

\* Illusion of unlimited fast memory.

eg: of library:

if we want to read a particular topic, collect related books, bring it to desk and start reading.

having several books on desk - saves time compared to having only one book & constantly having to go back to shelf to take another.

~ a pgm does not access all of its code or data at once.

Principle of locality - pgm accesses a relatively small portion of their address space at any instant of time.

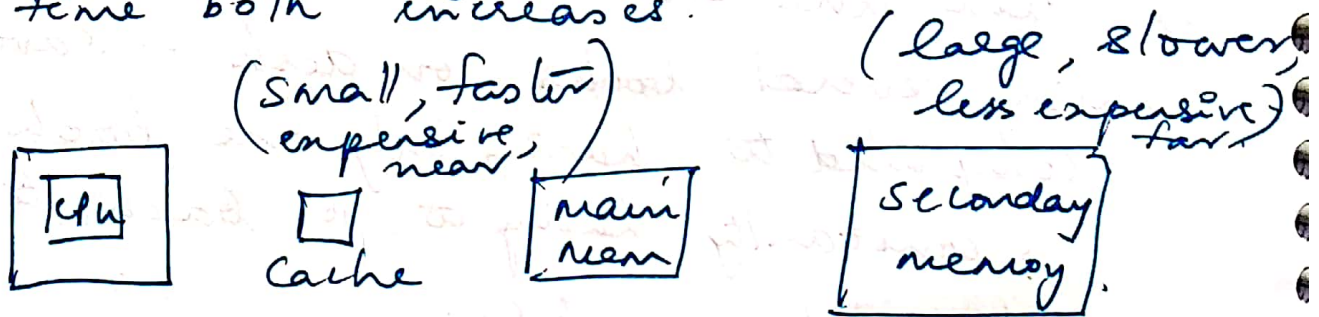
Two types: (in time)

① temporal locality - if an item is referenced, it will tend to be referred again soon. eg: loops.

(in space)  
② spatial locality - if an item is referenced, items whose addresses are close by will tend to be referenced soon. eg: arrays

By taking advantage of principle of locality, memory of a computer is implemented as a mem hierarchy.

— a structure that uses multiple levels of memories, as distance from CPU increases, size & access time both increases.



size — increases  
Cost — decreases  
speed — increases  
(access time)  
distance — increases

Mem hierarchy — multiple levels

Two level hierarchy { One focus — only on two levels.  
upper level — closer to CPU, smaller, faster, expensive  
than lower level.

\* Minimum unit of info that can either be present or not in 2-level hierarchy — block or line.



(81)

\* if data requested by CPU,  
appears in some block in upper  
level - hit.

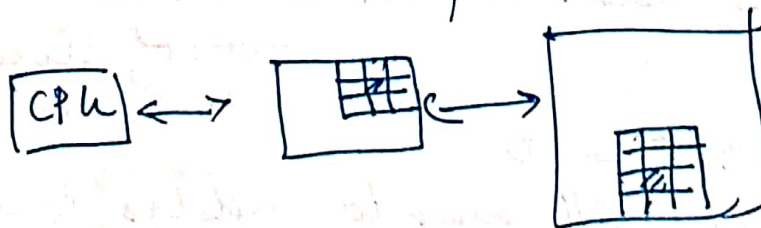
not found - miss.

hit ratio - fraction of mem  
(rate) accesses found in upper level  
(measure of part of mem  
hierarchy)  
Total hits / Total accesses.

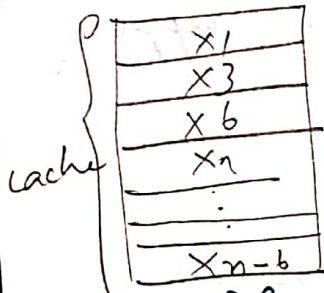
miss ratio  
(rate) - fraction of mem  
accesses not found in upper  
level.

