

COMPUTER SCIENCE

Computer Organization and Architecture

Cache Memory

Lecture_07



Vijay Agarwal sir





**TOPICS
TO BE
COVERED**

- o1 Mapping Techniques**
- o2 Replacement Algo & Updating Technique**

Cache Memory

① Memory org

CM

L.O	W.O
-----	-----

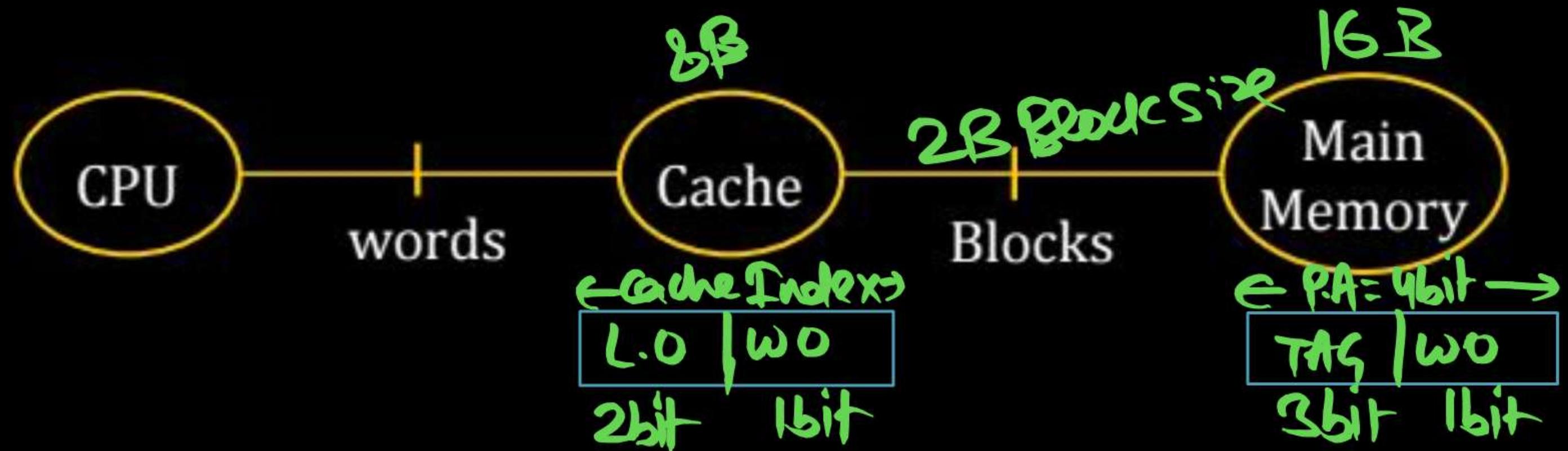
Main memory

TAG	W.O
-----	-----

② Mapping Technique: MM to CM

- ① Direct Mapping
- ② Set Associative Mapping
- ③ Fully associative Mapping.

Memory Organization



Mapping

The process of transfer the Data from Main Memory to Cache Memory is called mapping. There are 3 Type of Mapping Technique

- 1) Direct Mapping
- 2) Set Associative Mapping
- 3) Fully Associative Mapping

Mapping Function

- ❑ Because there are fewer cache lines than main memory blocks, an algorithm, is needed for mapping main memory blocks into cache lines.
- ❑ Three techniques can be used:

Direct

- The simplest technique
- Maps each block of main memory into only one possible cache lines.

Associative

- permits each main memory block to be loaded into any line of the cache.
- The cache control logic interprets a memory address simply as a Tag a word field
- To determine whether a block is in the cache, the cache control logic must simultaneously examine every line's Tag for a match

Set Associative

- A compromise that exhibits the strength of both the direct and associative approaches while reducing their disadvantage.

Direct Mapping

In this Direct Cache Controller interprets the CPU generated Request as follows:



$$\# \text{LINE} = \frac{\text{CM Size}}{\text{BLOCK Size}}$$

TAG = Physical Address - (Line offset + Word offset)

TAG Memory Size = #LINE's \times Tag bits (Depend on the mapping technique)

1) Direct Mapping

In this Technique mapping function is used to transfer the data from Main Memory to Cache Memory. The Mapping Function is

$$\text{Cache address} = \text{Main Memory request} \quad \text{MOD} \quad \# \text{ CM LINES}$$

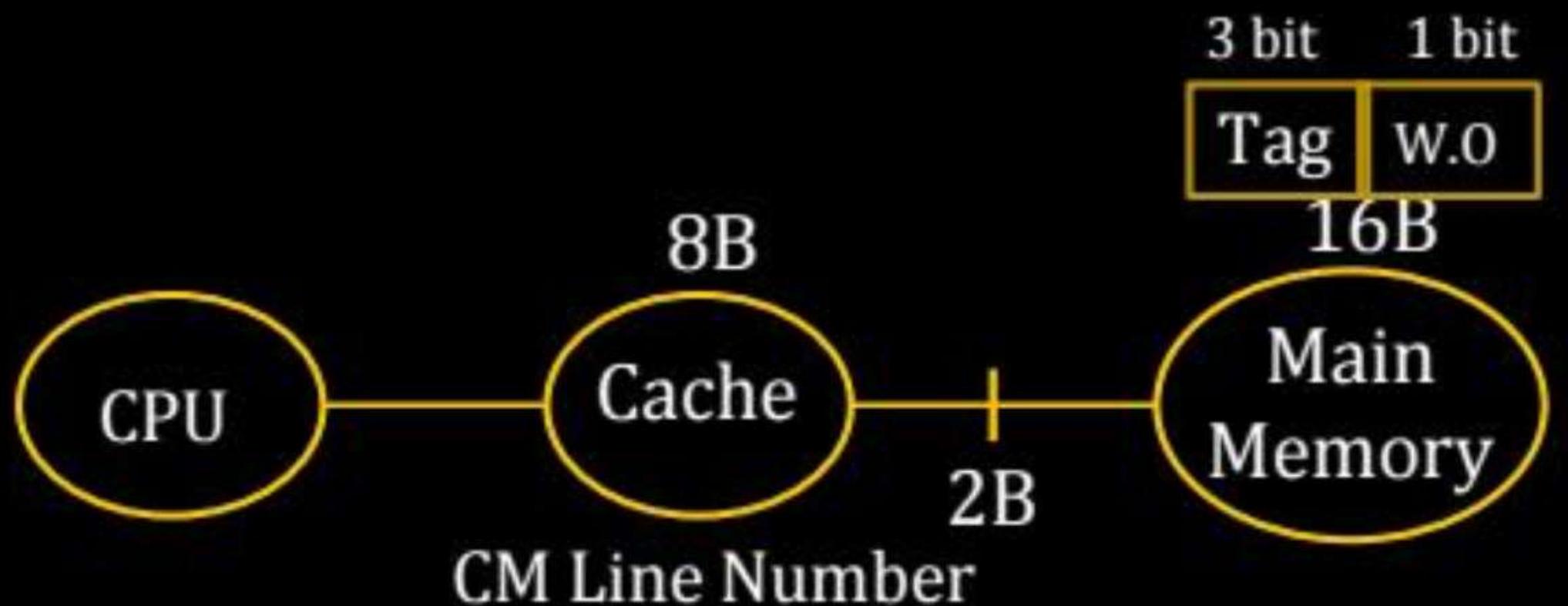
(Or)

$$K \text{ MOD } N = i$$

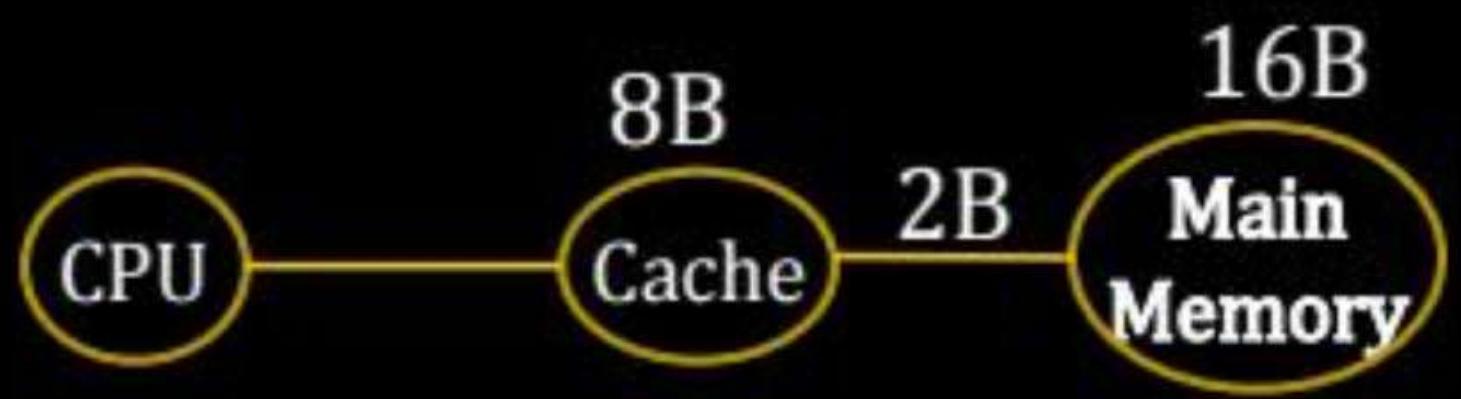
K: MM Block No.

N: # of Cache Line

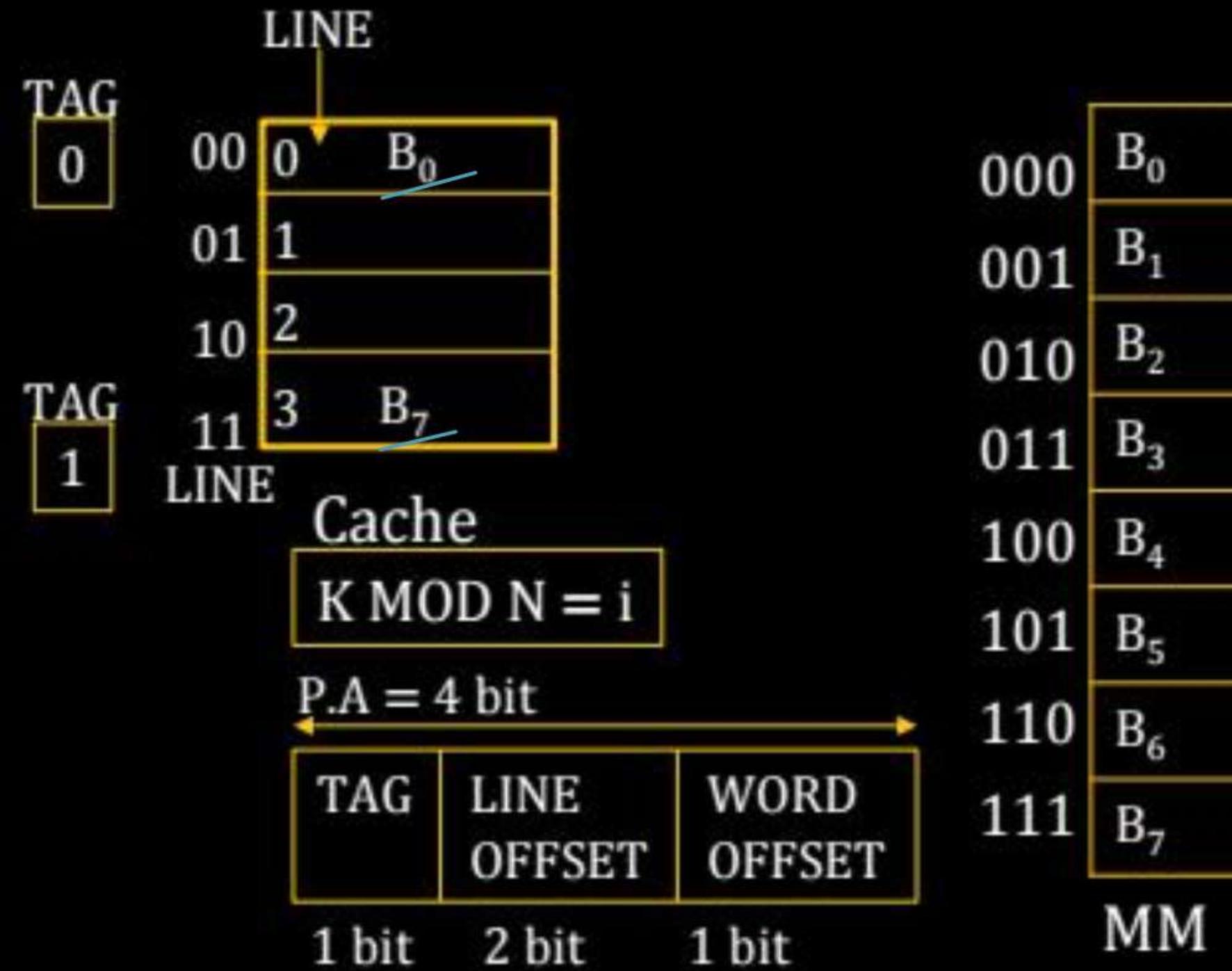
i: CM Line Number



	CM
$0 \bmod 4 = 0$	$B_0 \& B_4$: LINE 0
$1 \bmod 4 = 1$	$B_4 \& B_5$: LINE 1
$2 \bmod 4 = 2$	$B_2 \& B_6$: LINE 2
$3 \bmod 4 = 3$	$B_3 \& B_7$: LINE 3
$4 \bmod 4 = 0$	
$5 \bmod 4 = 1$	
$6 \bmod 4 = 2$	
$7 \bmod 4 = 3$	



