

Now, processor at once copies all the data items needed for a given piece of code rather than one at a time.

Now it is based on principle of locality

### ① temporal locality

if we access  $x$ , then there is higher probability of accessing  $x$  in near future (e.g. loops) ( $\because$  copy to cache)

### ② spatial locality

if we access  $x$ , then there is higher prob. of accessing data near to  $x$  in near future

### 4. Arrays

(If we copy items / subset of data so that if we need a data we need to search only cache (smaller) than  $1^{\circ}$  or  $2^{\circ}$  (larger))

floating point numbers not so.

Ex:  $x = -1.5 \times 10^{38}$

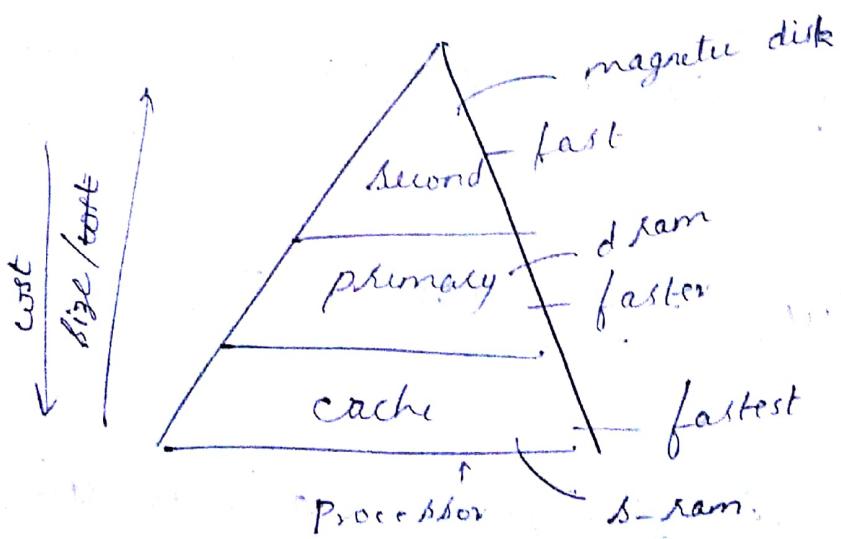
$y = 1.5 \times 10^{38}$

$z = 1.0$

I]  $x + (y+z) = x + 1.5 \times 10^{38}$   
= 0

II]  $(x+y)+z = 0 \times 10^{38} + z$   
= 0 + 1  
= 1

Memory hierarchy:



processor access data from cache

data copied from  $2^{\circ} \rightarrow 1^{\circ} \rightarrow$  cache  $\rightarrow$  processor

block/line: Minimum unit of info.  
that is either present or not present  
mem.

hit: A hit is occurred if a data item  
is found in <sup>given</sup> a level else its a  
miss.

If a hit is occurred it copies it to  
next level closer to processor

else

Searches in next level above it

hit ratio: % of memory access satisfied  
by that level

$$\text{if } \frac{\text{no. of hits in last } \times 100}{\text{no. of access in last}}$$