Hit time - time <sup>reqd</sup> to access a ~~low~~ level  
of mem hierarchy including  
time needed to determine whether  
the access is a hit or miss.

miss penalty - time reqd to fetch a block  
into a level of mem hierarchy from  
the lower level, including the time  
to access the block, transmit it from  
one level to other & insert it in the  
level that experienced the miss.



# Cache - level of mem hierarchy <sup>4/10</sup> CPU & mm.

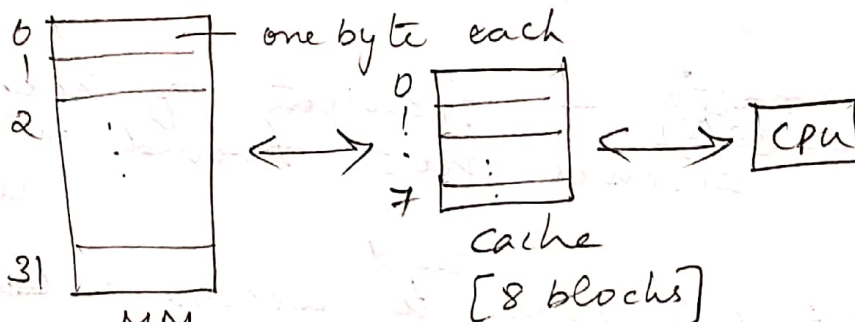


if a req  
for  $x_{n+1}$

← put into some free space → where to put

eg: if main memory has 32 blocks  
and cache is only 8 blocks.

which block in mm will go to  
which block in cache?



\* data to be processed by CPU has to  
brought to cache first.

various schemes → direct mapping  
→ set associative mapping  
→ fully associative.

⑦ Direct mapping.

$$\text{Cache addr} = \text{MM block addr} \% \text{no. of block in cache}$$

eg:  $2 \% 8 = 0$

M[2] will map to address 0 in cache

$12 \% 8 = 4$       M[12]      "      "      4 in cache



$M[0], M[8], M[16], M[24] \rightarrow$  cache address  $\textcircled{8}$   
 map to  
 $0 \% 8, 8 \% 8, 16 \% 8, 24 \% 8 = 0$

$M[1], M[9], M[17], M[25] \rightarrow 1$

$M[2], M[10], M[18], M[26] \rightarrow 2$

~~$M[23]$~~   
 $M[7], M[15], M[23], M[31] \rightarrow 7$

8 blocks in cache

In this eg. 4 mem references go to one entry in cache

$\therefore$  how do we know whether a requested word is in cache or not?

$\rightarrow$  tag bits.  $\rightarrow$  address info reqd to identify whether a word in the cache corresponds to requested word.

$\rightarrow$  upper portion of address, not used for index.

eg: for a 32 block main mem & 8 block cache

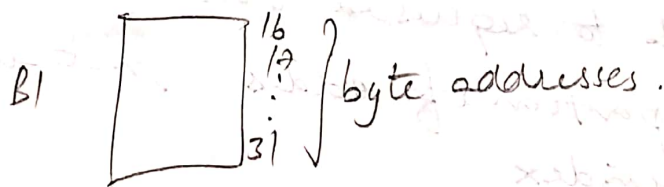
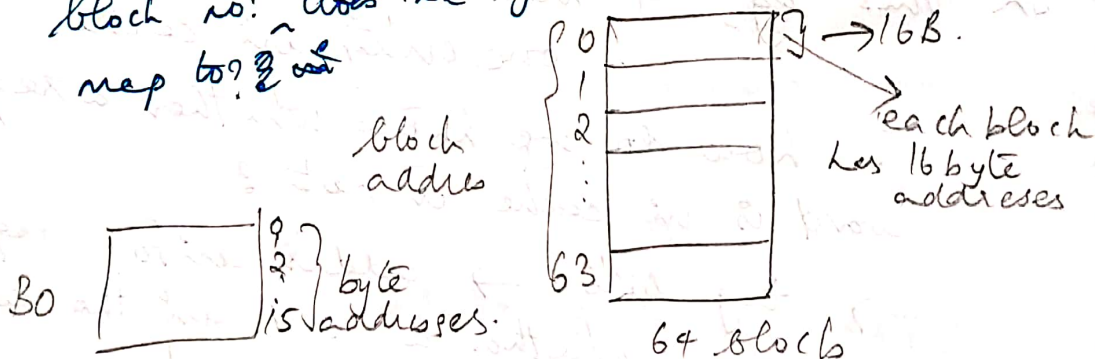
5 bits for 32	tag	index address		
			$\rightarrow$	for cache (8 blocks)
$M[0]$	00	0000		$M[1]$ - 00 001
$M[8]$	00	0000		$M[9]$ - 01 001
$M[16]$	10	0000		$M[17]$ - 10 001
$M[24]$	11	0000		$M[25]$ - 11 001
	diff	same - 0		diff same - 1
	for each	address		for each
	address			address

