

## chapter - 07

### MEMORY ORGANIZATION

Memory refers to any medium for data storage. It stores binary informations of program instructions as well as data.

#### 7.1. Memory Hierarchy

The purpose of all memory systems is storage and retrieval of data. But a wide variety of memories exists with contrasting performance and different characteristics. Some of the fundamental characteristics of memory are:

- \* read - write ability
- \* access time
- \* storage permanence and
- \* storage capability.

The memory hierarchy summarizes the memory types and their performance using pyramidal structure. The fig below shows the memory hierarchy.

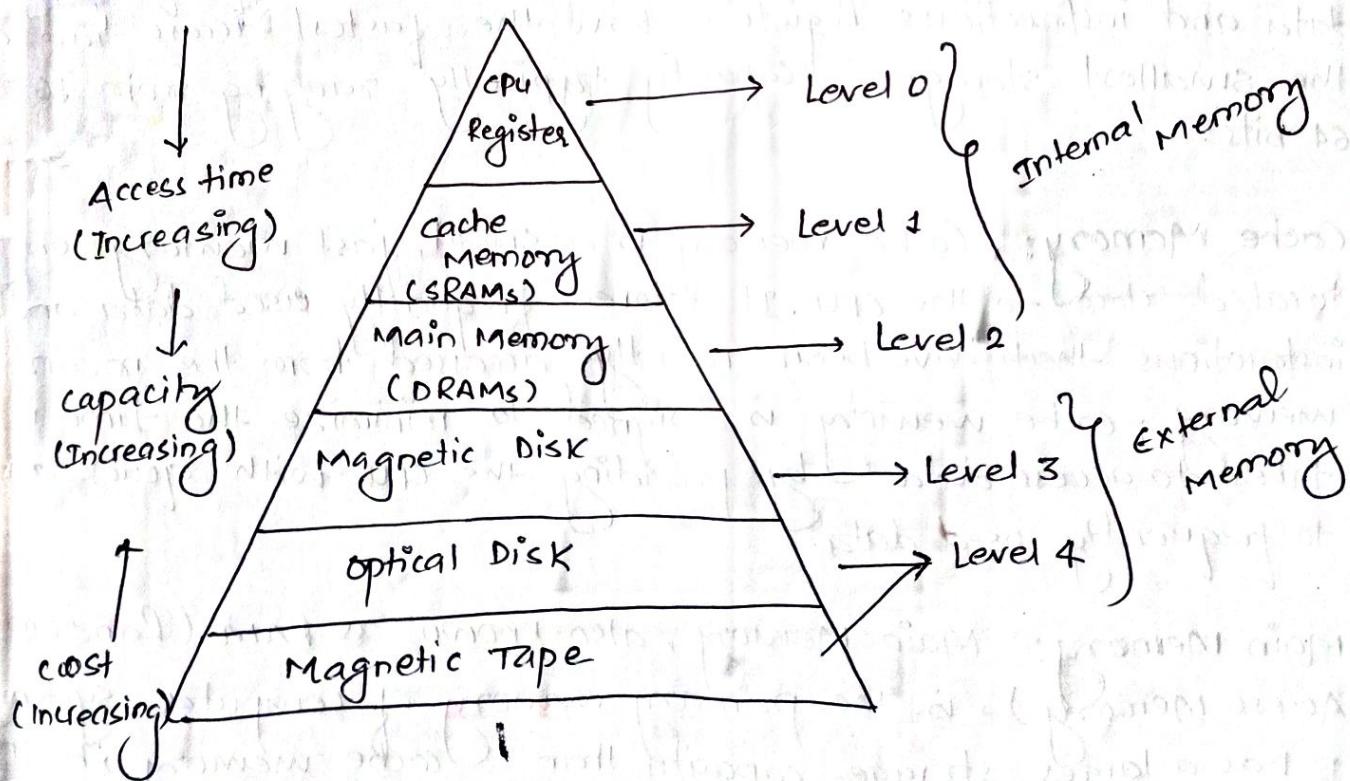


fig: Memory Hierarchy Design.

As one goes down the hierarchy

1. Distance from processor increases
2. Access time increases
3. Memory capacity increases
4. cost per bit decreases.

Hierarchy design is divided into 2 main types:

- a. External Memory or secondary memory: comprising of magnetic disk, optical disk, magnetic tape i.e. peripheral storage devices which are accessible by the processor via I/O modules.
- b. Internal Memory or Primary Memory: comprising of main memory, cache memory and CPU registers. This is directly accessible by the processor.

There are typically four levels of memory in a memory hierarchy

Registers: Registers are small, high speed memory units located in the CPU. They are used to store the most frequently used data and instructions. Registers have the fastest access time and the smallest storage capacity, typically ranging from 16 to 64 bits.

Cache Memory: Cache memory is a small, fast memory unit located close to the CPU. It stores frequently used data and instructions that have been recently accessed from the main memory. Cache memory is designed to minimize the time it takes to access data by providing the CPU with quick access to frequently used data.

