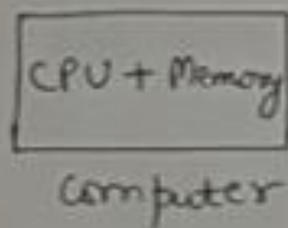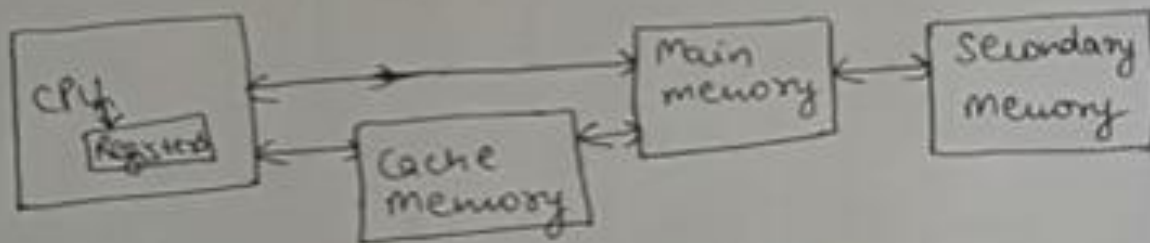# Memory Management

CPU + Memory

Computer

## 3 criteria to choose memory
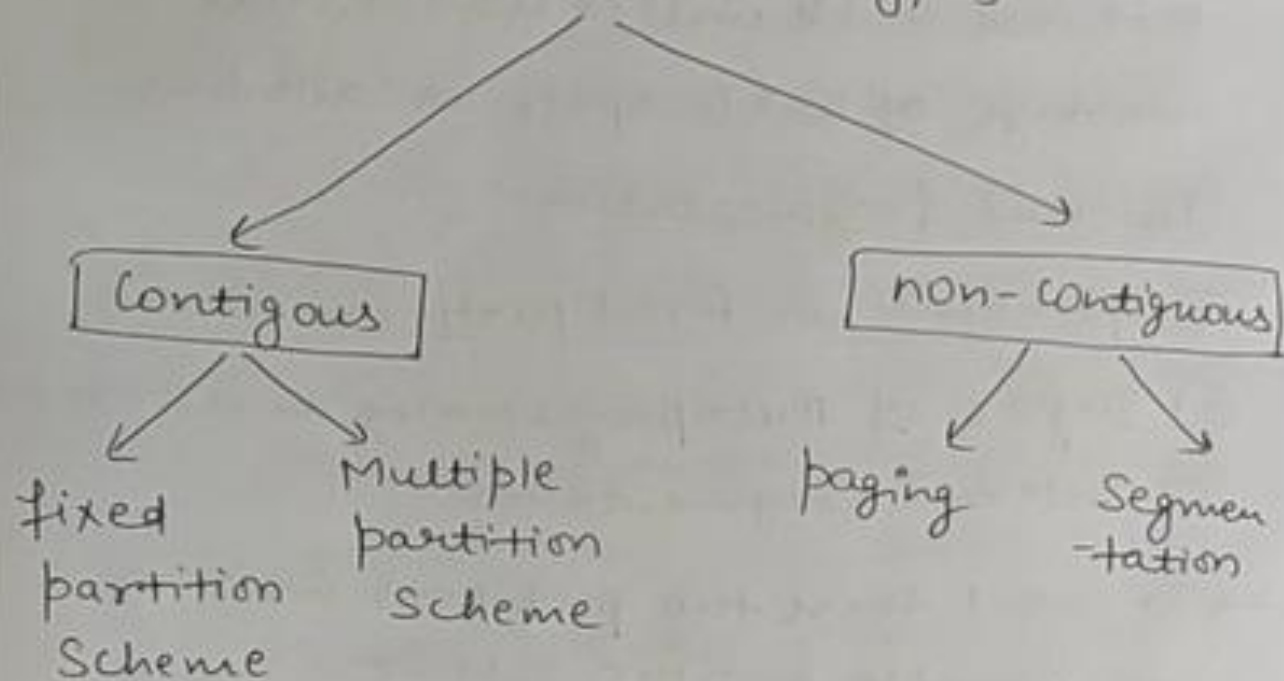
① Size  ② Access Time  ③ cost



**Registers** → It is a temporary storage area built in a CPU. Access time of Register is below 10ns, and registers have lowest capacity ie of few KB of words.

**Cache Memory** → It is a high speed memory. The purpose of cache memory is to store those programs that are repeatedly used or likely to be used in the near future.

# Memory Management Techniques

It is categorized into 2 types

```
                Memory Management Techniques
                  /                    \
           Contiguous              non-Contiguous
            /      \                  /        \
        fixed    Multiple         paging     Segmen
      partition  partition                   -tation
      Scheme     Scheme
```

→ **Fixed Partition Scheme**

- no. of partitions are fixed in memory.
- Size of each partition may be same or may be different.
- for eg:- we made 8 partitions in the memory. So, the no. 8 is fixed.
  In each partition only one process can be placed. In this case, we can place only 8 processes at a time. So, the degree of multiprogramming will be restricted.

If one of the partition is of size 50KB, and we have a process of size 20KB, in that case 30KB will be wasted. This wastage of 30KB space is called as Internal fragmentation.

Two problems in fixed partition scheme
① Degree of Multiprogramming is restricted
② Internal fragmentation.

→ To avoid these two problems we move to Variable partition scheme.
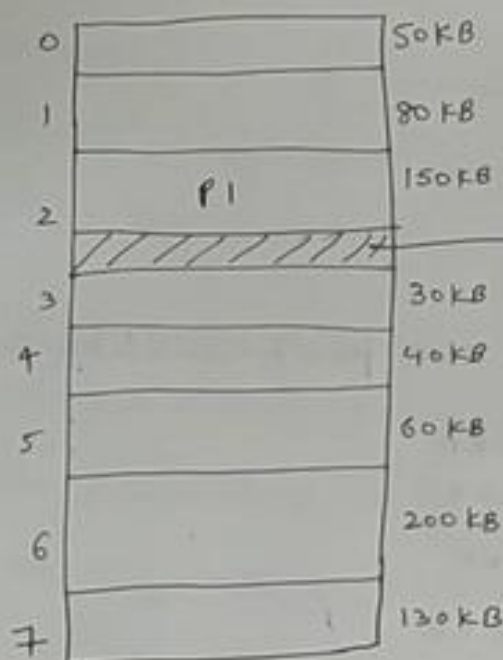
→ Variable partition Scheme
- In this, whenever the request of the process arrives accordingly partition is made in memory.
- for eg:- If a process of size 50KB arrives, a partition of 50KB is made in a memory.
- In this case, there is no restriction on degree of Multiprogramming, as we can make

as many partition as possible till the memory is free.

- Let us assume, after some various process are allocated and only 50 kB space is available in memory. And after sometime, a process of size 80 kB arrives, but we have only 50 kB space left, Therefore we cannot accomodate the process of size 80 kB. Hence, this 50 kB space is wasted and this wasted 50 kB space is known as <u>External fragmentation</u>

→ <u>Example of fixed partition</u>



| | | Process P1 of Size 100 KB arrives |
|---|---|---|
| 0 | 50KB | |
| 1 | 80 KB | |
| 2 | P1   150KB | |
| | → 50 KB is wasted |
| 3 | 30KB | and is known as |
| 4 | 40KB | <u>Internal fragment</u> |
| 5 | 60 KB | |
| 6 | 200 KB | |
| 7 | 130KB | |

⟶ Example of Variable partition

If the request is like this —:

$P1 = 50KB$

$P2 = 20KB$

$P3 = 40KB$

$P4 = 70KB$

$P5 = 100KB$

then we make the partitions like this —:

| | |
|---|---|
| P1 | 50KB |
| P2 | 20KB |
| P3 | 4+KB |
| P4 | 70KB |
| ⤷ external fragment | 60KB |

(only 60KB is left memory is full now)

Therefore, we can not allocate P5 of size 100KB.

→ Partition Allocation Methods

first fit

Best fit

worst fit

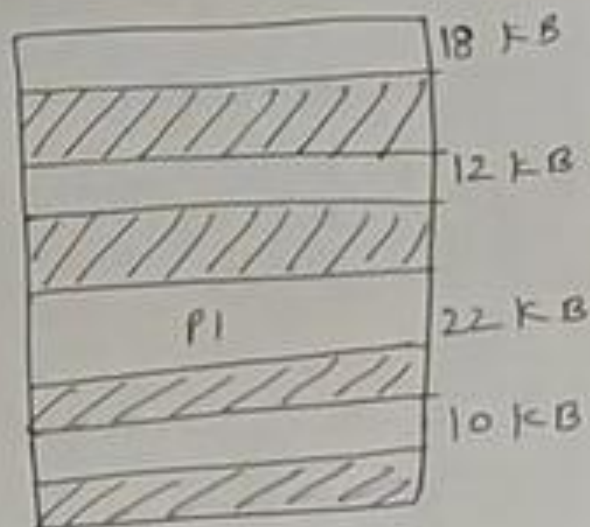| | | |
|---|---|---|
| 𝄞 P1 | | 100KB |
| P4 | | 20KB |
| | | 30KB |
| | | 5KB |
| | | 10KB |

P1 = 20KB

→ Partition Allocation Methods

① First fit → In this, allocate the process in a partition which is first sufficient partition from the top of the memory.

