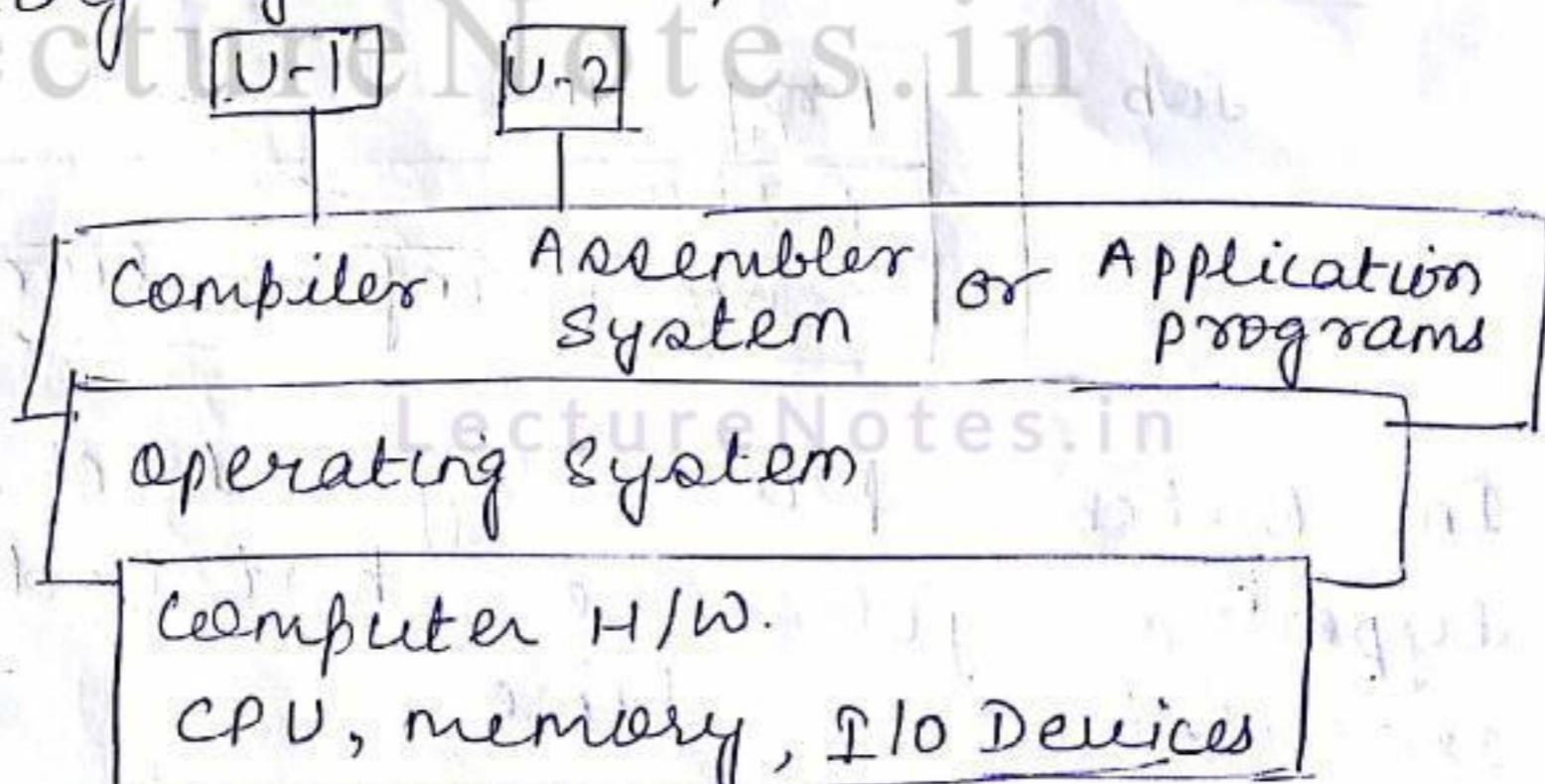


11/7/17

## Resource manager/ allocator.

- 1.) Introduction
- 2.) Process management :- storage and execution of program.
- 3.) process synchronization
- 4.) deadlock : tackling problems.
- 5.) Memory management
- 6.) File management.

- An Operating system(OS) is a set of program that act as an interface between the user and the computer hardware.
  - The purpose of OS is to provide an environment in which a user can execute programs in a convenient and efficient manner.
  - OS is a resource manager or Resource allocator. examples of resources are: CPU, memory, I/O devices
  - OS also behaves as a service provider.
- what operating system do?



### ① user-view

- PC: maximise work.
- mainframe (MFC) : maximise resource utilization

### ② System view

- Resource allocation

12/7/17

## → computer system operation

- i) initialize : All devices are connected or not  
ii,) OS will be loaded.

Kernel will be loaded in RAM.  
↳ (core part of processor)

Interrupt : unwanted signal

Interrupt vector : interrupts are stored in array.

LectureNotes.in

Normal execution



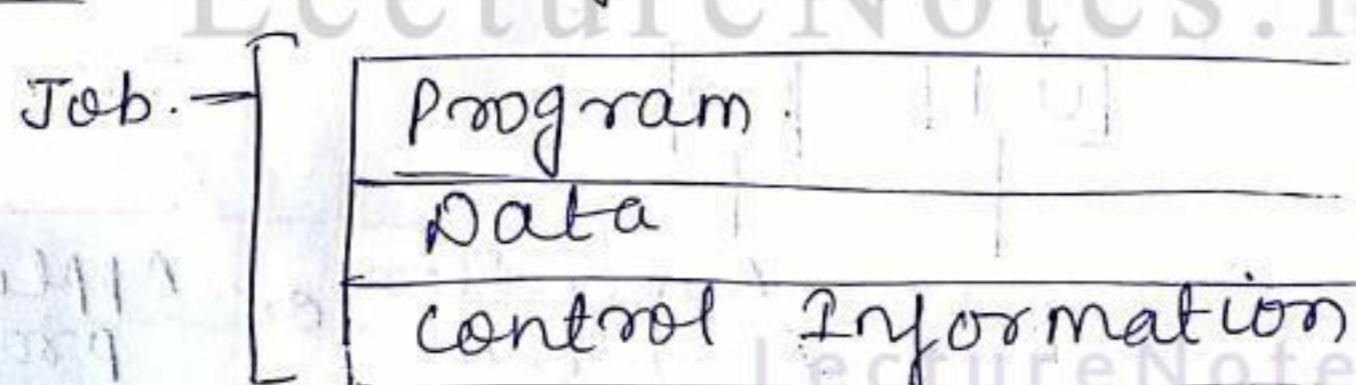
I.V.

interrupt service routine  
ISR  
I =  
}

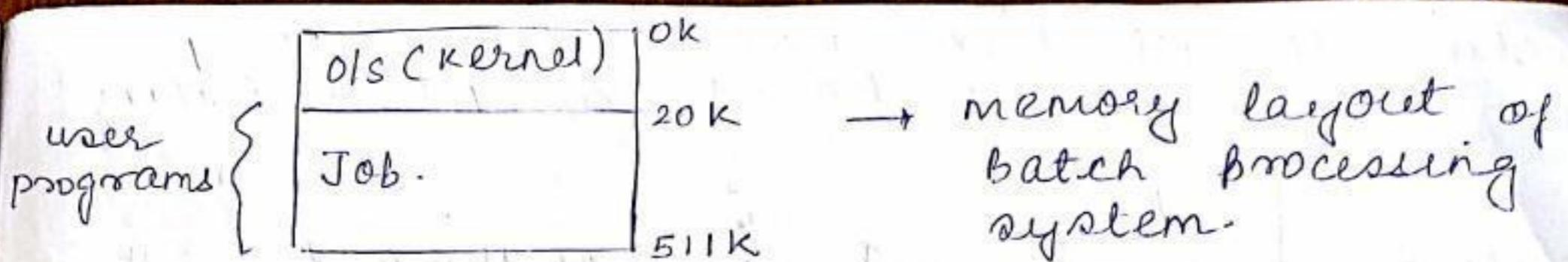
system call : Deviating from main program and executing interrupt and again coming back and executing main program is called system call.

## → Evolution of operating system.

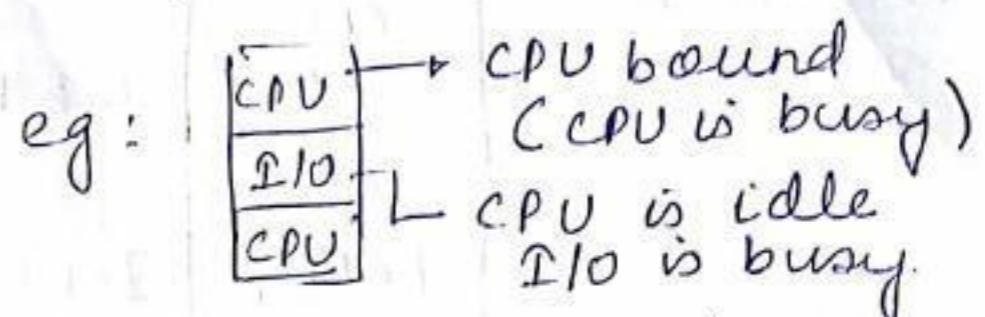
### i.) Batch processing system.



- In Batch processing technique similar types of jobs are batched together and executed at a time.
- The major task of OS was to transfer control automatically from one job to the next job.



Job  $\leftarrow$  CPU bound.  
Job  $\leftarrow$  I/O bound.



CPU should always be busy so as to increase performance.

### - Advantage of Batch processing system

i.) processing is faster

### - Disadvantage

i.) CPU is often idle.

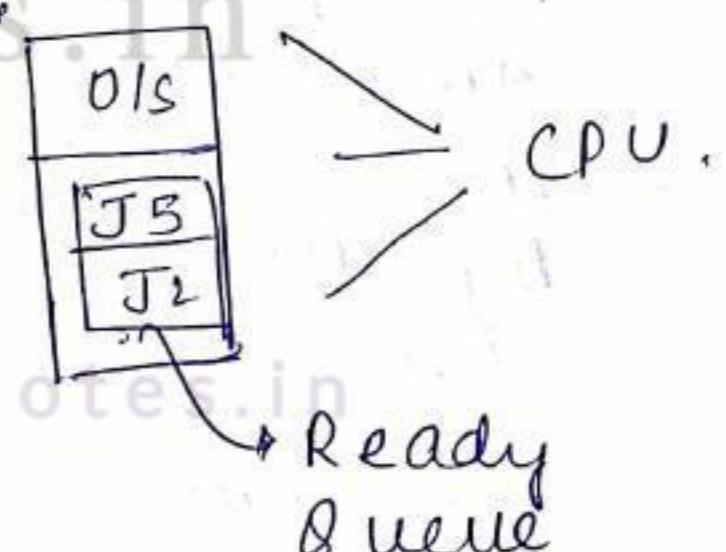
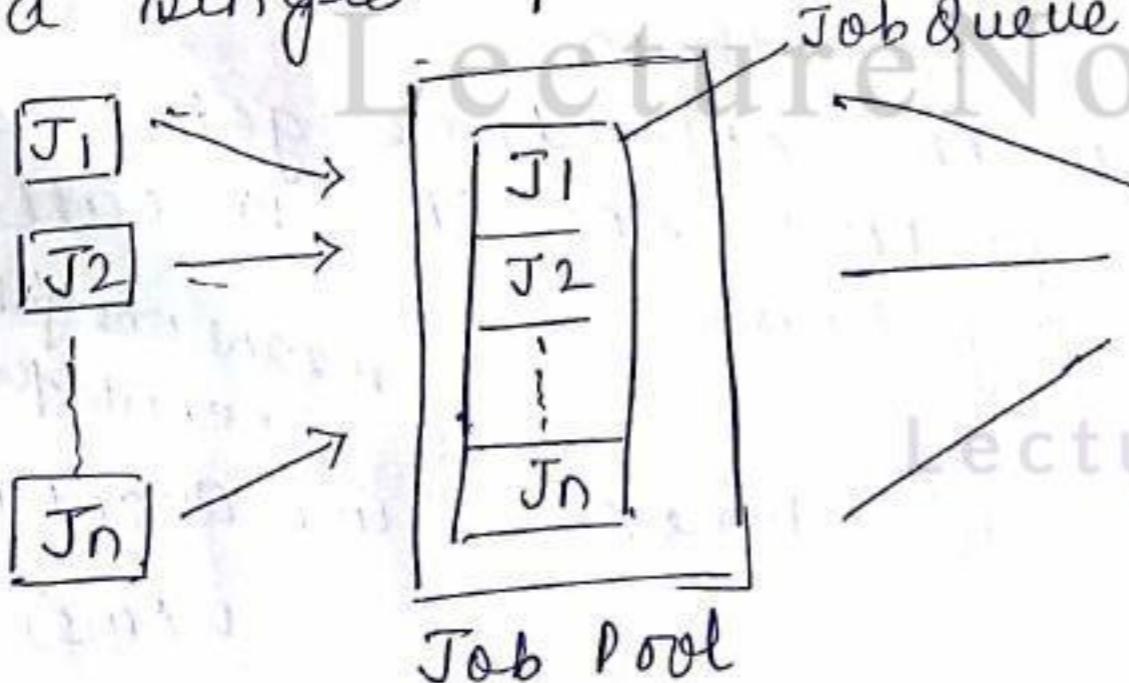
ii.) user cannot interact with the job when it is executing.

### ii) Multiprogramming System :-

To avoid (CPU is often idle) disadvantage this system was introduced.

Multiprogramming is a technique to execute multiple programs simultaneously by a single processor.

Job scheduling



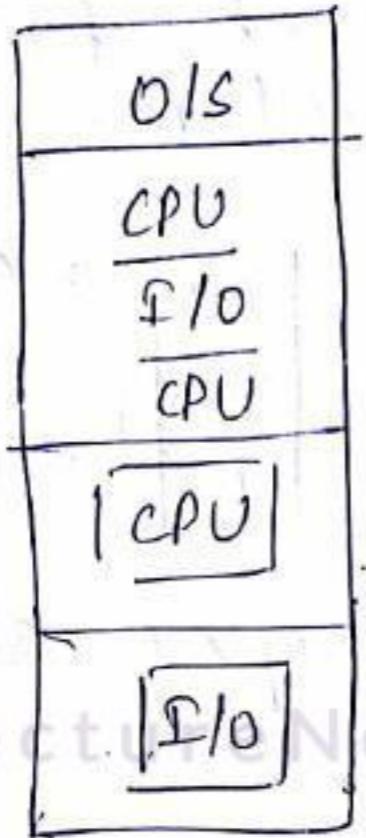
main memory

Ready Queue queue which contains jobs and those jobs are ready for execution. CPU can any time start its execution.

Job Queue : Here Jobs may or may not be ready.