Main Memory: Main Memory, also known as RAM (Random Access Memory), is the primary memory of computer system. It has a larger storage capacity than Cache memory, but it is slower. Main memory is used to store data and instructions that are currently in use by the CPU.

## Types of Main Memory:

- » **Static RAM**: It stores the binary information in flip flops and information remains valid until power is supplied. It has faster access time and is used in implementing cache memory.
- » **Dynamic RAM**: It stores the binary information as a charge on the capacitor. It requires refreshing circuitry to maintain the charge on the capacitors after few milliseconds. It contains more memory cells per unit area as compared to SRAM.

## Secondary Storage:

Secondary storage, such as hard disk drives (HDD) and solid state drives (SSD), is a non-volatile memory unit that has a larger storage capacity than the main memory. It is used to store data and instructions that are not currently in use by the CPU. Secondary storage has the lowest access time and is typically the least expensive type of memory in the memory hierarchy.

7.2

## Associative Memory:

An associative memory can be treated as memory unit whose saved information can be recognized for approach by content of the information itself instead of by an address or memory location.

Associative memory is also known as Content Addressable Memory (CAM).

Associative memory is accessed simultaneously and in parallel on the basis of data content rather by specific address or location.

The block diagram of associative memory is shown in fig. below:

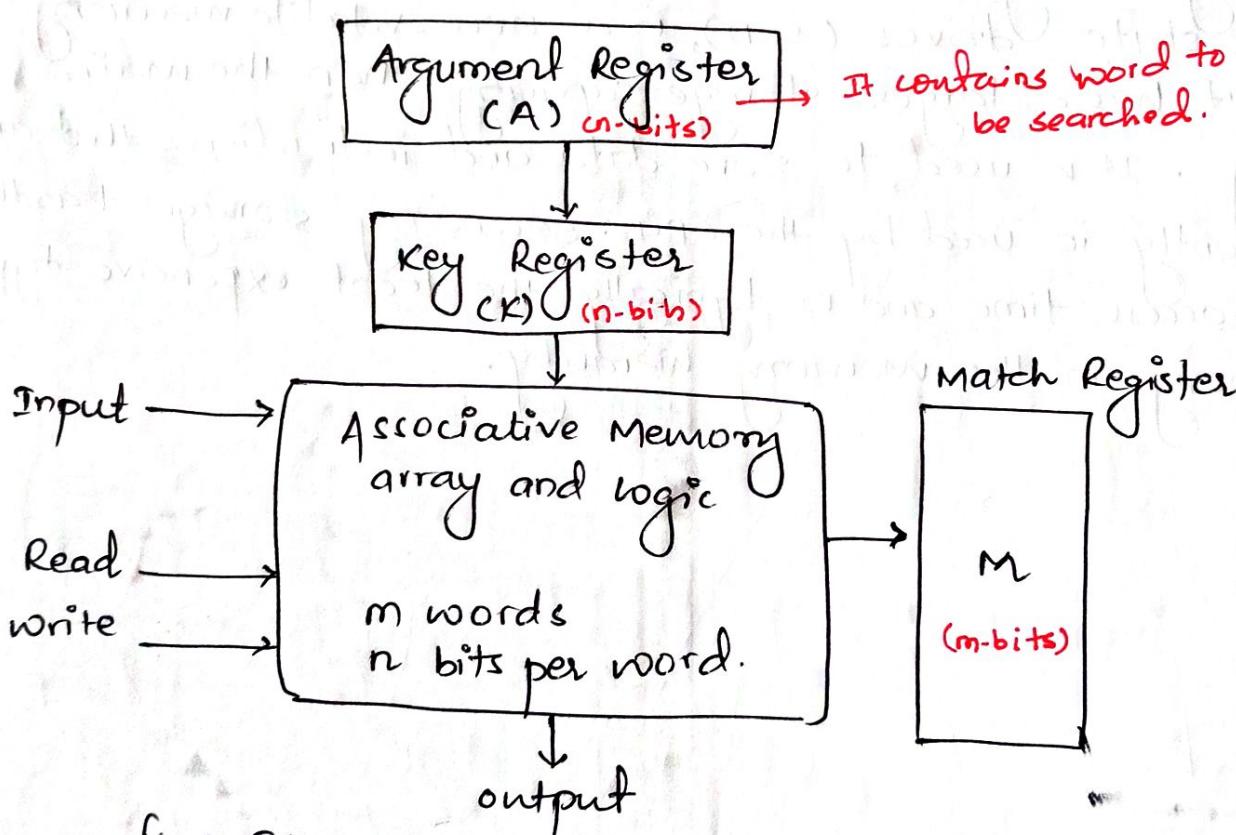


fig: Block diagram of Associative Memory.

