

bandwidth : is the Data transfer capacity of a computer Network in bits per second (Bps).

## Unit - III

Memory Interleaving : (memory Modules)

⇒ It is a Technique that divides memory into a number of modules such that successive words in the address space are placed in the different modules.

⇒ Memory Interleaving increases bandwidth by allowing simultaneous access to more than one chunk of memory.

⇒ This improves performance because the word = 2 bytes processor can transfer more information to / from memory in the same amount of time.

Divides system memory into multiple Blocks.

2 Methods → 2 way Interleaving

( 2 memory blocks are accessed at same time for Reading & writing op.)

4 way Interleaving

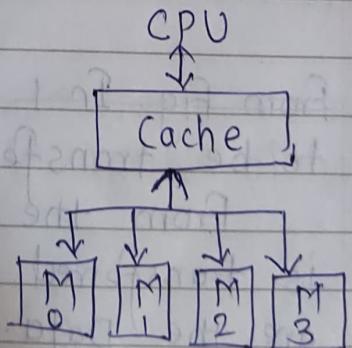
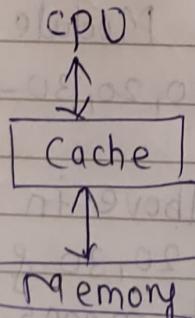
( 4 memory blocks are accessed at same time )

1 fork 1 fork

1 people

2 fork 2 fork

2 people



How it is done ?

By spreading memory addresses evenly across memory Modules.

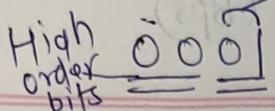


fig. Consecutive Word in a Module

0	0000	10				
1	0001	20	Module 00			Module 01
2	0010	30	00	10	00	50
3	0011	40	01	20	01	60
4	0100	50	10	30	10	70
5	0101	60	11	40	11	80
6	0110	70				
7	0111	80				
8	1000	90	Module 10		Module 11	
9	1001	100	00	90	00	130
10	1010	110	01	100	01	140
11	1011	120	10	110	10	150
12	1100	130	11	120	11	160
13	1101	140				
14	1110	150				
15	1111	160				

e.g.) 4-way interleaved memory  
4 modules:

Let's Assume 16 Data's to be Transferred to the 4 Module (here Module 00 be Module 1, Module 01 be Module 2, Module 10 be Module 3, Module 11 be Module 4).  
From Fig. In 1 Also 10, 20, 30---160 are the Data to be transferred.

From the fig. Above in Module 1, 10 [data] is transferred then 20, 30 & finally 40 which are the Data. That means the data are added consecutively in the Module till its max capacity.

Most significant bit (MSB) provides the Address of the Module.

Least significant bit (LSB) provides the Address of the Data (in the module).

27/9/22 → 43,

classmate

Module Data  
MSB LSB

For Example, to get  $g_0(\text{Data})$   $\overline{1000}$  will be provided by the processor.

Module 0 contains Data : 10, 20, 30, 40

Module 1 contains Data : 50, 60, 70, 80

Module 2 contains Data : 90, 100, 110, 120

Module 3 contains Data : 130, 140, 150, 160

Consecutive Word in Consecutive Module

	Module 00	Module 01	Module 10	Module 11
00	10	20	30	40
01	50	60	70	80
10	90	100	110	120
11	130	140	150	160

Least Significant Bit (LSB) : Address of Module

(MSB) : Address of Data

MSB      LSB  
10      00

\* A memory is a storage unit - This is a place to hold data & instructions.

Memory Location: Memory is divided into multiple small parts of fixed size & capable of holding info



## Hierarchy of Memory

A Memory unit is the collection of storage units or devices together. The memory unit stores the Binary information in the form of bits. Generally Memory storage is classified into 2 categories.

RAM  $\Rightarrow$  Volatile Memory: This loses its data, when power is switched off.

ROM  $\Rightarrow$  Non Volatile Memory: This is a permanent storage & does not lose any data when power is switched off.

$\Rightarrow$  The total memory capacity of a Computer can be visualized by hierarchy of components.

$\Rightarrow$  The memory hierarchy system consists of all storage devices contained in a computer system.

