# COMPUTER SCIENCE

## Computer Organization and Architecture

### Cache Memory

Lecture_02

Vijay Agarwal sir

Memory Hier.

Type of Access.

# Memory

Cache Miss ①

MM Miss (Page fault) ③

Words

Blocks

Pages

**CPU** reg ──── **Cache** ──── **MM** ──── **S M**

⑥

⑤

④

Managed by H/W

MM Hit (Page Hit)

Physical address

Hit Ratio = 1

Virtual / logical address

CO

Managed by OS

# Memory



CPU reg — Words — **Cache** — Blocks — MM — Pages — S M

Managed by H/W

Physical address

Virtual/ logical address

Mapping → CO

Paging → Managed by OS
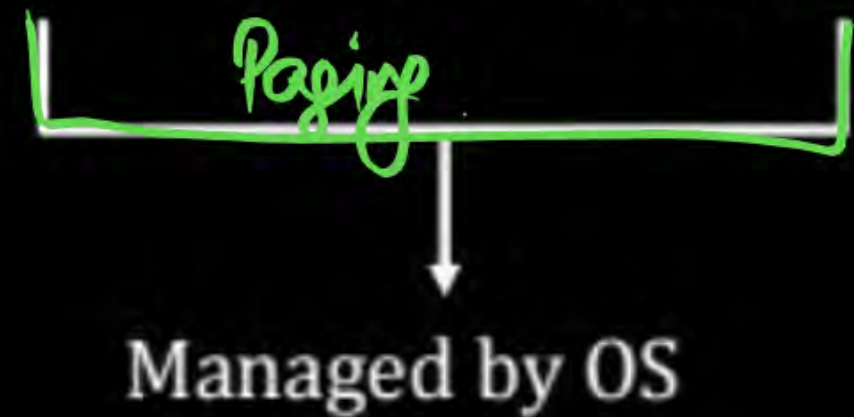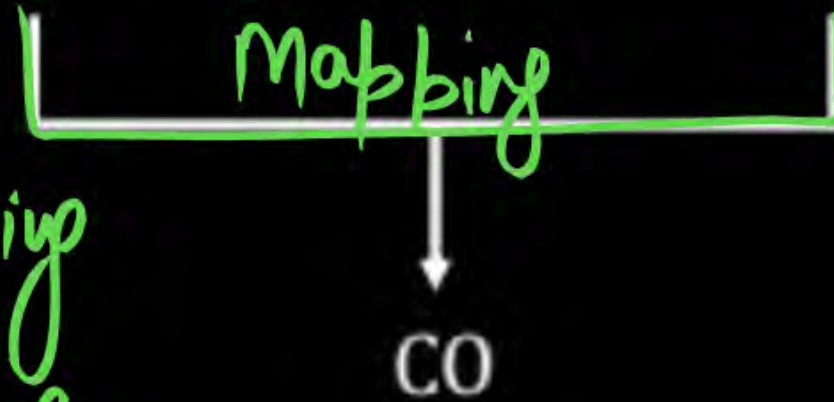
① Direct Mapping
② Set Associative Mapping
③ Fully Associative Mapping

