

## Write operation

Write operation uses the same addressing technique as read operation. The tag part and data is presented to cache for updation. Cache write operation in variable follows a cache miss during read operation.

## Address Mapping

When a tag address is presented to the cache, it must be quickly compared to the stored tags to determine whether a matching tag is currently assigned to the cache. The obvious approach of scanning all the tags in sequence is unacceptably slow. There are faster techniques for the same. They are,

- ① Direct mapping cache
- ② Fully associative cache
- ③ set associative cache.

Direct mapping cache: In this mapping each block of memory is mapped to a fixed slot of cache. For example if cache has four slots (sets) then the memory block 0 or 4 or 8 or 12 can be found in slot 0 (set 0), while 1 or 5 or 9 or 13 --- can be found in slot 1; 2 or 6 or 10 or 14 --- in slot 2 and 3 or 7 or 11 or 15 --- in slot 3. This can be mathematically defined as.

(20)

cache slot(set) number =  $\left[ \frac{\text{Block number of main memory}}{\text{Total number of slots in cache}} \right] \bmod \text{modulo operator}$

$\left[ \frac{\text{Total number of slots in cache}}{\text{Number of slots in cache}} \right]$

In this technique, it can easily be determined whether a memory block is in cache or not As shown in figure.

Number of main memory blocks = 16 (B<sub>0</sub> to B<sub>15</sub>)

Number of slots(sets) in cache memory = 4

Mapping of main memory block:

Suppose we want to find the cache memory set, in which the block number B<sub>11</sub> should be stored. This can be done by dividing 11 (Block Number) by 4 (Number of sets in cache memory).

$$\frac{11}{4} = 2 \frac{3}{4}, \text{ Quotient} = 2, \text{ gives group number or tag value}$$

Remainder = 3, gives the set number of cache data memory.

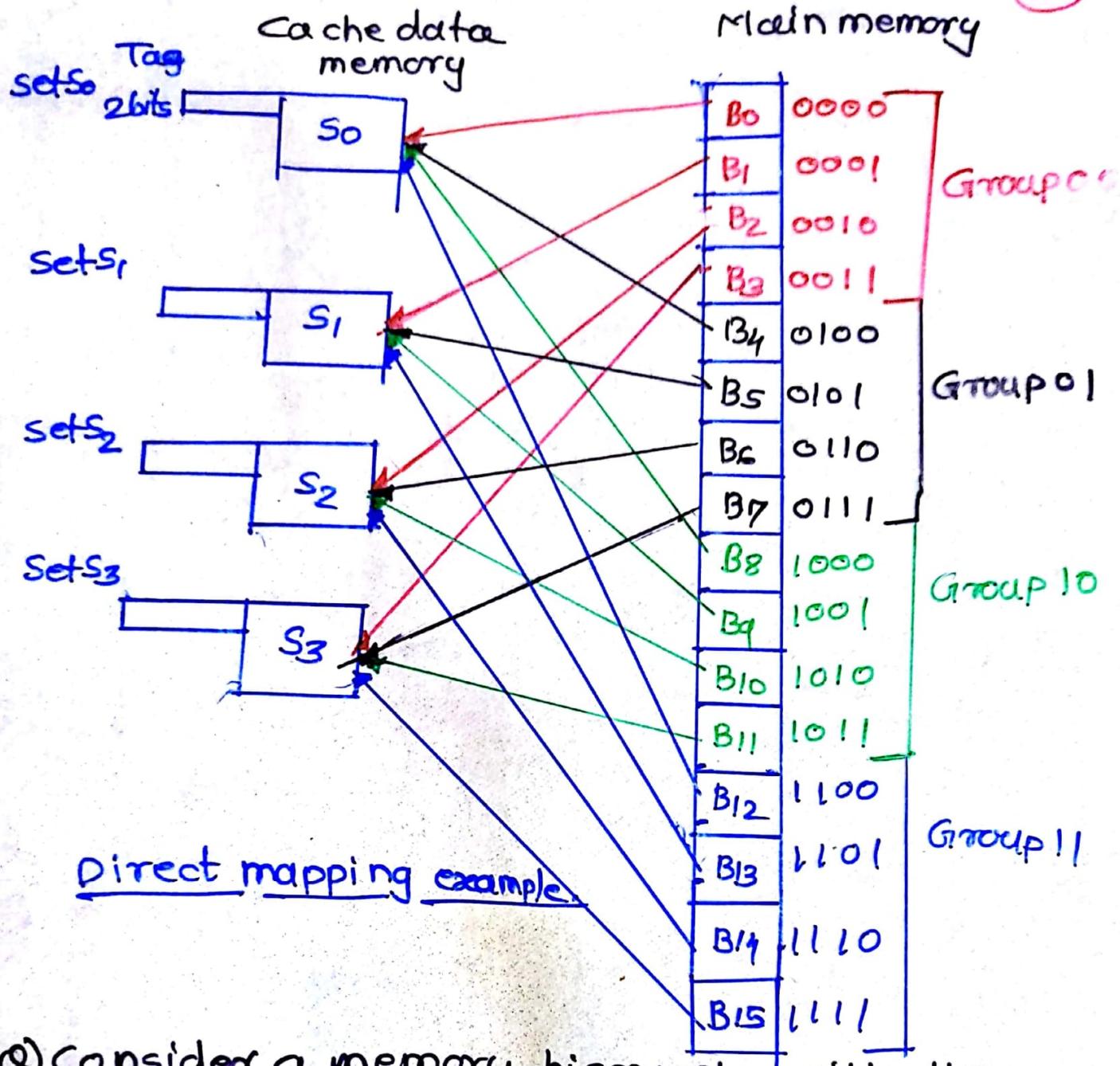
If the block is represented in binary then we can easily find the quotient and remainder.