Selecting the job from Ready Queue and processing it in CPU is based on CPU scheduling.

Execution starts



↑ CPU works this much.  
↓ Job 1 I/O.

Job 2

Job 3

comes here  
works till  
here then  
returns as  
CPU work is  
remaining in Job 1

### → Advantage

- efficient memory utilization as multiple jobs are stored.
- CPU is never idle.
- Throughput of the CPU may also increase.
- Waiting time is limited.

1817117

### (iii) Time-sharing operating system

- Time-sharing or multitasking is a logical extension of multiprogramming.
- Multiple jobs are executed by the CPU switching between them.
- In this method the CPU time gets shared by different programs so it is called time sharing system.

(means one job can execute for max. of 4ms)

Program

Execution time

Time Quantum

4(ms)

J<sub>1</sub>

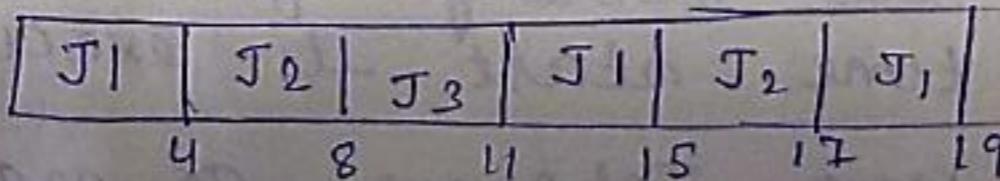
10.

J<sub>2</sub>

6.

J<sub>3</sub>

3



### - Advantage:

- i.) user can interact with the job when it is executing.
- ii.) the efficient use of CPU i.e. CPU utilization.

### IV) MULTIPROCESSOR SYSTEM or PARALLEL SYSTEM or

#### a) lightly coupled system:

In multiprocessor system have more than one processors will share computer bus, memory, peripheral devices etc. Then this type of system are called multiprocessor system.

#### Advantage

- i.) increased throughput (programs executed in particular time period: throughput).
- ii.) It is economical. { (e.g.: 5 CPU in 5 comp. & 5 CPU in one)
- iii.) increased reliability.

a.) The ability to continue providing service proportional to the level of surviving hardware

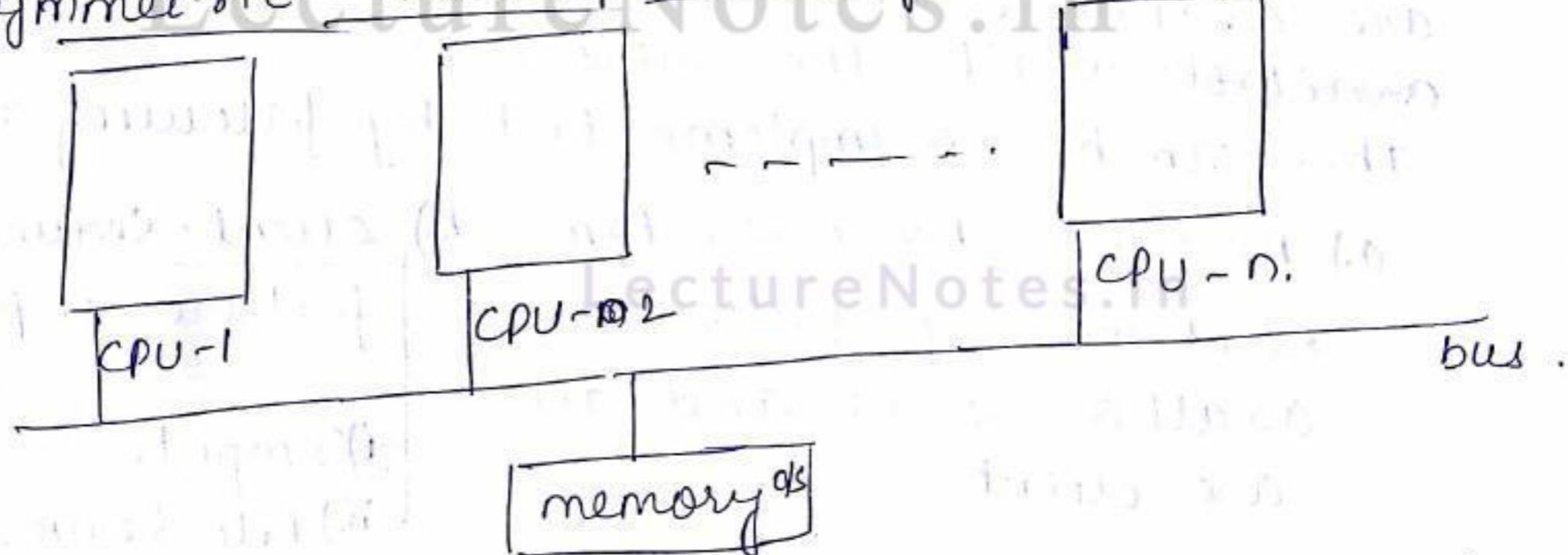
#### degradation

(graceful degradation  
max. no. of failure & in which system can operate efficiently.

- b) some systems go beyond the graceful degradation are called fault tolerance.

There are two types of multiprocessor system.

#### i) Symmetric multiprocessing (SMP)



- each processor runs a copy of the OS and these copies communicate with one another when needed
- In SMP all the processors are peers

## (ii) Asymmetric multiprocessing system:

- In this case each processor is assigned a specific task
- A master processor controls the whole system and other processors follow the master processor's instructions
- The master processor schedules and allocates the work to the slave processor.

## b) Distributed system or loosely coupled system.

In distributed systems, the processors cannot share memory, clock and peripheral devices. Each processor has its own local memory and they are communicated among themselves through high speed buses (LAN/WAN/Internet).

### Advantage:

- i) Resource sharing.
- ii) computation speedup.
- iii) Reliability.

1917117

## v) Network operating system.

A Network has a number of nodes and a node has many resources. When we want to use the resources, we need to know where the resources are available and connect that machine to avail the resources.

This can be implemented by following architectures.

### a) Peer-to-Peer System : b) client-server system.

No chance of failure as all are server and all are client.

failure is possible  
i) compute-server system.  
ii) file-server system.

## vi) Real-time operating system.

- A Real-time system is used when a rigid time requirement have been placed on the operation of a processor or the flow of data.

- A real-time system has well-defined fixed time constraint. Processing must be done within the defined constraint otherwise the system will fail.
- There are two types of real time system.
  - a) Hard real time system: This system guarantees that critical tasks are completed on time.
  - b) soft real time system: It means that only the precedence and the sequence for the task operations are defined.

## COMPONENTS OF OPERATING SYSTEM (18) pg-24.

### i) Process management:

program in execution is called a process.

program is passive entity whereas process is an active entity.

when a process is executed/created, a resource is allocated to it.

when execution is over the process releases all the resource..

One which manages all program and processes is done by process management.

### ii) Main-memory management

main memory: large collection of semiconductor cells capable of storing one bit of information either 0 or 1.

PA	1111
	1001
	0110
	1010

CPU generates operation.

CPU → logical address

memory management unit

↓  
physical address

In main memory.  
randomly we can  
access data. and time to  
access any address is same.

### iii) File management      file: collection of information.

Folder: collection of file. All folders make :superfolder type

file      character

numeric

Alphanumeric

Binary

20/7/17

#### 4. I/O system management:

{ I/O are electromechanical device. It cannot be connected directly due to speed mismatch.  
so some buffer/cache memory is put in low-speed device (eg: printer)

Here buffer is I/O subsystem.

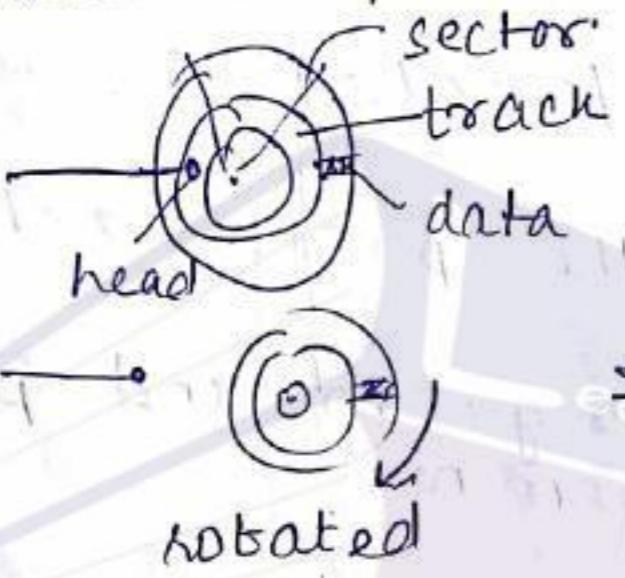
some predefined softwares / drivers are also installed sometimes.

#### 5. Secondary storage management:

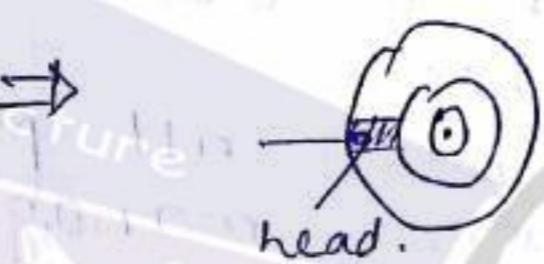
main m/m : volatile (power on → works).

80

disc is partitioned into tracks. It is further divided into sectors.



To read head is used.



and data is read.

fragmentation

{ free space one side  
data on other side  
defragmentation

#### 6. Network: LAN, WAN, Internet

program → program needs some resource that is not with me → then I need to access that resource from a resource system. This is provided in this Network system.

#### 7. Protection System

If two processes <sup>are running</sup> both should be protected from each other

2 types

i.) user process : we write program & execution

process

ii.) system process : when we access some program ~~from our~~ we need some system program. and if we access it then it becomes ~~system~~ system process.

## 8. Command Interpreter system:

Command Interpreter:

eg: in dos       $\text{dir} \rightarrow$  will show list of directory  
                ↳ this is a program to show list.  
↳ interface b/w user and OS

→ Operating System SERVICES. (LQ). Pg-54.  
OS provides services to user, system etc.

### 1. Program execution

program executes  $\xrightarrow{\text{normally}} \text{O/P}$   
 $\xrightarrow{\text{abnormally}} \text{error}$

### 2. I/O operation

while executing program we give input  
using keyboard, mouse, LAN, some file..

### 3. File System Manipulation

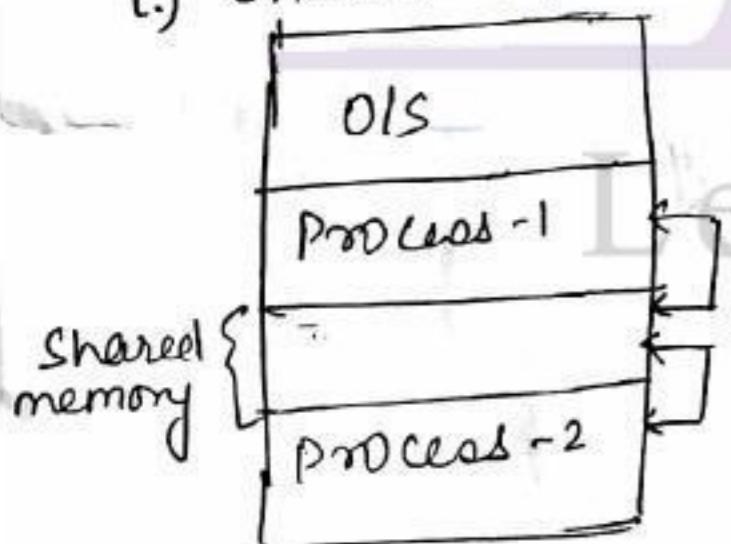
Reading data from file  $\rightarrow$  write/append etc in file.

### 4. Communication

2 processes are running  $\rightarrow$  They may communicate  
among themselves : i) shared m/m concept    ii) message  
passing  
Two process within system:

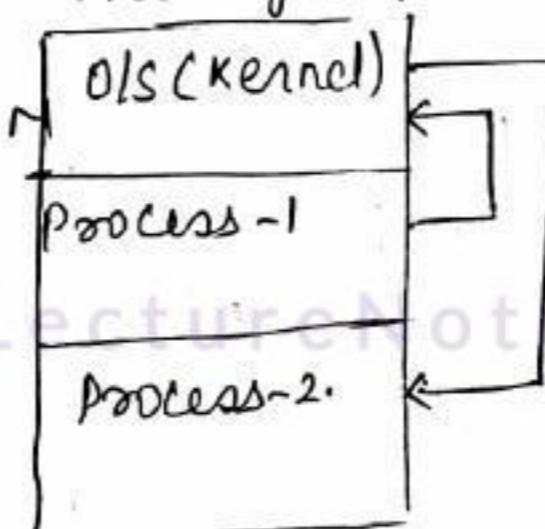
Two process in diff. system:

#### i) Shared m/m



both process can send or receive data

#### Message passing



P-1 P-2 commu  
P-1 ~~comm.~~  
with O.S then  
O.S communicates with P-2

### 5. Error detection

Errors generated by I/O, CPU, memory, program  
is handled by this.

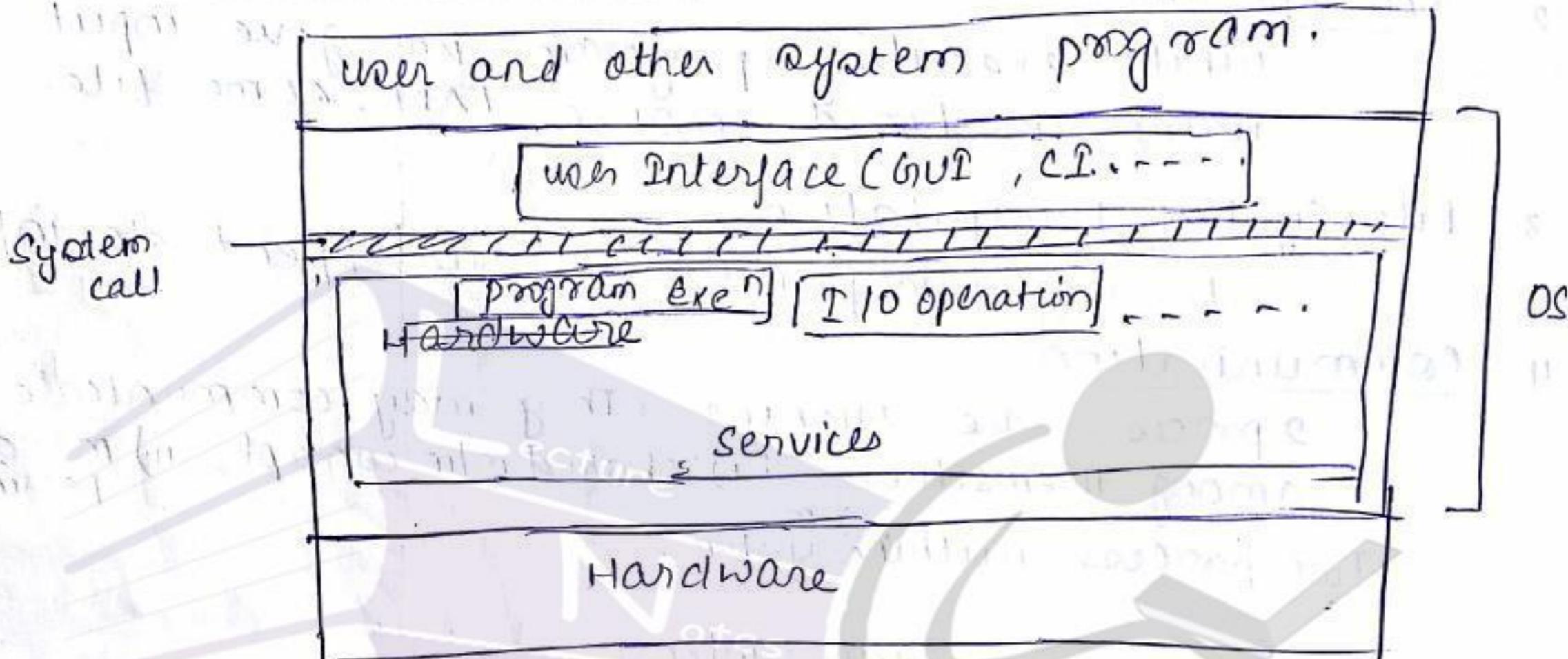
### 6. Resource Allocation

Allocating resources to process or a system at  
time of execution.