**Q.** Calculate the average Access time, when the CPU request for the memory 100 times, out of 100 times, 90 times hit & 10 Time miss. If time taken when there is a hit(Each hit) is 20ns & time taken when there is a Miss(Each Miss) is 150ns. ?
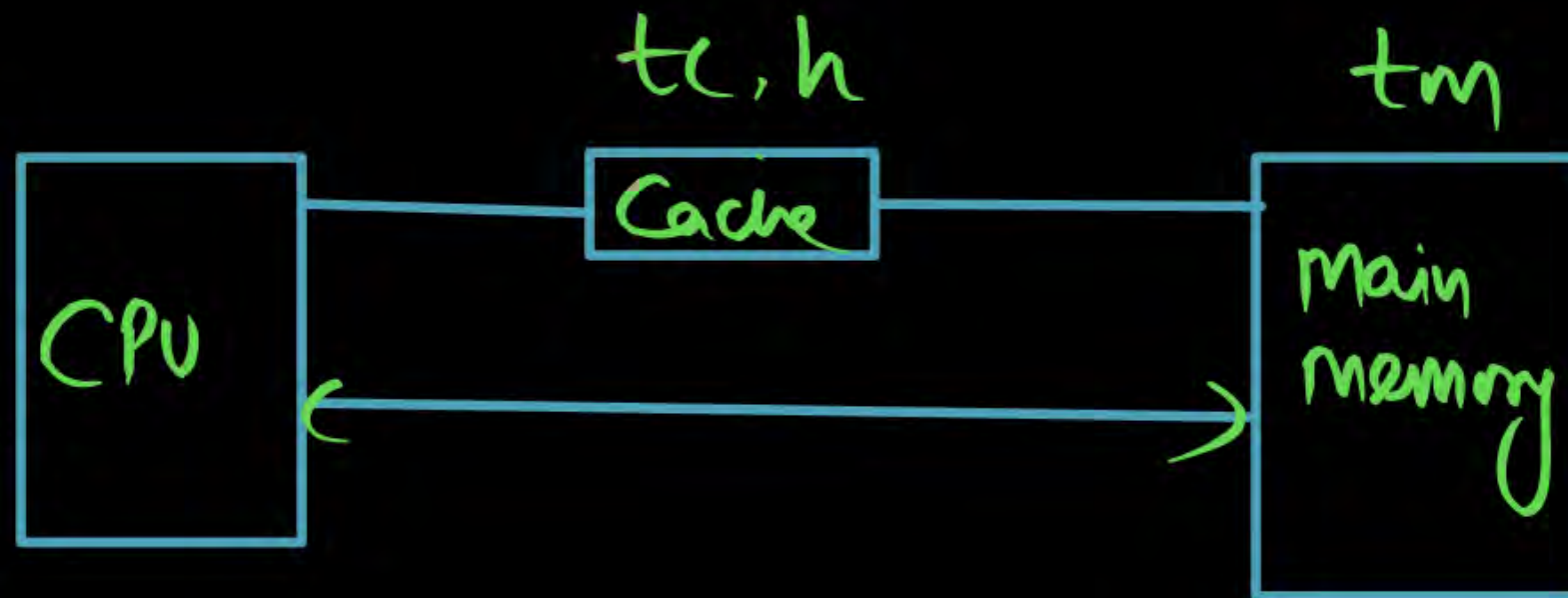
# Type of Memory Org

h: Cache Hit Ratio.

tc : Cache Access time

tm: MM Access time

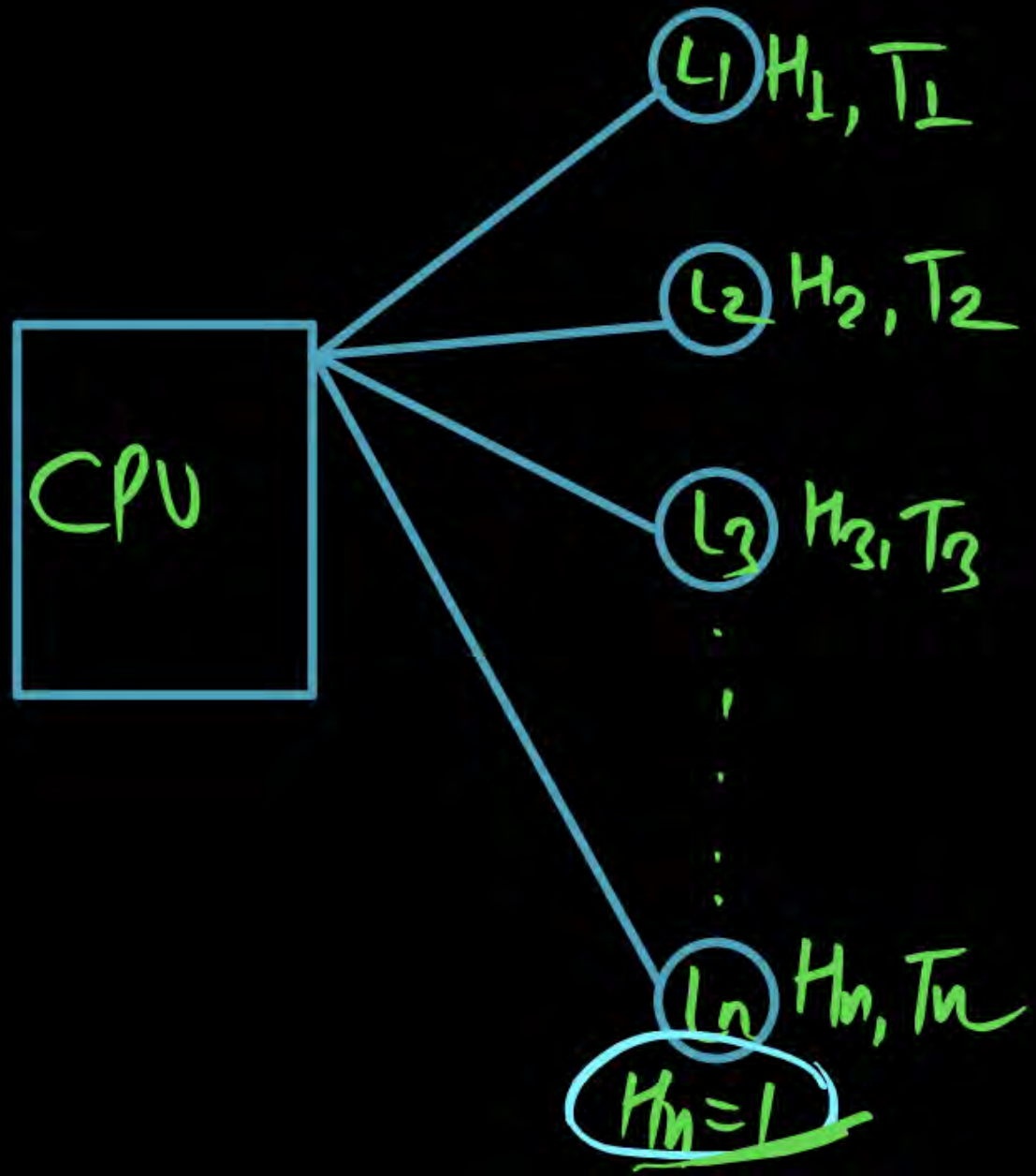1. Simultaneous Access Memory Org. : 2 level

