

main memory → RAM.

SENDOSS
FILA
JA3FAR
DOUSS

Chapter 5 S.Brahim

hamdan

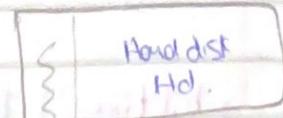
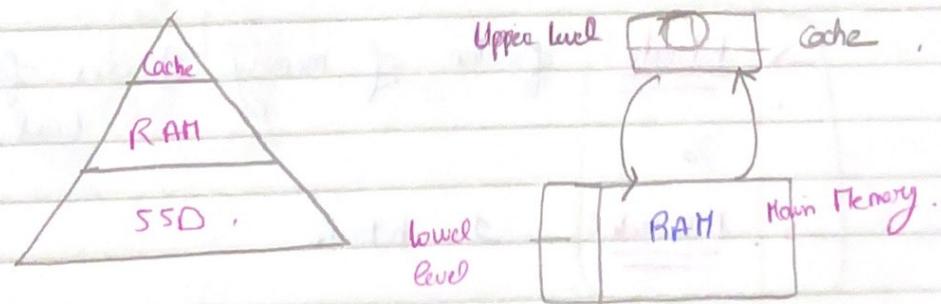
→ Principle of Locality: States that -the Program occurs a relative small portion of their address space at any instance in time -

Temporal locality:
instruction execute now
will most probably
be needed in future

Spacial locality:
code work
Sequentially.

→ Memory hierarchy:

CPU



→ block: is the minimum amount of information that can be transferred between 2 levels of memory.

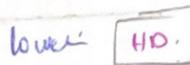
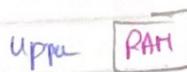
\rightarrow hit: If the data requested by the Processor is found in Upper Level its a ~~bad~~ hit
yes / No

~~yes / No
question~~

→ Miss

f The data requested by the processor
is found in Upper Level its a ~~bad~~ hit

Not a hit.



- hit rate
- ↓
%
- miss rate

fraction of memory Accesses found in upper level.

Miss rate

1 - nitrate

- hit time
- Time
- Miss penalty

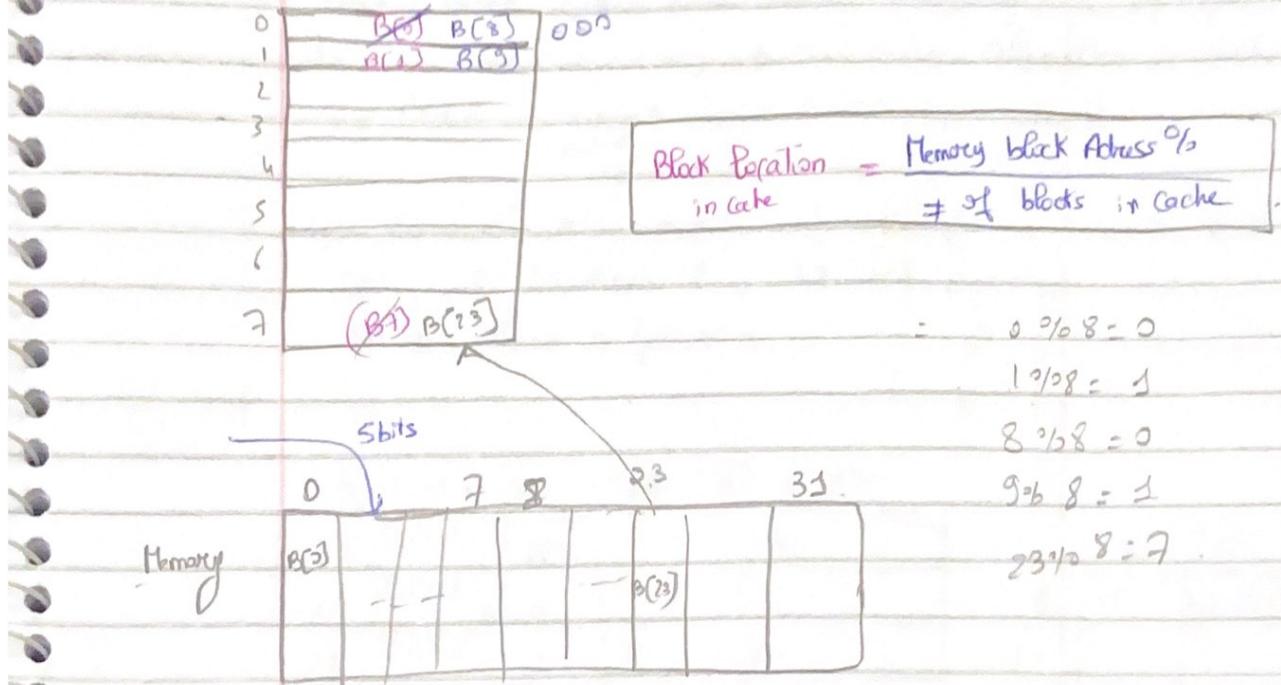
time to access upper level of memory hierarchy.