Fig 7.6 — example page — 477. (86)

valid bit — also in cache

— cache contains valid bit to indicate whether an entry contains a valid address or not.

Q) Consider a cache with 64 blocks and a block size of 16 B. what block no. does the byte address 1200 map to?



~ other blocks.

byte address 0 → 15 go to block 0 in cache.

$$\text{block address } 0 \rightarrow 0 \% 64 = 0$$

16 → 31 : block address 1

~ 32 → 47 " 2

$\lfloor \text{byte address} / \text{no. of bytes} \rfloor = \text{block no.}$

$$\therefore \lfloor \frac{1200}{16} \rfloor = 75$$



$$\text{cache address} = \text{block address} \% \text{no. of blocks} \\ = 75 \% 64 = \underline{\underline{11}}$$

Q) cache with 128 blocks & a block size of 8B. what block no. does the byte address 2004 map to?

$$\text{block no.} = \left\lfloor \frac{\text{byte address}}{\text{block size}} \right\rfloor = 250$$

$$\text{cache address} = 250 \% 128 \\ = \underline{\underline{122}}$$

$$\text{Size of single entry in cache} = 1(\text{vbit}) + \text{no. of tag bits} + \text{no. of bits for data.}$$

$$\text{tag} = \frac{\text{Total mem mem address} - \text{no. of bits for index + B0}}{\text{no. of bits for index + B0}}$$

$$\text{Total cache size} = \text{total no. of blocks} \times \text{size of single entry.}$$

$$\text{Total data size} = \text{total no. of blocks} \times \text{data / block.}$$

Q) For a cache with 32KB of data with and 8 word blocks, find the no. of blocks in cache?

$$= \frac{32 \times 2^{10} \text{ B}}{8 \times 4 \text{ B}} = \frac{2^{15}}{2^5} \text{ B} = \underline{\underline{2^{10}}}$$

- Q) Assume ~~32 bit~~ 1GB main mem with a cache of 1KB blocks and 4 word block size, find the no: of bits for tag, index and BO.

4 word block size

$$= 4 \times 4 = 16 = 2^4$$

$$1GB = 2^{31}$$

$$\therefore BO = \underline{\underline{4}}$$

$$Cache = 1K \text{ blocks} = 2^{10}$$

$$\therefore \text{index} = 10$$

$$\therefore \text{tag} = 31 - (10 + 4) = 7$$

- Q) cache with 8K blocks and 2 words / block find tag, index & BO for a 32 bit address

- Q) How many total bits are required for a direct mapped cache with 16KB of data and 4 word blocks assuming 32 bit address.

total bits = no: of blocks  $\times$  size of each entry.

