

# COMPUTER SCIENCE



## Computer Organization and Architecture

### Cache Memory

Lecture\_02

Vijay Agarwal sir





An orange diamond-shaped sign with a black border, mounted on a white pole. Below the sign is a construction barrier with two orange lights on top.

TOPICS  
TO BE  
COVERED

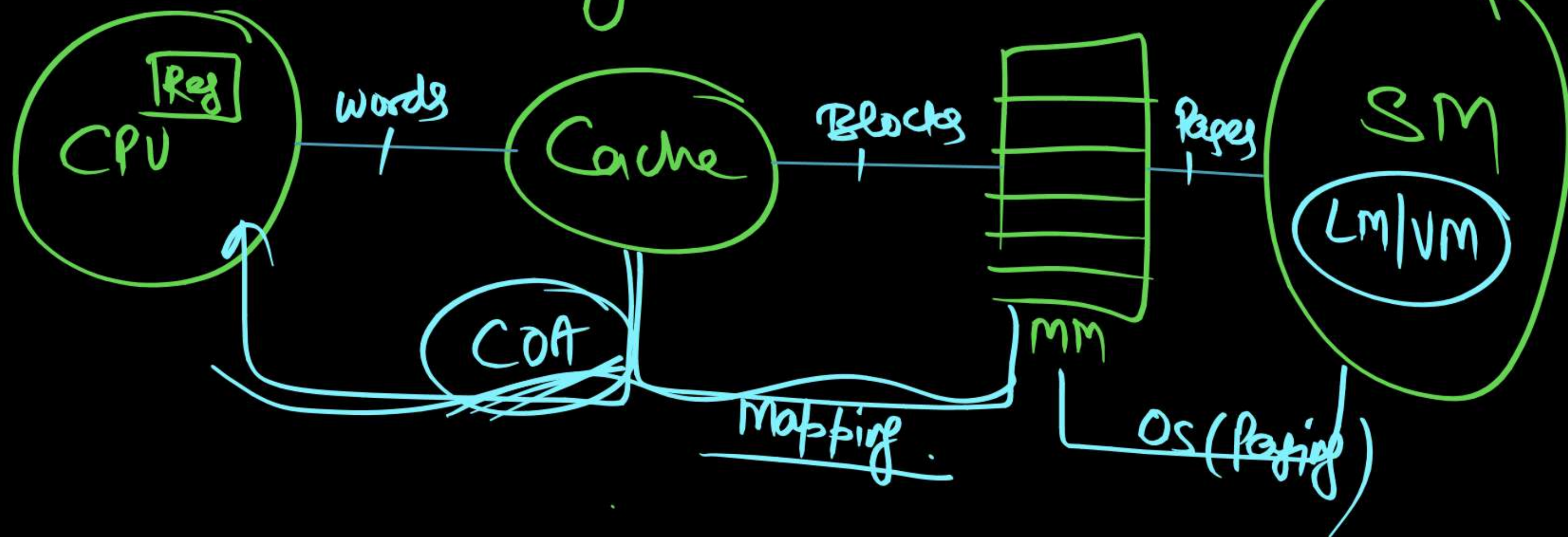
o1

Memory Access

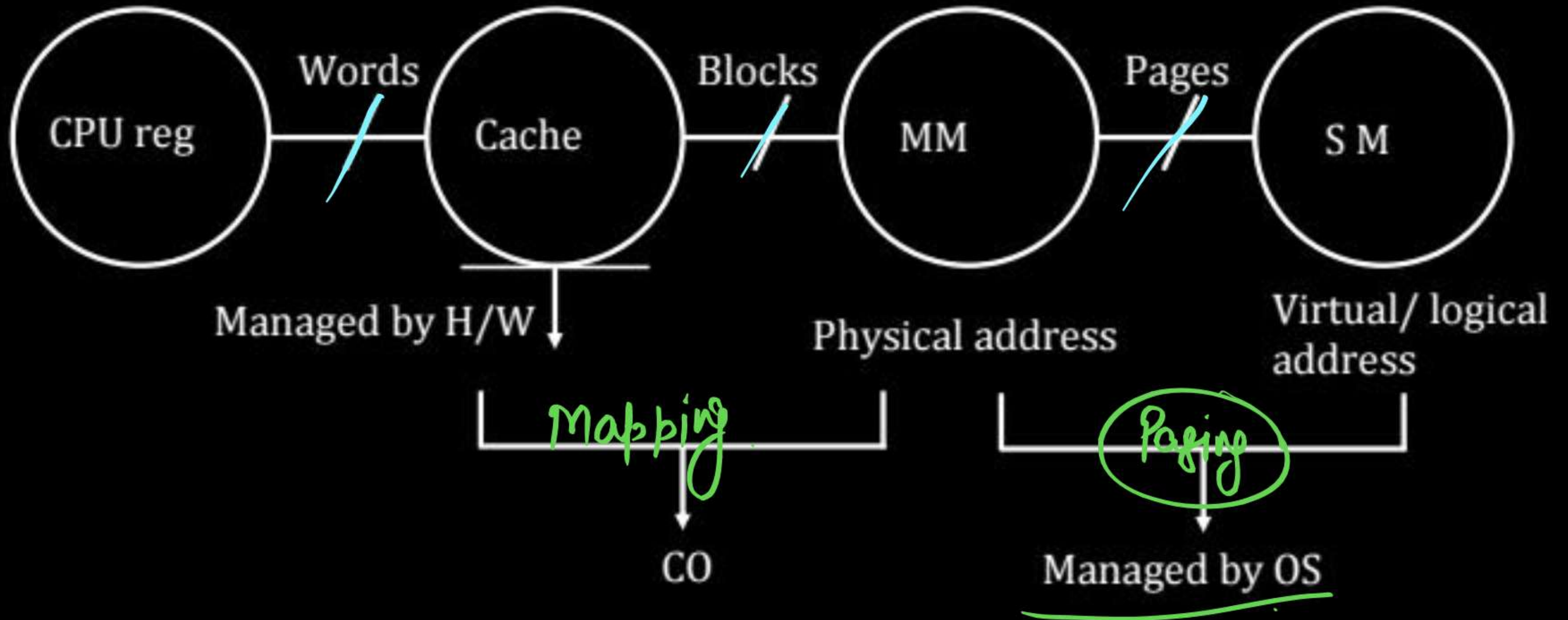
o2

Cache Memory

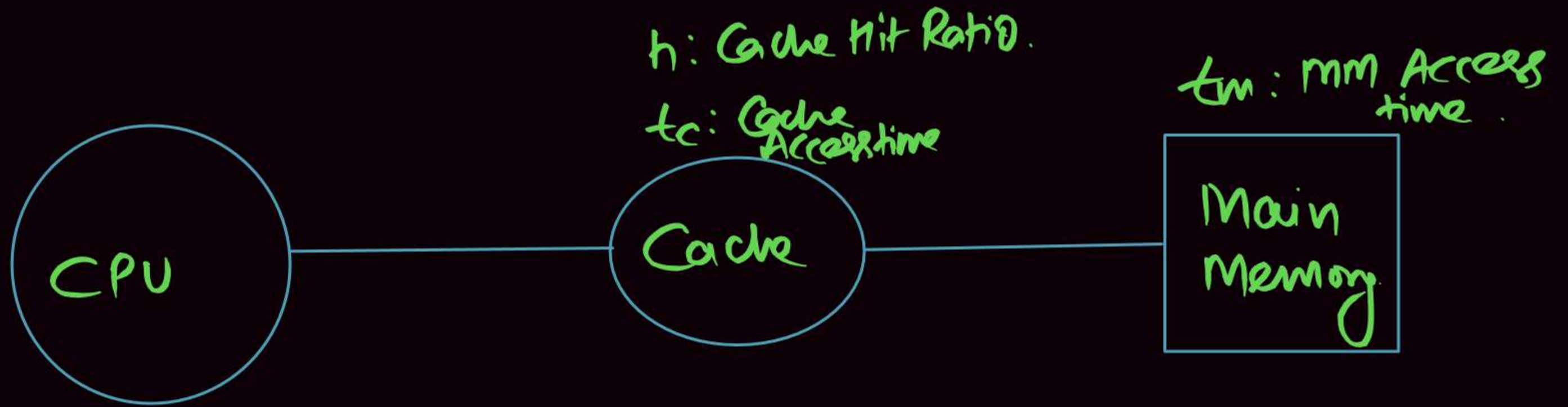
# Memory Hierarchy

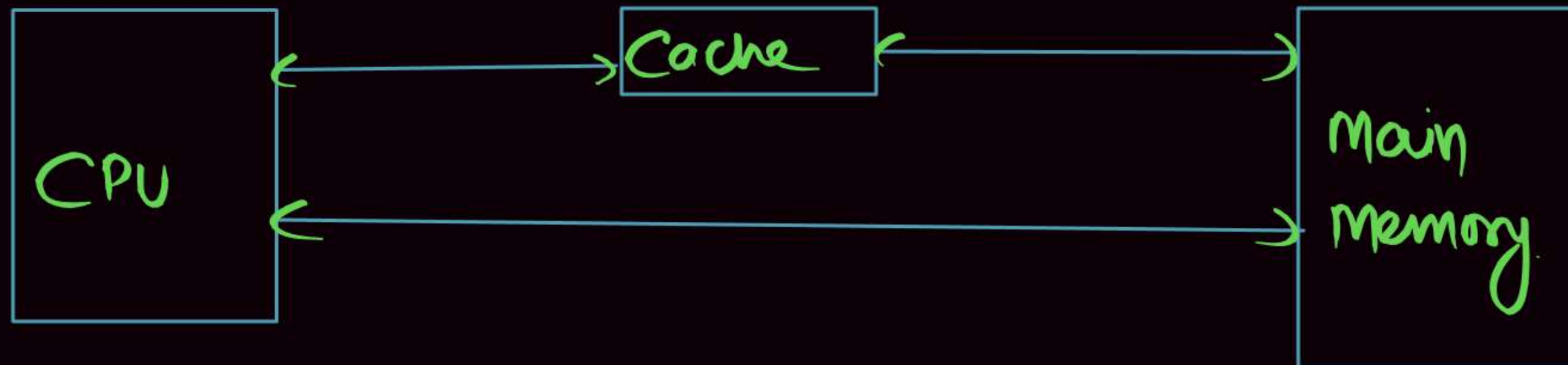


# Memory









Average Memory Access time :  $[T_{avg}]$

$$\boxed{T_{avg} = \underbrace{Hit}_H \times \text{Time taken by memory when there is a Hit} + (1 - H) \times \text{Time Taken by memory when there is a miss.}}$$

# Hit

$$\text{Hit Ratio } [H] = \frac{\text{Total \# Hit}}{\text{Total \# Access.}}$$





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. ?



Total CPU Request = 100  
# times hit = 90 ; Hit take time of 20ns  
# Time Miss = 10 When there is a miss then 150nsec

$$\text{Total time} = 90 \times 20 + 10 \times 150$$
$$= 1800 + 1500$$

$$\text{Total Time} = 3300 \text{ nsec}$$

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

$$\text{Total CPU Request} = 100$$

$$\text{Total \# Hit} = 90, \text{ Total Miss} = 10$$

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

[H]

$$\text{Miss Ratio} = \frac{10}{100} = 0.1$$