$$(11)_{10} = 10 \underbrace{11}_{\substack{\text{Quotient} \\ (\text{Tag value})}}$$

$\underbrace{11}_{\text{Remainder}}$

$\underbrace{11}_{(\text{set number of cache})}$

(2f)



Q) Consider a memory hierarchy with the following details.

Number of main memory blocks = 16  
 $\downarrow$  ————— cache —————  $\downarrow$  = 4

Block size = 4

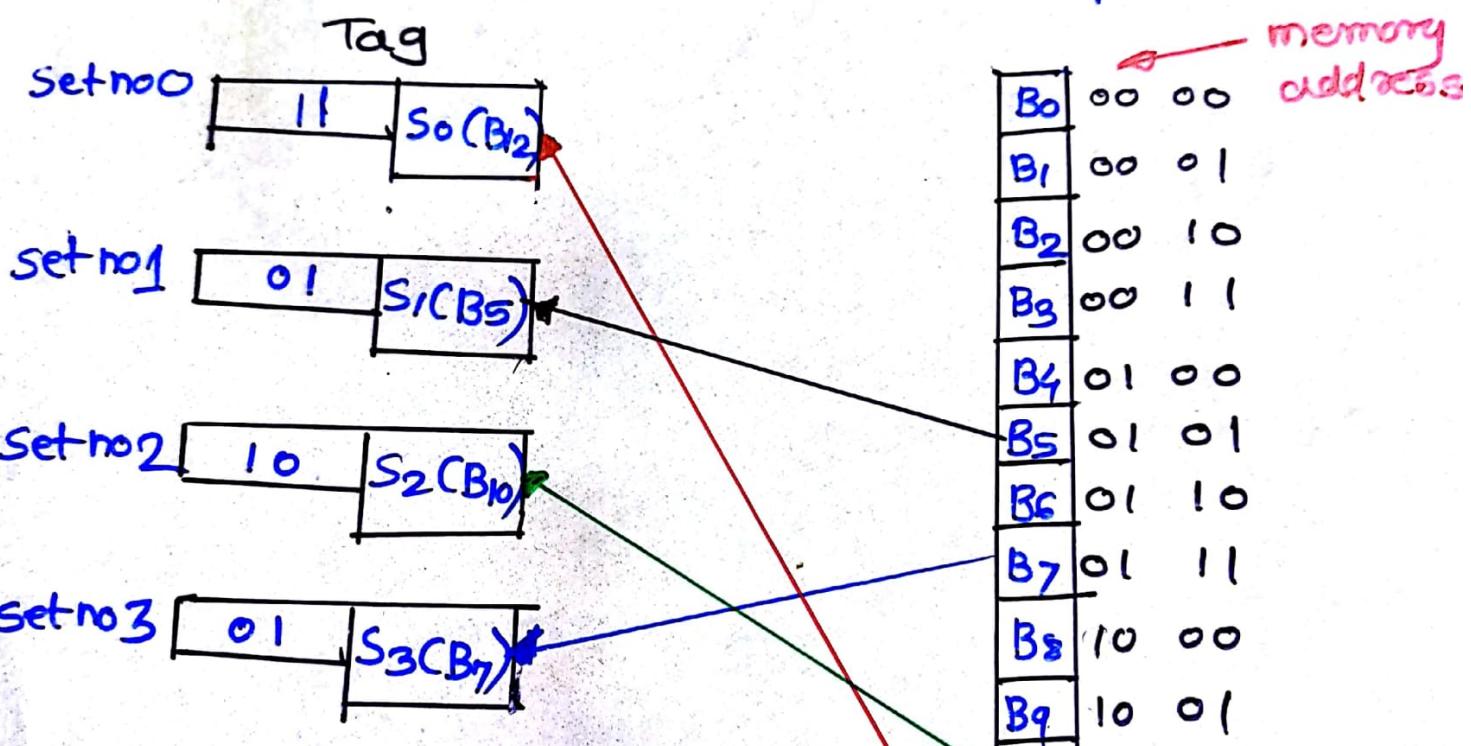
Direct mapping cache

Show the contents of tags if the following blocks are stored in cache memory  $B_{12}, B_7, B_5, B_{10}$

solution: We must use 4 bit address for a block as there are 16 blocks in memory.

Block no	Block address in memory	Tag	Set no
12	11 00	11	00
7	01 11	01	11
5	01 01	01	01
10	10 10	10	10

Tag value      Remainder (set number)



Direct mapping example with value of tag.

Tag      set no  
Main memory