P
W



Direct Mapping

MM Block

TAG LINE

0	00
---	----

 $B_0_{[000]}$

Direct Mapping

$$\frac{K \bmod N = i}{0 \bmod 4 = '0'}$$

CM LINE

LINE '0'

TAG LINE

1	11
---	----

 $B_7_{[111]}$

$$\frac{K \bmod N = i}{7 \bmod 4 = '3'}$$

LINE '3'

$$\text{Tag Memory Size} = \# \text{LINE's} \times \text{Tag bits}$$

Depends
On the
Mapping
technique

In the above example: # LINE = 4
Tag bit = 1 bit
(Direct Mapping)

$$\text{Tag Memory Size} = 4 \times 1 = 4 \text{ bits}$$

LOR [Locality of Reference]

Consider the following program

I₁: MOV r₀ [0000]

I₂: MOV r₁ [0001]

I₃: MOV r₂ [1000]

I₄: MOV r₃ [1001]

Disadvantage of Direct Mapping

I₁: MOV r₀ [0000] Bo

I₂: MOV r₁ [1000] By

I₃: MOV r₂ [0001] Bo

I₄: MOV r₃ [1001] By

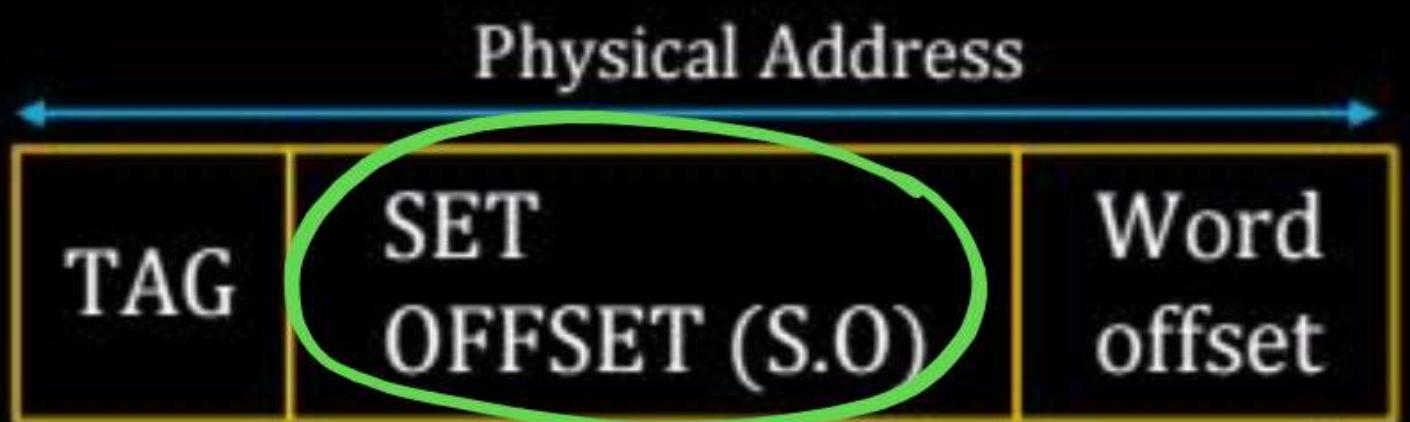
I₅: MOV r₄ [0000] Bo

I₆: MOV r₅ [1000] By

Conflict Miss Increase.

2) Set Associative Cache

SET associative cache controller, Interpreter the CPU generated request as follows:



$$\text{Word Offset} = \log_2 \text{Block Size}$$

$$\#SETS = \frac{\#Lines}{N\text{-way}}$$

$$\text{SET OFFSET} = \log_2 \#SETS$$

$$\text{TAG} = \text{Physical address} - (\text{S.O} + \text{W.O})$$

$$\#SET = \frac{\#LINE}{N\text{-Way}}$$

$$\#lines = \frac{\text{CM Size}}{\text{Block Size}}$$

Set Associative Mapping function

Cache

Set = MM Request **MOD** #SET is in Cache
Address

(OR)

K MOD S = i

k: MM Block Number
s: # Cache Set
i: Cache Set Number

Example1: # Line's = 16 & 2 way set Associative

$$\# \text{SET} = \frac{\# \text{LINES}}{\text{N-way}} = \frac{16}{2} = 8$$

#SET = 8 S = 8

K MODS = i

K MOD 8 = i

✓	✓	✓	✓	✓	✓	✓	✓
---	---	---	---	---	---	---	---

SET₀

SET₁

2

3

4

5

SET 6

SET 7

K MOD 8 = i

0 - 7

Example2:

#LINE = 16 & 4 way set Associative

$$\# \text{SET} = \frac{16}{4} \Rightarrow 4$$

S = 4

K MOD 4 = i

SET 0



SET 1

SET 2

SET 3

Cache

Example3:

#LINE = 16 & 8 way set Associative

$$\# \text{ SET} = \frac{16}{8} \Rightarrow 2 \quad S = 2$$

$$K \bmod 2 = i$$

SET 0

8 way

SET 1

8 way

Example4:

#LINE = 16 & 16 way set Associative

$$\# \text{ SET} = \frac{16}{16} \Rightarrow 1$$

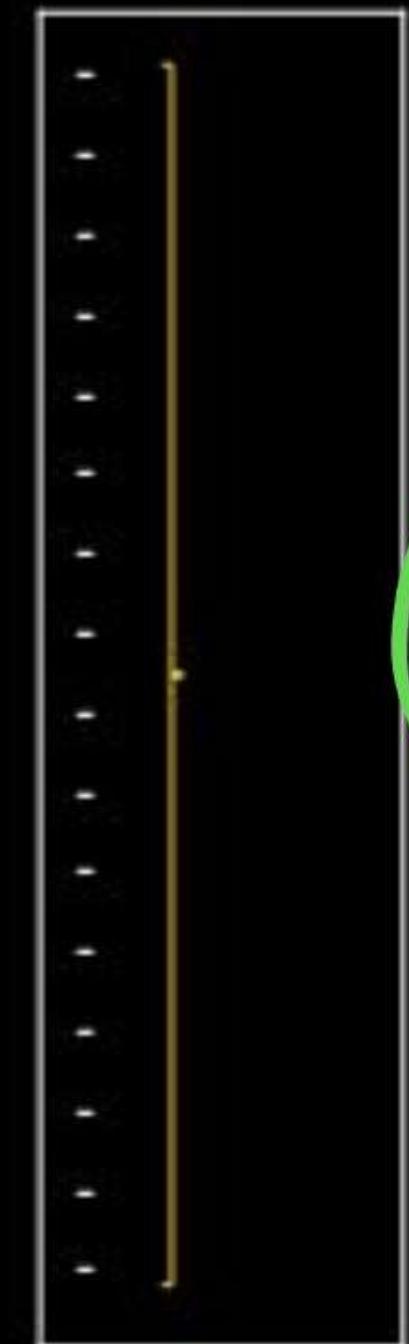
$$S = 1$$

$$K \bmod 1 = i$$

#LINE = 16

16 Way

$$S = 1$$



Fully
Associative
(Whole cache as a set)

Set Associative Mapping function

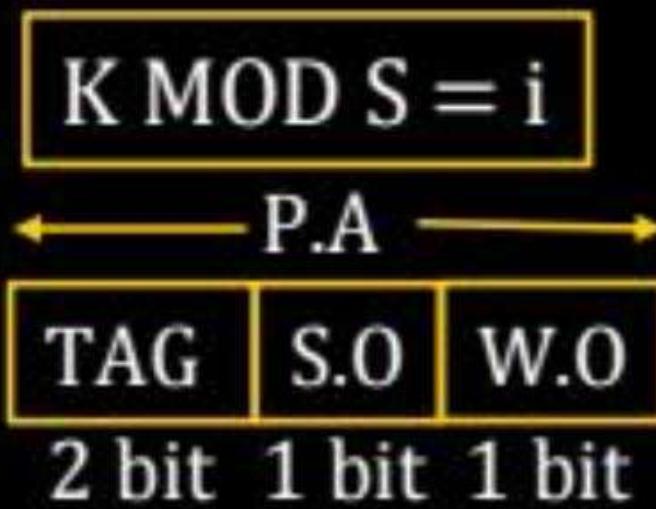
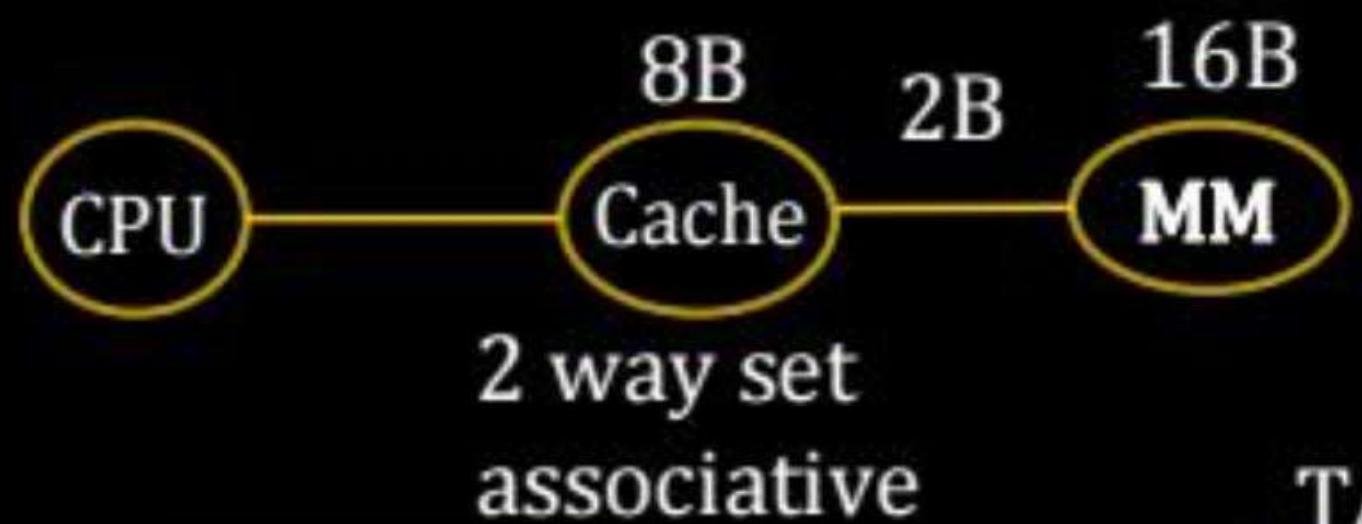
Cache

Set = MM Request MOD #SET is in Cache
Address

(OR)

K MOD S = i

k: MM Block Number
s: # Cache Set
i: Cache Set Number



$B_0[000] \rightarrow$

TAG	S.O
00	0

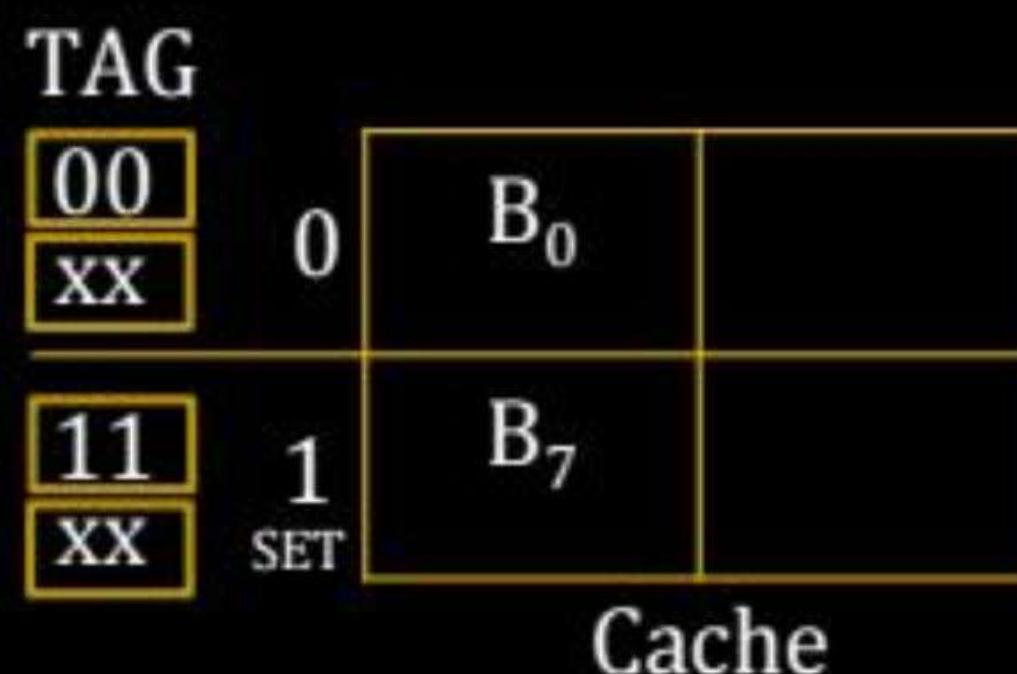
 → SET '0'