② Best fit → In this, allocate the process in a partition which is the smallest sufficient partition among the free available partition. To find the smallest sufficient partition it requires to search all the free partitions in the memory.

③ Worst fit → In this, allocate the process in a partition which is largest sufficient among the free available partition. To find the largest partition it requires to search all the free partitions in the memory.

→ Example of first fit          P1 = 20 KB
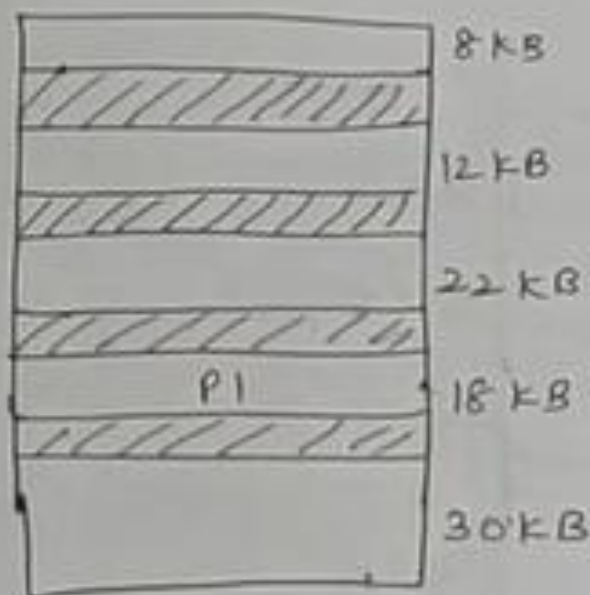


☐ → free partition

▨ → busy partition

Suppose, program P1 arrives of size
20 KB.

In case of first fit, P1 (20KB) will
be allocated in the ~~black~~ frame
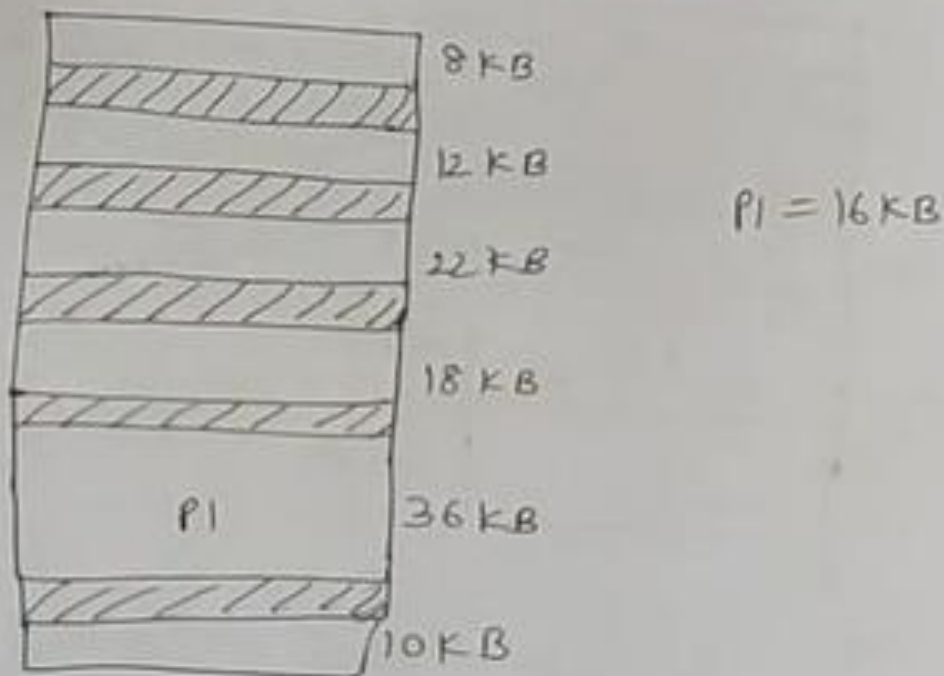of size 22KB.

→ Example of Best fit



P1 = 16 KB

□ → free partition

▨ → busy partition

Suppose, P1 of size 16 KB arrives,
In case of Best fit, P1 will be placed
at 18 KB partition.

→ **Example of worst fit**



8 KB
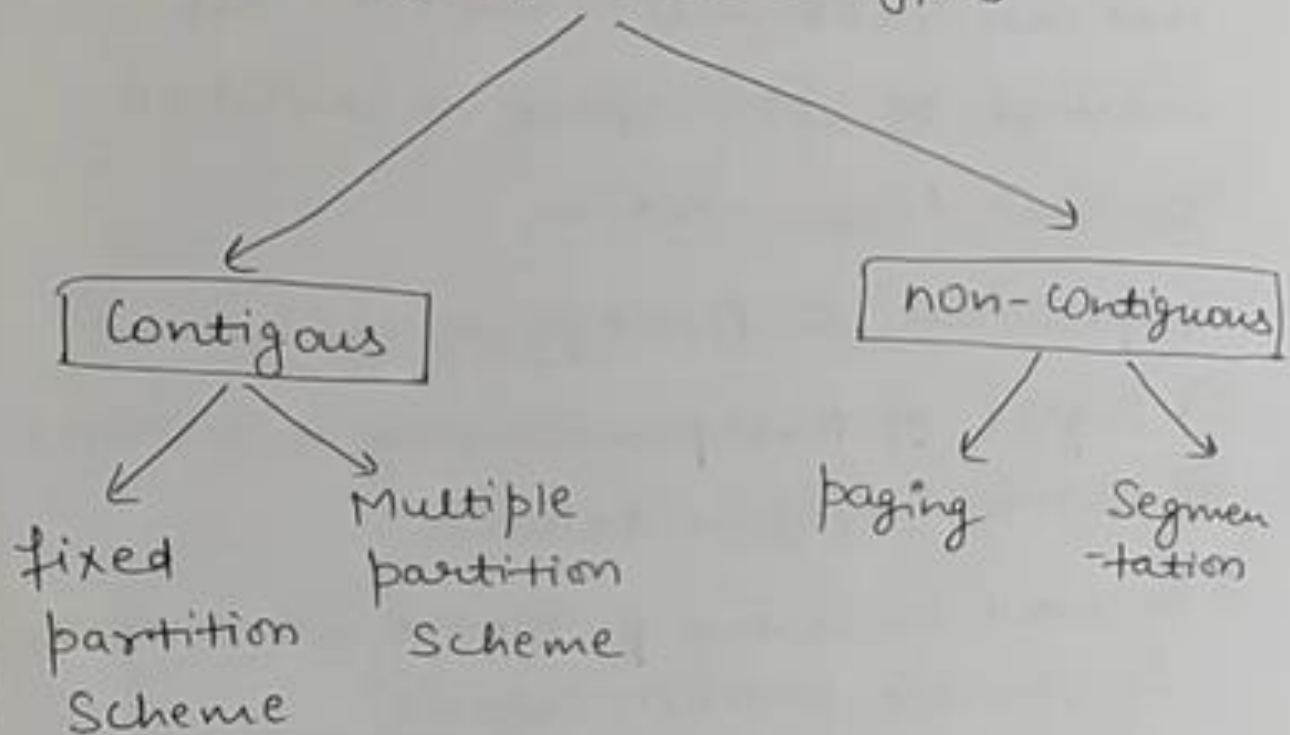
12 KB

22 KB

18 KB

PI   36 KB

10 KB

PI = 16 KB

▭ → free partition

▨ → Busy partition

Suppose, PI of size 16 KB arrives. In case of worst fit PI will be placed at 36 KB partition.