tc, h

tm

Cache

CPU

main memory

1 Word Access time Tavg.

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

# Type of Memory Org

1. Simultaneous Access Memory Org. $n$ Level



$$T_{avg} = H_1 T_1 + (1-H_1) H_2 T_2 + (1-H_1)(1-H_2) H_3 T_3$$

$$+ \ldots (1-H_1)(1-H_2)(1-H_3)\ldots(1-H_{n-1}) H_n T_n$$

$L_1$   $H_1, T_1$

$L_2$   $H_2, T_2$

CPU

$L_3$   $H_3, T_3$

$L_n$   $H_n, T_n$

$H_n = 1$

# Type of Memory Org

1. Simultaneous Access Memory Org.

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

Data transfer Rate
$$\left(\frac{\#\,Words}{sec}\right) = \frac{1}{T_{avg}}$$
(Performance)
(efficiency)

## Remember

$$\text{Avg Instn ET} = 5.51 \times 10^{-9} \text{ sec.}$$

$$1 \text{ Instn ET} = 5.51 \times 10^{-9} \text{ sec.}$$

In 1 Sec ⇒ How Many # Instn

$$\text{In 1 sec} = \frac{1}{5.51 \times 10^{-9}} \text{ Instn/sec}$$

$$\Rightarrow \frac{1}{5.51} \times 10^{9} \text{ Instn/sec}$$

$$\Rightarrow \frac{1000 \times 10^{6}}{5.51} \Rightarrow 181.4 \times 10^{6}$$

$$= 181.4 \text{ MIPS}$$

# Type of Memory Org
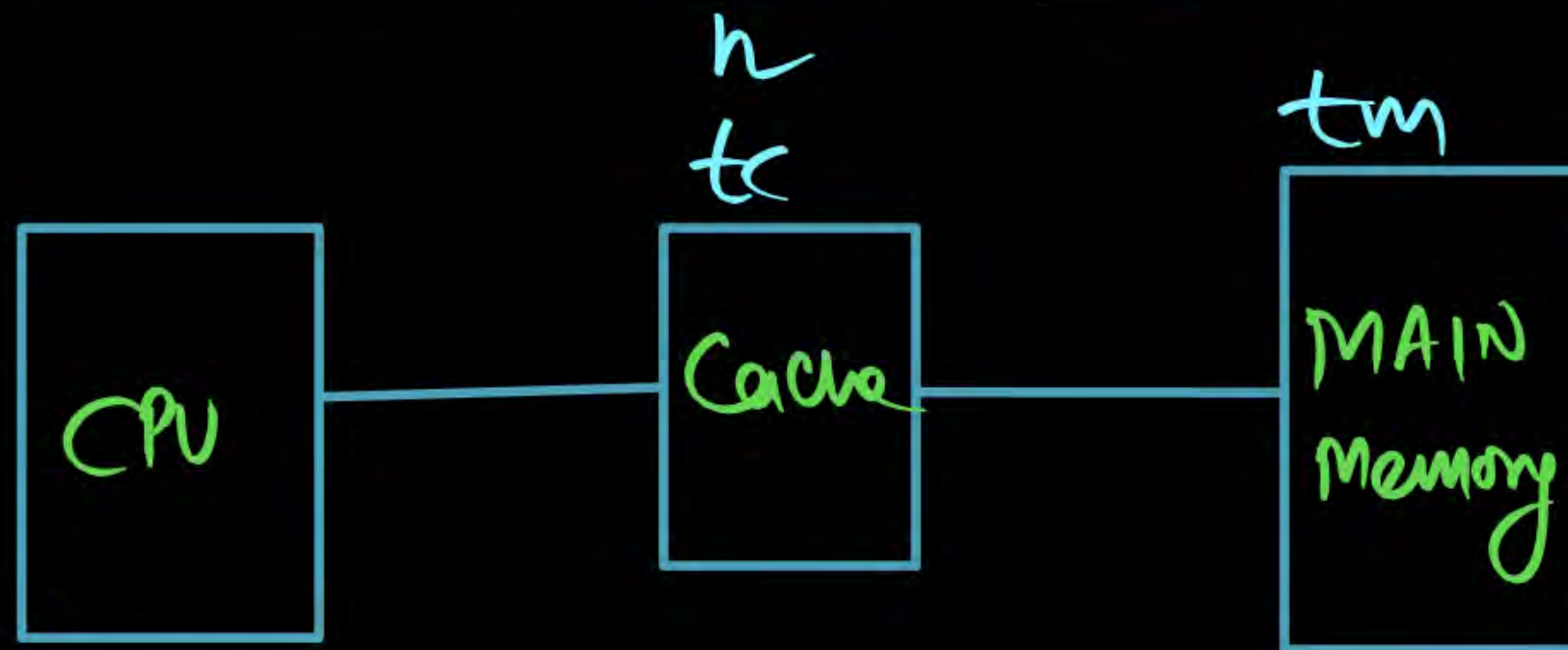
1. Simultaneous Access Memory Org.

2. <u>Hierarchical Access Memory Org.</u> : <u>2 Level</u>



$h$

$t_c$

$t_m$

CPU — Cache — MAIN Memory

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

<u>Hierarchical Access:</u>

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

# Hierarchical Access

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

$$\Rightarrow h\!\!\!/t_c + t_m - ht_m + t_c - h\!\!\!/t_c$$

$$\Rightarrow t_c + t_m - ht_m$$

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

**Q)** Hierarchal Access    Hit Ratio 80%,    $t_c = 20ns$.  MM = 100nsec

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

$\Rightarrow 0.8 \times 20 + (1-0.8)[100+20]$

$\Rightarrow 16 + 0.2(120)$

