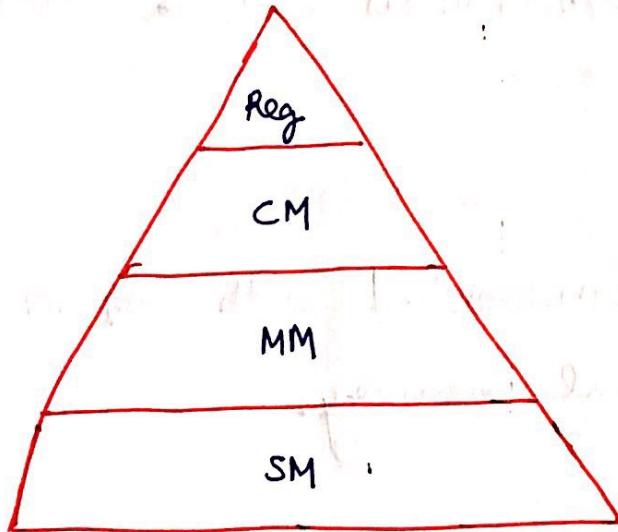


Memory hierarchy design :-

According to mem hierarchy design, system supported memory standards is as follows :-

- Level	: 1	2	3	4
- Name	: Register	cache memory	main memory	secondary memory
- Typical - size	: < 1KB	< 16MB	< 16GB	> 100GB
- Implementation:	: customized multiport CMOS	SRAM	DRAM	Magnetic
- Access time : (ns)	: (0.25 - 0.5)	(0.5-25)	(80-250)	(50,00,000)
- Bandwidth : (MB/sec)	: (20,000 - 1,00,000)	(5000, - 10,000)	(1000 - 5000)	(20 - 150)
- Managed by :	: Compiler	H/w	OS	OS
- Backed by :	: CM	MM	SM	Compact Disk (CD)

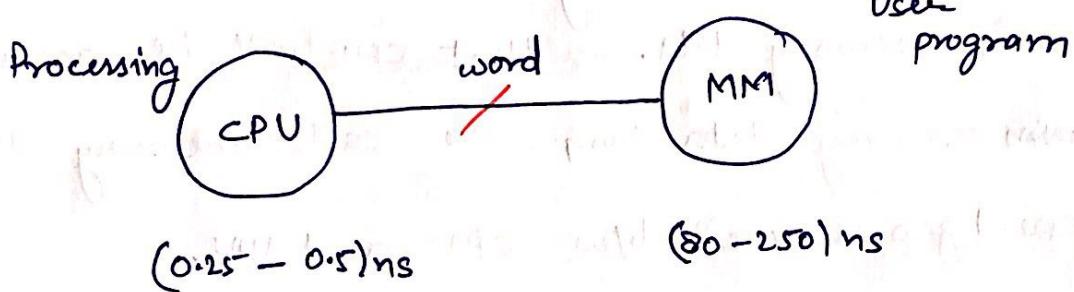


Memory hierarchy design

Bottom - top approach of a memory hierarchy states:

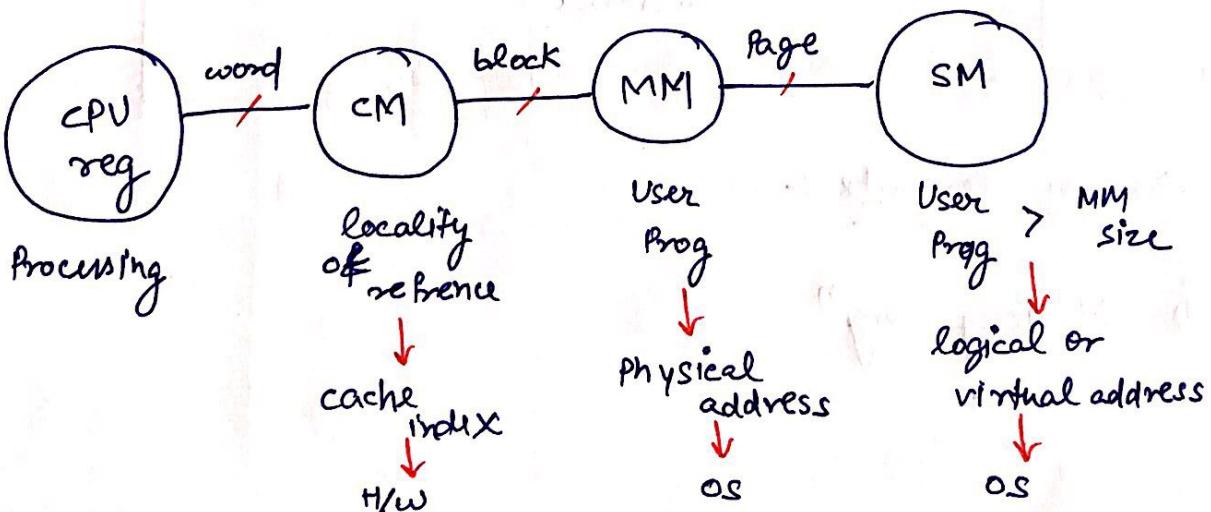
- Capacity \downarrow
- Access time \downarrow
- Performance \uparrow
- Cost/bit \uparrow
- Expensive

Sys without hierarchy design :-



Speed gap \uparrow b/w CPU and MM

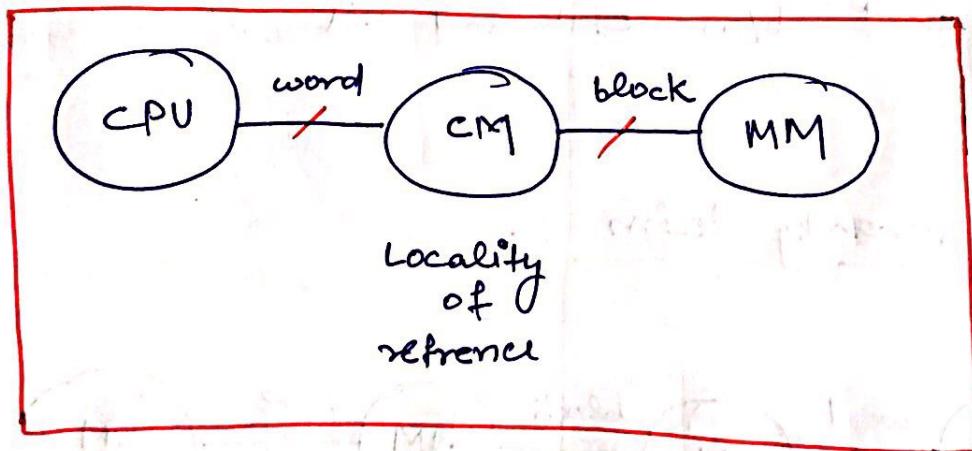
Sys with hierarchy design :-



- Main memory speed is synchronized with a CPU speed using cache memory.
- Main memory space is synchronized with any application program size using virtual memory.

Cache memory Design:-

- Cache memory acts as an intermediate memory b/w the CPU and the main memory, used to hold the image of a main memory data. So that CPU will be accessing the main memory data from the cache memory there - fore speed gap is syncⁿ b/w CPU and MM.



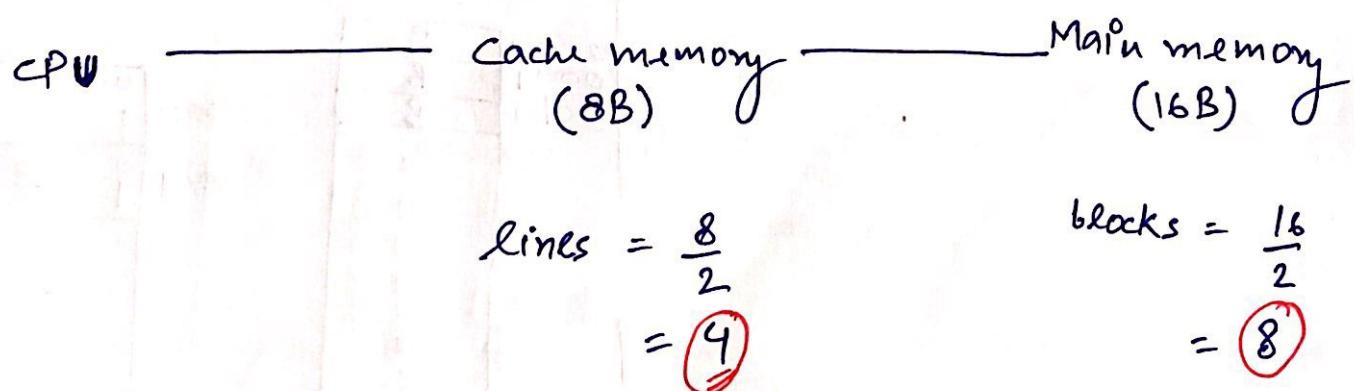
Design elements :-

- i) Memory orgⁿ

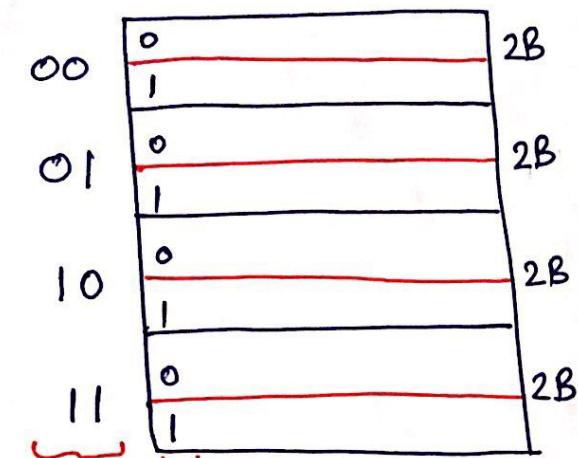
- (ii) Mapping techniques
- (iii) Replacement organization
- (iv) Updating techniques
- (v) Multi-level caches.

