## Efficiency of memory hierarchy :-

Let efficiency ($e$) be the factor by which $t_1$ differs from average time $t_{Avg}$.

i.e.   $e = t_1 / t_{avg}$

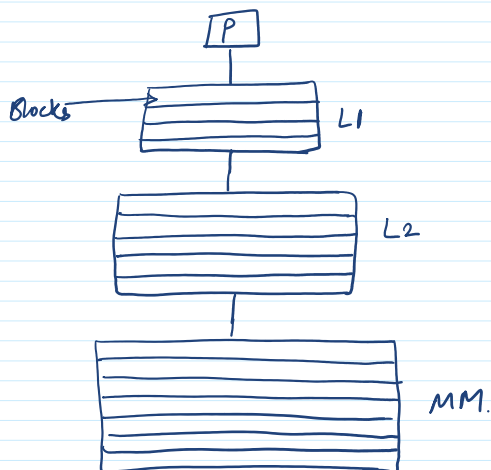Efficiency  $e = \dfrac{t_1}{H \cdot t_1 + (1-H) t_2}$  —①

Let $r = t_2 / t_1$
which is access time ratio of two levels.

putting in ①

$\Rightarrow e = \dfrac{1}{H + (1-H) r}$

## Speedup - gained by - Memory Hierarchy :-

$$S = \dfrac{t_2}{H \cdot t_1 + (1-H) t_2} = \dfrac{1}{H/r + (1-H)}$$

**Block** :- The smallest unit of information transferred between two levels.
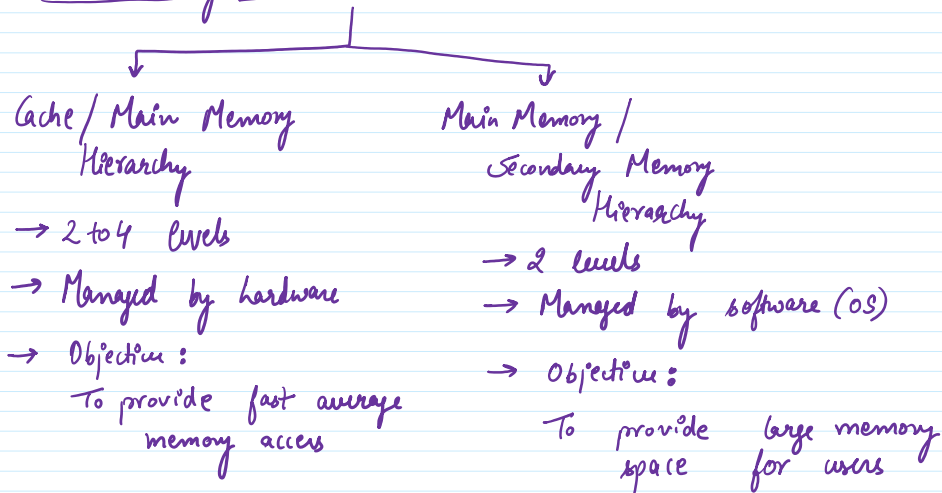


Block placement :- where is the block placed.
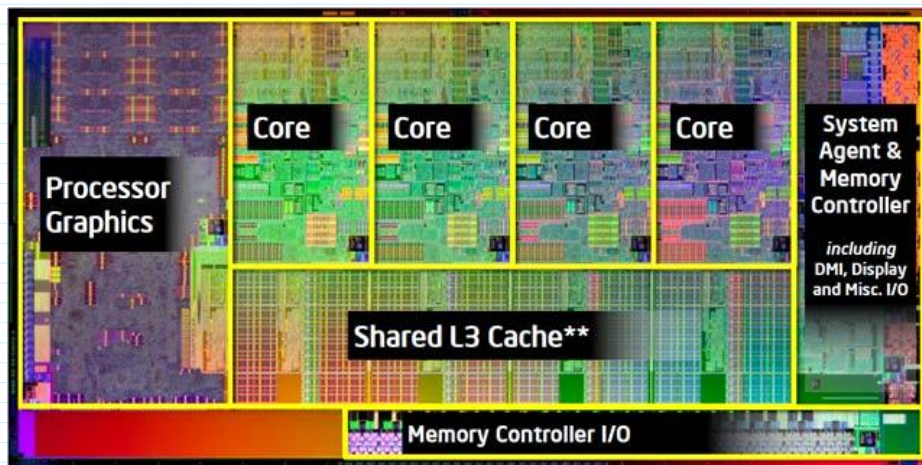
Block Identification :- Identify the block if
                        it is matched with
                        the upper level of the memory

Block replacement :- Which block is to be replaced
                     on a miss.


Common Memory Hierarchies :-

Cache / Main Memory                  Main Memory /
Hierarchy                            Secondary Memory
                                     Hierarchy

→ 2 to 4 levels                      → 2 levels

→ Managed by hardware                → Managed by software (OS)

→ Objection :                        → Objection :
  To provide fast average              To provide large memory
  memory access                        space for users


# Intel Core i7 - 4790K Processor :-



Core — A single processing unit within a CPU that
       can execute instructions.

4 cores — every core has its own L1 and L2 cache

          L3 is shared by all the cores.