✓	tag	data

$$\text{no: of blocks} = \frac{16 \text{ KB}}{4 \text{ word}} = \frac{16 \times 2^{10}}{4 \times 4 \text{ B}}$$

$$\text{index} = 10$$

$$= 2^{10} \text{ blocks}$$



size of each entry =  $v + \text{tag} + \text{data}$

$$4 \text{ word blocks} = 4 \times 4 \text{ B} = \cancel{4}^{16} 16 \text{ B} \\ = 2^4$$

$$\therefore 130 = \underline{\underline{4}}$$

$$\therefore \text{tag} = 32 - (10 + 4) = \cancel{26}^{18} \underline{\underline{18}}$$

we are row = 4 words (1 word = 32 bits)

$$\therefore \text{total bits} = 4 \times 32 = \underline{\underline{128}}$$

$$= 1 + 128 + 18 = 147$$

$$\therefore \text{total bits} = \text{size of each} \times \text{no. of blocks} \\ = \underline{\underline{147 \times 2^{10}}}$$

Q) same question as b4 with 64 kB  
of data & 8 word blocks  
assuming 30 bit address?

## Handling cache miss.

(90)

cache miss - A request for data from the cache that cannot be fulfilled because the data is not present in cache.

steps taken on an write cache miss.

- ① send the original PC value (of the instr not present in cache) (PC-4) to memory
- ② Instruct main mem to perform a read op for the requested instr & wait for the mem to complete its access. (wait is equivalent to stall).
- ③ Once main mem completes the access, write this instr to cache putting data from ~~data~~ mem in the data portion of cache entry, writing the upper bits of address (tag) into the tag field and turning valid bit on.
- ④ Restart instr exec at the 1st step, which will refetch the instr, this time finding it in cache.



## Handling writes.

- \* on a store instr, we have to write data both to cache as well as main mem., else if we write only to cache - cache & mem are said to be inconsistent.

various schemes.

### ① write through

\* simplest

\* update both mem & cache at the same time.

\* disadv: poor performance  
 coz - writing data to memory - long time -  $\approx 100 \text{ ns}$

### ② write buffer.

\* a queue that stores data while it is waiting to be written to mem.

\* after writing data to cache & buffer, processor can continue execution.

\* when write to mem is complete, write buffer is cleared.

\* if buffer is full when cpu reaches a write, cpu must stall until there is an empty pos in buffer.

### ③ write back

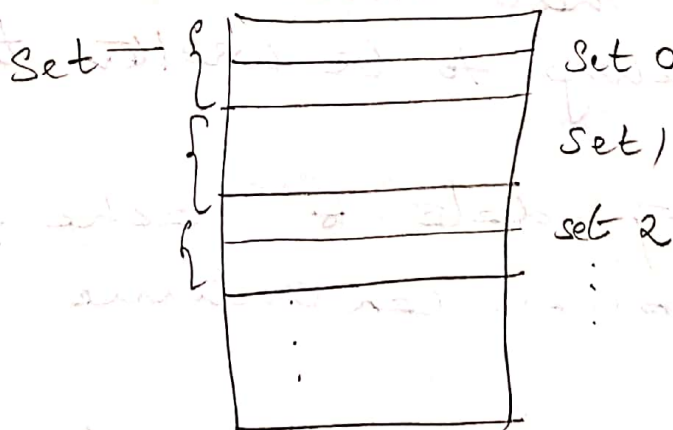
92

- \* handles writes by updating values only to the block in cache.
- \* modified block is written to mem only when it is replaced.
- \* good performance but complicated to implement compared to write through.

②

### Set Associative cache

- \* fixed no: of loc where each block can be placed.



2-way set associative.

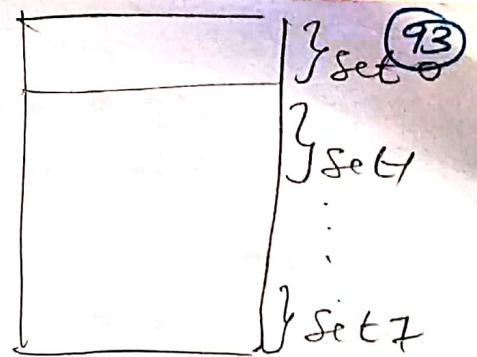
every set - 2 blocks

formula -  $\text{block no:} \% \text{ no: of sets in cache}$



eg: cache with 16 blocks  
for a 2-way set  
associative

16 blocks - 8 sets.