(i) Memory organization :-

- According to memory hierarchy design, data will be copied from the MM to CM in a form of "Blocks" so both of the memory are organized into blocks based on the block-size.
- No. of block in MM = MM size / Block-size.
- No. of lines in CM = CM size / Block-size.
- Let us consider the system which contains 8-byte cache memory & 16-byte main memory with 2B-block



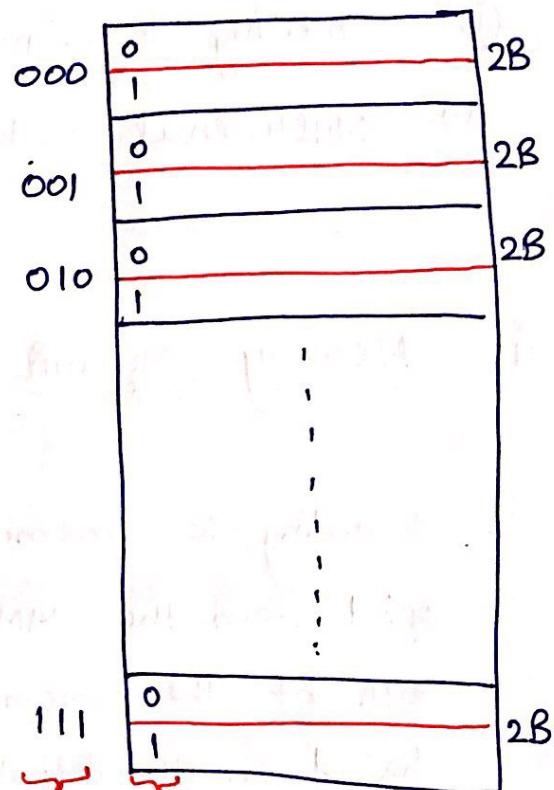
cache memory (8B)



line offset (LO)

word - offset (WO)

Main memory (16B)

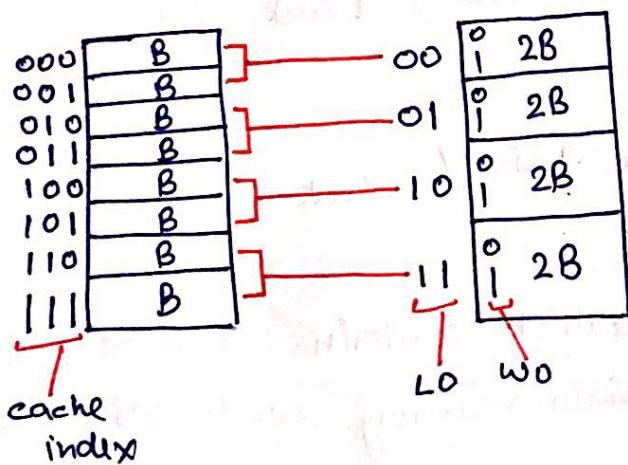


tag - offset (TO)

word - offset (WO)

Before orgⁿ

After orgⁿ



cache index

LO

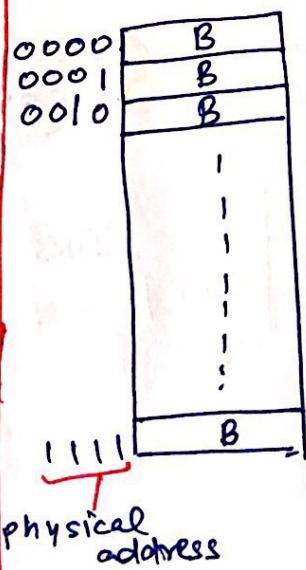
WO

Before orgⁿ

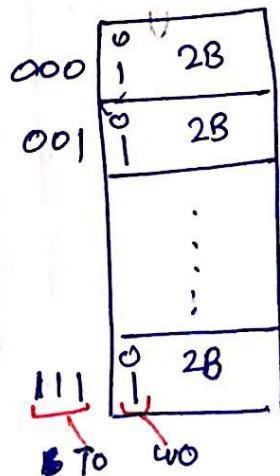
After orgⁿ

16B → MM

$$\text{Blocks} = \frac{16}{2} = 8$$

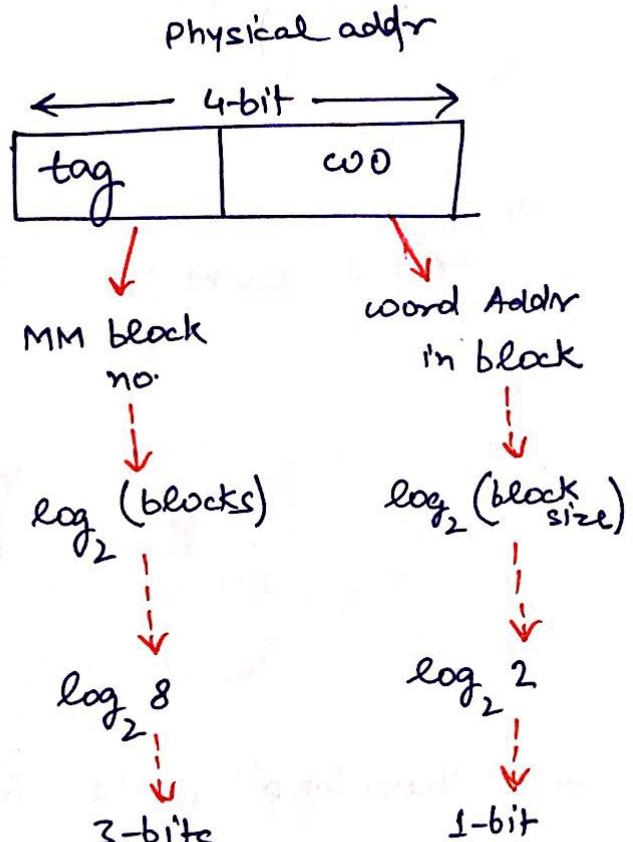
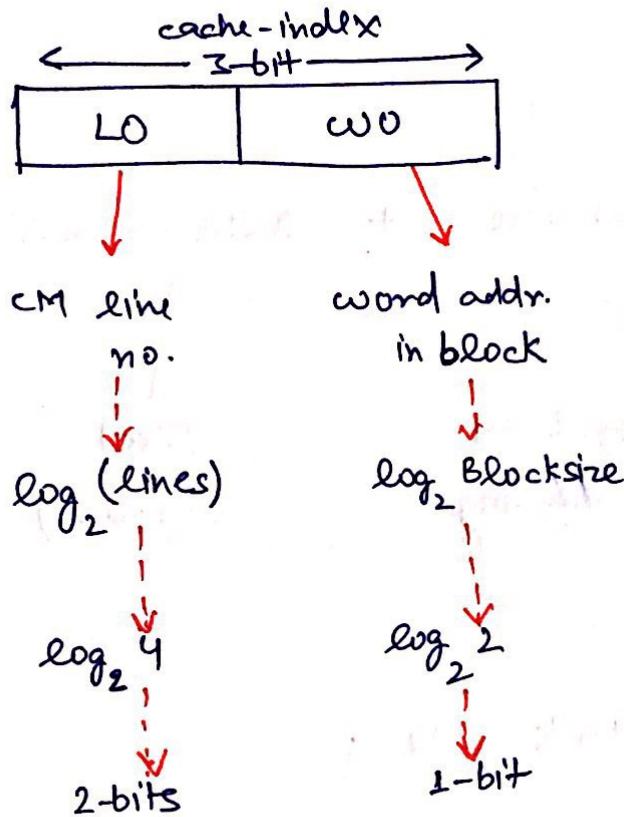


physical address



TO
WO

Address Formats :-



(ii) Mapping Techniques :-

- The process of copying the data from the main memory to cache memory is called as mapping.
 - In the mapping process, data-block is copied into a cache memory along with the physical address.
 - To hold the physical address, storage space is required in cache line known as tag space. Therefore,
- (Tag memory-size = lines in CM * tag space in the line)

$$(\text{Data-memory size} = \text{lines in CM} * \text{Block-size})$$

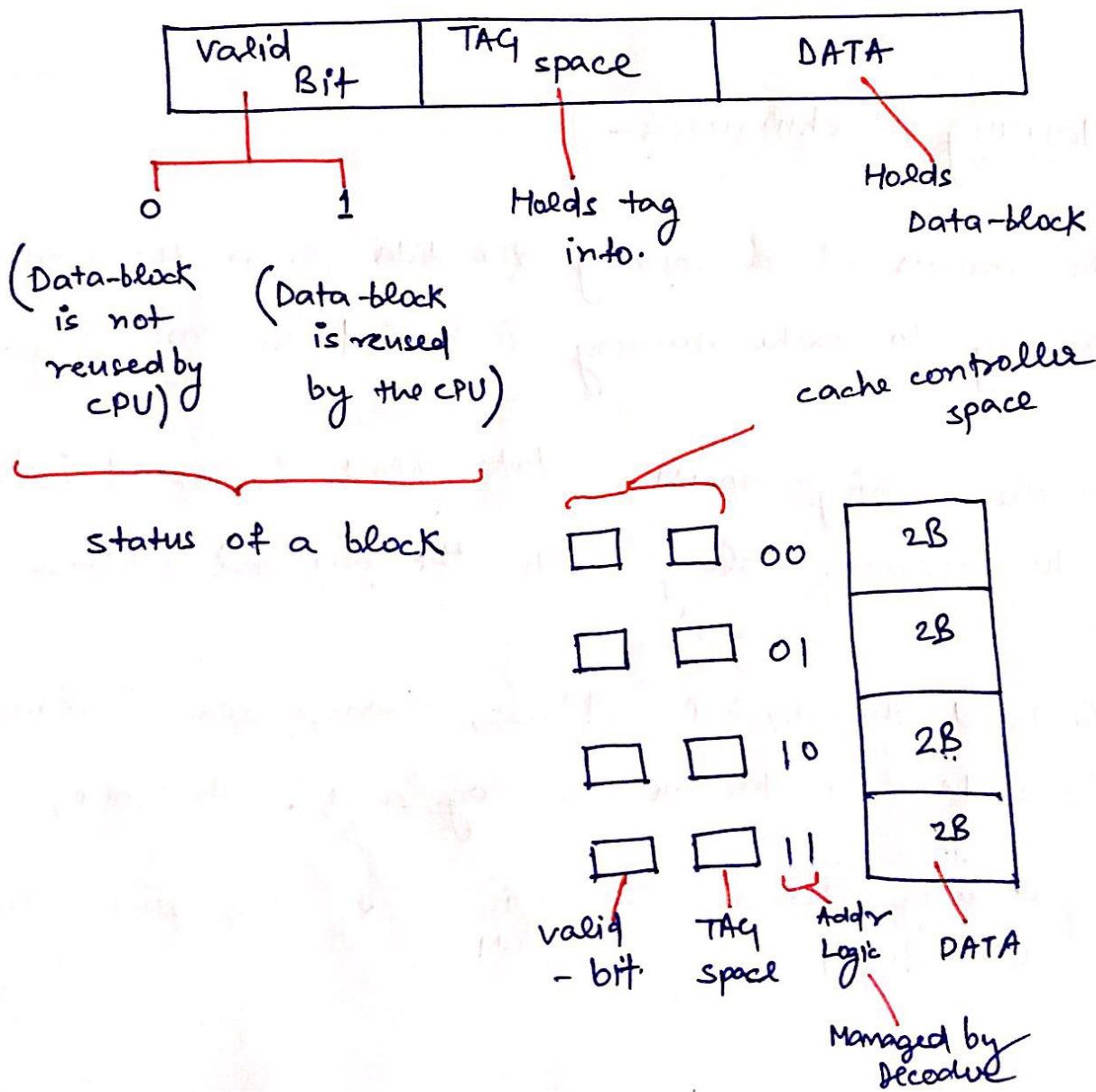
So,

$$\boxed{\text{Total cache} = \text{Tag mem size} + \text{Data-mem size}}$$

|
varies depends on type of mapping

|
Fixed (eg. 8B)

- Therefore cache line contents are :



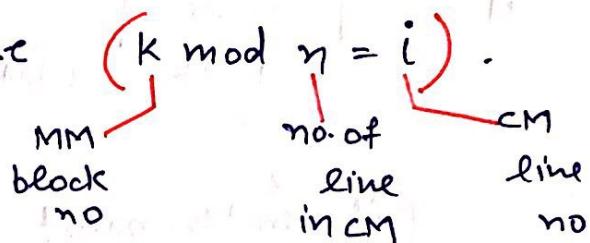
- Three kinds of mapping techniques used in cache design:-

(i) Direct mapping

(ii) Associative mapping

(iii) Set-Associative mapping

⇒ Direct - Mapping :- In this technique 'mod' func is used to map the data i.e $(k \bmod n = i)$.



