# COMPUTER SCIENCE

## Computer Organization and Architecture

### Cache Memory

**Lecture_01**

Vijay Agarwal sir

TOPICS TO BE COVERED
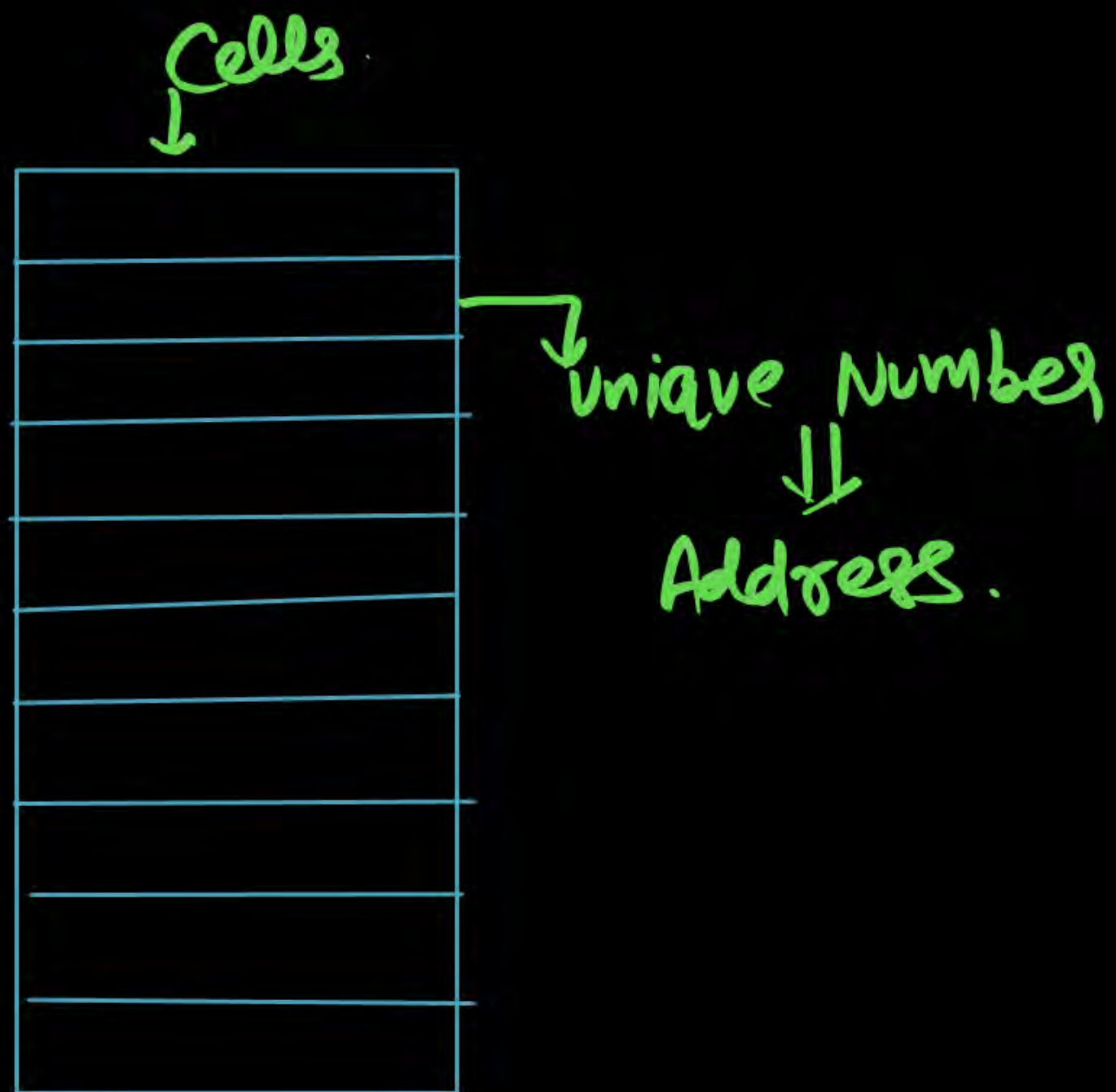
o1 Memory Hierarchy

o2 Cache Memory

① Introduction of COA

② Machine Instn & AM

③ Floating Point Representation

④ ALU Data Path & Control Unit

⑤ PIPE LINING.

⑥ CACHE Memory.

# Memory

Cells



Unique Number

$\Downarrow$

Address.

Q) Memory 256 KB then Address Size ?

Soln

256 KB.

$2.2^8 . 2^{10} \times 8\,bit$

$2^{18} \times 8\,bit$

Address = 18 bit.