OR

$$(1-H) = 1 - 0.9 = 0.1$$

$$T_{avg} = H \times \text{Time taken by memory for each hit} + (1-H) \times \text{Time taken by memory for each miss}$$
$$\Rightarrow 0.9 \times 20 + (1 - 0.9) \times 150$$
$$= 18 + 15 = 33 \text{ nsec}$$





Calculate the average Access time, when the CPU request for the memory 100 times, out of 400 times, 300 times hit & 100 times miss. If time taken when there is a hit (Each hit) is 20ns & time taken when there is a Miss (Each Miss) is 150ns. ?

$$\text{Total CPU Request} = 400$$

$$\# \text{ times Hit} = 300; \text{ Hit take time of } 20\text{ns}$$

$$\# \text{ Time Miss} = 100 \quad \text{When there is a miss then } 150\text{nsec}$$

$$\begin{aligned} \text{Total time} &= 300 \times 20 + 100 \times 150 \\ &= 6000 + 15000 \end{aligned}$$

$$\text{Total Time} = 21000 \text{ nsec.}$$

$$T_{avg} = \frac{21000}{400} = \frac{210}{4} = 52.5 \text{ nsec}$$

$$\text{Total CPU Request} = 400$$

$$\text{Total \# Hit} = 300, \text{ Total Miss} = 100$$

$$\text{Hit Ratio} = \frac{300}{400} = 0.75$$

[H]

$$\text{Miss Ratio} = \frac{100}{400} = 0.25$$

OR

$$(1-H) = 1 - 0.75 = 0.25$$

$$T_{avg} = H \times \text{Time taken by memory for each hit} + (1-H) \times \text{Time taken by memory for each miss}$$

$$\begin{aligned} &= 0.75 \times 20 + 0.25 \times 150 \\ &= 15 + 37.5 = 52.5 \text{ Avg} \end{aligned}$$

