

14/07/12

Computer Organization

83

Ref book:

COA by Morris Mano

COA by Paul Chowdary

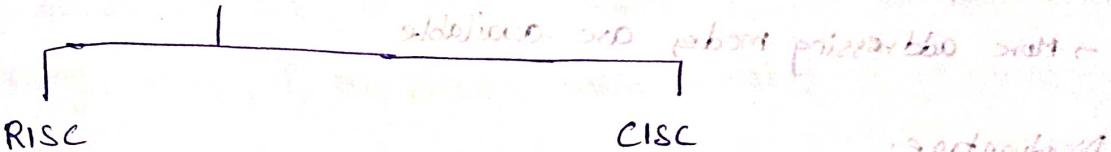
Computer Architecture is about

i) Instruction

ii) Addressing modes

iii) Data Format.

* Capacity of instruction



(Reduced Instruction Set Computer)

(Complex Instruction Set Computer)

RISC: Necessary condition for an instruction to be in RISC is

$$CPI \geq 1$$

not between either of (multiple) programs diff to execute

Clocks Per Instruction

should be same for both programs to be executed

since every instruction must be finished in single clock cycle,

RISC has ~~below Disadvantages~~

* limited capacity

* Fixed instruction size

* Limited addressing mode (as devices can be addressed only with direct addressing mode)

Advantages:

→ Due to fixed size instruction,

* It is simple to design compiler

* They are more compatible with instruction pipeline

→ RISC processors have more CPU registers

∴ Efficient parameter parsing (this because of the

fact that more registers provides register window)

→ RISC systems are used in Hardware Control Unit design
Programmer overhead is more in RISC eg Processor.

CISC:

- Has more instructions
- CPI ≥ 1
- Has More Capacity
- Instruction size is not fixed (i.e., variable length)
- More addressing modes are available

Disadvantage:

- * Complexity is more while designing a Compiler
- * Less compatible for instruction pipeline.
- * Less efficient works for parameter passing (\because No of CPU registers is less and hence no register windows)
- Most of the memory (registers) is mainly devoted for special purposes. Hence register window is not available.

Advantages:

- * Programmer overhead is less
- Let programmer overhead in RISC & CISC be denoted by P and compiler design complexity in RISC & CISC be denoted by R . Then we have below relations

$$P > R, R \ll S, P \ll S$$

$$P > R, Q \ll S$$

Addressing Modes:

It is about how we access operands or how operand reference is given in the instruction.

Flexibility or No of different addressing modes.

~~Ques~~ ∵ CISC is more flexible than RISC

Data Format:

It is about how we interpret binary strings

ASCII — uses 7 bits for Alpha Numeric

EBCDIC — uses 8 bits for Alpha Numeric

+ Extended Binary Coded Decimal ~~Interchange~~ for Information Interchange

→ ASCII uses 8th bit for parity through which can do self error

checking (Checking if the word is valid or not) bit given by 8th bit

→ In ASCII

W.A.

$$2^7 = 128$$

still we have left over combinations which can be left

for future use.

bit mapping bits in combination based program to fit in

→ Flynn Classification of Computer ~~Architecture~~

Bob Bob does not know what happens in architecture but

soothingly it is now that insights from serial and

Instructions &
Memory organization

Memory store for
program & Data

Flynn Classification

S	M	Instructions
Single	Multiple	
(Multicore) N ALUS SIMD (with diff memory modules) (vector processing)	n-CUs m-ALUs MIMD (Memory modules)	translating of stream 1 CU 1 Memory
Sequential processing	MISD 1 ALU	multiple ALU

→ Instruction is issued by CU ∵ multiple instructions → multiple CUs

→ Multiple data ⇒ More ALUs are required.

Von-Neumann Architecture:

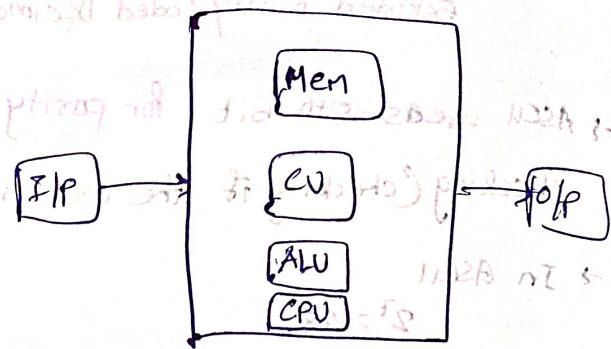
Here program and data are stored in same memory. This is called stored program concept.

i.e., program & data has to be stored before any processing.

Infact it is the first-memory based system. $112A \rightarrow 112B$
 $112B \rightarrow 112C$

→ This stored program concept supports inplace modification.

i.e., modifying the program or data itself while processing.



Harvard Architecture:

→ It is a memory based architecture in which program and data are stored separately.

This separation produces parallel access to both code and data and hence more efficient than Von-Neumann architecture.

Note: Since program and data are stored in different memory, they can be modified independently.

The disadvantage with Von-Neumann architecture is when program size is larger than memory (we cannot store all the program in the memory) and thus secondary memory came into the picture.

→ Since complete program is now not visible to primary memory, the added memory is called

Virtual Memory (Secondary memory)

Computer Organization

It deals with physical devices and their interconnections to improve performance.

→ performance: work done per unit time

Million Instructions Per Second (MIPS) is a measure of performance.



Computer Organization

(Q1) Performance gain of the machine M₁ over M₂ is 3. With respect to the following statements identify correct answer.

S₁: M₁ takes 3 times the time that of M₂.

S₂: M₁ increases performance by 200%.

S₃: M₁ is 3 times faster than M₂.

which are true?

- a) only S₁, S₂ b) S₂, S₃ c) S₁, S₃ d) All

(Q2) A system M₁ when operated with 1GHz clock takes 5 clocks per instruction on average. It runs the program of 1000 instructions.

i) what is total time to run the program.

$$CPI = 5$$

1000 instructions \Rightarrow 5000 clocks

Time period = 1 ns

$$\Rightarrow \text{total time} = 5000 \text{ ns}$$

$$= 5 \mu\text{s}$$

ii) What is its performance?

Performance = no. of instructions per second

5000 ns \rightarrow 1000 instruction

1 ns
(10^9 ns) \rightarrow (Performance) ?

$$\text{Performance} = \frac{1000}{5000} \times 10^9$$

$$= \frac{10^9}{5} = 200 \times 10^6 \text{ Inst/sec}$$

$$= 200 \text{ MIPS}$$

(iii) Its improved version M₂ uses 20% faster clock, what is performance of this new machine?

Sol:

$$\text{Performance of } M_2 = \text{Performance of } M_1 + 20\% \text{ (Performance of } M_1)$$

$$= 200 + 20\% (200)$$

$$= 240 \text{ MIPS}$$

(iv) Performance gain of M₂ over M₁ is _____

$$= \frac{\text{Performance of } M_2}{\text{Performance of } M_1} = \frac{240}{200} = 1.2$$

i.e., performance of M₂ is 1.2 times that of M₁

Note:

* $t = n * CPI * T$

t → execution time of program

n → no of instructions

CPI → clocks per instruction

T → Time period (per one clock)

* performance = Instructions per second

(Q) A processor P₁ takes 2 clocks per instruction and runs at 1GHz clock. Its improved version processor P₂ takes 2/3 of time of P₁ but increases CPI by 20%. What is clock rate of processor P₂?

Sol:

$$CPI_1 = 2 \quad CPI_2 = 2 + 0.2 \times 2 = 2.4$$

$$f_1 = 1 \text{ GHz} \quad f_2 = ?$$

$$\text{Performance of } P_2 \text{ over } P_1 = \frac{\text{Time of } P_1}{\text{Time of } P_2} = \frac{\text{throughput of } P_2}{\text{throughput of } P_1}$$

let P a program be with n instructions then

$$\frac{2}{3} (\text{time by } P_1) = \text{time by } P_2$$

$$\frac{2}{3} \times n \times CPI_1 \times T_1 = n \times CPI_2 \times T_2$$

$$\frac{2}{3} \times \frac{CPI_1}{f_1} = 1.2 \times \frac{CPI_1}{f_2}$$

$$f_2 = \frac{0.4}{1.2 \times f_1} = \frac{3}{2} \times 1.2 \times f_1$$

$$f_2 = \frac{0.8 \times 1}{0.8 \times 1.2} \Rightarrow f_2 = 1.8 \text{ GHz}$$

(Q4) A dual core ~~processor~~ machine runs the program containing 60% parallel portion. What is the performance gain of dual core processor compare to single core processor

Sol:

for single core processor

100% time

for dual core processor

$\frac{40\% + 60\%}{2}$

total % time, hence $40\% + 30\% = 70\%$

Performance of dual core over single core = $\frac{100}{70} = \frac{10}{7} \approx 1.4285$

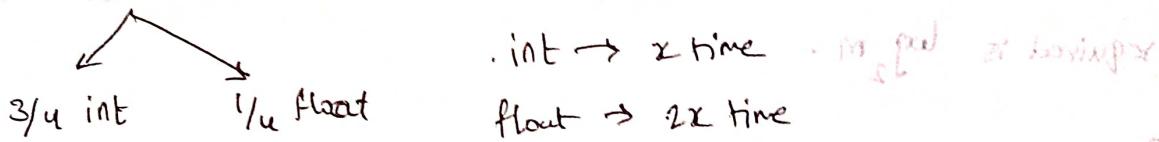
i.e., performance is increased by 42.85 %

(Q5) A novice processor takes twice the time to implement floating point operation compared to integer operation. However its improved new version takes only $\frac{2}{3}$ of time for floating operation. If the program is the mix of 3:4 integer and floating point operations respectively what is the speedup of new processor compared to novice one?