# Cache Memory :-

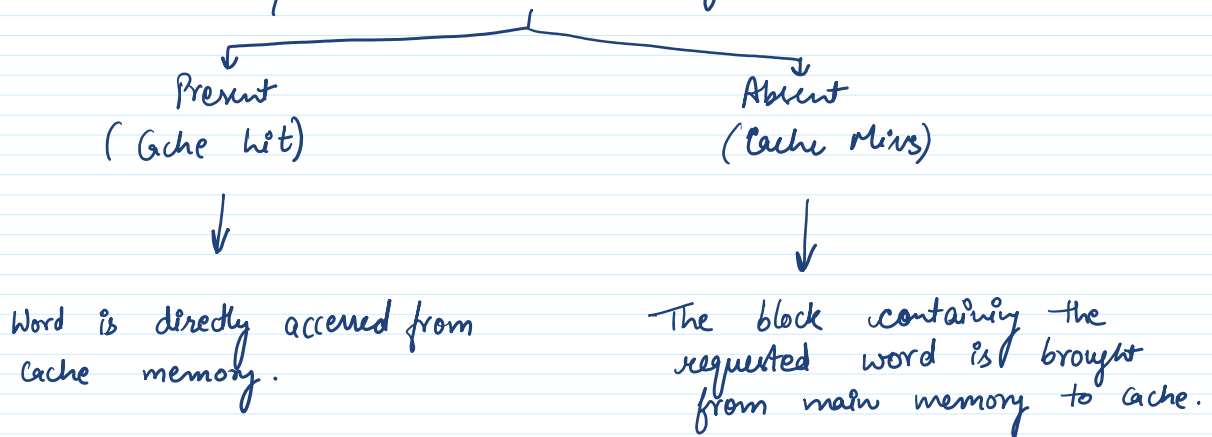→ A fast memory between processor and main memory.

Read/Write strategies, Block replacement, Mapping Techniques etc.

→ For fast execution, frequently used data and instructions are brought into the cache.

→ First time there is definitely a miss and then data is brought from lower levels into cache.

→ Cache Memory is logically divided into blocks/lines, where every block typically contains 8 to 256 bytes.

→ When the CPU wants to access a word in memory, a special hardware first checks whether it is present in cache memory.
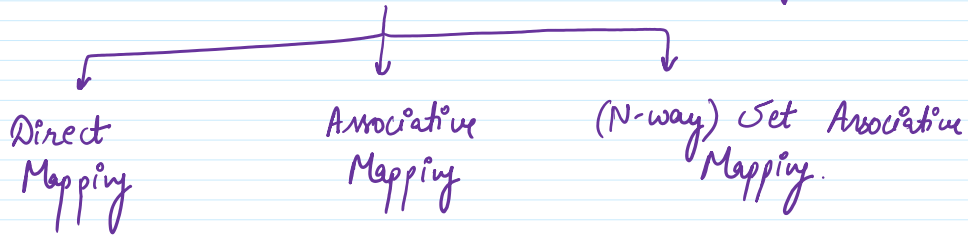
Present
(Cache Hit)

↓

Word is directly accessed from cache memory.

Absent
(Cache Miss)

↓

The block containing the requested word is brought from main memory to cache.

## Block Placement :—

Where can the block be placed in the cache?

→ determined by mapping algorithm.

Which main memory blocks can reside in which cache memory blocks.

only a small subset of main memory blocks can be held in cache memory.

Direct Mapping

Associative Mapping

(N-way) Set Associative Mapping.

A 2-Level Cache/Main Memory Hierarchy
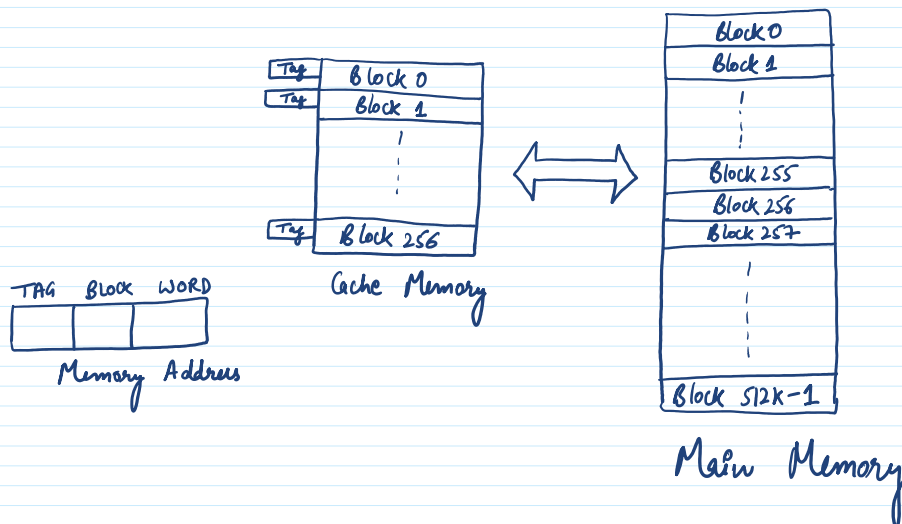
→ Cache Memory :—

256 blocks/lines of 32 words each
Total size = 8192 (8k) words.

→ Main Memory :—

Total size = 16 M words = $2^{24}$

$\Rightarrow$ 24-bit addressable

No. of 32-word blocks = $\dfrac{16M}{32}$ = 512K

## 1. Direct Mapping :—



| TAG | BLOCK | WORD |
|-----|-------|------|

Memory Address

Cache Memory

Main Memory

→ Each main memory block can be placed in only one block in the cache.

The mapping function is :—

$$\text{Cache Block} = \dfrac{(\text{Main Memory Block})}{(\text{No. of cache blocks})}$$

→     Direct Memory Address

| TAG | BLOCK | WORD |
|-----|-------|------|
| 11 | 8 | 5 |

32 Word Memory → $2^5$ ⇒ 5 bits for each word

256 blocks → $2^8$ → 8 bits to represent each block.

$$TAG = \frac{\text{\# of blocks in Main Memory}}{\text{\# of blocks in Cache Memory}}$$

TAG → which block of main memory is mapped to a particular block of cache memory.

$$TAG = \frac{512k}{256} = \frac{2^{19}}{2^8} = 2^{11}$$