$$2^1 = 2$$

$$2^2 = 4$$

$$2^3 = 8$$

$$2^4 = 16$$

$$2^5 = 32$$

$$2^6 = 64$$

$$2^7 = 128$$

$$2^8 = 256$$

$$2^9 = 512$$

$$2^{10} = 1024 \; [1k]$$

$$2^{10} = 1k \; (kilo)$$

$$2^{20} = 1M \; (mega)$$

$$2^{30} = 1G \; (Giga)$$

$$2^{40} = 1T \; (Tera)$$

$$2^{50} = 1P \; (Peta)$$

$$2^{60} = 1E \; (Exa)$$

$$1 \, msec = 10^{-3} \, sec$$

$$1 \, \mu sec = 10^{-6} \, sec$$

$$1 \, nsec = 10^{-9} \, sec.$$
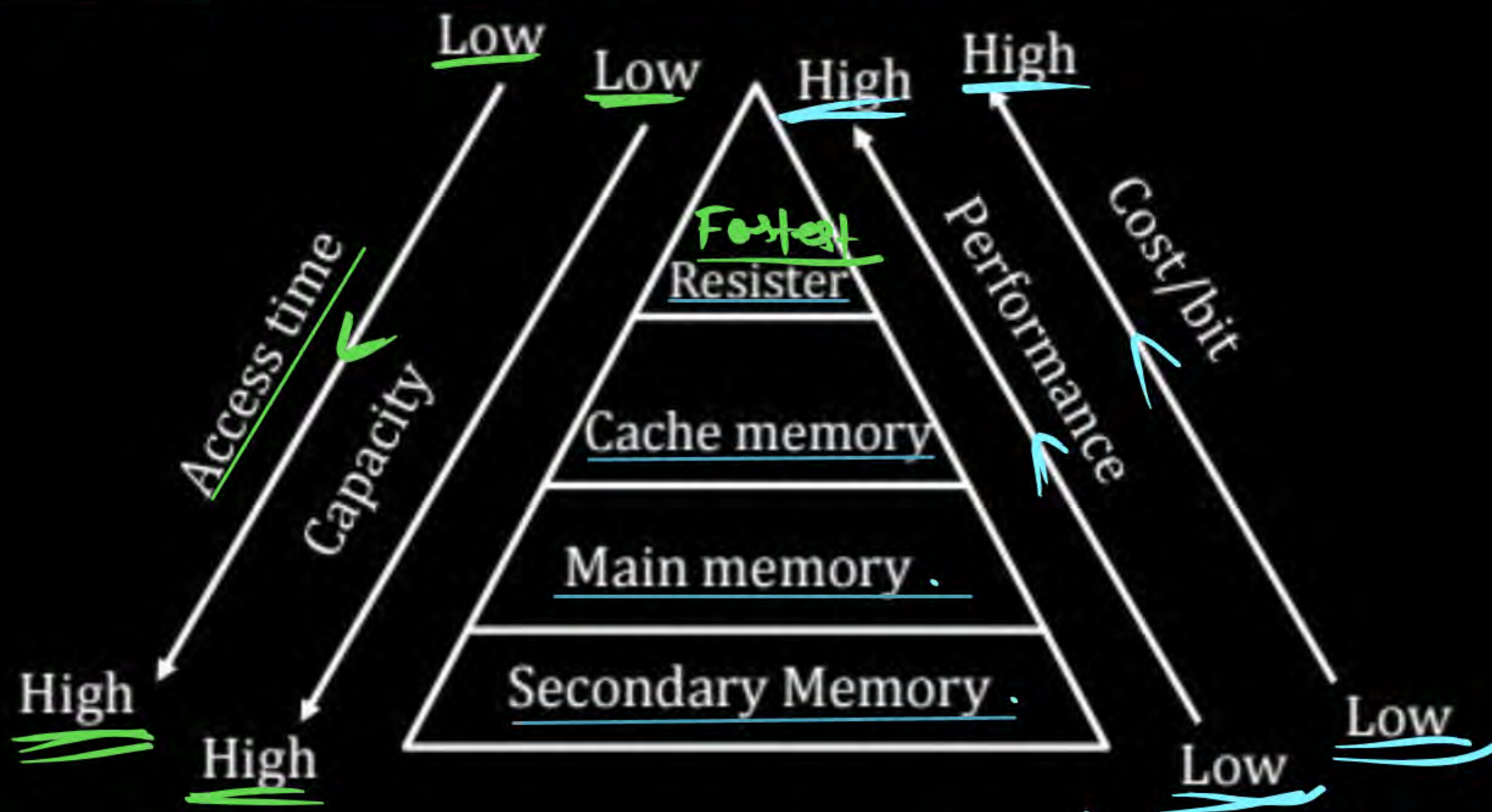
Word
Addressable
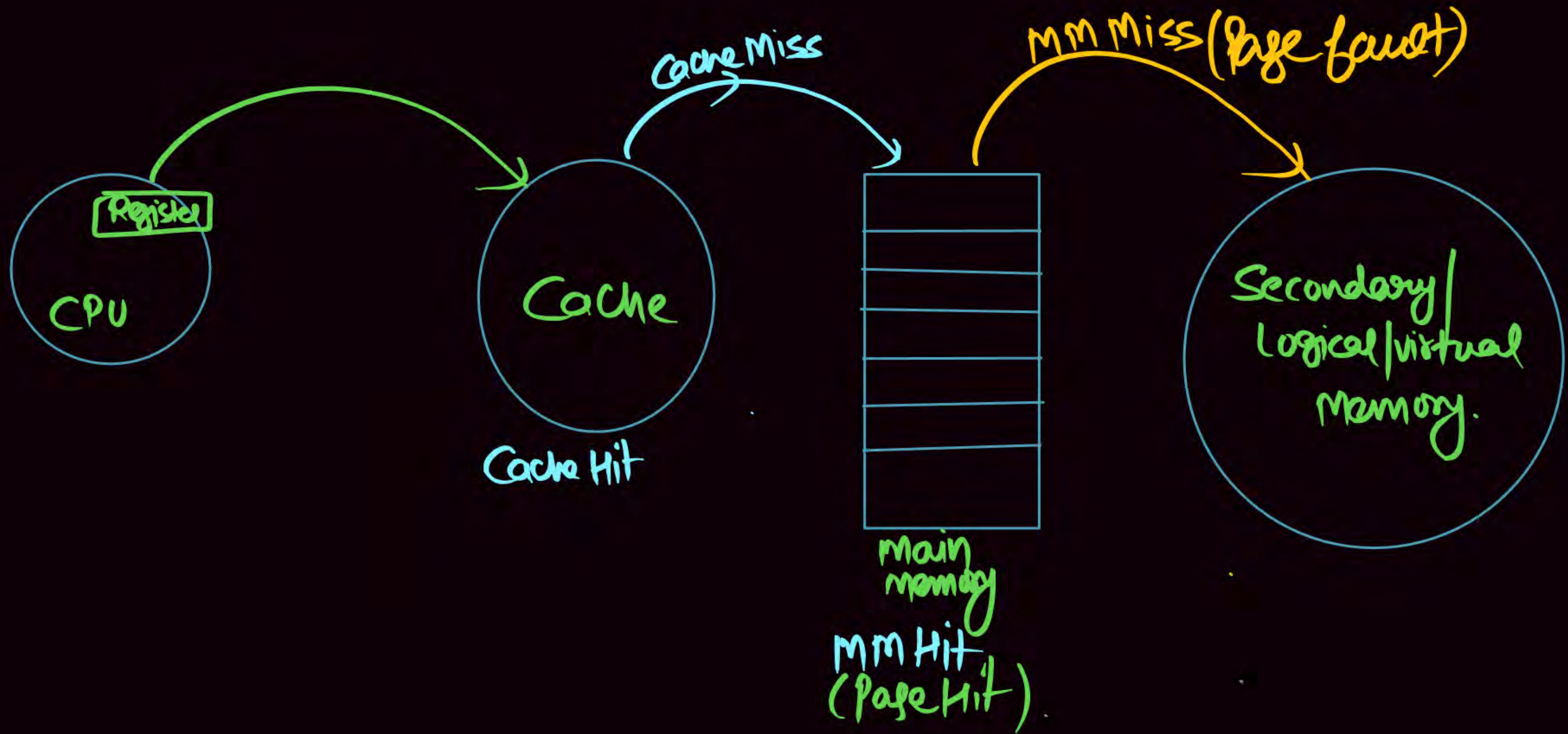
[eg 8 G Words]

Byte
Addressable

[eg 8 G Byte].