2 bit 1 bit

$B_7[111] \rightarrow$

TAG	S.O
11	1

 → SET '1'



P W

000	B_0
001	B_1
010	B_2
011	B_3
100	B_4
101	B_5
110	B_6
111	B_7

MM BLOCK Mapping Tech. CM SET

$$B_0[000] \xrightarrow{\begin{array}{l} K \bmod S = i \\ 0 \bmod 2 = '0' \end{array}} \text{SET '0'}$$

$$B_7[111] \xrightarrow{\begin{array}{l} K \bmod S = i \\ 7 \bmod 2 = '1' \end{array}} \text{SET '1'}$$

Block

$$\frac{\text{Tag Memory}}{\text{Size}} = \text{#SETS} \times \frac{\text{#LINES}}{\text{In a each set}} \times \text{Tag bits}$$

Example: # SET = 2 $\frac{\text{# LINE in}}{\text{Each set}} = 2$ TAG = 2 bit

$$\begin{aligned}\text{Tag Memory} &= 2 \times 2 \times 2 \\ \text{size} &= 8 \text{ bits}\end{aligned}$$

1) In direct Mapping

Hit Latency = Latency of Tag comparator

2) In Set Associative Mapping

Hit Latency = Latency of Tag comparator + Latency of Multiplexer

Set Associative Mapping:

$K \bmod S = i$

$$S = \frac{\# \text{LINES}}{N-\text{ways}}$$

From Mapping onwards

29 GATE QUESTION

+ Practice Question

Important Points About Set Associative



#LINE = L & N way set Associative

$$\# \text{SET}[S] = \frac{L}{N}$$

$$N=L$$

If $N=1$; Direct Mapping

If $N=L$; Fully Associative Mapping; ($S=1$) ie Only 1 Set

$$K \bmod S = i$$

k: MM Block Number
s: # Cache Set
i: Cache Set Number

1 way

$$\# \text{LINE} = 8$$

1 way

$$\# \text{SET} = \frac{8}{1} = 8$$

Direct Mapping.

$$\# \text{LINE} = 8$$

8 Way Set Associative

$$\# \text{SET} = \frac{8}{8} = 1 \quad [\text{Fully Association}]$$

Example4:

#LINE = 16 & 16 way set Associative

$$\# \text{ SET} = \frac{16}{16} \Rightarrow 1 \quad S = 1$$

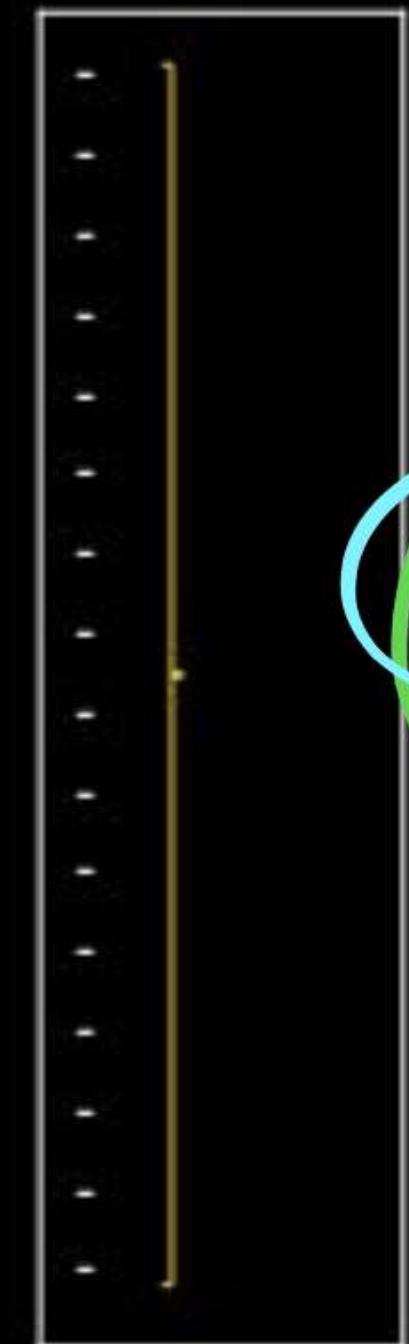
$$K \bmod 1 = i$$

L = N

#^(L) LINE = 16

(N) 16 Way

S = 1



Fully
Associative
(Whole cache as a set)

Important Points About Set Associative

#LINE = L & N way set Associative

$$\# \text{SET}[S] = \frac{L}{N}$$

If N=1; Direct Mapping

If N=L ; Associative Mapping; (S=1)ie Only 1 Set;

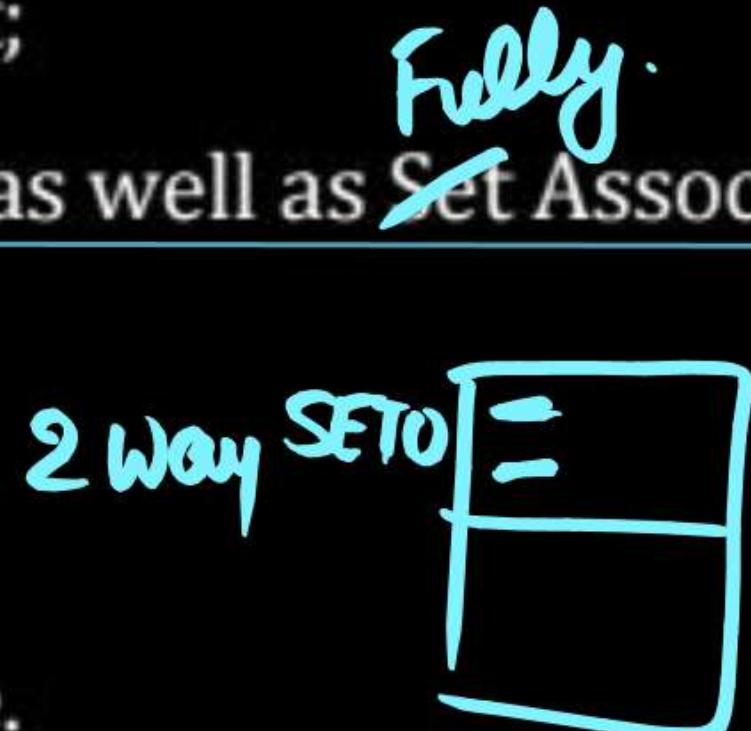
Set Associative Mapping Follows Direct Mapping as well as ~~Set Associative~~

Mapping. Using Mapping Function as:

$$K \bmod S = i$$

within the Set MM Block Can be placed anywhere.

- In the Set Associative Mapping Replacement Algorithm is required along with the Mapping function when cache set is Full.



Important Points About Set Associative

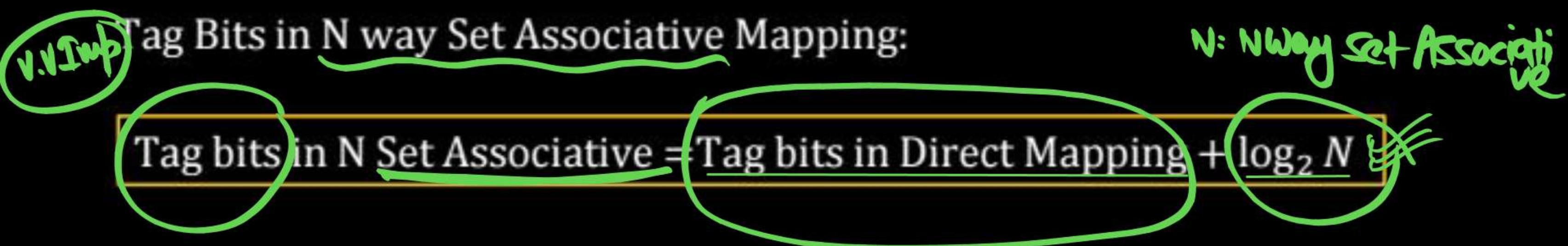


#LINE = L & N way set Associative

$$\# \text{SET}[S] = \frac{L}{N}$$

If N=1; Direct Mapping

If N=L ; Associative Mapping; (S=1)ie Only 1 Set.



Important Points About Set Associative

Tag bits in N Set Associative = Tag bits in Direct Mapping + $\log_2 N$

Eg. Consider a **Direct Mapping** if the size of cache memory is 512KB & Main Memory 512MB & Cache line size is 2KB then calculate the Number of bit Required for TAG? $\rightarrow 29\text{bit}$



$$\begin{aligned} \# \text{LINES} &= \frac{\text{CM Size}}{\text{Block Size}} \\ &= \frac{512\text{KB}}{2\text{KB}} = \frac{2^9}{2^{11}} = 2^8 \end{aligned}$$

Tag bit in Direct Mapping = 10bit

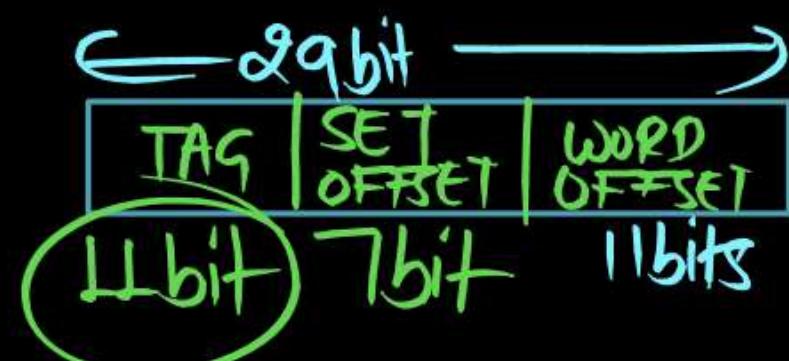
$$L.O = 8\text{bit}$$

Important Points About Set Associative

Tag bits in N Set Associative = Tag bits in Direct Mapping + $\log_2 N$

- Q1 Eg. Consider a **2-way set associative** if the size of cache memory is 512KB & Main Memory 512MB & Cache line size is 2KB then calculate the Number of bit Required for TAG?

2 Way Set Associative



$$\# \text{LINES} = 2^8$$

$$\# \text{SET} = \frac{2^8}{2} = 2^7 \text{ SET}$$

$$SO = 7 \text{ bit}$$

2 Way Set Associative

$$\begin{aligned} \text{Tag bit} &= \text{Direct mapping} + \log_2 1 \\ &= 10 + 1 \end{aligned}$$

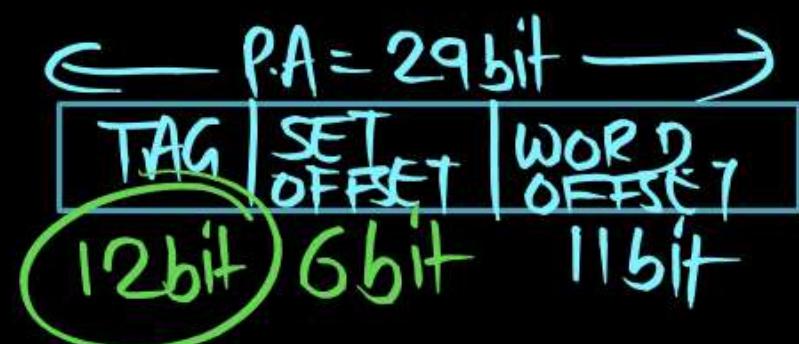
Tag = 11 bits

Important Points About Set Associative

Tag bits in N Set Associative = Tag bits in Direct Mapping + $\log_2 N$

- Q2 Eg. Consider a **4-way set associative** if the size of cache memory is 512KB & Main Memory 512MB & Cache line size is 2KB then calculate the Number of bit Required for TAG?

