$\boxed{T_{avg} = 4.5\,msec}$  when  Cache is  40MB

Case V :  when Cache Size = 50MB  then Miss = 0.3  $\boxed{Hit = 0.7}$

$T_{avg}$ = 0.7 × 1 + 0.3 ($\perp\perp$)

         = 0.7 + 3.3

$\boxed{T_{avg} = 4\,msec}$  when Cache is 50MB.