Sq:

let t_1 be time taken by 1st processor

let n be no. of instructions



$$t_1 = \frac{3}{4}x + \frac{1}{4}(2x) = \frac{11}{4}x$$

let t_2 be time taken by 2nd processor

int $\rightarrow x$ time

$$\text{float } \rightarrow \frac{2}{3}(2x) = \frac{4}{3}x$$

~~$t_2 = \frac{3}{4}x + \frac{4}{7} \cdot \frac{1}{3}x = \frac{9}{14}x + \frac{16}{21}x = \frac{25}{21}x$~~

gain of P_2 over P_1 =

~~$$= \frac{t_1}{t_2} = \frac{\frac{11}{4}x}{\frac{25}{21}x} = \frac{33}{25} = 1.32$$~~

Memory Organization:

Words & Bits:

- Memory is accessed with address bus & databus.
- words are addressed using address bus.
- The unit of addressing in memory is word which does not need to be 8-bit unit.
- The memory ^{in which successive addresses refer successive bytes} ~~with 8-bit words~~ is called byte addressable memory.
- If no of ~~address bits~~ to words are m , then no of address bits required is $\log_2 m$.

Eg :

256 word memory need address size of 8 bits.

- when memory word is addressed, then all bits of word are accessed in parallel using databus. (Normally databus size is of the word size)
- Address bus is always unidirectional (Address is given to memory but memory never gives any address)
- Data bus can be unidirectional or bidirectional

(read only)	(read & write)
ROM	RAM
- Control bus dictates whether to perform read operation or write operation on addressed memory word.
- Address bus, Data bus & control bus together are called system bus.

Note: (a)

Consider 64 MB memory. It is can be designed in following ways:

- 1. $64M \times 8$ bits (64 million words and 8-bits/word)

$$64M \text{ words} \rightarrow 2^6 \times 2^{20} \text{ words} = 2^{26} \text{ words}$$

∴ address bus is of 2^{26} bit size

8 bits/word \rightarrow data bits are 8 or data bus size ≥ 8

2. ~~32M~~ $32M \times 16$ bits (32 Million words and 16 bits/word)

$$\text{no of words} = 32 \times 2^{20} = 2^{25} \text{ words}$$

$$\therefore \text{size of address bus} = 2^{25}$$

$$\text{size of data bus} = 16$$

3. ~~16M~~ $16M \times 32$ bits (16 Million words and 32 bits/word)

$$\text{no of words} = 16 \times 2^{20} = 2^{24} \text{ words}$$

$$\therefore \text{size of address bus} = 2^{24}$$

$$\text{size of data bus} = 32$$

4. $8M \times 64$ bits (8 Million words and 64 bits/word)

5. ~~4M~~ $4M \times 128$ bits (4 Million words and 128 bits/word)

Ques
address

Q6) In a system design 4GB main memory is organized as 64 bit-words (Eight 8 bit words).

What will be the number of address bits required to access the memory.

Sol:

$$4\text{GB memory} \Rightarrow 4 \times 2^{30} \times 8 \text{ bits}$$

bus

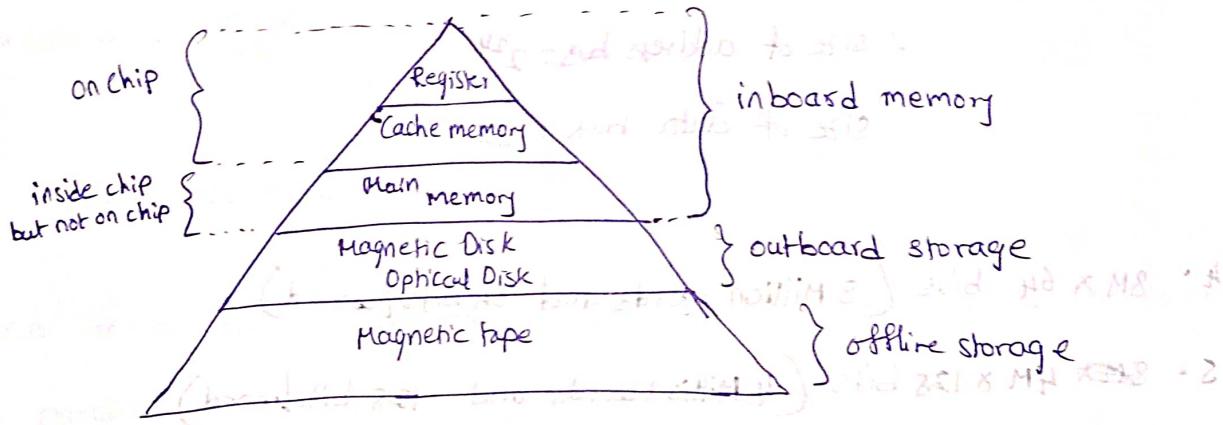
$$\text{each word is } 64 \text{ bits} = 2^6$$

$$\therefore \text{no of words} = \frac{4 \times 2^{30} \times 8}{2^6} = \frac{2^2 \times 2^{24} \times 2^3}{2^6} = \frac{1 \times 8 \times 16}{2^3} = 2^{29} \text{ words}$$

$$\therefore \text{no of address bits} = 29$$

Memory Hierarchy

- Memory hierarchy is arrangement of different types of memories of different sizes with the objective of reducing average access time.
- It ensures faster memories are accessed more frequently than slower memories.
- Hierarchy is logical arrangement (not physical) so the reference for accessing memory is given by word address.



16/07/20

- (Q7) Given the following design statements identify the most performing design

S1: Machine with 1GHz clock and CPI=5

S2: Machine with 2GHz clock and CPI=7

S3: Machine with 5GHz clock and CPI=10

- a) S1 b) S2 c) S3 d) All S1, S2, S3 have same

Sol:

$$t_1 = n \times 5 \times \frac{1}{10^9} = 0.5 \times 10^{-9}$$

$$t_2 = n \times 7 \times \frac{1}{2 \times 10^9} = 3.5 \times 10^{-9}$$

$$t_3 = n \times 10 \times \frac{1}{5 \times 10^9} = 2 \times 10^{-9}$$

$t_3 < t_2 < t_1$ \therefore S3 has highest performance

Ques. No. 1

→ It says that mostly 10% of the program uses 90% of the resources used by the program.

→ This 10% is called locality of reference.

→ The locality of reference is generally maintained in cache.

Ques. In a system design 20 bit address bus is used and 8 bit bidirectional data bus is maintained. Identify the correct among the following options.

a) System uses 1GB RAM

b) System uses 1MB RAM

c) System uses 1GB ROM

d) System uses 1MB ROM

Sol:

Given bidirectional data bus \Rightarrow RAM but not ROM

20 bit address bus \Rightarrow 2^{20} words

$$\text{Memory} = 2^{20} \times 8 \text{ bits}$$

$$= 2^{20} \text{ bytes} = 1 \text{ MB RAM}$$

\therefore opt b

→ Based on how close the memory to the processor the memories are assigned levels

register memory — level 0

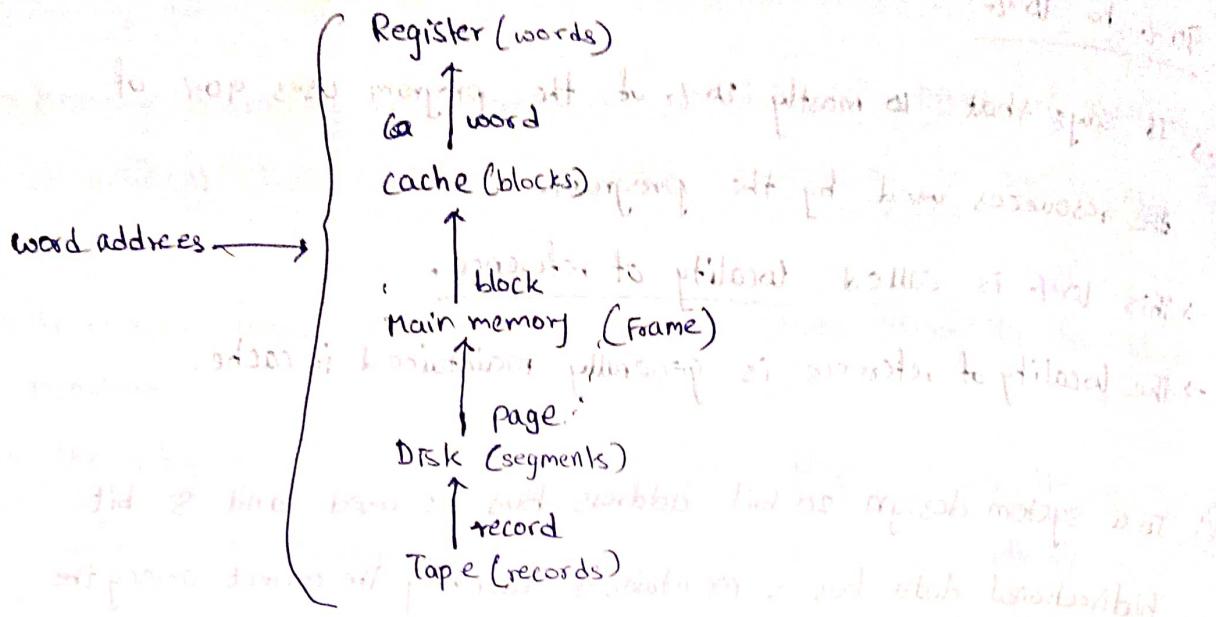
cache memory — level 1

Main memory — level 2

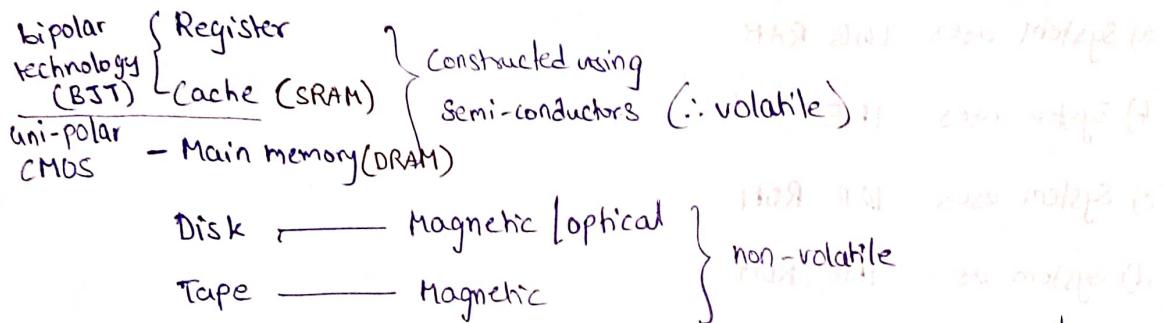
Disk Memory — level 3

Tape Memory — level 4

Readjustment in memory hierarchy



Behaviour:

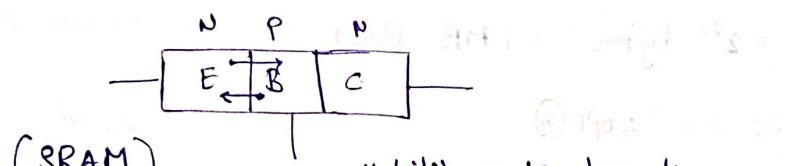


CMOS - Complementary Metal Oxide

Semiconductor set

BJT - Bi-Junction Transistor

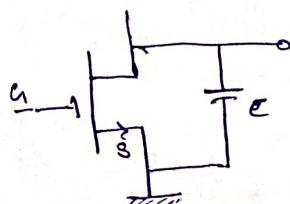
Bi Junction Transistor:



Mobility is in two directions

(Static RAM) and hence memory is faster

CMOS



Only one type of carriers are moving

Hence mobility will be less.

The charge across the capacitor can't stay infinitely.

∴ It is called Dynamic RAM (DRAM)

Before the content is lost the capacitor needs to be charged again.
This process is called Refreshing.

SRAM	DRAM
(i) Faster due to bipolar technology	(i) Slow due to unipolar technology
(ii) Volatile	(ii) Volatile
(iii) Less package density (∴ we need two carriers for each bit)	(iii) More Package Density (∴ we need one carrier for each bit)
(iv) Refreshing is not required	(iv) Refreshing is required regularly

Package density: no of bits that can be accommodated over given area

Properties of Memory Hierarchy:

1. In the hierarchy level 'i' memory is placed on top of level 'i+1' memory.
2. These levels are arranged such that
 $T_1 < T_2$; $S_1 < S_2$ & $C_1 > C_2$; $f_1 > f_2$
(time) (size) (cost) (cost/byt)e
(frequency of accessing)
3. Container i.e., Information in level 'i' \subset Information in level 'i+1'
Inclusion
4. Performance of memory hierarchy is given by hit ratio, h
(i.e., availability of referred information at referred level)
 $h \propto \text{size}$

$$h_i \propto \frac{1}{T_{avg}} \quad (\text{Avg access time})$$

5) The drawback is data inconsistency

i.e., same information will have different values at different memories

+ proper write operation would reduce this problem.

Note: ~~parallel operation of sub-memories~~

If 'h' is hit ratio then

If h_i is hit ratio of level i then

$$h_0 < h_1 < h_2 < h_3 < h_4 \quad \& \quad h_4 = 1$$

(for direct address and block unit)

→ Consider a two level memory system with

level $L_1: T_1, S_1, C_1$

level $L_2: T_2, S_2, C_2$

$S_T = \text{total size of memory} = S_1 + S_2$

$C_{avg} = \text{Avg cost of memory system} = \frac{\text{total cost}}{\text{total size (in bytes)}}$

$$= \frac{S_1 C_1 + S_2 C_2}{S_1 + S_2}$$

$$= \frac{\sum S_i C_i}{\sum S_i}$$

→ Consider two 2-level memory system with

'hit' based organization $\Rightarrow L_1: T_1 = 10 \text{ ns}, h_1 = 0.8$ and $L_2: T_2 = 100 \text{ ns}, h_2 = 1$

Find avg access time of external memory

Sol:

$$T_{avg} = 0.8 \times 10 \text{ ns} + (1 - 0.8) \times 100 \text{ ns}$$

$$= 8 \text{ ns} + 20 \text{ ns} = \underline{28 \text{ ns}}$$

(Q9) Given the following design statements identify most performing design

S1: 256MB cache, 8GB RAM (MM), 1TB ~~HDD~~ HDD

S2: 512MB cache, 4GB RAM (MM), 500GB HDD

S3: 1GB cache, 8GB RAM (MM), 500GB HDD

- a) S1 b) S2 c) S3 d) All have same performance.

Sol:

for we need to check which system has more cache, then we check main memory and then HDD

\therefore S3 gives better performance

\therefore opt C

(Q10) In a three level memory hierarchy the hit ratios are h_1, h_2, h_3 identify most appropriate one among the following.

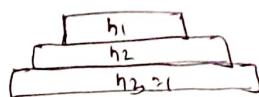
a) $h_1 > h_2 > h_3$ and $h_3 = 1$

b) $h_1 > h_2 < h_3$ and $h_3 = 1$

c) $h_1 < h_2 < h_3$ and $h_3 = 1$

d) $h_1 < h_2 < h_3$

Sol:



$\therefore h_1 < h_2 < h_3 \& h_3 = 1$

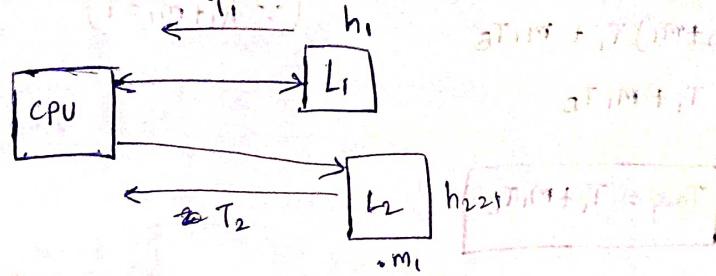
Average access time Calculation:

Case 1: (2-Level)

Processor is having bus for cache memory and main memory
(Processor is having access with both level)

Let $L_1: S_1, C_1, T_1$

$L_2: S_2, C_2, T_2$



m_1 is miss ratio of level 1

$$T_{avg} > h_1 T_1 + m_1 T_2$$

$$\Rightarrow T_{avg} = h_1 T_1 + (1-h_1) T_2$$

T_{avg} is always for one word (i.e., time/word)

∴ Data transfer rate/throughput = $\frac{1}{T_{avg}}$ words/sec

Case 2: Exploring Locality of reference)

i) Block size = one word

ii) Block size > one word

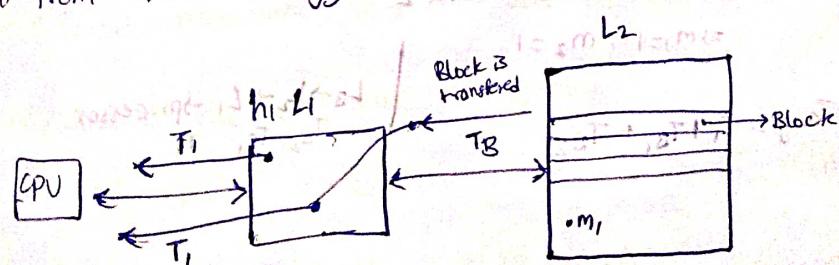
→ while solving problem if case is not mentioned we need to check for both the cases - however mostly in question case (ii) will be given

→ Processor is having access with faster memory; first a block of words are

→ If there is a miss at faster memory, first a block of words are moved from slower memory to faster memory and addressed word

is given to processor from faster memory

* word is always supplied from faster memory to CPU (it is never sent from slower memory)



$$T_{avg} = h_1 T_1 + m_1 [T_B + T_2]$$

$T_B \rightarrow$ time taken to move block from memory to cache

forming

then we

are h_1, h_2, h_3

memory

$$\begin{aligned} &= h_1 T_1 + M_1 T_B + M_1 T_1 \\ &= (h_1 + M_1) T_1 + M_1 T_B \\ &= T_1 + M_1 T_B \end{aligned}$$