$$M^{\text{ry}} \text{ miss ratio} = 100 - \text{hit ratio}$$

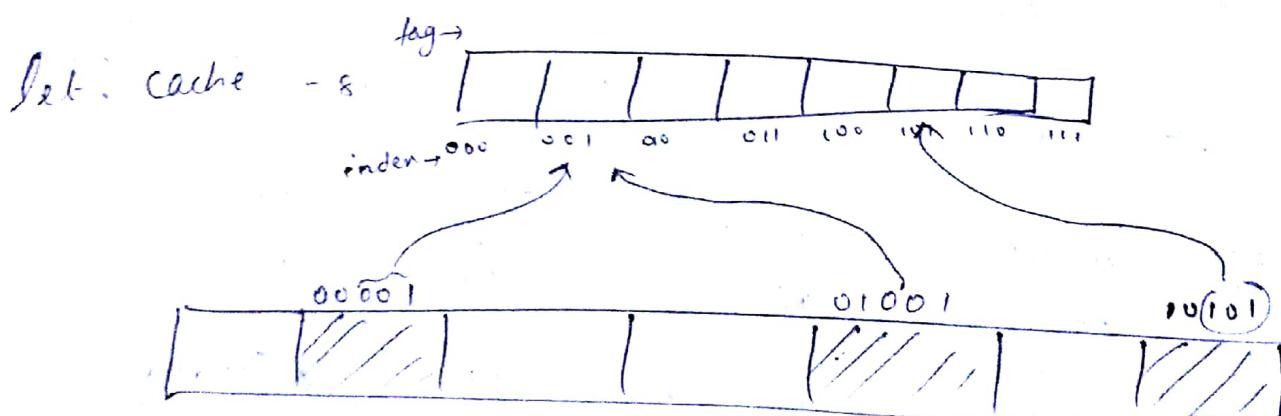
hit-time: time taken to provide a data  
to CPU including the time to

check if its a miss / hit

misperability: If its a miss time to copy lower to upper mem + time to deliver data to CPU

Now, in cache limited space is there so how do we copy all the elements needed for a program in there (duplication or overwrite may occur)

↳ how the processor can identify the data from cache?



• lower 3 bit as index & upper as tag  
here by using modulo on address

• if add data = 00001 we do  
modulo 8 to find index 001

Now to know if its lower bit of  $01001_0$   
00001 place upper bit in tag, associated  
that position in cache.

This is called Direct mapped cache

Now to see if its valid data or not we  
have a valid bit also associated with each  
cell.

If an address of a data is

$x_n x_{n-1} \dots x_0$  then if size of  
cache is  $2^n$ , then lower  $n$  bits is used  
to decide the index & rest as tag.

Index	valid	Tag	Data
000	X	10	M[0000]
001	N	11	M[1010]
010	Y	00	M[0001]
011	Y	00	M[0001]
100	.	.	
101	.	.	
110	Y	10	M[1010]
111	N	.	

If 10110 is needed by processor, it checks

index : 110 (taghe = 2<sup>(3)</sup>)

first it check 110, since von of tag & index  
& not same ∴ its a miss ∴ copy data  
of store address in corresponding index