- 7.) Accounting:  
to keep track of how many users use how many resources.
- 8.) Protection and security:  
protection: protecting one process to another  
security: whether or not you are authorized person to use system. To test it userid and password is used.

21/7/17

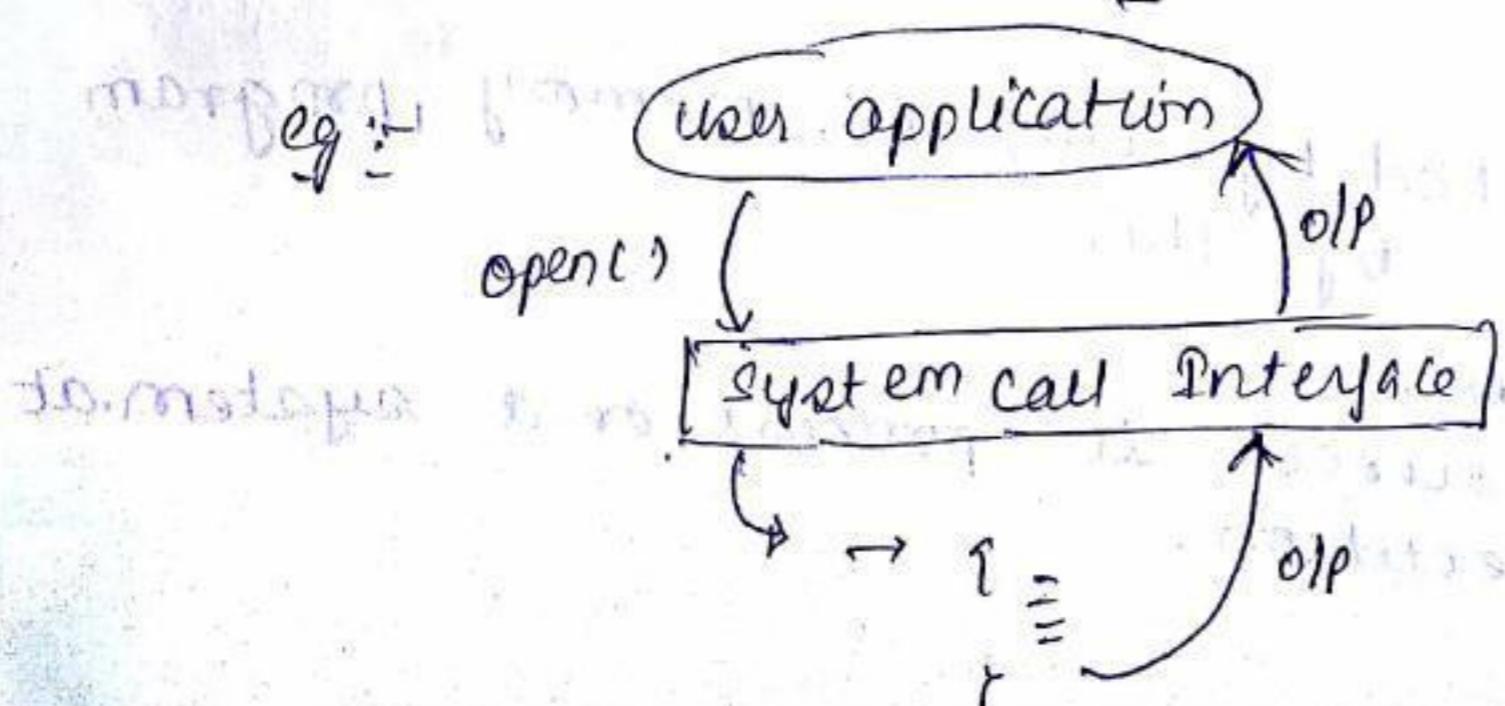
LectureNotes.in

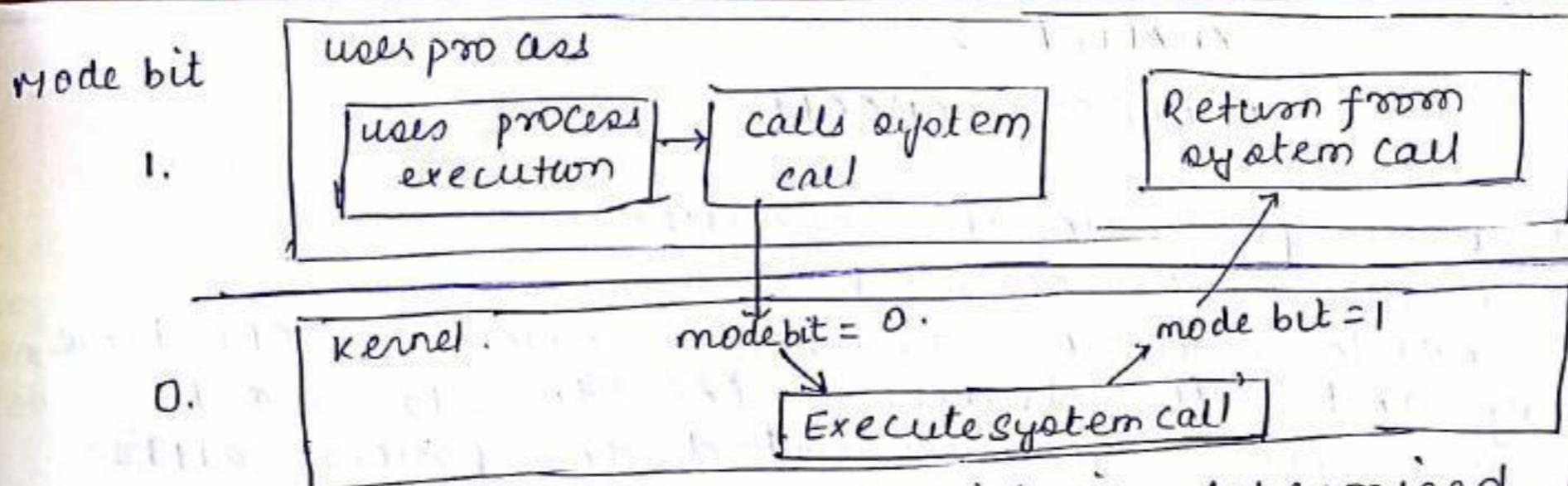


→ System call : provides the interface between the process and the OS. (Pg - 61).

Types of system call: there are 6 types :

- i.) process control
- ii.) file manipulation.
- iii.) device manipulation
- iv.) information maintenance
- v.) communications.
- vi.) protection. (Pg 65)





System runs in 2 modes which is determined by bit  
 if bit is 1 then system is in user mode  
 if bit is 0 then system is in kernel mode  
 mode bit is changed using interrupt and interrupt  
 is executed by system call.  
 This is how user and kernel process work without  
 hampering each other.



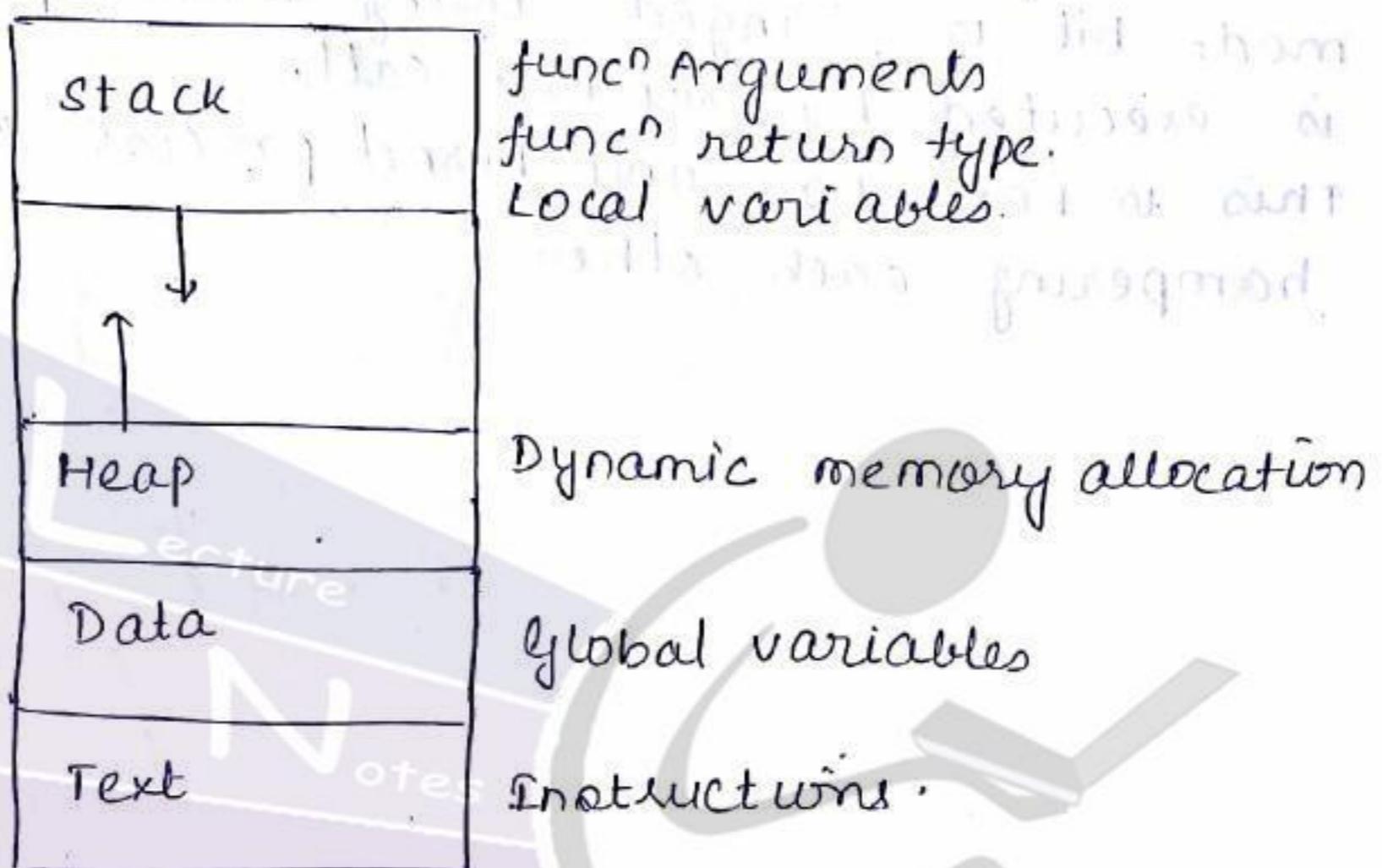
LectureNotes.in

21/7/17

## CHAPTER - 2 PROCESS CONCEPT

- Process is a program in execution.
- Process is a unit of work.
- Process needs certain resources such as CPU time, memory and I/O devices to do its work. These resources are allocated to process either at the time of creation or at the time of execution.
- Process : main m/m Program : Secondary m/m.

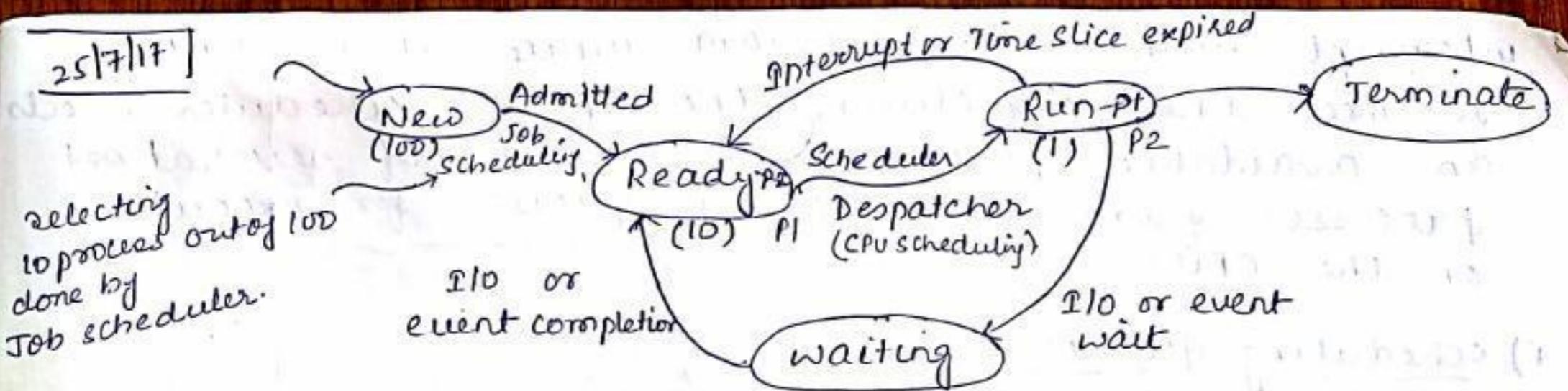
Internal structure of process :



### → Process state:

- When a process executes it changes the state
- The state of the process is determined by current activity of the process.
- The states are
  1. New : The process is being created.
  2. Running : The instructions are being executed.
  3. Waiting : The process is waiting for some event to occur.
  4. Ready : The process is waiting to be assigned to a processor.
  5. Terminate : The process has finished its execution.

$$1 \rightarrow 4 \rightarrow 2 \rightarrow 3 \rightarrow 4 \rightarrow 2 \rightarrow 5.$$



out of 10 processes CPU can run }  
 a single process at a time } done by CPU & scheduler  
 process: process / CPU scheduling

### [Working system of process state]

#### → Process control Block (PCB).

Each process is represented in the OS by its process control block or Task control block.