## Type of Memory Org

1. Simultaneous Access Memory Org.
2. Hierarchical Access Memory Org.

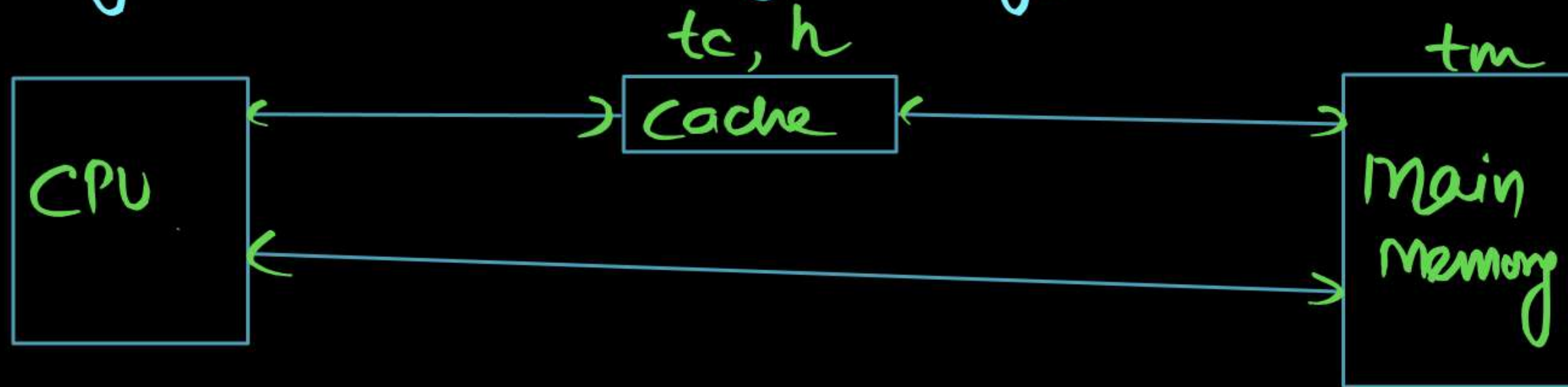


# Type of Memory Org



## 1. Simultaneous Access Memory Org.

(Both Memory Access Simultaneously/Parallely)



$t_c$ : Cache Access time

$h$ : Hit Ratio

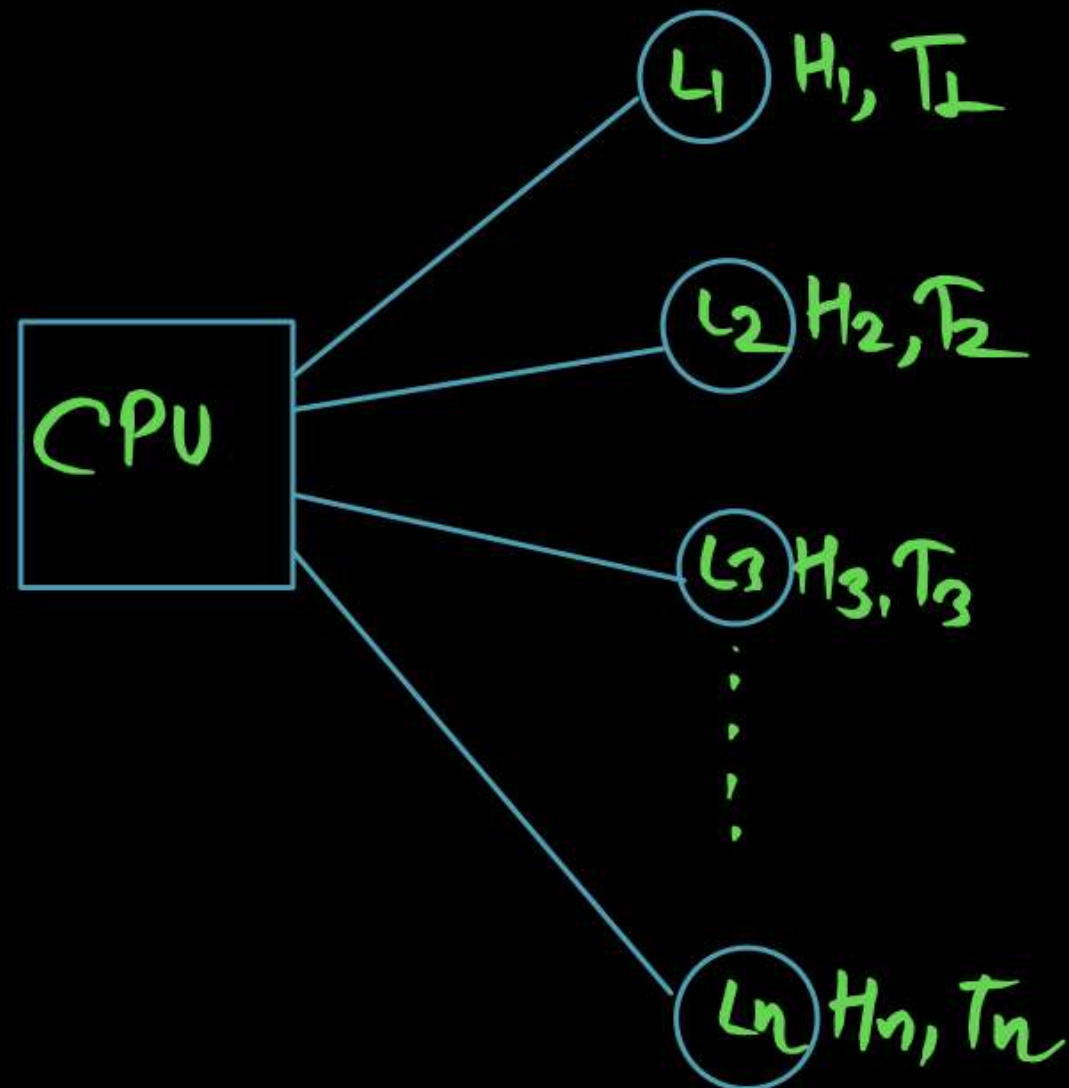
$t_m$ : Main Memory Access time.