# Memory Management Techniques

It is Categorized into 2 types

```
                    Memory Management
                    /              \
              Contiguous        non-Contiguous
              /      \            /         \
         fixed    Multiple    paging    Segmen
      partition   partition             -tation
       Scheme      Scheme
```
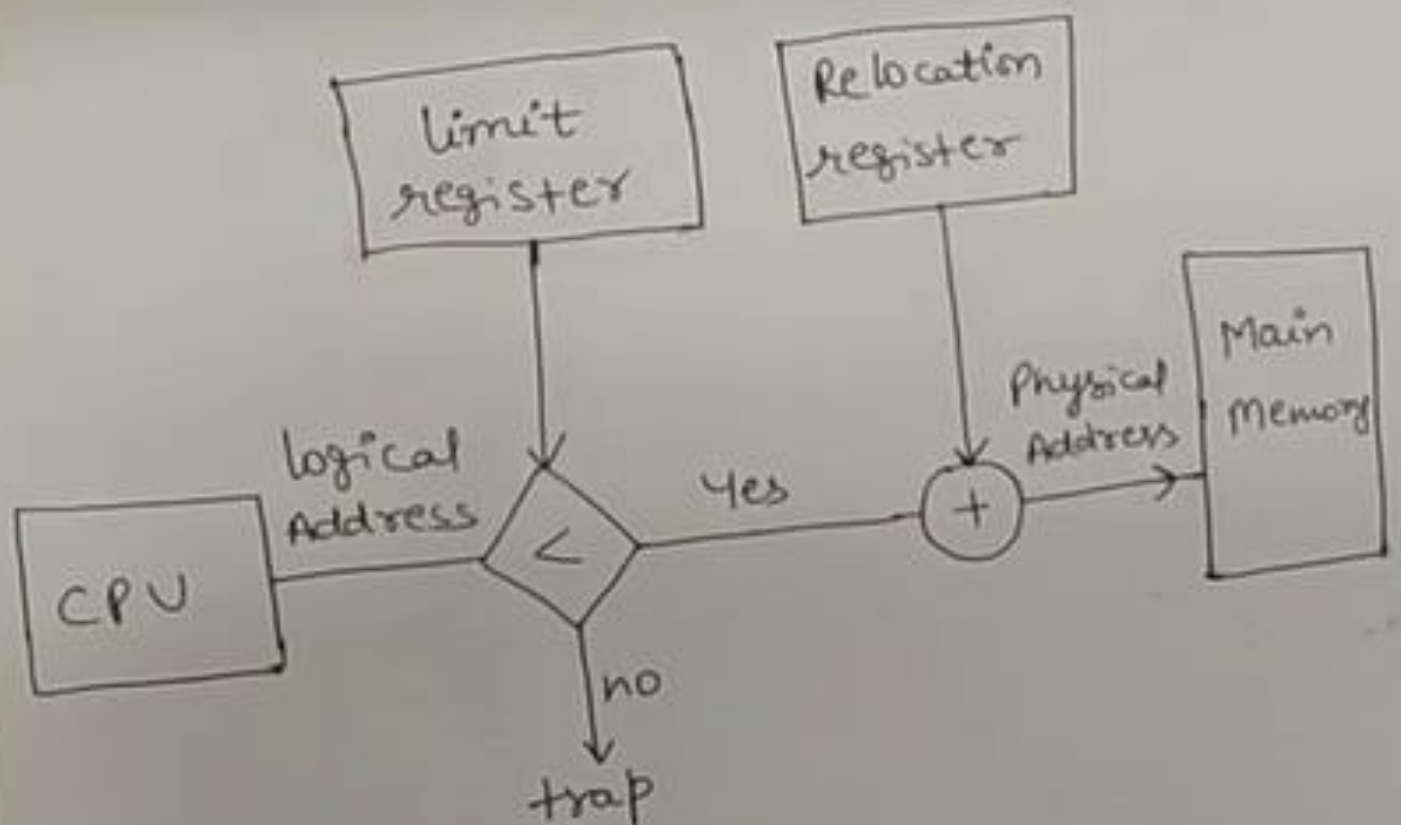
→ **Fixed partition Scheme**

- no. of partitions are fixed in memory.
- Size of each partition may be same or may be different.
- for eg:- we made 8 partitions in the memory. So, the no. 8 is fixed.
  In each partition only one process can be placed. In this case, we can place only 8 processes at a time. so, the degree of multiprogramming will be restricted.

→ Contigious Memory Allocation
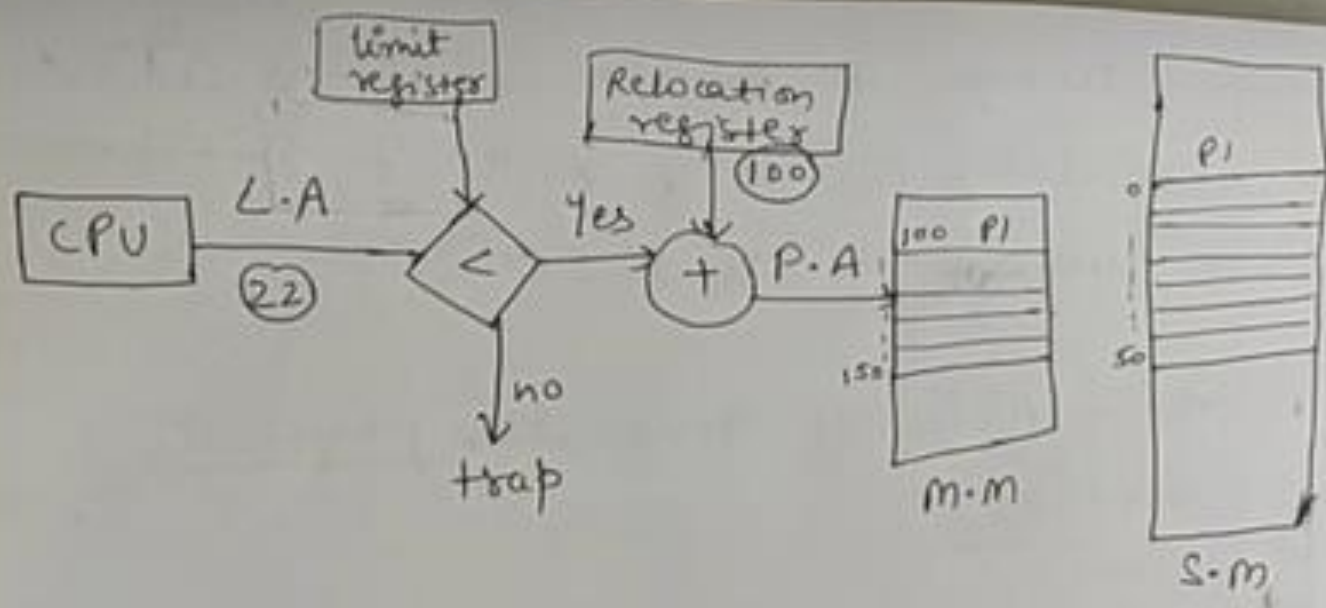
## Address Translation

- Logical Address to Physical Address



→ Two issues we face in Contiguous memory Allocation -:

① Fragmentation ⟨ → internal
              → external

② we have to translate logical Address to physical Address.

## Explanation of this above diagram

Suppose, we have a program P1 which contains 50 instructions in Secondary memory. CPU wants to execute instruction no. 22. CPU generates the logical address (LA) 22. And the programs which are to be executed will reside in the main memory. So, this P1 program is stored in the Main memory. Since main memory generates the Physical address and CPU generates the logical address, we have to translate LA to PA.

So, when CPU generates the logical Address 22, we will check the limit register. Limit register contains the size of the program, in our example the size of the program is 50. So, we will check $22 < 50$, if no it will go to trap, and if yes, then we will add the L·A (22) to relocation register to get the physical address.

The relocation register contains the base address of the program in main memory. ie in our example base address is 100. So, we add $100 + 22 \Rightarrow \boxed{122} \rightarrow$ which is the physical address.

As, we use contigious memory allocation the instructions are stored in a continuous manner. In this way we

Can translate Logical address (LA) to physical Address (P·A)

→ Note:

① Set of logical Addresses is Known as Logical Address space (L·A·S)
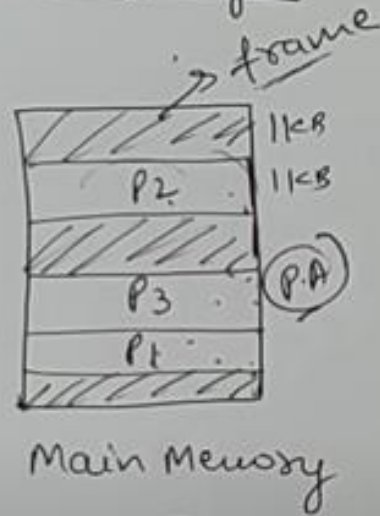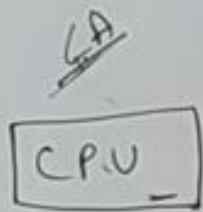
② Set of physical Addresses is known as Physical Address Space (P·A·S)
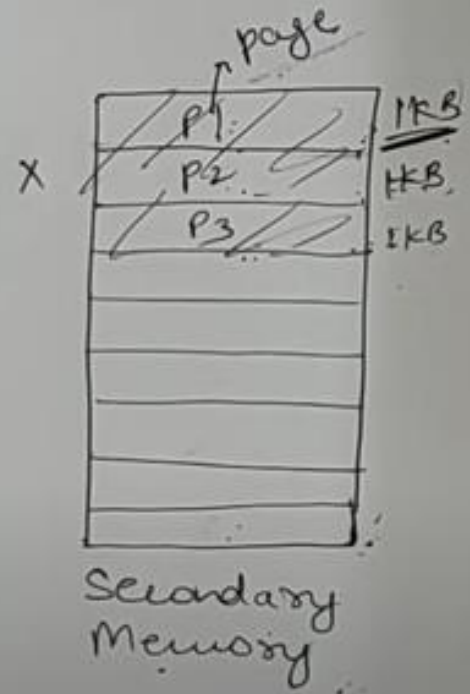
③ CPU generates Logical Address (LA)

④ Main memory Generates Physical Address (P·A).

