

$$\text{tag + valid} = 18 + 2 \text{ bits}$$

∴ total =

Fully associative mapping

If the cache has n blocks, data address can be mapped to any one of the block.

Use comparators to check tag bit of each block to that of address of data needed

$$\Rightarrow \left(\frac{\text{mem}}{\text{block no}} \right) \% (\text{no. of block in cache})$$

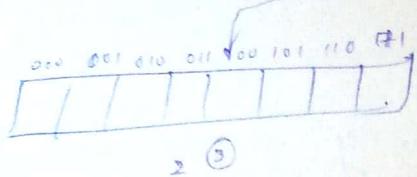
Set associative

n -way set associative \Rightarrow a set has n blocks

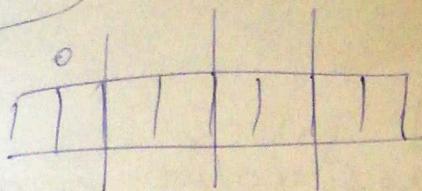
In this we need to identify the set & then use tag to find block within that set.

$$\Rightarrow \left(\frac{\text{mem}}{\text{block no}} \right) \% \left(\frac{\text{no. of sets}}{\text{in cache}} \right)$$

eg: mem. add = 12 = 1100



direct map

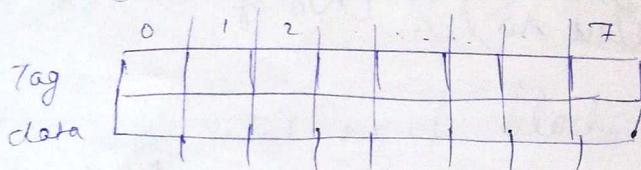


2-way set
associative

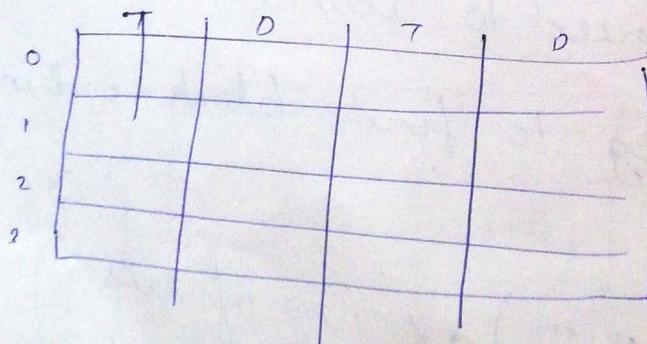
$$\text{no. of set} = \frac{8}{2} = 4$$

$$\therefore 12 \% 4 = 0$$

1-way set ass. with 8 block.



2-way - 8 block \Rightarrow 4 set



There are 3 cache each with 4 word block.

- 1 is fully associative
- 2 is 2 way ass
- 3 is direct

find no. of misses for each cache for

given seq of block address

0, 8, 0, 6, 8

In 2 way set, if a set $\boxed{x \mid y}$ has 2 data & another data is mapped to it, by least recent principle, the least recently used is replaced (e.g. if x is placed after y , y is replaced)

Add. mem block	hit/miss	Contents of cache block			
		0	1	2	3
0	M	m[0]			
8	M	m[8] <small>update</small>			
0	M	m[0]			
6	M	m[2]			m[6]
8	M	m[8] <small>update</small>			

in direct $0 \% \xrightarrow{\text{no of blocks}} 0$, for direct mapping

$$8 \% 4 = 0$$

$$6 \% 4 = 2$$

⇒ for 2 ways - set associative

addr.	hit/miss	set 0	set 1	set 2	set 3
0	M	m[0]			
8	M		m[8]		
0	H	m[0]			
6	M		m[6]		
8	M	m[8]			

replace least recently used
in m[8]
in m[0]

$$\text{no. of set} = 4/2 = 2$$

$$0 \% 2 = 0 \quad \begin{matrix} \text{no. of} \\ \text{block} \\ \text{no. of set} \end{matrix}$$

$$8 \% 2 = 0$$

$$6 \% 2 = 0$$

fully associative
in a set
place in any block of set 0

direct mapped
about set

⇒ for fully associative

addr.	hit/miss	0	1	2	3
0	M	m[0]			
8	M		m[8]		
0	H	m[0]			
6	M			m[6]	
8	H		m[8]		

We have to store
in such a way
to minimise
the misses

if they're
free block
store if there

* 8 block ≠ 1G block

$$\text{no. of set} = 4 + 8$$

Handling Cache misses

If an instruction is checked, if its a miss in cache, we check the address of that instruction ($= \underline{PC} - 4$) & already inc by 4
we need prev instr add
check in 1° mem.