Pointer	state
Process Number	
Program counter	
Registers	
Memory limits	
list of open files	
list of open devices	
⋮	

PCB.

pointer : holds address of each PCB.

state : PCB is in which state is specified here.

process number : every process is identified by its no.

program counter : holds address of next instruction to be executed.

Registers : eg: Accumulator : for holding temp. results ; GPRS R0-R7

, DR : holds current instruction.

ML limit : base limit --- etc.

eg. condition code / flags : status register, eg. numlock

ON  Numlock  
OFF  Numlock

#### → Process scheduling

- The objective of multiprogramming is to have some process running at all times to maximize the CPU utilization.

- The objective of time sharing system is to switch the CPU among processes so frequently that user can

interrupt with each program while it is running

- To meet these objectives, the process scheduler selects an available process from a set of several processes from the ready queue for execution on the CPU.

### i) scheduling queues:

{ Job queue (newly created process)

Ready queue

Device queue (every device maintains a queue)

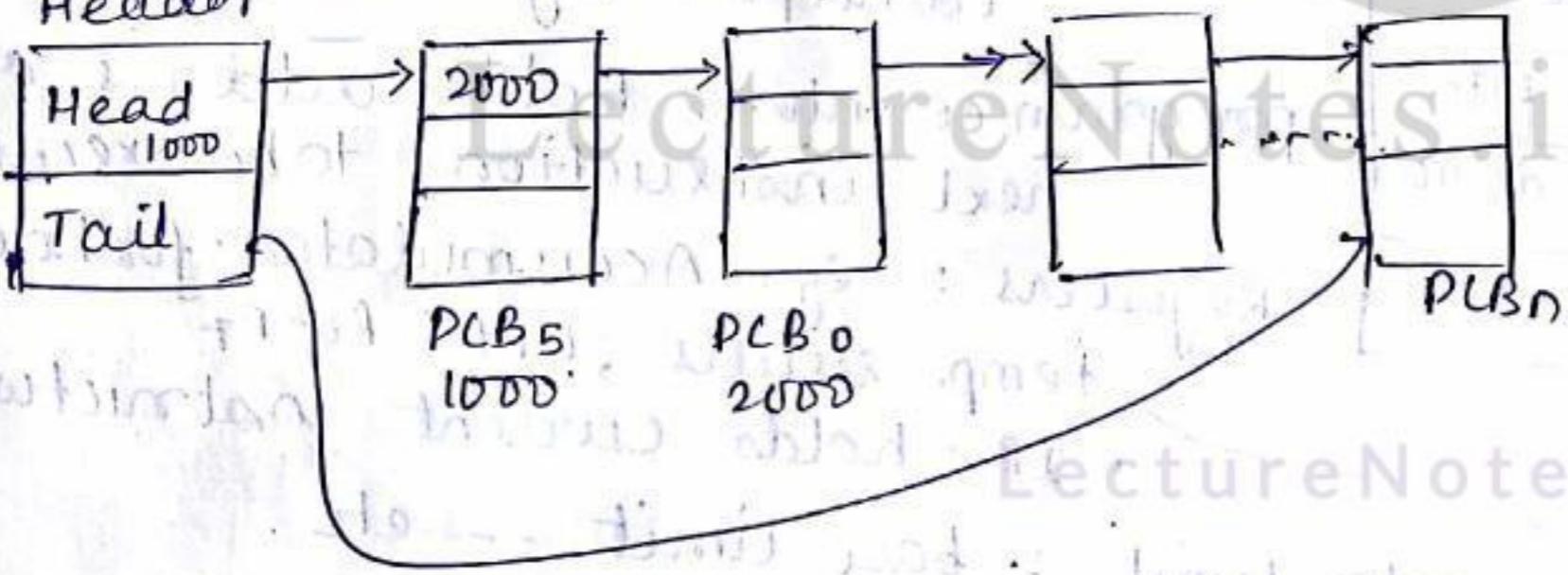
In process scheduling (Job queue is not used)

    └ Ready, Device queue.

- Job queue: When the process enters the system, they are put into a job queue. Job queue consists of all processes in the system.
- Ready queue: The processes that are residing in the main memory and are ready and waiting to execute are kept on a list called Ready queue.
- Device queue: The list of processes waiting for a particular I/O device is called device queue.

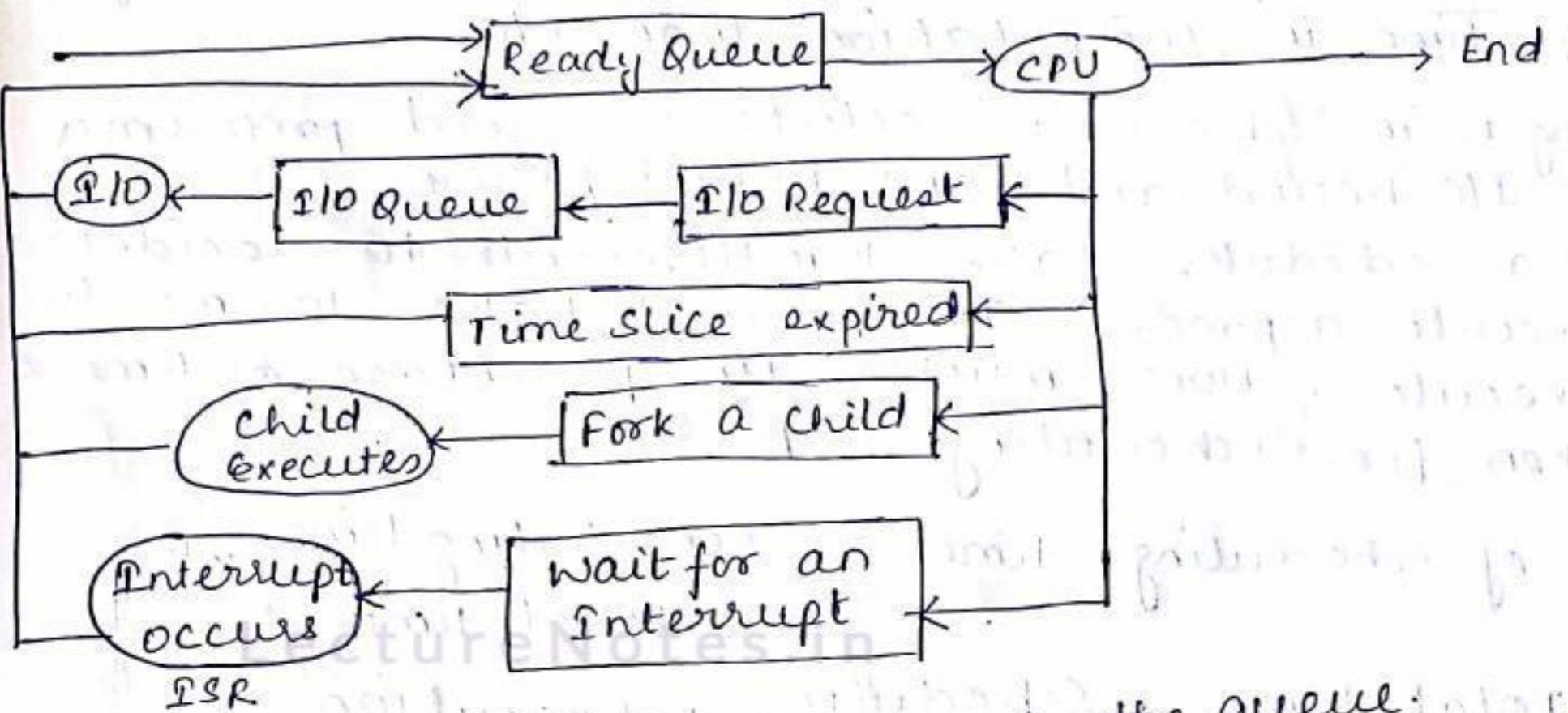
Every Queue is maintained in the form of linked list.

Header



26/7/17

→ Queuing Diagram Representation of process scheduling:



- each rectangular box represents the queue.
- Here two queues are used (Ready and Device).
- The circle represents the resources that serve the queues.

#### → Scheduler:

- The process migrates among various scheduling queues throughout its lifetime.
- The OS must select a process from these queues in some fashion. This selection process is carried out by a scheduler.
- There are 3 types of schedulers:
  - i.) Long-term scheduler or Job scheduler
  - ii.) short-term scheduler or CPU scheduler
  - iii.) Medium-term scheduler

#### i.) long-term or Job scheduler:

It selects a process from job pool and loads them into main memory (Ready Queue), for execution.

#### ii.) short-term

It selects a process from ready queue that are ready to execute and allocates the CPU to one of them.

- The number of processes in the memory is called degree of multiprogramming -
- I/O bound process is one that spends more of its time doing I/O then computations.

- A CPU bound process is one that spends more of its time in computation than I/O.

Note\* Long term scheduler selects a good process mix of I/O bound and CPU bound process.

Q: eg. If a scheduler takes 10 milliseconds to decide to execute a process. The process takes 100 ms to execute. How much time percentage of time is taken for scheduling a process?

$$\% \text{ of scheduling time} = \frac{\text{scheduling time}}{\text{Total time}} \times 100$$

$$\text{Total time} = \text{Scheduling time} + \text{Execution time}$$

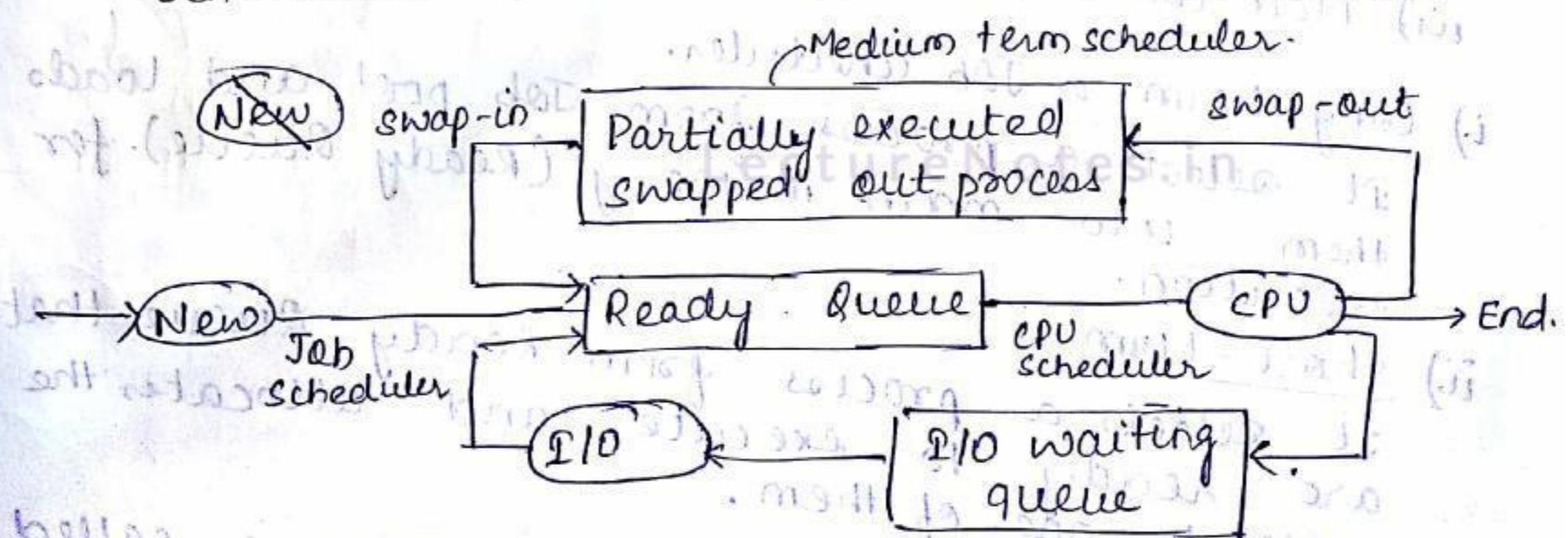
$$\therefore \% \text{ of scheduling time} = \frac{10}{110} \times 100$$

$$\frac{100}{11} = 9\%$$

$\therefore$  9% of total CPU time is wasted for scheduling a process.

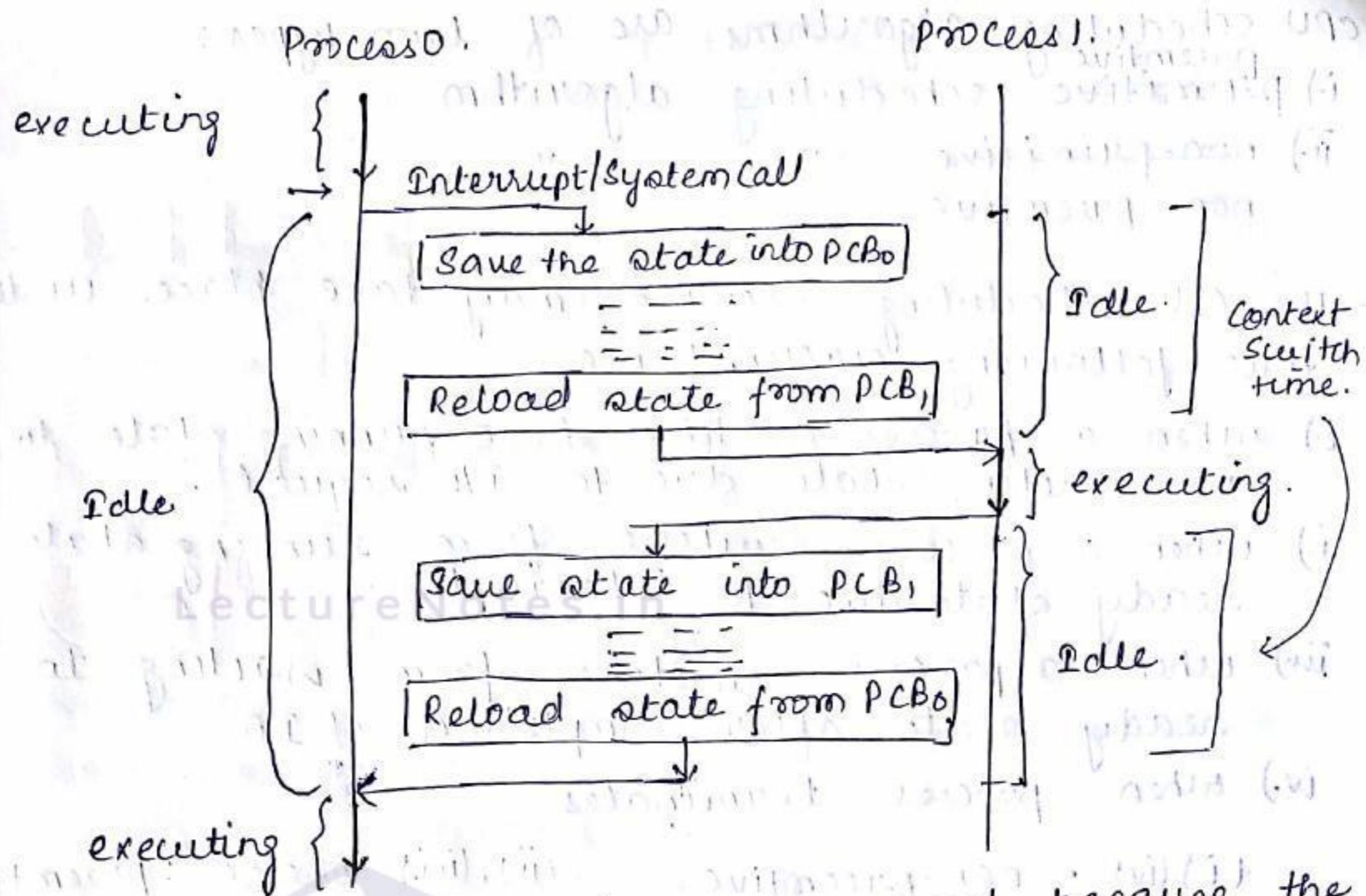
### iii) Medium term scheduler:

- In time sharing system, we are removing a running process from CPU to the ready queue after time slice expires. So removing a process from CPU to ready queue is done by a scheduler called medium term scheduler.