# Memory Hierarchy

❏ Hierarchy design organize the system supported memory into 4 levels to minimize the Accessing times. They are:

$$\text{Performance} \propto \frac{1}{\text{Execution Time}}$$

**CPU** — [Register]

**Cache** — Cache Hit

Cache Miss

MM Miss (page fault)

main memory — MM Hit (Page Hit)

Secondary/ Logical/virtual memory.

$$\boxed{\text{Hit Ratio} = \frac{\text{Number of Hit}}{\text{Total Number of Access}}}$$

(9) If Cache Hit Ratio 80%.

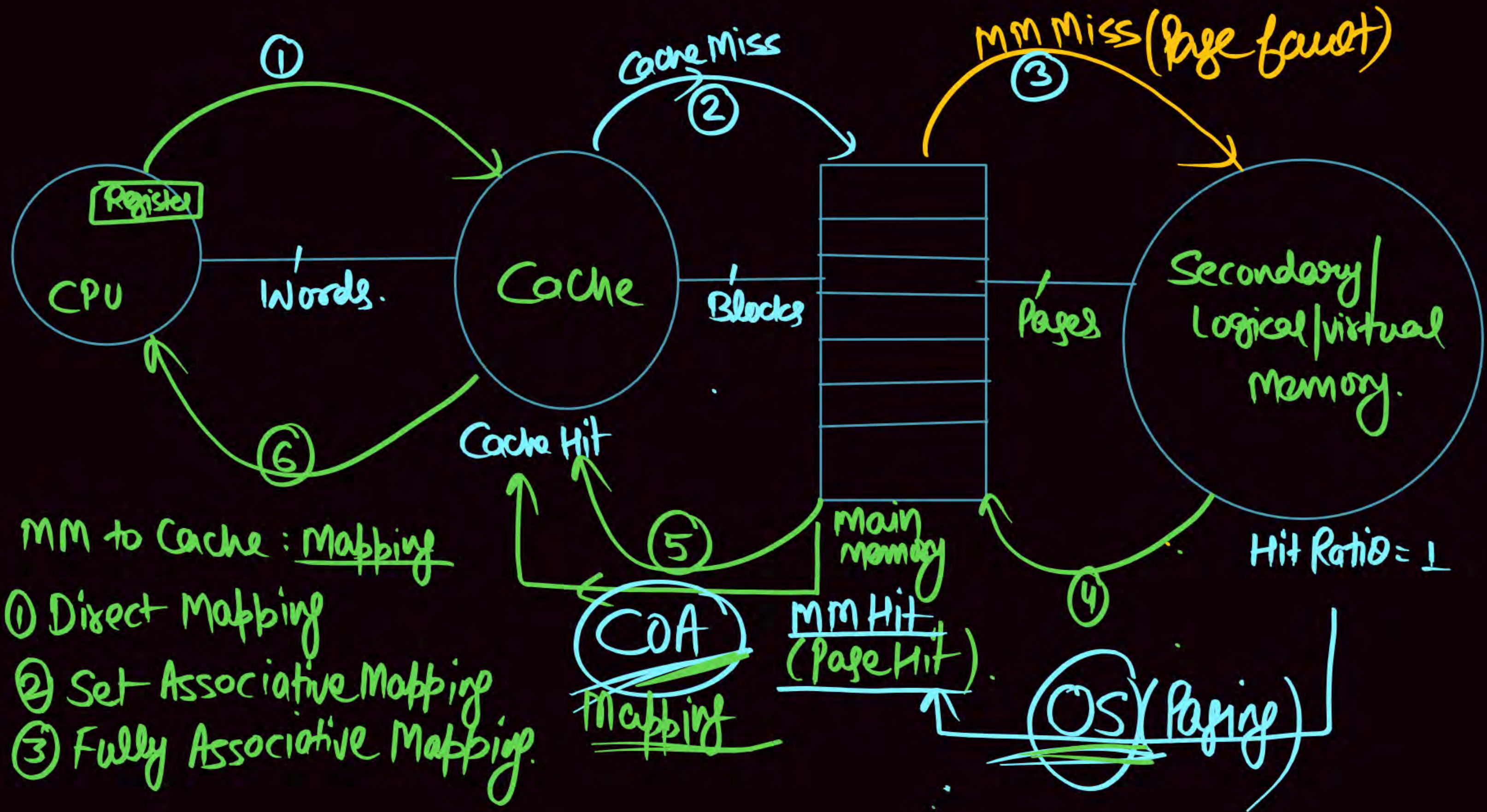    ⌐ ie out of 100, 80 Time there is Hit in Cache.

                          [Reference found in Cache].

CPU [Reg] — Words — Cache — Block. — MAIN memory — Pages. — SM|LM|VM

① (green arrow, CPU → Cache)

Register

CPU —— Words —— Cache

Cache Miss
②

MM Miss (page fault)
③

Secondary /
Logical / virtual
memory.

Cache —— Blocks —— [main memory table] —— Pages —— 

Cache Hit

⑤

COA
Mapping

main
memory

MM Hit
(Page Hit)

Hit Ratio = 1

④

OS (Paging)

⑥

MM to Cache : Mapping
① Direct Mapping
② Set Associative Mapping
③ Fully Associative Mapping.