# Non-Contiguous Memory Allocation

## Paging



frame → 1KB / P2 1KB / P3 (P.A) / Pt — Main Memory

CPU, LA

Program X (3KB)

page → P1 1KB, P2 1KB, P3 1KB — x — Secondary Memory
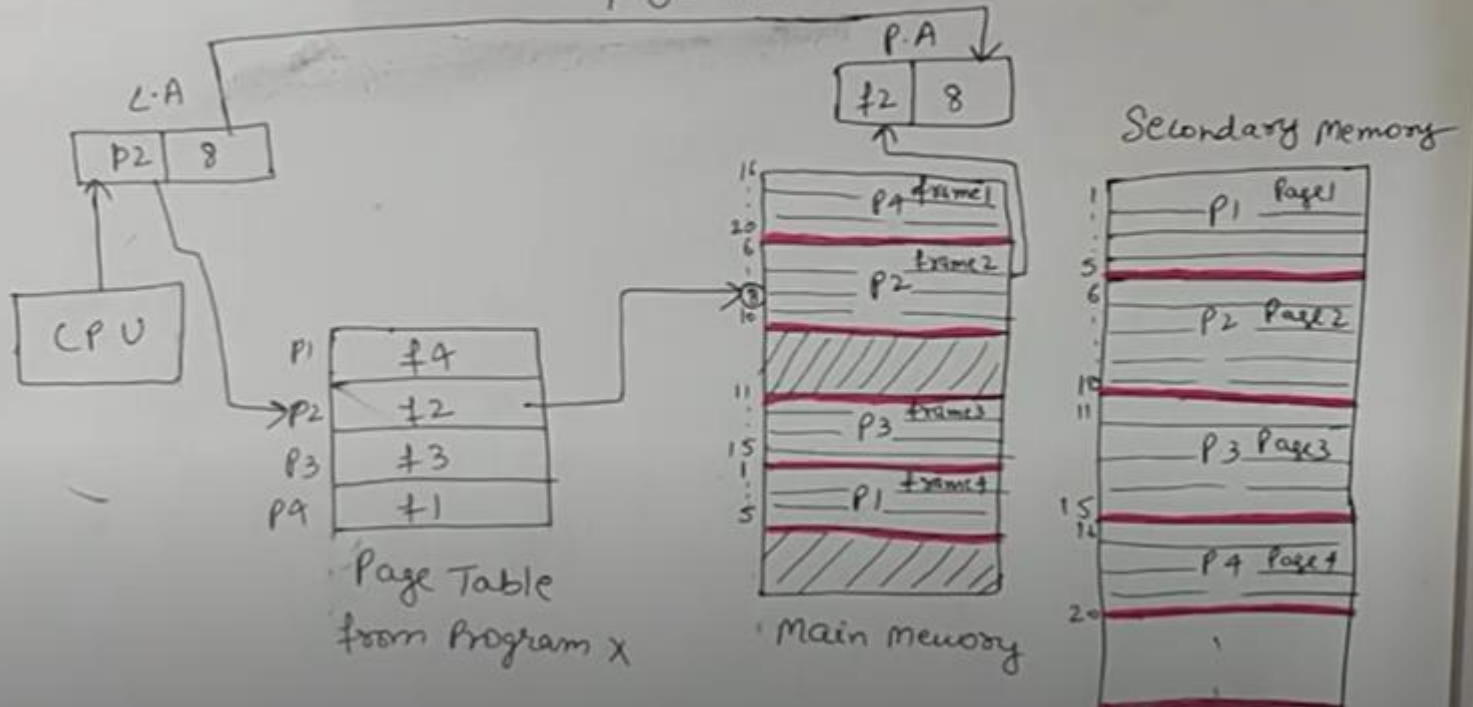
---

# Defination of paging

Paging is a memory Management scheme by which a computer stores and retrieves data from secondary storage for use in Main Memory. Operating System reads data from secondary storage in blocks called __pages__, all of which have identical size. And the Blocks (identical size) in Main Memory are called __frames__.
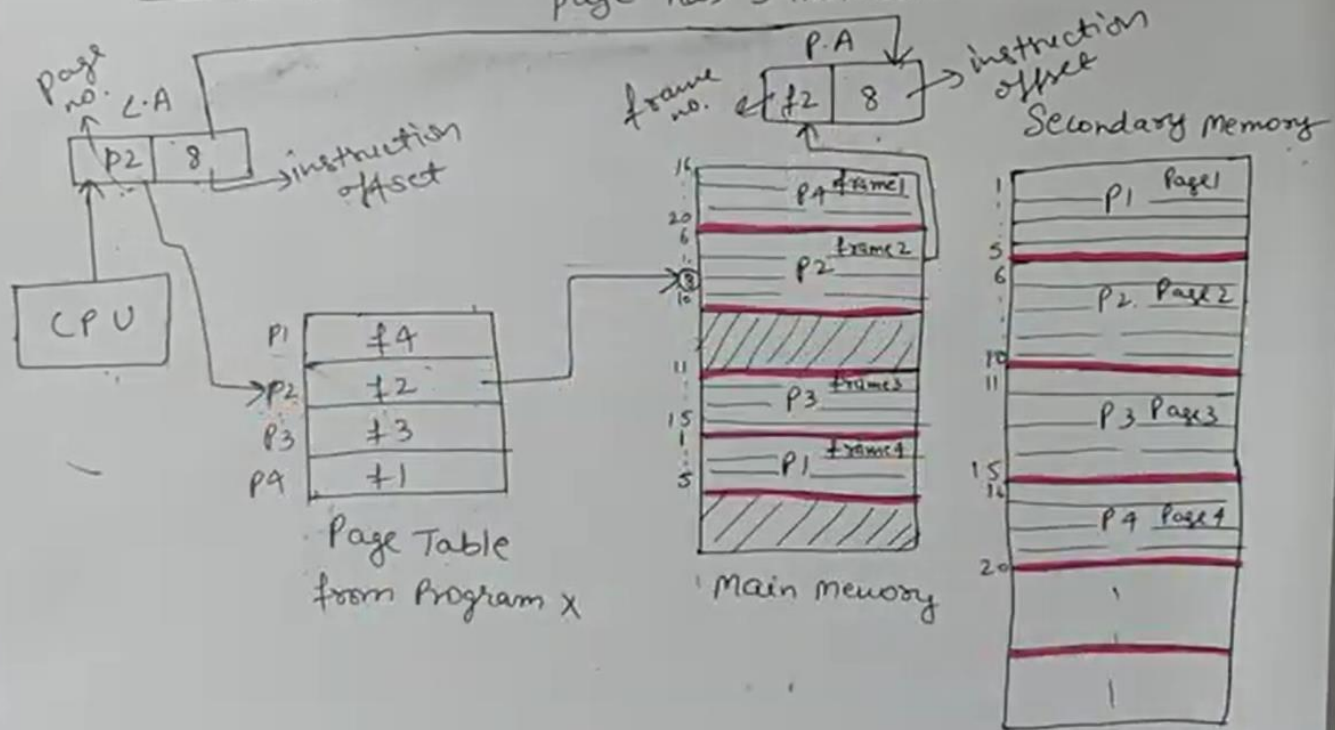
# Address Translation In Paging

Program X has 4 pages each page has 5 instructions.



Figure 1: Address translation in paging — page table for Program X maps to main memory and secondary memory.

L.A: P2 | 8

P.A: 72 | 8

## Page Table from Program X

| | |
|----|----|
| P1 | 74 |
| P2 | 72 |
| P3 | 73 |
| P4 | 71 |

CPU

Main memory frames: P4 frame1, P2 frame2, P3 frame3, P1 frame4

Secondary Memory:
- P1 Page1
- P2 Page2
- P3 Page3
- P4 Page4

---

# Address Translation In Paging

Program X has 4 pages each page has 5 instructions.



page no. L.A: P2 | 8 → instruction offset

frame no. → P.A: 72 | 8 → instruction offset

## Page Table from Program X

| | |
|----|----|
| P1 | 74 |
| P2 | 72 |
| P3 | 73 |
| P4 | 71 |

CPU

Main memory frames: P4 frame1, P2 frame2, P3 frame3, P1 frame4

Secondary Memory:
- P1 Page1
- P2 Page2
- P3 Page3
- P4 Page4

**Foreg.**

<u>Defination of Page Table</u> — It is the data structure stored in Main Memory that stores the mapping between Virtual addresses and Physical addresses.