Then we need 11010  
tag index → miss ∴ store

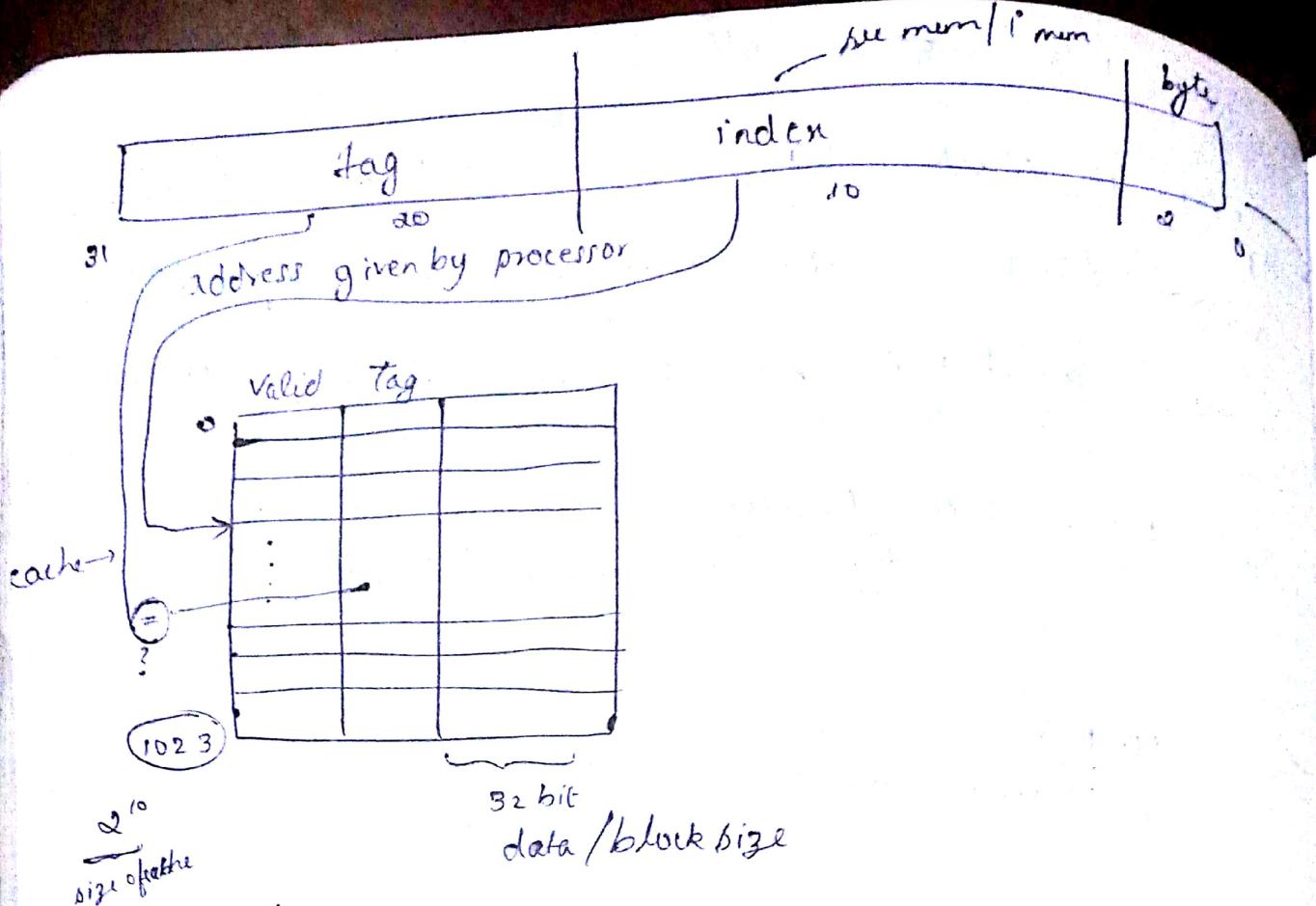
Then 10 000  
tag index → miss ∴ store

Then 00011  
tag index → miss ∴ store

Then 10 010  
tag index → miss [index same but tag not equal]

∴ store update tag

∴ if info is changed in the memory  
change valid bit of corresponding address as  
in the cache.



here we need only part of data stored in cache (earlier we accessed all the bits stored in cache) processor req all data stored in a block here if we need 1<sup>st</sup> byte of data, what to do?

Here we use 10 bits for index

Now how to access byte by byte.

we need 2 bits to know which portion of word you need to access (words has 4 bytes)

∴ index - 10  
byte portion - 2  
tag - 20 bits.

If processor gives a 32 bit address, we use its index part to find its location in cache. tag part to check if its the requested data that is stored in that location. & valid bit to check validity. If its a hit, then we use the byte part to know which byte portion of the word stored in that cache location is required by processor.

a) 32 bit byte addressed direct mapped cache  
of size  $2^k$  blocks, of block size  $2^m$  words.  
 find size of tag field. total no. of bits in direct mapped cache

$$\text{total} = 2^k \cdot 2^m \times 2^l$$

$$\text{tag field} = 32 - \underbrace{n}_{\text{index + byte}} - m - 32 - [n+m+2]$$

after to identify a word from  $2^m$  words

to find the byte position from each word

Store data

$$\text{total} = 2^n [\text{Block size} + \text{tag size} + \text{valid}]$$

$$= 2^n \left[ \underbrace{2^m \times 32}_{\text{bits}} + 32 - [n+m+2] + 1 \right]$$

Q) How many total bits & tag for direct mapped cache with 16 kB of data & 4 word block assuming 32 bit address