- ↳ It includes a memory array and logic for  $m$  words with  $n$  bits per word. The argument register A and key register K each have  $n$  bits, one for each bit of a word.
- ↳ The match register M has  $m$  bits, one for each memory word. Each word in memory is related in parallel to the content of the argument register.

The words that connect the bits of the argument register set an equivalent bit in the match register. After the matching process, those bits in the match register that have been set denote the fact that their equivalent words have been connected.

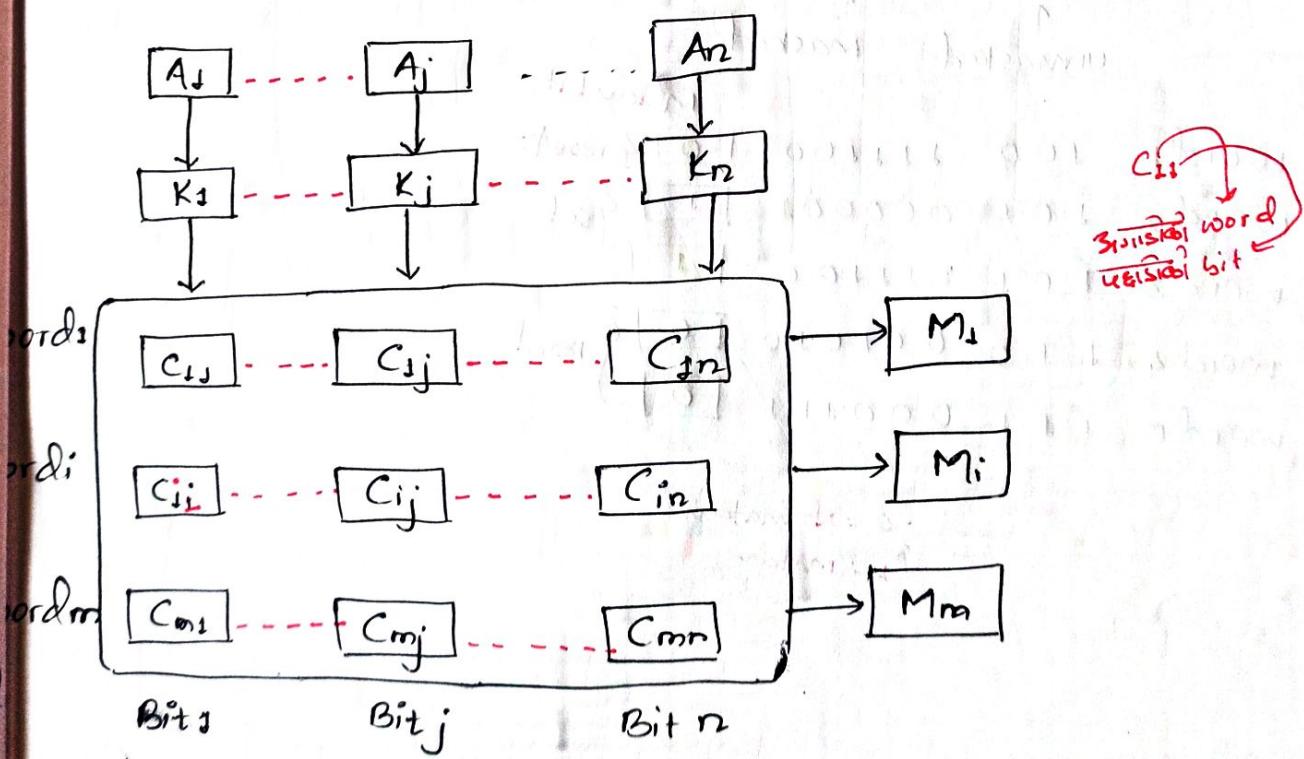
The key register supports a mask for selecting a specific field or key in the argument word. The whole argument is distinguished with each memory word if the key register includes all 1's.

Hence, there are only those bits in the argument that have 1's in their equivalent position of the key register are compared.

Therefore, the key gives a mask or seed recognizing a piece of data that determines how the reference to memory is created.

The following figure can define the relation between the memory array and the external registers in associative memory.

Associative Memory of  $m$  word,  
 $n$  cells per word



The cells in the array are considered by the letter C with two subscripts. The first subscript provides the word number and the second determines the bit position in the word. Therefore cell  $c_{ij}$  is the cell for bit  $j$  in word  $i$ .

A bit in the argument register is compared with all the bits in column  $j$  of the array supported that  $K_j = 1$ . ~~Therefore~~ cell  $c_{ij}$  is the cell for bit  $j$  in word  $i$ . ~~This~~ This is completed for all columns  $j = 1, 2, \dots, n$ .

If a match appears between all the unmasked bits of the argument and the bits in word  $i$ , the equivalent bit  $M_i$  in the match register is set to 1. If one or more unmasked bits of the argument and the word do not match,  $M_i$  is cleared to 0.

e.g:

A	101	111000
K	111	000000

↑                   ↑  
unmasked        masked.