Now CPU generates Logical address. The logical address consists of 2 parts — ① Page no. ② Instruction offset then CPU refers the Page Table and finds out which page no. is present at which frame no. and hence translate the Physical address. Physical address consists of two parts — ① frame no. ② Instruction Offset. foreg:- CPU wants to execute Page no. 1 and inside Page no. 1 instruction no. 5. f.e

| P1 | 5 |
|----|---|

L·A

---

P1 is page no. and 5 is the instruction offset. now CPU refers Main memory and finds out say P1 is present at frame no. 4 (f4).
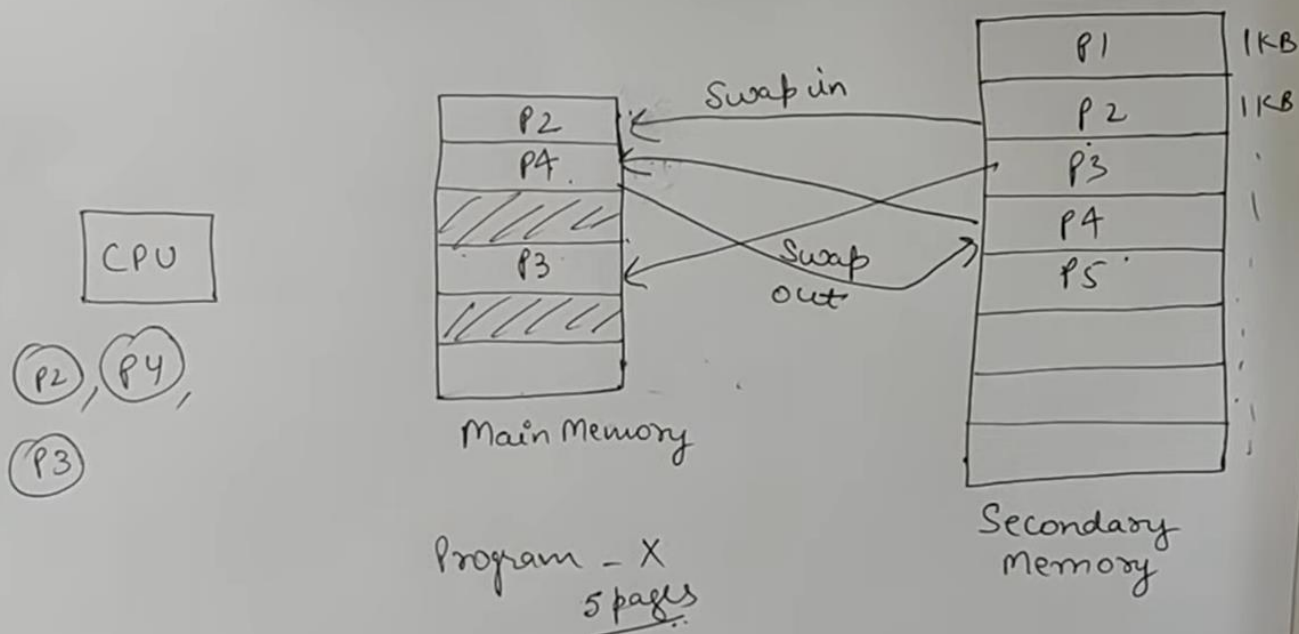
| f4 | 5 |
|----|---|

P·A

The instruction offset is same 5. CPU goes to frame no. 4 and picks 5th instruction and executes it and gives the output. This is how we translate L·A to P·A in paging.

# Demand Paging    Swapping



CPU

(P2), (P4),

(P3)

Main Memory

Program - X
5 pages

Swap in

Swap out

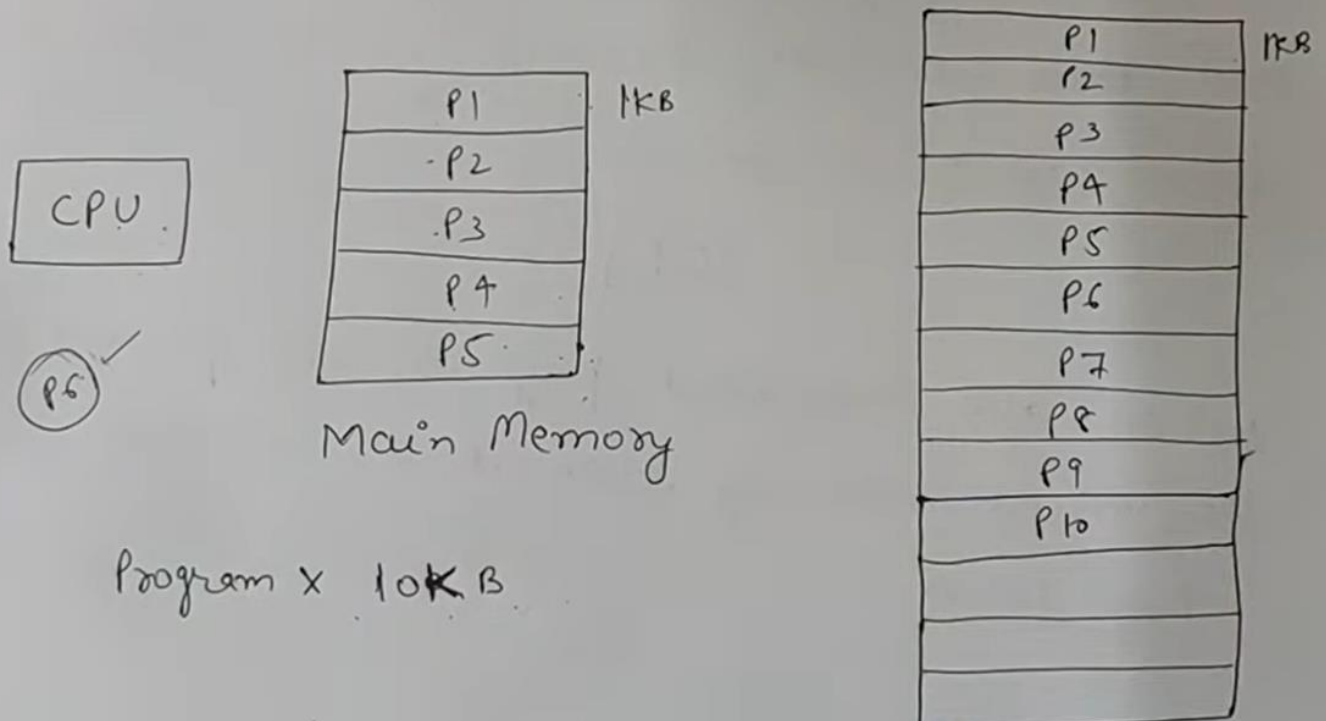| P1 | 1KB |
| P2 | 1KB |
| P3 | |
| P4 | |
| P5 | |

Secondary Memory

---

# Demand Paging

The Program resides in the secondary Memory, and pages are loaded in main Memory only on demand not in advance.