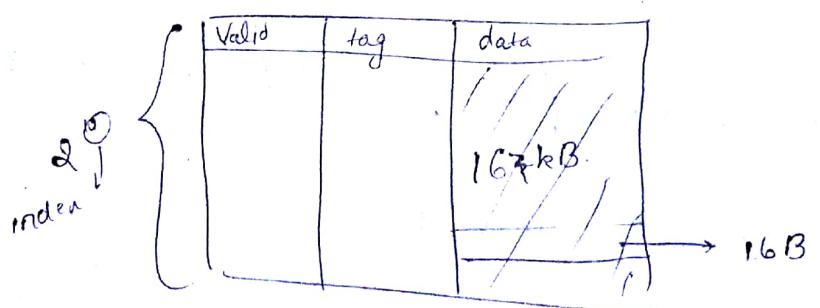
$$\text{A. tot. block size} = 16 \text{ kB}$$

$$= \underbrace{\text{size of 1 block}}_{9 \times 4 \text{ byte}} \times \text{total no. of blocks}$$

$$\Rightarrow \text{total no. of block lines} = \frac{16 \text{ kB}}{16 \text{ byte}}$$

$$= \cancel{16 \text{ kB}}$$

$$= 2^{10}$$



$$\therefore \text{tag} = 32 - [10 + 2 + 2 \text{ for word } - 2 \text{ for byte}]$$

$$= 32 - [14] = 18$$

$$\therefore 1 \text{ line} = 1 + 18 + 4 \times 32 = 19 + 128$$

$$\text{total} = 2^{10} \times 147 \text{ bits} // = \underline{\underline{147 \text{ kB}}}$$

~~tag + valid~~ = 18+2 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

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

## Set associative

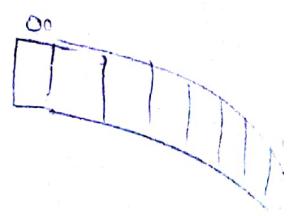
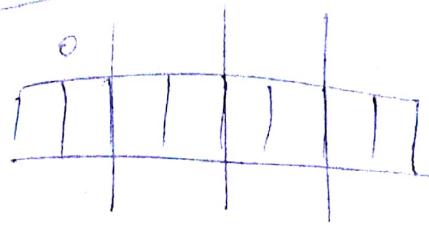
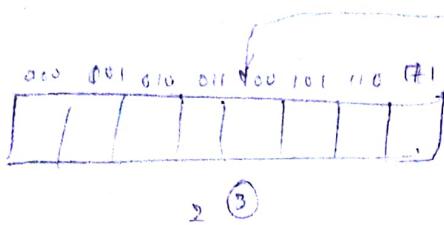
$n$ -way set associative  $\Rightarrow$  a set has  $n$  block

first we need to identify the set & then use tag to find block within that

set.

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

e.g: mem. add = 12 = 1100



direct map

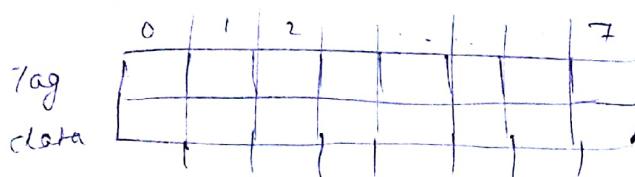
2-way set

associative

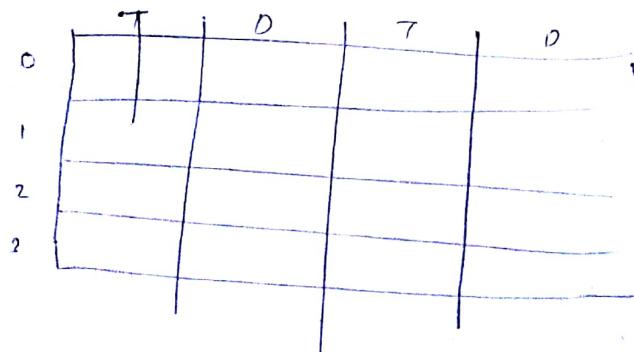
$$\text{no. of set} = \frac{8}{2} = 4$$

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

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	$m[8]$	update		
0	M	$m[0]$			
6	M	$m[2]$			$m[6]$
8	M	$m[8]$	update		

in which  $0 \% \stackrel{\text{no of block}}{=} 0$   $\rightarrow$  for direct mapping

$$3 \% 4 = 0$$

$$6 \% 4 = 2$$

→ for ~~ways~~ - set associative

add	hit/miss	set 0	set 1	set 2	set 3	
8	M	M[0]				
9	M		M[8]			
0	H	M[0]				
6	M		M[6]			→ replace least recently used
8	M	M[8]				M[6]