4 way set associative.



$$\# \text{LINES} = 2^8$$

$$\# \text{SET} = \frac{2^8}{2^2} = 2^6 \text{ SET}$$

$$S.O = 6 \text{ bit}$$

4 Way Set Associative

$$\text{Tag bit} = 10 + \log_2 4$$

$$= 10 + 2$$

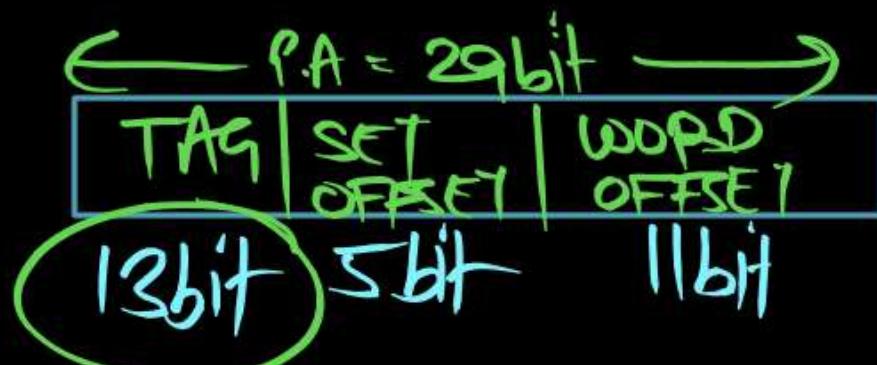
$$= \underline{12 \text{ bit}}$$

Important Points About Set Associative

Tag bits in N Set Associative = Tag bits in Direct Mapping + $\log_2 N$

- Q3 Eg. Consider a **8-way set associative** if the size of cache memory is 512KB & Main Memory 512MB & Cache line size is 2KB then calculate the Number of bit Required for TAG?

8 way Set Associative



$$\begin{aligned} \# \text{LINE} &= 2^8 \\ \# \text{SET} &= \frac{2^8}{2^3} = 2^5 \text{ Set} \\ S0 &= 5 \text{ bit} \end{aligned}$$

8 Way Set Associative

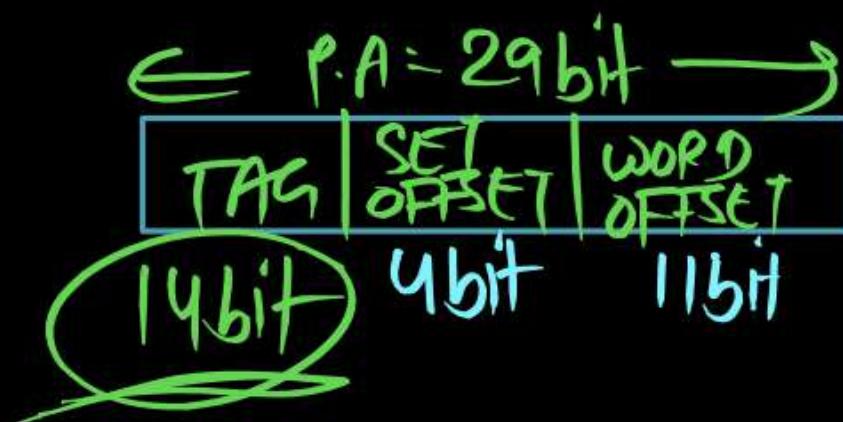
$$\begin{aligned} \text{Tag bit} &= 10 + \log_2 8 \\ &= 10 + 3 \\ \text{Tag} &= 13 \text{ bit} \end{aligned}$$

Important Points About Set Associative

Tag bits in N Set Associative = Tag bits in Direct Mapping + $\log_2 N$

- Q4 Eg. Consider a **16-way set associative** if the size of cache memory is 512KB & Main Memory 512MB & Cache line size is 2KB then calculate the Number of bit Required for TAG?

16 way Set Associative



$$\# \text{LINE} = 2^8$$
$$\# \text{SET} = 2^4 \rightarrow 2^4 \text{ SET}$$

$$S.O = 4 \text{ bit}$$

16 Way Sel Associative

$$\begin{aligned}\text{Tag bit} &= 10 + \log_2 16 \\ &= 10 + 4\end{aligned}$$

Tag = 14 bit

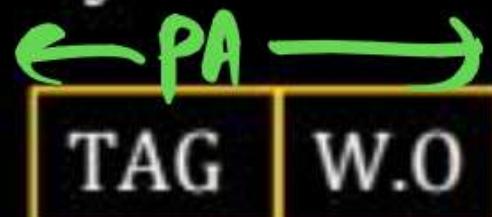
③ Associative Mapping:

Till
71 Pipeline
10 Approach in Cache
29 Today Homework

110 PYQ's.

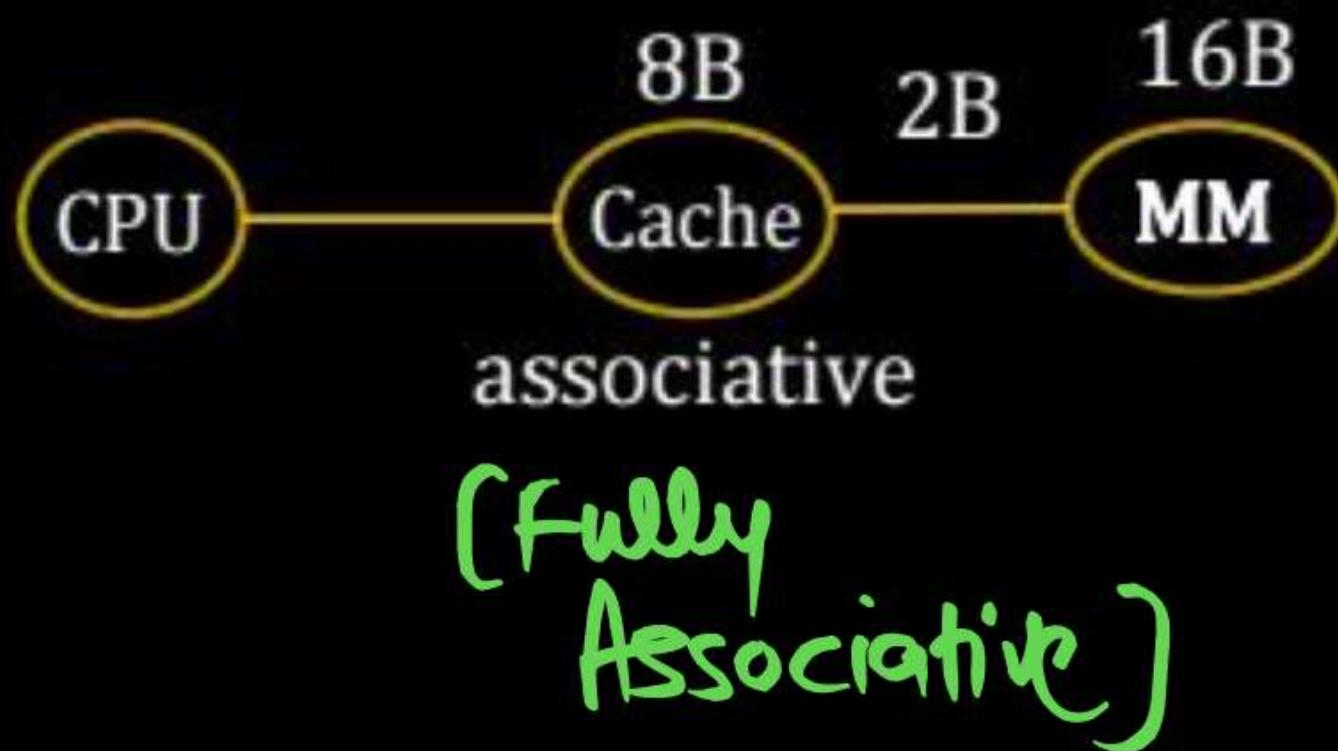
3) Associative Mapping

- In this no mapping function is used to transfer the data from MM to CM.
- Any Block of Main Memory can be placed Anywhere in Cache Memory.
- No Conflict Miss;
- This Cache is designed without address called as Content Addressable Memory.
- In Associative Mapping More Tag bits is required & More Tag Memory size.
- In Associative Mapping More Hardware Required, Expensive(N Tag Comparator, Here N is Number of Lines., For Each line Comparator Required)
- In Associative Cache Design, Counter Sequence is used to map the Data, means any Main Memory Block any cache memory line in a sequence.
- In Associative Mapping Physical Address is Interpreted as:

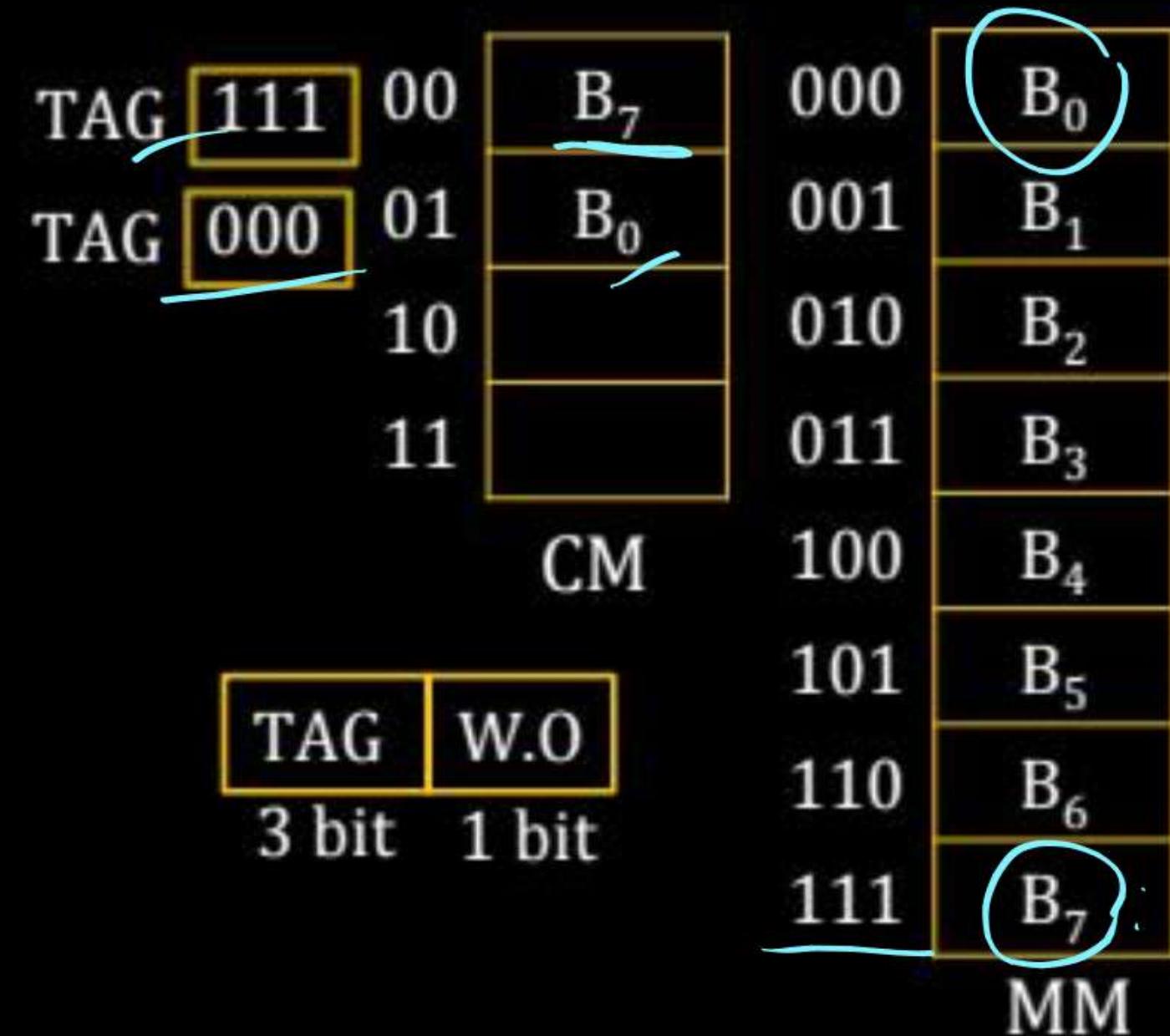


3) Associative Mapping

- In this no mapping function is used to transfer the data from MM to CM.
- No Conflict Miss;



4CMLINE \Rightarrow 4 Tag Comparator Required



MM Block Associative Mapping

$B_7[111]$ No mapping Function

$B_0[000]$ No mapping Function

$$\boxed{\text{Tag Memory Size} = \# \text{LINES} \times \text{Tag bits}}$$

Example: #LINE = 4 & Tag bits = 3

$$\boxed{\text{Tag Memory Size} = 4 \times 3 = 12 \text{ bits}}$$

CM Line



Any line

Any line

① In Direct Mapping

$$\text{Tag Memory Size} = 4 \times 1 = \underline{\underline{4 \text{ bit}}}$$

② In 2 Way Set Associative

$$\text{Tag Memory Size} = 4 \times 2^{\frac{(\text{Topbit})}{2}} = \underline{\underline{8 \text{ bit}}}$$

③ In Fully Associative

$$\text{Tag Memory Size} = 4 \times 3^{\frac{(\text{Topbit})}{3}} = \underline{\underline{12 \text{ bit}}}$$

In Direct Mapping

$$\text{Tag bit} = \frac{\text{MM Size}}{\text{CM Size}}$$

Today I told
Directly in
Set Associative mapping

$$\text{Top bit} = \# \text{Top bit in Direct Mapping} + \log_2 N$$

Directly calculate
Tag bit

N: Nway
Set
Associative

L Mixing type of
Question

Q.