$\Rightarrow 16 + 24$

$$\boxed{T_{avg} = 40 \, nsec} \quad \underline{Ans}$$

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

$\Rightarrow 20 + (1-0.8)100$

$= 20 + (0.2)100$

$\Rightarrow 20 + 20 \Rightarrow \boxed{40 nsec}$

---

$T_{avg}$.

$t_c + t_m$

$(1-0.8)[20+100]$

$\Rightarrow 0.8 \times 20 + 0.2[20+100]$

$\Rightarrow \boxed{0.8 \times 20 + 0.2(20)} + 0.2(100)$

$\Rightarrow 20 + 0.2(100)$

$\boxed{40 nsec} \quad \underline{Ans}$

$h = 0.8$

$(1-h) = 1 - 0.8$

$\quad = 0.2$

## Type of Memory Org

2. Hierarchical Access Memory Org. **3 Level**

$(1-h_2)$: Miss Ratio in level $(M_2)$

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

$$T_{avg} = h_1 t_1 + M_1 h_2 (t_2 + t_1) + M_1 M_2 h_3 (t_3 + t_2 + t_1)$$

$M_1$: Miss Ratio in Level 1

$M_2$: Miss Ratio in Level 2.

# Type of Memory Org

2. Hierarchical Access Memory Org.   **3 Level**

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

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

(Performance)
(Data transfer Rate)

# Type of Memory Org

**2. Hierarchical Access Memory Org.** : *n Level*



$$T_{avg} = H_1 T_1 + (1-H_1) H_2 (T_2 + T_1) + (1-H_1)(1-H_2) H_3 (T_3 + T_2 + T_1)$$
$$+ \ldots (1-H_1)(1-H_2)(1-H_3) \ldots (1-H_{n-1}) H_n (T_n + T_{n-1} + \ldots T_3 + T_2 + T_1)$$

**Q.** Calculate the average Access time with the cache access time 1ns, and main memory access time 100ns, Hit ratio 90%?
Using Hierarchical Access?

$$11 \text{ nsec} \quad \text{Ans}$$

**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. Let level 1 Access time is 20ns, What is the hit ratio? Using simultaneous Access org?

**Q.** Consider a system with 2 levels. Level 1 Access time is 20ns Level 2 Access time $= 150$ns $T_{avg} = 30$ using simultaneous Access.

(i) What is the Hit Ratio?