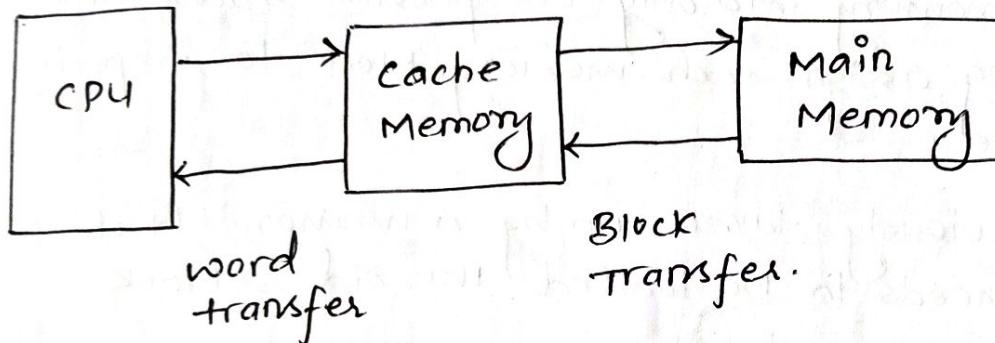
			Match bit
word <sub>1</sub>	100	111000	0 {reset}
word <sub>2</sub>	101	000001	1 {set}
word <sub>3</sub>	101	111100	1
word <sub>4</sub>	110	001010	0 {reset}
word <sub>5</sub>	111	000111	0

; 0 = No match  
1 = Match

### 7.3 Cache Memory and Cache Mapping Technique.

The main memory is slower than processor. A designer would like to build a main memory as fast as the processor but that would significantly increase the cost. A direct and frequent access to the main memory would slow down the processor operation. To address this problem, a faster memory whose speed is same as that of the processor is employed which is known as cache.

This memory resides within the processor chip and contains a copy of block of main memory so that processor can access contents faster. Cache memory comprises of SRAM cells which are expensive but faster to access.



When CPU needs to access some memory, cache is examined. If the desired word is available in the cache it indicates cache hit and processing continues. If it is not available in cache, it indicates cache miss and block of words including the desired word is transferred into cache from main memory.

Cache connects to the processor via address, data and control bus. The address and data bus are connected to the address and data buffer as well.

$$\text{Hit ratio (H)} = \frac{\text{hit}}{(\text{hit} + \text{miss})} = \frac{\text{no. of hits}}{\text{Total accesses}}$$

$$\text{Miss ratio (M)} = \frac{\text{miss}}{(\text{hit} + \text{miss})} = \frac{\text{no. of miss}}{\text{Total accesses}} = 1 - \text{hit ratio (H)}$$

## #. Cache Mapping Technique

Basically cache mapping refers to a technique using which the content present in the main memory is brought into the memory of the cache.

There are three different types of mapping used for the purpose of cache memory which is as follows:

1. Direct Mapping
2. Associative Mapping
3. set-Associative Mapping.

### 1. Direct Mapping:

- ↳ The simplest technique, known as direct mapping, maps each block of main memory into only one possible cache line or In direct mapping, assign each memory block to a specific line in the cache.
- ↳ If a line is previously taken up by a memory block when a new block needs to be loaded, the old block is trashed.
- ↳ An address is split into two parts index field and tag field. The cache is used to store the tag field when the rest is stored in the main memory. Direct mapping performance is directly proportional to the hit ratio.

$$i = j \bmod m$$

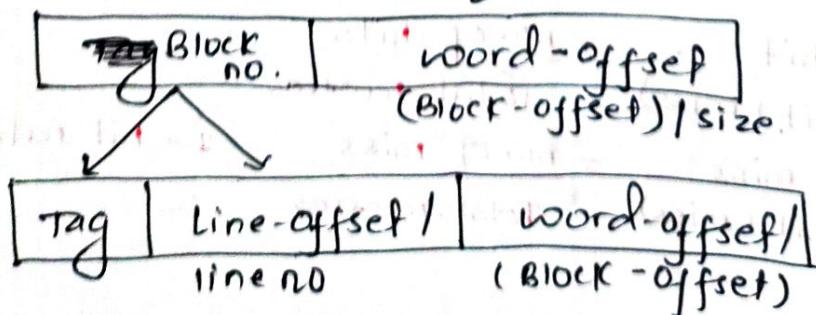
where

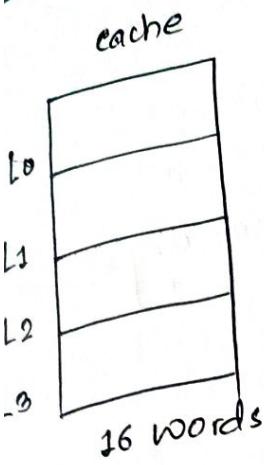
i = cache line number

j = main memory block number

m = number of lines in the cache.