The time to replace a block in the upper level
with the corresponding block from lower level
+ the time to deliver the block to the
processor.

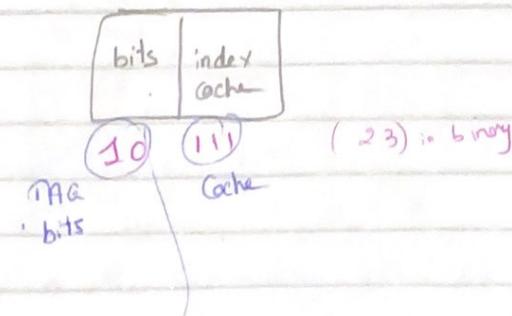
Ibrahim hamdan is the joker of the world.

Che

Direct Mapped cache



PAG



Valid	TAC	Data
Y	10	$B[7]$
N	0	
N	0	

Direct Mapped caches works by bringing a block from
Memory & Applying the modulo operator on it (% # blocks
include)

To

Direct Mapped Cache

Each block taken from memory goes into a specific place in cache according to the below location in cache
 $= \text{Block address } \% \# \text{ of blocks in cache}$

Ex: Given the below memory reference for a 1-word block
 Cache - 8 locations & a 32 word memory.
 After these references.

Address in Cache	TAG	Hit or Miss	v	TAG	Date
22 $22 \% 8 = 6$ 10110	10	Miss	0	1	10 MEM[16]
26 $26 \% 8 = 2$ 01010	11	Miss	1	0	
22 $22 \% 8 = 6$ 10110	10	Hit	2	1	11 10 MEM[26]
26 $26 \% 8 = 2$ 11010	11	Hit	3	1	00 MEM[03]
16 $16 \% 8 = 0$ 10000	10	Miss	4	0	
3 $3 \% 8 = 3$ 00011	00	Miss	5	0	
16 $16 \% 8 = 0$ 10000	10	Hit	6	1	10 MEM[22]
18 $18 \% 8 = 2$ 10010	10	Miss	7	0	
16 $16 \% 8 = 0$ 10000	10	Hit			
26 $26 \% 8 = 2$ 11010	11	Miss			

TAG is the remaining bits of Memory address After taking out bits related to cache addressing

Ex: Address \rightarrow 8 bits
 Cache \rightarrow 3 bits TAG = 5 bits

How many blocks did we replace?

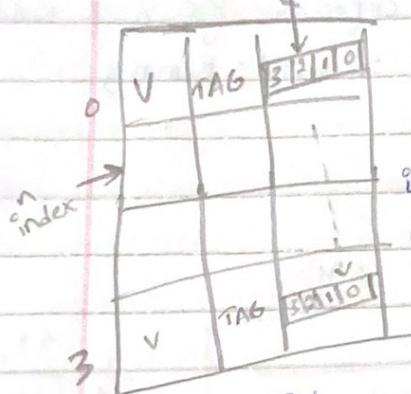
We replaced 2 blocks

- ① MEM [26] by MEM [18]
- ② MEM [8] by MEM [26]