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

# Type of Memory Org



## 1. Simultaneous Access Memory Org.



In the Simultaneous Access CPU is communicating All Level of Memory Directly. [All the levels of memory Directly connected to the CPU] But follow the sequence.

- When there is a Miss in  $L_1$  & Hit in Level  $L_2$  Memory then Directly Data is transferred from Level 2 memory to CPU without copying into Level 1 Memory.

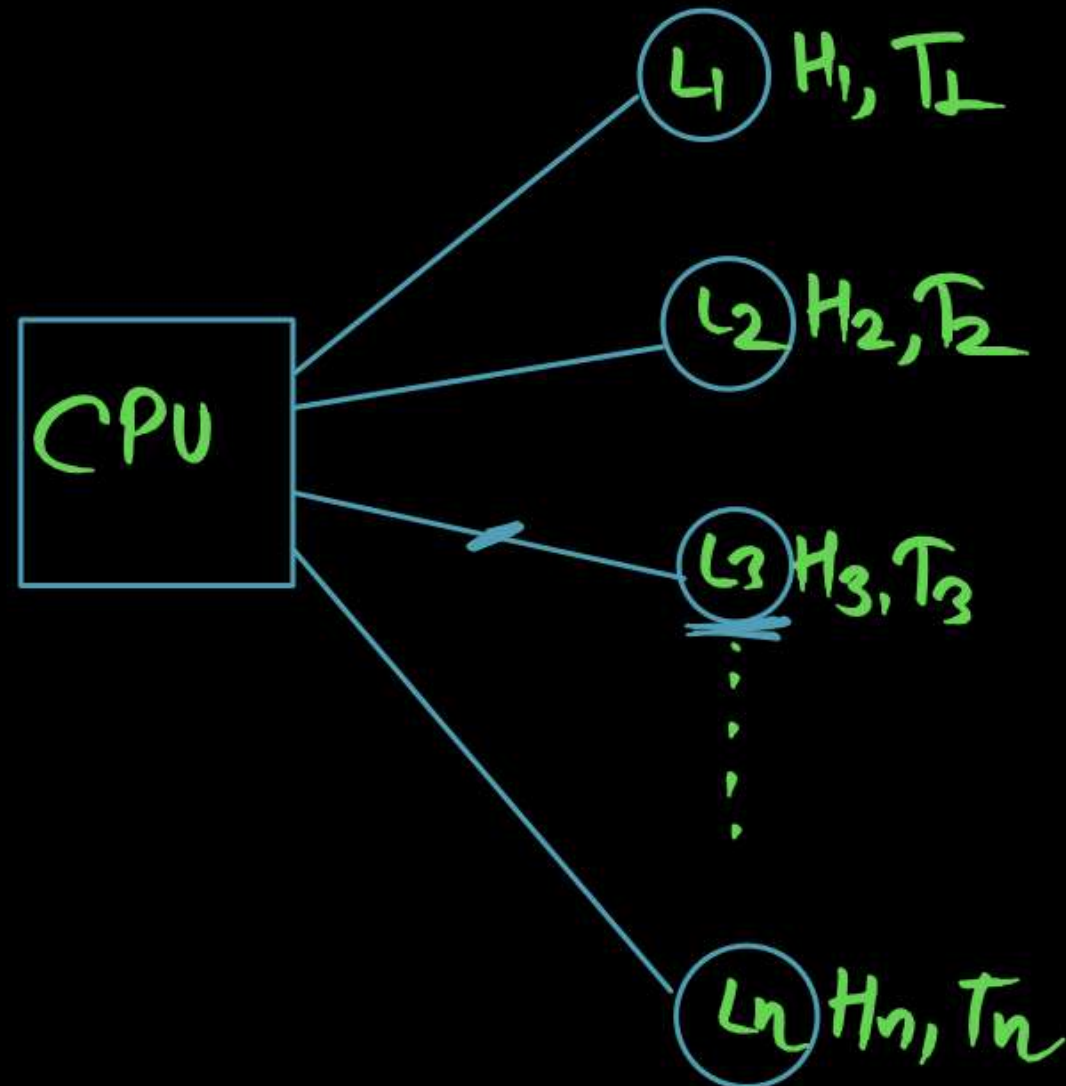


# Type of Memory Org



## 1. Simultaneous Access Memory Org.

- When there is a Miss in Level 1, Level 2 ( $L_2$ ) but Hit in Level 3 memory then Directly Data given from  $L_3$  memory to CPU Without Copying into Level 1 & Level 2 Memory.