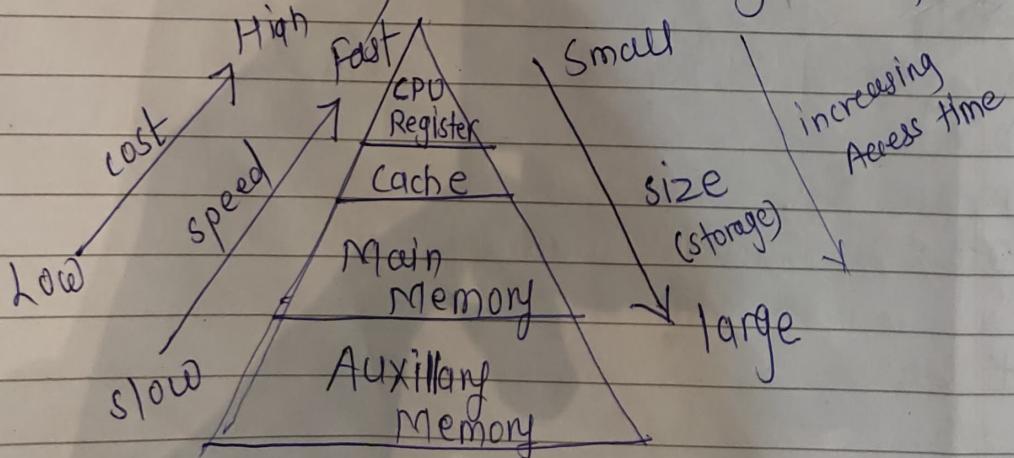
External Memory or Secondary Memory

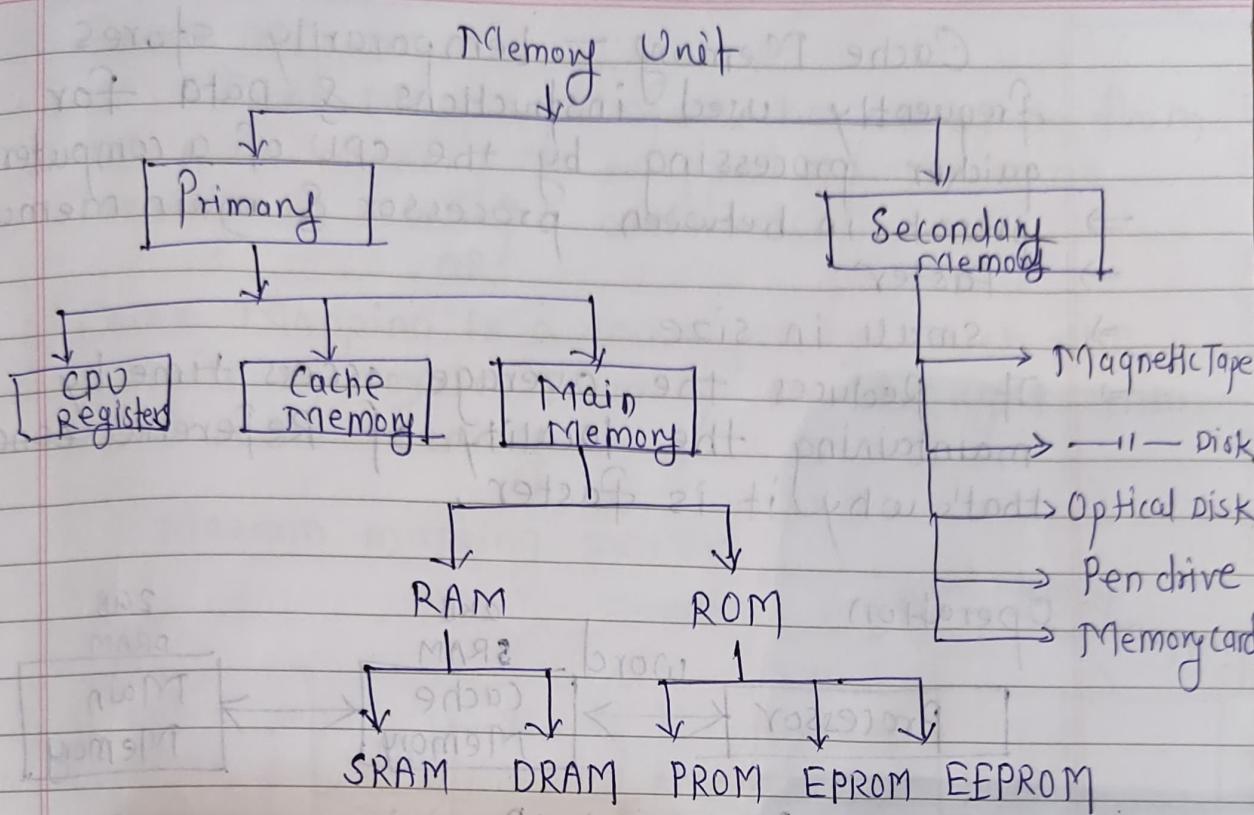
Magnetic Disk, Optical Disk, Magnetic Tape i.e. peripheral storage devices which are accessible by the processor via I/O Module.

Internal Memory or Primary Memory

Main Memory, Cache Memory & CPU Registers.

This is Directly Accessible by the processor.





### Secondary (Auxillary) Memory

Magnetic Tape (seq. Access)

Magnetic Disk → HDD  
FDD

Optical Disk → CD, DVD, BD, HVD

Flash Memory → Memory card  
Pendrive

PROM Programmable Read only Memory

EPROM Erasable

EEPROM Electrically erasable

*(Note: PROM, EPROM, EEPROM are collectively known as pocket money.)*

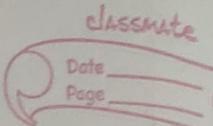
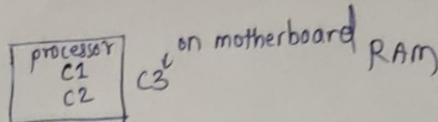
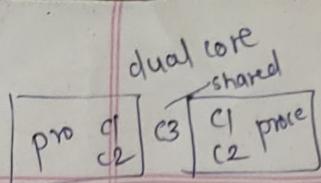
cupboard

Bank

Cache mem

RAM

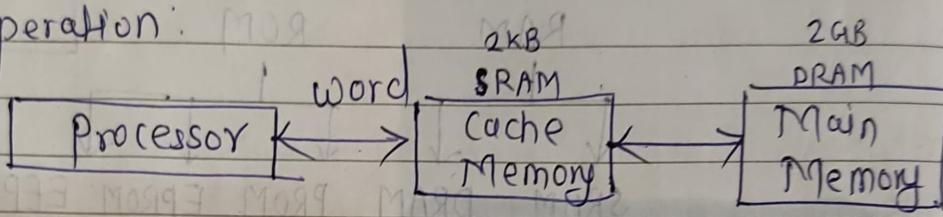
Hard disk



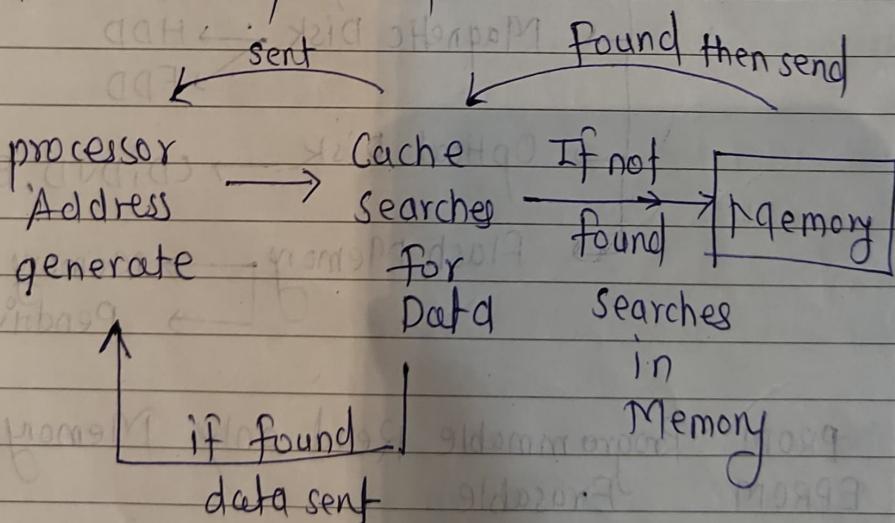
**Cache Memory** — temporarily stores frequently used instructions & data for quicker processing by the CPU of a computer.