Main  
memory  
  
cache  
memory





## Main Memory

W <sub>0</sub> W <sub>1</sub> W <sub>2</sub> W <sub>3</sub>	B <sub>0</sub>
W <sub>4</sub> W <sub>5</sub> W <sub>6</sub> W <sub>7</sub>	B <sub>1</sub>
W <sub>8</sub> W <sub>9</sub> W <sub>10</sub> W <sub>11</sub>	B <sub>2</sub>
W <sub>12</sub> W <sub>13</sub> W <sub>14</sub> W <sub>15</sub>	B <sub>3</sub>
W <sub>16</sub> W <sub>17</sub> W <sub>18</sub> W <sub>19</sub>	B <sub>4</sub>
⋮	⋮
W <sub>24</sub> W <sub>25</sub> W <sub>26</sub> W <sub>27</sub>	B <sub>31</sub>

128 words

(1 byte = 1 word)

one block contains 4 words.

$$\text{Total No of blocks} = \frac{128}{4} = 32 \text{ Blocks.}$$

Line size = Block size

$$\text{No. of lines} = \frac{32}{4} = 4 \text{ lines.}$$

i.e have,

$$i = j \bmod m$$

Let  $j=0$  and  $m=4$

$j=1$  and  $m=4$

$j=2$  and  $m=4$

$j=3$  and  $m=4$

$j=4$  and  $m=4$

Then,  $i = 0 \bmod 4 = 0$  (i.e. Line no zero)

Then,  $i = 1 \bmod 4 = 1$  (i.e. L<sub>1</sub>)

Then,  $i = 2 \bmod 4 = 2$  (i.e. L<sub>2</sub>)

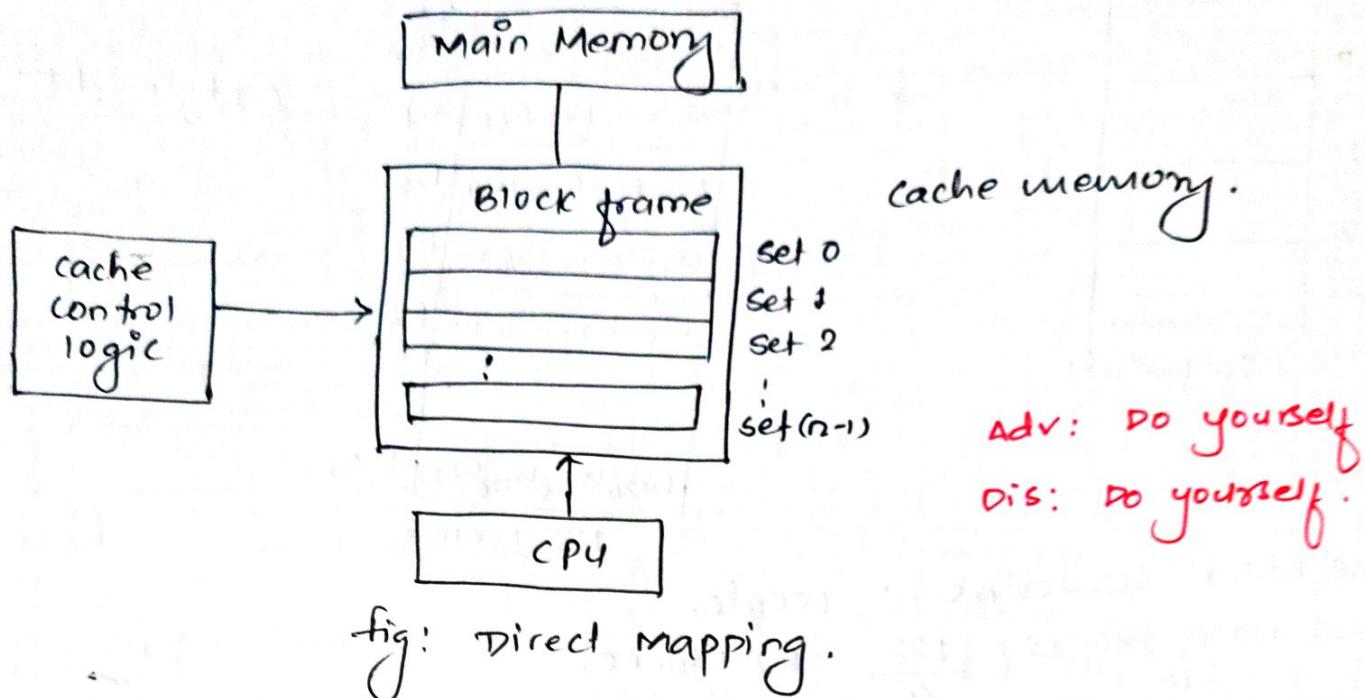
Then,  $i = 3 \bmod 4 = 3$  (i.e. L<sub>3</sub>)

Then,  $i = 4 \bmod 4 = 0$  (i.e. L<sub>0</sub>)