Handling Cache writes

After changing value of say variable n , if we change only the value of n in cache (copied from its original loc in 1° mem) during fetching the cache will have modified value & 1° mem will have old value causing an inconsistency.

To avoid this, we can use

(i) Write through:

"processor writes to cache & at the same time it writes to memory." takes more time This causes more delay to processor

(ii) Buffered write

processor writes to buffer instead of writing to \sim mem (reducing the time to write to \sim mem) \because mem writes / by processor at buff is faster
update itself from buffer

Now if buffer is full, it creates a stall

we can avoid this by:

- (i) increase buffer size
- (ii) keep multiple queues/buffer

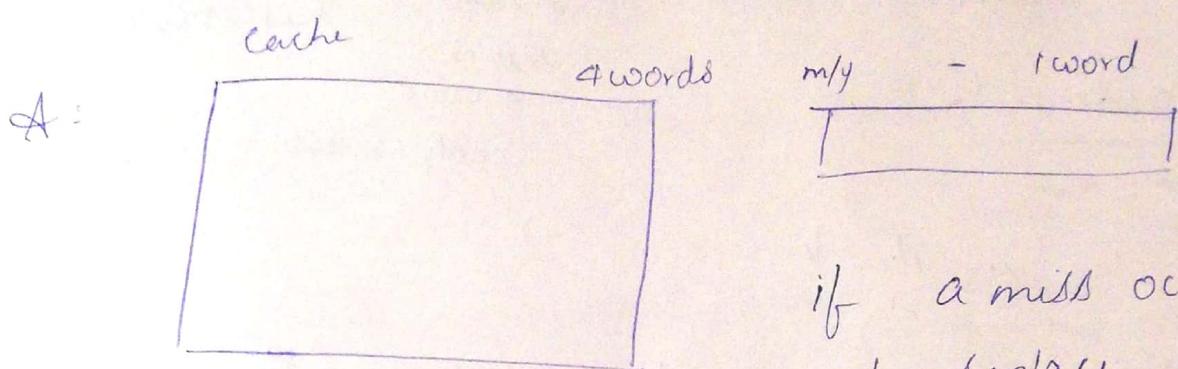
(iii) Write Back/copy back

Update \sim mem only if a replacement occurs to corresponding add in cache or write copy back add data in cache when power is off.

DESIGNING M/Y SYSTEMS

Q. Assume 1 mem bus clk cycle is needed to send address. 15 mem bus clk cycles for each d-ham access initiated
 \sim mem

1 mem bus clk cycle to send a word of data
 If we have a cache block of 4 words &
 word byte wide bank of 8 banks. Calculate
 miss penalty in terms of mem. bus
 clk cycle?



if a miss occur, we have
 to replace all the 4 words
 in cache (not that 1
 word alone)

If miss happens

- send add to lower mem $\rightarrow 1$
 - access 4 words in 1' mem $\rightarrow 15 \times 4$
 - ~~wpy back data~~ $\rightarrow 1 \times 4$ word
- $15 \times 4 + 1 = 65$

if we can access the 4 words seq. if we provide only the starting address

if we have to provide add. of each of 4 word
 in cache, then

$$4 \times 1 + 15 \times 4 + 1 \times 4 = 68$$

Now, if we \downarrow cache size, miss penalty \downarrow
but miss rate \uparrow .

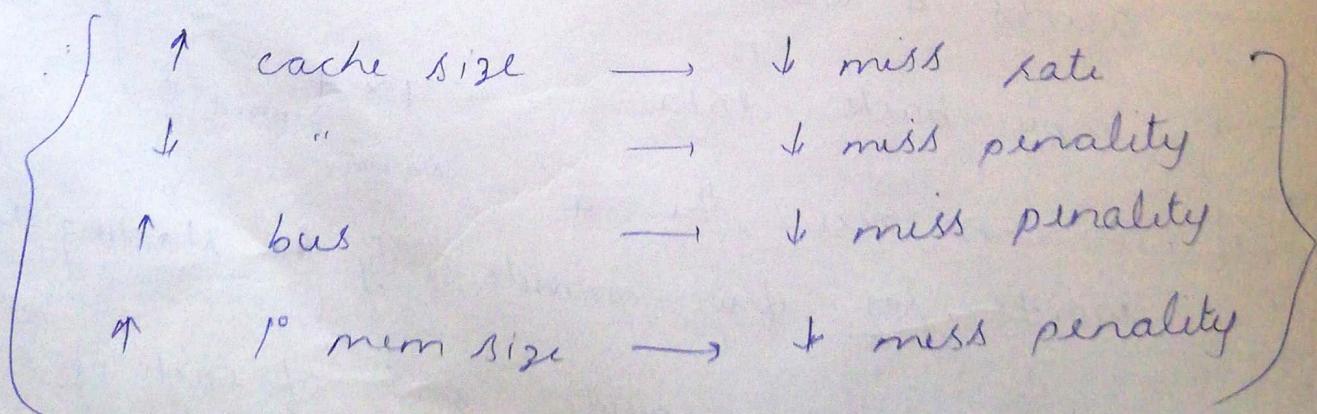
but if we \uparrow size of 1° mem (instead
of \downarrow cache size), then

