

Computer Architecture

Refers to the design, structure & organization of computer system

- It defines how processor, memory & I/O devices communicate & work together.
- it includes instruction sets, M/M hierarchy, control mechanism

Technology trends - focusses on how components like processor, memory, storage and network are improving with time

Ex

Faster processor
improvement in wireless communication like
5G

- AI, cloud computing
smaller components

Fast CPU with more cores

less power consumption

Moore's Law - proposed by Gordon Moore in 1965

- It states that the number of transistors on a microchip doubles approx every 2 years.
- This leads to exponential growth in computing power while size and price reduces.

Application - More powerful processors, better graphics, Faster networks, AI, cloud computing.

Parallel computers - are designed to

process multiple tasks simultaneously to improve speed & efficiency

SISD - Single instruction single data

- Traditional computers
- Where 1 instruction operates on single data element at a time

SIMD → single instruction multiple data

- Used in GPU
- Single instruction operates on multiple data elements

MISD - Multiple instruction single data

Used in Fault tolerant systems

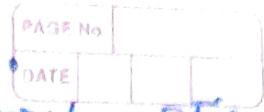
- Multiple instruction operate on same data

MIMD - "M" multiple data

- Used in multiple core processor
- Each processor works independently on its own instruction and data

core; no processor is an independent processing unit within CPU

e.g. dual core \rightarrow can handle 2 tasks efficiently
dual core \Rightarrow " 4 task



Think of cores as workers in factory
single core can do only 1 job at once

Myopic view of C.A

- it is a narrow perspective that focuses on improving the processor speed without considering factors like

System Cost

Energy consumption

Reliability, maintainability

A balanced view is required to ensure practical system designs

used to

Andahl's law - Shows the limit of

speedup by parallelizing only a part of program

Speedup - ratio of time taken to run program on a single processor to the time on multiple processor

$$S = \frac{T_1}{T_n} \leftarrow \begin{array}{l} \text{single processor (Sequential)} \\ \text{N processor (Parallel)} \end{array}$$

measures how much faster a task becomes when parallel techniques are applied.

$S=1$ No improvement

$S=N$ Perfect speedup

$S < N \rightarrow$ Prohibitive

$$\text{Efficiency} = E = \frac{\text{Speedup}}{\text{no of processor}}$$

Measuring how well the processes are divided in a OS system

Formula for speed up

ordinary law

$$\text{Speed Up} = \frac{1}{(1-P) + \frac{P}{N}} \quad \begin{matrix} \rightarrow \text{ // Fraction} \\ \downarrow \qquad \qquad \qquad \downarrow \\ \text{Sequential} \qquad \qquad \qquad \text{Fraction} \end{matrix} \quad \begin{matrix} \leftarrow \text{ something } \\ \text{Faster} \end{matrix}$$

$P \rightarrow$ fraction of program that can be parallelized

$N \rightarrow$ number of processor

e.g. if 80% can be // then $P = 0.8$

e.g. if 40% of is floating pt operation
and factor is 4 find speed up

Aptitude

Formula

Sequential Parallel

$$\therefore \frac{1}{1} = 1$$

60% 40%

$$\frac{1}{1 - 0.4 + \frac{0.4}{4}} = 0.7$$

now \sqrt{factor}

$$0.6 + 0.1 \cdot 4$$

$$\text{Earlier it was } \frac{1}{0.7}$$

if we have 1 processor

$$T_1 = S + P$$

Serial (S) Parallel (P)
↓ time to execute serial part ↓ time to execute "n=1"