consider a cache consisting of 256 blocks of  $2^4$  words each for a total of 4096(4k) words and assume that the main memory is addressable by 16 bit address and it consists of 4k blocks. How many bits are there in each of the TAG, BLOCK/SET and WORD fields for direct mapping cache.

### Solution

i) Main memory size =  $2^{16}$  bytes [Main memory is addressable by a 16 bit address]

ii) Cache memory size =  $256 \times 16 = 2^8 \times 2^4 = 2^{12}$

$\frac{\text{No of blocks}}{\text{Cache memory size}}$        $\frac{\text{No of words}}{\text{in a block.}}$

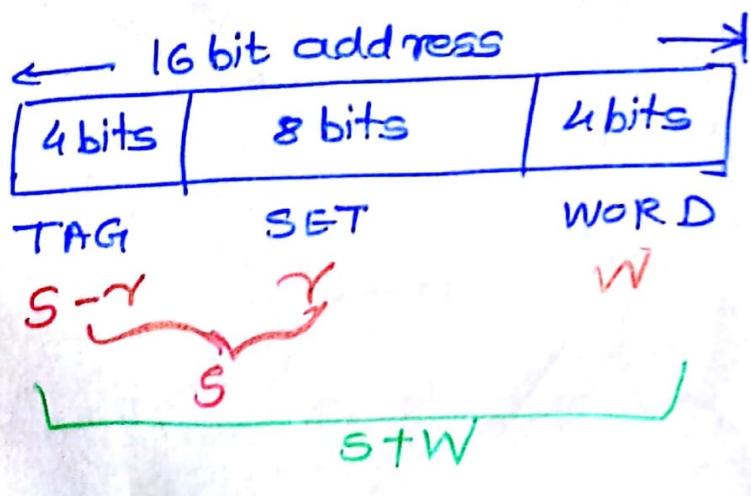
Main memory size       $= \frac{2^{16}}{2^{12}} = 2^4$ , therefore the  
cache memory size      size of TAG in bits = 4

since the block containing  $16(2^4)$  words - 4 bits are required to address a word.