Consider a 64KB Direct Mapped Cache organized into a 64 word blocks. Word length of the CPU is 32bits. Main Memory 4GB. In the Cache Controller is comprising of 1 Valid bit & 1 Update bits.

the calculate the number of bit required for

- (i) P.A (32bit)
- (ii) TAG (16bit)
- (iii) L.O (8bit)
- (iv) W.O (8bit)
- (v) #LINES
- (vi) TAG Memory Size.**

2^8 Lines
(256 Lines)

Direct Mapped

Cache Controller: 1 Valid bit + 1 Updated bit

$$\text{Cache Size} = 64 \text{ KB} [2^{16} \text{ B}]$$

$$\text{Block Size} = 64 \text{ words}$$

$$[L] \text{ Word length (size)} = 32 \text{ bit}$$

$$\text{Main Memory} = 4 \text{ GB} [2^{32} \text{ B}]$$

$$\text{Block Size} = 64 \text{ words}$$

$$[L] \text{ Word Size} = 32 \text{ bit} = 4 \text{ Byte}$$

$$\text{Block Size} = \frac{64 \times 4 \text{ Byte}}{2^8} \Rightarrow 2^6 \times 2^2 \text{ B} \\ \Rightarrow 2^8 \text{ Byte}$$

$$\text{Word offset} = 8 \text{ bit}$$

$$\leftarrow \text{PA} = 32 \text{ bit} \rightarrow$$



$$16 \text{ bit} \quad 8 \text{ bit}$$

$$\# \text{LINE} = \frac{\text{Cache Size}}{\text{Block Size}}$$

$$\Rightarrow \frac{2^{16} \text{ B}}{2^8 \text{ B}}$$

$$2^8 \text{ B}$$

$$2^8 \text{ lines}$$

$$\text{Tag entry size} = \text{Tag bit} + \text{extra bits (if given)} \\ L.O = 8 \text{ bit}$$

$$\Rightarrow 16 + 1 + 1$$

$$\text{Tag entry size} = 18 \text{ bit}$$

$$\frac{\text{Tag memory size}}{\text{size}} = \# \text{LINE} \times \text{Tag bits}$$

$$= 2^8 \times 18$$

$$= 256 \times 18 \text{ bit}$$

$$= 4608 \text{ bits}$$

Q.

P
W

Consider fully associative cache consists 8 Block & MM contains 128 Block & Request made by the CPU:
119, 84, 37, 0, 16, 0, 84, 120, 121, 93, 37, 0, 43, 39, 47, 48.
Calculate # of compulsory & Capacity miss?



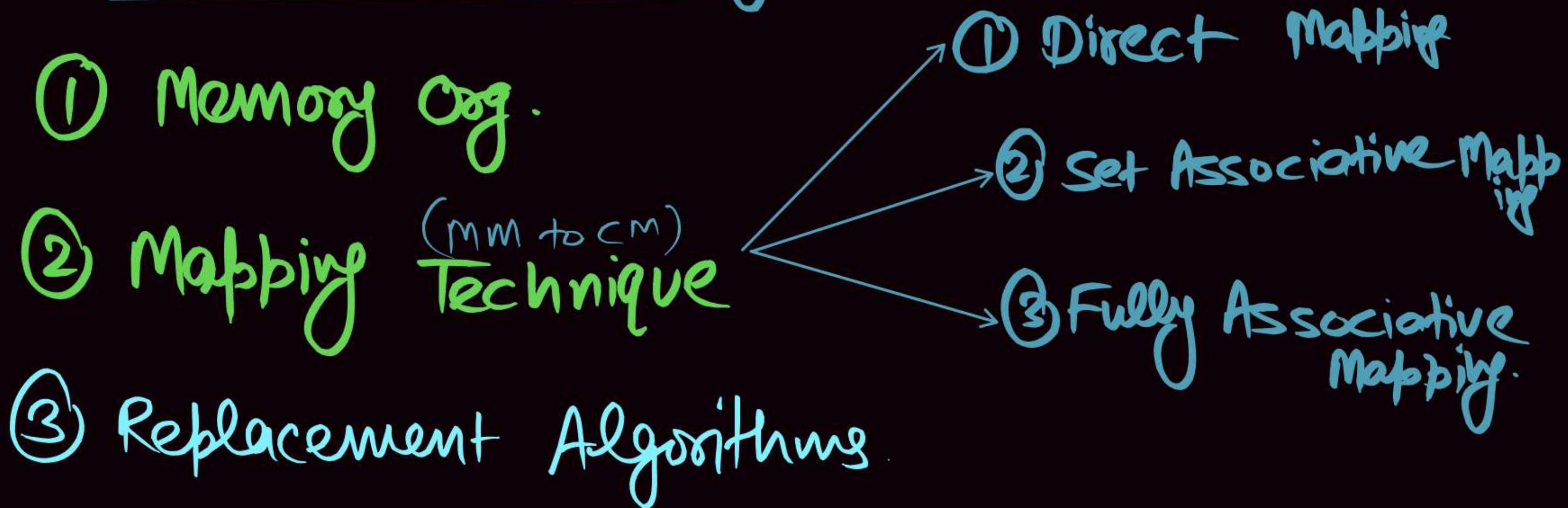
0	119 43
1	84 39
2	37 47
3	0 48
4	16
5	120
6	121
7	93

Hit: 0, 84, 37, 0

#Hits: 4

Total Miss = 12

Cache Work On Locality of Reference.



Replacement Algorithm

When Cache is Full, then replacement algorithm are required to replace the exist cache block with new block.

In the CM design 3 type of replacement algorithm is used.

1) Random Algorithm

✓ 2) FIFO Replacement

✓ 3) LRU Replacement [Least Recently Used]

In the random algorithm, any cache block can be replaced based on the random selection.

In the Cache There are 3 Type of Miss.

- ① Compulsory | First Reference | Cold Start Miss.
- ② Capacity Miss
- ③ Conflict Miss | Collision Miss.

Types of Misses

In the CM design 3 types of misses are present.

- 1) Compulsory miss - (Cold start miss / first reference miss)

This miss will occur when the very first reference to the cache itself
a miss.

- 2) Capacity Miss - This miss will occur when cache is full.

- 3) Conflict Miss (Collision miss / reference miss)

This miss will occur when the too many blocks are placed into same
cache line or same cache SET.

Q.1

Consider 4 block cache memory (initially empty) with the following MM block references.

7, 8, 10, 15, 7, 8, 16, 7, 8, 10

Identify the Hit Ratio using

(i) FIFO

(ii) LRU

(iii) Direct Mapped cache

(iv) 2 - way Set Associative with LRU

Number of Cache Line = ^(Block) 4

MM Block Request : 7, 8, 10, 15, 7, 8, 16, 7, 8, 10

↳ Total Request(Access) = 10

Q.1

Consider 4 block cache memory (initially empty) with the following MM block references.

7, 8, 10, 15, 7, 8, 16, 7, 8, 10

Identify the Hit Ratio using

(i) FIFO [FIRST IN FIRST OUT]

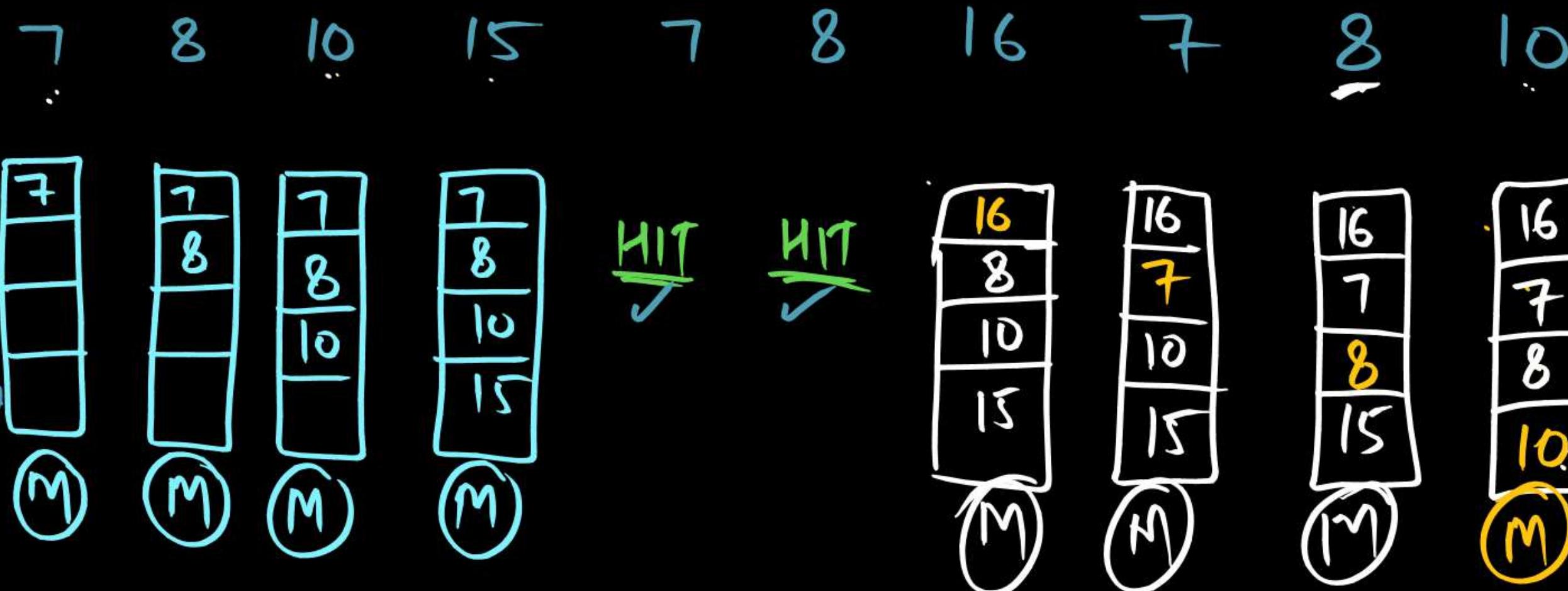
M: Miss

H: Hit

$$\begin{aligned} \# \text{ HIT} &= 2 \\ \# \text{ MISS} &= 8 \end{aligned}$$

$$\text{Hit Ratio} = \frac{2}{10}$$

Ans



Q.1

Consider 4 block cache memory (initially empty) with the following MM block references.

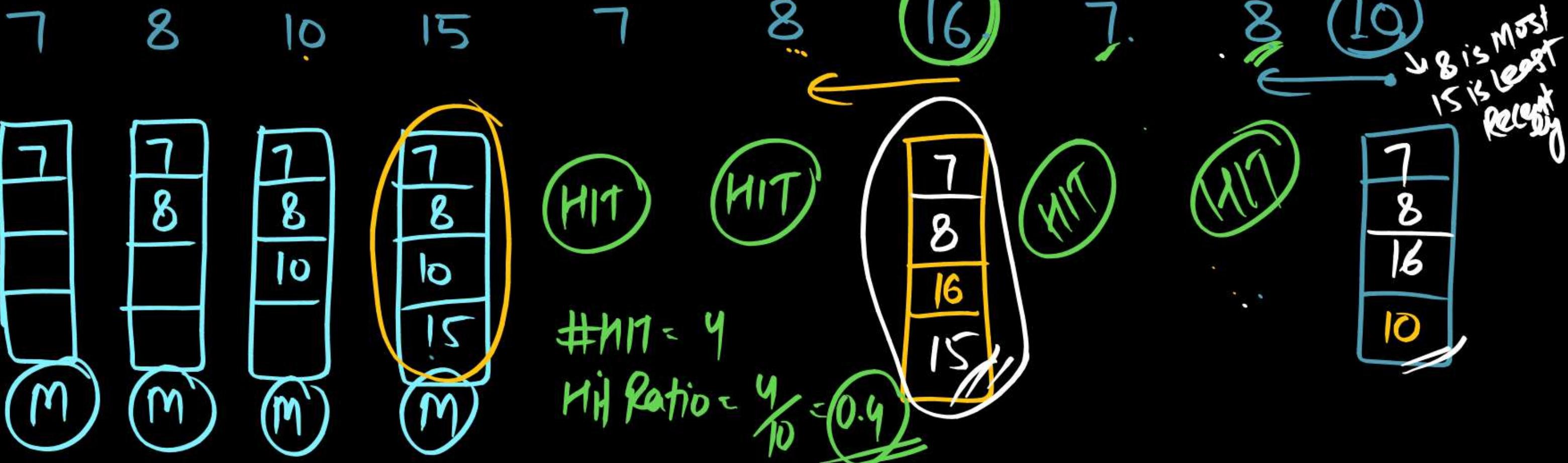
7, 8, 10, 15, 7, 8, 16, 7, 8, 10

Identify the Hit Ratio using

(ii) LRU [Least Recently USED]

P
W

8 is Most Recently Used
10 is the Least Recently Used



Q.1 Consider 4 block cache memory (initially empty) with the following MM block references.