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

$$0 \% 2 = 0 \text{ set}$$

$$8 \% 2 = 0$$

$$6 \% 2 = 0$$

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

direct mapped  
about set

→ for fully associative

add	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 minimize  
the misses

if they  
free block  
store it there

\* 8 block ≠ 16 block

$$\text{no of set} = 4 / 8$$

## Handling Cache miss

If an instruction is checked, if it's a miss in cache, we check the address of that instruction ( $= PC - 4$ )  $\phi$   
already inc by 4  
we need prev instr  
add  
check in 1<sup>o</sup> 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<sup>o</sup> mem) during fetching the cache will have modified value  $\phi$  1<sup>o</sup> mem and will have old value causing an inconsistency.

To avoid this, we can use

### (i) Write through:

processor writes to cache  $\phi$  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  $r^{\text{mem}}$  (reducing the time to write to  $r^{\text{mem}}$ ) by processor as buff is faster, mem writes / update itself from buffer

Now if buffer is full, it creates a st

we can avoid this by:

(i) increase buffer size

(ii) keep multiple queues/buffer

### (iii) Write Back/copy back

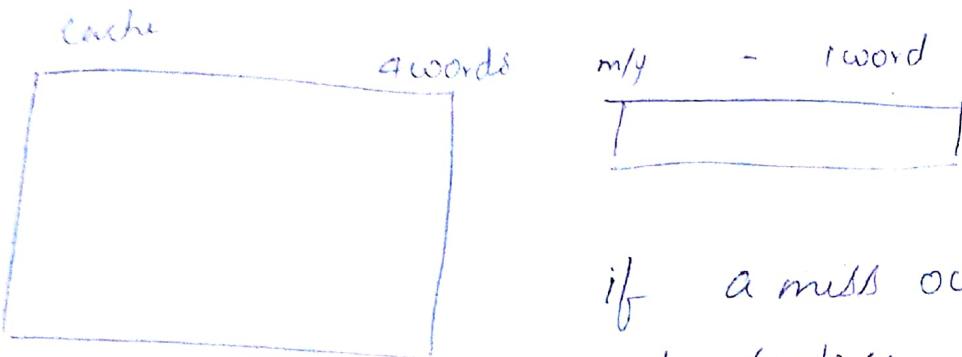
Update  $r^{\text{mem}}$  only if a replacement occurs to corresponding add. in cache or while copy back add data is in cache when power is off

### DESIGNING MY SYSTEMS

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

$r^{\text{mem}}$

From last cycle to send a word of data  
 If we have a cache block of 4 words of  
 1 word byte wide bank of 4-ways Calculate  
 miss penalty in terms of mem. bus  
 blk cycle?



If a miss occurs, 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

if we can access the 4  
 words beg. 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  $\uparrow$   
but miss rate  $\uparrow$ .

but if we  $\uparrow$  size of  $1^{\circ}$  mem (instead  
of  $\downarrow$  cache size), then