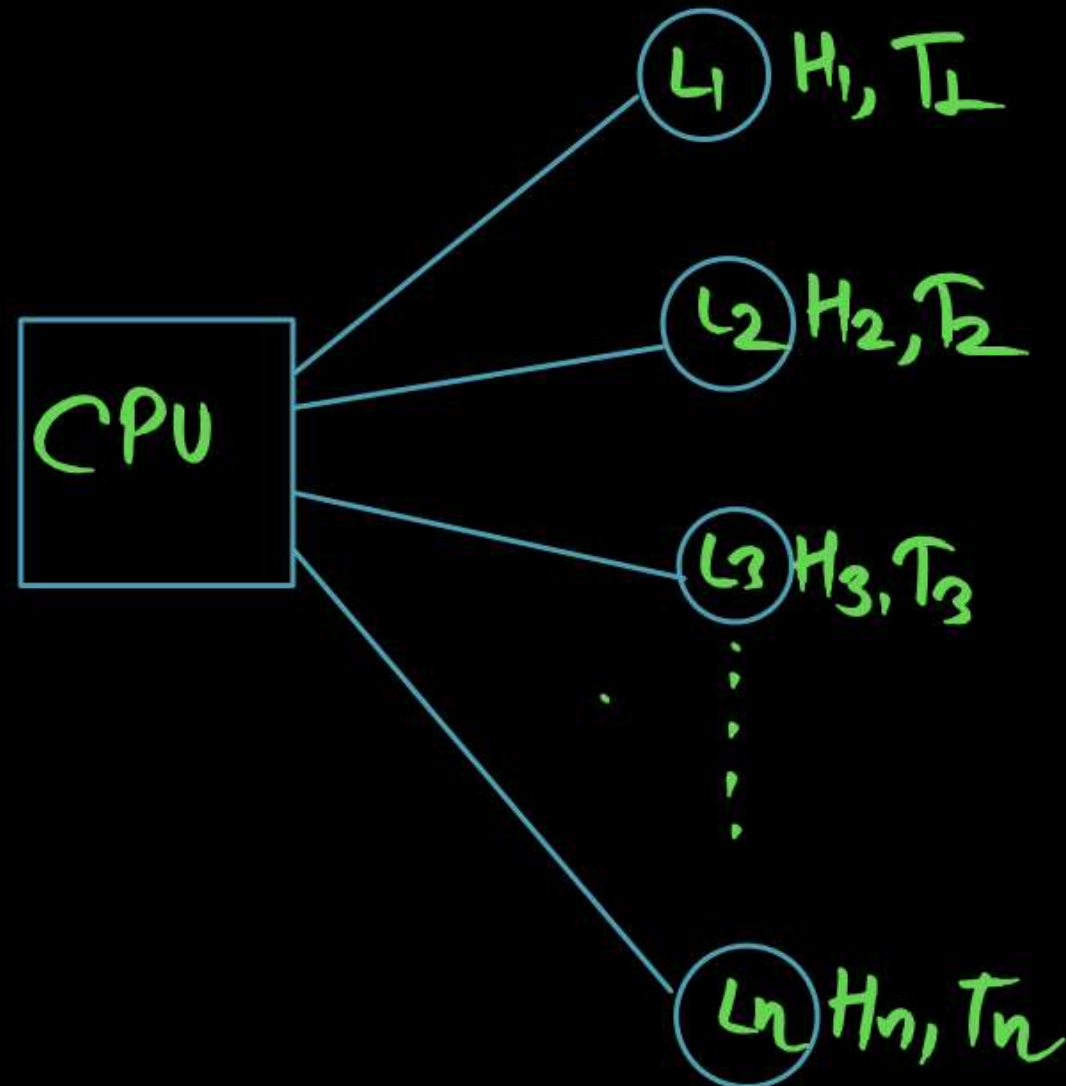
# Type of Memory Org



## 1. Simultaneous Access Memory Org.

Here  $H_1, H_2, H_3, \dots, H_n$  are Hit Ratio &

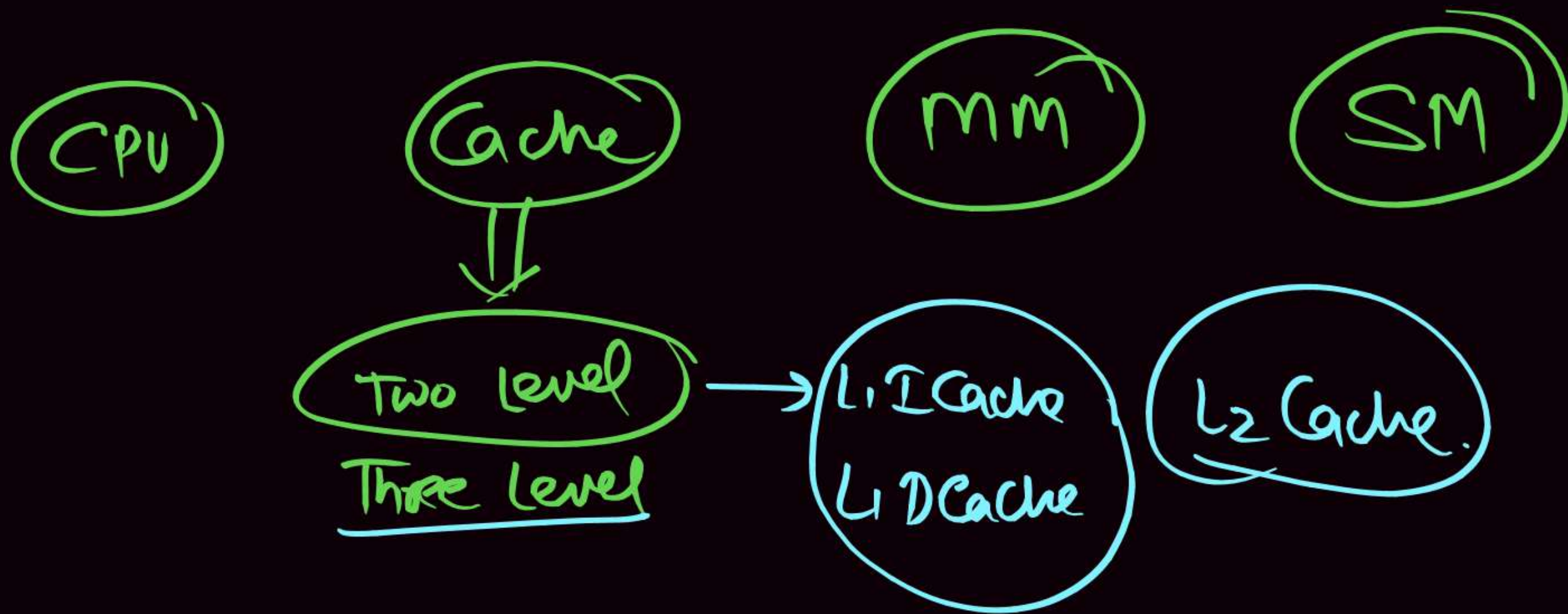
$T_1, T_2, T_3$  &  $T_n$  are Access time of Respective memory.



Time Required to Access (Read/Write) 1 word from the memory is called  $T_{avg}$ .