## Defination of Demand Paging → It is a type of swapping in which pages are copied from secondary memory to main memory when they are needed.
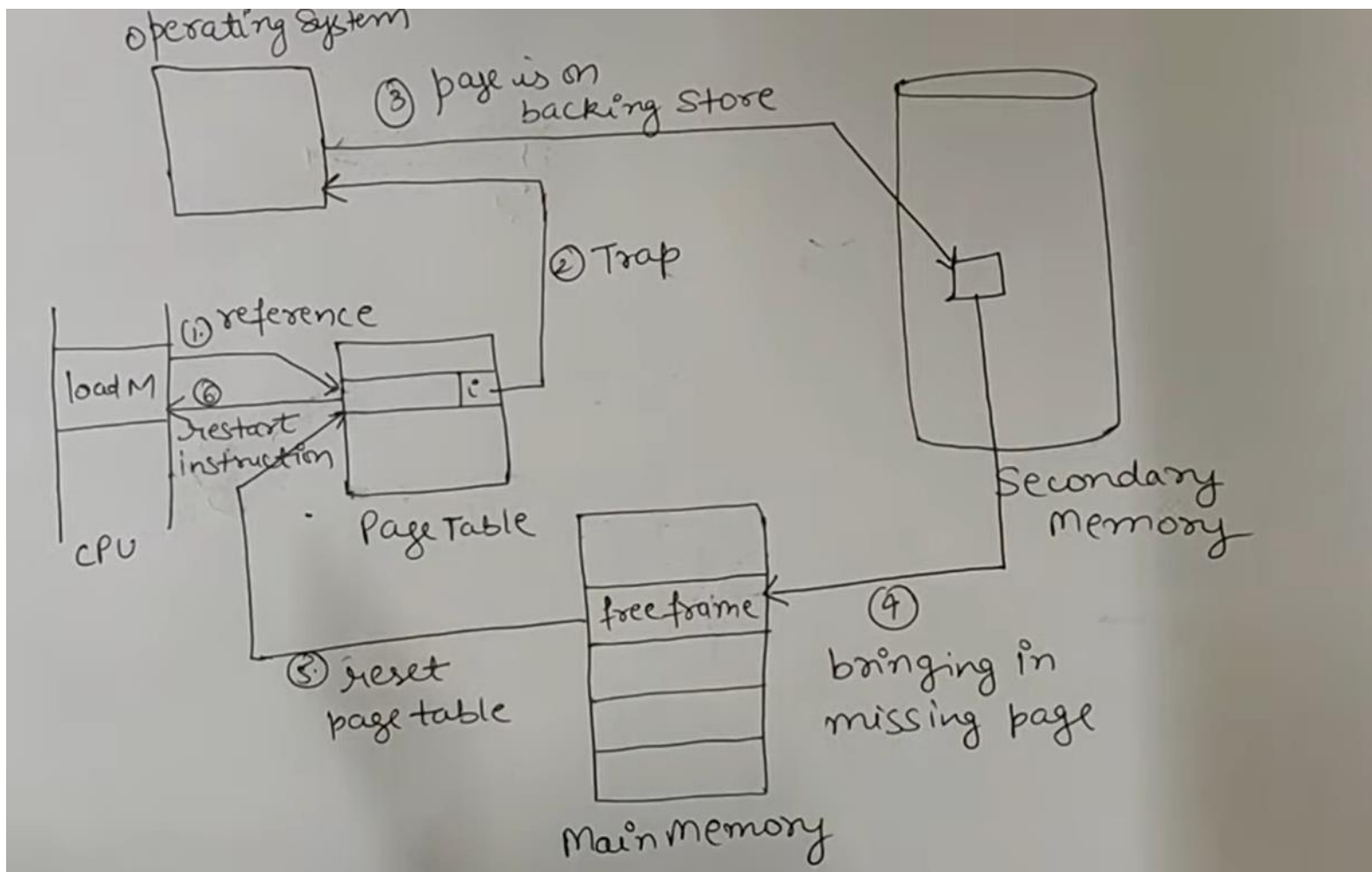
# Page fault



**Definition of Page fault** → when a processor wants to access a page, and that particular page is not ~~cuter~~ currently present in the main memory, then page fault occurs.
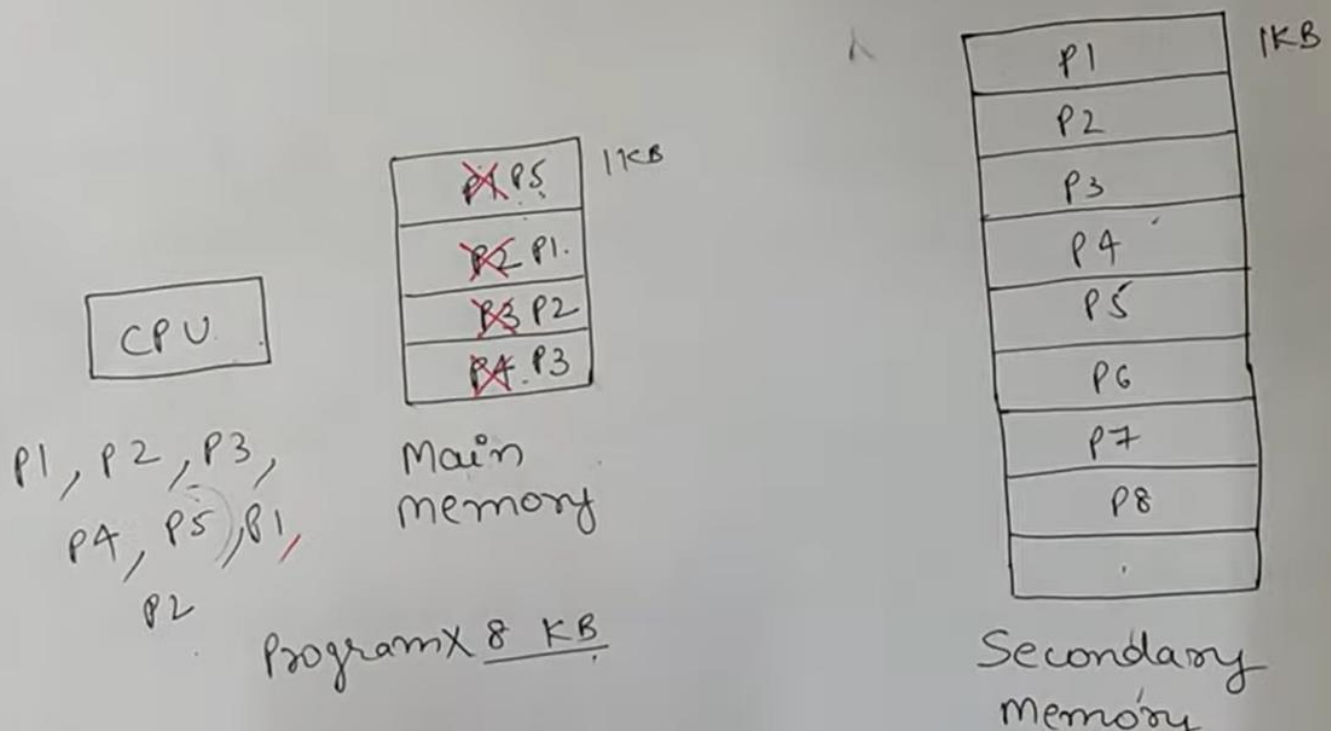
## OR

when the demanded page is not present in the main memory then it is called as Page fault.
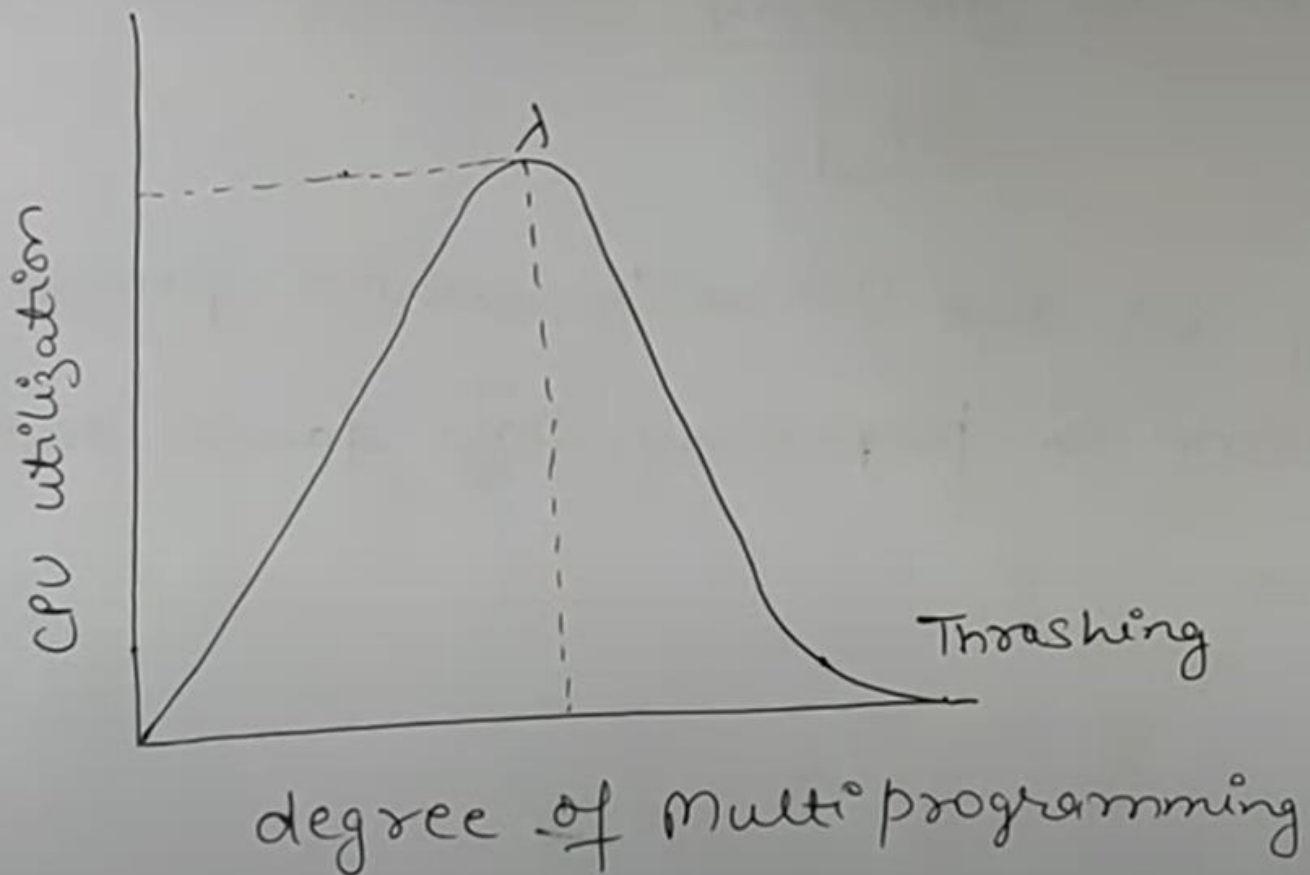
operating System