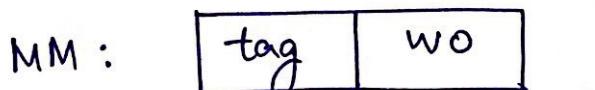
$0 \bmod 4 = 0$
$1 \bmod 4 = 1$
$2 \bmod 4 = 2$
$3 \bmod 4 = 3$
$4 \bmod 4 = 0$
$5 \bmod 4 = 1$
$6 \bmod 4 = 2$
$7 \bmod 4 = 3$

0	0, m, 2m
1	1, m+1, 2m+1
⋮	⋮
$(m-1)$	$(m-1), 2m-1, 3m-1$

(row = line-status)

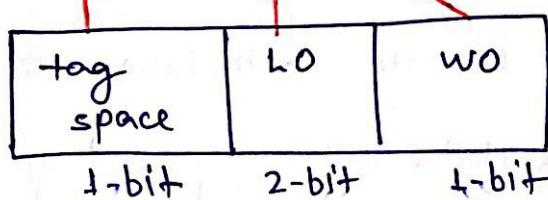
(col = CM - status)

- Mapping funcⁿ shows the relationship b/w the block no. in MM and line no. in CM. Therefore address binding is :



3-bit 1-bit

Direct :
CM



1-bit 2-bit 4-bit

(if denotes no. of
blocks assigned to given
line possible single at a time)

$$2^1 = 2 \text{ blocks possible}$$

i.e. $(0, 4) \rightarrow \text{line } 0$

Tag Directory size = lines * tag-size in line

$$= 4 * 1\text{-bit}$$

$$= \underline{\text{4-bits}}$$

$$\begin{aligned} \text{Total cache} &= \text{Tag cache} + \text{Data cache} \\ &= 4\text{-bits} + 8\text{-bytes} \\ &= \underline{68\text{-bits}} \end{aligned}$$

MM

Direct CM

Tag



Tag	LO
-----	----

$B_0(000)$

$B_4(100)$

$B_1(001)$

$B_5(101)$

$B_2(010)$

$B_6(110)$

$B_3(011)$

$B_7(111)$



0/1	00
-----	----

B_0 is present

B_4 is present

line 0

0/1	01
-----	----

line 1

0/1	10
-----	----

line 2

0/1	11
-----	----

line 3

Runtime analysis :-

Assume CPU generates MM ref : $B_0, B_7, B_0, B_0, B_7, B_4, B_0, B_4, B_0, B_4$.

Initially CM is empty

Accessing :-

(mapping)

MM - Block

$$K \bmod \eta = i$$

CM Line

opⁿ

B₀ — Miss

$$0 \bmod 4 = 0$$

line 0

compulsory
-miss

B₇ — Miss

$$7 \bmod 4 = 3$$

line 3

compulsory
-miss

B₀ — Hit

—

— (000) hit

B₀ — Hit

—

— (000) hit

B₇ — Hit

—

— (000) hit

B₄ — ~~Hit~~ Miss

$$4 \bmod 4 = 0$$

line 0

compulsory
miss
+ replace

B₀ — Miss

$$0 \bmod 4 = 0$$

line 0

conflict
+ miss
+ replace

B₄ — Miss

$$4 \bmod 4 = 0$$

line 0

"

B₀ — Miss

$$0 \bmod 4 = 0$$

line 0

"

B₄ — Miss

$$4 \bmod 4 = 0$$

line 0

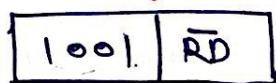
"

$$\text{Hit - ratio} = \frac{\text{Total Hit}}{\text{Total access}} = \frac{3}{10} = (0.3)$$

Access Sequence :-

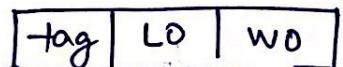
M/C Instⁿ : MOV r₀, [1001]

CPU generates mem req

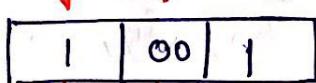


Direct CM

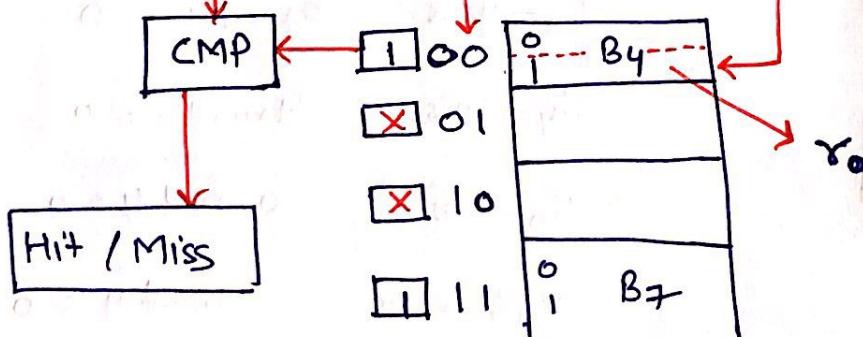
Addr Format



1-bit 2-bit 1-bit



Addr Decoder



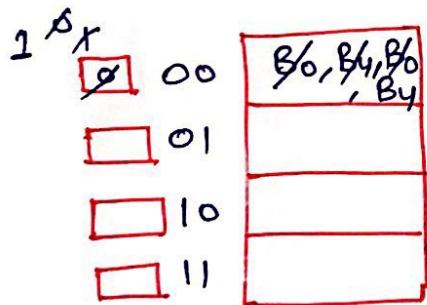
Limitations :-

- i) In the direct cache design, each line is capable to hold only 1-block (tag) at a time therefore no. of conflict misses will be increases.
- ii) When the CPU frequently prefers the multiple blocks which are mapped into a same cache line. Then the cache blocks are continuously swapping. So, hit ratio will be falling down, this phenomenon is known as 'Thrashing'.

= Thrashing -

Assume MM Block ref: B_0, B_4, B_0, B_4

C_M : empty (initially)



B_0 - miss $0 \bmod 4 = 0$

B_4 - miss $4 \bmod 4 = 0$

B_0 - miss $0 \bmod 4 = 0$

B_4 - miss $4 \bmod 4 = 0$

$$\text{Hit ratio} = \frac{0}{4} = 0$$

$H = 0 \rightarrow \text{Thrashing}$

- For preventing from Thrashing alternative mapping is required i.e. Associative mapping.

⇒ Associative mapping :- It is designed without address called as content - addressable memory.

- In this cache design, sequence func' (counter logic) is used to map the data so any MM block can be mapped into any CM line in a sequence. Therefore address binding is :-

MM Addr :

tag	wo
3-bits	1-bit

Associative CM :

tag	wo
3-bits	1-bit
(tag-space)	

$$\text{Tag Directory size} = \frac{\text{lines}}{\text{in CM}} * \frac{\text{tag space}}{\text{in the line}}$$

$$= 4 * 3 \text{ bits}$$

$$= \underline{12 \text{- bits}}$$

$$\begin{aligned} \text{so, Total cache} &= \text{Tag cache} + \text{Data cache} \\ &= 12 \text{- bits} + 8 \text{ Bytes} \\ &= \underline{76 \text{ bits}} \end{aligned}$$

Assume /
MM Block ref: $B_0, B_4, B_0, B_4, B_3, B_7, B_3, B_7, B_0, B_2$

cM: empty

MM Block ref

seq - funcⁿ

B_0 - Miss

Any line in a sequence

B_4 - Miss

B_0 - Hit

B_4 - Hit

B_3 - Miss

B_7 - Miss

B_3 - Hit

B_7 - Hit

B_0 - Hit

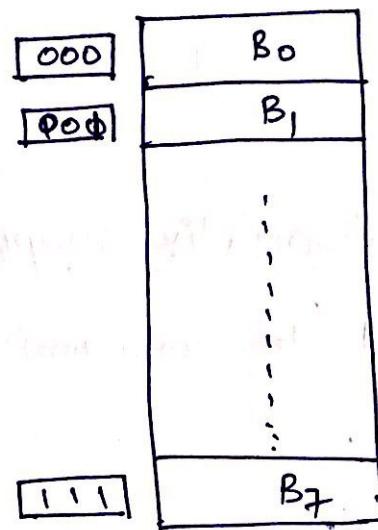
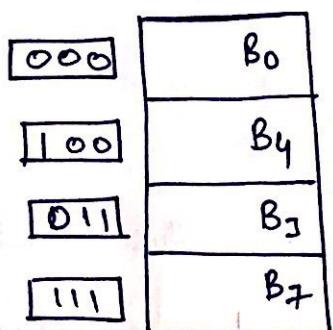
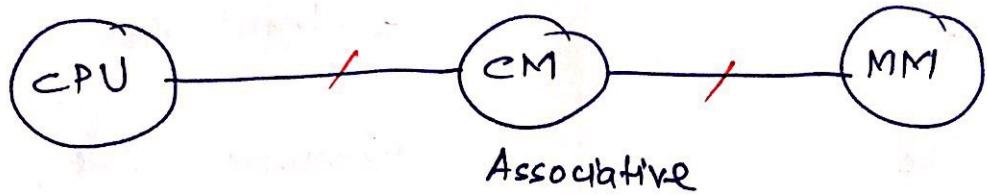
B_2 - Miss



replacement policy → [FIFO
LRU]

* FIFO & LRU policies are required to replace the data when the cache is FULL.