word Addressing \rightarrow 1 word output

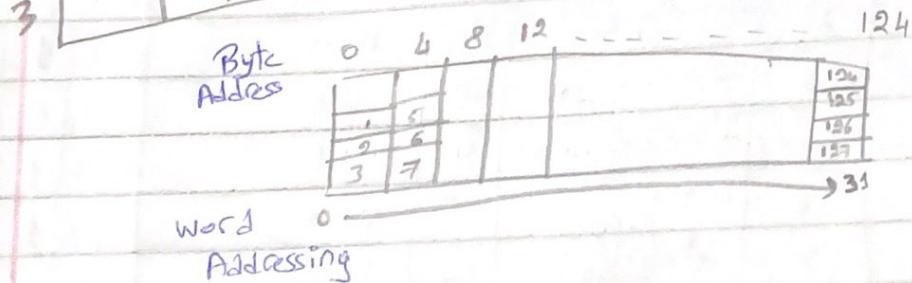
byte Addressing \rightarrow 1 word

byte offset \rightarrow 1 byte



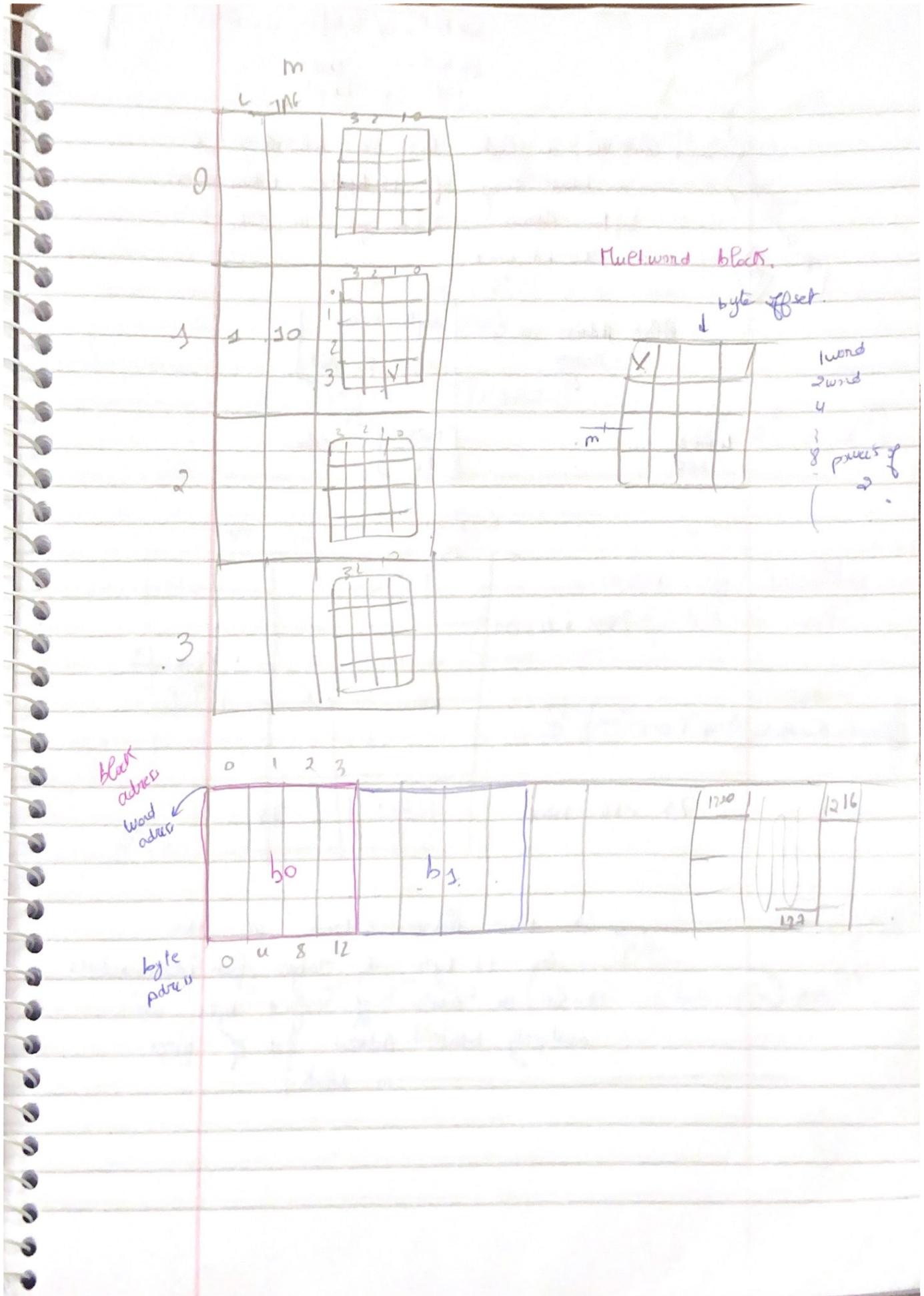
1 word \rightarrow 4 bytes

in case of 32 bit word we have byte offset to be 2 bits.



Given that we changed the Memory to byte Address & Kept it 32-words, what the muche of bits needed to Address the memory?