- ⇒ present in between processor & main memory
- ⇒ faster
- ⇒ small in size
- ⇒ It reduces the average access time by maintaining the locality of reference concept that's why it is faster.

**Operation:**



- ⇒ Transfer of data b/w main memory & cache memory is block by block
- ⇒ Transfer of data b/w cache & processor is word by word.



cache hit → Address / data is found in cache memory  
 cache miss → —————— is not found.

The process of transfer the Data from main memory to cache Memory is called as mapping

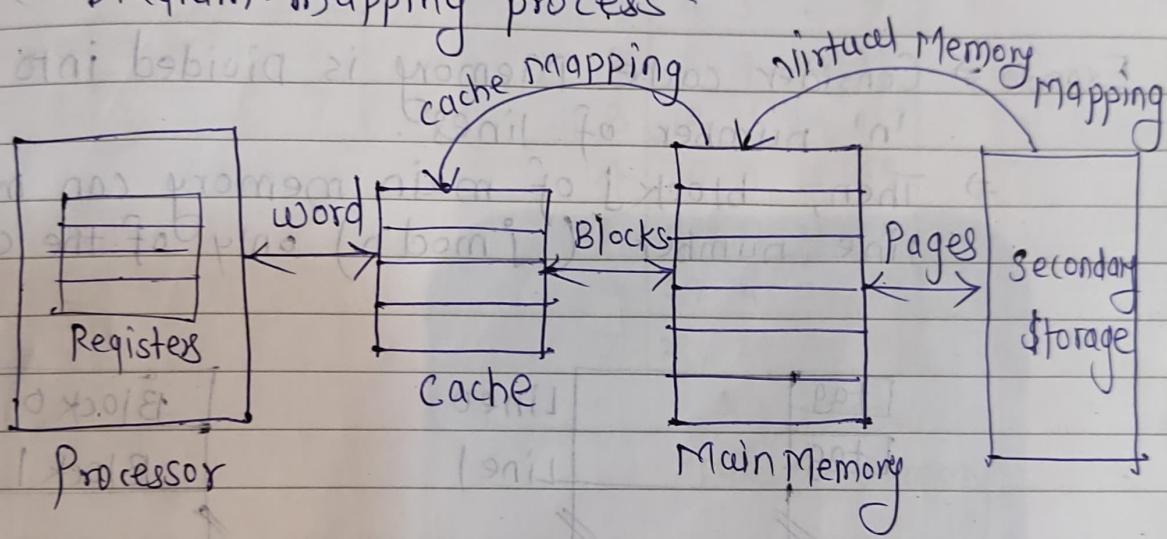
### Cache Mapping

Cache Mapping defines how a block from the main memory is mapped to the cache memory in case of a cache miss.

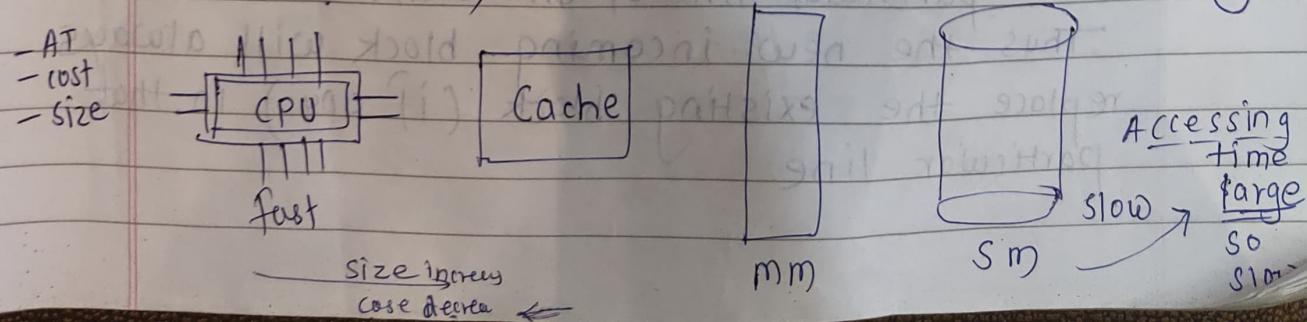
OR

Cache Mapping is a technique by which the contents of main Memory are brought into the cache Memory.

Diagram mapping process -



- ⇒ Main Memory is divided into equal size partitions called as blocks or frames.
- ⇒ Cache Memory is divided into partitions having same size as that of blocks called as lines.
- ⇒ During cache Mapping, block of main Memory is simply copied to the cache & the block is not actually brought from main Memory.



# Cache mapping Tech.

CLASSMATE

Date \_\_\_\_\_

Page \_\_\_\_\_

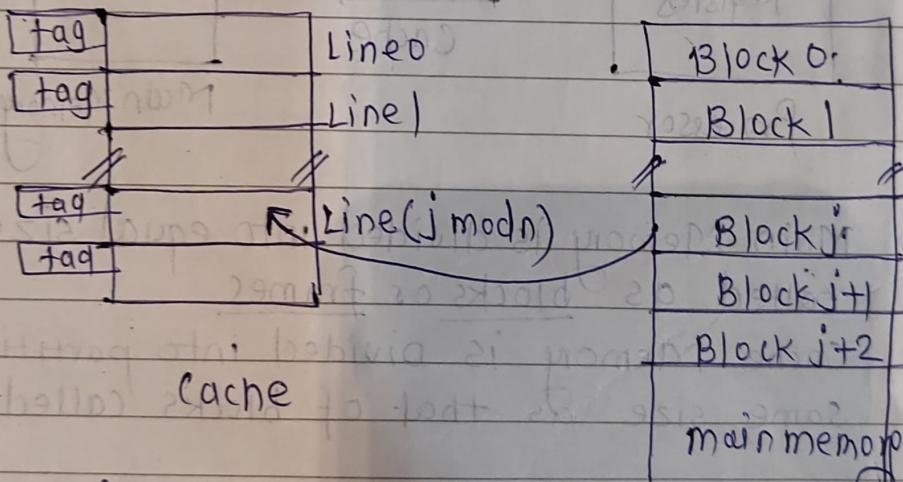
↓  
Direct      fully Associative      K-way Set Association

- 1) Direct Mapping : In Direct Mapping
- ⇒ A particular block of main memory can map only a particular line of the cache
  - ⇒ The line number of cache to which a particular block can map is given by

$$\text{Cache line Number} = \left( (\text{Main Memory Block Address}) \bmod n \right) \text{ mod } \left( \frac{\text{Number of lines in cache}}{\text{Number of lines in main memory}} \right)$$

Example ⇒ Consider cache memory is divided into 'n' number of lines.