### Context Switch

Switching the CPU from one process to another process requires saving the state of the old process and loading the saved state for the new process. This switching is known as context switch.



context switch time is pure overhead because the system does no useful work while switching.

## → CPU / process scheduling Algorithms

27/7/17

- The CPU scheduling algorithm determines how the CPU will be allocated to the process.
  - CPU scheduler: The CPU scheduler selects a process from several available processes in memory (Ready Queue) that are ready to execute and allocates the CPU to one of them.
  - Dispatcher: The Dispatcher is a module that gives control of the CPU to the process selected by the CPU scheduler. The main function of the scheduler is switching the CPU from one process to another process.
- The Dispatcher should be fast because it is invoked during each and every process switch. The time it takes for the dispatcher to stop one process and start another process is known as dispatch latency/delay.

CPU scheduling algorithms are of two types:

- i.) ~~preemptive~~ scheduling algorithm.
- ii.) ~~non-preemptive~~ " "

non-preemptive

- The CPU scheduling decision may take place under four following circumstances.

- i.) when a process switches from running state to the waiting state due to I/O request.
- ii.) when a process switches from running state to ready state due to interrupt.
- iii.) when a process switches from waiting to ready state after completion of I/O.
- iv.) when process terminates.

(i),(iv) : non-preemptive. (ii),(iii) scheme : preemptive

- In preemptive scheduling, the CPU can release the process even in the middle of the execution.
- In non-preemptive scheduling, once the CPU has been allotted to a process, the process keeps the CPU until it releases the CPU either by terminating or by switching to the waiting state.

→ Scheduling criteria :

• The criteria determines the best algorithm among several best algorithm.

i.) CPU utilization : It is the percentage of time that the processor is busy.

ii.) Throughput : no. of processes completed per time unit is called throughput.

iii.) Turnaround time : The interval from the time of submission of process to the time of completion is called turn around time.

Turnaround time = time spent in Job queue + time spent in Ready queue + execution time of process  
+ waiting at waiting set + doing I/O.  
Turnaround time = finish time - arrival time

iv) waiting time: It is the amount of time that a process spends waiting in the ready queue.

$$\text{Waiting time} = \text{Starting time} - \text{Arrival time}$$

v) Response time: The time interval from submission of a request until the first response is known as Response time.

$$\text{Response time} = \text{First response} - \text{Arrival time}$$

i) ii) should be T      iii) iv) v) should be +

→ Algorithm

i) First come first serve

Q. Let a set of processes P1 to P5 arrives in following order

Arrival time      Burst time

P.NO.	AT	BT
1	0	4
2	1	3
3	2	1
4	3	2
5	4	5

Find the schedule

Grantt

P1	P2	P3	P4	P5
0	4	7	8	10

FR : First response

P NO.	AT	BT	CT	TAT CT-AT	WT TAT-BT	RT FR-AT
1	0	4	4	4	0	0
2	1	3	7	6	3	3
3	2	1	8	6	5	5
4	3	2	10	7	5	5
5	4	5	15	11	6	6

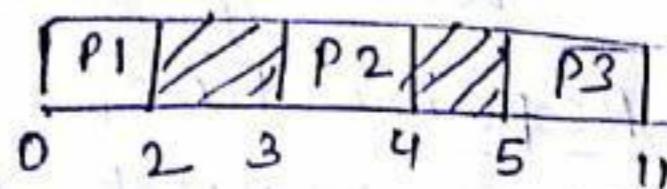
$$\text{Avg. Turn around time} = (4+6+6+7+11)/5 =$$

$$\text{Avg WT} = (0+3+5+5+6)/5 =$$

$$\text{Avg RT} = (0+3+5+5+6)/5 =$$

28/7/17

P NO.	AT	BT
1	0	2
2	3	1
3	5	6

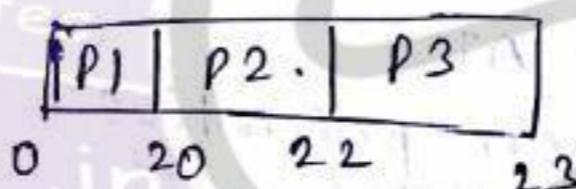


P NO.	AT	BT	CT	TAT = CT - AT	WT = TAT - BT	RT = FR - AT
1	0	2	2	2	0.	0
2	3	1	4	1	0.	0
3	5	6	11	6	0	0

Here  $WT=0$  means CPU is waiting for process whereas in previous eg. there was waiting time means processes were waiting for CPU.

- convey effect: all other process wait for one big process to get off the CPU. This effect results in lower utilization  
 eg { If arrival time is not given all process enter the system at time = 0 }

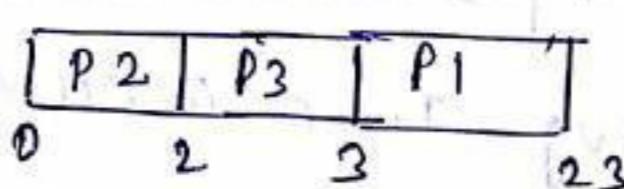
P NO.	AT	BT
1	0	20
2	1	12
3	1	15



P NO.	AT	BT	CT	TAT = CT - AT	WT = TAT - BT	RT = FR - AT
1	0	20	20	20	0.	0.
2	1	2	22	21	19	19
3	1	15	23	22	21	21

Avg. TAT =  $\frac{63}{3}$  Avg. WT =  $\frac{40}{3}$

P NO.	AT	BT
1	1	20
2	0	2
3	0	1



P NO.	AT	BT	CT	TAT = CT - AT	WT = TAT - BT	RT = FR - AT
1	1	20	23	22	2.	.
2	0	2	2	2	0.	.
3	0	1	3	3	2	.

$$\text{Avg TAT} = \frac{27}{3}$$

$$WT = \frac{4}{3}$$

whenever you are executing process you should not execute big process first so as to keep waiting time less.

convey effect : increase in time - - -

### → ii) shortest Job First (SJF)

- The process having smallest execution time will be assigned to the processor next.
- If 2 processes having same execution p time then you apply first come first serve to break the tie  
Here: This algo works according to burst time(less), mode: non-preemptive.

P.NO.	AT	BT
1	1	7
2	2	5
3	3	1
4	4	2
5	5	8

use SJF algo to find

avg. waiting time , - - -

0 | P1 | P3 | P4 | P2 | P5 |

1 | 8 | 9 | 11 | 16 | 24 |

TAT = CT - AT

WT = TAT - BT

RT = FR - AT

P.NO.	AT	BT	CT	TAT	WT	RT
1	1	7	8	7	0	0
2	2	5	16	14	9	9
3	3	1	9	6	5	5
4	4	2	11	7	5	5
5	5	8	24	19	11	11

$$\text{Avg. WT} = \frac{30}{5} = 6 \quad \text{Avg. TAT} = \frac{53}{5} = 10.6$$

$$\text{Avg. RT} = 10.6$$

### shortest Job first with prediction of burst/execution time

Advantage of SJF.

- maximum throughput
- min. avg. waiting time and turn around time

Disadvantage of SJF

- It leads to starvation (long process waits due to many small ones)
- It is not implementable because burst time of the process cannot be known in advance in real situations

The solution is SJF with prediction / predicted burst time .

## Prediction Technique

↓  
Static

- process size
- process type

Dynamic.

→ simple average

→ exponential average

11/8/17

→ Process size (Byte)

$$\begin{aligned} P_{old} &= 200 \text{ KB} \Rightarrow 20 \text{ units} \\ P_{new} &= 201 \Rightarrow 20 \text{ units} \end{aligned} \quad \left. \begin{array}{l} \text{not advisable to use.} \\ \end{array} \right\}$$

→ Process type

↓  
Kernel process      user process

e.g.: scheduler  
dispatcher

(3-5) units

Interactive process  
ex: games  
(5-8)

foreground process  
(10-15) units  
of time

Background process  
(15-20) units

Dynamic

• simple average

let  $n$  processes starting from  $(p_1, \dots, p_n)$

Let  $t_i$  be the actual burst time

$\tau_i$  denotes predicted burst time of  $i$  process.

$$\text{so. } \boxed{\tau_{i+1} = \frac{1}{n} \sum_{i=1}^n t_i}$$

• exponential average:

$$\tau_{n+1} = \alpha t_n + (1-\alpha) \tau_n \quad \text{--- (1)}, \quad 0 \leq \alpha \leq 1$$

predicted time of  
( $n+1$ )th process

$\alpha$ : const : smoothing factor.

$$\tau_n = \alpha t_{n-1} + (1-\alpha) \tau_{n-1} \quad \text{--- (2)}$$

so, input  $\tau_0$  in (1) contributing  $\alpha t_1$  at 2. in (2).

$$T_{n+1} = \alpha t_n + (1-\alpha)[\alpha t_{n-1} + (1-\alpha)T_{n-1}]$$

$$T_{n+1} = \alpha t_n + (1-\alpha)\alpha t_{n-1} + (1-\alpha)^2 T_{n-1}$$

$$T_{n+1} = \alpha t_n + (1-\alpha)\alpha t_{n-1} + \dots + (1-\alpha)^{n-1}\alpha t_1 + (1-\alpha)^n T_1$$

ex  $\alpha = 0.5$   $T_1 = 10$  Actual BT ( $P_1, P_2, P_3, P_4$ ) = (4, 8, 6, 7)

$$T_5 = ?$$

$$T_2 = \alpha t_1 + (1-\alpha)T_1$$

$$\geq 0.5 \times 4 + 0.5 * 10$$

$$= 7$$

$$T_3 = \alpha t_2 + (1-\alpha)T_2$$

$$= 0.5 \times 8 + 0.5 \times 7 = 7.5$$

$$T_4 = \alpha t_3 + (1-\alpha)T_3$$

$$= 0.5 \times 6 + 0.5 \times 7.5$$

$$= 3 + 3.75$$

$$= 6.75$$

$$T_5 = 6.75$$

→ shortest remaining time first scheduling time algorithm

- In this algorithm CPU scheduler always selects the process that has the shortest remaining processing time.
- when a new process enters the ready queue, the short term scheduler (CPU scheduler) compares the remaining time of executing process with the newly entered process.
- If the new process has least execution time, the scheduler selects that process for execution. Otherwise continues with the old process execution.

P. NO.	AT	BT	CT	TAT = CT - AT	WT = TAT - BT
1	0	7.50	19	19 - 0 = 19	19 - 7.50 = 11.50
2	1.50	5.40	13	13 - 1.50 = 11.50	11.50 - 5.40 = 6.10
3	2.50	2.10	6	6 - 2.50 = 3.50	3.50 - 2.10 = 1.40
4	3	1.0	4	4 - 3 = 1	1 - 1 = 0
5	4	2	9	9 - 4 = 5	5 - 2 = 3
6	5	1.0	7	7 - 5 = 2	2 - 1.0 = 1

schedule

P1	P2	P3	P4	P3	P3	P6	P5	P2	PL
0	1	2	3	4	5	6	7	9	13 . 19.

when tie apply fcfs

criteria : shortest remaining time

mode : preemptive

for completion time move right to left

#### → Priority scheduling algorithm

- A priority is associated with each process and the CPU is allocated to the process with the highest priority
- equal priority processes are scheduled in fcfs order (first come first serve).
- priority can be defined internally or externally.  
nothing given: smallest no. → more priority

218117

#### II. Preemptive priority scheduling algorithm.

P. NO.	Priority	AT	BT	CT	TAT = CT - AT	WT = TAT - BT	RT = FR - AT
1	2 (low)	0	4.3	25	25 - 0 = 25	25 - 4.3 = 20.7	0
2	4	1	2.1	22	22 - 1 = 21	21 - 2.1 = 18.9	0
3	6	2	3.2	21	21 - 2 = 19	19 - 3.2 = 15.8	0
4	10	3	5.3	12	12 - 3 = 9	9 - 5.3 = 3.7	0
5	8	4	1.5	19	19 - 4 = 15	15 - 1.5 = 13.5	0
6	12 (high)	5	4	9	9 - 5 = 4	4 - 4 = 0	0
7	9	6	6	18	18 - 6 = 12	12 - 6 = 6	0

P1	P2	P3	P4	P6.	P4	P7	P5	P3
0	1	2	3	9.5	9	12	18	19

Avg. WT = ? Avg. TAT = ? Avg. RT = ?

$$\text{Throughput} = \frac{\text{Total job}}{\text{Total time}} = \frac{7}{25}$$

### → Non-preemptive priority scheduling :-

P.NO.	Priority 2 (low) 4 6 10 8 12 (high) 9	AT 0 1 2 3 4 3 6	BT 4 2 3 5 1 4 6	CT 4 25 23 9 20 13 19	TAT 4 24 21 6 16 8 13	WT 0 22 18 1 15 4 7	RT 0 22 18 1 15 4 7
1	2 (low)	0	4	4	4	0	0
2	4	1	2	25	24	22	22
3	6	2	3	23	21	18	18
4	10	3	5	9	6	1	1
5	8	4	1	20	16	15	15
6	12 (high)	3	4	13	8	4	4
7	9	6	6	19	13	7	7

schedule

P1	P4	P6	P7	P5	P3	P2
0	4	9	13	19	20	23 25

Avg. WT = ?