(ii) If the Hit Ratio is made to 100% then what is the Access time of $L_1$ & $L_2$ Memory?

→ Hit Ratio Not Effect levels Access time Only $T_{avg}$ change.

$T_1 = 20$ns

$T_2 = 150$ns

$T_{avg} = 30$ns

**Simultaneous**

$T_{avg} = ht_1 + (1-h)t_2$

$30 = h*20 + (1-h)150$

$30 = 20H + 150 - 150H$

$120 = 130H$

$H = \dfrac{120}{130} = 0.923$

$\Rightarrow \boxed{92.3\%}$

Ans

(ii) If Hit Ratio Made to 100% then Access time of $L_1$ & $L_2$ Memory.

$\boxed{\begin{array}{l} T_1 = 20\text{ns} \\ T_2 = 150\text{ns} \end{array}}$ $L_1$ & $L_2$ Access Remain Same.

$\boxed{h=1}$ $\boxed{(1-h)}$ $= 0$

If Hit Ratio = 100%

$T_{avg} = ht_1 + (1-h)t_2$

$\Rightarrow 1 \times 20 + (1-h)150$

$T_{avg} = 20$ns

Only $T_{avg}$ change.

**Q.** If the above Question if $T_{avg}$ is increased by 10% then what is % of change in Hit Ratio?

$T_1 = 20 \quad T_2 = 150, \qquad T_{avg} = 30 \Rightarrow$ Now $T_{avg}$ increased by $10\% \Rightarrow 30 + 10\% \text{ of } 30$

$$\Rightarrow 30 + 3$$

$$\boxed{T_{avg_{new}} = 33\,msec} \quad = 33\,nsec$$

$$T_{avg_{new}} = h * t_1 + (1-h)t_2$$

$$33 = h * 20 + (1-h)150$$

$$33 = 20h + 150 - 150h$$

$$130h = 117$$

$$H = \frac{117}{130}$$

$$\boxed{H = 0.9}$$

$\boxed{\text{Now Hit Ratio} = 90\%}$

Previous Hit Ratio $= 92.3\%$

Hit Ratio Decreased by

$\boxed{2.3\%} \downarrow$

$92.3 - 90 = 2.3\%$

**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 _14 nsec_ Ans

$$T_{avg} = h*5 + (1-h)*50$$

$$= 0.8 \times 5 + (1-0.8) \times 50$$

$$= 4 + 10$$

$$\boxed{T_{avg} = 14 \, nsec} \quad \underline{Ans}$$

# NAT

(P W)

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?

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

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

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

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

$\Rightarrow 71.42 \times 10^{6}$ Words/sec

$\Rightarrow 71.42$ millions word per sec.

$\Rightarrow 72$ Millions word per Sec.

72 Millions word Per Sec

(ii) Bandwidth ⇒ Word Size = 8 bits

⇒ 72 million word Per Sec ⇒ $72 \times 10^6$ word Per Sec.

⇒ $72 \times 10^6 \times$ (8 bits) Per Sec ⇒ $72 \times 8$ M bits/sec

⇒ 576. M bits/sec ⇒ $\dfrac{576 \text{ Mbit}}{8 \text{ bits}}$ Byte [1 Byte = 8 bit]

OR

⇒ 72 MBps Ans
   Byte

Cache Work on Locality of Reference.

# Locality of Reference [LOR]

Accessing the Higher Level of Memory Data from the Level 1 Memory (Cache Memory) is Called Locality of Reference. (Faster Memory)