$$\text{miss penalty} = 1 + \underbrace{15}_{\substack{\text{d-lam.} \\ \text{size is} \\ 4 \text{ word}}} + \underbrace{4 \times 1}_{\substack{\text{bus = 1 byte} \\ \therefore \text{only 1 access}}} = 20$$

in it \downarrow

Now if we \uparrow bus size,

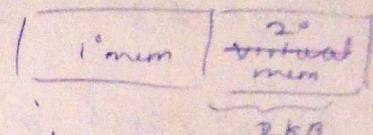
$$\text{miss penalty} = \underbrace{1}_{\substack{\text{send} \\ \text{one add}}} + \underbrace{15}_{\substack{\text{only 1} \\ \text{access}}} + \underbrace{1}_{\substack{\text{only 1} \\ \text{access} \\ \text{to} \\ \text{return} \\ \text{data}}} = 17$$



Virtual memory

If RAM is having less by miss rate can be higher along with miss penalty.

Hence, we attach a part of 2° mem
virtual to 1° mem so that processor has a larger view. This is called virtual mem



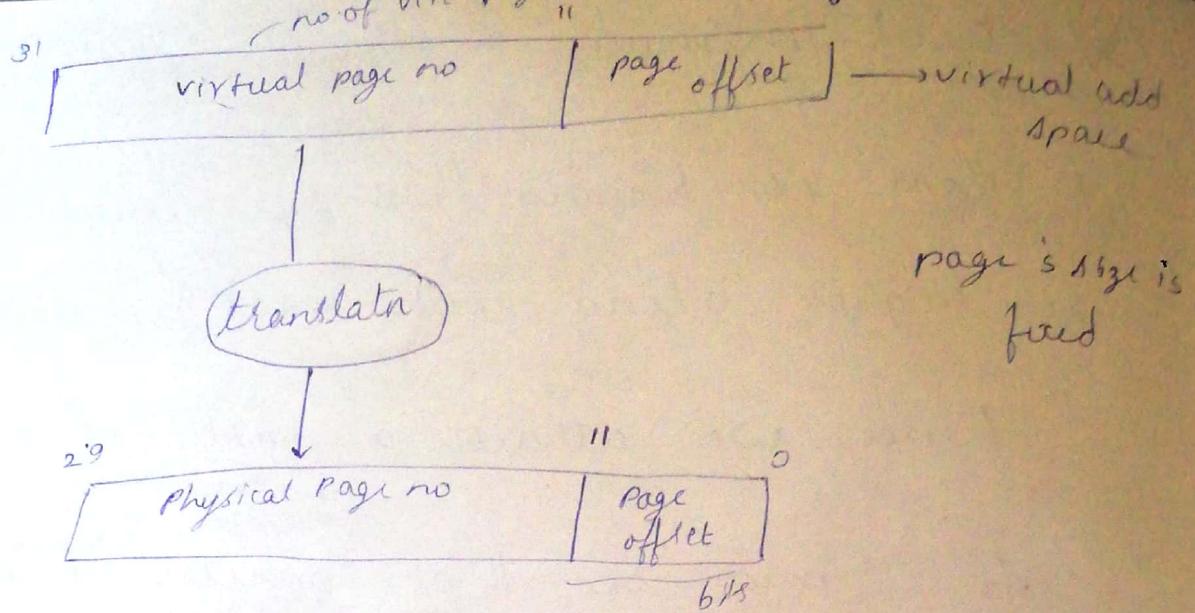
Also, virtual mem allows sharing of programs.