($\because h_1 + M_1 = 1$)

$$\therefore \boxed{T_{avg} = T_1 + M_1 T_B}$$

Block size can be one word or n words

if block size = 1 word

$$\Rightarrow T_B \geq T_2$$

if block size = n words

$$\Rightarrow T_B = n * T_2$$

Extension for multiple levels

$$T_{avg} = T_1 + m_1 T_{B_2} + m_1 m_2 T_{B_3} + m_1 m_2 m_3 T_{B_4} + \dots$$

for 3 levels private stores

$$T_{avg} \geq T_1 + m_1 T_{B_1} + m_1 m_2 T_{B_2}$$

with condition (i) and (ii)

i) hit at level 1: $h_1 = 1 \Rightarrow M_1 = 0$ thus private store is not required

$$T_{avg} = T_1 \quad (\text{Level gives the word})$$

blocks before bno promote instead of private words most bottom

ii) miss at level 1 & hit at level 2

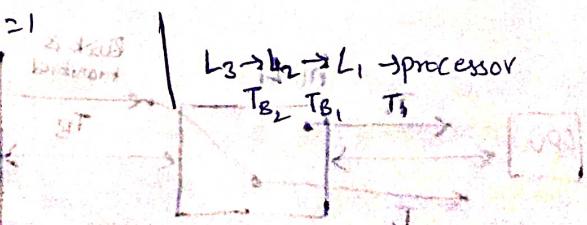
$$\Rightarrow m_1 = 1, h_2 = 1 \Rightarrow m_2 = 0$$

$$T_{avg} \geq T_1 + T_B,$$

iii) miss at level 1 & miss at level 2

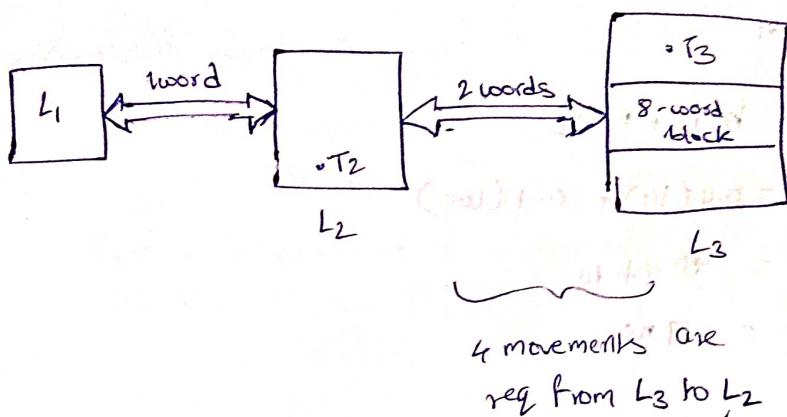
$$\Rightarrow m_1 = 1, m_2 = 1$$

$$T_{avg} = T_1 + T_{B_1} + T_{B_2}$$



→ Locality of reference will reduce penalty (extratime)
 But it reduces access time for future references (not for current ref)

Placement time Consideration:



This is done as follows:

1st one word is moved from L3 to L2.

once the previous word is made available in L2 then only start next two word unit.

Placing a word from data bus to memory is called placement time.

In the previous two cases we didn't consider the placement time.

T₃ - Read from L₃ → Now data is in the bus (not in L₂)

T₂ - Then place the ~~the~~ word in the bus in L₂

∴ Time taken for one word unit to move from L₃ to L₂ is T₃ + T₂

- Q11) In a 2-level system access time of level 1 is 10 times faster than L₂. L₁ maintains 90% hit and its access time is 10ns

i) what is the ~~average~~ access time

ii) If average access time is allowed at most 40ns, what is the minimum hit ratio at level 1?

Sol:

$T_1 \geq \frac{T_2}{10}$ (minimum value of first stage broadcast time)

$$T_1 = 10 \text{ ns} \Rightarrow T_2 = 100 \text{ ns}$$

$$h_1 = 0.9$$

$$m_1 = 0.1$$

$$T_{avg} = h_1 T_1 + m_1 T_2$$

$$= 0.9(10) + 0.1(100)$$

$$= 19 \text{ ns}$$

$$\therefore \boxed{19 \text{ ns}}$$

Case ii: If the processor has access only to faster memory then
 $\boxed{\text{fast}} \rightarrow \text{no cache or SMT}$

$$T_{avg} = h_1 T_1 + m_1 (T_1 + T_2)$$

$$= 0.9(10) + 0.1(10+100)$$

$$= 9 + 11$$

$$= 20 \text{ ns}$$

(ii) $T_{avg} \leq u_0$

$$\therefore h_1 T_1 + m_1 T_2 \leq u_0$$

$$h_1(10) + (1-h_1)(100) \leq u_0$$

$$10h_1 + 100 - 100h_1 \leq u_0$$

$$-90h_1 \leq -60$$

$$90h_1 \geq 60$$

$$h_1 \geq \frac{2}{3}$$

most strict constraint $\Rightarrow \frac{2}{3} \rightarrow \text{SMT access}$ $\boxed{\text{no cache or SMT}}$

Case iii: $\boxed{\text{SMT access in last slot}}$ \Rightarrow minimum broadcast time $\geq 10 \text{ ns}$
 $\therefore T_{avg} \leq u_0$

$$T_{avg} = h_1 T_1 + m_1 (T_1 + T_2) \leq u_0$$

$$\Rightarrow 10h_1 + (1-h_1)(110) \leq u_0$$

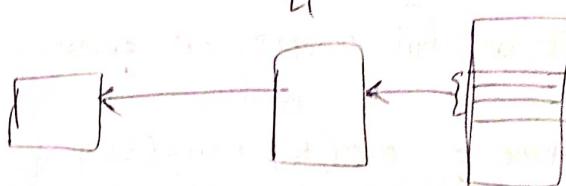
$$10h_1 + 110 - 100h_1 \leq u_0$$

$$-100h_1 \leq -70 \Rightarrow$$

$$\boxed{h_1 \geq 0.7}$$

27/08/20

- Q12 In a 2-level memory, level1 access time is 20ns and level2 access time is 200 ns. 90% of references present in L1. If there is a miss a 4-word block must be moved into L1 and addressed word is given to CPU later. what is avg access time?

Sol:

$$\begin{aligned} \text{avg access time} &= 0.9(20) + 0.1\left(\frac{4 \times 200 + 20}{4}\right) \\ &\quad \text{Block transfer to L1} \\ &= 18 + 0.1(820) \\ &= 18 + 82 = 100 \text{ ns} \end{aligned}$$

$$\text{Data transfer rate} = \frac{1}{\text{avg access time}} \text{ words/sec}$$

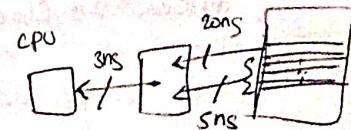
$$= \frac{1}{100 \text{ ns}} = \frac{10^9}{100} = 10^7 \text{ words/sec}$$

- Q13 A 2-level memory will have 3ns when hit occurs at level1. 94% of references are found in level1. If there is a miss at level1 the system takes 20ns to bring the first word from L2 to L1. Afterwards it takes 5 ns for each of the remaining words. Let block size is 32 words. Find avg access time.

Sol:

$$0.94(3) + 0.06(20 + 31 \times 5 + 3)$$

$$2.82 + 0.06(178) = 2.82 + 10.68 = 13.5 \text{ ns}$$



Q14 In a memory system 50ns time is taken if there is a miss at faster memory and 5 ns for a hit with faster memory. If 80% references are present in faster memory, find the avg access time.

Here individual memory access times are not given. Only times when hit or fail occurs, is given.

$$\therefore \text{avg access time} = 0.8(5) + 0.2(50)$$

$$= 4 + 10$$

$$= 14 \text{ ns}$$

Also performance or Data transfer rate

$$\text{Data Transfer Rate} = \frac{1}{14 \text{ ns}}$$

$$= \frac{100}{14} \times 10^7 \text{ words/sec}$$

$$= 7.14 \times 10^7 \text{ words/sec}$$

Q15 In above question, if word size is 8 bits what is the minimum bandwidth required for the bus.

Sol:

$$\text{Bandwidth} = (7.14 \times 10^7) \times 8 \text{ bits/sec}$$

$$= 57.12 \text{ MbPS}$$

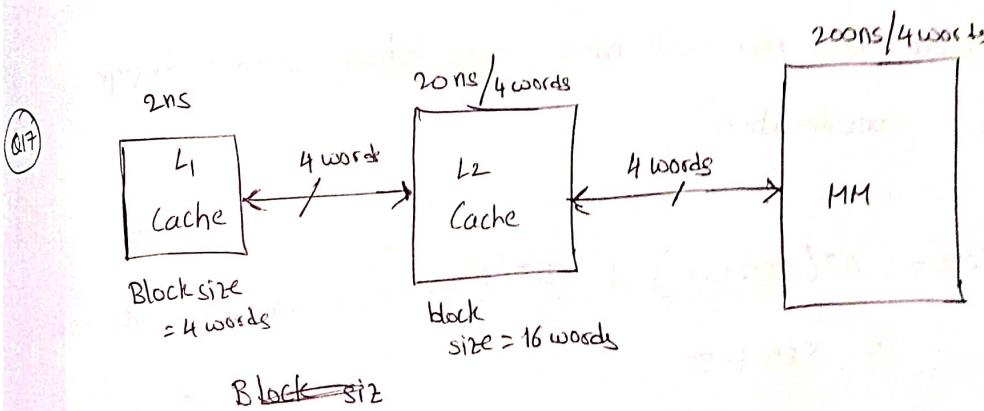
Q16 A memory system takes 2 clocks for cache hit to a read operation and 10 clocks for miss at cache for the same operation. For the write hit it takes 5 clocks and for miss it takes 20 clocks. If there are 60% read operations what is average access time? (take hit ratio as 80%)

Sol:

$$\begin{aligned}\text{avg access time for read} &= 0.8(2) + 0.2(10) \\ &= 1.6 + 2 \\ &= 3.6 \text{ clocks}\end{aligned}$$

$$\begin{aligned}\text{avg access time for write} &= 0.8(5) + 0.2(20) \\ &= 4 + 4 \\ &= 8 \text{ clocks}\end{aligned}$$

$$\begin{aligned}\text{total avg access time} &= 0.6(3.6) + 0.4(8) \\ &= 2.16 + 3.2 \\ &= 5.36 \text{ clocks}\end{aligned}$$



when there is a miss at L1 cache & hit at L2 cache a block

is transferred from L2 to L1. what is the time taken to transfer.

- a) 2ns b) 20 ns c) 22 ns d) 88 ns

Sol:

Since L2 is sending the block

block size = 16 words

4 words are sent at a time.

\therefore we need to move 4 moves.