- ⇒ Then, block  $j$  of main memory can map to line number  $(j \bmod n)$  only of the cache



Need of Replacement Algorithm : -

There is no need of any Replacement Algorithm because a main memory block can map only to a particular line of the cache.

Thus the new incoming block will always replace the existing block (if any) in that particular line.

In Direct Mapping, the physical Address is divided as

Tag	Line number	Block	Line offset

Block Number.

block	0	1	2	3
1	4	5	6	7
2	8	9	10	11
3	12	13	14	15
4	16	17	18	19
5	20	21	22	23
6	24	25	26	27
7	28	29	30	31

8 Blocks

8 words.

32 words

Main Memory

1 = word = 1 byte

1 Block size = 4 words

No. of Blocks =  $32/4 = 8$  blocks - words/size of block

No. of cache line =  $8/4 = 2$  lines

Cache line size = Block size

Cash line no. =  $(MM \text{ Block}^{\text{no}}) \bmod (\text{No. of lines in cache})$

$$\bullet = (0) \bmod (2)$$

= 0 Remainder  $\rightarrow$  Block 0  $\Rightarrow$  line 0

$$\bullet = (1) \bmod 2 \rightarrow$$

= 1 Remainder

$$\bullet = (2) \bmod 2 = 0^{\text{th}} \text{ line}$$

$$\bullet = (3) \bmod 2 = 1$$

## Physical Address

Tag	Line no.	Block size
2 bits	1 bit	2 bits
blocks	Block Number	
$2^3 = 8$		

physical Address      32 words       $2^{(5)} = 5 \text{ bits}$

$$\text{Block size} = 4 \text{ words} = 2^{(2)} = 2 \text{ bits}$$

$$\text{Line nos.} = 2 = 2^{(1)} = 1 \text{ bit}$$

$$\text{Total Blocks} = 8 \text{ blocks} = 2^{(3)}$$

e.g.

Bo

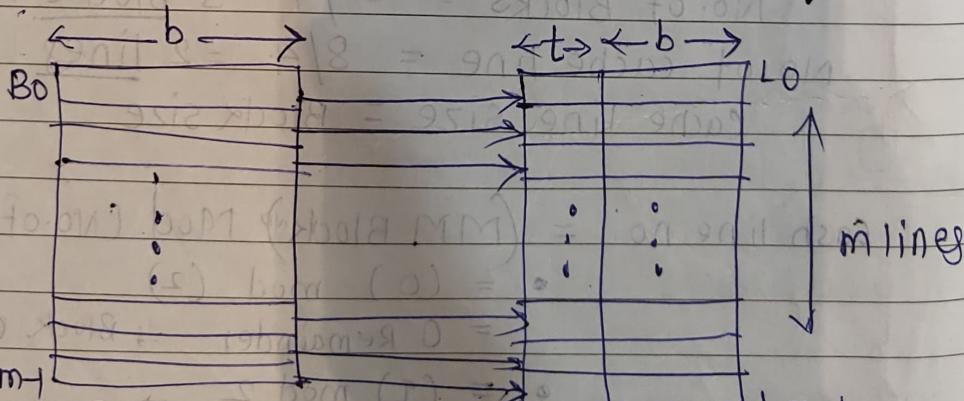
0	0	1	2	3	00	0	0	0	0	v
1	4	5	6	7	BB	1	0	0	1	

Phy Add.

e.g. 000001, hite.g. 010001, miss

Direct

Mapping



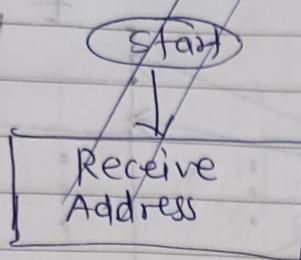
first m blocks  
of mm  
(equal to size of cache)

$b$  = length of Block in bits  
 $t$  = length of tag in bits

classmate  
Date \_\_\_\_\_

physical Address - refers to a memory address or the location of a memory cell in the main memory.