Since the cache contains  $256(2^8)$  blocks, 8 bits are required to address a set.

Main memory address

5 bits specifying  $2^5$  blocks in main memory (block address in main memory)



Address length =  $(S+W)$  bits.

(24)

Number of addressable units =  $2^{S+W}$  words or

Block size = line size =  $2^W$  words or bytes.

Number of blocks in main memory =  $\frac{2^{S+W}}{2^W} = 2^S$

Number of lines in cache =  $m = 2^r / 2^w$

size of tag =  $(S-r)$  bits.

$S$  bits specify  $2^S$  blocks of main memo

① Main memory = 16MB =  $2^4 \times 2^{20}$  = 24 address lines  
 cache memory = 64KB =  $2^6 \times 2^{10}$  = 16 bits (24 bits)

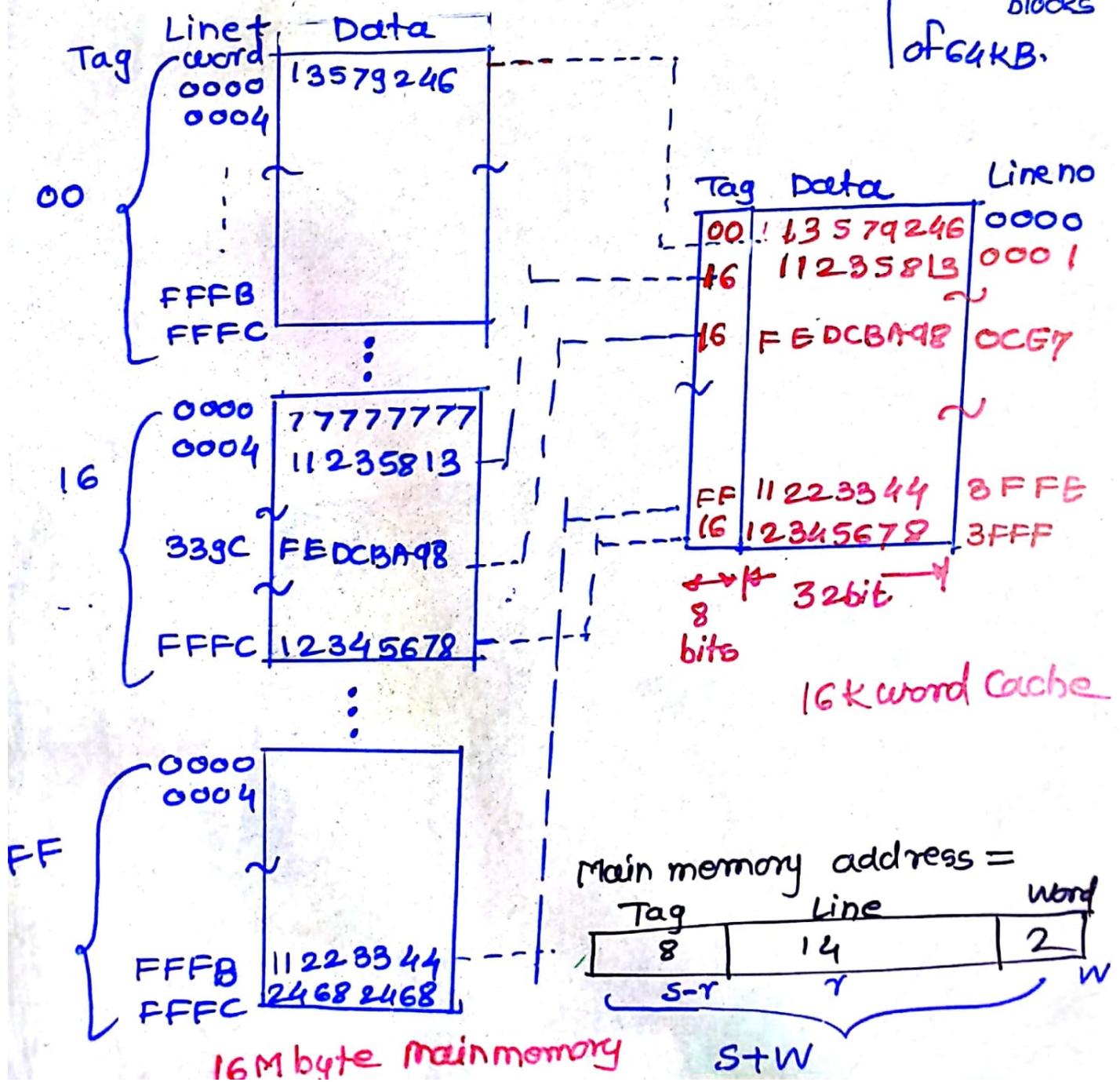
No of words in each block = 04 = 2 bits required to select a word