## Cache Memory

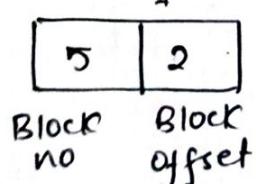
L <sub>0</sub>	B <sub>0</sub> B <sub>4</sub> B <sub>8</sub> B <sub>12</sub> ... B <sub>28</sub>
L <sub>1</sub>	B <sub>1</sub> B <sub>5</sub> B <sub>9</sub> B <sub>13</sub> ... B <sub>29</sub>
L <sub>2</sub>	B <sub>2</sub> B <sub>6</sub> B <sub>10</sub> B <sub>14</sub> ... B <sub>30</sub>
L <sub>3</sub>	B <sub>3</sub> B <sub>7</sub> B <sub>11</sub> B <sub>15</sub> ... B <sub>31</sub>

When we store main memory address content in cache, we also store the tag. To determine if desired main memory content is available in the cache or not, we go to the cache address indicated by the index and we compare the tag available there with the desired tag.



## 2. Associative Mapping

- In this type of mapping, associative memory is used to store the content and addresses of the memory word. Any block can go to any line of the cache. This means that the word id bits are used to identify which word in the block is needed, but the tag becomes all of the remain bits. This enables the placement of any word at any place in the cache memory.
- It is considered to be the fastest and most flexible mapping form. In associative mapping, the index bits are zero.
- We have, the physical address (PA) is of 7 bit (i.e.  $2^7 = 128$  no



- In cache memory, there is no need of line-offset block. Therefore our tag is the block no.

$$\text{i.e. } \text{Tag} = \text{Block no}$$

This shows that  $2^5 = 32$  (0-32) blocks can appear in any line in cache.

$L_0$	$B_0$	$B_1$	$B_2$	$B_3$	$B_4$	---	$B_{31}$
$L_1$	$B_0$	$B_1$	$B_2$	$B_3$	$B_4$	---	$B_{31}$
$L_2$	$B_0$	$B_1$	$B_2$	$B_3$	$B_4$	---	$B_{31}$
$L_3$	$B_0$	$B_1$	$B_2$	$B_3$	$B_4$	---	$B_{31}$

### Advantage:

- conflict miss issue is resolved
- hit ratio/rate is increased
- hit ratio/rate is increased

### Disadvantage:

Primary disadvantage is that the comparison time is increased.

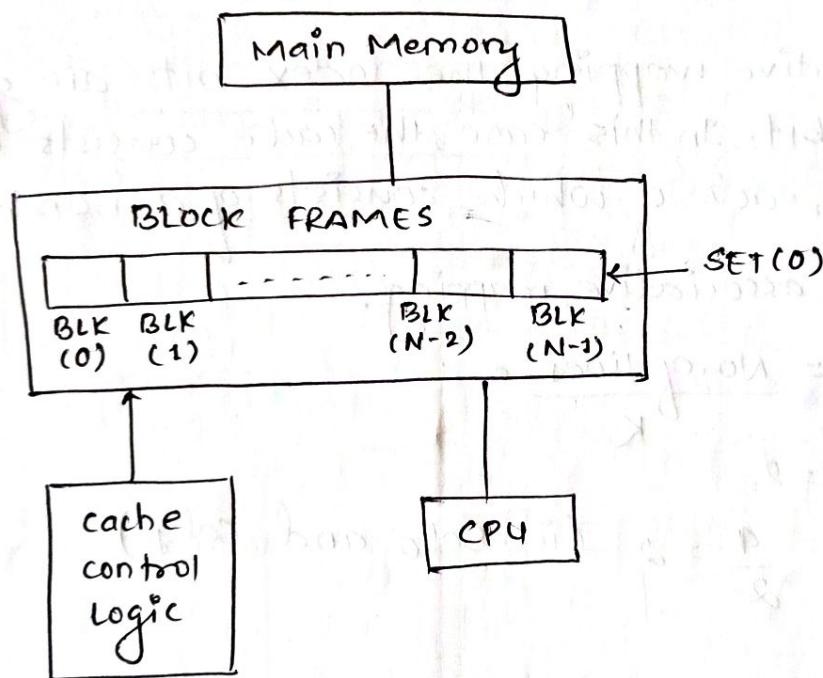


fig: Associative mapping - structure.