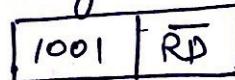
$$(\text{Hit-Ratio} = \frac{5}{10} = 0.5)$$



Accessing Sequence -

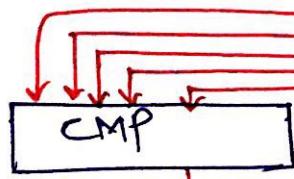
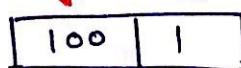
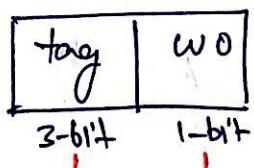
M/C-instrⁿ : MOV r₀, [100]

Memory Request

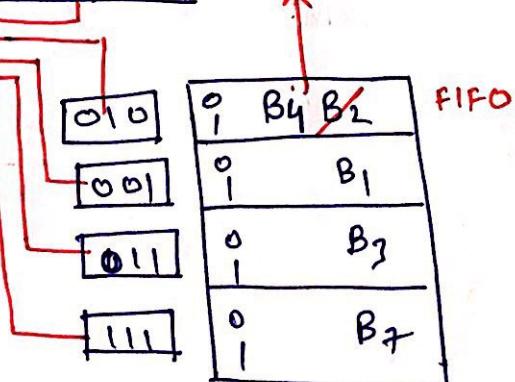


Assoc. CM

Adddr Format:



Hit/Miss
(Miss)



Advantage:-

Fastest cache

Disadvantage:-

Expensive cache, required alternative i.e Set-Associative.

⇒ Set - Associative mapping / cache :- This cache structure is used to compromise the direct and associative cache designs.

- In this cache design, lines are grouped into sets to accommodate multiple blocks in the set.

- Number of sets , $(S = \frac{N}{P\text{-way}})$

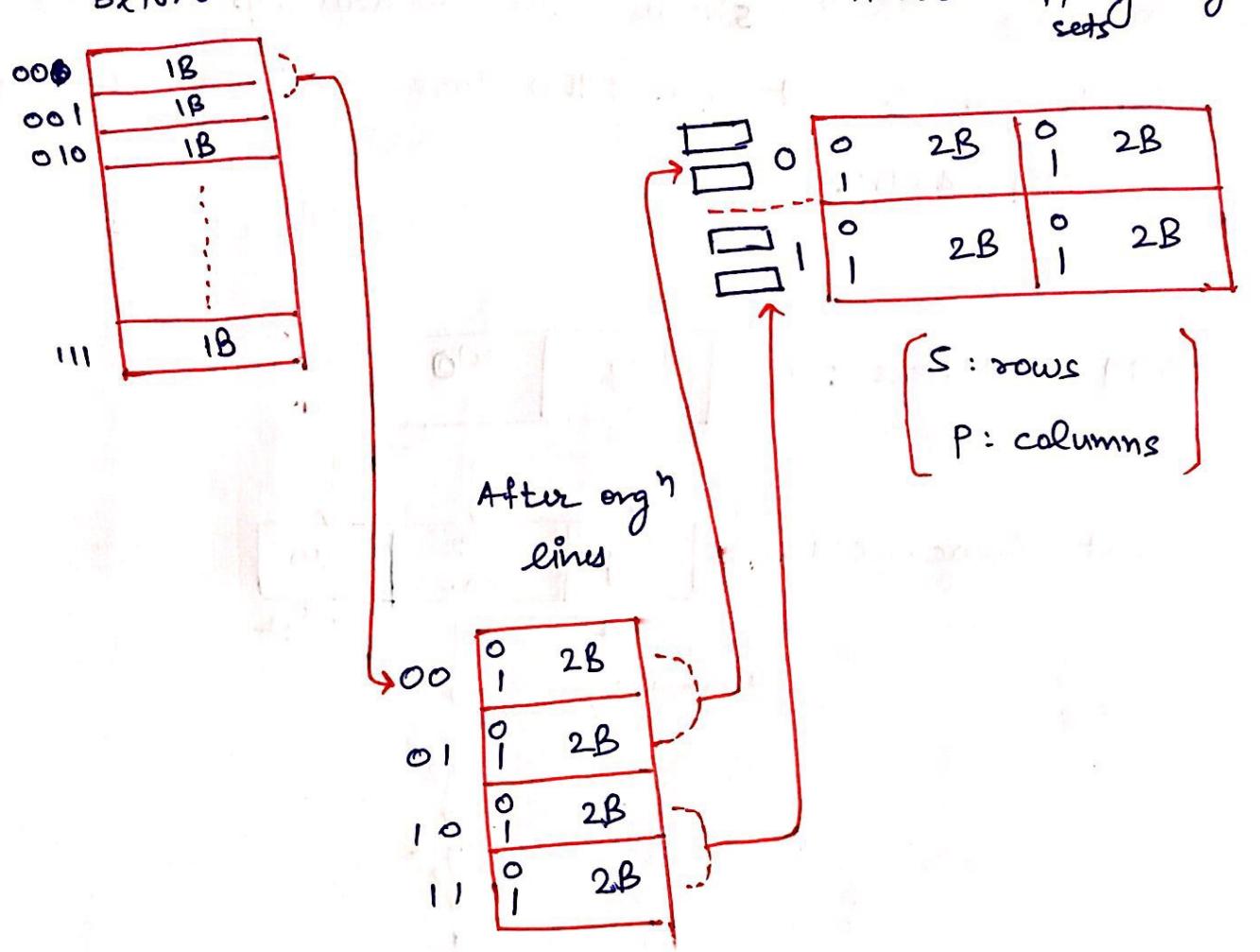
where N is lines in CM

P is line in a set

Eg:- consider 2-way set associative cache of 8-byte with 2-byte block

$$\text{lines} = \frac{8}{2} = 4$$

$$\text{sets} = \frac{N}{P} = \frac{4}{2} = 2$$



- In this design 'mod' funcn is used to map the data

i.e $(K \bmod S = i)$.

$\begin{matrix} \text{MM} \\ \text{Block} \\ \text{no} \end{matrix}$ $\begin{matrix} \text{sets} \\ \text{in} \\ \text{CM} \end{matrix}$ $\begin{matrix} \text{CM} \\ \text{set} \\ \text{number} \end{matrix}$

$$0 \bmod 2 = 0 \quad s_1 \quad 6 \bmod 2 = 0 \quad s_1$$

$$1 \bmod 2 = 1 \quad s_2 \quad 7 \bmod 2 = 1 \quad s_2$$

$$2 \bmod 2 = 0 \quad s_1$$

$$3 \bmod 2 = 1 \quad s_2$$

$$4 \bmod 2 = 0 \quad s_1$$

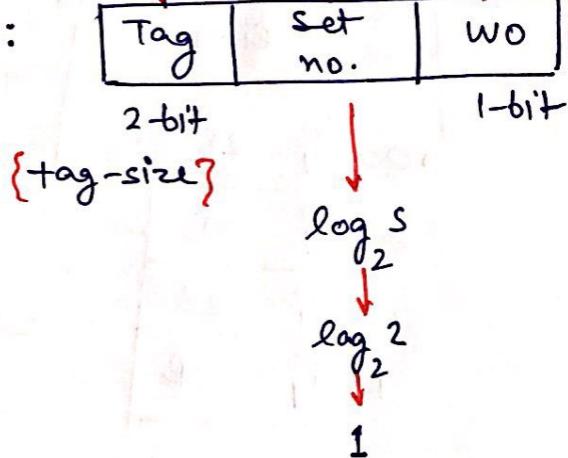
$$5 \bmod 2 = 1 \quad s_2$$

- mapping func' shows the relationship b/w the block no. & set no. Therefore address binding is MM Addr.

MM Address :



Set - Assoc. CM :



Tag - Directory size = $\underset{\text{in CM}}{\text{lines}} * \underset{\text{in the line}}{\text{tag-space}}$

$$= N * 2\text{-bits}$$

$$= (P * S) * 2$$

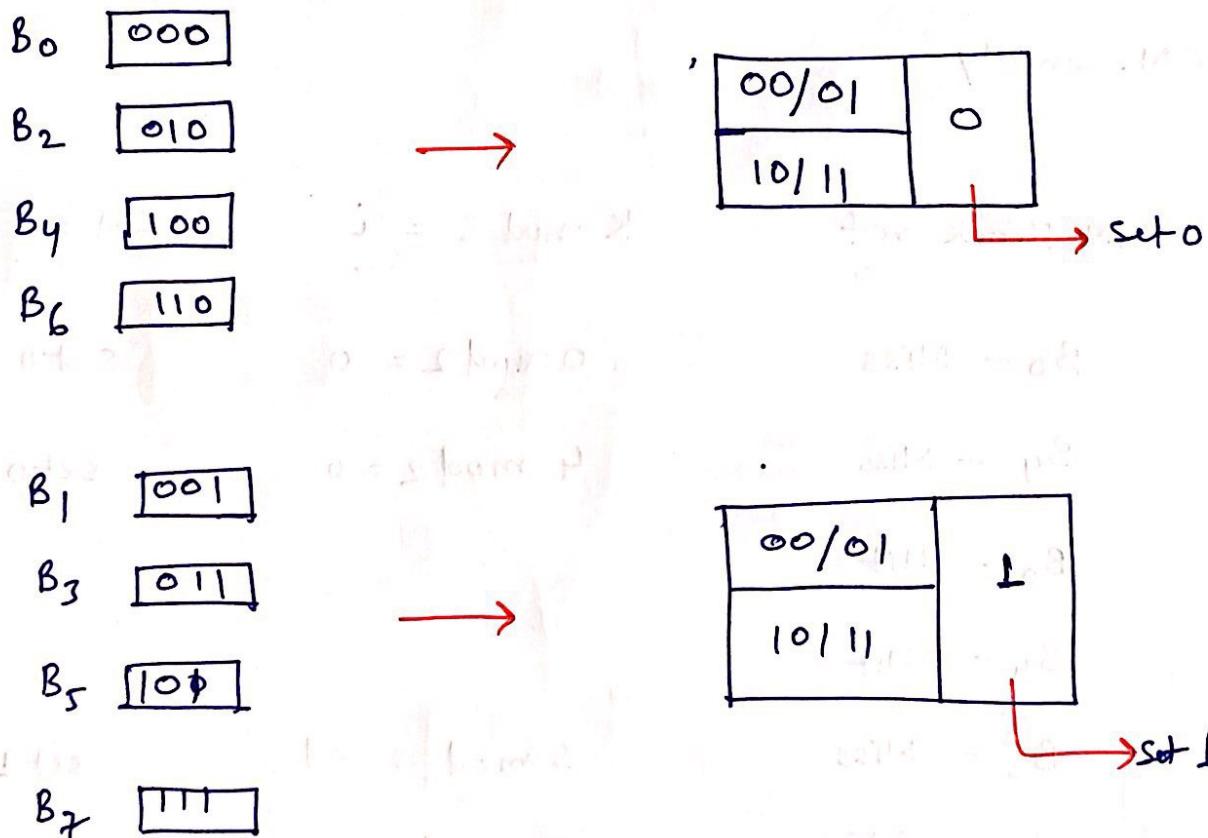
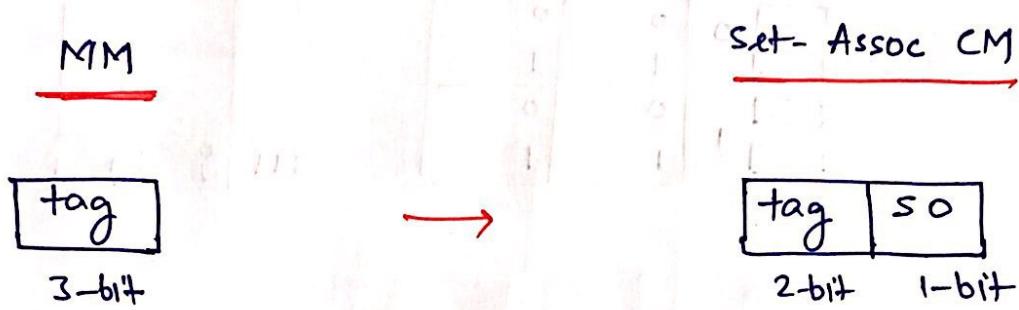
$$= (2 \times 2) \times 2$$