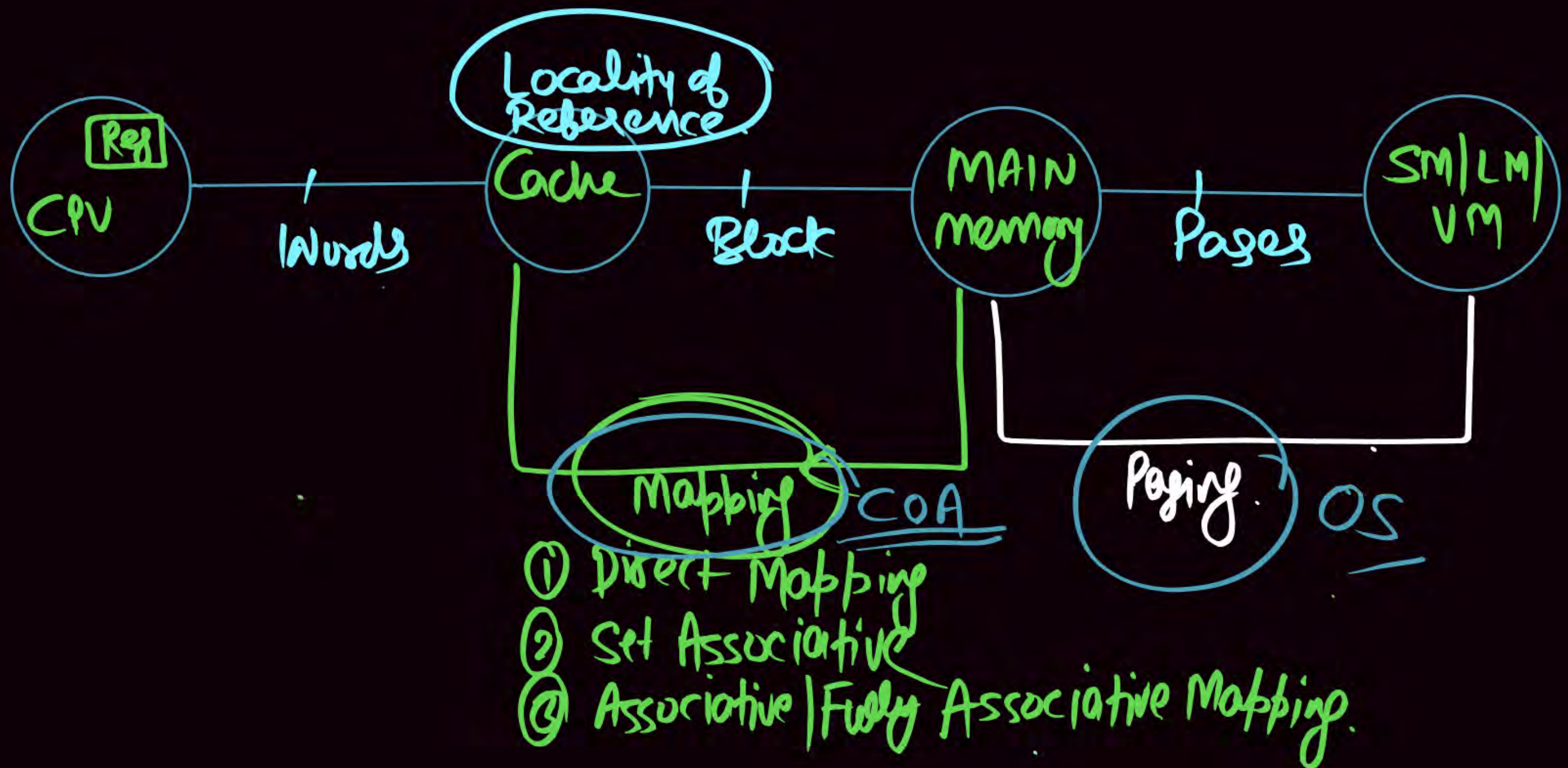# Memory

- CPU generate Request initally Refer to the Cache.

- If the Reference [Respective Data] found in the Cache that is Called Cache Hit [operation is Called Hit] then Respective Data give from Cache to CPU in the form of Words.

- If the Reference is Not found in the Cache, that is known as Cache Miss, then Reference forwarded to Main Memory.

# Memory

- If the Reference found in the Main Memory then its Called MM Hit (or) Page Hit. then Respective Data given from Main Memory to Cache in the form of blocks, & Cache to CPU in the form of words.

- If the Reference is Not found in the Main Memory that is Called MM Miss (or) Page fault then Reference forwarded to secondary memory.

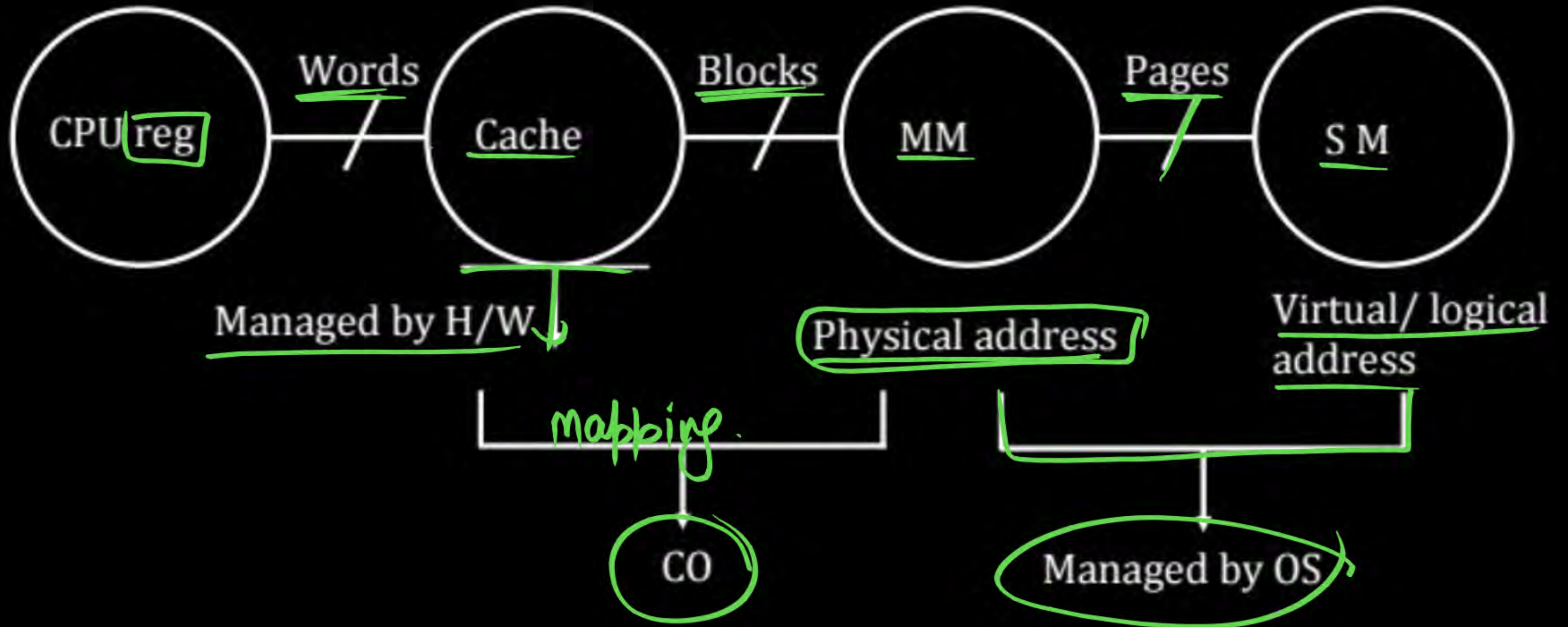(Note) Secondary Memory is the Last Level of Memory in which Hit Ratio always '1'.

# Memory

So Respective Data is transfered 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.

· The Process of transfer the Data from Main memory to Cache Memory is Called 'Mapping'.