### 3. Set-Associative Mapping

- This form of mapping is an enhanced form of direct mapping where the drawbacks of direct mapping are removed. Set associative addresses the problem of possible trashing in the direct mapping method. It does this by saying that instead of having exactly one line that a block can map to in the cache. We will group a few lines together creating a set.
- ↳ also known as K-way set associative.
  - ↳ Then a block in memory can map to any one of the lines of a specific set.
  - ↳ Set-associative mapping allows each word that is present in the cache can have two or more words in the main memory for the same index address.
  - ↳ In set associative mapping the index bits are given by the set offset bit. In this case, the cache consists of a number of sets, each of which consists of a number of lines.
  - ↳ for 2-way set associative mapping.

$$\text{No of set} = \frac{\text{No. of lines}}{K}$$

for  $K=2$ ,

$$\text{No of set} = \frac{4}{2} = 2 \quad (\text{ie. Set}_0 \text{ and Set}_1)$$

$L_0$	$B_0 \ B_2 \ B_4 \dots \ B_{30}$
$L_1$	$B_0 \ B_2 \ B_4 \dots \ B_{30}$
$\dots$	$\dots$
$L_2$	$B_1 \ B_3 \ B_5 \dots \ B_{31}$
$L_3$	$B_1 \ B_3 \ B_5 \dots \ B_{31}$

formula:

$$i = j \bmod n$$

$i$  = Line number

$j$  = Block number

$n$  = no. of set.

e.g.:

$$0 \bmod 2 = 0 \quad (\text{Set}_0)$$

$$1 \bmod 2 = 1 \quad (\text{Set}_1)$$

$$2 \bmod 2 = 0 \quad (\text{Set}_0)$$

$$3 \bmod 2 = 1 \quad (\text{Set}_1)$$

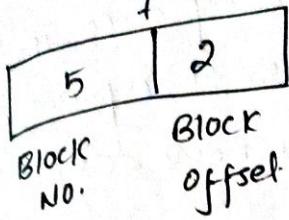
Note:  $B_0, B_2, B_4, \dots, B_{30}$  line 0 or line 1

मुख्य लाइन में स्थान संवधि।

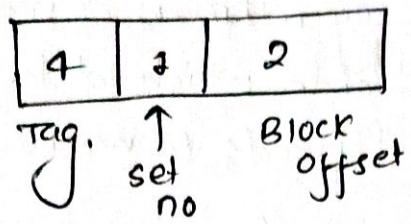
यहाँ तीसरी लाइन 2 या लाइन 3

में स्थान संवधि।

main memory



cache memory.



(ie.  $2^1 = 2$  = 80 and s1)

Here, tag bit = 4 =  $2^4 = 16$  combination.

i.e. B0, B2, B4, ..., B30 (16 Blocks)

OR,

B1, B3, B5, ..., B31 (16 Blocks)

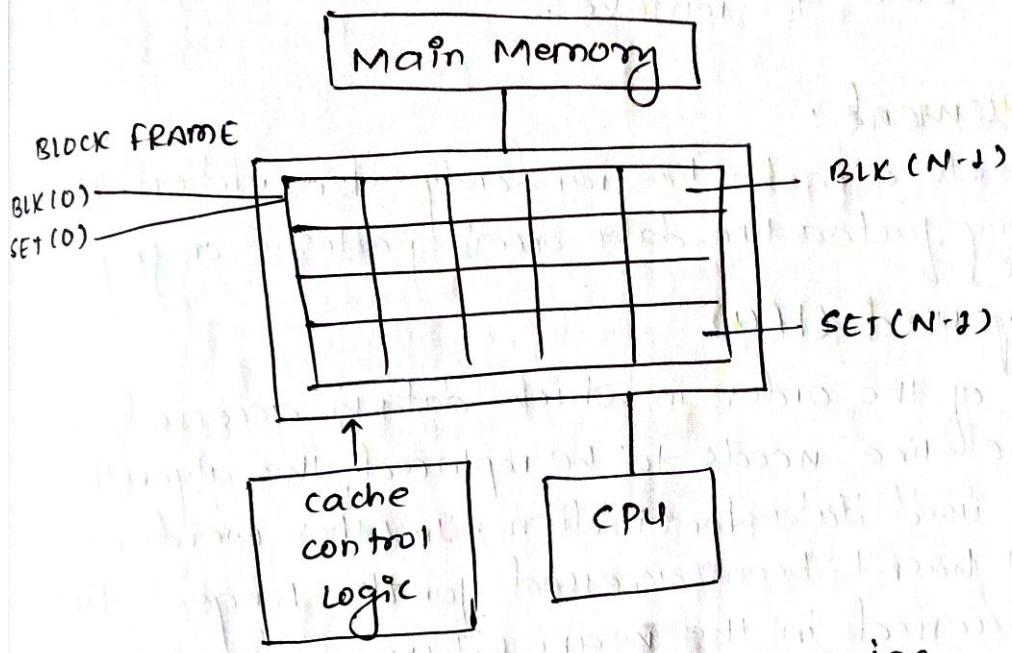


fig: set - Associative mapping.

### Advantage:

The comparison time is reduced in the K-way set associative mapping as compared to the fully associative mapping.

### Disadvantage:

It involves the factor of conflict miss.

## 7.4. Cache Replacement Algorithm:

Cache replacement algorithm is a method used by a computer's cache memory to decide which data to remove or replace when the cache is full and needs to make room for new data.