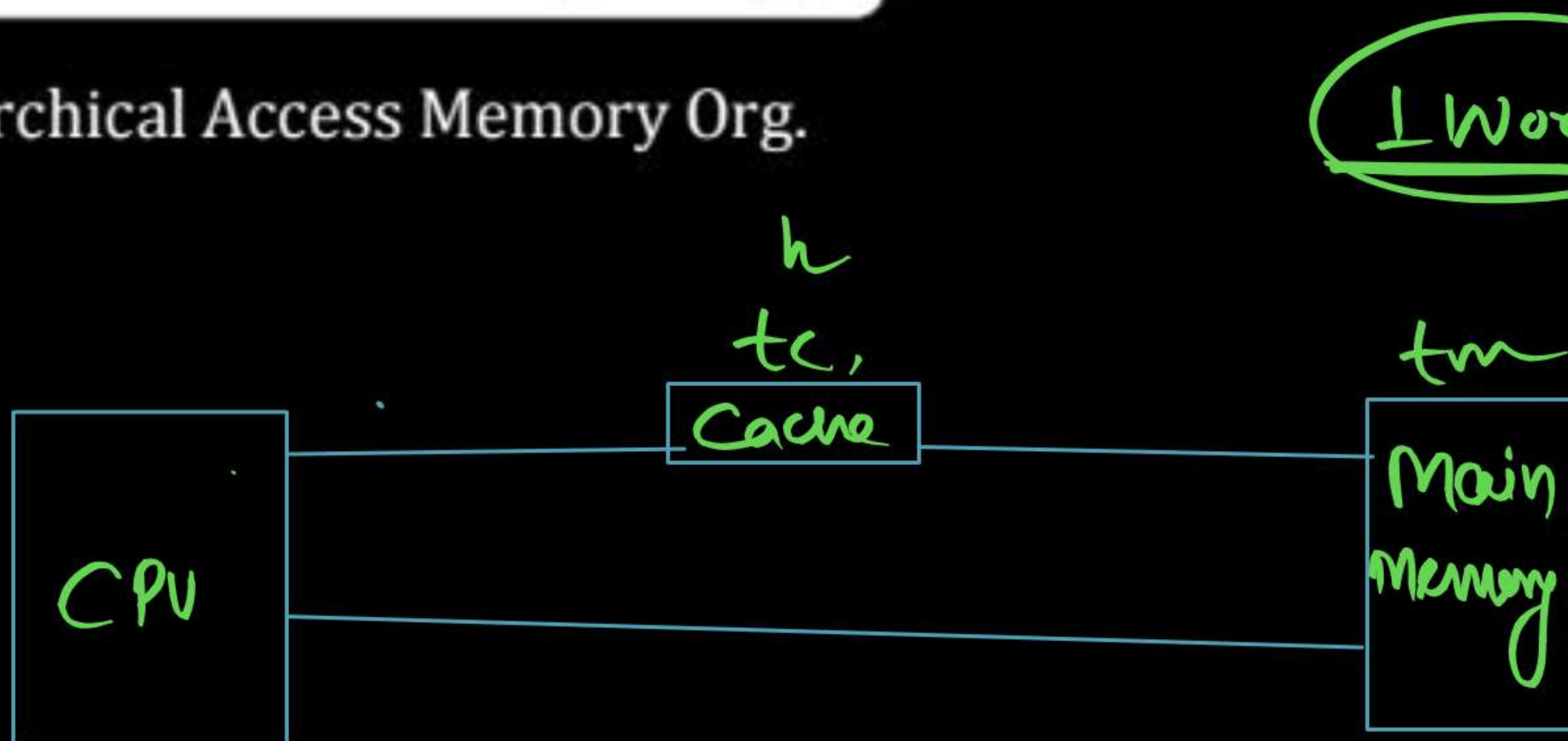
$$T_{avg} = H_1 T_1 + (1 - H_1) H_2 T_2 + (1 - H_1)(1 - H_2) H_3 T_3 + \dots + (1 - H_1)(1 - H_2)(1 - H_3) \dots (1 - H_{n-1}) H_n T_n$$

$H_n = 1$  Last Level Hit Ratio = 1



# Type of Memory Org

## 2. Hierarchical Access Memory Org.



1 Word

Block Size  
Assume ↓  
32 Byte /  
32 Words.



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

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

EMAT  
[ $T_{avg}$ ]  
1 Word Access  
time

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

in Hierarchical.

⑧ Hit Ratio = 80%  $t_c = 20ns$

$t_m = 100ns$

Tag Using Hierarchical Access?