7, 8, 10, 15, 7, 8, 16, 7, 8, 10

Identify the Hit Ratio using

(iii) Direct Mapped cache

	$k \bmod N = j$	$k \bmod 4 = j$	
7	$\frac{7}{4} = 1$	<u>Line No</u>	
8	$\frac{8}{4} = 2$		# HIT < 3
10	$\frac{10}{4} = 2$		
15	$\frac{15}{4} = 3$		
7	$\frac{7}{4} = 1$	M	
8	$\frac{8}{4} = 2$	M	
16	$\frac{16}{4} = 4$	M	
7	$\frac{7}{4} = 1$	HIT	
8	$\frac{8}{4} = 2$	HIT	
16	$\frac{16}{4} = 4$	HIT	
7	$\frac{7}{4} = 1$	HIT	
8	$\frac{8}{4} = 2$	HIT	
10	$\frac{10}{4} = 2$	HIT	

LINE 0	8	16	8
LINE 1			
LINE 2	10		
LINE 3	7	15	7

CM

Q.1

Consider 4 block cache memory (initially empty) with the following MM block references.

7, 8, 10, 15, 7, 8, 16, 7, 8, 10

Identify the Hit Ratio using

(iii) Direct Mapped cache

$$\begin{array}{l} \text{MODUS} \\ \text{K MOD N = } j \\ \text{K MOD 4 = } j \\ \text{Line No} \\ \text{comp 7} \\ \text{comp 8} \\ \text{comp 10} \\ \text{conflict 15} \\ \text{conflict 7} \\ \text{comp 8} \\ \text{comp 16} \\ \text{8} \\ \text{10} \\ \text{14} \end{array} = \begin{array}{l} 3 \\ 0 \\ 2 \\ 3 \\ 3 \\ 3 \\ 0 \\ 0 \\ 3 \\ 0 \\ 2 \end{array} \begin{array}{l} M \\ M \end{array}$$

HIT = 3

Hit Ratio = $\frac{3}{10}$

LINE 0	8	16	8
LINE 1	Empty		
LINE 2	10		
LINE 3	7	15	7

CM

Q.1 Consider 4 block cache memory (initially empty) with the

following MM block references. $\#SET = \frac{\#LINE}{N\text{-ways}} = \frac{4}{2} = 2$ Set

7, 8, 10, 15, 7, 8, 16, 7, 8, 10

k MODS Identify the Hit Ratio using

S = 2

MOD 2

K MOD 2 = i [Set Number]

$$7 \div 2 = 1 \quad M$$

$$8 \div 2 = 0 \quad M$$

$$10 \div 2 = 0 \quad M$$

$$15 \div 2 = 1 \quad M$$

$7 \div 2 = 1$ (Set No L) \rightarrow HIT

$8 \div 2 = 0$ \rightarrow HIT

$16 \div 2 = 0$ [Set No '0' But Set is Full, So Apply LRU].

$7 \div 2 = 1$ \rightarrow HIT

$8 \div 2 = 0$ \rightarrow HIT

$10 \div 2 = 0$ \rightarrow miss, set full Apply LRU

#HIT = 4
#MISS = 6
Hit Ratio = $\frac{4}{10} = 0.4$

SET 0	8	1	10	2 way LRU
SET 1	7	15		

2 way set associative with LRU

Q.1 Consider 4 block cache memory (initially empty) with the following MM block references. $\#SET = \frac{\#LINE}{N\text{-ways}} = \frac{4}{2} = 2$ set

7, 8, 10, 15, 7, 8, 16, 7, 8, 10

k_{MODS} Identify the Hit Ratio using

$$S=2$$

Mod 2

(iv) 2 - way Set Associative with LRU

$$k_{MOD2} = i[\text{Set Number}]$$

$$7 \div 2 = 1 \quad M$$

$$8 \div 2 = 0 \quad M$$

$$10 \div 2 = 0 \quad M$$

$$15 \div 2 = 1 \quad M$$

$$7 \div 2 = 1 \quad (\text{Set No.}) \rightarrow \text{HIT}$$

$$8 \div 2 = 0 \quad \rightarrow \text{HIT}$$

16 $\div 2 = 0$ [Set No. '0' But Set is Full, So Apply LRU].

$$7 \div 2 = 1 \rightarrow \text{HIT}$$

$$8 \div 2 = 0 \rightarrow \text{HIT}$$

10 $\div 2 = 0$ \rightarrow miss, Set Full Apply LRU

#HIT = 4.
#MISS = 6
Hit Ratio = $\frac{4}{10} = 0.4$

SET 0

SET 1

8	1	10	15
7			

2 way
+ LRU

2 way set
Associative with LRU

Q.

Consider a small two-way set-associative cache memory, consisting of 4 blocks. For choosing the block to be replaced, use the least recently used (LRU) scheme. The number of cache misses for the following sequence of block addresses is 8, 12, 0, 12, 8

$$K \bmod 2 = j$$

(a) 2

(b) 3

(c) 4

(d) 5

$$K \bmod 2 = j$$

SET Number

$8 \bmod 2 = 0$	M
$12 \bmod 2 = 0$	M
$0 \bmod 2 = 0$	Miss <small>Cache Full Apply LRU</small>
$12 \bmod 2 = 0$	HIT
$8 \bmod 2 = 0$	Miss <small>Cache set '0' is full, Apply LRU</small>

#Miss = 4

[GATE - 2004]

$$CM Lines = 4$$

2 way set associative

$$\# SET = \frac{4}{2} = 2$$



2 Way Set
+ LRU

Types of Misses

In the CM design 3 types of misses are present.

- 1) Compulsory miss - (Cold start miss / first reference miss)

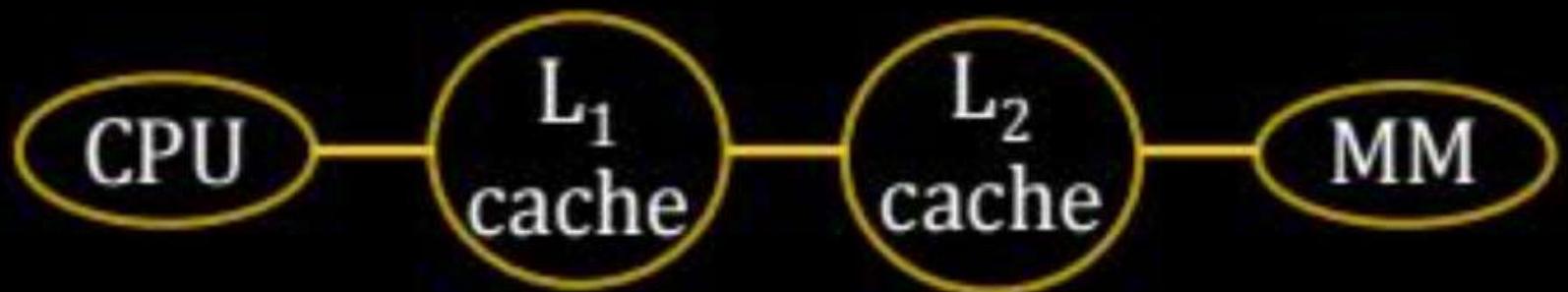
This miss will occur when the very first reference to the cache itself a miss.

- 2) Capacity Miss - This miss will occur when cache is full.
- 3) Conflict Miss (Collision miss / reference miss)

This miss will occur when the too many blocks are placed into same cache line or same cache SET.

Multi level cache

- ❑ To reduce the miss penalty multi-level caches are used in the system design.
- ❑ The number of cycles required to transfer the data from higher levels to L_1 due to miss operation is called as miss penalty



Local Miss Rate [LMR]

Global Miss Rate [GMR]

$$\text{Local Miss Rate} = \frac{\text{\#misses in the cache}}{\text{\# accesses to that cache}}$$

$$\text{Global Miss Rate} = \frac{\text{\#misses in the cache}}{\text{Total \#CPU generated reference}}$$

Average access time of the memory is always calculated in terms of Hit time, miss rate and miss penalty is as follows:

$$T_{avg} = \text{Hit time } L_1 + (\text{Miss Rate } L_1 * \text{Miss penalty } L_1)$$

$$\text{Miss penalty } L_1 = \text{Hit time } L_2 + (\text{Miss rate } L_2 * \text{Miss penalty } L_2)$$

$$\text{Miss penalty } L_2 = \text{MM Access Time}$$

Types of Misses

In the CM design 3 types of misses are present.

- 1) Compulsory miss - (Cold start miss / first reference miss)

This miss will occur when the very first reference to the cache itself a miss.

- 2) Capacity Miss - This miss will occur when cache is full.
- 3) Conflict Miss (Collision miss / reference miss)

This miss will occur when the too many blocks are placed into same cache line or same cache SET.

- Q.1** Consider a Direct Mapping if the size of Cache memory is 512KB & Main Memory 512 KB & Cache line size (Block) is 64KB the calculate the number of bit required for
- (i) P.A
 - (ii) TAG
 - (iii) B.O
 - (iv) W.O
 - (v) #LINES
 - (vi) TAG Memory Size.

Q.2

Consider a Direct Mapping, Cache size = 64 byte, Line Size = 8

Byte. MM = 256 Byte then #bits for P.A, TAG, L.O, W.O Tag
memory size?

Q.3 Consider a Direct Mapping, Cache size = 128 KB, Line Size = 64

Byte. Main Memory is 1MB then what is the line number of physical address $(ABCDE)_{16}$?

Q.4

Consider a 2-way set associative if the size of cache memory is 512KB & Main Memory 512MB & Cache line size is 64KB then calculate the Number of bit Required for

Q.5

Consider a 2-way set associative Cache Size = 256 KB, Line size = $\frac{P}{W}$,
32 Byte, MM = 1MB, then what is the set number of Physical
address $(ABCDE)_{16}$?

An 8-way set associative cache of size 64 KB (1 KB = 1024 bytes) is used in a system with 32-bit address. The address is sub-divided into TAG, INDEX, and BLOCK OFFSET.

The number of bits in the TAG is _____.

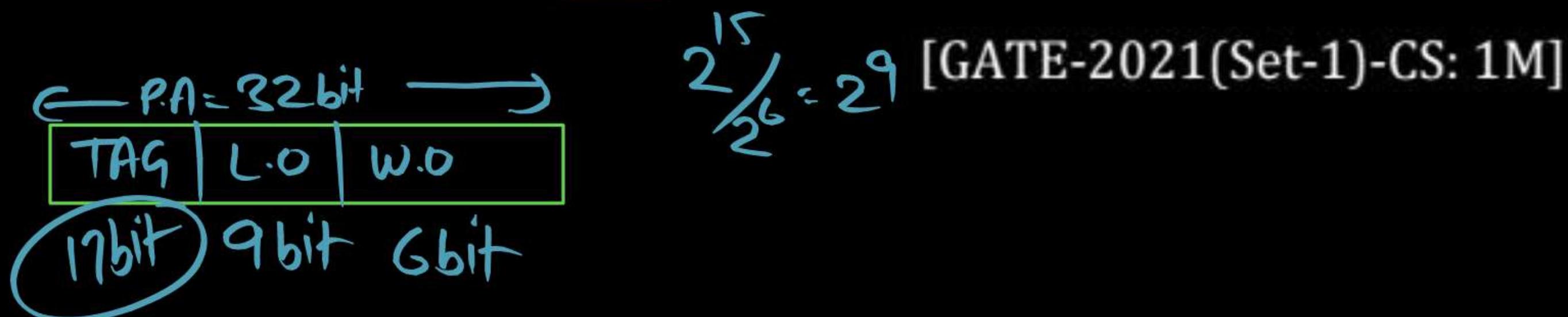
[GATE-2023-CS: 2M]

Consider a computer system with a byte-addressable primary memory of size 2^{32} bytes. Assume the computer system has a direct-mapped cache of size 32 KB ($1 \text{ KB} = 2^{10}$ bytes), and each cache block is of size 64 bytes.

The size of the tag field is 17 bits.