$$T_n = \frac{S + P}{n}$$
$$= S + \frac{P}{n}$$

assuming $T_1 = 1$

$$1 = S + P$$

$$S = 1 - P$$

Replacing in T_n

$$T_n = (1 - P) + \frac{P}{n}$$

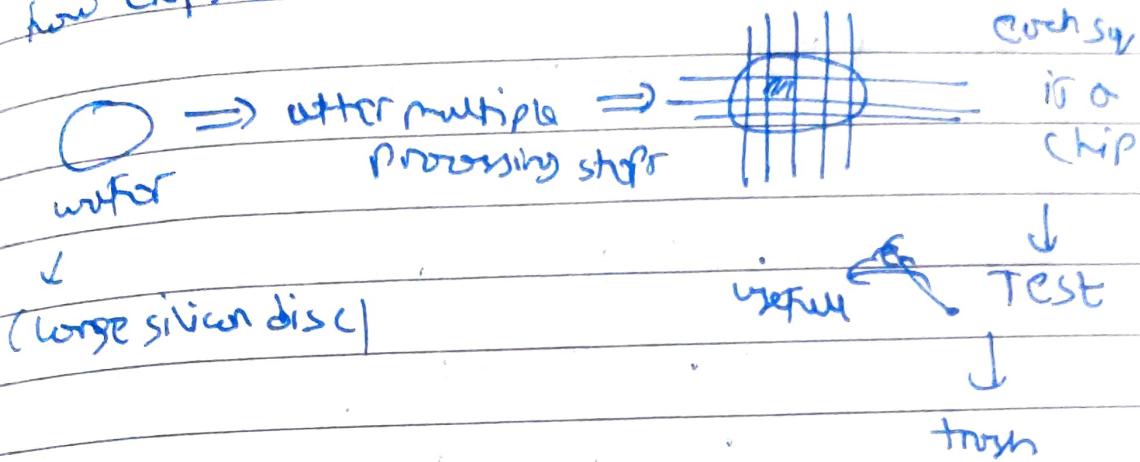
$$\text{Speed up} = \frac{T_1}{T_n} = \frac{1}{(1 - P) + \frac{P}{n}}$$

Die = chip

PAGE No.	
DATE	

Fabrication cost

how chips are made



Fabrication yield - refer to the % of functional chips produced upon the total no of chips

$$\text{yield} = \frac{\text{Number of good chips}}{\text{Total no of chips}} \times 100$$

e.g. Cost of wafer = ₹ 2000

Number of dies per wafer = 1000

$$\text{yield} = 80\%$$

i.e. only 800 good die

$$\text{cost per chip} = \frac{2000}{800} = 2.5$$

$$\text{Cost per chip} = \frac{\text{Cost of wafer}}{\text{Total no of good chips}}$$

$$\text{yield} = e^{-\text{defect density}} \rightarrow \text{defect density}$$

Eg defect density = $0.08/\text{cm}^2$

wafer area = 150cm^2

No. of wafers per run = 500

Cost per wafer = ₹ 2000

Dies per wafer = 1000

$$\text{yield} = e^{-\text{defect density} \times \text{area of die}} \rightarrow$$

$$e^{-0.08 \times 0.15}$$

area of
wafer
 $\frac{150}{1000}$

$$\approx 0.988$$

Yield = 98.8 %

Cost of production

$$\text{Cost} = \text{No. of wafers} \times \text{cost per wafer}$$

$$200 \times 500$$

Cost per die =

No. of functional die

$$= \frac{98.8}{100} \times 1000 = 988$$

$$\text{Cost per diis} = \frac{\text{unit cost}}{\text{instructions}} = \frac{2000}{488} = 2.024$$

CPI → clockcycles per instruction

It is a performance metric that tells us the average number of clockcycles required to execute a single instruction in a program

Lower CPI → Faster execution

Higher CPI → Slower execution

$$\text{CPI} = \frac{\text{Total clockcycles}}{\text{instruction count}}$$

total cycles taken
 to execute
 prog

$$\text{CPU Execution Time} =$$

$$\boxed{\text{C.P.I} \times \text{instruction count} \times \text{clockcycle time}}$$