Cache replacement algorithms help determine which data should be evicted from the cache to make space for new data. These algorithms aim to maximize the cache hit rate and minimize cache misses, as cache hits result in faster access time.

↳ Different cache replacement algorithms use various strategies to decide which data to remove.

### a) Random Replacement:

This algorithm selects a cache line randomly for eviction. It does not consider any factors like data access patterns or frequency.

### b) Least Recently Used (LRU)

LRU keeps track of the order in which data is accessed in the cache. When cache line needs to be replaced, the algorithm selects the least recently used data for eviction. In other words, it assumes that the data that hasn't been accessed for the longest time is less likely to be accessed in the near future.

LRU cache.

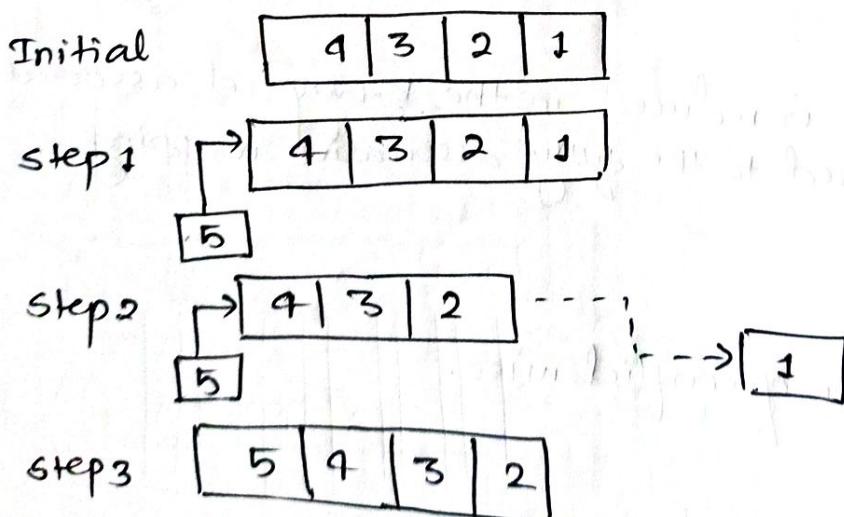


fig: LRU cache

### c) Least frequently used (LFU)

the LFU cache replacement algorithm is designed to remove or replace data from a cache based on the frequency of its access. It assumes that data items with lower access frequencies are less likely to be accessed in the near future, so it prioritizes evicting those items.

The LFU algorithm maintains a counter for each data item in the cache, which keeps track of how many times that item has been accessed. When the cache becomes full and needs to make room for new data, the LFU algorithm identifies the data items with the lowest access frequency and removes it from the cache.

### d) FIFO (first in first out)

FIFO replaces the data that has been in the cache the longest. It works on the principle of a queue, where the first data to enter the cache is the first one to be evicted when necessary.

Objects are added to the queue and are evicted with the same order. Even though it provides a simple and low-cost method to manage the cache but even the most used objects are eventually evicted when they're old enough.

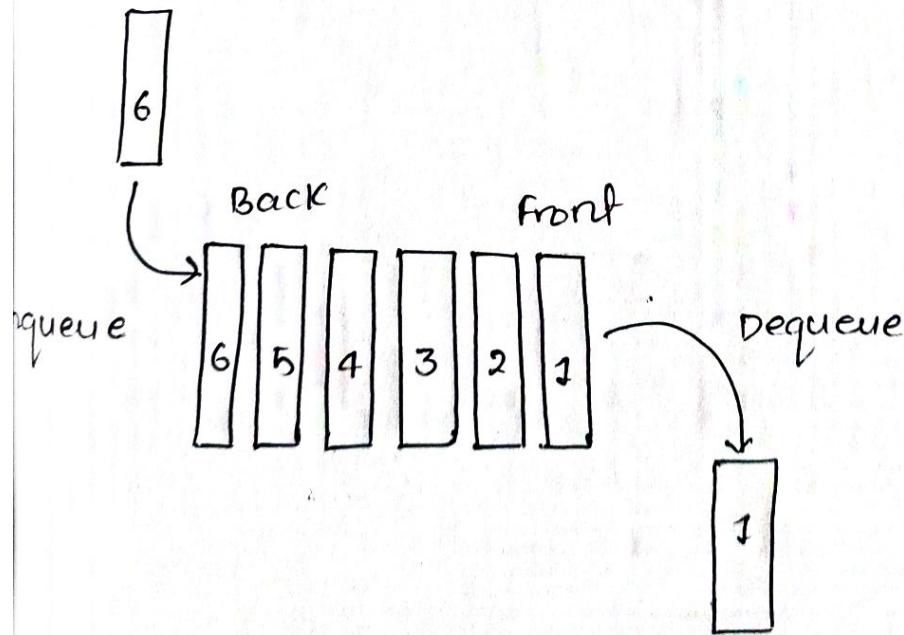


fig: FIFO Technique.