③ page is on backing store

② Trap

① reference

load M

⑥ restart instruction

CPU

Page Table

i

Secondary Memory

free frame

⑤ reset page table

④ bringing in missing page

Main memory

# Thrashing



CPU

P1, P2, P3,
P4, P5, P1,
P2

Program X 8 KB

X P5    1 KB
X P1.
X P2
X P3

Main
memory

| P1 | 1KB |
| P2 | |
| P3 | |
| P4 | |
| P5 | |
| P6 | |
| P7 | |
| P8 | |

Secondary
Memory

① In the initial stage when we increase the degree of Multiprogramming, CPU utilization is very high upto 1.

② Further, if we increase the degree of Multipr-ogramming, CPU utilization is drastically falling down.

③ This situation in the system is **Thrashing**.

④ Thrashing degrades the performance of System.

⑤ There can be a situation when main memory is full of pages that are accessed frequently. A page fault will occur if the required page is not present.

⑥ In order to make space for swapping in the required page, on of the frequently accessed page is swapped out.

⑦ Soon, the swapped out page is required for execution and this again results in



Degree of Multiprogramming: Number of pages are being added.

Defination of Thrashing → high paging activity is known as Thrashing.

$$\boxed{OR}$$

Thrashing is the Situation where process spends more time in processing page faults rather than execution.

## Virtual Memory

- Virtual Memory gives an illusion to the programmer that programs of larger size than actual physical memory can be executed.

- Virtual memory doesn't really exists, but the part of secondary memory are made as Virtual memory.

- Defination of Virtual Memory
  Virtual Memory is a feature in O.S, where large programs can store themselves in form of pages
  while their exe...

pages or portions of processes are loaded into the main memory.

Virtual Memory can be implemented by using -:

Demand Paging

Demand Segmentation

---

<u>Race Condition Defination</u>

When more thane one processes are executing the same code or resource or any shared Variable in that condition there is a possibility that the output or ~~bad~~ the Value of that Shared Variable is wrong, So for that all the proceses doing race to say that ~~so~~ my output is correct, this condition is known as Race condition.

# Critical Section in O.S

do

{

entry Section → controls the entry into C.S and gets a lock on resources.
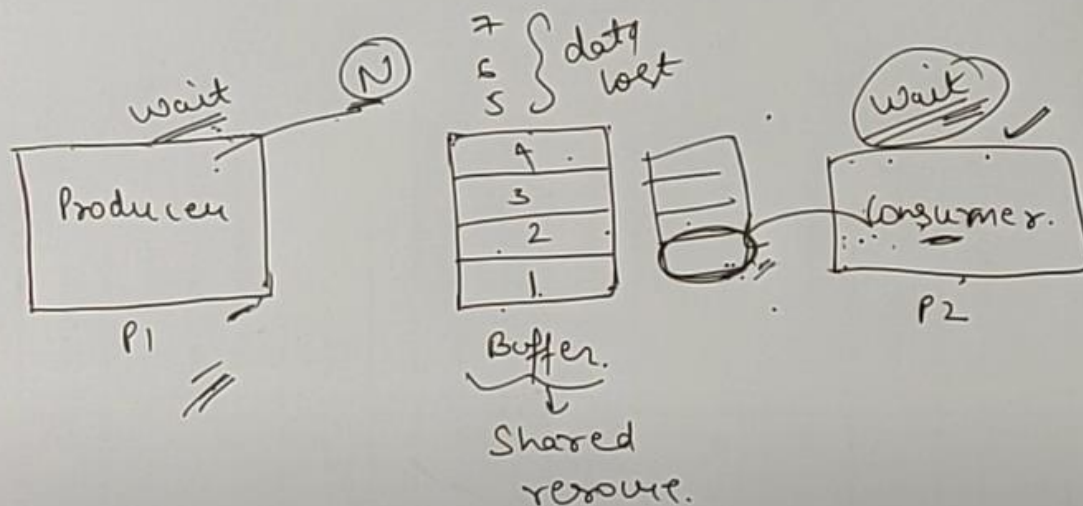
critical Section —

exit Section → removes the lock from the resources.

Remainder Section

} while (True);

---

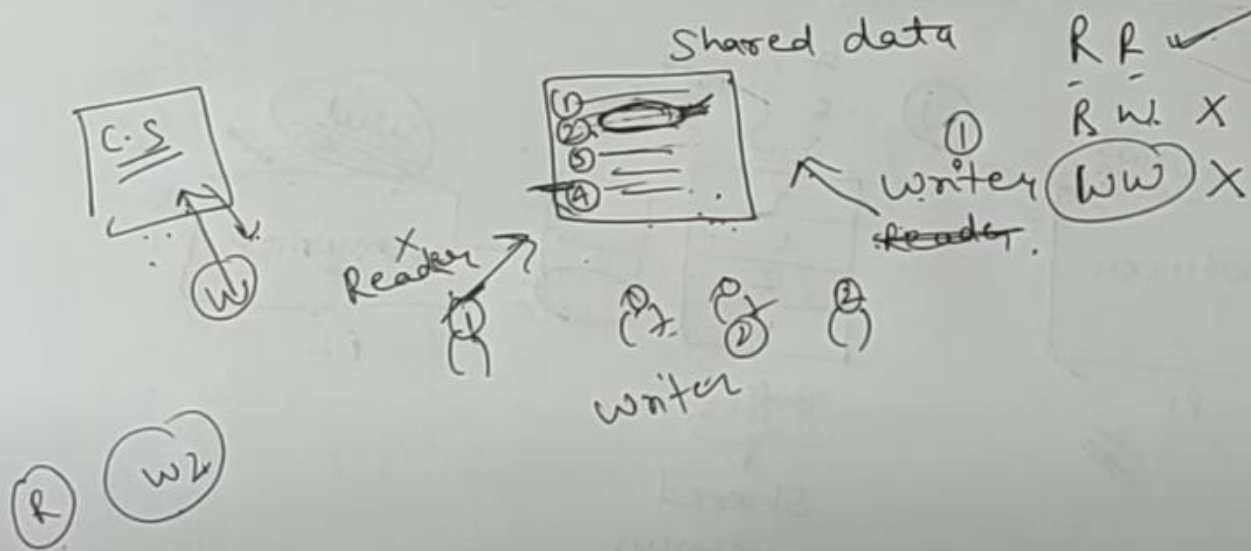## Classical Problem of Synchronization

### Producer-Consumer / Bounded Buffer Problem
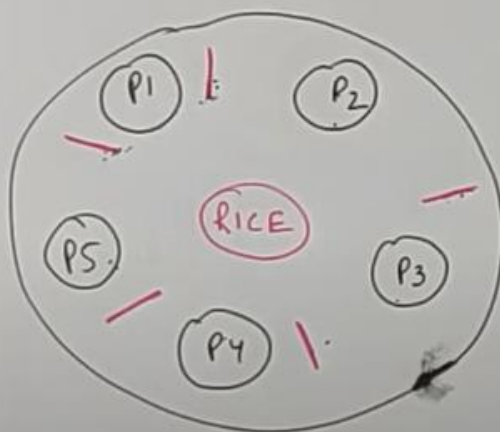
Producer/Consumer problem is solved using semaphores.