### Cache Read operation.

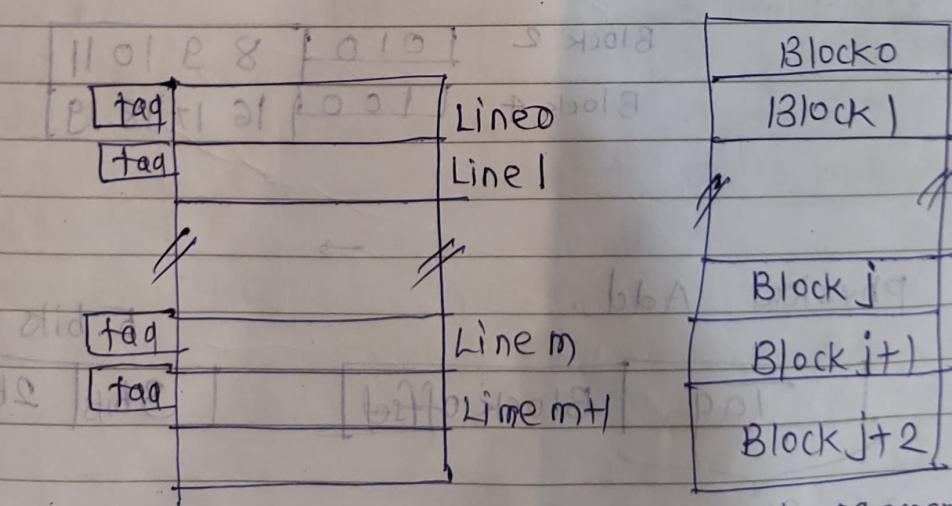


### Fully Associative Mapping:

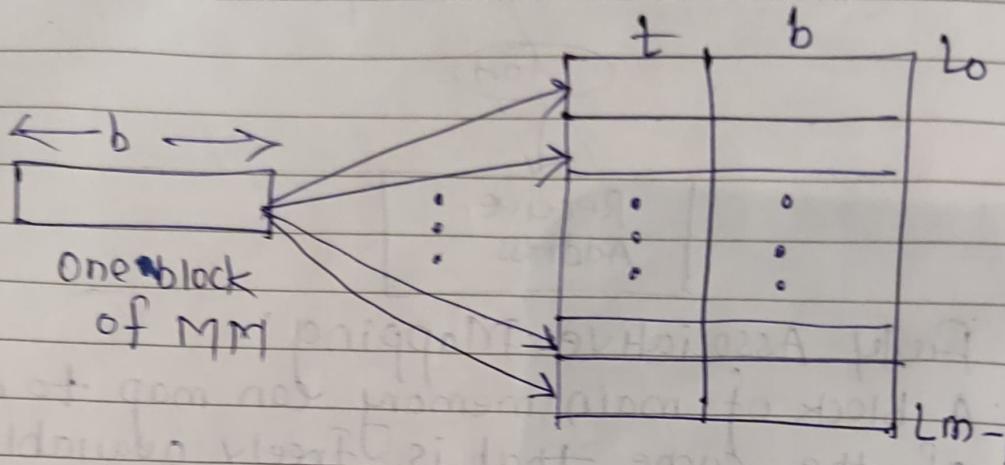
- ⇒ A block of main memory can map to any line of the cache that is freely available at that moment.
- ⇒ This makes Fully Associative mapping more flexible than Direct Mapping.
- ⇒ All the lines of cache are freely available.
- ⇒ Thus, any block of main memory can map to any line of the cache.
- ⇒ If all the cache lines occupied, then one of the existing blocks will have to be replaced.

### Division of physical Address

Block Number/Tag	Block/Line offset
------------------	-------------------



## Associative mapping



Cache Memory:

*(8 words)*

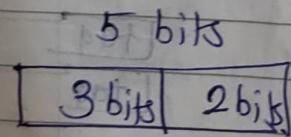
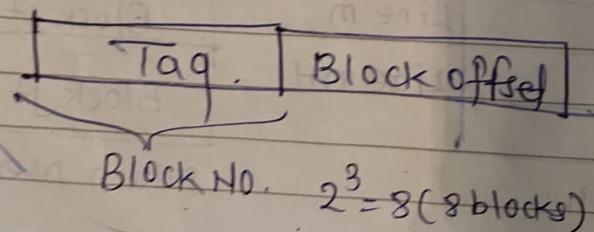
0	0 1 2 3	Block offset
1	4 5 6 7	00 0th
2	8 9 10 11	01 1st
3	12 13 14 15	10 2nd
4	16 17 18 19	11 3rd
5	20 21 22 23	
6	24 25 26 27	
7	28 29 30 31	

→ Phy. Address      ↓ block's content  
 e.g. 0 1 0 | 1 1      hit      32 words.

Block 2      Block 4

0 1 0	8 9 10 11	initial cache state
1 0 0	16 17 18 19	

Physical Add.



## Set Associative Mapping

0	0
1	
2	1
3	SAPT

16 words

0	0	1	2	3
1	4	8	12	16
2	B	9	10	11
3	SAPT			
4				
5				
6				
7				
8				
9				
10				
11				
12				
13				
14				
15	60	61	62	63

64 words

$$BS = 4 \text{ words}$$

$$NCL = 16$$

$$NB_S = \frac{64}{4} = 16 \text{ Blocks}$$

$$\frac{64}{4} = 16 \rightarrow \underline{\text{Blocks}}$$

$$2^8 = 2 \times 2 \times 2 \times 2 \times 2 \times 2 \times 2 = 64$$

### 2-way Set Associative

#### Physical Address

No. of sets	Set	Block NO.	Tag	set no.	Block offset
			set 0	set 1	set 2
0 mod 2 $\Rightarrow$ 0			000	0 1 2 3	B0
1 mod 2 $\Rightarrow$ 1			001	8 9 10 11	B2
2 mod 2 $\Rightarrow$ 0			000	4 5 6 7	B1
3 mod 2 $\Rightarrow$ 1			001	12 13 14 15	B3
4 mod 2 $\Rightarrow$ 0					

Initial state

0 0000	0 0 1 0 1 0	hit
1 0001	0 1 0 1 0 0	miss
2 0010		
3 0011		

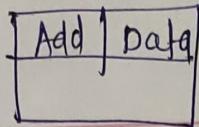
(82)