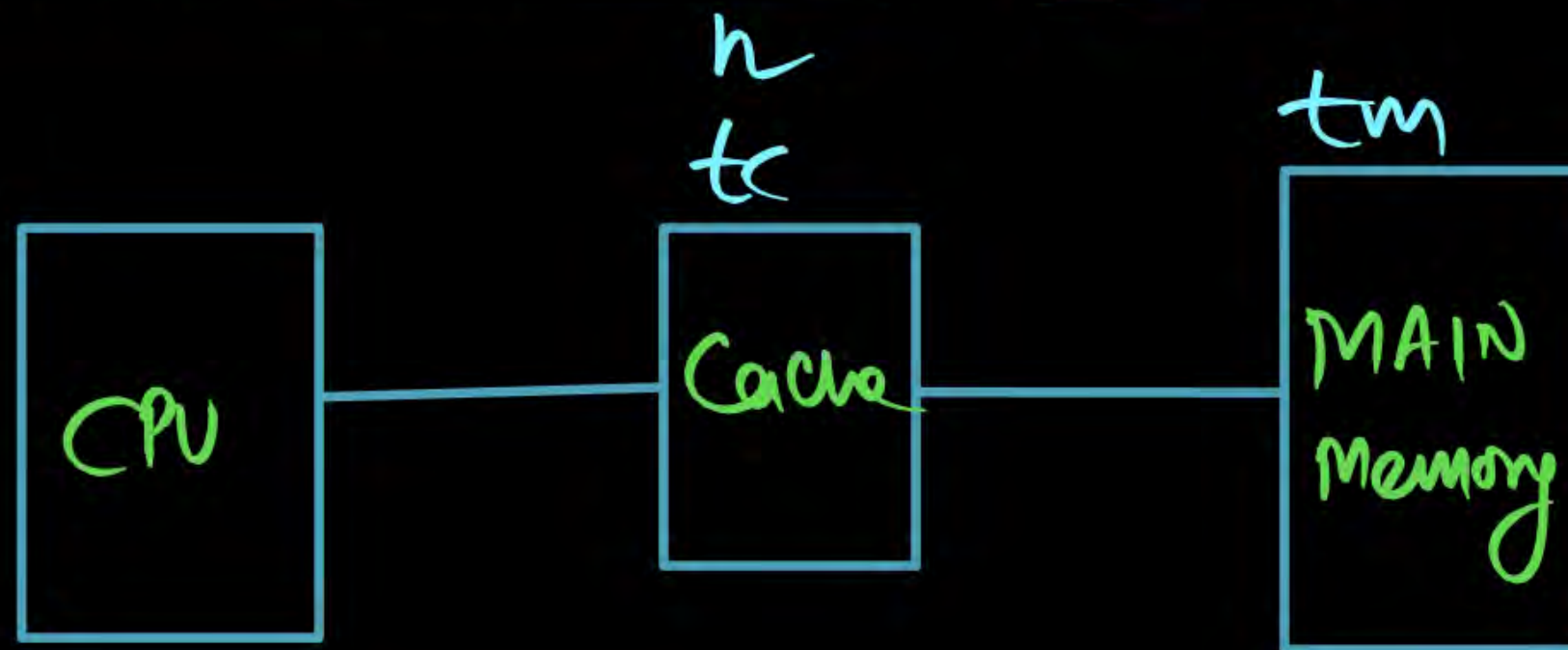# Locality of Reference [LOR]

① Temporal LOR

② Spatial LOR

## 2. Hierarchical Access Memory Org. : 2 Level



CPU ── Cache ── MAIN Memory

$h$
$t_c$
$t_m$

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

### Hierarchical Access:

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

# Memory

Cache Miss

MM Miss (Page fault)

① ② ③

Words — Blocks — Pages

**CPU** reg — **Cache** — MM — S M

⑥ Locality of Reference ⑤ MM Hit (Page Hit) ④

Managed by H/W

Physical address

Virtual / logical address

Hit Ratio = 1

CO

Managed by OS

## Non Technical eg

**LOR:** Medical Store

Only 1 Time go to Medical Store Buy
One Patta of
(strips)
So Next time if Same Tablet 10/20 Tablet.
Require then No Need to go
Medical Store, take it from Home.

# Locality of Reference:

CPU —— / Word. —— Cache —— / Block —— MAIN Memory.

1 Block Size = 32 Byte / 32 Words

If Miss in Cache then 1 Complete (32 Word / 32 Byte) Block is Transferred From Main Memory to Cache. the Respective Word (which is Requested / Demanded by CPU) give from Cache Memory to CPU.

( Next Reference there is 0 Cache Hit ).

# Locality of Reference:

CPU — | Word. — Cache — | Block — MAIN Memory.

1 Block Size = 32 Byte / 32 Words

⊗ Here 1 Block Size = 32 Word / 32 Byte & CPU Require Only 1 Word. ?

(Sol^n) In this Process 1 Complete Block of 32 Words / 32 Byte transfered from Main Memory to Cache Memory. then Demanded (Requested) 1 Word transfered from Cache Memory to CPU.

(Note) So in the Next time When CPU Request Some Word (or) Adjacent Word then that Request Available in Cache. [Cache Hit].

# Locality of Reference [LOR]

# LOR [Locality of Reference]

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

*Faster Memory [Cache Memory]*

(1) Temporal LOR

*eg (32 Word)*

*Cache*

$M(2000)$ $[n]$ *Memory Location*
*Word Available in $B_5$.*

(2) Spatial LOR.