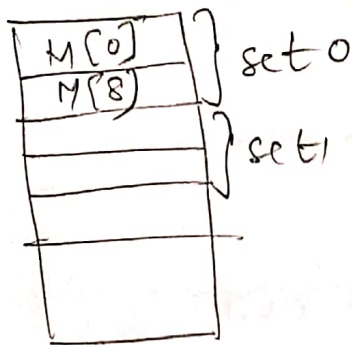


16 blocks.

if main mem with 64 blocks.

cache address = block address % no. of sets

$$\left. \begin{array}{l} M[0] \\ M[8] \\ M[16] \\ \vdots \\ M[56] \end{array} \right\} \% 8 = 0 \rightarrow \text{will go to set 0.}$$



if  $M[0]$  comes first, then  $M[8]$ , no need to replace  $M[0]$ , it will go to the next slot in set 0.  
only if another address, say  $M[16]$  comes again, we need to go for replacement.

[For replacement - any replacement policy. eg: LRU least recently used]

Four-way set associative - 4 slots / 4 blocks/set.

8-way - 8 blocks /set.

fully set associative - all blocks in one set.

→ any main block address can go to any block in cache.

# Floating point representation.

IEEE 754 single precision format. (32 bits)



Q) -0.75

in binary :- 0.11

scientific notation :  $-0.11 \times 2^0$

normalised sci not :  $-1.1 \times 2^{-1}$

gen representation for a single precision no. 75 :

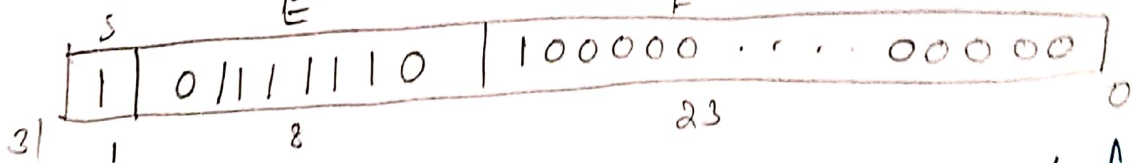
$$(-1)^S \times (1 + F) \times 2^{(E - \text{bias})}$$

= actual/exp

$$(E - \text{bias}) = -1$$
$$\therefore E = \underline{126}$$

$126 - 127$

$$(-1)^1 \times (1 + .10000 \dots 00) \times 2^{126 - 127}$$



- single precision binary representation

Q) 25.4

binary : 11001.01100

sci not :  $11001.01100 \times 2^0$



normalized not :

$$1.1001011 \times 2^4$$

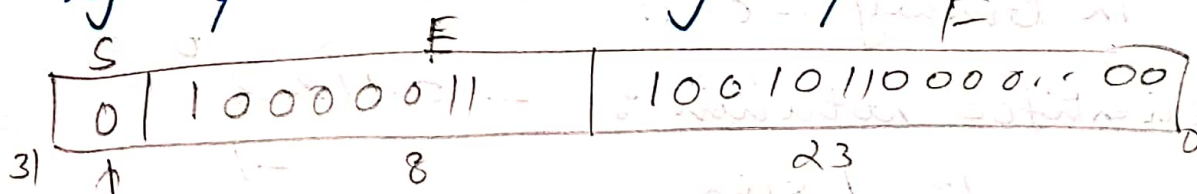
gen. rep: for single precision

(E-127)

$$(-1)^S \times (1 + F) \times 2$$

$$(-1)^0 \times (1 + 0.1001011) \times 2^{(131-127)}$$

single precision binary expr:



Double precision: (64 bits)  
(bias-1023)