Here we need to consider placement time also

i.e., 20 ns to ~~put~~ put it on the bus and 2 ns for placement time.

$$\therefore \text{total time} = 22 * 4 = 88 \text{ ns}$$

Q18 In above question, if there is a miss at both L₁ & L₂ caches the block must be moved from MM to L₂ cache and the concerned block must be moved from L₂ cache to L₁ cache. what is the total block transfer time.

Sol:

- a) 222 ns b) 888 ns c) 902 ns d) 968 ns

Sol:

total block transfer time

= block transfer time from MM to L₂ Cache

+

block transfer time from L₂ cache to L₁ cache.

The cache closer to MM will have same block sizes (design consideration)

$$\begin{aligned}\therefore \text{req time} &= 4 * (200 + 20) + 4 * (20 + 2) \\ &= 880 + 88 \\ &= 968\end{aligned}$$

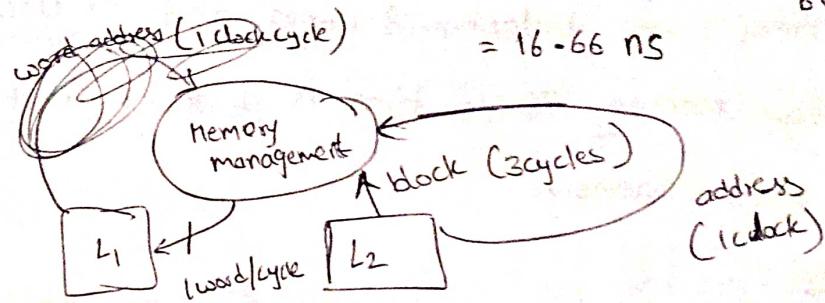
Q19 Certain processor operation with 60MHz clock and using 2-level memory. The block size is 8 words and each word contains 4 bytes. When there is a miss at level 1, the internal management system takes 1 clock cycle to accept first word address of the block and pull all remaining words of the block in 3 clock cycles from level 2, later it handover all the words to level 1 at the rate of 1 word per cycle. How many total clocks are needed to serve miss request by this processor and how much time it takes?

Sol:

$$\text{freq} = 60 \text{ MHz}$$

$$\text{time period of clock} = \frac{1}{60 \times 10^6} = \frac{1}{6 \times 10^7} = \frac{100}{6} \times 10^{-9} \text{ (or)} \frac{1}{60} \mu\text{s}$$

$$= 16.66 \text{ ns}$$



∴ total no of clocks = $1 + 3 + (8 * 1)$
 \downarrow
 address

$$= 12 \text{ clocks}$$

$$\text{time req} = 12 + \frac{1}{60} \mu\text{s} = \frac{1}{5} \mu\text{s}$$

$$= 0.2 \mu\text{s}$$

(Q20) In above question, the minimum bandwidth required for the bus connecting L2 and L1 is _____ 10^6 bytes/sec

Sol:

to transfer one block (8 words
 $= 32 \text{ bytes}) = 0.2 \mu\text{s}$

$$32 \text{ bytes} \longrightarrow 0.2 \mu\text{s}$$

$$160 \text{ bytes} \longleftarrow 1 \mu\text{s}$$

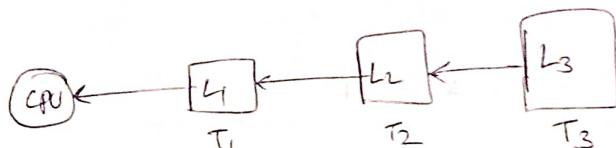
$$\Rightarrow 1 \text{ s} \longrightarrow 160 \times 10^6 \text{ bytes}$$

$$\therefore 160 \times 10^6 \text{ bytes/sec}$$

$$\therefore 160 \text{ MBPS}$$

(Q21) A processor design uses 3 level memory. The miss rate of level 1 memory is twice to that of level 2 memory. The access times are 1 clock, 8 clocks and 18 clocks respectively. If average memory access time is 2 clocks what are miss rates at L₁ and L₂ memory.

So:



$$t = T_1 + m_1 T_2 + m_1 m_2 T_3$$

$$t = T_1 + (2m) T_2 + (2m)(m) T_3$$

$$2 = 1 + 2m(8) + 2m^2(18)$$

$$\Rightarrow 36m^2 + 16m - 1 = 0$$

$$36m^2 + 18m - 2m - 1 = 0$$

~~$2m + 1 = 0$~~

$$18m(2m+1) - 1(2m+1) = 0$$

$$\Rightarrow m = -1/2 \quad \cancel{+ m} \quad \cancel{+ 1/18}$$

∴

$$\therefore 2m, m$$

$$\frac{2}{18}, \frac{1}{18}$$

$$\text{i.e., } \frac{1}{9}, \frac{1}{18}$$

$$\therefore 0.111 \text{ & } 0.056$$

28/08/20

Cache Memory

→ Cache is the smallest & fastest component in memory hierarchy
It bridges speed miss match b/w fastest processor to slowest
memory.

→ It maintains locality of reference thereby reduces average
access time.

→ Cache and Main memory are logically divided into fixed size
blocks.

→ the total no of blocks in MM are more compared to that of
no of blocks in the cache.

→ Hence it is required to know which main memory block
is to be placed and where in the cache.

for this we use Address mapping techniques.

→ Direct mapping is the simplest but results conflict problems.

→ Associative mapping resolves conflict problems but require
more hardware.

→ Set associative mapping is the combination of direct mapping &
associative mapping. (∴ it requires less hardware and less
conflict problems)

→ Tag bits: These bits are used for cache memory blocks.

They denote which MM block among the possible ones,
for that cache position present at that time.

→ Write operation results cache coherence problem (i.e., same
variable has different value in cache from memory)

→ proper updation techniques will reduce / resolve the cache coherence problem.

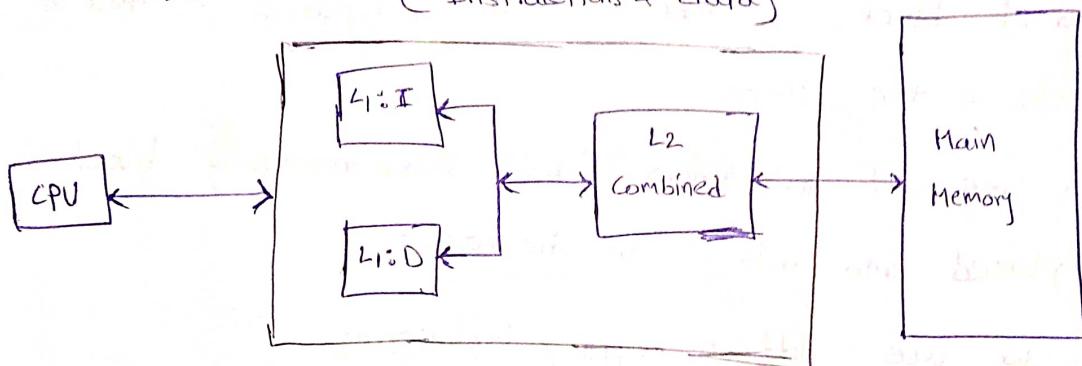
→ Block replacement techniques are aimed to minimize the no of block movement b/w cache & MM.
(e.g.: FIFO, LRU, Direct mapped cache)

Cache Design

① L₁ - Instruction cache (Holds only instructions)

L₂ - Data cache (Holds only data)

L₂ - Combined (Instructions & data)



* If CPU needs an instruction first L₁:I is checked, then L₂ and later main memory

• Only if CPU needs data then

L₁:D \longleftrightarrow L₂ \longleftrightarrow MM

So while accessing instruction, data cache (D-cache) is not involved. Only while accessing data, instruction cache is not involved.

② Physically indexed cache (while solving ques, this is ~~not~~ taken by default)

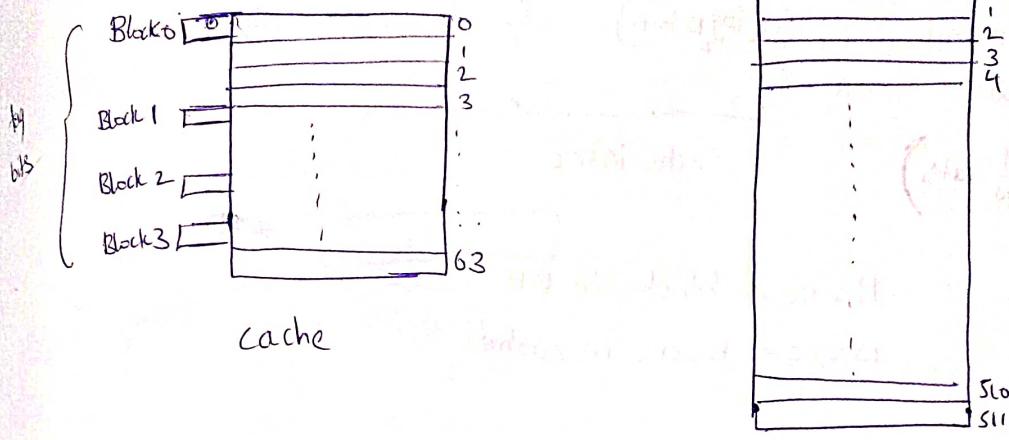
→ Main memory address is used in address mapping.
(Physical address)

→ Cache is referred using physical address.

3) Logical Indexed Cache / Virtually Indexed Cache:

- * This is TLB (Translation look ahead Buffer)
- * Virtual address is used to refer the cache.

Eg: Cache size: 64 words ; Main memory size: 512 words;
Block size: 16 words;



- no of blocks in main memory,

$$M = \frac{512}{16} = 32 \text{ blocks}$$

no of blocks in cache memory, $N = \frac{64}{16} = 4 \text{ blocks}$

Every block has tag bits.