OR

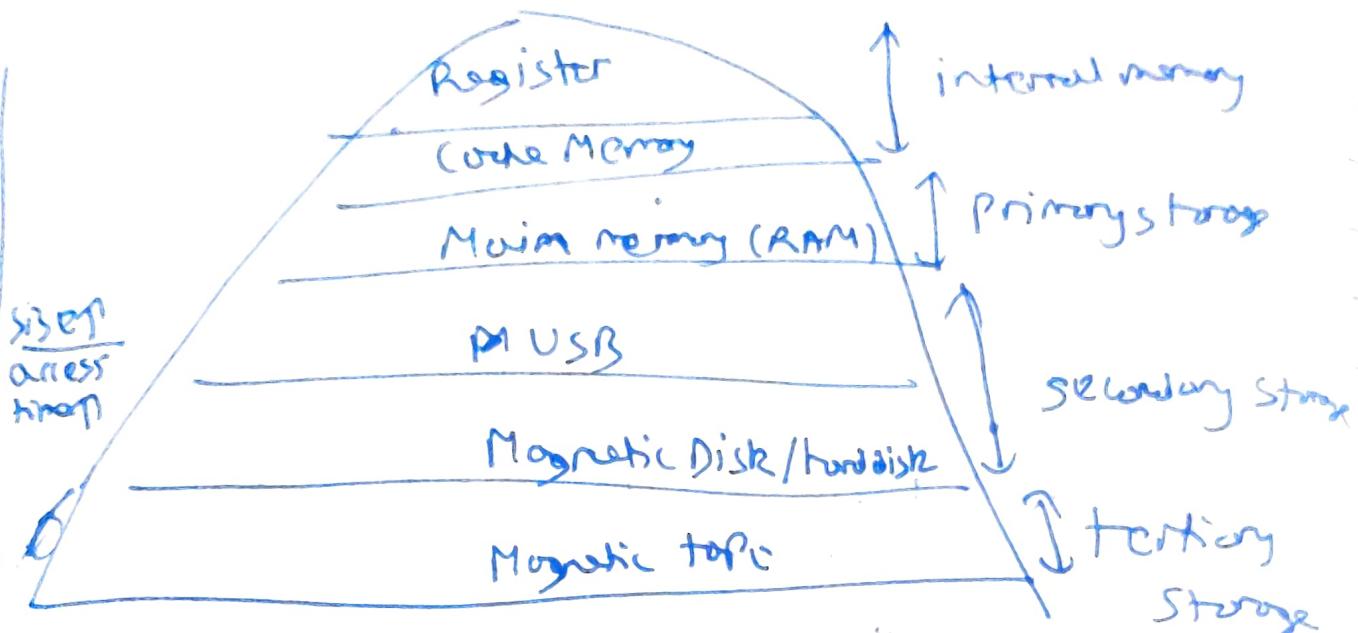
$$\text{C.P.I} \times \frac{\text{instruction count}}{\text{clockrate}}$$

in Hz ←

clockrate (freq)

$$\text{Clock cycle time} = \frac{1}{\text{clock rate}}$$

MM hierarchy



Locality of reference

temporal locality : Recently accessed data will be used again

spatial locality : Data near recently accessed data will be accessed soon

Cohesive - In multiprocessor system multiple Cores have copies of same data

Core coherence ensures that all processor see the most recent copy

Cache - Small, fast memory located between CPU & main memory
It stores copies of frequently used data so that CPU can access them quickly instead of going to main memory

W+P
↳ built using SRAM (Static RAM)
↳ works on principle of locality of reference

levels of cache L1 - closest to CPU, Fastest
L2, slightly slower
L3 - longer, slower than L1/L2
but faster than RAM

Reduce CPU access time

Mapping techniques

- Direct mapping Fully organizing memory

Set associative mapping

Unit for memory is word

Cache is divided in lines

A block from main memory is placed in one line in cache

basic wrpt

From 2 cache in Main memory and cache in both system and CPU have access time. That is mapping techniques

when CPU want data .

- 1) it checks cache \rightarrow if the block is there
 - if yes \rightarrow cache hit
 - If no \rightarrow cache miss (Fetch from RAM)
- 2) if block not in cache , CPU decides where in cache to place it.
This is called cache mapping.

Cachemapping decides how RAM blocks are placed in cache lines , so CPU knows where to look when it needs data.

Direct Mapping

Each block of RAM can go to only fixed line in cache

- Simple , fast

~~Disadvantage~~ 2 frequently used blocks may collide

Fully Associative mapping

- Any block or row can go to any line cache line

1 word = 1 Byte

Line size = Work size

PAGE No.	
DATE	

CPU searches all lines in /I to find block

Advantage - No conflict

Dis - Slower searching (expensive hardware)

Key Set Association - Cache is divided into sets and each set has a few lines

A block or ram is mapped to one set but inside that set it can go to any line

Adv - Better than direct and convolution

dis - slightly more complex

	Direct	Fully associative	Set associative
Flexibility	Low	High	medium
Hardware cost	Low	High	med
Conflict	High	v. low	low
Speed (Latency)	Fast	Slow	balanced
Hit rate	Lowest (due to high conflict)	Highest	medium
Suitability	suited for simple systems (microcontroller)	suitable when hit rate is critical (e.g. GPU)	Used in general purpose CPU like laptop, server