Page: block is called a page

Pagefault: miss is called pagefault

Virtual add: add corresponding to virtual space

Address mapping / translation: translation of virtual add to a 1° mem add using s/w & h/w



$$\text{Page size} = 2^{12} = 4\text{kB}$$

$$\text{No. of pages in physical mem} = 2^{18}$$

No. of pages in virtual mem = (2^{20}) - large block

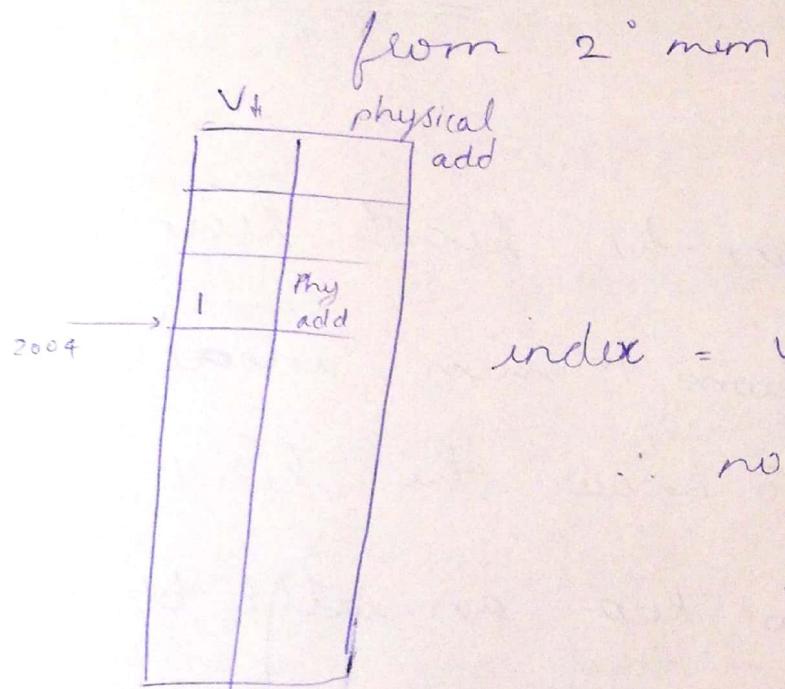
Now, we need to know whether a page is present in 1^o mem / 2^o mem of virtual mem

∴ we need a page table with each cell having . . .

- iii) valid bit $V=0 \Rightarrow$ page fault a page in 2^o
- $V=1 \Rightarrow$ page in 1^o mem

Once we know the virtual mem add, we use it as an index to search the

page table. If $v=1$, the add phy. add corresponding to it shows page is in 1° mem if $v=0$, \Rightarrow page in 2° mem, then we use phy. add. adjacent to that index to search 2° mem & update page table by moving data to 1° mem



index = virtual page no

\therefore no. of index in pagetable

$$= 2^{20}$$

each program has its own page table & page table e.g. store the starting add of page table.

⇒ Another type of mapping is called segmentation

⇒ Now, when we have to replace some value in 1^o mem when a page fault occur in pagetable, we use copy back mechanism. (value in 2^o mem moved to add. in 1^o mem & the value at add n is moved to a space called swap space in 2^o mem.)

⇒ If we replace the least recently used data from 1^o mem incase of a page fault. To know the least recently used block, we keep an extra bit called reference bit R

⇒ When a data is used we set R=1 but after a time stamp (eg 6ms) its value is cleared saying that it is not a recently used data.

In case of cache, we need to copy back in case of replacement of if data in cache was updated. keep dirty bit = 1 if cache was updated frequently & hence we have to copy back data in cache \neq data in corresponding 1^o mem.

or before replacing the value in cache in case of miss (if dirty bit = 0, no copy back just replace)

11th in case of virtual mem. space in page table dirty bit is set as 1 if data in that table is changed. If its 1, then we have to copy back the data in that cell to corresponding byte in 2^o mem before replacing it.

Ques
Page table resides in physical mem.

Translation look aside buffer
(TLB)

→ page table in cache

In 1^o mem.

	V	D	R	phy add.
1	1	0	1	
1	1	0	0	
0	0	0	0	
0	0	1	1	
1	1	0	0	

We keep cache to get faster access to page table. Hence we keep only the entries of the recently used index in cache

∴ TLB

V	D	R	phy add.
✓	0	1	

→ recently used index in page table of 1^o

Cache TLB can store may be max of 1 or 2 entries. If cache bi miss occur in this case we replace the recently used
~~if the 2 index is passed to~~

TLB miss \Rightarrow entry miss not that
it's a page fault. That may be
available in page table. Access page table
modify if its page fault & then replace
the least recently used entry in TLB

{
1st TLB if miss goto page table update TLB
after modifying page table in case of page fault}