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

$$\Rightarrow .80 \times 20 + (1-0.8)(100+20)$$

$$= .80 \times 20 + .20[120]$$

$$16 + 24 = 40ns$$

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

$$= 20 + (1-0.8) \times 100$$

$$\Rightarrow 20 + (0.2 \times 100)$$

$$= 20 + 20 = 40ns$$

Hit ( $h$ ) = 0.80  
(m) Miss  $(1-h) = 0.20$

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

$$0.8 \times 20 + .20 \times 20 + .20 \times 100$$

$$\frac{100\%}{1} \times 20 + \frac{(1-h)}{tm} \times 100$$

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

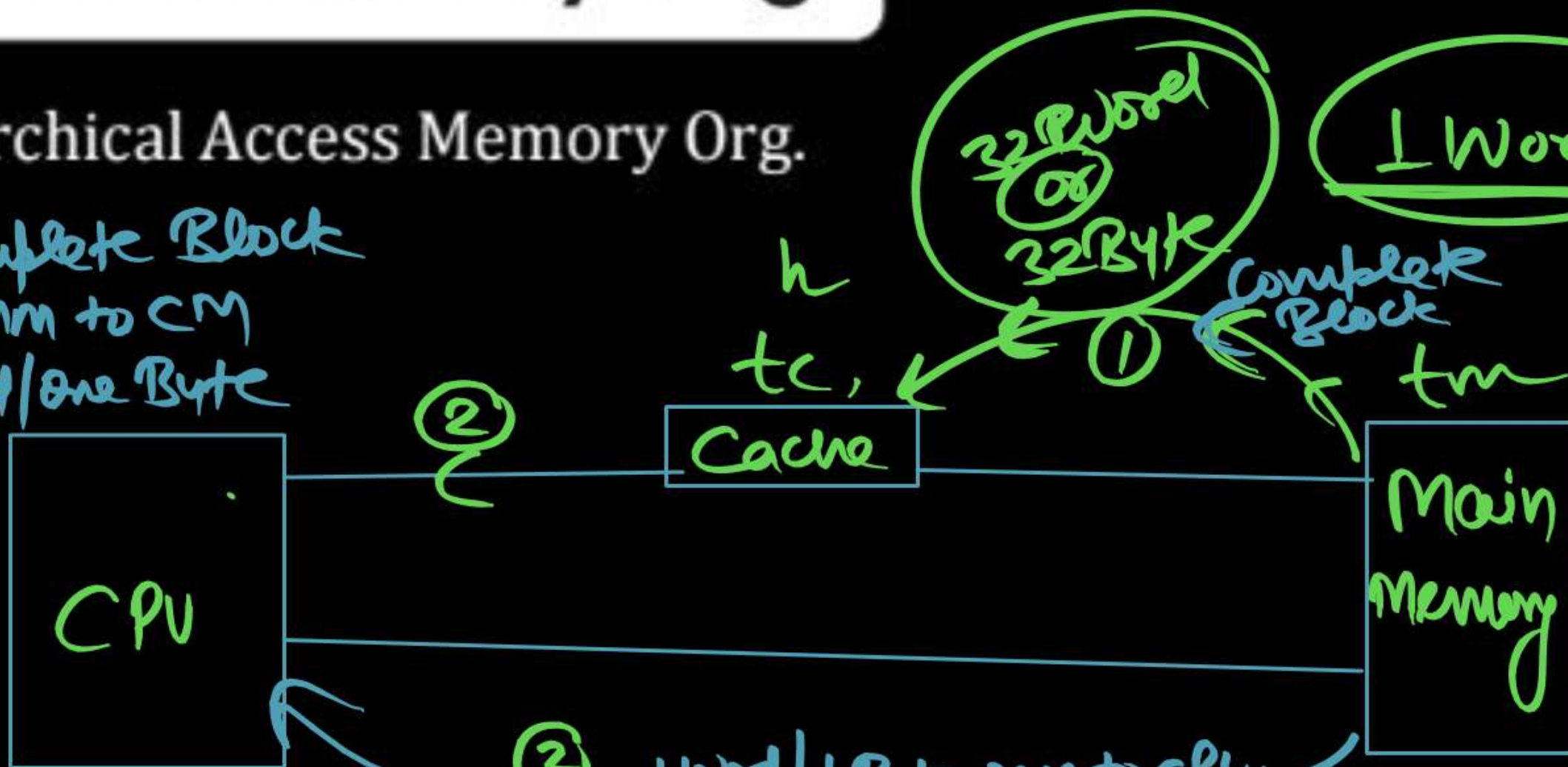
$$T_{avg} = t_c + (1-h)t_m$$



# Type of Memory Org

## 2. Hierarchical Access Memory Org.

- ① One Complete Block from MM to CM
- ② One Word/one Byte from MM to CPU



② 1 Word/1 Byte MM to CPU.  
In Next Reference Due to L.D.R  
there is a cache Hit.

Block  
Size  
Assume ↓  
32 Byte/  
32 Words.



② Block Size is 32 Byte / 32 Words But CPU Require only 1 Word then ?

Sw<sup>n</sup> In this Process Complete Block (32 Byte / 32 Words) Copied / Transferred from MM to Cache & that Respective 1 Word given to CPU [whichever CPU Demanded] [Locality of Reference]

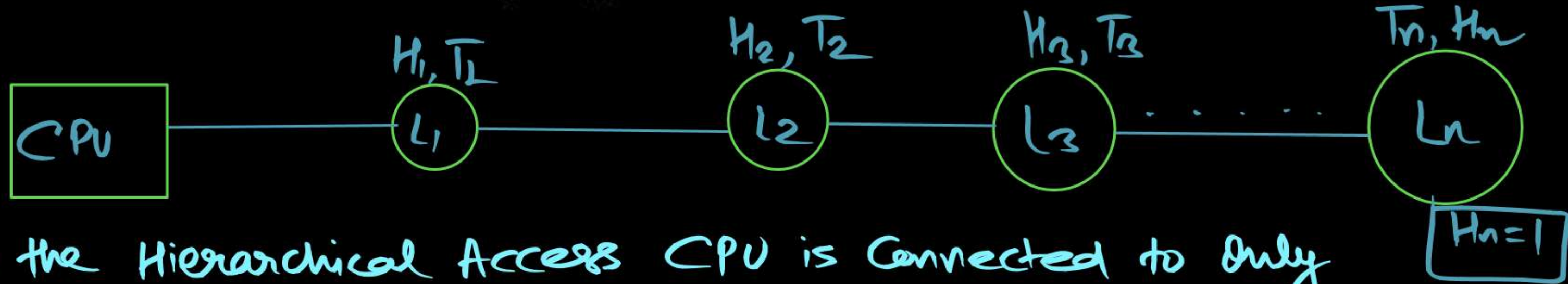
Advantage So in the Next time when CPU Request either same word or Adjacent word then that Request we find in Cache [Cache Hit]



# Type of Memory Org



## 2. Hierarchical Access Memory Org.



In the Hierarchical Access CPU is Connected to only Level 1 Memory.

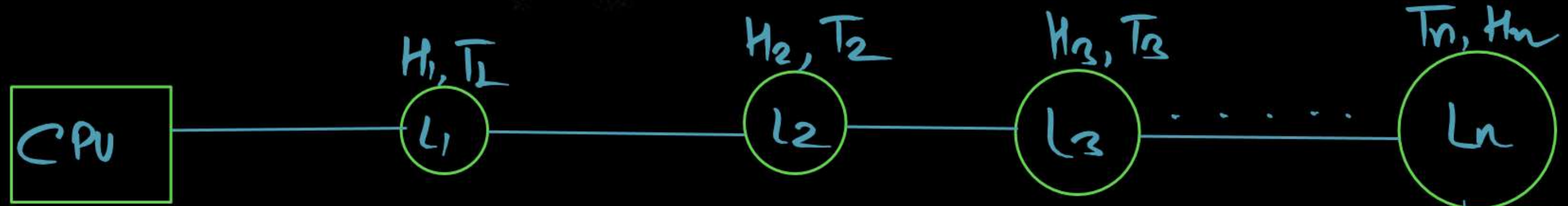
When there is a Miss in Level 1 Memory, & operation hit in Level 2 memory then firstly Data Moves (copied) from L<sub>2</sub> Memory to L<sub>1</sub> Memory & then Level 1 memory to CPU.



# Type of Memory Org



## 2. Hierarchical Access Memory Org.

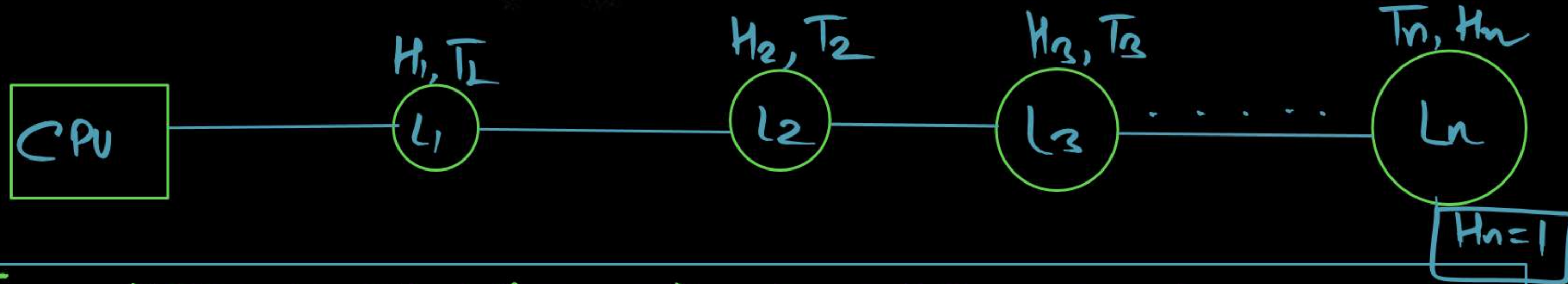


If there is a Miss in Level L<sub>1</sub> Memory & Level L<sub>2</sub> Memory But H<sub>n</sub>=1 Hit in Level 3 Memory then first Data Copied from L<sub>3</sub> Memory to L<sub>2</sub> Memory & Level 2 Memory to L<sub>1</sub> Memory then L<sub>1</sub> Memory to CPU.