1 word = 1 byte

size of block = size of cache line

K way set associativity

↳ means no of lines in a set

e.g. if $k=2, 4, 8$

means that divide lines in

e.g. cache has 16 lines

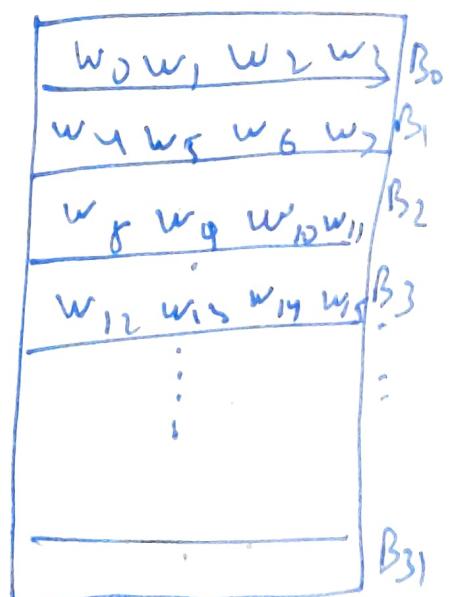
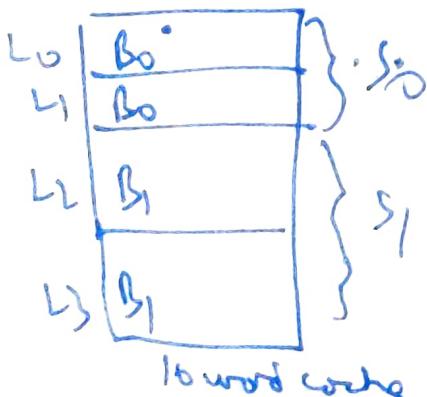
2 way set associativity

so no of sets = 8

each with 2 lines

e.g. 16 word cache 128 word MIM

each block contains 4 words



$$\text{so no of blocks} = \frac{128}{4} = 32$$

Lines \Rightarrow slots

PAGE No.	
DATE	

block size = 4 word

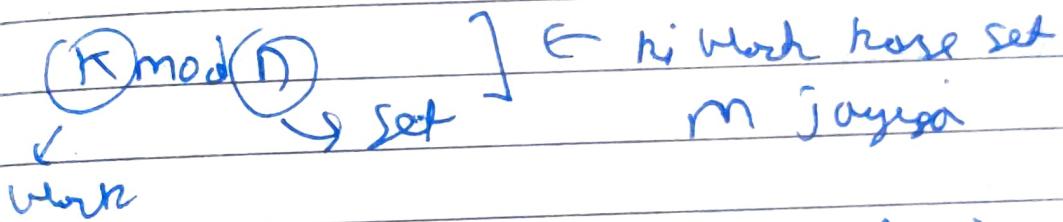
so line size also 4 word

$$\therefore \text{so no of lines} = \frac{16}{4} = 4$$

if 2 way associate ie 2 lines per set

Kmodn

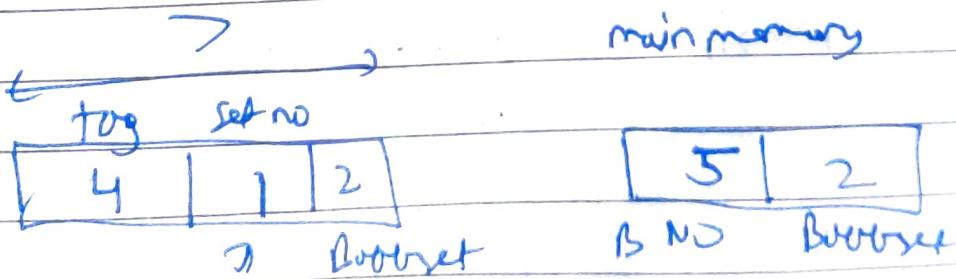
$$\text{no of set} = \frac{\text{Total lines}}{2} = \frac{4}{2} = 2 \text{ sets}$$



eg 0 mod 2 B₀ in set0
 1 mod 2 B₁ in set1
 2 mod 2 B₂ in set0

and set m Work

Hi LRU replacement Associativity concept



C) tag 2 Set 0 bit

0010 0.00

in Work.

Now in set 0 it will look for tag in any line it can get it

Other than offset rest is Works

$$\text{offset} = 2^{\text{bit}} \\ \text{blocksize} = 2^{12} = 4096$$

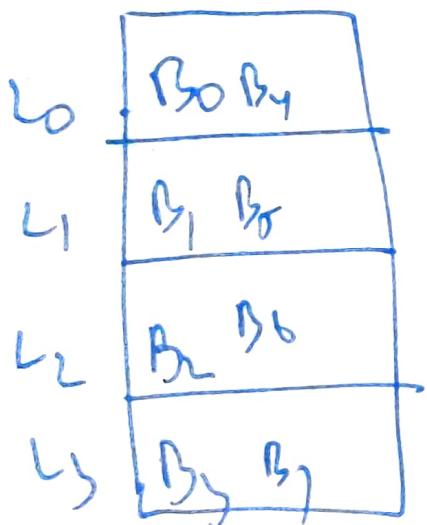
$$P-A = 1011\ 0110\ 1100 \\ \text{Tag index offset}$$