# Memory



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

Managed by H/W

Physical address

Virtual/ logical address

mapping

CO

Managed by OS

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

$$T_{avg} = \text{Hit} \atop [H] \times \begin{bmatrix} \text{Time taken by} \\ \text{Memory when} \\ \text{there is a Hit} \end{bmatrix} + (1-H) \begin{bmatrix} \text{Time taken by} \\ \text{memory when} \\ \text{there is a Miss} \end{bmatrix}$$

Q.1 Consider CPU generate 100 Request. Out of 100 ⇒ 90 times Hit & 10 times Miss. When there is Hit then time taken is 20nsec, & when there is a Miss then time taken is 150nsec. Calculate the $T_{avg}$?

**Soln**

Total CPU Request = 100

Hit = 90 Times
Miss = 10 Times

Hit will takes = 20nsec
Miss will takes = 150nsec

Total Time = 90 × 20 + 10 × 150

= 1800 + 1500

100 Request

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

$$T_{avg} = \frac{3300}{100} = \boxed{33ns \; Ans}$$

---

Total CPU Request = 100.

Hit = 90 times   Miss = 10 Times

$$\frac{10}{100} = 0.1$$

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

**or**

Miss Ratio = $\frac{1 - 0.9}{} \; (1-H)$
= 0.1

$T_{avg}$ = H × Time Taken when there is a Hit + (1-H) [Time taken when there is a Miss]

⇒ 0.9 × 20 + 0.1 [150]

⇒ 18 + 15

$$\boxed{T_{avg} = 33nsec} \quad \underline{Ans}$$

**Q.2** Consider CPU generate 400 Request. out of 300 times Hit 100 times Miss. when there is Hit then time taken is 20nsec, & when there is a Miss then time taken is 150nsec. Calculate the $T_{avg}$?

**Soln**

Total CPU Request = 400

$$Hit = 300$$
$$Miss = 100$$

Hit will takes = 20nsec
Miss will takes = 150nsec

Total Time = $300 \times 20 + 100 \times 150$
$$= 6000 + 15000$$

But 400 Req:

$$\boxed{Total\ Time = 21000\ ns}$$

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

---

Total CPU Request = 400

$$Hit = 300 \qquad Miss = 100$$

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

$$\frac{100}{400} = 0.25$$

**or**

$$\frac{Miss}{Ratio} = (1-H) = (1-0.75)$$
$$= 0.25$$

$T_{avg}$ = H × Time Taken when there is a Hit $+ (1-H)$ [Time taken when there is a Miss]

$$\Rightarrow .75 \times 20 + .25 \times 150$$

$$\Rightarrow 15 + 37.5$$

$$\boxed{T_{avg} = 52.5\ ns} \quad \underline{Ans}$$

CPU ⟷ Cache ⟷ MM ⟷ SM

Cache
- L₁ Cache
  - I-Cache [Instruction Cache]
  - D-Cache [Data Cache]
- L₂ Cache
- L₃ Cache

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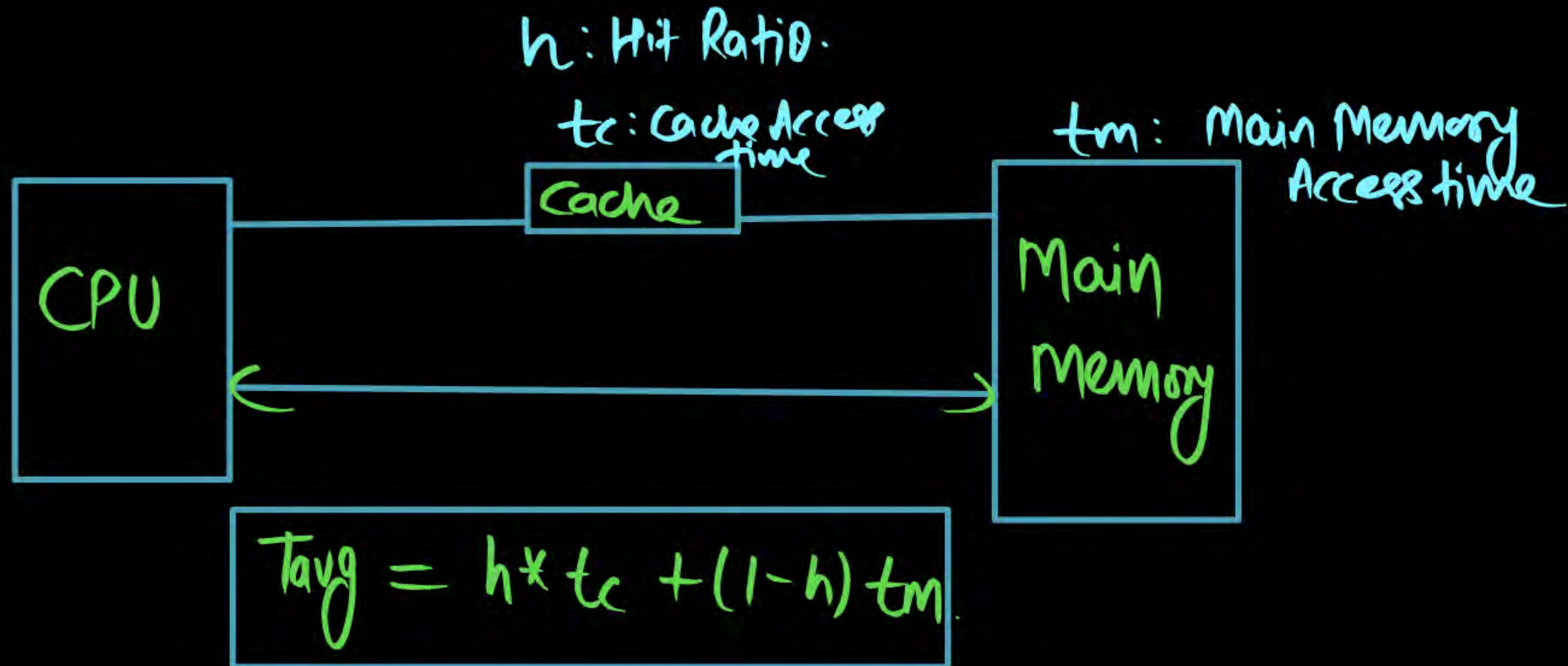
1. Simultaneous Access Memory Org.
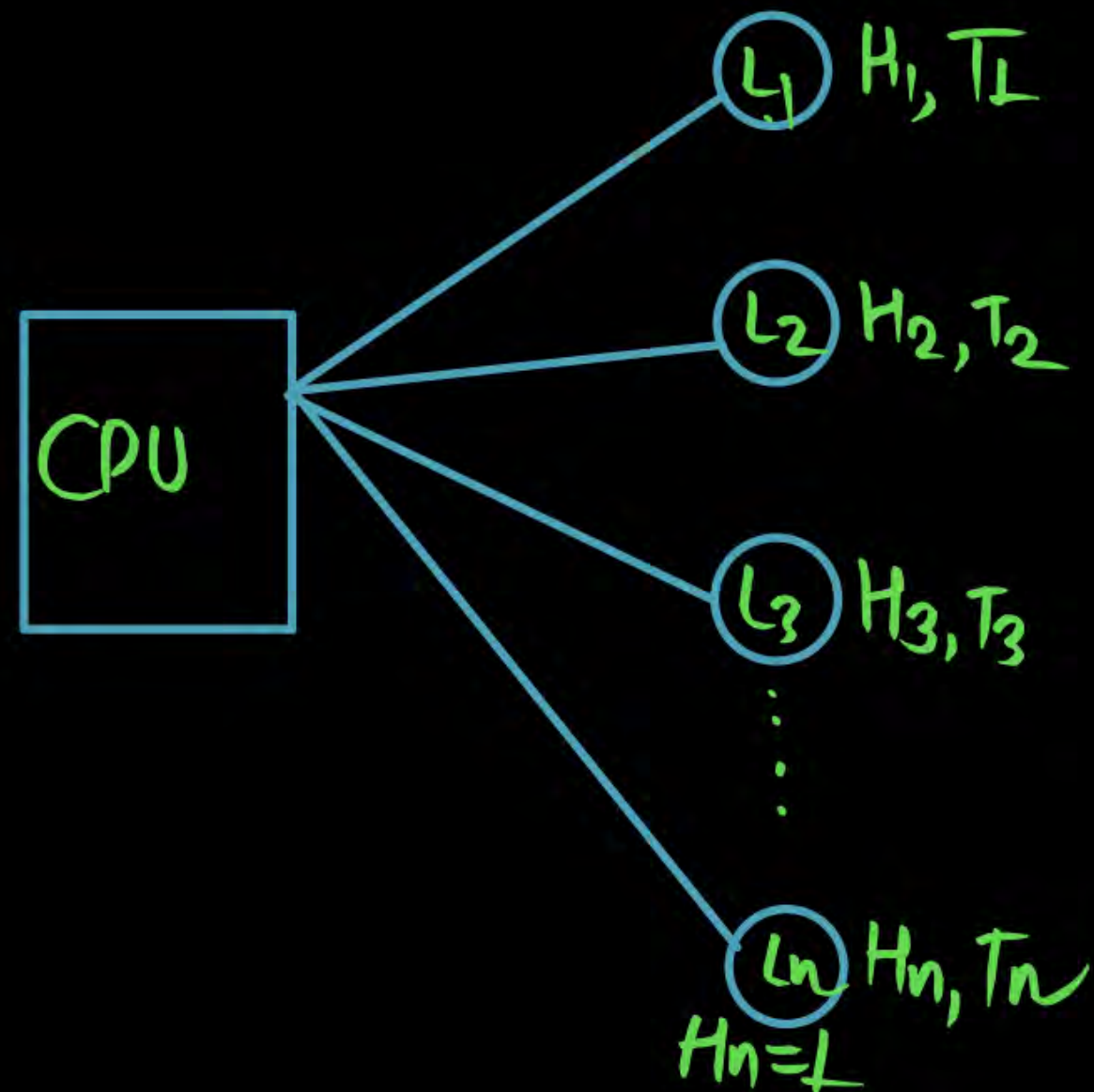
2. Hierarchical Access Memory Org.