$$\frac{\text{Main memory size}}{\text{cache memory size}} = \frac{2^{24}}{2^8} = 2^8$$

∴ Tag size is of 8 bits.

To represent set/line = 14 bits.

Total memory is divided into  $\frac{16\text{MB}}{64\text{KB}} = 256$  blocks of 64KB.



## Advantages of direct mapping

- ① Simple to implement
- ② Required normal SRAM

## Drawback of direct mapping

If a program happens to reference words repeatedly from two different blocks that map into the same line, then the blocks will be continually swapped in the cache, and the hit ratio will be low (a phenomenon known as thrashing).

## Associative Mapping

Associative mapping overcomes the disadvantage of direct mapping by permitting each memory block to be decided <sup>main</sup> into any line/set of the cache. In this case the cache control logic interprets a memory address simply as a tag and a word field.

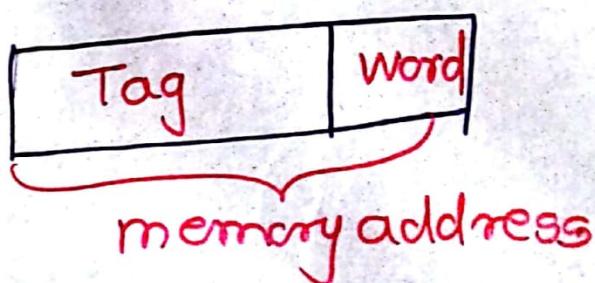
The tag field uniquely identifies a block of main memory. To determine whether the block is in the cache, the cache control logic must simultaneously examine every line's tag for a match.

Note that no field in the address corresponds to line number so that the number of lines in the cache is not determined by the address format.