in line 0110 i.e. 6

Check the tag 1011 if fit then block is Tag + index
1011 0110 =

Example works
4 words

Direct mapping

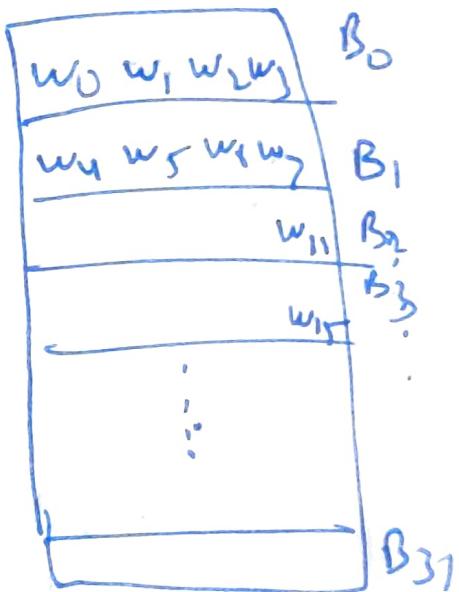


Cache size 16 word

line size = block size \geq 4 word

\Rightarrow 4 words line

$$\frac{\text{cache size}}{\text{line size}} = \frac{16}{4} = 4$$



128 words

To block size \geq

4 word

$$\Rightarrow \text{Total block} = \frac{128}{4} = 32$$

128

\rightarrow

0-127

$$\text{associativity} = 2^7 = 128$$

no of
bits
 $= 7$

Ek time pe ek hi work in a line

PAGE No.		
DATE		

Direct Mapping

Formula : $\text{work} \mod n$

\downarrow
worknumber

→ ~~for no mapping~~

Ki konsa block konsi line mijayega

work 0 → 0 mod 4 → 0th line

Main m/m →

how to get?

Physical address

size of m/m = 128

12 bits

Physical

workno	work offset
(PA) 7 bits	→ each work in 1 line

• 5 bit 2 bit



key

to get 4 we need
2 bit

e.g. if CPU says 10th word

0001010 ie I need 10th word

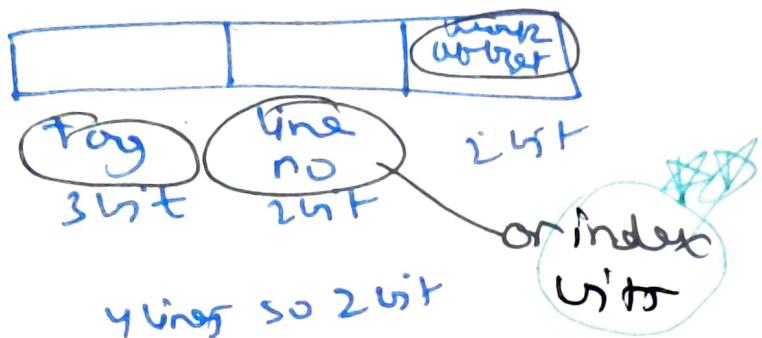
0001010
 \downarrow \downarrow
5 2

it is 2nd block 2nd word

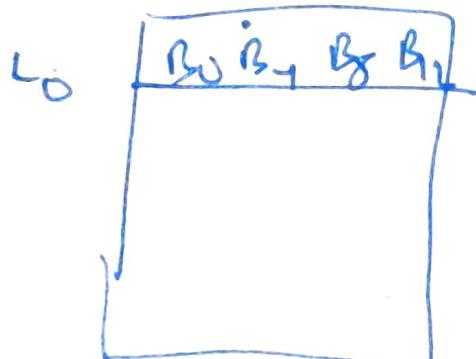
Why P.A assumes size of memory in cache? Ans

now for cache

The size remains same 7 bits



Tag tells that hi line in random work



Block number = Tag + line no

or lines

Tag lines no to offset

001 00 00

~~16 8 4~~ 10
16 8 4 21

ie 10th word is in

16 block lines

see Tag & line no for

		offset	
001	01	00	w_{20}
001	01	01	w_{21}
001	01	10	w_{22}
001	01	11	w_{23}

At a time can have only 1 word out of many to check which word is present see the possible words in line I right now