# Type of Memory Org



## 2. Hierarchical Access Memory Org.



$$\begin{aligned} T_{avg} = & H_1 T_1 + (1 - H_1) H_2 (T_2 + T_1) + (1 - H_1) (1 - H_2) H_3 (T_3 + T_2 + T_1) \\ & \dots + (1 - H_1) (1 - H_2) (1 - H_3) \dots (1 - H_{n-1}) H_n (T_n + T_{n-1} + \dots + T_3 + T_2 + T_1) \end{aligned}$$



# Cache Work on Locality of Reference

- ① Temporal LOR
- ② Spatial LOR.

Note

If in a Question Mention the keyword 'Hierarchical Access' or 'Levels' of memory Accessing Hierarchical then Using Hierarchical Access.



Q.

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



Sol<sup>n</sup>

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

$$\Rightarrow .90 \times 1 + (1-0.9) (100+1)$$

$$= 0.9 + 0.1(101)$$

$$= 0.9 + 10.1$$

$$= \underline{11.0 \text{ ns}} \quad \underline{\text{Avg}}$$

Q.

PW

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?

$$T_1 = 20$$

$$H_2 = 1$$

Assume Access time of  
 $L_1 \Rightarrow T_1$

$L_2 \Rightarrow T_2$

$$5 = \frac{P_{L_1}}{P_{L_2}} \Rightarrow \frac{1/T_1}{1/T_2}$$

$$5 = \frac{T_2}{T_1}$$

$$T_1 = T_{avg} - 10$$

$$T_{avg} = T_1 + 10$$

$$T_2 = 5T_1$$

Here Given  $T_1 = 20$

$$T_2 = 5 \times T_1 = 5 \times 20 \Rightarrow T_2 = 100 \text{ nsec}$$

$$T_{avg} = T_1 + 10 \Rightarrow 20 + 10$$
$$T_{avg} = 30 \text{ nsec}$$

$$T_{avg} = HT_1 + (1-H)T_2$$

$$30 = H \times 20 + (1-H)100$$

$$30 = 20H + 100 - 100H$$

$$70 = 80H$$

$$H = \frac{70}{80} = 0.875 \text{ Ans}$$

$$87.5\% \text{ Ans}$$



$$\text{Performance} \propto \frac{1}{\text{ET}}$$

Ram

10 Hours

SHYAM

5 Hours

Same Work

SHYAM Performance  
is Fast.

Q.



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

$$T_1 = 20 \text{ nsec}$$

$$T_2 = 150 \text{ nsec}$$

$$T_{avg} = 30 \text{ nsec}$$

Simultaneous Access

$$T_{avg} = h \times T_1 + (1-h) T_2$$

$$30 = h \times 20 + (1-h) 150$$

$$30 = 20h + 150 - 150h$$

$$130h = 120$$

$$\text{Hit Ratio} = 92.33\% \text{ Ans}$$

$$H = \frac{120}{130}$$

$$H = \frac{12}{13}$$

$$H = 0.9233 \text{ Ans}$$

(ii) If Hit Ratio = 100%. then  
What's time of  $L_1$  &  $L_2$  memory

$$\left. \begin{array}{l} (T_1) L_1 = 20 \text{ ns} \\ (T_2) L_2 = 150 \text{ ns} \end{array} \right\} \text{Remain Same.}$$

Note Hit Ratio does not Effect Level Access time  
it effect  $T_{avg}$

$$T_{avg} = h \times t_1 + (1-h) t_2$$
$$= 1 \times 20 + (1-1) 150$$

$$T_{avg} = 20 \text{ ns} \text{ Ans}$$



Q.

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



$$T_1 = 20 \text{ ns}$$
$$T_2 = 150 \text{ ns}$$

$$T_{avg} = 30 \text{ ns}$$

Increased  
by 10%

$$T_{avg \text{ new}} = 30 + 10\% \text{ of } 30 \Rightarrow 30 + 3$$

$$T_{avg \text{ new}} = 33 \text{ nsec.}$$

Simultaneous Access

$$T_{avg \text{ new}} = h_{\text{new}} * t_1 + (1 - h_{\text{new}}) t_2$$

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

$$33 = 20H + 150 - 150H$$

$$117 = 130H$$

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

$$H = 0.9$$

Hit Ratio = 90%

& Earlier was 92.33%

Hit Ratio  
Decreased

by 2.33%  
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 \_\_\_\_\_. [GATE - 2015]

$$\text{Hit Ratio} = 80\% = 0.8$$

$$\text{When Hit time taken} = 5\text{ns}$$

$$\text{Miss Ratio} = (1 - 0.8) = 0.2$$

$$\text{When Miss time taken} = 50\text{ns}$$

$$T_{\text{avg}} = 0.8 \times 5 + 0.2 \times 50$$
$$4 + 10$$

$$T_{\text{avg}} = 14\text{ns}$$
 Ans



# LOR [Locality of Reference]

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

(Faster)

(1) Temporal LOR

$$R \leftarrow \underline{M(2000)}$$

Add: [2000]

A Location Wallah word  
which is available in  
Block No B5.

(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.

(x+1) location wallah word

which is available in  
Block No B<sub>5</sub>.

↓

Cache Hit

Bread  
milk }  $\Rightarrow$  1 Portion



# Types of Cache



- 1) Unified Cache: Instruction & Data both are placed in Same Cache.
- 2) Split Cache: This Cache logically Divide into two parts
  - (i) Instruction Cache [I - cache]
  - (ii) Data Cache [D- cache]
- 3) Multilevel Cache:



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

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



**THANK  
YOU!**