# Type of Memory Org

**1. Simultaneous Access Memory Org. :** (Both Memory Access Simultaneously / Parallely)

$h$ : Hit Ratio.

$t_c$ : Cache Access time

$t_m$ : Main Memory Access time

CPU

Cache

Main Memory

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

# Type of Memory Org

**P W**

## 1. Simultaneous Access Memory Org.

```
        (L₁) H₁, T₁

        (L₂) H₂, T₂

CPU

        (L₃) H₃, T₃
          ⋮
        (Lₙ) Hₙ, Tₙ
      Hₙ = 1
```
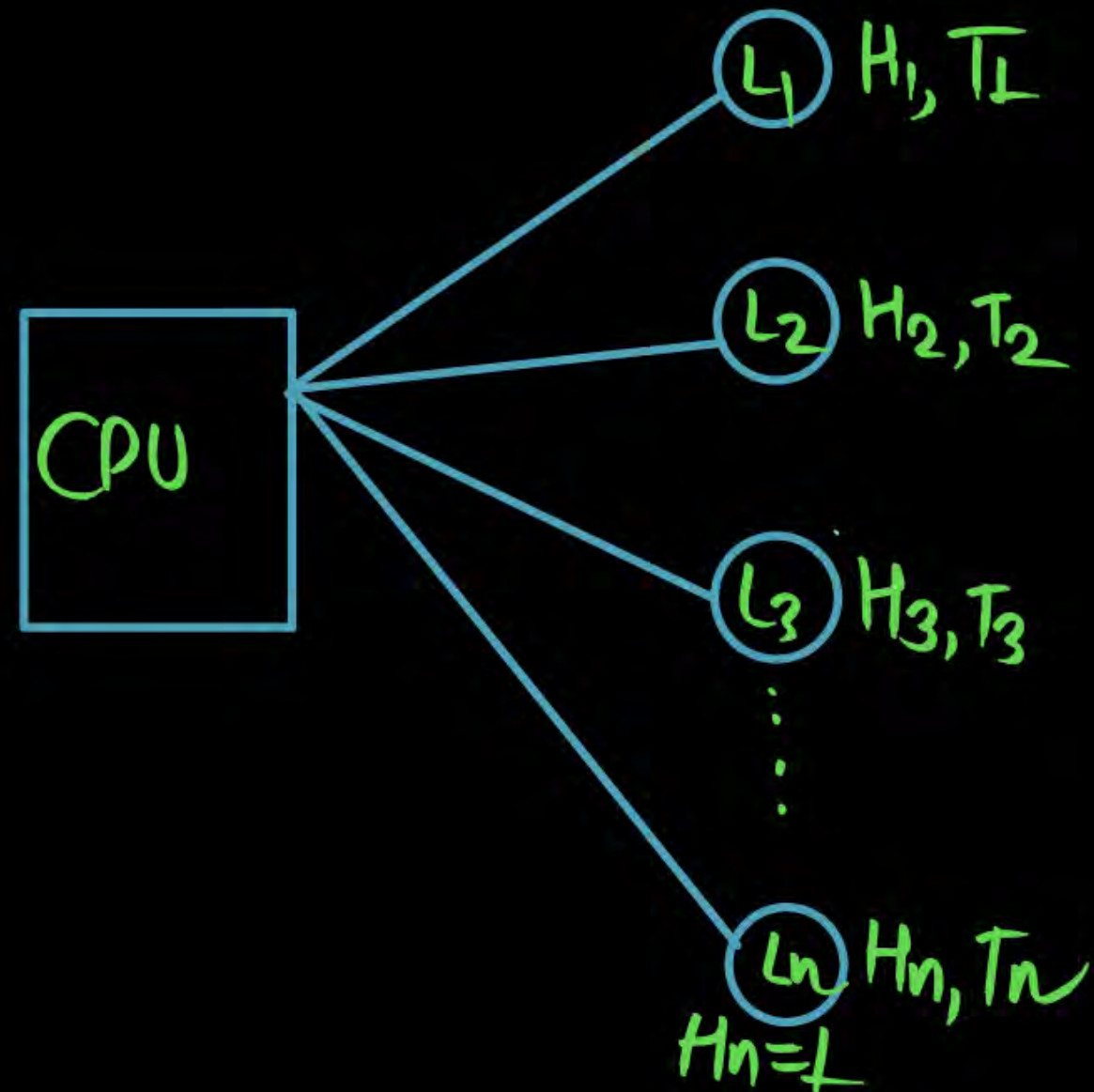
In the Simultaneous Access org All the level of Memory is Directly connected to CPU. But follow in Sequence (Acces in a Sequence)