$B_5 \leftarrow M[2000]$

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

Or

$M(2000)$ *location wallah word Again & Again*

Same data which access again and again then that type of data stored in Temporal LOR.
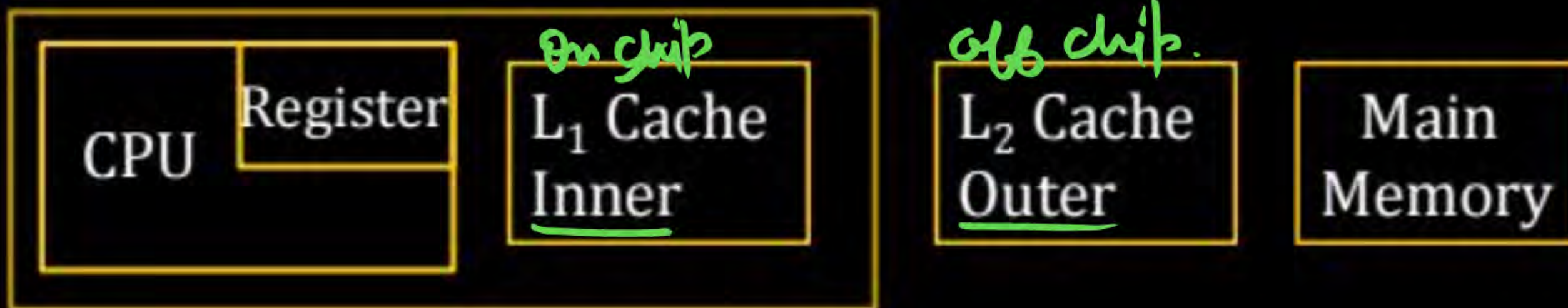
## LOR [Locality of Reference]

(2) Spatial LOR means adjacent word in the same block is

referenced by the CPU in a sequence.

Memory $(n+1)$ Wallah Word
Location

Which is Available in B5.

Cache hit

# Types of Cache

1) <u>Unified Cache</u>:   <u>Instruction & Data both are placed in Same Cache.</u>

2) <u>Split Cache</u>:   <u>This Cache logically Divide into two parts</u>

                         (i) Instruction Cache [<u>I</u> - cache]

                         (ii) Data Cache [<u>D</u>- cache]

3) <u>Multilevel Cache</u>:

| CPU | Register | On chip<br>$L_1$ Cache<br><u>Inner</u> | off chip.<br>$L_2$ Cache<br><u>Outer</u> | Main<br>Memory |
|---|---|---|---|---|

Size $L_1 < L_2$
Speed $L_1 > L_2$

# 2 Level of Memory (If Locality of Reference Included)



CASE
① Block Size = 1 word
CASE Ⅱ Block Size = n words.

CPU — words — L₁ memory (CM) — Blocks — L₂ memory (MAIN memory) — Block

CPU Always Access the Data from the faster Memory (Level1 | Cache Memory). If there is Miss in Level1 (Cache) Memory then One Complete Block is transfered From Level2 (Slow) Memory to Level1 [Cache] Memory & addressed word (which is Requested | Demanded by CPU) that Respective word given from Cache Memory to CPU.

# 2 Level of Memory (If Locality of Reference Included) PW



CPU — Words. — $L_1$ Memory (CM) — Blocks — $L_2$ Memory / Block — MAIN memory.

Case I: Block Size: 1 Word

CASE II: Block Size = n words

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

① CASE I: Block Size = 1 Word

$$TB = T_2$$ Ans

② CASE II: Block Size = n Words

$$TB = n * T_2$$ Ans

# 2 Level of Memory (If Locality of Reference Included)



CPU --- Words. --- L1 memory (CM) --- Blocks --- L2 Memory (MAIN memory)

Block

$t_1$

$T_B$

# 2 Level of Memory (If Locality of Reference Included)



CPU — Words. → L1 memory (CM) — Blocks → L2 Memory (MAIN memory)

## If. 2 Level Memory.

$$h_2 = 1$$

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

(OR)

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

## If Locality of Reference Included

$(T_B = n * T_2)$

$$T_{avg} = h\,t_1 + (1-h)\left[T_B + T_1\right]$$

(OR)

$$T_{avg} = t_1 + (1-h)\,T_B.$$

**If 3 Level Memory**

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

Hit Ratio of
$h_3 = 1$ [last level memory always]

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

**If Locality of Reference included.**

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

TB₁ → $TB_1$

TB₂ → $TB_2$

CPU

1 Word
$(T_1)$
CPU Request

$L_1$

Block

$L_2$

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? $\boxed{h_3 = 1}$

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

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

(iii) If there is a miss in level1 and Level 2 & hit in level3 $\boxed{h_1 \& h_2 = 0}$ $\boxed{h_3 = 1}$

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

(OR)