128 bytes \rightarrow 7 bits to Address.



$$\begin{array}{rcl}
 32 \text{ bit} & \rightarrow & \text{byte} \rightarrow 2 \text{ bit} \\
 64 \text{ bits} & \rightarrow & 3 \text{ bit} \\
 128 & \rightarrow & 4 \text{ bit}
 \end{array}$$

→ Consider a Cache with 64 blocks &
a block size of 16 bytes when does
byte Address 1200 map in Cache
32 bit word

$\boxed{\text{X5}}$

$\boxed{\text{X5}}$

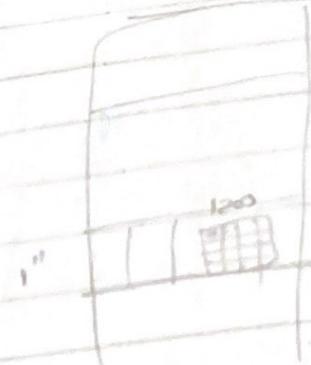
$\boxed{\text{X5}}$

Block Address in Memory = $\left[\begin{array}{l} \text{byte Address} \\ \text{bytes per block} \end{array} \right]$

$$6 \text{ word block} = \left[\begin{array}{l} 1200 \\ 16 \end{array} \right] \text{ floor}$$

$$= 75$$

$$75 \% 64 = 31$$

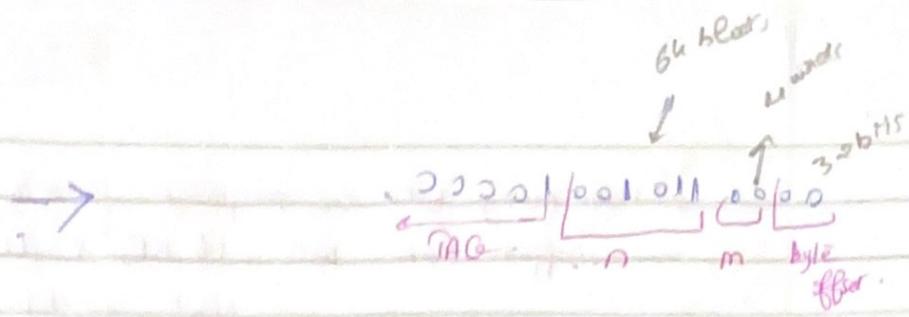


$$\therefore 75 \times 16 = 1200$$

$$\left(\frac{1200}{16} \right) = 75$$

→ The block Address contains in this case 16 bytes if range from 1200 → 1215

To set the Address of first byte we multiply block Address \times # of bytes in block.