$2^6 = 64 = 6 \text{ bits phy Add.}$

hit  
many sets  
 $2^4 = 16$

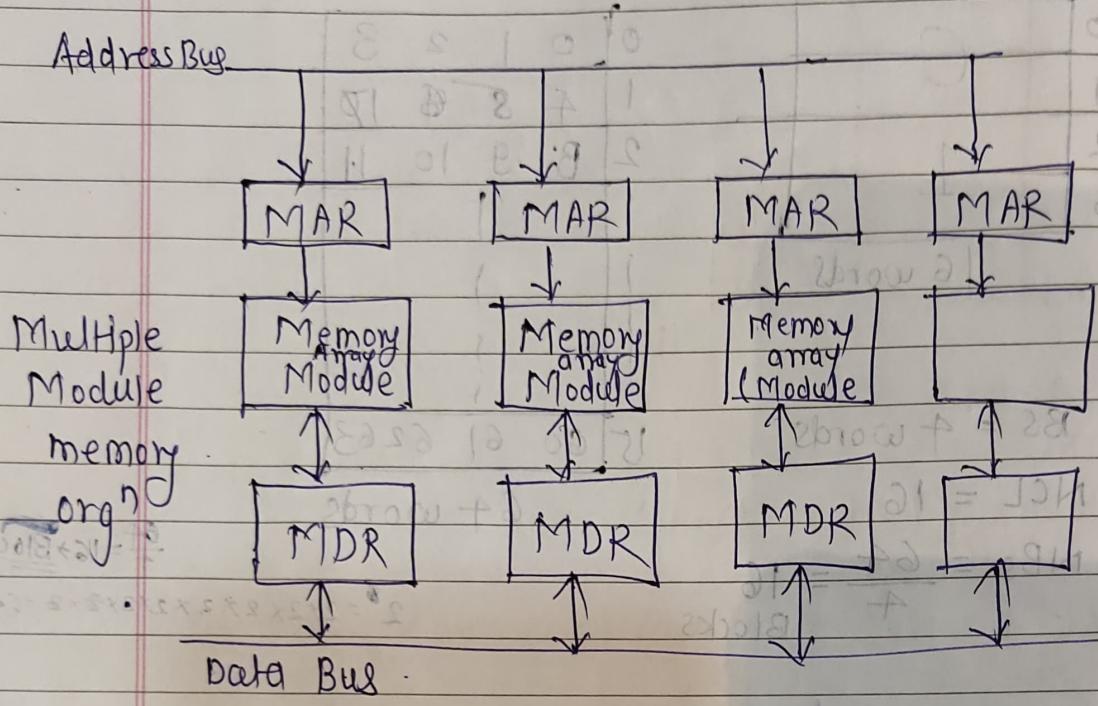
1 bit  
2 bits

Shanu Kuttan CSE classes



M.M.  
Memory Module is a memory array together with its address & Data registers.

Address Bus

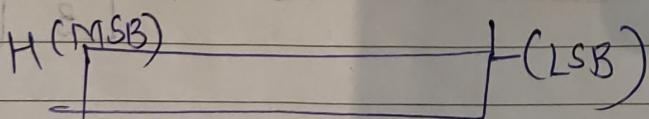


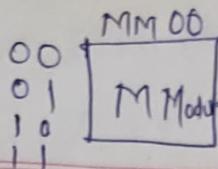
MAR Receives info<sup>m</sup> from a common Address Bus  
MDR communicate with a bidirectional data bus

System Performance is enhanced because read & write Activity ~~occurs~~ <sup>occurs</sup> simultaneously across the multiple modules in a similar fashion.

There are 2 - Address formats for Memory Interleaving the Address space

- ① High-Order Interleaving : High-order bits as Module Address & the lower order bits as the word address within each module



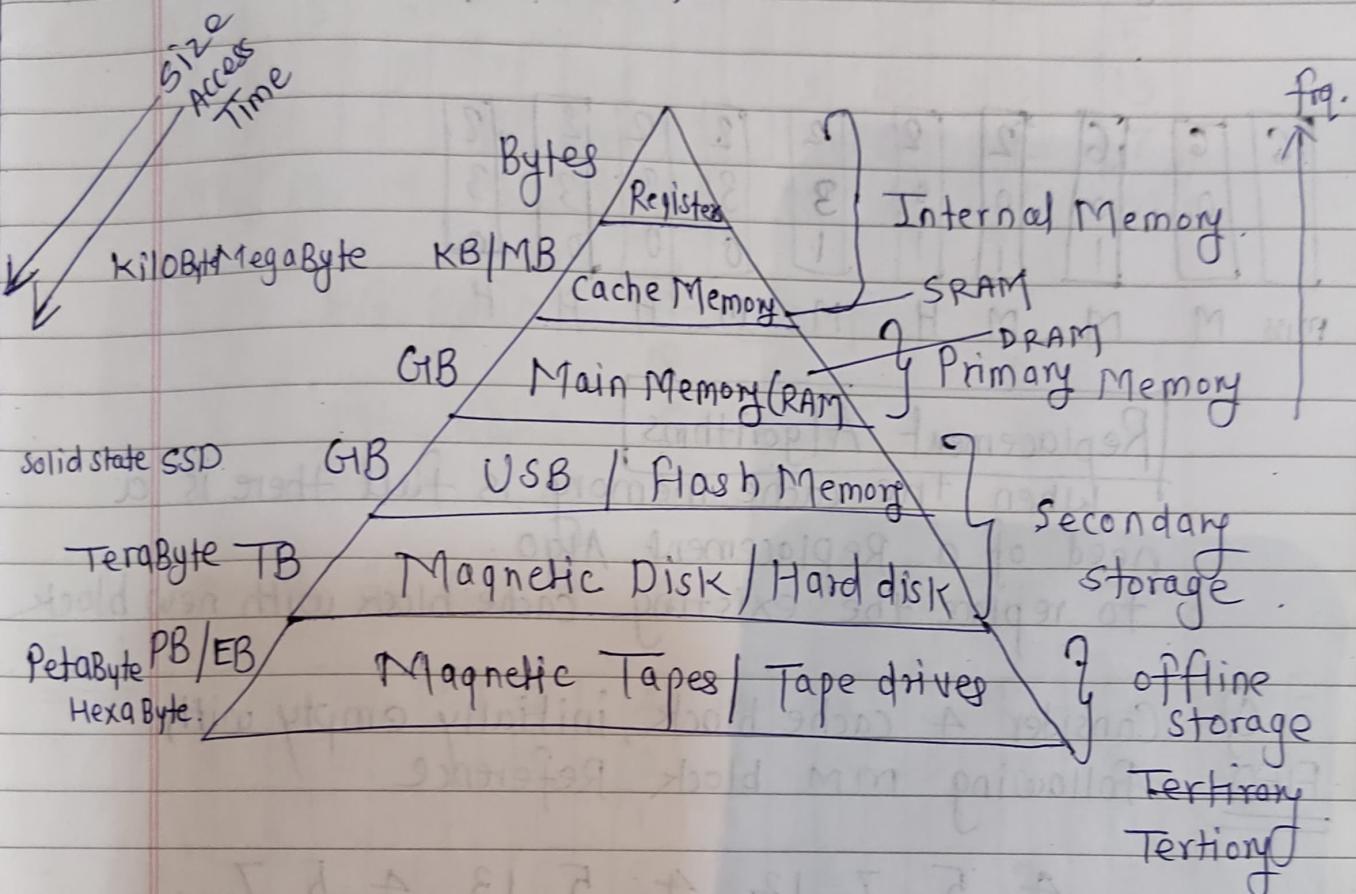


classmate

Date \_\_\_\_\_

Page \_\_\_\_\_

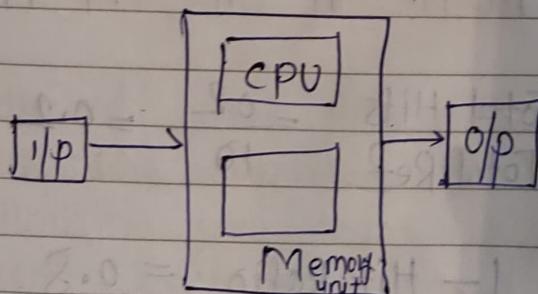
Low-order Interleaving: Low-order bits as module Address & the high-order bits as the word address



## Memory Organization:

Dif. betw Organization: - Peripherals Organized.