Here we have four tag such tag

Direct Mapping:

→ Block 'k' of MM has to be placed in $(k \bmod N)$ cache block

$N \rightarrow$ no of blocks in the cache

Eg: $M=32$, $N=4$;

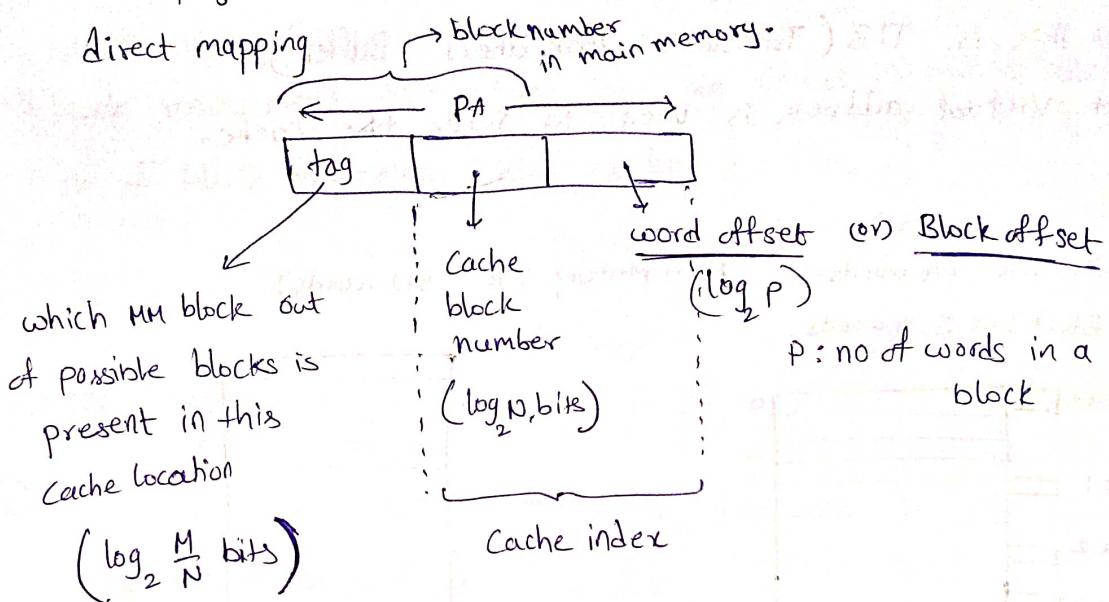
Now consider 7th MM block

Now $7 \bmod 4 = 3$

∴ 7th block of MM has to be placed in 3rd block of cache.

→ Here, physical address is partitioned into 8 fields with

direct mapping



M : no of blocks in MM

N : no of blocks in cache

→ The blocks $(0, 4, 8, \dots, 28)$ of MM are competing for blocks of the cache.

Tag	000	001	010	011	100	101	110	111	Aug bits	Cache
0	4	8	12	16	20	-24	28			
1	5	9	13	17	21	25	29			
2	6	10	14	18	22	26	30			
3	7	11	15	19	23	27	31			

* 8 blocks of MM are competing for each block of cache

$$\therefore \text{Size of tag bits} = \log_2 8 = 3 \text{ bits}$$

Ex: If the content in tags is B_1, B_2, B_3, B_4 find the blocks of MM that are present in the cache.

Sol:

$$B_1 - 000 \Rightarrow 0^{\text{th}} \text{ block}$$

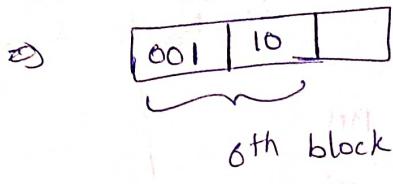
$$B_2 - 100 \quad \text{(crossed out)}$$

The most significant 5 bits will be

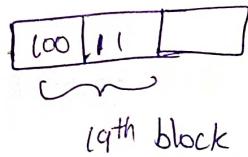
$$\begin{array}{r} 100 \quad 01 \\ \downarrow \quad \downarrow \\ \text{tag} \quad \therefore B_1 \Rightarrow \text{remainder} = 1 \end{array}$$

$$\therefore (10001)_2 = 17$$

$B_2 - 001$



$B_3 - 100$



- (Q2) ~~In a direct mapped cache, the no of tag bits for a cache block is 5 - what does it denote?~~

Sol:

→ Tag bits denote no of MM block competing for a cache block.

→ since we have 5 bits in a tag
there must be 2^5 MM blocks competing for a cache block.

It means for every block of cache we have 32 MM blocks.

i.e., If M is no of MM blocks and N is no of cache memory

blocks then

$$M/N = 32$$

Note:

- The MM blocks having same $(k \bmod N)$ value can not be placed together in cache inspite of cache blocks being free (this is known as conflict problem)

→ Direct mapping requires least number of TAG bits and require only one tag comparator.
 (tag comparator is used to check whether the addressed word is present in the cache or not)

$$\text{no of tag bits} = \log_2 \frac{\text{MM capacity}}{\text{Cache capacity}} = \log_2 \frac{M}{N}$$

$M \rightarrow$ no of blocks in MM

$N \rightarrow$ no of blocks in cache

(Q23) Main memory size = 8 MB

Cache memory size = 16 kB

Block size = 128 bytes

Direct mapping is used. Find no of tag bits.

Sol:

$$\log_2 \frac{8 \text{ MB}}{16 \text{ kB}} = \log_2 \frac{2^{23}}{2^{14}} = \log_2 2^9 = 9 \text{ tag bits.}$$

(Q24) MM size: 4 GB

Block size: 32 bytes

word size: 16 bits = 2 bytes

no of lines in cache = 512 lines

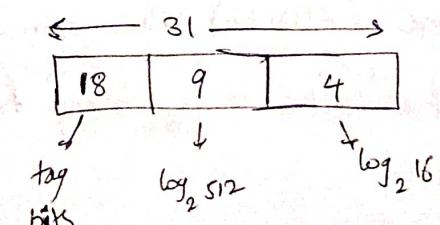
Assuming direct mapped cache

Find size of physical address & size of Cache memory.

Sol:

$$\text{size of cache memory} = 512 * 32 = 2^{14} = 16 \text{ kB}$$

PA:



block size = 16 words

MM size = 2^{32} bytes

= 2^{31} words

(Q5) In a system design, direct mapped cache is having 2^n words and main memory is having 'm' bit physical address. The block contain 2^m words and each word contain 2^p bytes.

The no of tag bits are

Sol:

$$\text{no of tag bit} = \log_2 \frac{\text{MM Capacity}}{\text{Cache capacity}}$$

$$= \log_2 \frac{2^M}{2^n}$$

$\therefore m-n$ tag bits

(Q6) byte addressable MM - 2^{20} bytes

Direct mapped cache - 2^{12} blocks

block size - 16 bytes

find tag, cache line address for the hexa decimal address

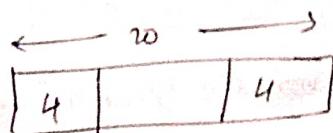
"E201F"

Sol:

$$\text{no of tag bits} = \log \frac{2^{20}}{2^{12+4}} = 4$$

PA : 20 bits

∅ block offset : 4 bits



\therefore no of bits for cache line = $20 - 4 - 4 = 12$

\therefore In hexa decimal

PA:

1	3	1
---	---	---

$\therefore E, 201$

Cache Controller (or) Cache directory controller (or) tag Controller

- The physical address is first sent to cache controller. It will check whether the addressed word is present in cache or not.
- If present, the word is accessed from cache.
This is called cache hit.
- If not present, the block is moved from MM to cache.
This is called cache miss.
- Cache controller maintains one word for each cache block. Each of its word contain tag information of that cache block and status bits of the block.
- Status bits indicate the status of the block in cache
 - Eg: valid bit
 - use bit/s (recently used or not)
 - Modified
- No of words in cache controller = no of blocks in the cache

Status bits:

valid bit:

whether that block in cache has scope or not
(i.e., live / dead)

Used bits:

It denotes that ~~not~~ the block is most recently used or not. (used in LRU replacement technique).

Modified bit:

It denotes whether the block is ~~most~~ modified or not. (Dirty or clean). This information is used with write Back updation.

A direct mapped 32 KB cache maintains 32 byte blocks. Physical address is 32-bits. Cache controller maintains one valid bit, 2 use bits and 1 modified bit along with the tag information for each cache block. What is the size of cache controller (in bits)?

So:

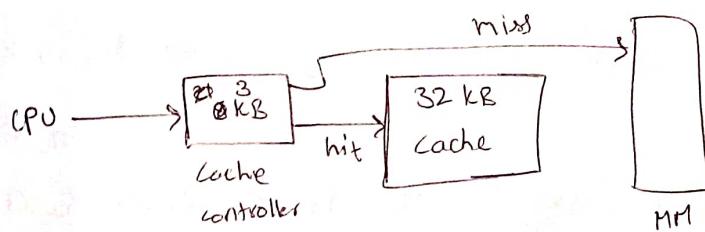
$$\text{no of tag bits} = \log_2 \frac{2^{32}}{2^5} = 17$$

$$\text{no of blocks in cache} = \frac{2^{15}}{2^5} = 2^{10}$$

~~in each~~ size of a word in cache controller = $17 + 4 = 21$

$$\therefore \text{size of cache controller} = 2^{10} \times 21 \text{ bits}$$

$$= 21 \text{ k bits} \approx 5 \text{ KB}$$

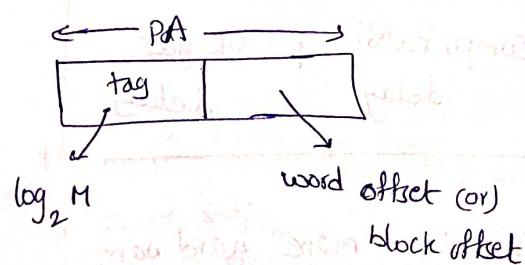


2. Associate Mapping:

→ Any block of the MM can be placed in any block of the cache.