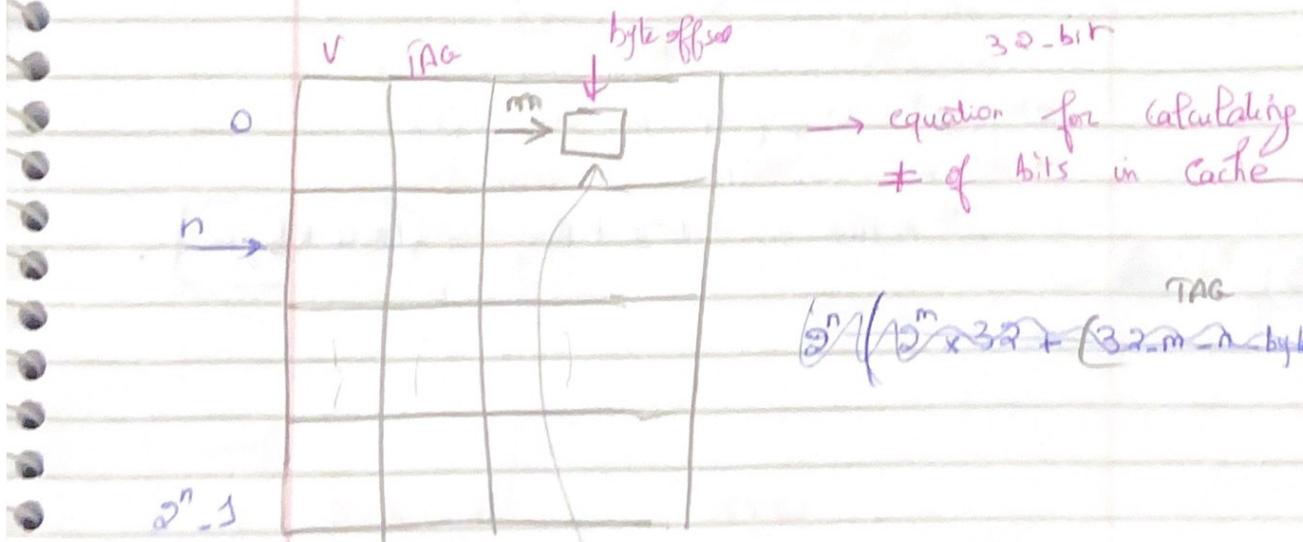


$$2^{14} = 2^{10} \times 2^4$$

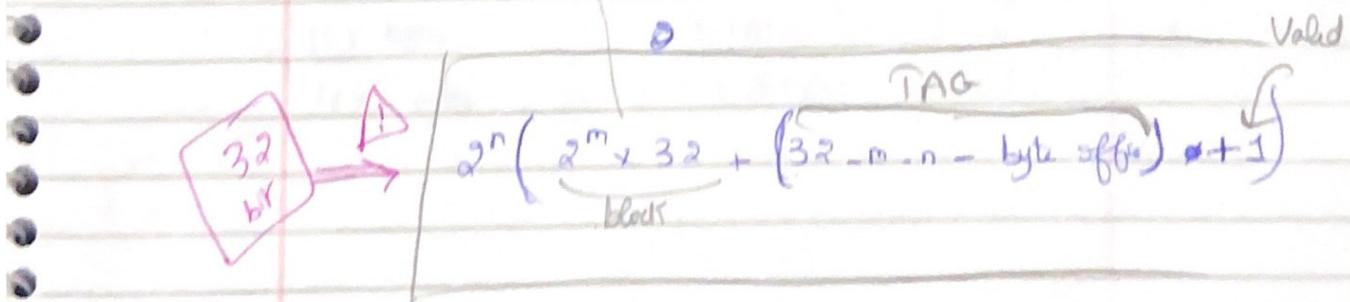
= 16 Kbytes

$$\frac{16\text{ Kbytes}}{16\text{ bytes}} = (1\text{ Kbytes})$$

$$2^6, 2^6 = (2^{10})$$



$$2^n(2^m \times 32 + (32-m-n \text{ bytes}))$$



$$2^n(2^m \times 32 + (32-m-n \text{ bytes})) + 1$$

→ How many bits are required for a direct mapped Cache with 16 KB 8 (1 word \Rightarrow 16 bytes) block 32 bit address $\left(2^4\right)$

byte offset = ?

$m = ?$

$$\frac{16 \text{ Kbytes}}{16 \text{ bytes}} = 1 \text{ K block}$$

$$2^{10} \Rightarrow n = 10$$

$$2^{10} \left(2^7 \times 3^2 + (3^2 - 10 - 9 - 2) + 1 \right)$$

$$= 2^{10} \times 147 = 157 \text{ Kbits} = 18.4 \text{ Kbytes}$$

$2^m \rightarrow$ word.

Cache sizes:

32 bits

$$2^n \left(32 \times 2^m + (32 - n - m) \underbrace{(2)}_{\text{1 for each block.}} + 1 \right)$$

nr. of blocks

block size in bits

byte offset

Valid bit

64 bits

$$2^n \left(64 \times 2^m + (64 - n - m) \underbrace{(2)}_{\substack{\text{Data} \\ \text{bytes}}} + 1 \right)$$

Data

Address

8 bytes in 64 bits

etc

64 blocks

\Rightarrow words

n (blocks)

6

m (word)

1

$$m = \cancel{2^6} = 2^1 - 2$$

32 bits data:

$$= \frac{16 \text{ K bytes}}{4 \times 4 \text{ bytes}}$$

$$\begin{cases} \text{Word: } 4 \\ 2^2 \end{cases}$$

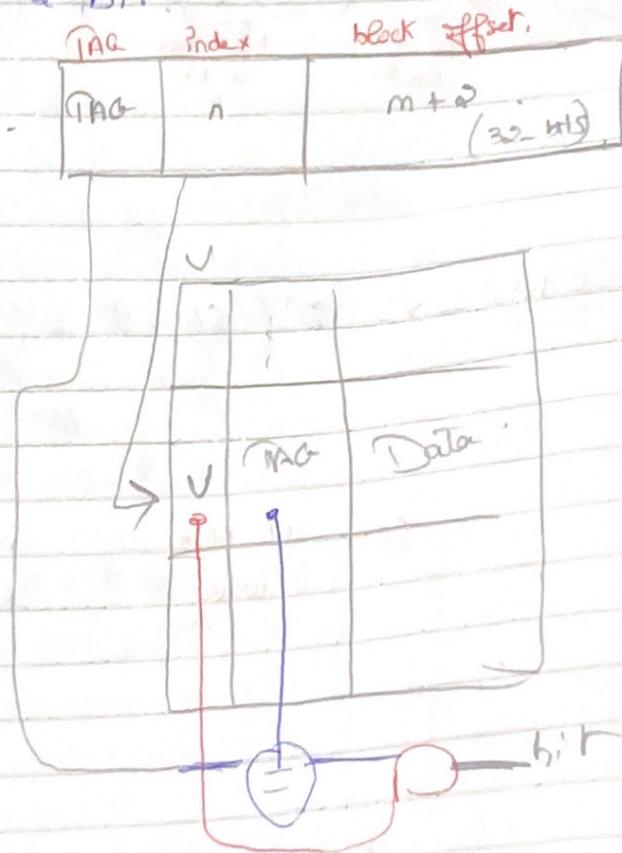
64 bits

$$= \frac{16 \text{ KB}}{4 \times 8 \text{ bytes}}$$

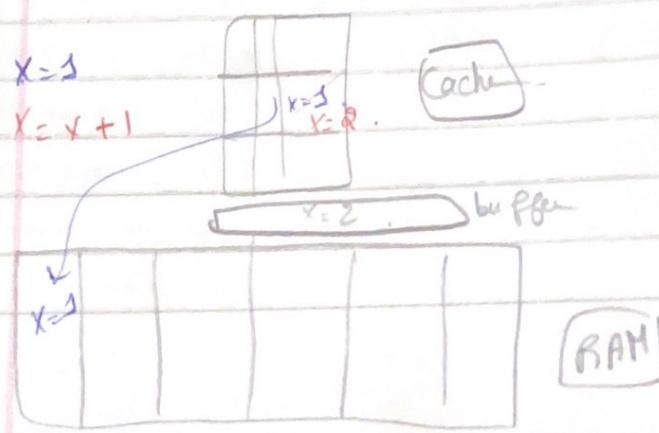
$$\left. \begin{array}{l} m \rightarrow \text{Same} \\ n \rightarrow 1 \end{array} \right\}$$

Cache

→ One comparator is needed to compare TAG bits found in address with the one pointed by the index in cache, if valid bit is one then we have a hit.



→ Write Policies?



→ Write-through
Processor ask for X

- ① it reads from RAM
- ② Change the value, now we have inconsistency between RAM & Cache.
- ③ One changed, it will also update RAM.
This means each write will take time to update RAM.

→ Why write-through? Without buffer is not a good idea ??

~~ex:~~ Suppose we have 10% store CPI without cache miss is 1.

Time to write to RAM is 100 cycles

What would be the new CPI with a write-through Cache with No buffer

$$CPI = 1 + 10\% \times 100 \text{ cycles} = 3.1$$

→ Write buffer will not be beneficial if the rate the Processor generates writes (SW) is more than the rate, the buffer can empty its content in RAM.

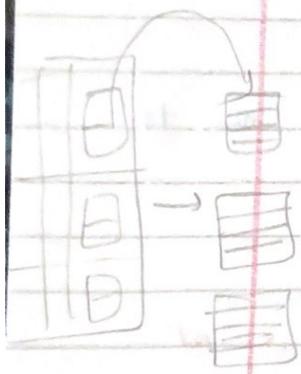
→ When buffer is full & we have an SW → The Processor will STOP.

block size? → Rate \rightarrow pen \rightarrow

→ blocksize \rightarrow miss Rate \rightarrow pen \rightarrow

→ Write back; instead of writing to RAM write each SW, & keep the content changed in cache until a block will replace the current content.

a Dirty is used to indicate if a block has ~~has~~ its content modified.



→ Assume you have a fixed Cache Size ex. 64 blocks & it is a good idea to increase block size.

+ if we increase block size this means more words in cache block, based on spatial locality, this will decrease miss rate however this will increase miss penalty.

→ increase block size will lead to less locations in cache to be searched & higher competition between blocks to contend for block replacement will be high this making miss rate higher.

Write allocate \rightarrow cache
Cache
Cache Full.

→ Write Allocate? No write \rightarrow cache
Cache Full.
Cache Full.

Consider a miss in a Write-through Cache

in this case \Rightarrow fetch the block from RAM.

Allocate it in Cache & Perform the write.

No write allocate - the Processor write an entire block of data in memory but not in cache.

* This is useful when O.S. needs to zeros a page in memory.

⇒ Measuring Cache Performance technique

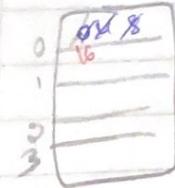
- * Reduce miss rate to reduce the probability that two different memory blocks contend to the same cache location.

- * Reduce miss penalty by using multi level cache

CPU time is divided into clock cycles that the CPU spends executing the program & clock cycles that CPU spends waiting for the memory system.