Architecture



→ K/b, mouse

→ Mouse

Memory.

Practically possible

① FIFO

② LRU (Least Recently Used)

③ Optimal → study purpose

When processor initiates,

$$\underline{42}, \underline{27} =$$

Reference string: 6 0 1 2 0 6 0 1 2 0 6 0 1 2 0 3 0 4 2 3 0 3 2 1 2 0 1 6 0 1

[6]	[6]	[6]	[2]	[2]	[2]	[2]	[2]	[2]	[2]	[2]
0	0	0	0	3	3	3	3	3	3	3
1	1	1	1	1	0	0	0	0	0	0
miss	M	M	M	H	M	m	H	H	H	H

### Replacement Algorithms

When the cache memory is full there is a need of a Replacement Algo.

To replace the existing cache block with new block.

FIFO

Consider 4 cache block, initially empty with following mm block Reference.

4 5 7 12 4 5 13 4 5 7

[4]	[4]	[4]	[4]	[4]	[4]	[4]	[13]	[13]	[13]	[13]
5	5	5	5	5	5	5	5	4	4	4
7	7	7	7	7	7	7	7	5	4	5
12	12	12	12	12	12	12	12	12	12	12
M	M	M	N	N	H	H	M	M	M	M

$$\text{Hit Ratio} = \frac{\text{Total Hits}}{\text{Total Ref}} = \frac{02}{10} = 0.2$$

$$\text{Miss Ratio} = 1 - \text{Hit Ratio} = 0.8$$

$$\frac{8}{10} = 0.8$$

page fault = miss

classmate

Date \_\_\_\_\_

Page \_\_\_\_\_

4 5 7 12 4 5 13 4 5 7

LRU Least Recently Used

4	4	4	4	4	4	4	4	4	4
5	5	5	5	5	5	5	5	5	5
7	7	7	7	7	7	13	13	13	13
12	12	12	12	12	12	12	12	12	7
M	M	M	M	H	H	M	H	H	M

$$\text{Hit Ratio} = \frac{7}{10} = 0.7$$

$$\text{Miss Ratio} = \frac{6}{10} = 0.6$$

LRU 0 2 1 6 4 0 0 1 0 3 1 2 1

Ex. 2

0	0	0	0	4	4	4	4	4	4	2	2
2	2	2	2	2	0	0	0	0	6	0	0
1	1	1	1	1	1	1	1	1	1	1	1
6	6	6	6	6	6	6	6	3	3	3	3
M	M	M	M	M	M	H	H	M	H	M	H

$$\text{Hit Ratio} = \frac{8}{12} = 0.67$$

$$\text{miss / page fault Ratio} = \frac{4}{12} = 0.33$$

discard least Recently accessed content

A LRU cache organizes items in order of use allowing you to quickly identify which item hasn't been used for the longest amount of time

EX 2  
FIFO

0 2 1 6 4 0 0 0 3 1 2 | 0 3 1 2 | 0 3 1 2 |

0	0	0	0	4	4	4	4	4	4	2	2
2	2	2	2	2	0	0	0	0	0	0	0
1	1	1	1	1	1	1	1	1	1	1	1
6	6	6	6	6	6	6	6	3	3	3	3
M	M	M	M	M	M	M	H	M	M	M	H

FIFO

7 0 1 2 0 3 0 4 2 3

cache memory has 4 lines

7	7	7	7	7	3	3	3	3	3	3	3
0	0	0	0	0	0	0	4	4	4	4	4
1	1	1	1	1	1	1	1	1	1	1	1
2	2	2	2	2	2	2	2	2	2	2	2
M	M	M	M	H	M	H	M	H	H	H	H

6 misses.

The block which has entered first in the main be replaced first.

This can lead to a problem known as "Belady's Anomaly" it starts that if we increase the no. of lines in cache memory the cache miss will increase

## LRU (Least Recently Used)

The page which was not used for the largest period of time in the past will get reported first. (looking backward in time)  
Let we have a sequence

e.g. 7 0 1 2 0 3 0 4 2 3

Cache memory has 3 lines

7	7	7	2	2	2	2	4	4	4
0	0	0	0	0	0	0	0	0	3
1	1	1	1	3	3	3	2	2	2
M	M	M	H	M	H	M	M	M	M

7	7	7	7	3	3	3	3	3	3
0	0	0	0	0	0	0	0	0	0
1	1	1	1	4	4	4	4	4	4
2	2	2	2	2	2	2	2	2	2

With  
4 cache  
lines

M	M	M	M	H	M	H	M	H	H
---	---	---	---	---	---	---	---	---	---

Optimal Approach: This The page which will not be used for the largest period of time in future reference will be replaced first.

⇒ This optimal Algorithm will provide the best performance but this is difficult to implement as it Requires the future Knowledge of pages which is not possible.

⇒ Used for comparison studies.

⇒ The use of this algo. guarantees the lowest possible page fault rate for a fixed rate of frames.

FIFO

When there is a need for page replacement, FIFO swaps out the page at the front of the queue that is the page which has been in the memory for longest.

Consider the page Reference String of size 12: 1, 2, 3, 4, 5, 1, 3, 1, 6, 3, 2, 3

with frame size (i.e. maximum 4 pages in a frame)

FIFO

Advantages: Simple & easy to implement  
Dis: Poor Performance

Doesn't consider the frequency of use or last used time, simply replaces the oldest page.

Suffers from Belady's Anomaly (i.e. more page faults when we increase the no. of page frames).

1	1	1	1	5	5	5	5	5	5	2	2
2	2	2	2	1	1	1	1	1	1	1	1
3	3	3	3	3	3	3	3	3	3	3	3
4	4	4	4	4	4	4	4	4	4	6	6
M	M	M	M	M	H	H	H	H	H	M	H

5	5	2	2
1	1	1	1
6	6	6	6
4	3	3	3

Total pg. fault = 09

LRU

1, 2, 3, 4, 5, 1, 3, 1, 6, 3, 2, 3

1	1	1	1	5	5	5	5	5	5	2	2
2	2	2	2	1	1	1	1	1	1	1	1
3	3	3	3	3	3	3	3	3	3	3	3
4	4	4	4	4	4	6	6	6	6	6	6
M	M	M	M	M	H	H	M	H	M	H	

Total page Fault = 8

In LRU whenever page Replacement happens, the page which has not been used for the longest amount of time is replaced.