# Classical Problem of Synchronization

## Reader Writer Problem

Shared data

R R ✓
R W X
writer W W X
Reader

C.S

Reader

writer

W2

R



# Classical Problem of Synchronization

## Dining philosphers Problem.

thinking ✓
eating

Solution

max 2 person's
Can eat at
the same
time

P1
P2
RICE
P5
P3
P4

## Semaphores in Operating System

Semaphore is an integer variable that solves the critical section problem. After initialization, it can only be accessed by two atomic operations —:
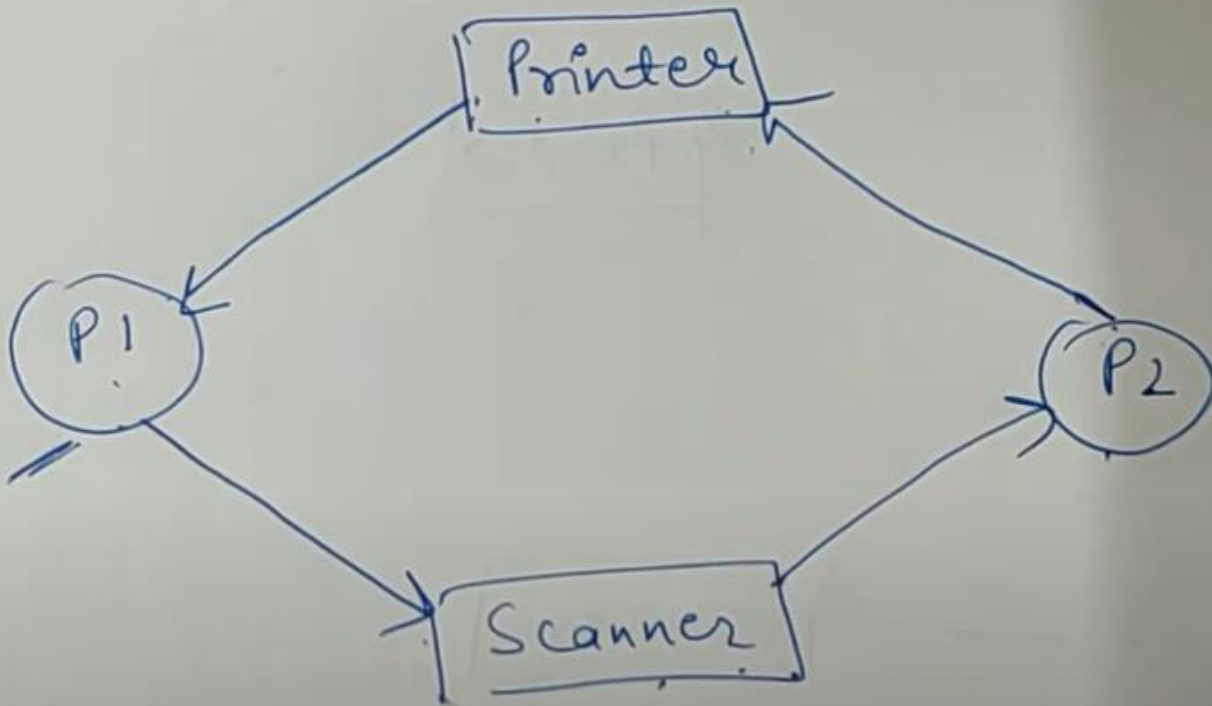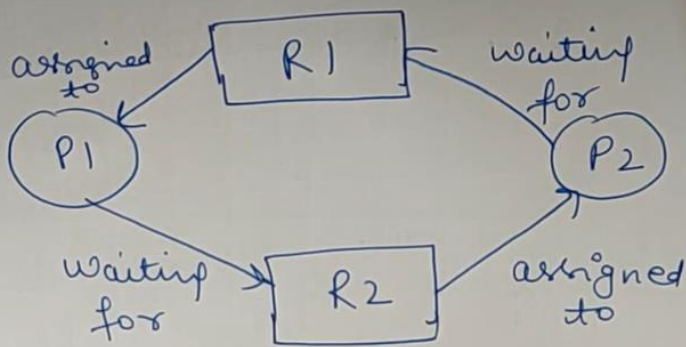
① wait (s)

    { while (s<=0);

        S = S-1;

    }

② signal (s)

    { S = S+1;

    }

## Deadlock in Operating System

assigned to · R1 · waiting for

P1 · P2

waiting for · R2 · assigned to

**Defination of Deadlock** → Deadlock is a situation where a set of processes are blocked because each process is holding a resource and waiting for another resource acquired by some other process.

→ **Four Necessary conditions for deadlock to occur**

① **Mutal Exclusion** → one or more than one resource are non sharable (only one process can use it).

②. **Hold and wait** → A process is holding at least one resource and waiting for another resource.

③ **No Preemption** → A resource cannot be taken from a process unless the process releases that resource.

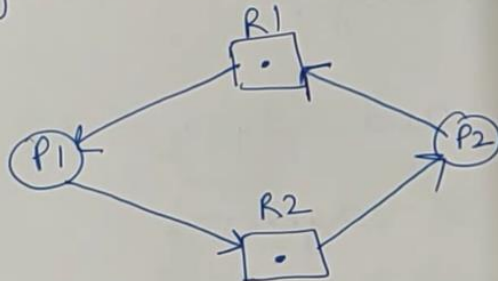④. **circular wait** → A set of processes are waiting for

# Resource Allocation Graph (RAG)

- Using RAG, deadlock can be easily detected.
- It is a Graph that represents the state of a system pictorially.
- It has vertices and edges.
- It has two types of edges —:
  ① Request edge    $P_i \longrightarrow R_j$
  ② Assignment edge    $R_j \longrightarrow P_i$

---

# Rules to detect deadlock using RAG

**Rule-1** → In RAG, where all the resources are Single instance —:

ⓐ If a cycle is being formed, then the System is in a deadlock State.



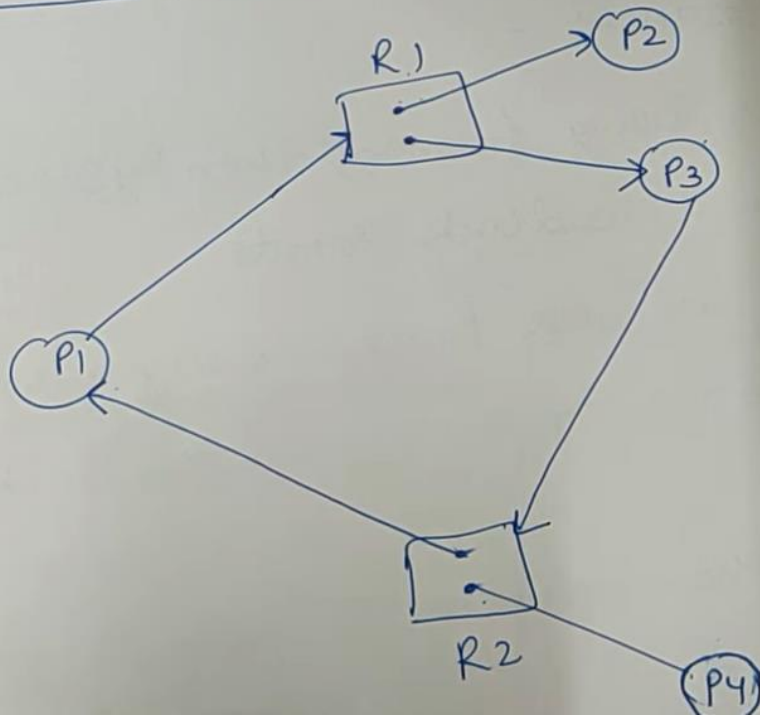ⓑ If no cycle is being formed, then the Syste is not in a deadlock State.

**Rule -2** → In RAG, where all the resources are multiple instances -:

ⓐ If a cycle is being formed, then System may de may not be in a deadlock state.

In this case we use Banker's algorithm to confirm whether a System is in deadlock state or not.

ⓑ If no cycle, then System is not in a deadlor State.

RAG, with multiple instances but no deadlock

# Banker's Algorithm (Deadlock Avoidance Algo)

Total A = 10, B = 5, C = 7

| Process | Allocation | | | Max need | | | Available | | | Remaining need | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | A | B | C | A | B | C | A | B | C | A | B | C |
| P1 | 0 | 1 | 0 | 7 | 5 | 3 | 3 3 2 / 2 0 0 | | | 7 4 3 P1 | | |
| P2 | 2 | 0 | 0 | 3 | 2 | 2 | 5 3 2 / 2 1 1 | | | 1 2 2 P2 | | |
| P3 | 3 | 0 | 2 | 9 | 0 | 2 | 7 4 3 / 10 3 2 | | | 6 0 0 P3 | | |
| P4 | 2 | 1 | 1 | 4 | 2 | 2 | 7 4 5 / 0 10 | | | 2 1 1 P4 | | |
| P5 | 0 | 0 | 2 | 5 | 3 | 3 | 7 5 5 | | | 5 3 1 P5 | | |
| | 7 | 2 | 5 | | | | | | | | | |

Safe seq

P2 → P4 → P5 → P1 → P3

Remaining need = Max need – allocation; for all

A1 (Available 1) = Total (A, B, C) – Allocation (sum (A, dim=1), sum (B, dim=1), sum (C, dim=1)); #dim 1 means column

Available for P1 = A1

Now execute the processing by looking at the remaining need and focusing on the available column;