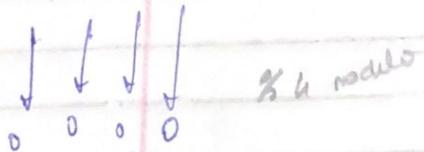
Cache Performance

- miss rate: Cache design reduce number of blocks contending for the same location



- miss penalty: Multi-level cache

0, 4, 8, 16, 0 4, 8, 16.



→ CPU Time = CPU Execution cycles + Memory stall cycles

Time waiting for
Cache to retrieve
Data from Memory

→ Memory stall cycles = Read stall + Write stall

$$\text{Read stall} = \frac{\text{Reads}}{\text{Program}} \times \text{miss Rate} \times \text{miss Penalty}$$

Same

→ Perfect cache is cache with no misses.

$$\text{Write stall} = \frac{\text{Writes}}{\text{Program}} \times \text{miss Rate} \times \text{miss Penalty}$$

+ write buffer stalls
(negligible)

D-cache → Data memory stalls = $\frac{\text{Memory access}}{\text{Prog.}}$ $\times \text{miss rate} \times \text{miss penalty}$

I-cache → Instruction memory stall = $\frac{\text{Instruction}}{\text{Prog.}}$ $\times \text{misses/instruction} \times \text{miss penalty}$

Assume that the miss rate of an instruction cache is (2%) & miss rate of data cache is (4%).

I-cache D-cache If a processor has a CPI of $\frac{2}{3}$ without memory stalls & miss penalty is 100 cycles.

$$CPI = 2$$

Determine how much faster a processor would run with a perfect cache that never misses. Assume the frequency of

All loads & stores is 36% $\xrightarrow{\text{All}} \text{w/sw}$.

$$\rightarrow \text{Instruction cache stalls} = \frac{1}{I} \times 2\% \times 100 \text{ cycles} \\ = [2I \text{ cycles}]$$

$$\rightarrow \text{Data cache stalls} = \frac{1}{I} \times 36\% \times 0.04 \times 100 \\ = [1.44I \text{ cycles}]$$

$$\rightarrow \text{CPI}_{\text{normal with cache}} = \text{perfect} + \frac{\text{CPI}}{I} + \frac{\text{CPI}}{I} \\ = [2 + \frac{2I}{I} + \frac{1.44I}{I}] \\ = [5.44]$$

$$\% \text{ faster} = \frac{5.44}{2} = 12.72 \text{ Time}$$

→ Perfect $\rightarrow \begin{pmatrix} 3 \\ 0 \\ 0 \\ 0 \end{pmatrix}$
miss

Δ → If we speed up the Processor
 $CPI \rightarrow 3.$

→ What will be the effect on the overall system

$$\begin{aligned} CPI \text{ with cache} &= 1 + 2 + 1.44 \\ &= 4.44 \end{aligned}$$

slower than Perfect Cache.

$$\% \text{ slower} = \frac{4.44}{1} = 4.44$$

CPI
 $5.44 \geq 4.44$

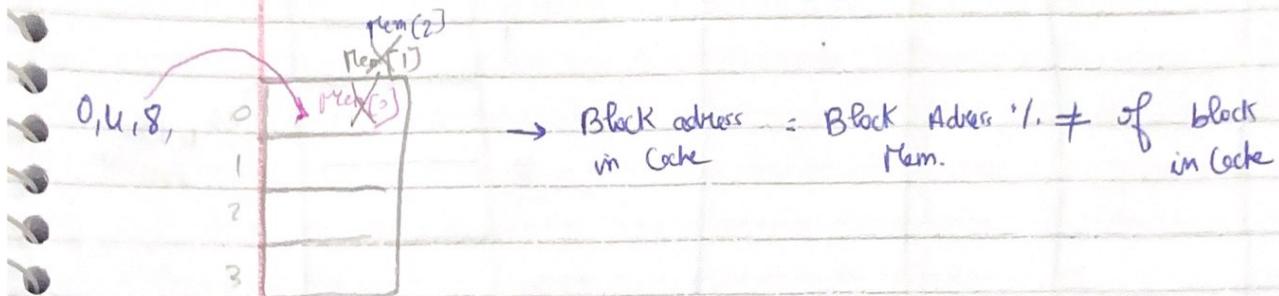
a.
100% miss
10% f. miss
80% hits
10% d. miss

AMAT ? Average Memory Access Time?