" TAT = ?

" RT = ?

Throughput = ?

Note\* In non-preemptive scheduling waiting time is equal to response time

- The major problem with priority scheduling algo is indefinite blocking or starvation.
- Starvation means only higher priority processes are executing whereas lower priority processes are waiting for the CPU for long period of time.
- A solution to this problem: aging
- aging is a technique to gradually increase the priority of the processes that are waiting in the system for longer period of time.

### # ROUND ROBIN ALGORITHM (RR)

ROUND ROBIN ALGORITHM (RR) is designed for time-sharing systems.

- The RR algorithm is similar to FCFS algo, but the preemption is added to switch b/w the processes.
- A small time unit called time quantum or time slice is defined by OS.
- In this algorithm the CPU switches b/w the processes when the time quantum expires.

	P.NO.	AT	BT	CT	TAT	WT	RT
TQ = 2	1	0	4 2.	8	8.	4.	
=	2	1	5 3 1	18	17	12	preemptive
	3	2	2	10	6.	4	
	4	3	1	9	6.	5	
	5	4	6 4 2.	21	17	11	
	6	6	3 4.	19	13.	10.	

Queue: P1 P2 P3 P1 P4 P5 P2 P6 P5 P2 P6 P5

Schedule:

	P1	P2	P3	P1	P4	P5	P2	P6	P5	P2	P6	P5.
	0	2	4	6.	8	9	11	13	15	17	18	20

3 | 8 | 17

Q	P.NO.	AT	BT	CT	TAT = CT - AT	WT = TAT - BT	RT = FR - AT
TQ=3	1	5	5 2	10	32	27	10.
=	2	4	6 3.	27	23	17	5.
	3	3	7 4 1	33	30.	23.	3.
	4	1	9 8 3.	30	29	20.	0.
	5	2	2	6.	4	2.	2.
	6	6.	3	21	15	12	12.

Queue: P4 P5 P3 P2 P4 P1 P6 P3 P2 P4 P1 P3 |

	P4	P5	P3	P2	P4	P1	P6	P3	P2	P4	P1	P3
	0	1	4.	6.	9	12	15	18	21	24	27	30

time quantum increase  $\rightarrow$  fcfs

time quantum decrease  $\rightarrow$  in ms  $\rightarrow$  then process just jumps from one to another.

- The performance of round robin algorithm depends on the size of time quantum.
- If the time quantum is very large, the RR algorithm behaves same as fcfs algorithm.
- If time quantum is very small (1ms), the RR algorithm is called processor sharing.

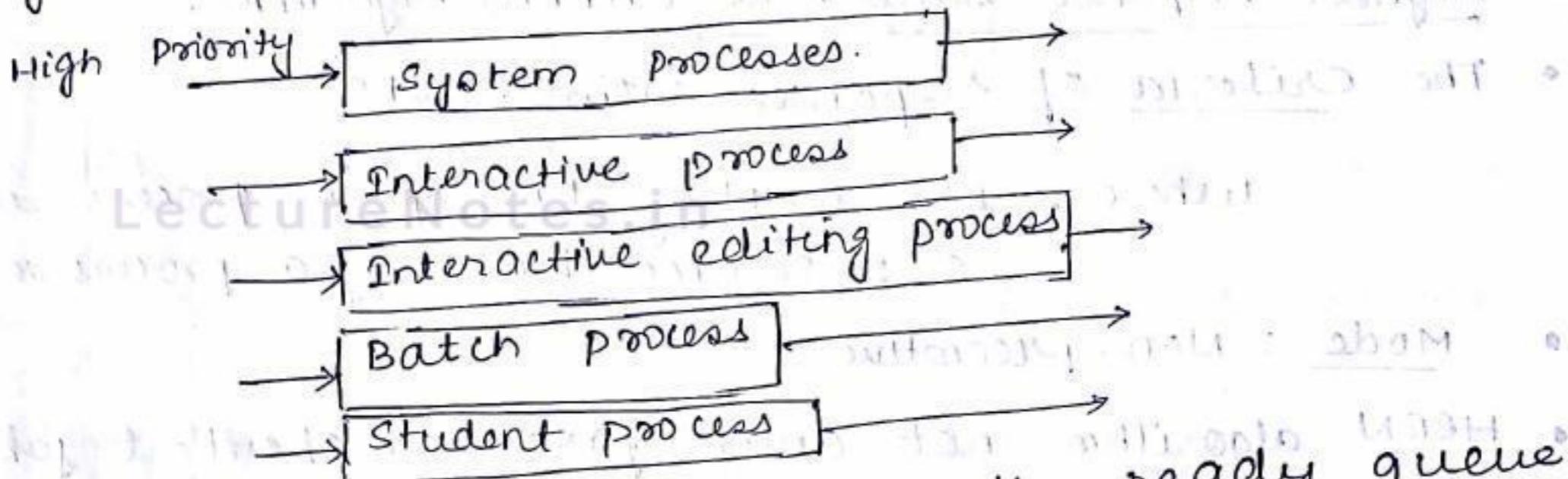
Process time Time Quantum context switch.

avg = 10.

- smaller the time quantum, context switch will be more.

## # Multilevel Queue scheduling Algorithm:

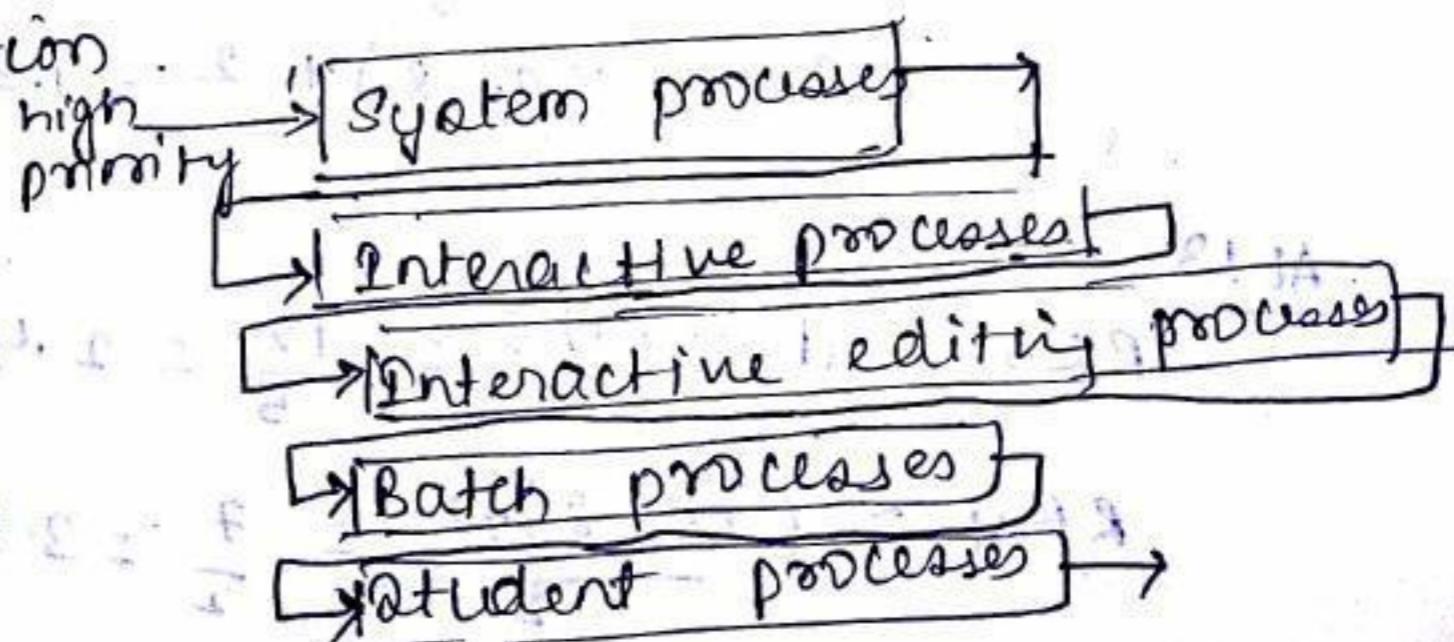
- Processes are classified into different groups such as - foreground - background etc.
- These processes have different response time so they need different scheduling schemes.



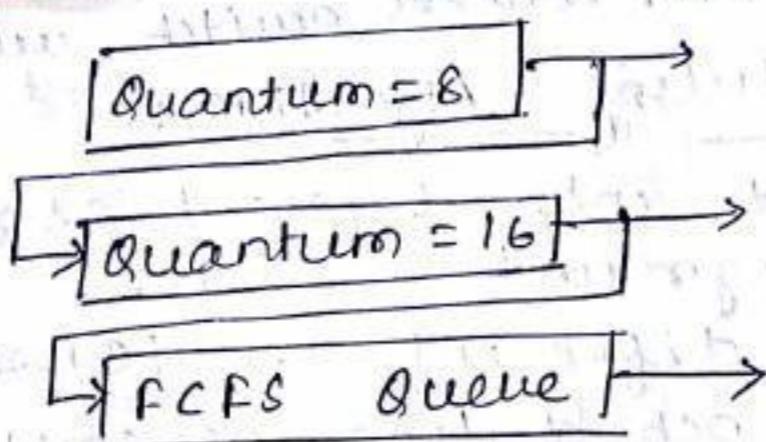
- In multilevel queue scheduling, the ready queue is divided into several separate queues.
- The processes are permanently assigned to 1 queue based on the priority of the process.
- Each queue has its own scheduling algorithms.

## # Multilevel feedback queue-scheduling algorithm:

- In multilevel queue scheduling, the processes are permanently assigned to a queue on entry to the system.
- Processes do not move between the Queue (Multilevel queue).
- In multilevel feedback Queue scheduling allows a process to move between queue.
- If a process uses too much CPU time, it will move to the lower priority queue.
- It prevents starvation.



eg:-



## # Highest Response ratio next (HRRN) Algorithm:

- The criteria of response ratio =  $\frac{W+S}{S}$   
where,  $W$  = waiting time for a process so far  
 $S$  = service time of a process or BT
- Mode : Non-preemptive
- HRRN algorithm not only favours shortest job but also limits the waiting time of the longer jobs.

4/8/17

P.NO.	AT	BT	CT	TAT = CT - AT	WT = TAT - BT	RT = PR - AT
0	0	3	3	3	0	0
1	2	6	9	7	1	1
2	4	4	13	9	5	5
3	6	5	15	9	9	9
4	8	2	20	12	10	5

schedule		P0	P1	P2	P3	P4
0	3	9	13	18	20	15

At 9

$$RR_2 = \frac{(9-4)+4}{4} = 2.25$$

$$RR_3 = \frac{(9-6)+5}{5} = 1.6$$

$$RR_4 = \frac{(9-8)+2}{2} = 1.5$$

Once whose response is more will go for execution.  
i.e P2.

At 13

$$RR_3 = \frac{(13-6)+5}{5} = \frac{12}{5} = 2.4$$

$$RR_4 = \frac{(13-8)+2}{2} = \frac{7}{2} = 3.5$$

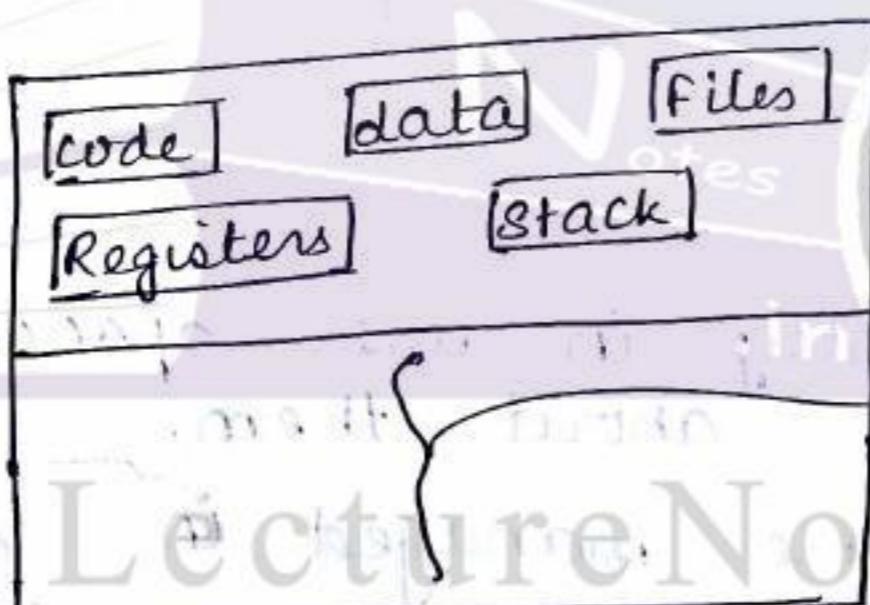
P4 will be executed.

## SJF schedule

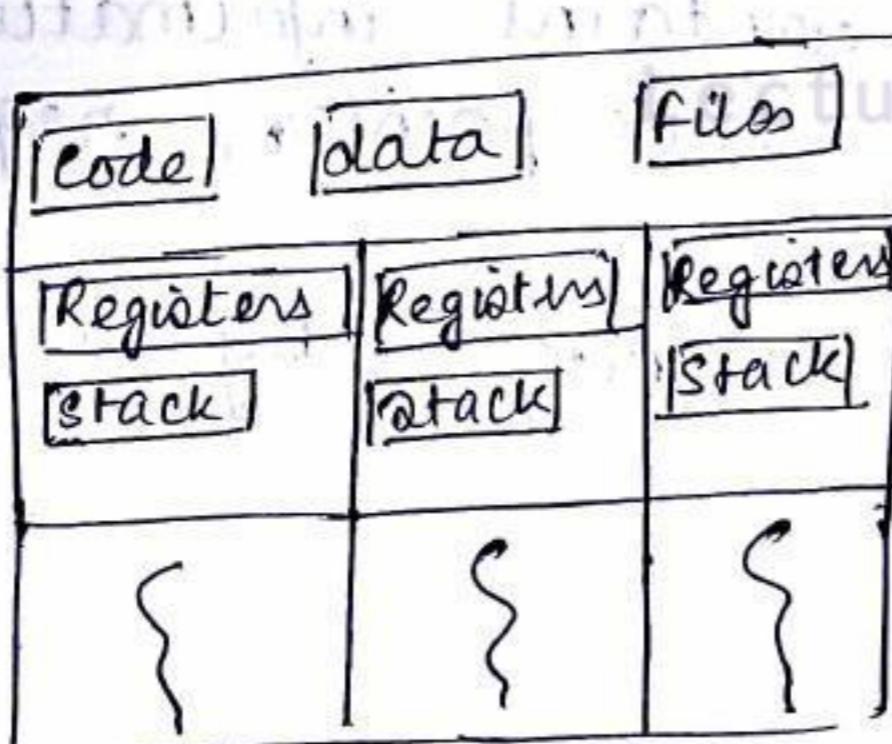
	P0	P1	P2	P3	
AT	0	3	9	11	15
PNO.	0	1	2	3	4
BT	10	10	10	10	10
	0	2	4	6	8

## → THREAD

- Process is divided into number of smaller task called Thread.
- A Thread sometimes called light weight process (LWP).
- It is a basic unit of CPU utilization.
- Thread has Thread id, program counter (PC), Register set and stack etc...
- NO. of threads within a process execute at a time called multithreading.



single Thread Process



Multithreading process.