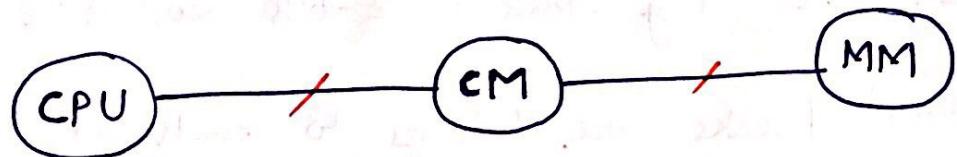
$$= \underline{\underline{8\text{-bits}}}$$

Type	Tag - Directory size
Direct	4-bits
Associative	12-bits
Set - Associative	8-bits

NOTE:- In this tag-space = 2-bit's so, 2^n i.e $2^2 = 4$

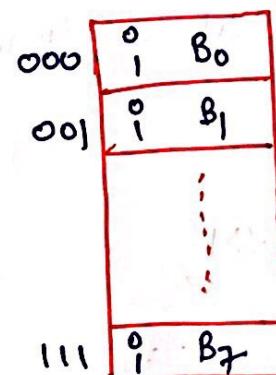
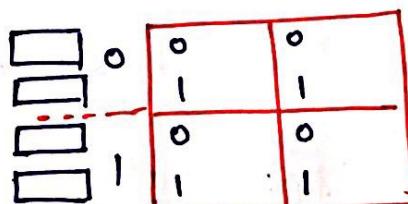
different blocks are belongs to each set but only 2 at a time.





2-way set

associative
CM



Assume,

MM Block ref: B0, B4, B0, B4, B3, B7, B3, B7, B0, B2

CM: empty

MM Block ref

$k \bmod s = i$

CM set

B0 - Miss

$$0 \bmod 2 = 0$$

set 0

B4 - Miss

$$4 \bmod 2 = 0$$

set 0

B0 - Hit

B4 - Hit

B3 - Miss

$$3 \bmod 2 = 1$$

set 1

B7 - Miss

$$7 \bmod 2 = 1$$

set 1

B3 - Hit

B7 - Hit

B0 - Hit

B4 - Miss

$$2 \bmod 2 = 0$$

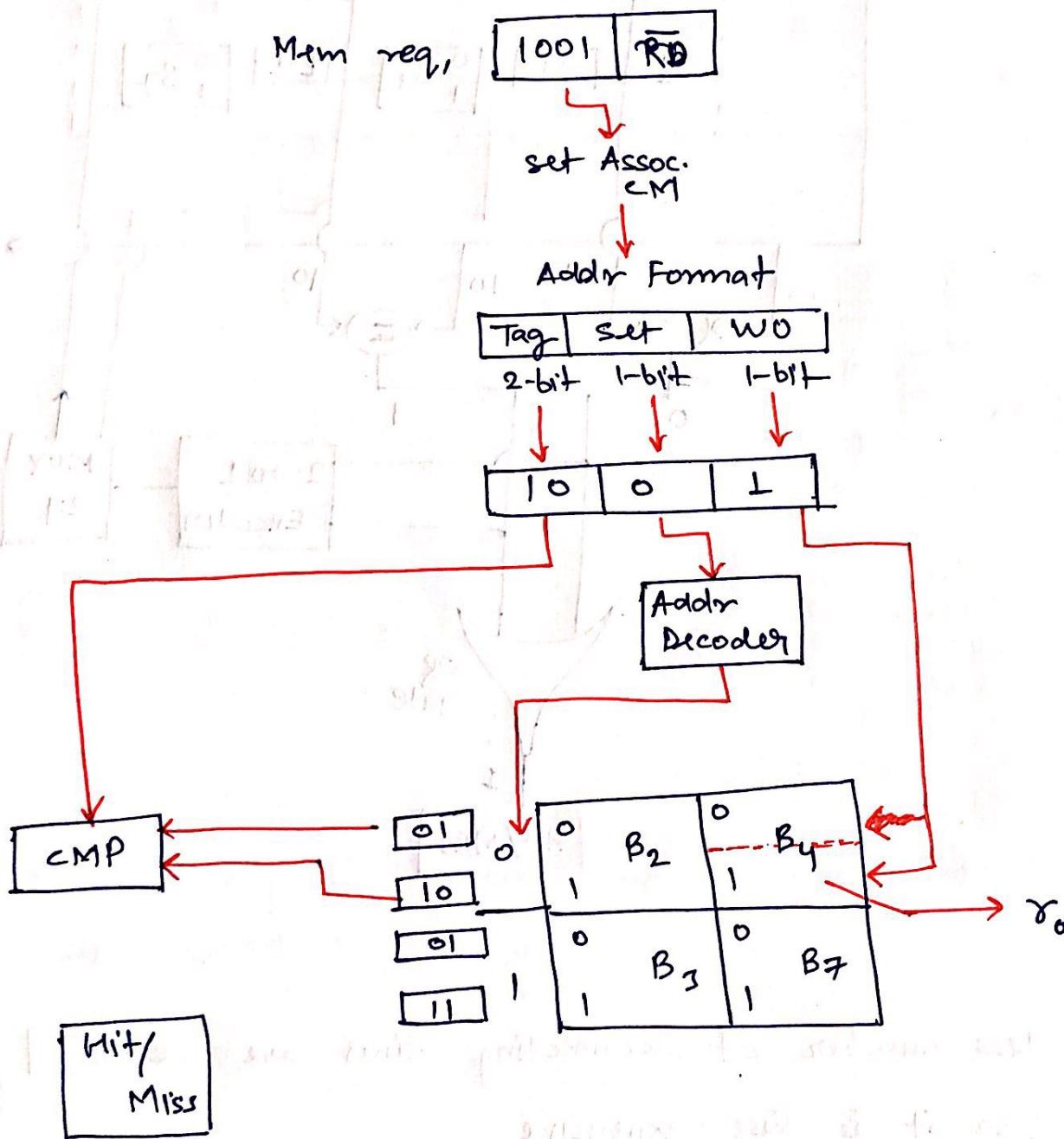
set 0 - FIFO: B0
LRU: B4
B2

$$(\text{Hit-Ratio } H = \frac{5}{10} = 0.5)$$

NOTE:- FIFO & LRU policies are used to replace the data when the set is full.

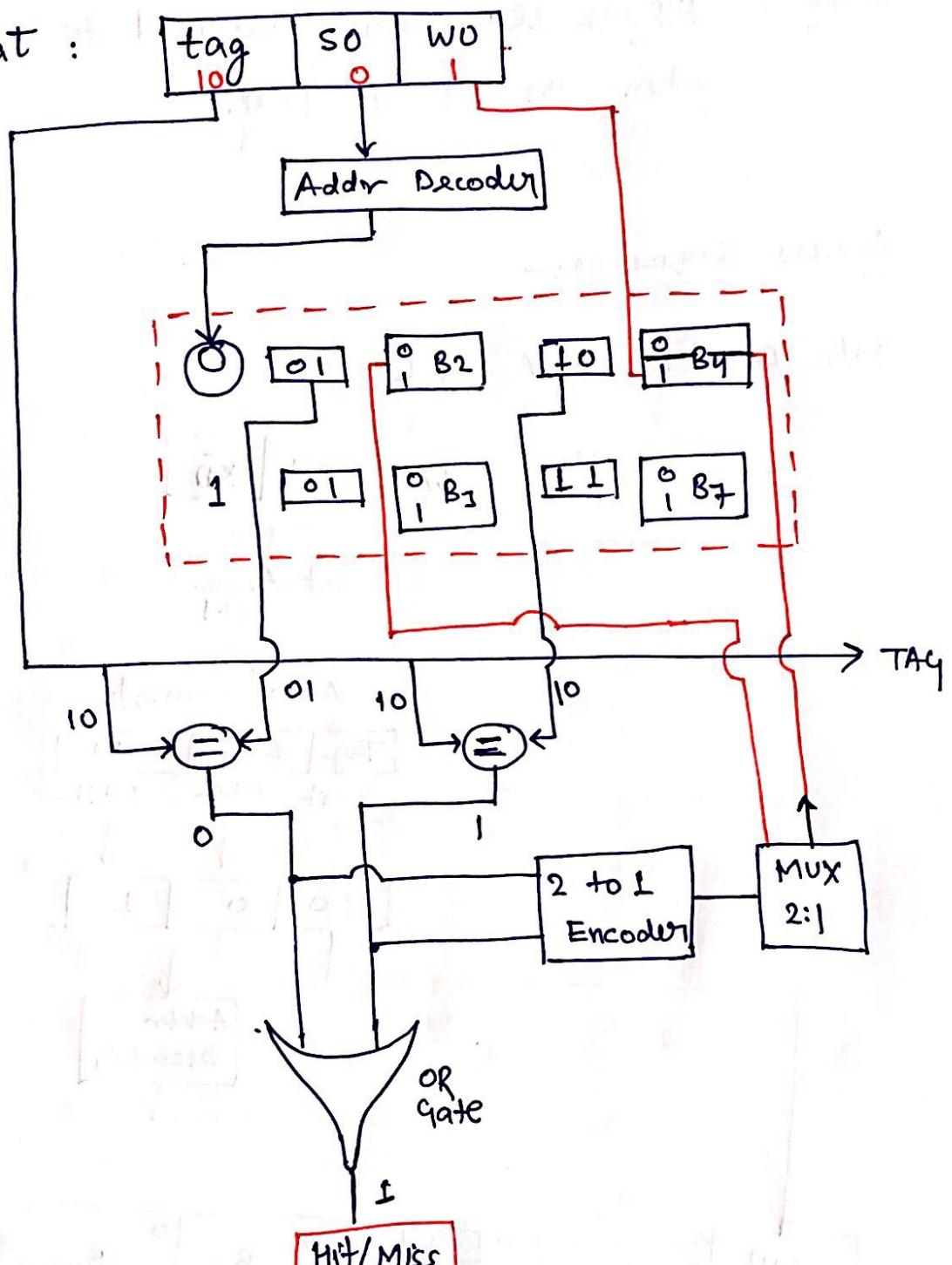
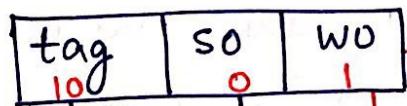
Access Sequence:-