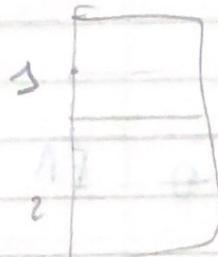
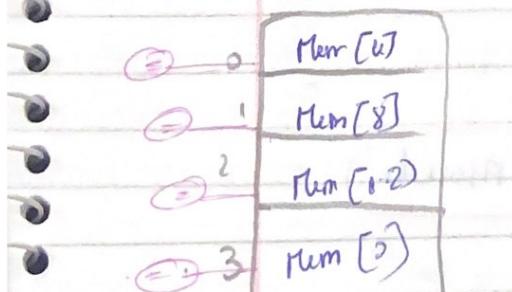
→ AMAT = time for a hit + Miss rate \times Miss Penalty.

Direct Mapped Cache:

- Every block has a specific location
- very fast in terms of addressing
- miss Rate $\uparrow\uparrow$



→ instead of every block going into a specific location why not ALL blocks can go in ANY location.

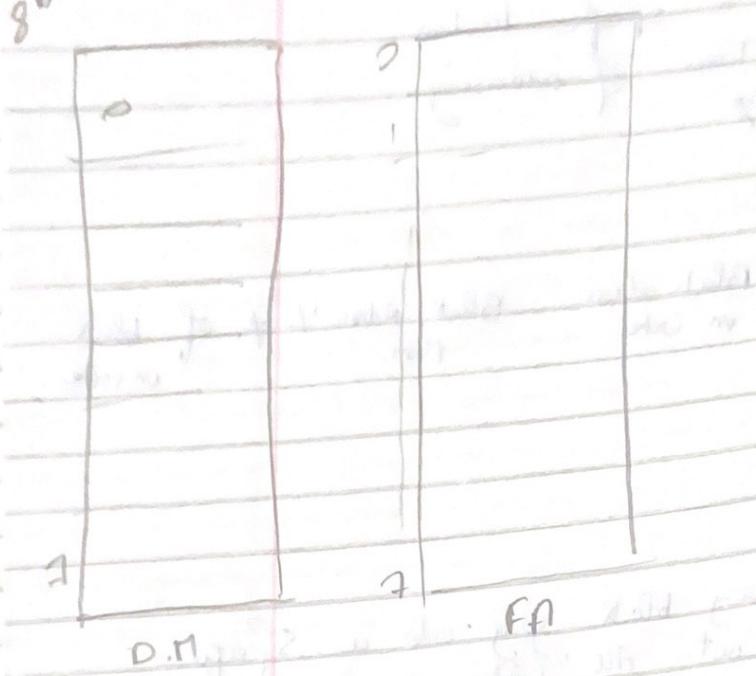


8 way → fully
 \leftrightarrow direct mapped

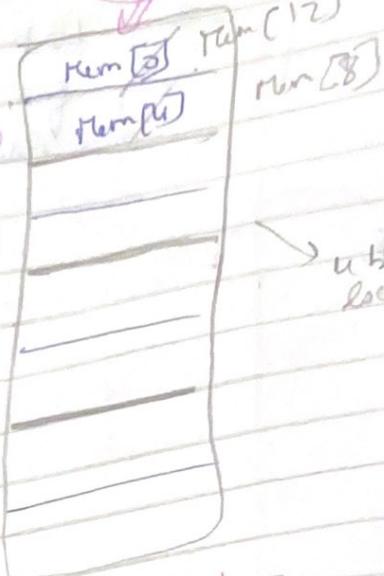
$0, 4, 12, 8$

8 blocks

1 word block?

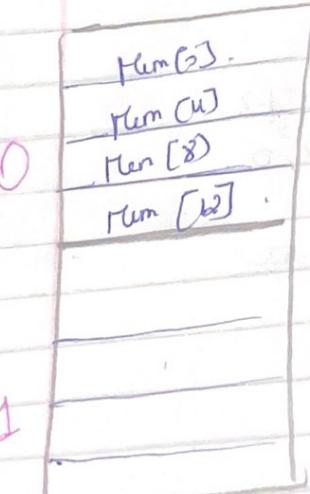


✓ 2 way set associative



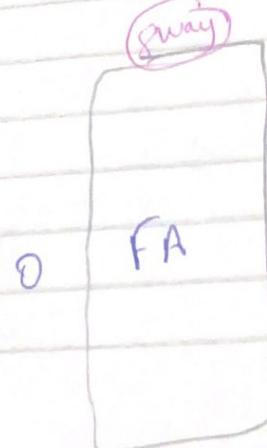
→ block locator

Address in cache = Address in memory / 4



⇒ 4 way set associative

↓
 block locator



(8 way)