Ans

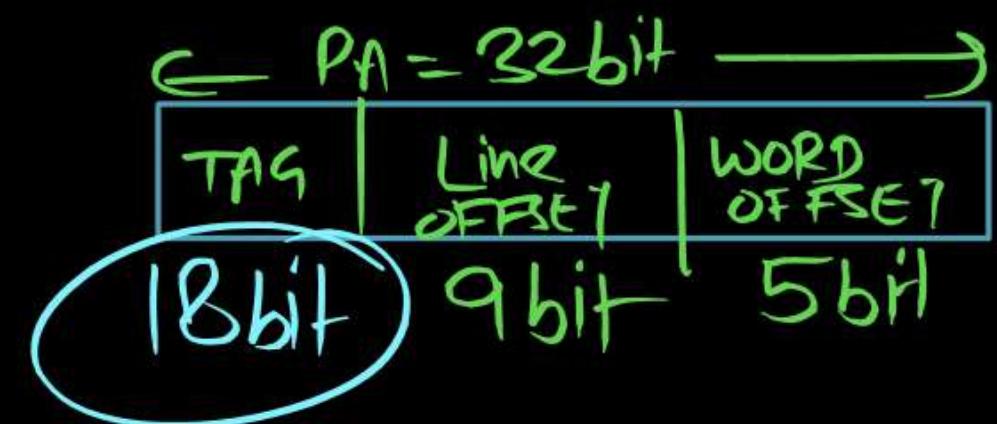
6



Consider a machine with a byte addressable main memory of 2^{32} bytes divided into blocks of size 32 bytes. Assume that a direct mapped cache having 512 cache lines is used with this machine. The size of the tag field in bits is ____.

\downarrow
 2^9

[GATE-2017(Set-1)-CS: 2M]



$$\text{Tag} = 32 - (9 + 5)$$

$$\text{Tag} = 18 \text{ bit}$$

MCQ

Q.4

Consider a machine with a byte addressable main memory of 2^{20} bytes, block size of 16 bytes and a direct mapped cache having 2^{12} cache lines. Let the addresses of two consecutive bytes in main memory be $(E201F)_{16}$ and $(E2020)_{16}$. What are the tag and cache line address (in hex) for main memory address $(E201F)_{16}$?

[GATE-2015(Set-3)-CS: 1M]

A E, 201

Ans(A)

B F, 201

PA = 20bit

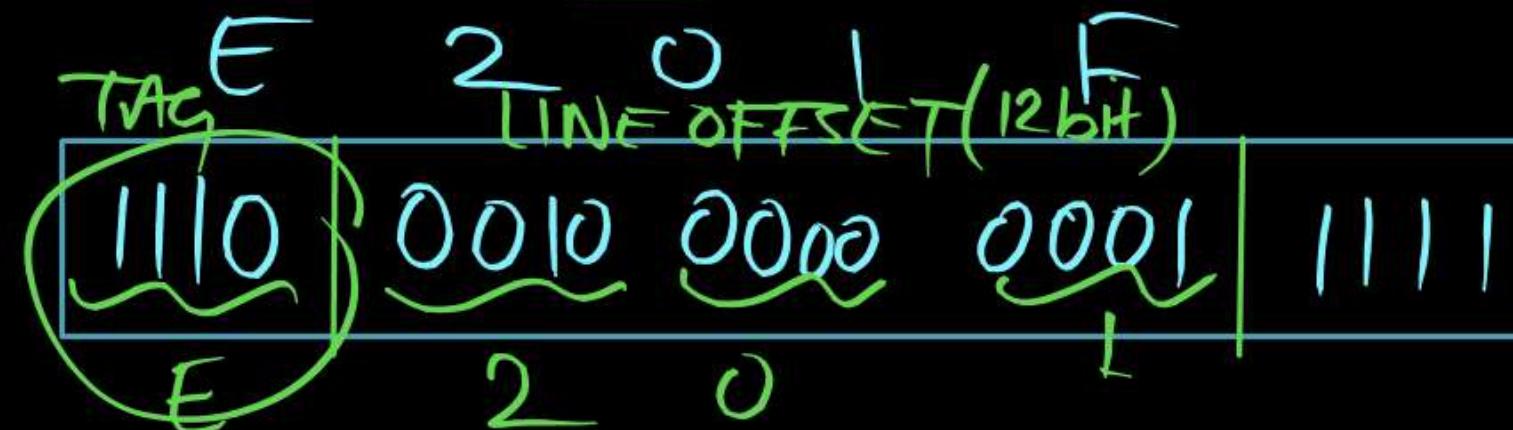
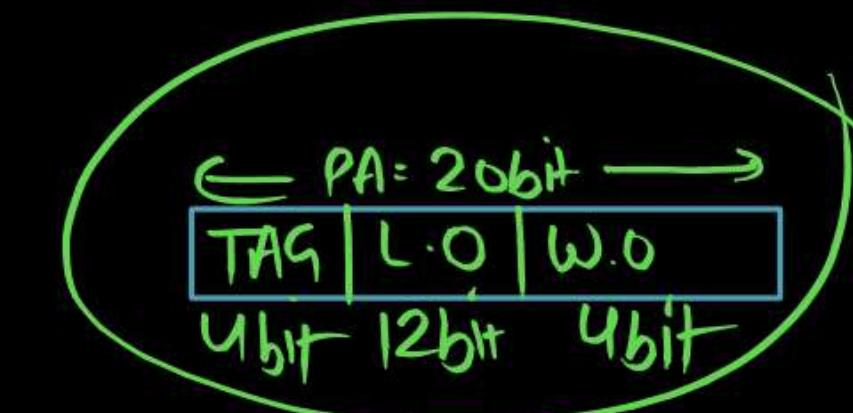
C E, E20

W.D = 4bits

TAG: 4bit

D 2, 01F

L.O = 12bit



Q.5

[Common Data for this and next question]

u marks

A computer has a 256 Kbyte, 4-way set associative. Write back data cache with block size of 32 Bytes. The processor sends 32 bit address to the cache controller. Each cache tag directory entry contains, in addition to address tag, 2 valid bits, 1 modified bit and 1 replacement bit.

2^5

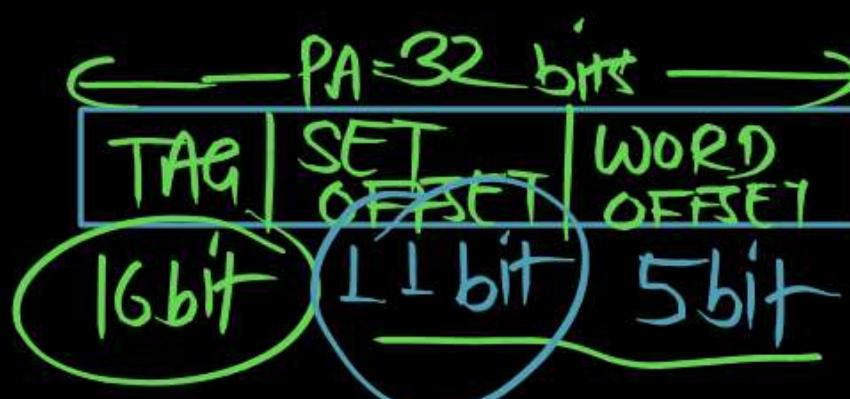
The number of bits in the tag field of an address is

(a) 11

(b) 14

(c) 16

(d) 27



$$\# \text{LINE} = \frac{\text{CMSize}}{\text{BlockSize}}$$

[GATE - 2012: 2 Marks]

$$= \frac{256 \text{ KB}}{32 \text{ B}} = \frac{2^8}{2^5} = 2^{13} \text{ LINES}$$

$$\# \text{SET} = \frac{2^{13}}{2^2} = 2^{11} \text{ SET}$$

S.O = 11 bit

Q.6

[Common Data from previous question]

A computer has a 256 Kbyte, 4-way set associative. Write back data cache with block size of 32 Bytes. The processor sends 32 bit address to the cache controller. Each cache tag directory entry contains, in addition to address tag, 2 valid bits, 1 modified bit and 1 replacement bit.

[GATE - 2012: 2 Marks]

The size of the cache tag directory is

- (a) 160 Kbits
- (b) 136 Kbits
- (c) 40 Kbits
- (d) 32 Kbits

$$\frac{\text{Tag entry}}{\text{Size}} = \frac{\text{Tag bit}}{\text{bit}} + \frac{\text{extra}}{\text{bit}} = 16 + 2 + 1 + 1 = 20 \text{ bit}$$

$$\frac{\text{Tag Memory}}{\text{Size}} = \# \text{LINES} \times \frac{\text{Tag Entry}}{\text{Size}}$$

$$\Rightarrow 2^{13} \times 20 \text{ bits}$$

$$\Rightarrow 8 \times 20 \times 2^{10} \text{ bit}$$

$$= 160 \text{ Kbits}$$

4 way Set Assoc.

$$\frac{\text{Tag Memory}}{\text{Size}} = \frac{\# \text{SET}}{\text{Set}} \times \frac{\text{LINES per SET}}{\text{Set}} \times \frac{\text{Tag entry}}{\text{Set}}$$

$$= 2^11 \times 4 \times 20 \text{ bit}$$

$$\Rightarrow 2 \times 4 \times 20 \times 2^{10} \text{ bit}$$

$$= 160 \text{ Kbits} \quad \underline{\text{Ans}}$$

An 8KB direct-mapped write back cache is organized as multiple blocks, each of size 32 bytes. The processor generates 32-bit addresses. The cache controller maintains the tag information for each cache block comprising of the following.

1Valid bit

1Modified bit

As many bits as the minimum needed to identify the memory block mapped in the cache.

What is the total size of memory needed at the cache controller to store meta-data (tags) for the cache?

[GATE-2011-CS: 2M]

A 4864 bits

B 6144 bits

C 6656 bits

D 5376 bits

Physical address = 32 bit

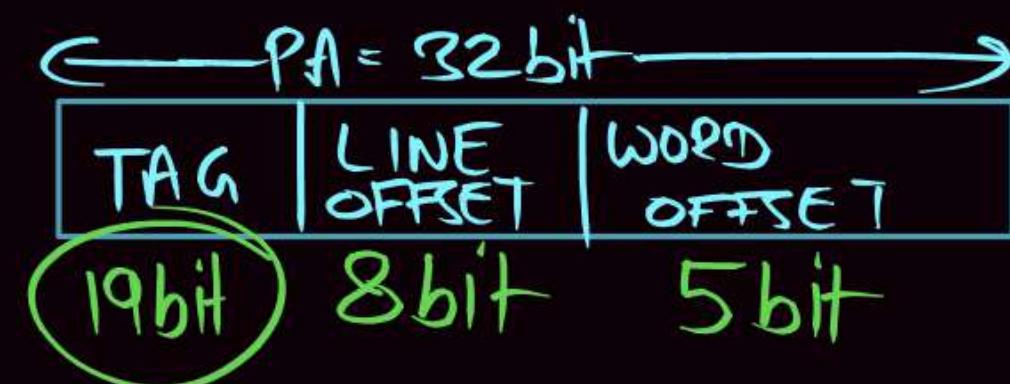
CM Size = 8KB = 2^{13} Byte

Block Size = 32B = 2^5 Byte

Direct Mapped Cache

$$\# \text{LINES} = \frac{2^{13}}{2^5} = 2^8 \text{ Lines}$$

L.O = 8 bit



$$\begin{aligned}\text{Tagentry Size} &= \text{Tag bit} + 1 + 1 \\ &= 19 + 1 + 1\end{aligned}$$

$$\text{Tagentry} = 21 \text{ bit}$$

$$\begin{aligned}\text{Tag Memory Array Size} &= \# \text{LINES} \times \text{Tagentry Size} \\ &= 2^8 \times 21 \\ &= 256 \times 21 \\ &= 5376 \text{ bits } \underline{\text{Ans}}\end{aligned}$$

Q.8**Q.9**

Common Data for next two questions:

Consider a computer with a 4-ways set-associative mapped cache of the following characteristics: a total of 1 MB of main memory, a word size of 1 byte, a block size of 128 words and a cache size of 8 KB.

While accessing the memory location 0C795H by the CPU, the contents of the TAG field of the corresponding cache line is

[GATE-2008-CS: 2M]

A 000011000

B 110001111

C 00011000

D 110010101

0 C 7 9 S

The number of bits in the TAG, SET and WORD fields, respectively are:

[GATE-2008-CS: 2M]

A 7, 6, 7

B 8, 5, 7

C 8, 6, 6

D 9, 4, 7

Q.10

Consider a 4-way set associative cache consisting of 128 lines with a line size of 64 words. The CPU generates a 20-bit address of a word in main memory. The number of bits in the TAG, SET and WORD fields are respectively.

P
W

- (a) 9, 6, 5
- (b) 7, 7, 6
- (c) 7, 5, 8
- (d) 9, 5, 6

[GATE - 2007]

A certain processor uses a fully associative cache of size 16 kB. The cache block size is 16 bytes. Assume that the main memory is byte addressable and uses a 32-bit address. How many bits are required for the Tag and the Index fields respectively in the addresses generated by the processor?

[GATE-2019-CS: 1M]

A 24-bits and 0-bits

B 28-bits and 4-bits

C 24-bits and 4-bits

D 28-bits and 0-bits

The width of the physical address on a machine is 40 bits. The width of the tag field In a 512 KB 8-way set associative cache is _____ bits.