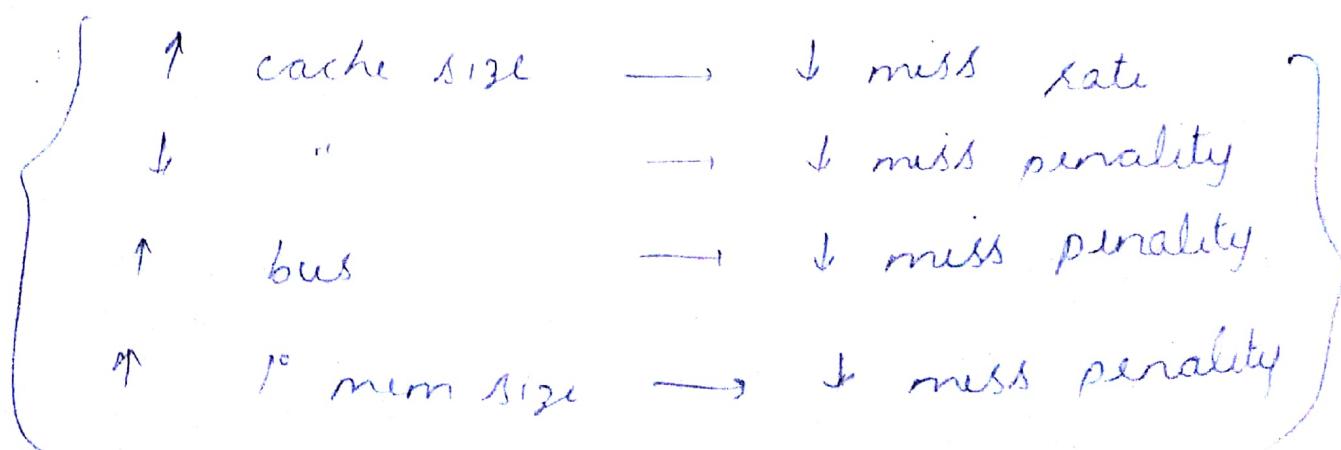
$$\text{miss penalty} = \underbrace{1}_{\text{d-ham}} + \underbrace{15}_{\text{bus size}} + \underbrace{4 \times 1}_{\text{4 word}} = 20$$

$\therefore$  only 1 access

so it  $\downarrow$

Now if we  $\uparrow$  bus size,

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

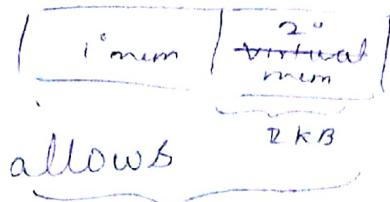


## Virtual memory.

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

hence, we attach a part of 2<sup>o</sup> mem  
virtual  
to 1<sup>o</sup> mem so that processor has a larger view. This is called virtual mem.

Also, virtual mem. allows sharing of programs.