M/c Instⁿ: MOV m, [1001]



~~composition~~ ~~composition~~ circuit :-

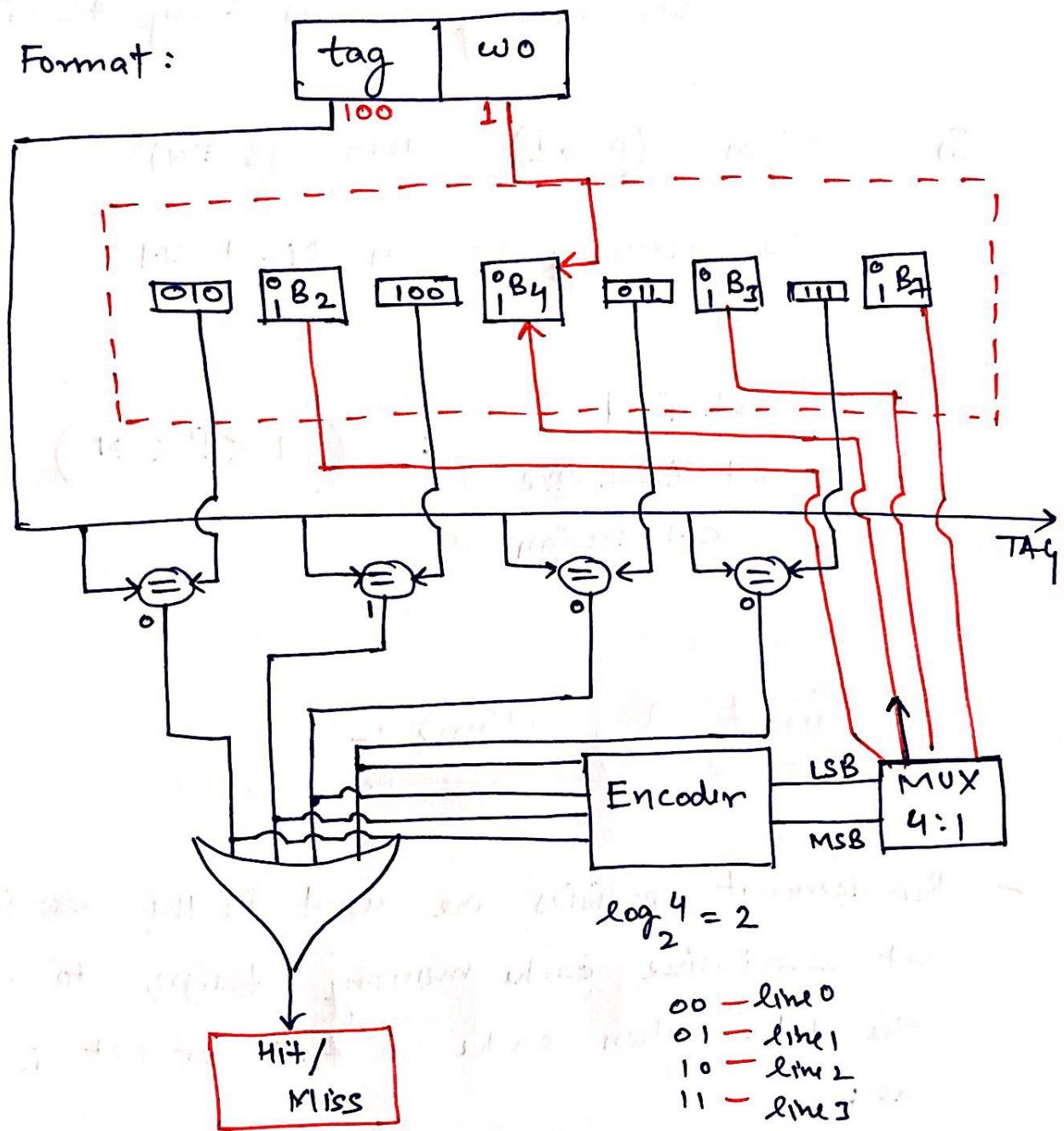
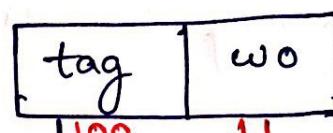
Addr Format :



* Less number of connecting lines are required
so it is less expensive

Comparision ckt of Associative EM :-

Addr Format:



- * More no. of CMP ckt and more no. data & control lines are required so it is expensive.

NOTE :- (i) When $(P \neq N)$ then $(S \approx 1)$
so, memory becomes fully Associative CM

(ii) when $(P \rightarrow 1)$ then $(S \approx N)$

so, memory becomes Direct CM

so, optimized set Associative CM Design : $(1 < P < N)$

Replacement Algorithms :-

- Replacement policies are used in the associative & set associative cache memory design, to replace the data when cache is full or set is full name as :
 - (i) FIFO
 - (ii) LRU
- In FIFO, replace the block which is ~~present~~ enters in the cache longest time without reference. firstly.
- In LRU, replace the block which is present in the cache longest time without reference.

Q. Consider 32KB direct cache organized into a 256 B blocks. Cache memory data is subset of a 2^{28} B storage space. In the cache controller, each tag comprising of 1 valid bit and 1 update bit. What is the tag directory in the cache controller?

$$CM \text{ size} = 32 \text{ KB}$$

Direct CM

$$N (\text{no. of lines in CM}) = \frac{CM \text{ size}}{\text{Block-size}}$$

$$\text{Block-size} = 256 \text{ B}$$

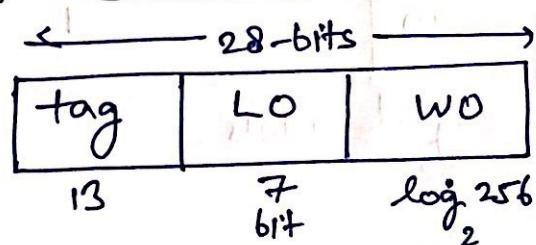
$$MM \text{ size} = 2^{28} \text{ B}$$

$$N = \frac{32 \text{ KB}}{256 \text{ B}}$$

$$\text{Physical addr size} = 28\text{-bits}$$

$$N = 2^7$$

So, Addr Format



$$\text{Tag directory size} = N * \text{tag-space}$$

$$= 2^7 (13 + 1 + 1)$$

$$= \boxed{2^7 * 15 \text{ bits}}$$

Q: Consider 64-bit hypothetical CPU which supports 8KB direct mapped cache org'n into 64 words blocks. To which cache line (in decimal) the following MM data is mapped. 0x FC3D6D

$$CM \text{ size} = 8KB = 2^{13} B$$

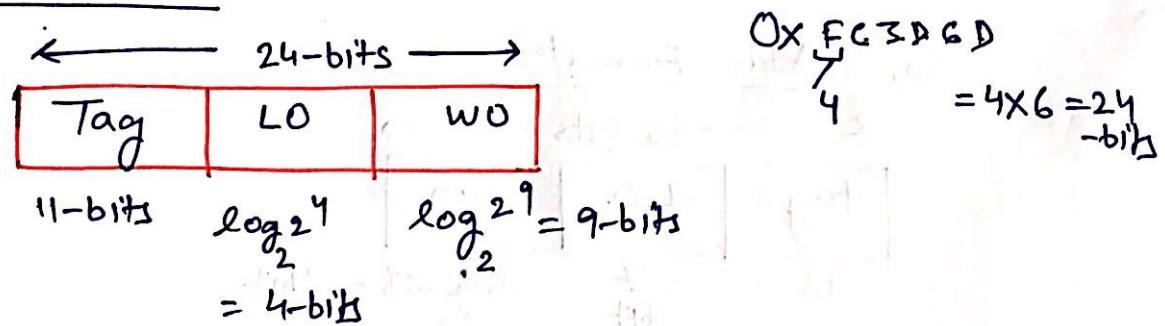
Direct mapping -

$$\text{Block size} = 64 \text{ w}$$

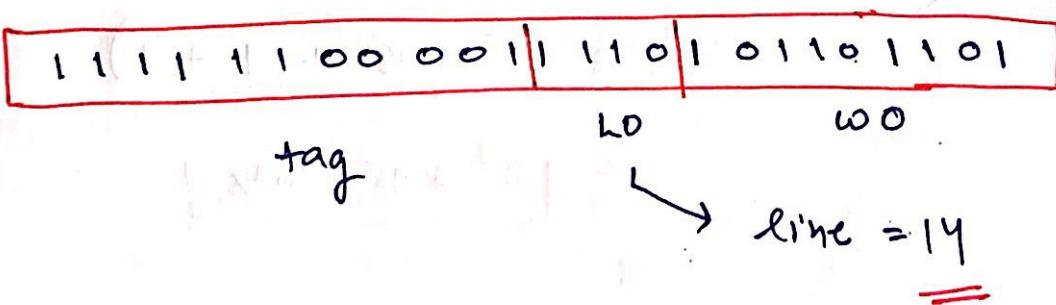
$$= 64 \times 8B = 2^9 B$$

$$\text{No. of lines (N)} = \frac{2^{13}}{2^9} = 2^4 \text{ lines}$$

Addr-format:-



FC3D6D



Q. Consider 32-bit CPU supports 8-way set assoc. cache of 64KW memory, organized into a 256 B blocks. CPU generates 32-bit physical address to access the data. In the cache controller each tag is consist of 1 valid bit, 1 update bit, 2-replace bits. What is the size of tag directory.