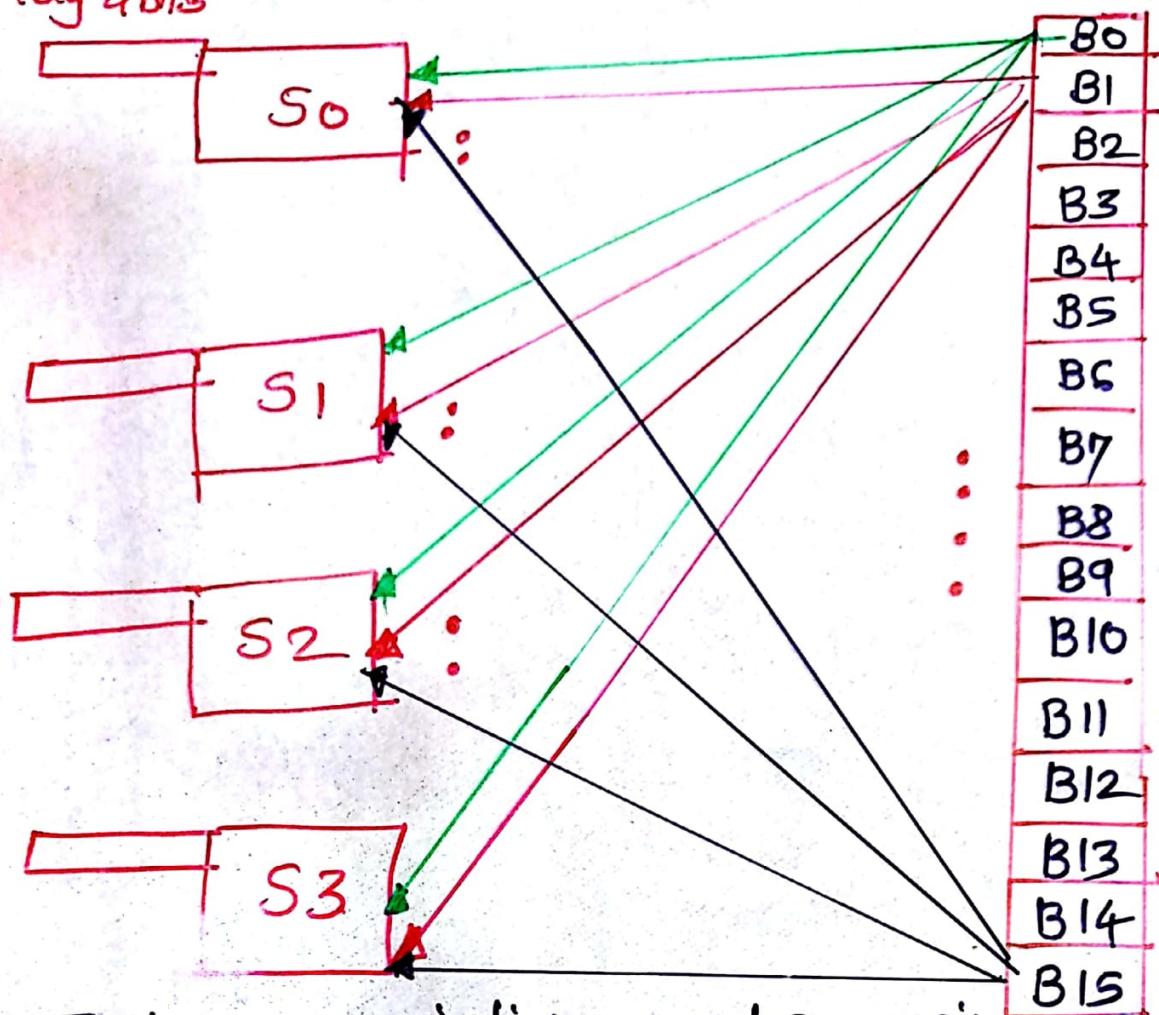
(27)

At the start of memory access, the incoming tag is compared simultaneously to all the tags in the cache's tag memory. If a match (cache hit) occurs, a match indicating signal triggers the cache to service the requested memory access. A no-match signal identifies a cache miss and the memory access requested is forwarded to main memory for service. A cache block containing the target address is then sent from main memory to the cache.

Each block of main memory can be placed in any of the available cache block frame (slot, set). Because of this, an  $s$ -bit tag is needed in each cache block, where  $s$  is the number of bits required to address a main memory block.



Tag 4 bits



Fully associative cache mapping,

Every block is mapped to any of the four block frames identified by the tag.

- Advantages:
- 1) Very high hit Ratio.
  - 2) Allows better block replacement strategy with reduced contention

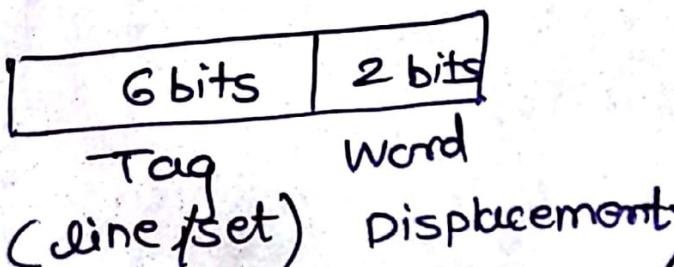
- Disadvantages:
- 1) Associative memory is very costly
  - 2) Significant increase in tag length
  - 3) Scheme is comparatively complex.