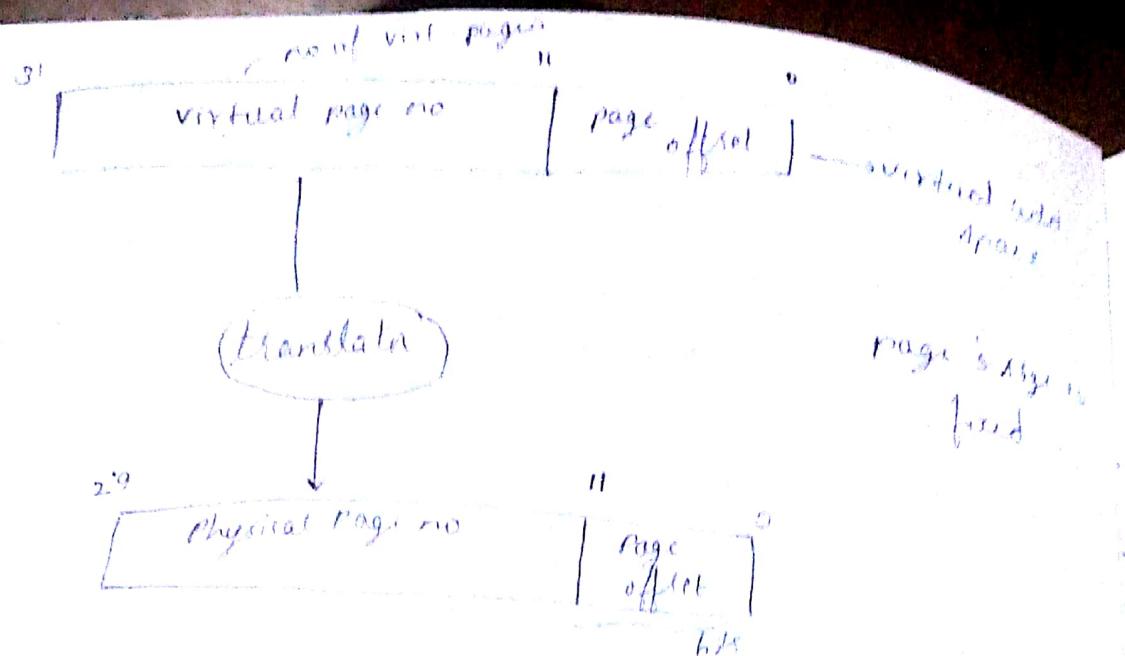


page: block is called a page

Pagefault: miss is called pagefault

virtual add: add corresponding to virtual space

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



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

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

No. of pages in virtual memory =  $(2^{20})$  - 1 pages  
block

Now, we need to know whether a page is present in  $1^{\circ}$  mem / 2<sup>18</sup> pages of virtual memory.

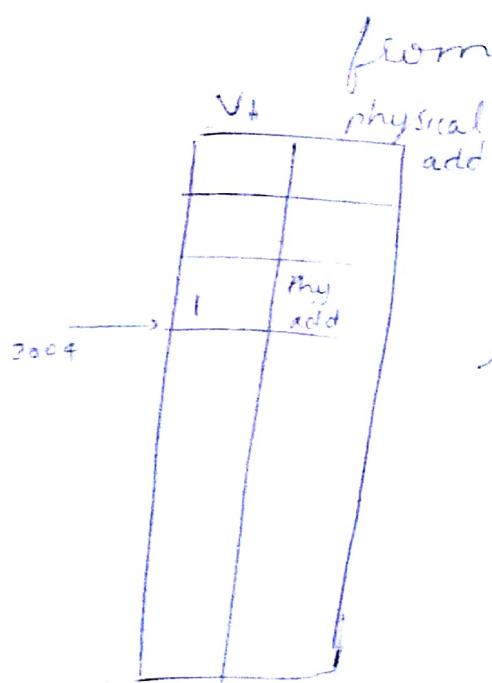
$\therefore$  we need a page table with each cell having .

if valid bit  $V=0 \Rightarrow$  page fault or page not present

$V=1 \Rightarrow$  page in  $1^{\circ}$  mem

Once we know the virtual memory, we use it as an index to search the virtual pages.

page table. If  $v=1$ , the addr phy add corresponding to it shows page is in  $1^{\text{st}}$  mem if  $v=0$ ,  $\Rightarrow$  page in  $2^{\text{nd}}$  mem. Then we use phy add adjacent to that index to search  $2^{\text{nd}}$  mem & update page table by moving data to  $1^{\text{st}}$  mem from  $2^{\text{nd}}$  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.

$\Rightarrow$  Another type of mapping is called segmentation

⇒ Now, when we have to replace value in 1<sup>st</sup> mem when a page fault occurs in page table, we use copy back mechanism. Value in 2<sup>nd</sup> mem moved to 1<sup>st</sup> mem & the value at add n is moved to a space called swap space in 2<sup>nd</sup> mem.

⇒ If we replace the least recently used data from 1<sup>st</sup> mem because of 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. key

dirty bit = 1 if cache was updated

hence we have to copy back

data in cache & data in corresponding main

or before replacing the value in cache

in case of miss (if dirty bit = 0, no copy back just replace)

III) 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 main before replacing it swapping.

Rank

Page table resides in physical mem.

Translation look aside buffer  
(TLB)

→ page table in cache

for program.

V	D	R	phy add.
1	0	1	
1	0	0	
0	0	0	
0	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.
✓	n	r	

→ recently used index in page table of

Cache TLB can store may be max of

1 or 2 entries. If cache be miss case 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

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