· 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 transfered from Level 2 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 Memory & Hit in Level 3 Memory then Directly Data is transfered from Level 3 Memory to CPU Without Copying into Level 1 & Level 2 Memory.
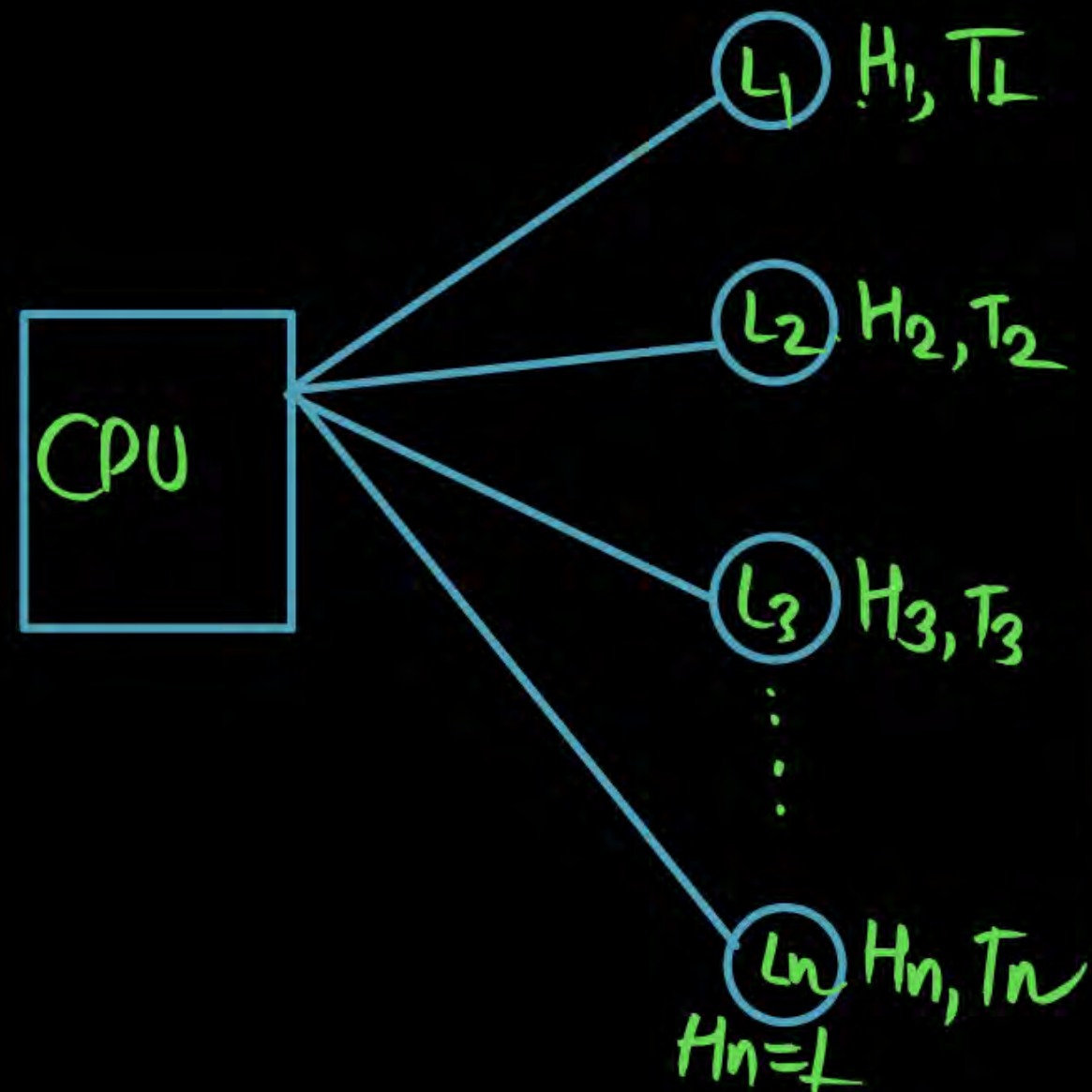
CPU

$L_1$ ) $H_1, T_1$

$L_2$ ) $H_2, T_2$

$L_3$ ) $H_3, T_3$

$L_n$ ) $H_n, T_n$

$H_n = 1$

# Type of Memory Org

$h_i$ : Hit Ratio of Level $i$

$(1-H_i)$ : Miss Ratio of Level $i$.

## 1. Simultaneous Access Memory Org.

Here $H_1, H_2, H_3 \ldots H_n$ are hit Ration of

$T_1, T_2, T_3 \ldots T_n$ are the Access of Respective level Memory.