- i) process takes more time to create, execute and complete its operation
- ii) It takes more time to switch b/w processes
- iii) Process communication is more difficult
- iv) Processes are loosely coupled
- v) It requires more no. of resources to execute
- vi) Processes are not suitable for parallel activities

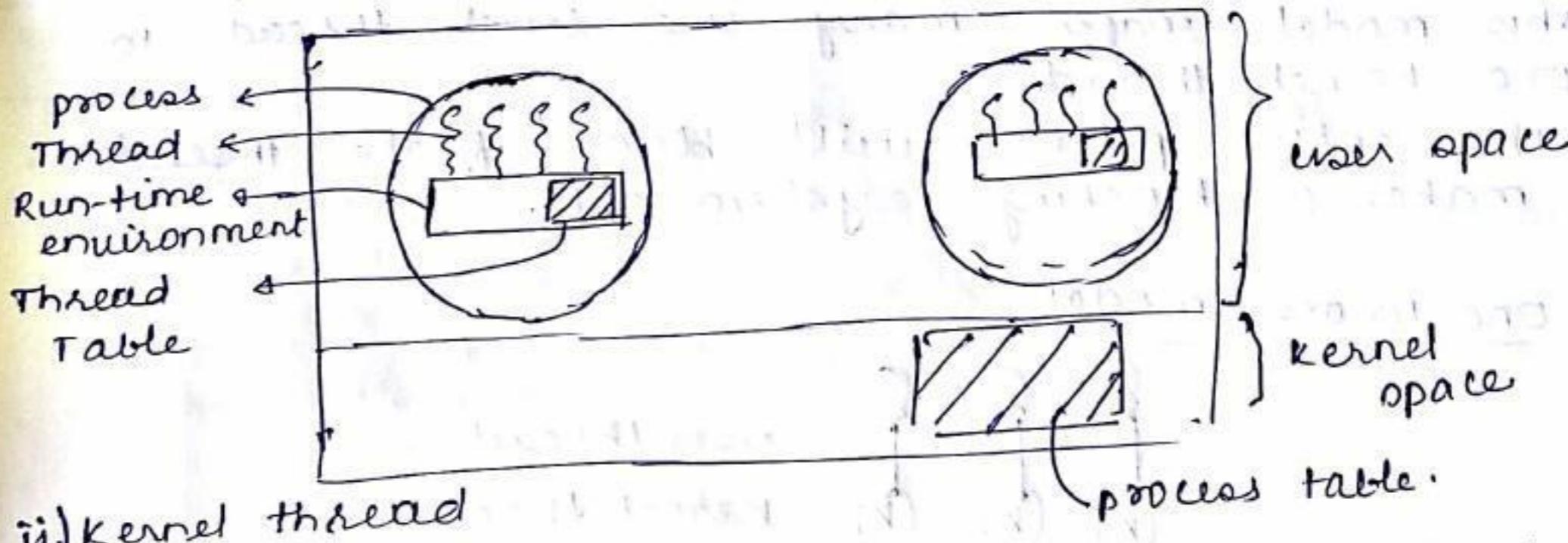
• There are 2 types of thread.

- i) user thread
- ii) kernel thread

#### → i) user thread.

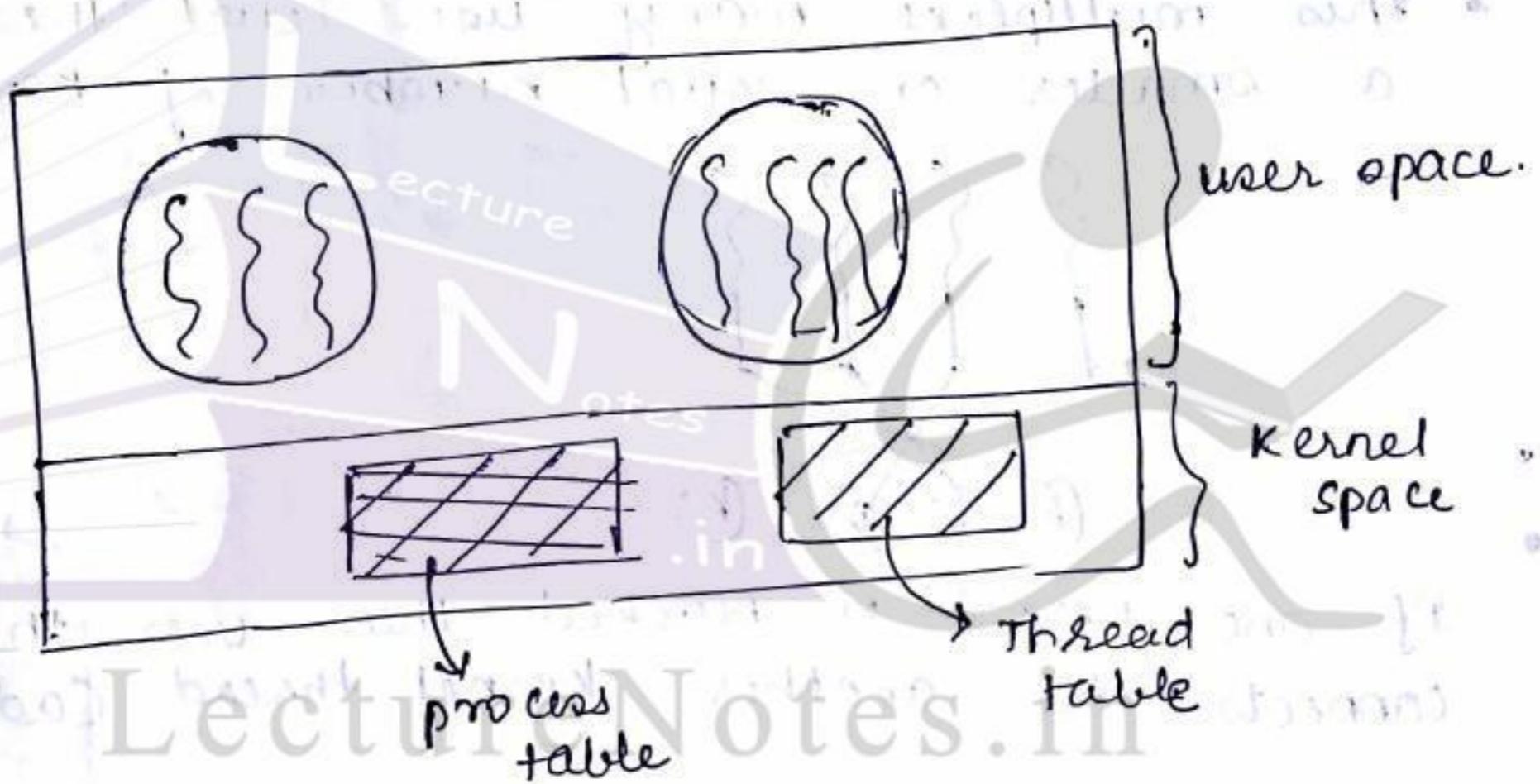
- are loaded entirely in user space, kernel knows nothing about them.
- when threads are managed in user space the process needs its own private thread table.
- The thread table contains information about program counter, stack pointers, registers and state (of thread) etc.
- Thread table is managed by a run-time environment

- ii) Thread takes less time to create, execute and complete its operation.
- ii) It takes less time to switch b/w threads.
  - iii) Thread communication is easier.
  - iv) Threads are tightly coupled
  - v) It requires fewer no. of resources to execute, i.e. it is light weight
  - vi) Threads are suitable for parallel activities.



### ii) Kernel thread

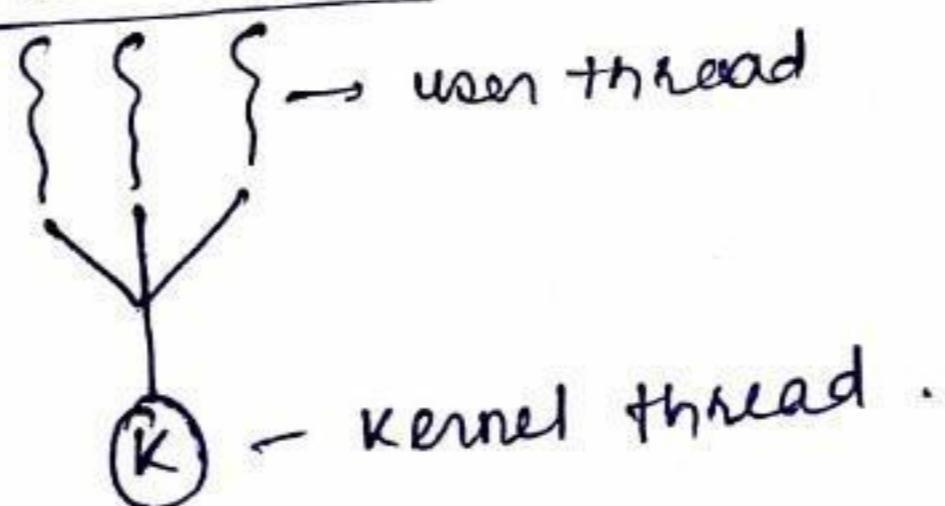
- In kernel thread, kernel does total work of thread management
- There is no thread table in each process
- Kernel has a thread table that keeps track of all the threads in the system.



### → Multi-threading model

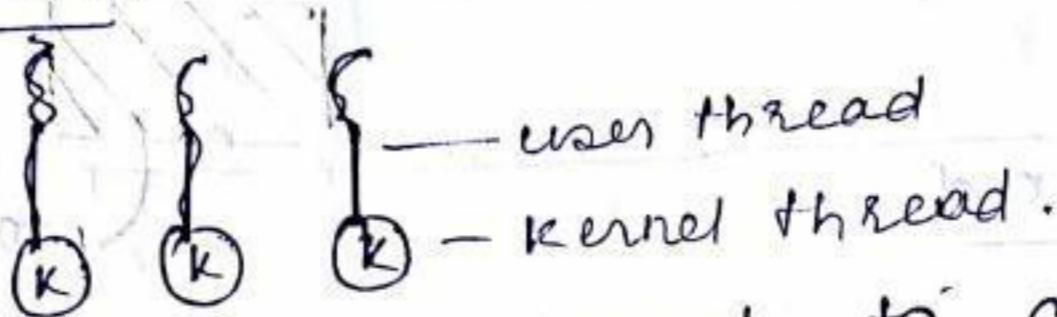
- many system provide support for both user and kernel thread.
- Three types of multithreading implementation

#### i.) many to one model



- This model maps many user level threads to one kernel thread.
- The entire process makes a blocking system call. will block if the thread

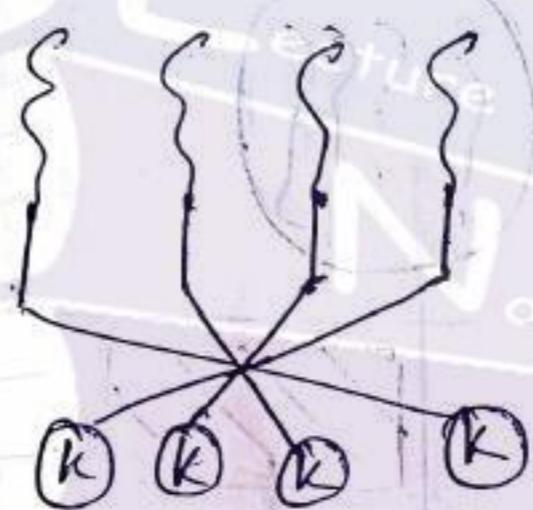
### ii) one to one model.



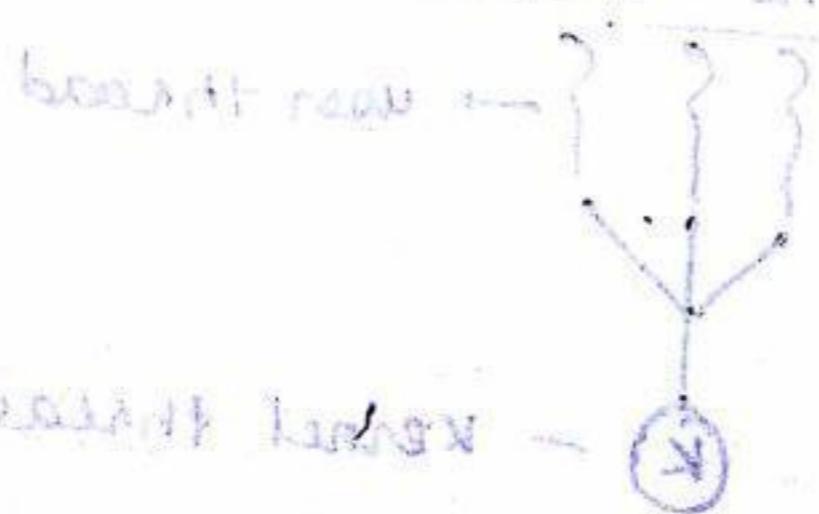
- In this maps each user thread to a kernel thread.
- one blocked yet others work. (advantage)

### iii) many to many model

- This multiplexes many user level threads to a smaller or equal number of kernel threads.