Advantages : Efficient

Doesn't suffer from Belady's Anomaly.

DisAdvantages : Complex Implementation.

Expensive

Requires HW support.

**Optimal**

In this algorithm, pages are replaced which would not be used for the longest duration of time in the future i.e. the pages in the memory which are going to be referred farthest in the future are replaced.

T. miss  
6

Advantages : Easy to implement.

Simple Data structures are used.  
Highly efficient.

Disadvantages : Requires Future Knowledge of the program.  
→ Time consuming.

1, 2, 3, 4, 5, 1, 3, 1, 6, 3, 2, 3

1	1	1	1	+	1	1	1	6	6	6	6
2	2	2	2	2	2	2	2	2	2	2	2
3	3	3	3	3	3	3	3	3	3	3	3
4	5	5	5	5	5	5	5	5	5	5	5
m	m	m	m	H	H	H	m	H	H	H	H

Total page fault = 06

**Optimal**

4 5 7 12 4 5 13 4 5 7

4	4	4	4	4	+	4	4	4	4	4	4
5	5	5	5	5	5	5	5	5	5	5	5
7	7	7	7	7	7	7	7	7	7	7	7
12	12	12	13	13	13	13	13	13	13	13	13
m	m	m	H	H	m	H	H	H	H	H	H

Total page faults = 05

~~Ex APO~~ Consider page Reference string 1, 3, 0, 3, 5, 6, 3  
 with 3 page frames. Find number of page faults for FIFO replacement policies.

1	1	1	1	1	5	5	5
3	3	3	3	3	6	6	6
0	0	0	0	0	0	3	3
M	M	M	H	M	M	M	Total page fault = 6

LRV      1, 3, 0, 3, 5, 6, 3

1	1	1	8	1	8	5	8	5	5
0	0	3	3	3	3	3	3	3	1
M	M	M	H	O	O	O	6	6	6
M	M	M	H	M	M	M	H	Total page fault = 5	

optimal    1, 3, 0, 3, 5, 6, 3

1	1	1	1	1	5	5	5
0	3	3	3	3	3	3	3
+	8	0	0	0	6	6	6
M	M	M	H	M	M	H	

FIFO    1    1    1    1    1    1    1    1    6    6    6

4 frame    3    2    2    2    2    2    2    2    2    2    2

          3    3    3    3    3    3    3    3    3    3

          4    4    4    4    4    4    4    4    4

          5    5    5    5    5    5    5    5    5

          M    M    M    M    H    H    H    M    H

6, 0, 5, 2, 0

1, 2, 3, 4, 1, 2, 5, 1, 2, 3, 4, 5

3 frames	1	1	1	4	4	4	5	5	5	3	3	3
	2	2	2	1	1	1	1	1	1	4	4	
	3	3	3	2	2	2	2	2	2	2	5	
	m	m	m	m	m	m	H	H	m	m	m	

3 2 1 0 3 2 4 3 2 1 0 4

FIFO

3	3	3	0	0	0	4	4	4	4	4	4	4
2	2	2	3	3	3	3	3	3	3	3	3	3
1	1	1	1	2	2	2	2	2	2	2	2	2
m	m	m	m	m	m	m	H	H	m	m	m	m
									4	4	4	4
									1	1	1	1
									2	0	0	0
									m	m	H	

3	3	3	3	3	3	4	4	4	4	0	0	0
2	2	2	2	2	2	3	3	3	3	3	3	4
1	1	1	1	1	1	1	2	2	2	2	2	2
0	0	0	0	0	0	0	0	0	1	1	1	1
m	m	m	m	H	H	m	m	m	m	m	m	m

Belady's Anomaly

A cache's write policy is the behavior of a cache while performing a write operation.

Hif : write thru, back

Miss : write-allocate, no write allocate

A device for storing digital information that is fabricated by using integrated circuit technology is known as semiconductor memory.

Thus semiconductor devices are preferred as primary memory.

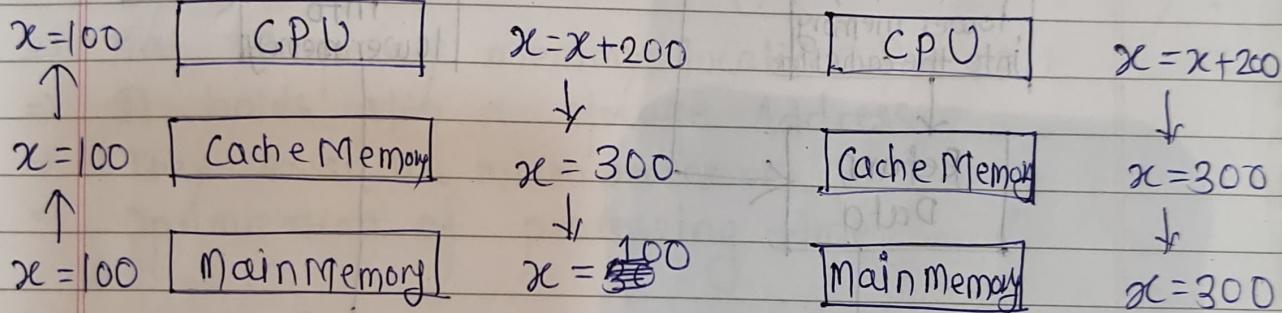
These are 2 main types of SM technology.

non-volatile (ROM)

Volatile (RAM) →

Cache write Policies:

Before      Write Back      After      Write Through      After



- Data is updated only in cache memory
  - simple method to implement.
  - When a cache block is not required then selected cache block is written back to main memory to UPDATE the main memory.
  - commonly used method
  - UnReliable Method
  - No Data Redundancy
  - No wastage of time as data occurs only in cache memory.
- Data is updated both in Cache Memory & main Memory
  - complex method to implement.
  - When cache block is not required there is no need to written back to main memory.
  - Unpopular Method
  - Reliable Method
  - Data Redundancy
  - wastage of time as each data updated both in cache & main memory