→ Thus there is no conflict miss problem with associative mapping.

→ Physical address will be made two partitions.

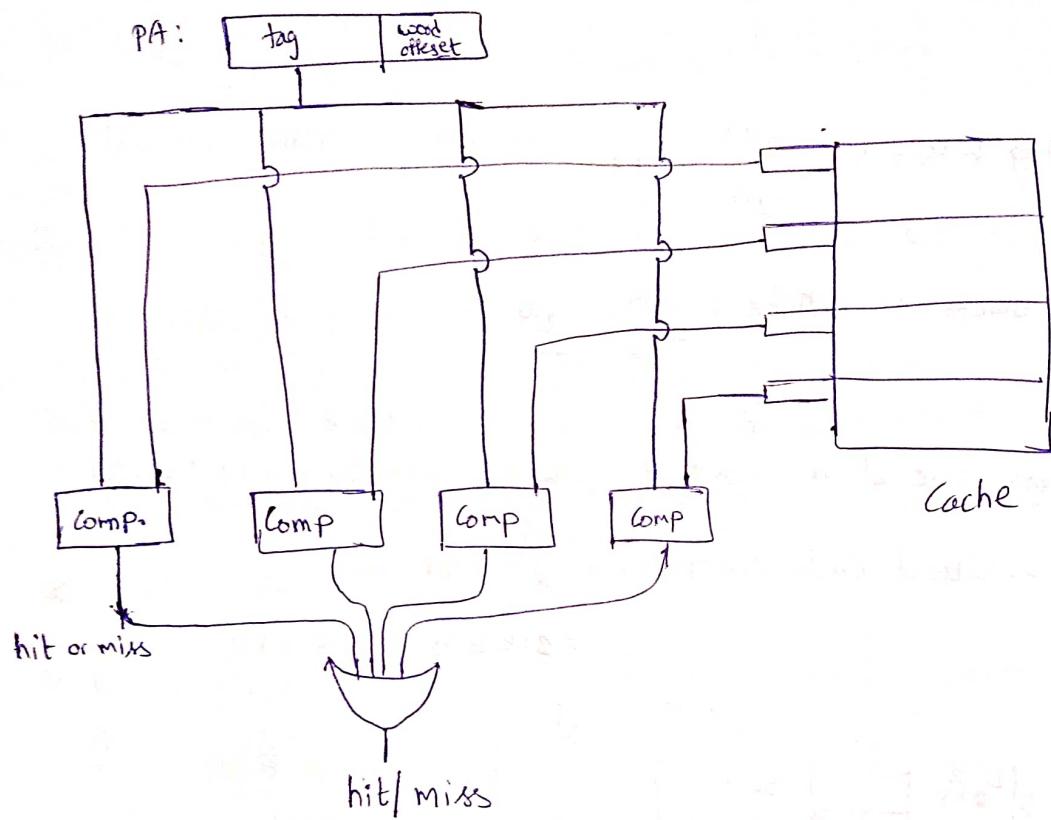


M: no of blocks in
main memory.

→ Thus, associative mapping has more tag bits and requires N-tage comparators.

(i.e., Each cache block has to be checked to know if a hit or a miss has occurred)

→ Thus, size of cache controller is also more (\because more tag bits)



Hit latency:

→ the amount of time required, after providing address, to known whether a cache hit or miss has occurred is called hit latency.

∴ In ~~set~~ associative mapping,

$$\text{hit latency} = \frac{\text{comparator delay}}{\text{delay}} + \frac{\text{OR gate delay}}{\text{delay}}$$

→ Associative mapping requires more hardware.

Associative Mapping

here cache blocks are divided into sets (logical)

Eg: 2-way set associative \rightarrow 2 cache blocks / set

k-way set associative \rightarrow k cache blocks / set

thus we will have $M/N/k$ sets

Now,

$k=1 \Rightarrow$ direct mapping (1-way set associative)

$k=N \Rightarrow$ Associate mapping (N-way " ")

set associate mapping follow direct mapping as well as associate mapping.

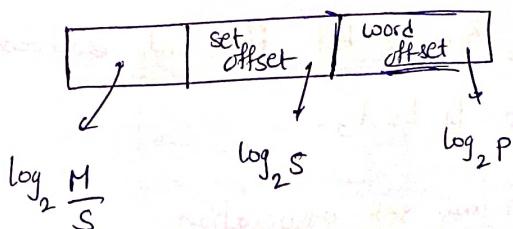
let no of sets, $S = \frac{N}{K}$

Now block 'x' of MM is placed in

$(x \bmod S)$ set of cache

within the set the block may be place anywhere.

Physical address:



M: no of blocks in MM

S: no of sets

p: word offset or block offset.

tag bits in k-way set associative $= \log_2 M/S$

$=$ tag bits in direct mapping + $\log_2 k$

Proof:

$$\text{no of tag bits} = \log_2 M/S \quad S = N/k$$

$$= \log_2 M / \left(\frac{N}{k}\right)$$

$$= \log_2 \frac{Mk}{N}$$

$$= \log_2 \frac{M}{N} + \log_2 k$$

$$= \text{no of tag bits in direct} + \log_2 k$$

mapping

→ To know hit/miss for a main memory reference only one cache set has to be searched. However all ~~blocks~~ blocks of that set have to be searched.

$$\therefore \text{no of tag comparators} = k \quad (1 \leq k \leq N)$$

(k-way set association)

→ Although conflict problem is present in set associative it is less. Also the hardware is ~~of~~ of medium cost.

30/08/20

Q27 Given the following designs for 32-bit physical address and 32kB cache with respective tag bits t_1, t_2, t_3 .

S1: 32 byte blocks with 2 way set association.

S2: 16 byte blocks with 4 way set association

S3: 8 byte blocks with 8 way set association

Find value of expression $2^*t_2 - t_3 - t_1$

Sol:

$$t_1 = \log_2 \frac{2^{32}}{2^{15}} + \log_2 2 = 17 + 1 = 18$$

$$t_2 = \log_2 \frac{2^{32}}{2^{15}} + \log_2 4$$

$$= 17 + 2 = 19$$

$$t_3 = 17 + 3 = 20$$

$$\therefore 2(19) - 18 = 28$$

$$= 0$$

1MB cache with 64 byte blocks is organized with 16-way set association. The virtual memory is having 8KB pages. Minimum no of page colors required such that synonym pages mapped to same cache sets will have same color.

- a) 1 b) 2 c) 4 d) 8

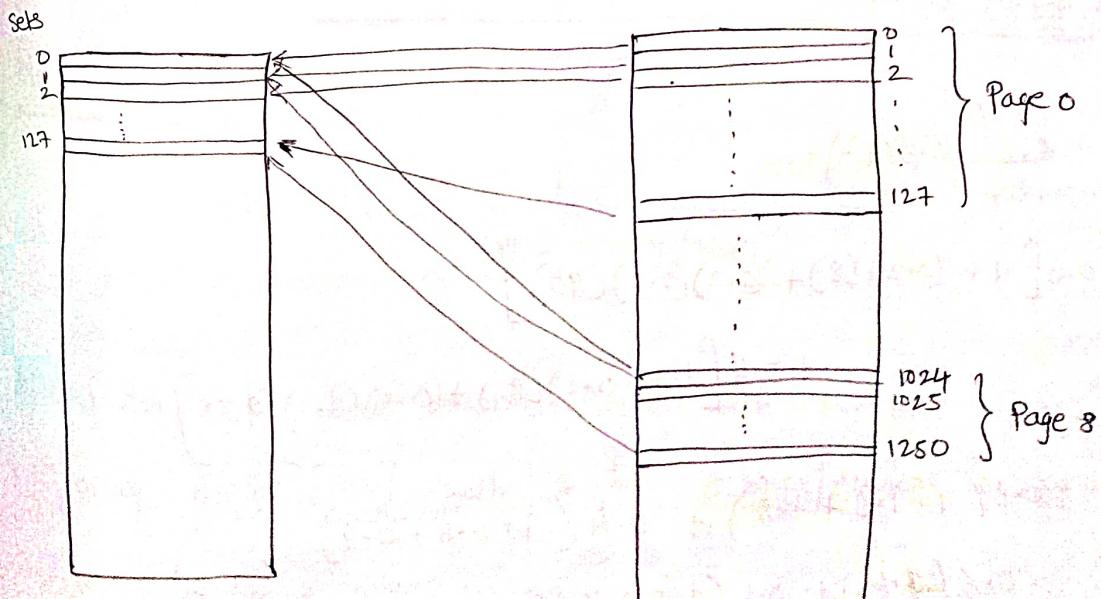
Sol:

$$\text{no of cache blocks} = \frac{2^{20}}{2^6} = 2^{14}$$

$$\text{no of sets} = \frac{2^{14}}{16} = 2^{10}$$

$$\text{size of each set} = \frac{2^{20}}{2^{10}} = 2^{10} = 1\text{KB}$$

$$\text{no of blocks per page} = \frac{2^{13}}{2^6} = 2^7 \text{ blocks}$$



$\therefore 8$ colors

Types of misses in cache:

Conflict miss:

Miss that occurs even after having free blocks.

Capacity miss: Cache association ↑ \Rightarrow conflict miss ↓

Miss that occurs when all blocks are full.

Cache size ↑ \Rightarrow capacity miss ↓

Compulsory miss:

Miss that occurs for the first time when the block is referred.