the shorter way or work with some fixed address as in M/M

P.A \rightarrow actual address of data in memory

PAGE NO.	DATE

In Fully associative

P.A

some ex

5	2
---	---

block offset
no

Now since in cache any block can go to any line so line no is not required

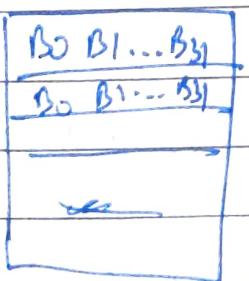
P.A at cache \rightarrow

5	2
---	---

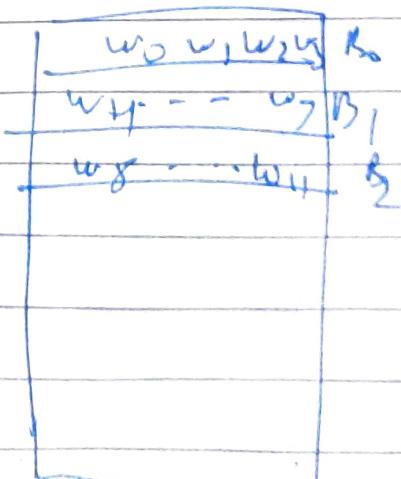
\nearrow tag offset

tells us block no

i.e. in a line any block can be present
as 32 block so 5 bits
means any block



tag address for this



16 word
 \nearrow tag offset
00100 00
01
10
11

is block 415
here
128 word

To convert decimal
in binary

First

Decimal to hex

then write it in binary

Iron law of performance - states that the performance of a program depends on 3 factors :

1) no of instructions

3) clock cycle time

2) Clock cycle per instruction

$$\text{CPU Execution time} = \frac{\text{Instruction count}}{\text{Clock cycle time}} \times \text{C.P.I}$$



To improve the performance / make prog faster

- Reduce instruction count
- Reduce C.P.I \rightarrow better architecture (pipelining)
- Reduce clock cycle time \rightarrow faster hardware

$$1KB = 1024 \text{ bytes}$$

$$1MB = 1024 KB$$

Just For understanding

cpu \rightarrow page frame
(in page frame change) \rightarrow maintains hi result from both

Virtual memory - memory management technique used by modern O.S that gives the illusion of a large, continuous main memory (RAM), even if the actual physical m/m is smaller.

it allows prog to believe it has access to large block of m/m by using secondary storage or on extension of m/m

- implemented using Paging or Segmentation

How virtual m/m Abstract & Extend Physical m/m

- Abstraction means hiding the complexity and presenting a logical view

In V.M system

O.S hides the real physical m/m layout and each process sees a simple clean continuous block of m/m starting from 0 even though in reality it is scattered in RAM

This is called address space abstraction

Physical mem is limited i.e 8 GB or 16 GB
but programs may need more space than this

V.M extends physical memory by using a
portion of secondary memory or backup
M/M

How it works

Only active progs (currently used) are kept in ram
inactive progs are moved to the disk
(not currently used)

when a prog needs them they are brought into
RAM

There is a page fault mechanism using page fault unit
which they are brought in ram

Fundamental concept of V.M

- i) V.M provides address space abstraction
- Each program is given its own virtual address

These virtual address are mapped to physical address in RAM by the Memory management unit using a structure called page table.

If any page is required that's not in MM then page fault and using page fault service mechanism it is brought into MM.

Components of V.M system

- 1) Page table - is a data structure that is maintained by OS that is used to map virtual address to their physical address in the main memory.
 - keeps track of which virtual page is stored in which frame in main MM

- 2) P.T.E - Each entry in the page table is called page table entry

- it contains info about virtual page line
- Frame number
 - if the page is in memory
 - valid / invalid bit
 - dirty bit
 - program modification bit
 - prefetch list

Translation Lookaside Buffer - small fast cache
inside the CPU that stores recent
Virtual to physical address translation

- Accessing the P.T for every translation is slow
- T.L.B reduces the overhead by storing the translation for frequently used pageframe entries

CPU generates V.A

M.M.U checks T.L.B if found, T.L.B hit Fast access
if not T.L.B miss Page table consulted

How Page table Maps

- Virtual addr - Page No, offset
- M.M.U uses Page no to find the corresponding Frame Number in ~~store~~ the Page frame
- Frame number + offset gives P.A

Advantage of TLB

- Improves / reduces memory access time
- Better CPU utilization
- improved throughput
- Reduces the time required for address translation
- Stores Frequently used Page table entries translation