$$1W = 2 \text{ 32-bit} = 4B$$

$$\# \text{ set} = \frac{\text{set-size (P)}}{\text{set}} = 8 \quad (\text{set - associative mapping})$$

$$CM \text{ size} = 64 \text{ KW}$$

$$= 64K \times 4B$$

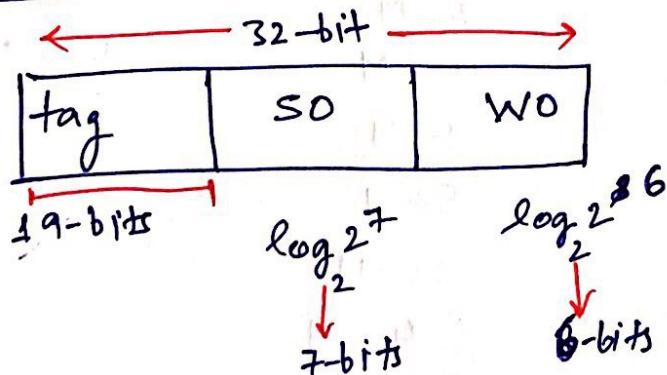
$$= 2^{18} B$$

$$\text{Block size} = 256 B = \underline{2^8 B} \quad \underline{2^6 W}$$

$$\text{No. of lines (N)} = \frac{2^{18}}{2^8} = 2^{10} \text{ lines}$$

$$\text{No. of set (S)} = \frac{N}{P} = \frac{2^{10}}{2^3} = 2^7$$

Addr Format:-



$$\begin{aligned}
 \text{Tag-directory size} &= N * (\text{tag-space} + 1 \text{ valid bit} \\
 &\quad + 1 \text{ update bit} + \\
 &\quad 2 \text{ replace bits}) \\
 &= N * (19 + 1 + 1 + 2) \\
 &= 2^{10} * 23 \text{ bits} \\
 &= \underline{\underline{23K \text{ bits}}}
 \end{aligned}$$

Q. Consider 4 block cache memory, initially empty with following main memory references 4, 5, 7, 12, 4, 5, 13, 4, 5, 7. Identify the hit ratio using -

- | | |
|----------|------------------------------------|
| (a) FIFO | (c) Direct Mapped CM |
| (b) LRU | (d) 2-way set Association with LRU |

(a) FIFO sequence -

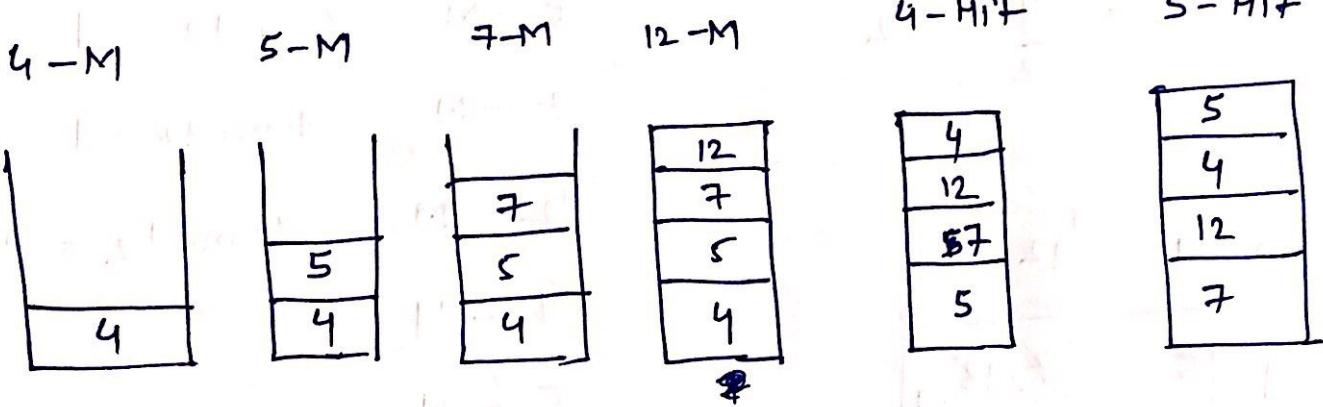
4	13
5	4
7	5
12	7

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

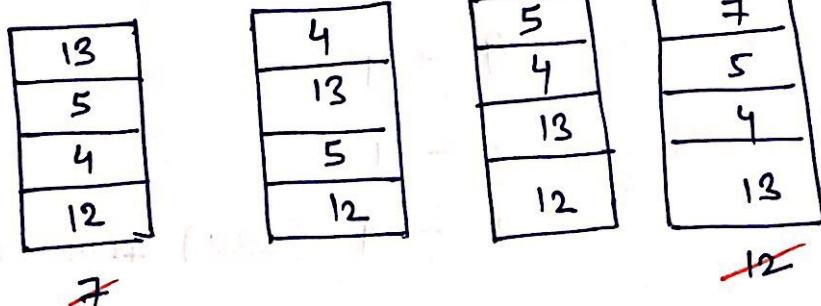
$$H = \frac{2}{10}$$

$$H = 0.2$$

(b) LRU sequence



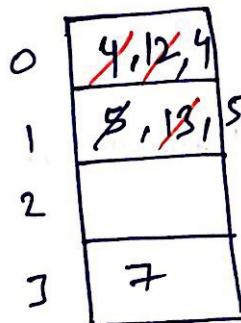
13-M 4-H 5-H 7-M



$$\text{Hit ratio} = \frac{4}{10}$$

$$H = 0.4$$

(c) Direct mapping :-



$$4-M \quad 4 \bmod 4 = 0$$

$$5-M \quad 5 \bmod 4 = 1$$

$$7-M \quad 7 \bmod 4 = 3$$

$$12-M \quad 12 \bmod 4 = 0$$

$$\text{Hit ratio} = \frac{3}{10}$$

$$4-M \quad 4 \bmod 4 = 0$$

$$5-H \quad 5$$

$$13-M \quad 13 \bmod 4 = 1$$

$$H = 0.3$$

$$4-H$$

$$5-M \quad 5 \bmod 4 = 1$$

$$7-H$$

$$(d) \quad N = 4$$

$$S = \frac{4}{2} = 2$$

0	4	12
1	5	7, 13, 7

$$K \bmod S = i$$

$$\text{Hit ratio} = \frac{4}{10}$$

$$H = 0.4$$

$$4 - M \quad 4 \bmod 2 = 0$$

$$5 - M \quad 5 \bmod 2 = 1$$

$$7 - M \quad 7 \bmod 2 = 1$$

$$12 - M \quad 12 \bmod 2 = 0$$

$$4 - H$$

$$5 - H$$

$$13 - M \quad (\text{LRU}) \quad 13 \bmod 2 = 1$$

$$4 - H$$

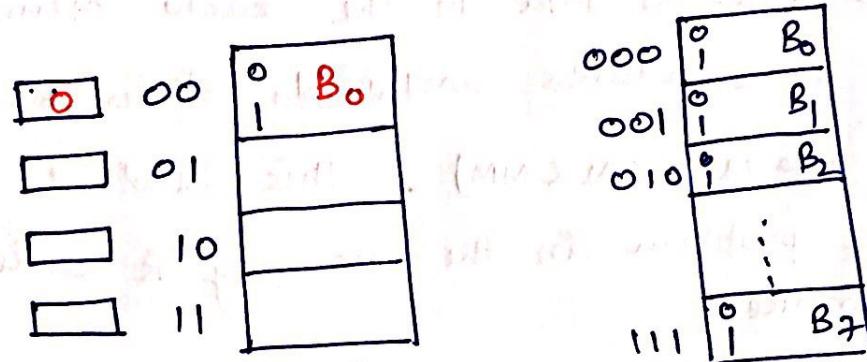
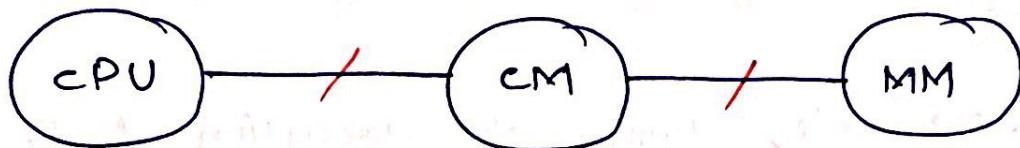
$$5 - H$$

$$7 - M \quad (\text{LRU}) \quad 7 \bmod 2 = 1$$

Updating Techniques :-

- According to memory hierarchy design, CPU performs the READ & WRITE opⁿ only on a cache memory.
- In this regard, CPU always checks the availability of data in the cache memory. When it is available, then opⁿ become READ hit or WRITE hit otherwise opⁿ becomes READ miss or WRITE miss.

- When miss occurred then respective block is copied from the MM to CM called as READ allocate or WRITE allocate. After the allocation, CPU performs READ & WRITE opⁿ only on a cache memory.



Prog:-

I₁: MOV r₀, [0000] B₀ - 0th byte \rightarrow r₀

I₂: ADD r₀, #23 r₀ \leftarrow r₀ + 23

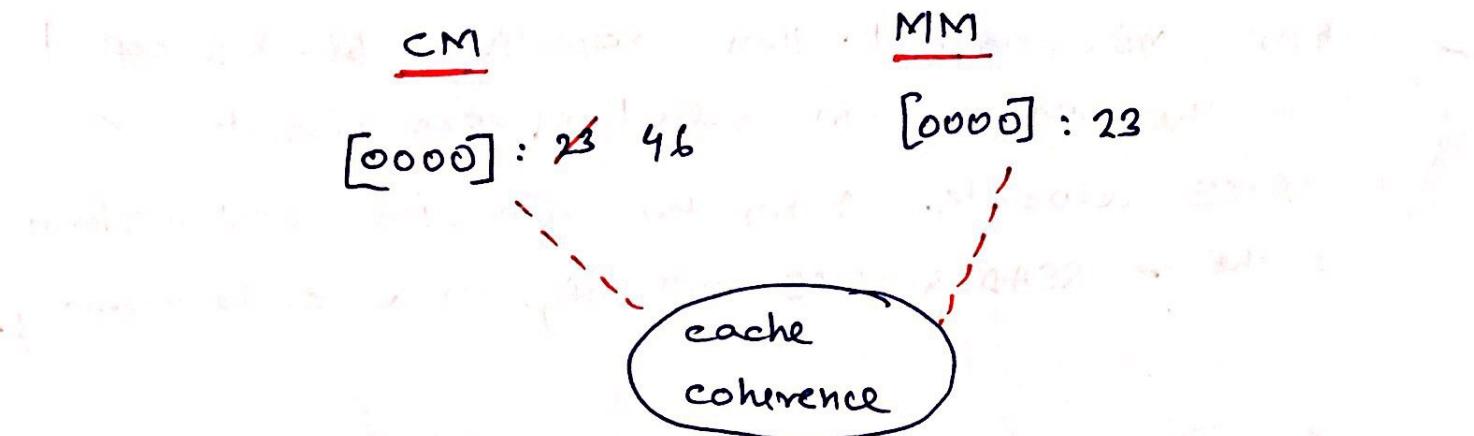
I₃: MOV [0000], r₀ B₀ - 0th byte \leftarrow r₀