Cache Block size ↑ \Rightarrow compulsory miss ↓

- Q29) In a memory hierarchy the hit ratios and access times are given below

Memory	Access time	hit
L1 - Instruction cache	1 ns	0.8
L1 - Data cache	1 ns	0.9
L2 - Combined cache	8 ns	0.9
MM	90 ns	1.0

If there are 60% instruction fetch operations and 40% data fetch operations what is the avg access time

Sol:

$$0.6 \left[0.8(1) + \dots \right]$$

$$0.6 \left[1 + (0.2)(8) + (0.2)(0.1)(90) \right]$$

$$+ 0.4 \left[1 + (0.1)(8) + (0.1)(0.1)(90) \right]$$

$$= 0.6 (1 + 1.6 + 1.8) + 0.4 (1 + 0.8 + 0.9)$$

$$= 0.6 (4.4) + 0.4 (2.7) = 3.72$$

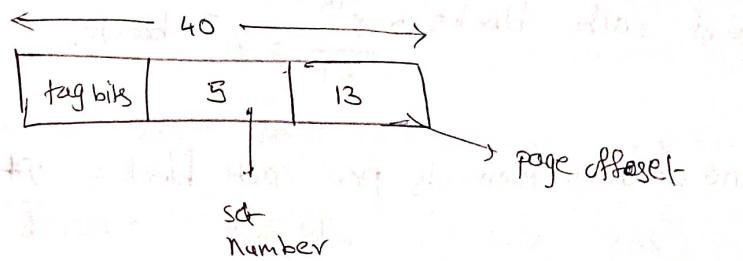
Q) 40 bit virtual address system of 8KB pages maintains 128 word entry TLB. The TLB is having 32 sets and uses 4-way set associate mapping. Number of TAG bits of TLB is _____

Sol:

$$VAS = 2^{40} \text{ bytes } 2^{40} \text{ words}$$

$$PS = 8\text{KB} = 2^{13} \text{ bytes}$$

VA:



$$\text{Tag bits} = 40 - 5 - 13 = 22$$

Array:

* Initially cache contains no data relate to array and array is present in main memory.

* If array is stored in row major order and accessed in row major order then the performance is ~~not~~ optimal. Try column major order.

* If array contains fairly large no of blocks ~~than~~ compared to no of blocks in the cache, then

$$\text{miss ratio} = \frac{1}{(\text{no of elements / block})} \quad (\text{irrespective of the mapping technique})$$

A CPU has 32 kB direct mapped cache with 128 byte blocks. A

2-D array $A[512][512]$ with ~~size~~ 8 bytes/element is stored in MM.

Consider following program

for($i=0$; $i < 512$; $i++$)

 for($j=0$; $j < 512$; $j++$)

$x += A[i][j];$

Assume i, j, k are present in processor registers.

i) what is the miss ratio?

Sol:

$$\text{no of cache blocks} = \frac{2^{15}}{2^7} = 2^8 \text{ blocks}$$

$$\text{no of array elements per cache block} = \frac{2^{14}}{8} = 2^{14-3} = 16 \text{ elements}$$

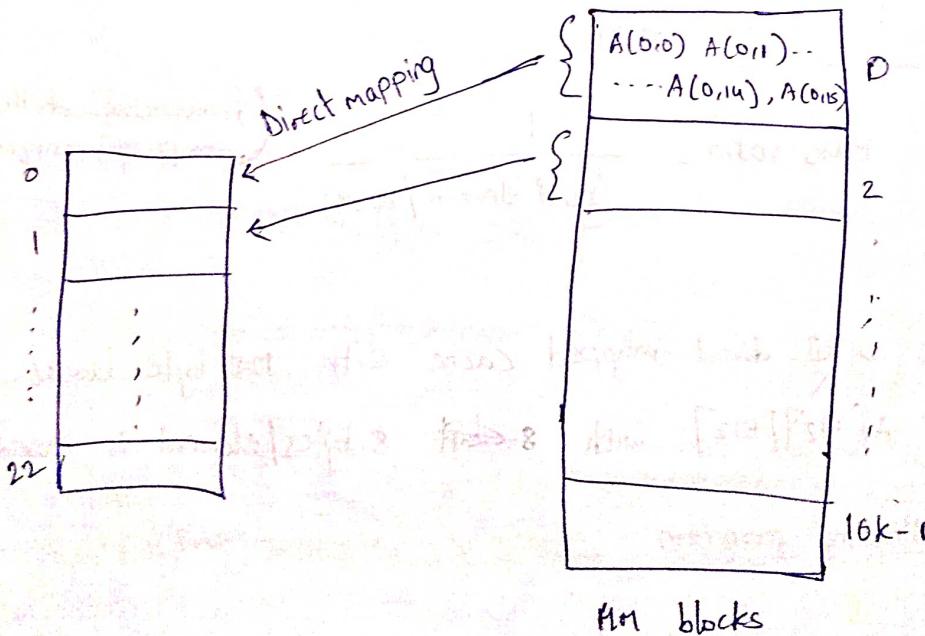
$$\text{size of array} = 512 * 512 * \frac{8}{2} = 2^{21} \text{ bytes}$$

$$\text{no of array blocks in array} = \frac{2^{21}}{2^7} = 2^{14} = 16k \text{ blocks}$$

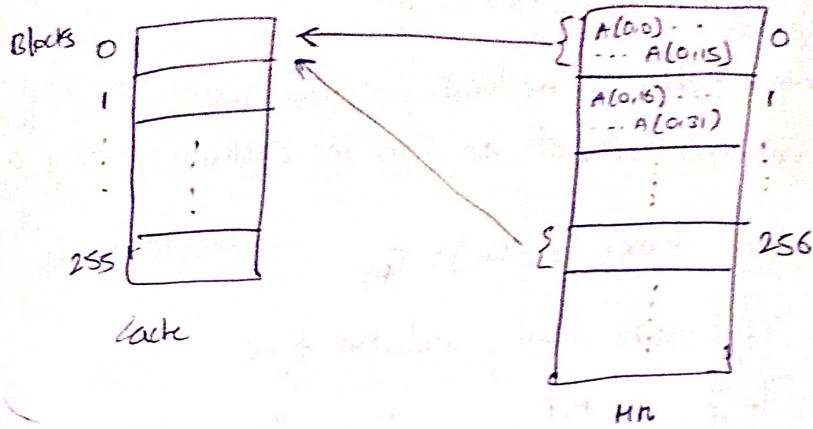
Since we have 16 elements per block and the elements are in row major order, the 1st element of the block generates a miss and the references of subsequent 15 elements will be a hit.

$$\therefore \text{miss ratio} = 1/16$$

$$\Rightarrow \text{hit ratio} = 15/16$$



iii) which array element replaces $A(0,0)$ in cache memory?



the first element is 256th block of array replaces $A(0,0)$

no of element in 256 blocks ($0, 1, 2, \dots, 255$)

$$\begin{aligned}
 & 256 \times 16 \\
 & \hookrightarrow \text{Ele/block} \\
 & = 2^{12} \text{ elem} \\
 & \therefore 2^{12} \text{ th element} \\
 & \text{if row number} = \frac{2^{12}}{512} = 8 \\
 & \quad \left. \begin{array}{l} \text{each row has 512 ele} \\ \text{one row} \rightarrow \frac{512 * 8}{27} = 2^5 \text{ blocks} \\ 2^5 \text{ row ends in} \rightarrow 2^6 \text{ blocks} \\ 4^{\text{th}} \text{ row ends in} \rightarrow 2^7 \text{ blocks} \\ 8^{\text{th}} \text{ row ends in} \rightarrow 2^8 \text{ blocks} \end{array} \right. \\
 & \therefore A(8,0)
 \end{aligned}$$

\therefore In general for this problem

$A(i,j)$ is replaced by $A(i+8,j)$

Cache Updations

* write operations ~~result~~ in cache coherence problem.

* so proper updation techniques are needed.

* we consider this problem in

i) uni.processor system

write through

write back

ii) Multi processor system

Cache snooping protocols

Cache directory based protocols.

Write through updation (uniprocessor system)

* Cache & MM are updated simultaneously

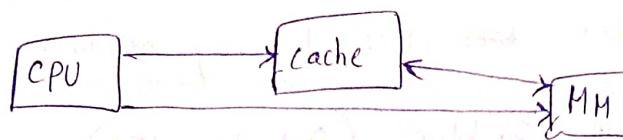
* Here cache coherence problem is resolved

i.e., all the time contents of an item in cache & memory are same.

$$\therefore \text{time for updation} = \max(T_c, T_m) = T_{wm}$$

T_c - Cache memory updation time

T_m - MM " "



Advantage:

* For less ~~no of~~ ~~with~~ write-~~trapping~~, it gives better performance

Disadvantage:

* Worst performance for more no of writes.

Write-Back Updation:

→ Here MM is updated only when the concerned block is taken out of the cache.

→ As long ~~as~~ the block is present in the cache, the update is not visible to MM.

→ Hence cache coherence ~~is~~ problem is still present but with less magnitude

→ For each cache block a modified bit is used to ~~know~~ whether the block is ~~not~~ updated or not.

$$\text{time for updation} = h * T_c + m \left[\underbrace{(T_{Bnew} + T_c)}_{\substack{\text{clean slot} \\ (\text{modified bit}=0)}} p + \underbrace{(T_{Bold} + T_{Bnew} + T_c)}_{\substack{\text{dirty slot} \\ (\text{modified bit}=1)}} (1-p) \right]$$

h : hit m : miss ratio

T_{Bnew} : Moving block from MM to Cache

T_{Bold} : Cache to MM

T_c : write operation

p : probability of clean slot
(modified bit=0)
(or)
empty block

Dirty slot
(modified bit=1)