The Time to Required to Access (Read/Write) 1 word from Memory



$(L_1) \, H_1, T_1$

$(L_2) \, H_2, T_2$

CPU

$(L_3) \, H_3, T_3$

$\vdots$

$(L_n) \, H_n, T_n$

$H_n = 1$

$$T_{avg} = h_1 t_1 + (1-h_1) h_2 t_2 + (1-h_1)(1-h_2) h_3 t_3 + \ldots$$
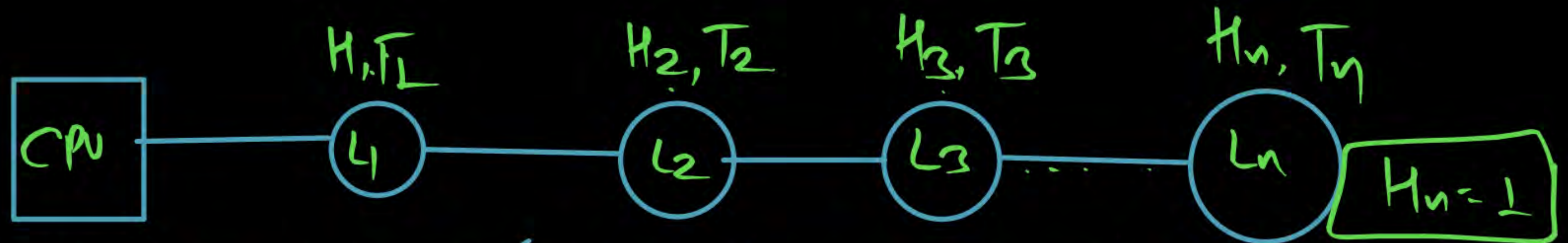$$\ldots (1-h_1)(1-h_2)(1-h_3) \ldots (1-h_{n-1}) H_n \, t_n \, .$$

$\boxed{H_n = 1}$ Last Level hit Ratio = 1
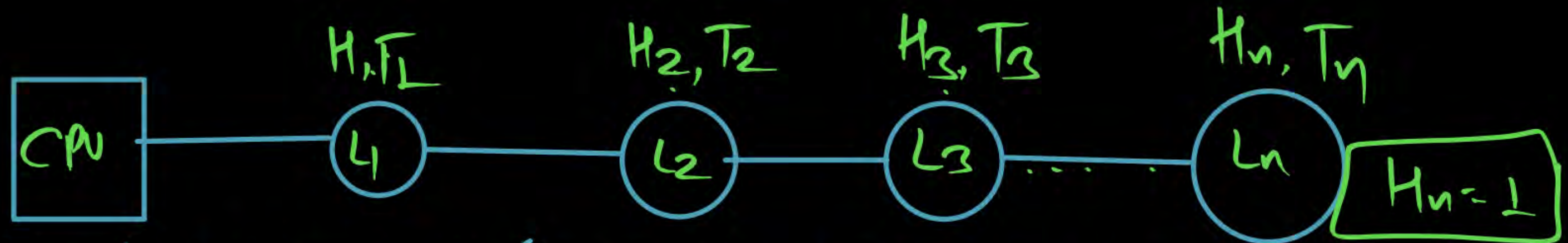
## 2. Hierarchical Access Memory Org.



· In the Hierarchical Access CPU is Communication with only Level 1 Memory.

· If there is a Miss in Level 1 Memory & Hit in Level 2 Memory then first Data is transfered from Level 2 Memory to Level 1 Memory then Level 1 Memory to CPU.

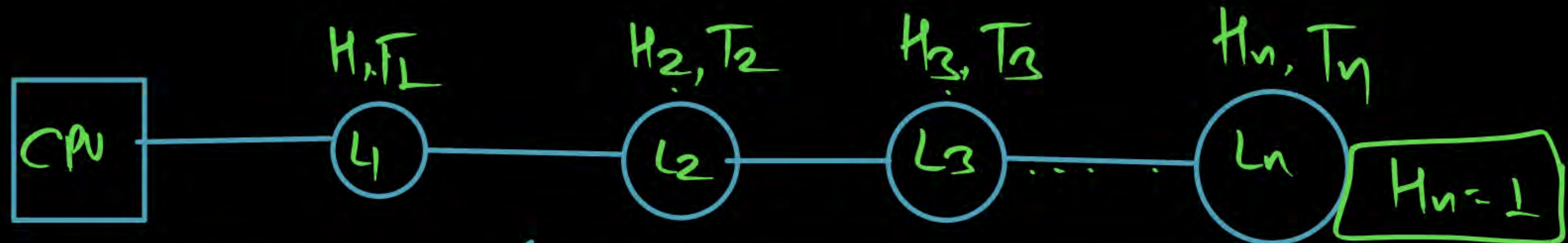# Type of Memory Org

## 2. Hierarchical Access Memory Org.

$$H_1, T_1 \qquad H_2, T_2 \qquad H_3, T_3 \qquad H_n, T_n$$

CPU — $L_1$ — $L_2$ — $L_3$ ..... $L_n$ — $H_{n-1}$

When there is a Miss in Level 1 & level 2 memory & Hit in level 3 memory then firstly Data is transfered from Level 3 $(L_3)$ memory to Level 2 $(L_2)$ memory then Level 2 $(L_2)$ memory to Level 1 $(L_1)$ Memory then from Level 1 $(L_1)$ Memory to CPU.

# Type of Memory Org

## 2. Hierarchical Access Memory Org.



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

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

$\boxed{H_n = 1}$ last level Hit Ratio $= 1$.

(Note) If in a Question Metioned the keyword 'Hierarchical Access' or Level of Access or Hierarchical Meaning 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?

**Soln.**

$$t_c = 1 \text{ nsec} \qquad h = 90\%$$
$$t_m = 100 \text{ nsec} \qquad h = 0.9$$

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

$$\Rightarrow 0.9 \times 1 + (1-0.9)(100+1)$$
$$= 0.9 + (0.1)(101)$$
$$= 0.9 + 10.1$$
$$= 11 \text{ nsec} \quad \text{Ans}$$

**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?

$$\boxed{\text{Perform}_{\text{ance}} \propto \frac{1}{ET}} \quad L_1 \Rightarrow T_L$$

$$L_2 \Rightarrow T_2$$

(Sol^n)

$$5 = \frac{PL_1}{PL_2} = \frac{1/T_1}{1/T_2}$$

$$5 = \frac{T_2}{T_1} \qquad T_1 = T_{avg} - 10$$

$$\boxed{T_2 = 5 * T_L} \qquad \boxed{T_{avg} = T_L + 10}$$

$$T_2 = 5 \times 20 \qquad T_1 = 20ns$$

$$\boxed{T_2 = 100 nsec} \qquad T_{avg} = T_1 + 10 \Rightarrow 20 + 10 \boxed{T_1 = 30 nsec}$$

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

$$30 = H * 20 + (1-H) 100$$

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

$$80H = 70$$

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

$$\Rightarrow 87.5\% \ \underline{Ans}$$

**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 mode to 100% then what is the Access time of $L_1$ & $L_2$ Memory?

**PART II**

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

**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 **[GATE – 2015]**

Hit Ratio = 80 = 0.8

Miss Ratio = (1- 0.8) = 0.2

When there is Hit
    Time taken = 5

When there is a
    Miss time taken = 50

$T_{avg}$ = 0.8 × 5 + 0.2 × 50

= 4 + 10

= 14 nsec Ans

Ans (14)