Execu:-

I₁: READ miss; READ allocate ; r₀ = 23
 $(0 \bmod 4 = 0)$

I₂: r₀ + 23 \rightarrow r₀ ; r₀ = 46

I₃: WRITE hit, CM: [0000] : 23 46



- In the above code, during the execution of I_3 instⁿ CPU updates the data in the cache memory only therefore same address contains different value at different places (CM & MM) . This kind of data-inconsistency problem in the memory is called as cache coherence.
- Coherence causes data loss described below -

$I_4:$ $\text{MOV } r_0, [1000] ;$ $B_4 - 0^{\text{th}}$ byte $\rightarrow r_0$

$I_5:$

$I_6:$

$I_7:$

$I_8:$ $\text{MOV } r_0, [0000] ;$ $B_0 - 0^{\text{th}}$ byte $\rightarrow r_0$

Execution:-

$I_4:$ READ miss; READ - allocate

$$K \bmod N = i$$

$$4 \bmod 4 = 0$$

$$\text{CM : } B_0 \quad B_4$$

CM: By \rightarrow updated data is lost 'Ig' Instn

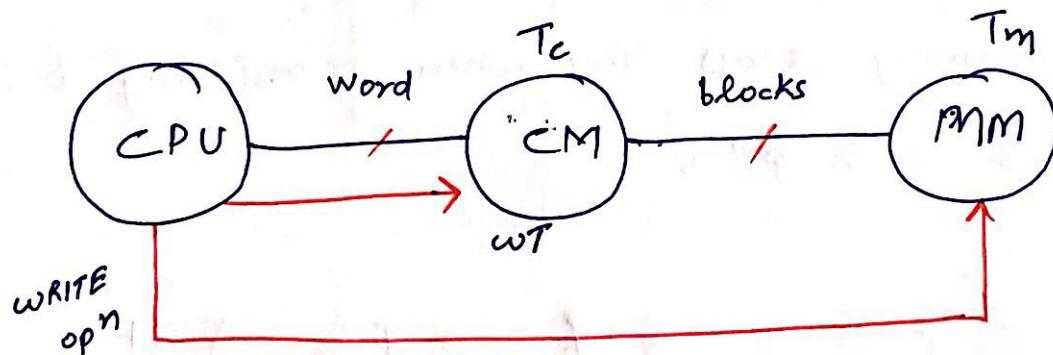
incorrectly access the old value from the block 'B₀'. Data Loss occurred.

- To handle the above problem updating techniques are used in the cache design. These are of two kinds:

Write Through
protocol

write back
protocol

- (i) Write Through protocol:- In this protocol, CPU performs the simultaneous write opⁿ in both CM and MM so no coherence. Inclusion is always ~~is~~ success. Inclusion means lower level memory data is always subset of higher level memory data.



simultaneous
WRITE opⁿ

$$(T_w) = \text{Max} \left(\begin{array}{l} \text{word update time in CM} \\ \text{word update time in MM} \end{array} \right)$$

- Time required to READ one-word data from the memory is called as Read - cycle time T_{avg} .

ce

$$T_{avg} = H_r T_c + (1-H_r)(T_m + T_c)$$

↓ ↓ ↓ ↓ ↓
Read Hit Read Data Read miss Read allocate Read data

- Time required to WRITE 1-word data into a memory is called as write - cycle time T_{avg_write} .

$$T_{avg_write} = H_w T_w + (1-H_w)(T_m + T_w)$$

↓ ↓ ↓ ↓ ↓
Write Hit Write simultaneous Write miss Write allocate Simultaneous write

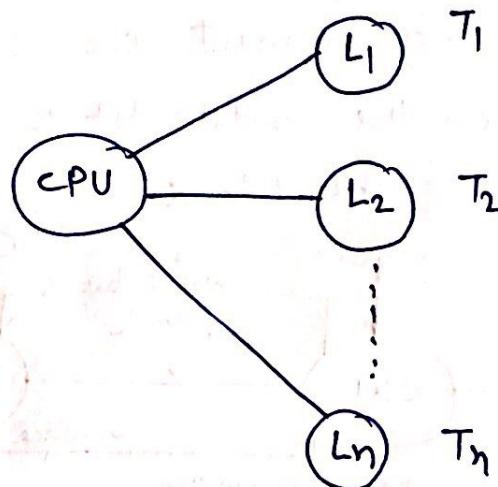
- Avg. memory access time when considering both READ & WRITE opⁿ is

$$T_{avg_wt} = (f_{read} * T_{avg_read}) + (f_{write} * T_{avg_write})$$

$$\gamma_{WT} = \frac{1}{T_{avg}} \text{ words/sec}$$

NOTE:- when the write-through protocol is used in a simultaneous access memory orgⁿ then H_w always becomes 1 because data block is not needed in the first level (L1 memory). implemented

$$H_w = 1$$



$$H_w = 1$$

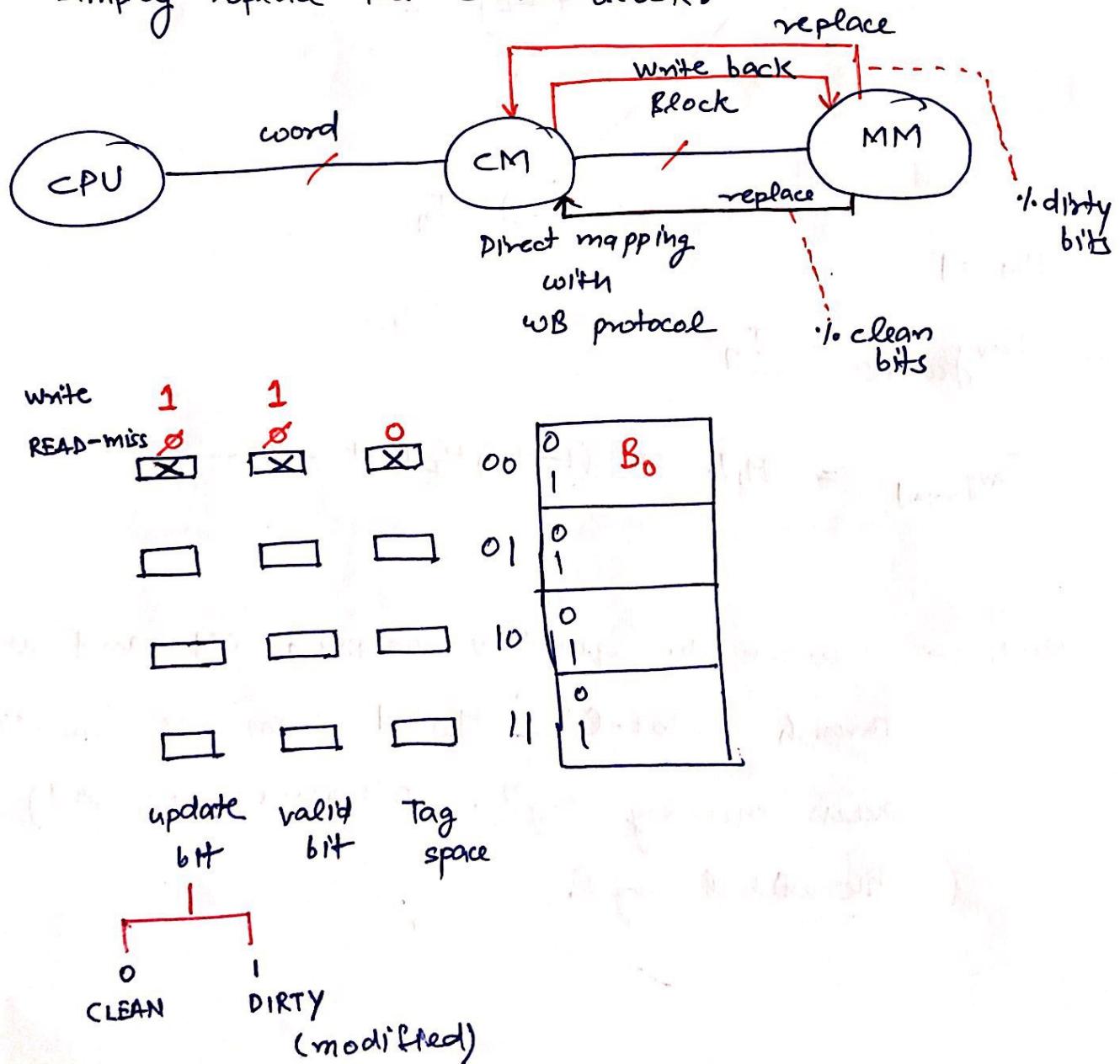
$$T_{avg\ write} = T_n$$

$$T_{avg\ read} = H_1 T_1 + (1-H_1) H_2 T_2 + \dots$$

NOTE:- when the question containing CM and write-through protocol & $H_w = 1$ then use simultaneous access memory orgⁿ. otherwise ($H_w \neq 1$), use hierarchical orgⁿ.

(ii) Write-back protocol :- In the write back cache design each line contain 1-bit storage space called as update bit, use to hold the status of a updation. This bit is set, when the cache block is updated otherwise, reset.

- In this protocol, CPU performs the WRITE opⁿ only in the cache memory . So, coherence present . It doesn't causes the data loss bcz CPU reads the status of update bit before replacement i.e CPU writeback the cache block when the update bit is 1 otherwise simply replace the cache block.



- Read cycle-time T_{avg} _{read}

$$T_{\text{avg, read}} = H_r T_c + (1 - H_r) \left[\begin{array}{l} \text{.1. dirty} \\ \text{- bits} \\ | \\ T_m + T_m + T_c \\ | \\ \text{write} \\ \text{back} \end{array} \right] + \begin{array}{l} \text{.1. clean} \\ \text{- bits} \\ | \\ T_m + T_c \\ | \\ \text{(replace} \\ \text{only)} \end{array} \right]$$

- Write - cycle time T_{avg_write} ;

$$\text{Write - cycle time} = T_{\text{arg write}} ;$$

write back
 write allocate
 write data

$$T_{\text{arg write}} = H_w T_c + (1-H_w) \left[\begin{array}{l} \% \text{ dirty bits} \\ \text{write miss} \end{array} \right] + \% \text{ clean bits} (T_m + T_c)$$

write hit
 write data
 write allocate
 write data

Avg. memory access time when considering both read & write op's is

$$T_{avg_{WB}} = \left(F_{read} * T_{avg_read} \right) + \left(F_{write} * T_{avg_write} \right)$$

$$\eta_{wB} = \frac{1}{T \arg_{wB}} \text{ words/sec}$$