(29)

Consider a cache consisting of 16 words. Each block consists of 4 words. Size of main memory is 256 bytes. Find number of bits in each of TAG, Block/SET and WORD for associative mapped cache!

$$\text{Main memory} = 256 \text{ bytes} = 2^8$$

$$\text{address bits} = 8$$



set 0	111100	8A
	8B	
	90	
	9A	
Set 1	101000	4A
	4B	
	4C	
Set 2	000000	5B
	A5	
	F1	
	B6	
	92	
Set 3	000001	81
	7B	
	6A	
	85	

Tag ↑

Cache  
data  
memory

A5	00000000
F1	00000001
B6	00000010
92	00000011
81	00000100
7B	00000101
6A	00000110
85	00000111
⋮	⋮
4A	10100000
4B	10100001
4C	10100010
5B	10100011
5B	10100100
59	10100101
85	10100110
43	10100111
⋮	⋮
8A	11110000
8B	11110001
90	11110010
9A	11110011
9B	11110100
9C	11110101
75	11110110
50	11110111

Block Data Main Memory

## Set Associative Cache

Set associative cache tries to take advantage of both direct mapping and associative mapping.

- ① Direct mapped cache requires normal SRAM.
- ② Associative mapping offers high hit-ratio.  
Set associative cache is implemented using normal SRAM and provides a high hit ratio.

If properly designed, this cache may offer the best performance-cost ratio.  
Most high performance computer system are based on this.

- ③ In a K associative cache, the m cache blocks are divided into  $\frac{m}{k}$  sets. ( $m = \text{no of frames in cache}$ )
- ④ Each set contains k blocks.
- ⑤ A memory block is mapped to a set using direct mapping.
- ⑥ If each set is identified by d bits then  $2^d = V$
- ⑦ If s bits are required to access a main memory block then tag length  
 $\boxed{\text{tag length} = s-d}$

- can be stored (3)
- ⑧ A block of main memory in any of the  $s$  frames inside the addressed set (associative mapping)

Example: Memory consists of  $n = 16$  blocks  
Two way set associative cache is used  
 $K=2$

$\therefore V = 4$  (Number of sets in cache)

$m = 8$  (Number of frames in memory)

$$\text{Tag size} = \log_2 \frac{\text{Number of blocks in main memory}}{\text{Number of sets in cache}}$$

$$= \log_2 \frac{16}{4} = \log_2 4 = \underline{2 \text{ bits.}}$$

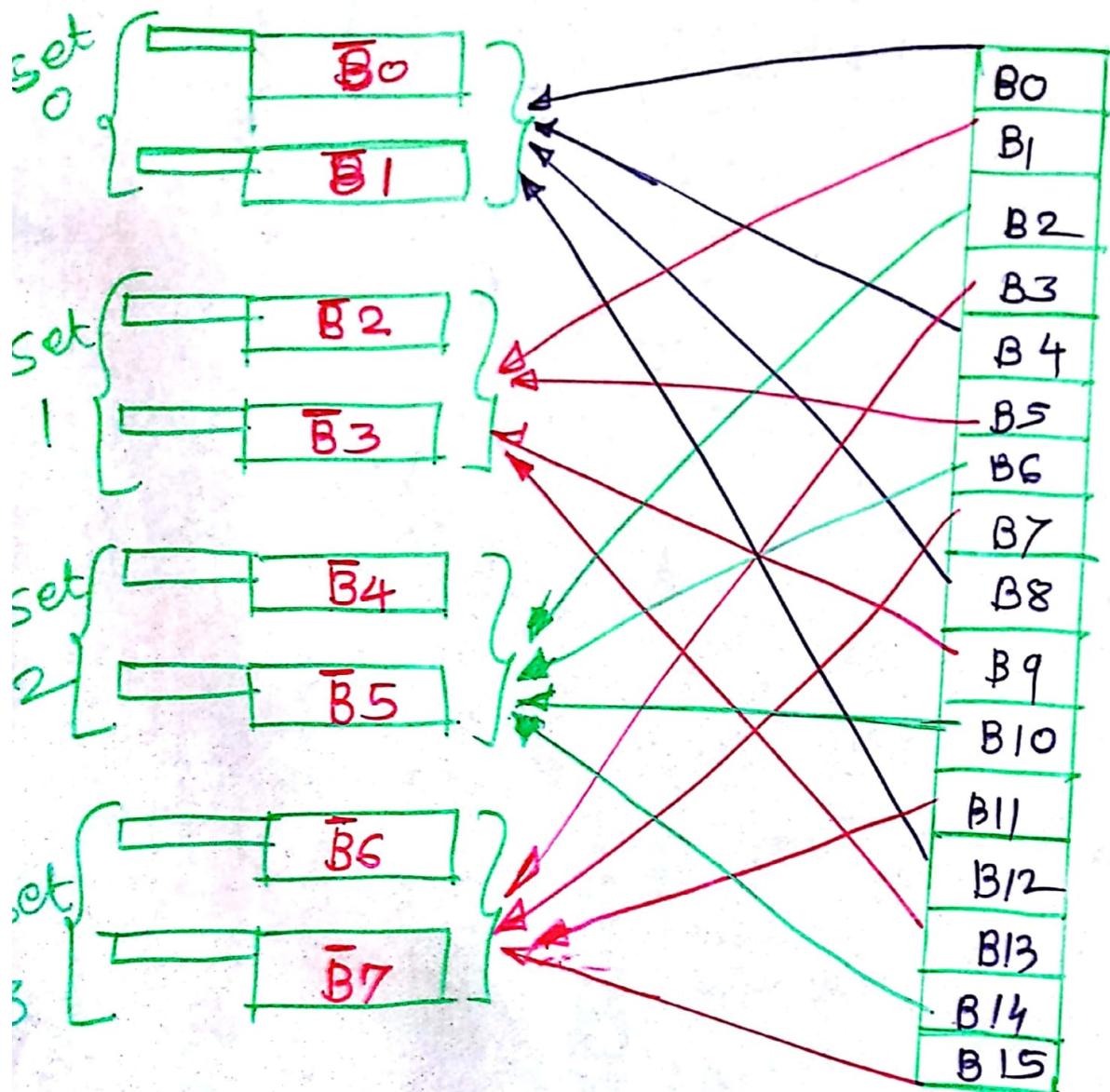
Set length ( $\in$  No of bits to address a set)