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

∴

$$\boxed{\overline{T}avg = h_1 t_1 + (1-h_1) h_2 (t_2 + t_1) + (1-h_1)(1-h_2) \atop (t_3 + t_2 + t_1)}$$

$$\underset{h_2}{}$$ ← $c_{91}$

**If locality of Reference Included.**

$$\boxed{\overline{T}avg = h_1 t_1 + (1-h_1) [TB_1 + T_I] + (1-h_1)(1-h_2) [TB_2 + TB_L + T_I)}$$

$$\underset{h_3}{}$$

**Q.1**

**Soln 1**

$h_1 = 100. \Rightarrow h_1 = 1$

$(1-h_1) = 0.$

$0 * n = 0$

$$\boxed{\overline{T}avg = T_I}$$

$h_1 = 100 \Rightarrow h_1 = 1$

$(1-h_1) = 0$

$\Rightarrow 0 * n = 0$

$$\boxed{\overline{T}avg = T_L}$$

**Soln 2**

$h_1 = 0 \quad \& \quad h_2 = 100 \Rightarrow \boxed{h_2 = 1}$

$$\boxed{1 - h_2 = 0}$$

$$\boxed{\overline{T}avg = T_2 + T_L} \quad \text{Ans}$$

**Soln 2**

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

$$\boxed{\overline{T}avg = TB_L + T_L} \quad \text{Ans.}$$

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

(Soln)  $h_1 = 0$  $h_2 = 0$  $\boxed{h_3 = 1}$

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

$$T_{avg} = h_1 t_1 + (1-h_1) h_2 \left[ TB_1 + T_1 \right] + (1-h_1)(1-h_2) \overset{h_3}{\left[ TB_2 + TB_1 + T_1 \right)}$$

(Soln)  $h_1 = 0$  $h_2 = 0$  $h_3 = 1$

$$\boxed{T_{avg} = TB_2 + TB_1 + T_1} \underline{Avg}$$

**Q.** 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?

$h_1 = 90\% \Rightarrow h_1 = 0.9$

$t_1 = 30 nsec$

$t_2 = 250 ns/word$

$$T_{avg} = T_1 + (1-h) TB_L$$

$TB_1 = 4 \times 250 \, ns/word$
$\phantom{TB_1 = 4 \times} word$

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

$$T_{avg} = h_1 t_1 + (1-h_1)(TB_1 + T_1)$$

$TB_1 = 1000 nsec$

(OR)

$\Rightarrow 30 + 100$

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

$\Rightarrow 27 + 103$

$\phantom{\Rightarrow} = 130 nsec$

$$T_{avg} = 130 nsec \quad \underline{Ans}$$

**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 _**14 nsec**_ **[GATE – 2015]**

Already Done.

## NAT

(P W)

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 __14nse__

(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?

Ans:

(i) 71.4 ⇒ 72 Million word Per Sec.

already Done in Today class Kindly
Refer.

Soln
(ii) 72 MBps.

**Q.** 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 _____.

$h = 0.8, \quad t_c = 10\,nsec$

miss Penalty $(t_m) = 100\,nsec$

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

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

$= 8 + 0.2(110)$

$= 8 + 22$

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

(OR)

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

$\Rightarrow 10 + (1-0.8)\,100$

$= 10 + 20$

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

Now optimization: $t_{c\,new} = 15\,nsec.$

tm & tavg Remain Same (Not Changed)

$Tavg_{new} = h * tc + (1-h)(tm+tc)$

$30 \Rightarrow h * 15 + (1-h)(100+15)$

$30 = 15H + 115 - 115H$

$85 = 100H$

$H = \dfrac{85}{100}$

$\boxed{H = 0.85}$ Ans

$Tavg_{new} = t_{c\,new} + (1-h)\,tm$

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

$30 = 15 + 100 - 100H$

$85 = 100H$

$H = \dfrac{85}{100} = 0.85$

$\boxed{H = 0.85}$ Ans

# NAT

(P)(W)

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 ____.

[GATE-2020]

Block Size = 256 Byte.
1 word Size = 64 bit ≈ 8 Byte.

$$\text{Number of Words} = \frac{256\,\text{Byte}}{8\,\text{Byte}} = \boxed{32\,\text{Words}}$$

1 Word takes = 20 nsec.

Remaing 31 Word = 5ns

$$T_{avg} = \underset{time}{h \times Cache} + \underset{miss}{(1-h)} \underset{20}{[Cache + MM} + 31(5)]$$

1st 1 word

$$T_{avg} = 0.94 \times 3 + (1-0.94) [3 + 20 + 31(5)]$$

$$\Rightarrow 2.82 + (0.06)(3 + 20 + 155)$$

$$\boxed{T_{avg} = 13.5\,nsec.} \quad \underline{Ans}$$

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 _____$\times 10^6$ bytes/sec.                   [GATE-2019-CS: 2M]