[GATE-2016(Set-2)-CS: 2M]

A 4-way set-associative cache memory unit with a capacity of 16 KB is built using a block size of 8 words. The word length is 32 bits ^{4B/w}. The size of the physical address space is 4 GB. The number of bits for the TAG field is ____.

[GATE-2014(Set-2)-CS: 1M]

A cache memory unit with capacity of N words and block size of B Words is to be designed. If it is designed as a direct mapped cache, the length of the TAG field is 10 bits. If the cache unit is now designed as a 16-way set-associative cache, the length of the TAG field is _____ bits.

[GATE-2017(Set-1)-CS: 2M]

Q.15

The main memory of a computer has 2^m blocks while the cache has 2^c blocks. If the cache uses the set associative mapping scheme with 2 blocks per set, then block k of the main memory maps to the set.

[GATE - 1999]

- (a) $(k \bmod m)$ of the cache
- (b) $(k \bmod c)$ of the cache
- (c) $(k \bmod 2^c)$ of the cache
- (d) $(k \bmod 2^m)$ of the cache

The size of the physical address space of a processor is 2^P bytes. The word length is 2^W bytes. The capacity of cache memory is 2^N bytes. The size of each cache block is 2^M words. For a K-way set-associative cache memory, the length (in number of bits) of the tag field is

[GATE-2018-CS: 2M]

A $P - N - \log_2 K$

B $P - N + \log_2 K$

C $P - N - M - W - \log_2 K$

D $P - N - M - W + \log_2 K$

A computer system with a word length of 32 bits has a 16 MB byte-addressable main memory and a 64 KB, 4-way set associative cache memory with a block size of 256 bytes. Consider the following four physical addresses represented in hexadecimal notation.

$A_1 = 0 \times 42C8A4, A_2 = 0 \times 546888, A_3 = 0 \times 6A289C, A_4 = 0 \times 5E4880$

Which one of the following is TRUE?

[GATE-2020-CS: 2M]

- A A1 and A3 are mapped to the same cache set.
- B A2 and A3 are mapped to the same cache set.
- C A3 and A4 are mapped to the same cache set.
- D A1 and A4 are mapped to different cache sets.

Consider a set-associative cache of size 2 kb ($1\text{ KB} = 2^{10}\text{ bytes}$) with cache block size of 64 bytes. Assume that the cache is byte - addressable and a 32-bit address is used for accessing the cache. If the width of the tag field is 22 bits, the associativity of the cache is ____.

[GATE-2021(set-2)-CS: 1M]

Consider a 4-way set associative cache (initially empty) with total 16 cache blocks. The main memory consists of 256 blocks and the request for memory blocks is in the following order:

0, 255, 1, 4, 3, 8, 133, 159, 216, 129, 63, 8, 48, 32, 73, 92, 155

$k \bmod 16$

Which one of the following memory block will NOT be in cache if LRU replacement policy is used?

[GATE-2009-CS: 2M]

A 3

B 8

C 129

D 216

Consider a 2-way set associative cache with 256 blocks and uses LRU replacement. Initially the cache is empty. Conflict misses are those misses which occur due to contention of multiple blocks for the same cache set. Compulsory misses occur due to first time access to the block. The following sequence of accesses to memory blocks (0, 128, 256, 128, 0, 128, 256, 128, 1, 129, 257, 129, 1, 129, 257, 129) is repeated 10 times. The number of conflict misses experienced by the cache is _____.

[GATE-2017(Set-1)-CS: 2M]

— : Compulsory

— : Conflict

If the associativity of a processor cache is doubled while keeping the capacity and block size unchanged, which one of the following is guaranteed to be NOT affected?

[GATE-2014(Set-2)-CS: 2M]

- A Width of tag comparator
- B Width of set index decoder
- C Width of way selection multiplexer
- D Width of processor to main memory data bus

Consider a two-level cache hierarchy with L_1 and L_2 caches. An application incurs 1.4 memory accesses per instruction on average. For this application, the miss rate of L_1 cache is 0.1; the L_2 cache experiences on average, 7 misses per 1000 instructions. The miss rate of L_2 expressed correct to two decimal places is _____.

[GATE-2017(Set-1)-CS: 1M]

In a two-level cache system, the access times of L_1 and L_2 caches are 1 and 8 clock cycles, respectively. The miss penalty from the L_2 cache to main memory is 18 clock cycles. The miss rate of L_1 cache is twice that of L_2 . The average memory access time (AMAT) of this cache system is 2 cycles. The miss rates of L_1 and L_2 respectively are:

[GATE-2017(Set-2)-CS: 2M]

- A 0.111 and 0.056
- B 0.056 and 0.111
- C 0.0892 and 0.1784
- D 0.1784 and 0.0892

Q.24



Common Data for next two questions:

Consider a machine a 2-way set associative data cache of size 64Kbytes and block size 16 bytes. The cache is managed using 32 bit virtual addresses and the page size is 4 Kbytes. A program to be run on this machine begins as follows:

Double ARR [1024] [1024]

Int i, j;

```
/* Initialize array ARR to 0.0 */  
for (i = 0; i < 1024; i++)  
for (j = 0; j < 1024; j++)  
ARR [i] [j] = 0.0;
```

The size of double 8 bytes. Array ARR is in memory starting at the beginning of virtual page 0xFF000 and stored in row major order. The cache is initially empty and no pre-fetching is done. The only data memory references made by the program are those to array ARR.

MCQ

The total size of the tags in the cache directory is

[GATE-2008-CS: 2M]

- A 32 kbits
- B 34 kbits
- C 64 kbits
- D 68 kbits

Q.25

4 Marks

[Common Data for this and next question]

Consider two cache organization. The first one is 32 KB 2-way set associative with 32-byte block size. The second one is of the same size but direct mapped. The size of an address is 32 bits in both cases. A 2-to-1 multiplexer has latency of 0.6 ns while a k-bit comparator has a latency of $k/10$ ns. The hit latency of the set associative organization is h_1 while that of the direct mapped one is h_2 . The value of h_1 is

- (a) 2.4 ns (b) 2.3 ns (c) 1.8 ns (d) 1.7 ns

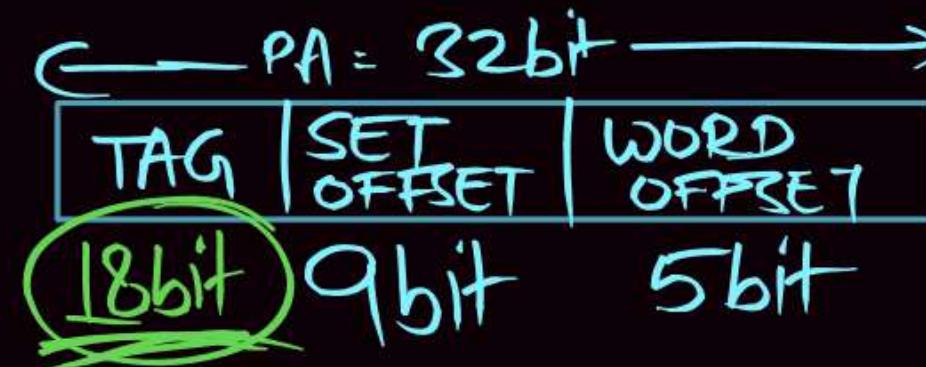
[GATE - 2006: 2 Marks]

P
W

2 way Set Associative

Cache = 32KB, 2 way set ass.

Block size = 32B, P.A = 32 bit



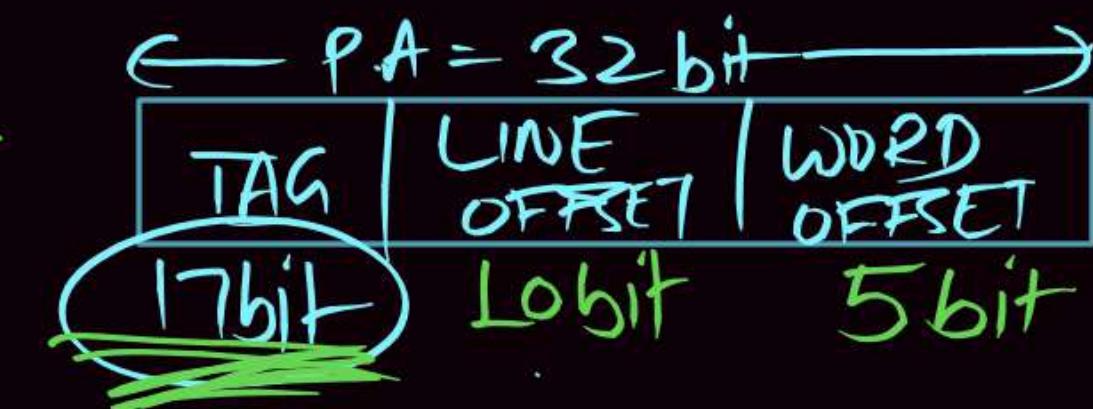
$$\#SET = \frac{\#LINES}{N\text{-Ways}} \Rightarrow \frac{2^{10}}{2^1} = 2^9 \text{ Set}$$

Direct Mapping

Cache Size = 32KB, Direct mapping

Block size = 32B (2^5) P.A = 32 bit

Direct



$$\#LINE = \frac{CMSize}{Blocksize}$$

$$\Rightarrow \frac{32KB}{32B} = \frac{2^{15}}{2^5}$$

$$L.O = 10bit \leftarrow 10$$

2 way Set Associative

$$\text{Hit Latency} = \frac{\text{Latency of Tag Combinator}}{10} + \frac{\text{Latency of Multiplexer}}{10}$$

$$= \frac{K}{10} + 0.6$$

$$= \frac{18}{10} + 0.6$$

$$= 1.8 + 0.6$$

$$h_L = 2.4 \text{ ns} \quad \text{Ans}$$

Direct Mapping

$$\text{Hit latency} = \text{Latency of Tag Combinator}$$

$$= \frac{K}{10} \quad (K: \# \text{ of Tag bits})$$

$$= \frac{17}{10}$$

$$h_2 = 1.7 \text{ nsec} \quad \text{Ans}$$

Q.26

[Common Data from previous question]

P
W

(Qn no 25)

Consider two cache organization. The first one is 32 KB 2-way set associative with 32-byte block size. The second one is of the same size but direct mapped. The size of an address is 32 bits in both cases. A 2-to-1 multiplexer has latency of 0.6 ns while a k-bit comparator has a latency of $k/10$ ns. The hit latency of the set associative organization is h_1 while that of the direct mapped one is h_2 . The value of h_2 is

- (a) 2.4 ns (b) 2.3 ns (c) 1.8 ns (d) 1.7 ns

[GATE - 2006: 2 Marks]

Q.27

A computer system has a level - 1 instruction cache (1-cache), a level-1 data cache(D-cache) and a level-2 cache(L2-cache) with the following specifications.

	Capacity	Mapping method	Block size
1-cache	4K words	Direct mapping	4 words
D-cache	4K words	2-way set associative mapping	4 words
L2-cache	64K words	4-way set associative mapping	16 words

Capacity mapping method block size
 1-cache 4K words direct mapping 4 words
 D-cache 4 k words 2 way set associative mapping 4 words L2-cache
 64K words 4-way set associative mapping 16 words. The length of the physical address of a word in the main memory is 30 bits. The capacity of the tag memory in the 1-cache, D-cache & L2-cache is. Respectively,

- (a) $1K \times 18\text{-bit}$, $1K \times 19\text{-bit}$, $4K \times 16\text{-bit}$
- (b) $1K \times 16\text{-bit}$, $1K \times 19\text{-bit}$, $4K \times 18\text{-bit}$
- (c) $1K \times 16\text{-bit}$, $512 \times 18\text{-bit}$, $1K \times 16\text{-bit}$
- (d) $1K \times 16\text{-bit}$, $512 \times 18\text{-bit}$, $1K \times 18\text{-bit}$

Q.28

Consider a small two-way set-associative cache memory, consisting of 4 blocks. For choosing the block to be replaced, use the least recently used (LRU) scheme. The number of cache misses for the following sequence of block addresses is 8, 12, 0, 12, 8

- (a) 2
- (b) 3
- (c) 4
- (d) 5

[GATE - 2004]

Q.29

Consider a system with 2 KB direct mapped data cache with a block size of 64 bytes. The system has a physical address space of 64 KB and a word length of 16 bits. During the execution of a program, four data words P, Q, R and S are accessed in that order 10 times (i.e., PQRSPQRS....) Hence, there are 40 accesses to data cache altogether. Assume that the data cache is initially empty and no other data words are accessed by the program. The addresses of the first bytes of P, Q, R and S are 0xA248, 0xC28A, 0xCA8A and 0xA262, respectively. For the execution of the above program, which of the following statements is/are TRUE with respect to the data cache?

[2022: MSQ 2M]

- A** Every access to S is a hit.
- B** Once P is brought to the cache it is never evicted.
- C** At the end of the execution only R and S reside in the cache.
- D** Every access to R evicts Q from the cache.

**THANK
YOU!**