$$= \log_2 (\text{Number of sets in cache}) = \log_2 4$$

$$\text{word length} = \log_2 (\text{Number of words in a block}) = \underline{2 \text{ bits}}$$

$$= \log_2 4 = \underline{2 \text{ bits.}}$$

Tag	set	word
2 bit	2 bit	2 bit



Mapping cache blocks in a two-way associative cache with four sets.

(2)

Design a 2-way set associative memory with the following details.

Cache consists of 32 words.

Each block consists of 4 words.

Size of the main memory = 256 words

Find the number of bits in each of TAG, Block / SET and word.

### Solution

① No of bits in main memory address

$$= \log_2 (\text{size of main memory})$$

$$= \log_2 256 = \log_2 2^8 = \underline{\underline{8 \text{ bits}}}$$

② No of bits required to address

main memory blocks =  $\log_2 (\text{No of blocks in main memory})$

$$= \log_2 \left( \frac{256}{4} \right) = \log_2 64$$

$$= \underline{\underline{6 \text{ bits}}} = \underline{\underline{S}}$$

③ Number of sets in cache memory

$$N = \frac{32}{4} = 4$$

~~words~~  $\rightarrow 4 \times 2$  per block  $\uparrow$  2 way

④ No of bits required to address a set

$$= \log_2 4 = \underline{\underline{2}} = \underline{\underline{d}}$$

⑤ Tag Size = S - d = 6 - 2 = 4

4 bit Tag	2 bits set	n 2 bits word
--------------	---------------	------------------

## Comparision b/w associative & set <sup>(Q4)</sup>

- Set associative mapping is a compromise b/w direct and associative mapping
- Set associative mapping preserves the strength of both associative and direct mapping. It also reduces their disadvantages.
- Associative mapping requires content addressable memory which is very costly.
- Set associative mapping can be implemented using normal SRAM.
- To reduce the contention among the blocks, getting mapped to same set of cache. A set consists of  $k$  blocks in a  $k$ -way set associative mapping. Inside a set, a mapped block can be stored in any of the frames of the set.
- Mapping inside the set in a cache can be considered as associative mapping.
- Associative mapping offers high hit ratio. Set - II - is a compromise b/w the hit ratio of direct and associative mapping.
- cost of set associative mapping  $\approx$  cost of direct mapping
- hit ratio of set associative mapping  $\approx$  hit ratio of associative mapping