- If one kernel is blocked then user thread connects to another kernel thread. (advantage)



8/8/17

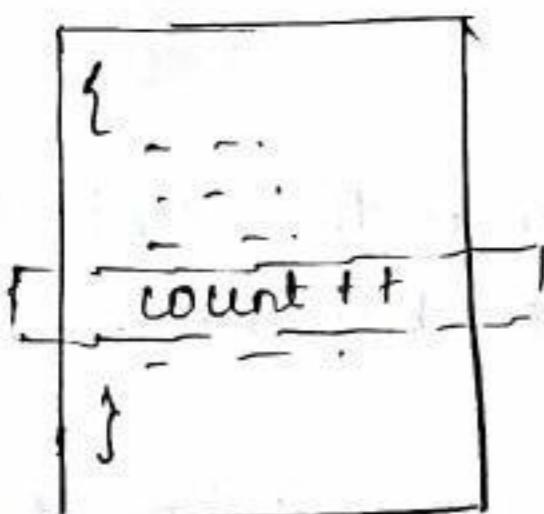
#. PROCESS SYNCHRONIZATION

- when two or more co-operating processes execute a shared variable, a shared record or a shared file simultaneously, then there may be chances of data inconsistency. This leads to a problem called race-condition.
- To avoid this problem process synchronization is needed.
- There are two kinds to synchronization:
  - Data-access synchronization: It ensures that shared data does not get loose consistency, when they are updated by several co-operating processes.
  - control-synchronization: It ensures that co-operating processes perform their actions in desired order.
- Race-condition
  - A situation where several processes access and manipulate the same data concurrently and the outcome of the execution depends on the particular order in which the access takes place.
  - A race condition occurs when multiple processes read or write data items so that the final result depends on the order of the execution of instructions in the multiple processes.

global shared variable

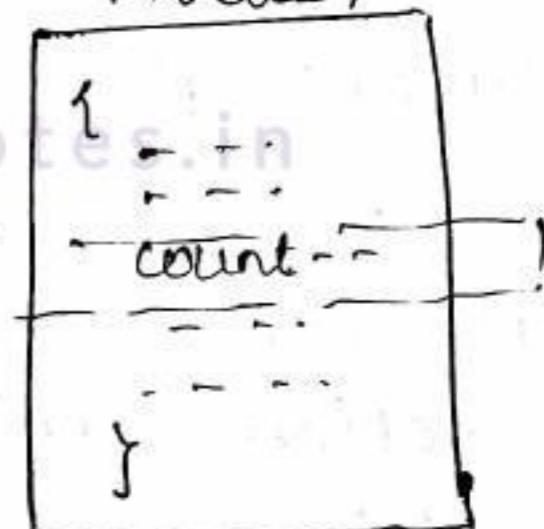
count = 5

Process 0



critical section

Process 1



expected O/P: 5 ✓

if they execute simultaneously 4, 6. X

$R1 \leftarrow count$	5
$R1 \leftarrow R1 + 1$	6
$count \leftarrow R1$	6
$R2 \leftarrow count$	6
$R2 \leftarrow R2 - 1$	5
$count \leftarrow R2$	5

context switch

O/P : 5  
✓

$R1 \leftarrow count$	5
$R2 \leftarrow count$	5
$R2 \leftarrow R2 - 1$	4
$count \leftarrow R2$	4
$R1 \leftarrow R1 + 1$	6
$count \leftarrow R1$	6

context switch

O/P : 6

$R2 \leftarrow count$	5
$R1 \leftarrow count$	5
$R1 \leftarrow R1 + 1$	6
$count \leftarrow R1$	6
$R2 \leftarrow R2 - 1$	4
$count \leftarrow R2$	4

context switch

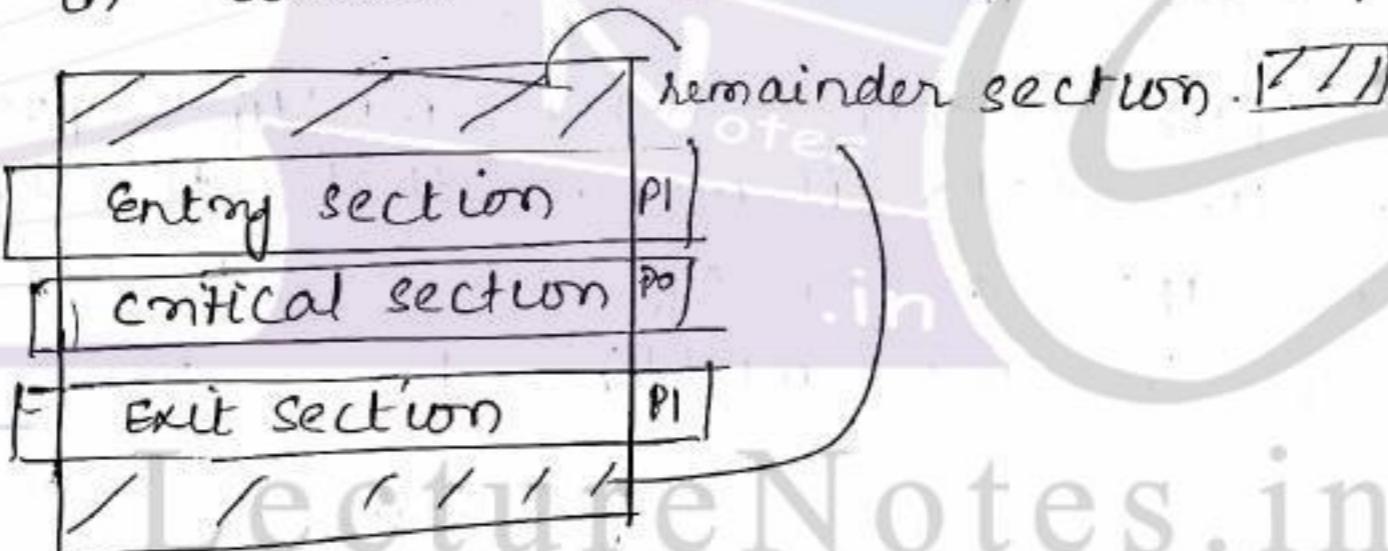
O/P : 4

X

Race condition

{ Due to two co-operating process which have a shared value }

- Race condition can be prevented by synchronization.
- Synchronization ensures that only one process at a time manipulate the critical data.
- The outcome depends on the
- A situation where the several processes access and manipulate the critical data is called critical section.
- The outcome of the critical section depends on the order in which the access takes place.



- critical section
  - only one process <sup>can</sup> execute here at a time
- entry section
  - It allows only <sup>one</sup> process to enter into the critical section and rest processes are blocked.
- exit-section
  - Once a process completes its critical section, it enters into the exit section.
  - In exit section, it allows another process to enter into the critical section.

- remainder section
  - The rest of the codes are called remainder section
- A solution to the critical section problem must satisfy
  - i) mutual exclusion
  - ii) progress
  - iii) bounded-waiting
- i.) mutual exclusion
  - It allows only one process to execute in critical section at a given instant of time.
- ii.) progress
  - a) when no process is in critical section, any process that requests entry into the critical section must be permitted within finite amount of time.
  - b) A decision who will enter into the critical section must be taken by considering only those processes who wants to enter into the critical section and these decision must be taken in finite time.
- iii.) bounded-waiting
  - there is an upper bound on the number of times a process enters the critical section, while another process is waiting.

9/8/17

If the above sections are satisfied then it is said to be synchronized. (process)

#### → LOCKS AND UNLOCKS

- All critical section problems use locking and unlocking. Suppose I have a shared variable.

shared variable  
 int count = 5;  
 lock\_t L;

Process 0

<code>i =</code>	
<code>lock(L);</code>	entry-section
<code>c++;</code>	cs.
<code>unlock(L);</code>	
exit section	

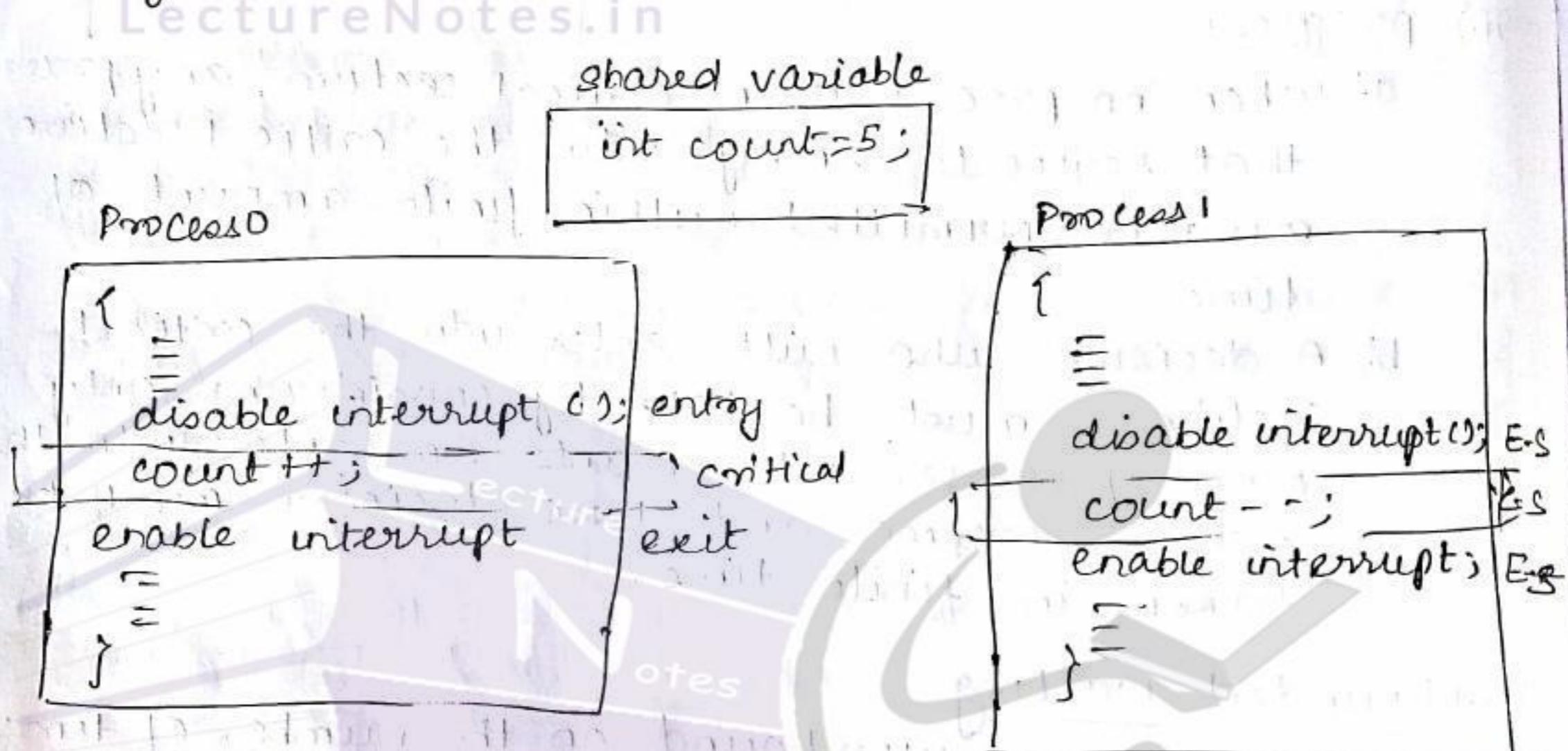
Process 1

<code>i =</code>	
<code>lock(L);</code>	
<code>count--;</code>	
<code>unlock(L);</code>	

- lock(L): acquire lock(L) exclusively  
only one process with L can access the critical section
- unlock(L): release exclusive access to lock(L)  
permitting other processes to access the critical section.

Note\*: Locking and unlocking should be designed in such a way that the three requirements of critical section should be satisfied.

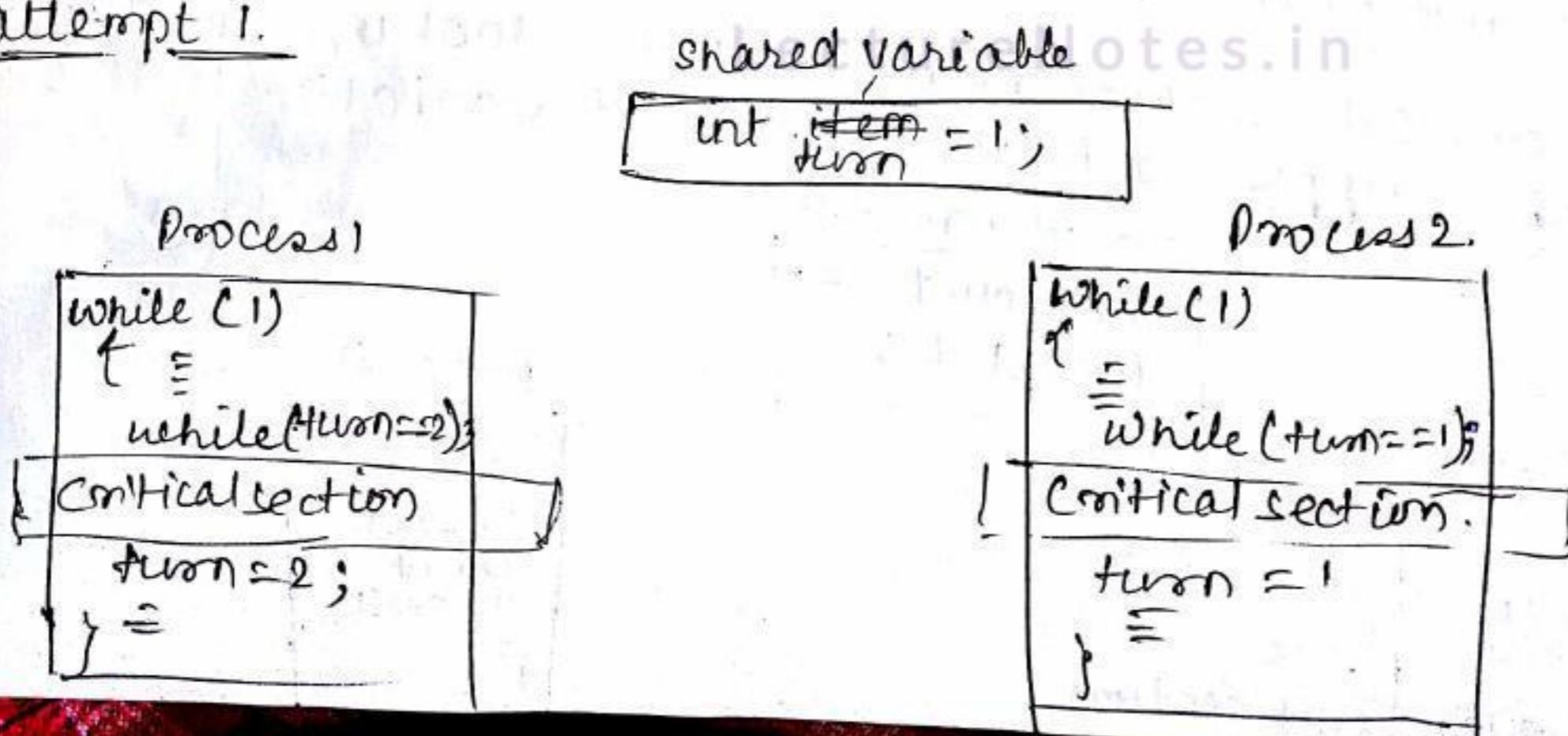
→ using interrupt



- i.) It is simple: when interrupts are disabled, context switch will not happen
- ii.) It requires privilege: user processes generally cannot disable the interrupts.

→ software solutions {for this critical section}

attempt 1.



$\left\{ \begin{array}{l} \text{while(cond)}; \\ \quad \left\{ \begin{array}{l} \text{statements} \\ \quad \left\{ \begin{array}{l} \text{if cond true then nothing executes} \\ \text{if " false then statements execute} \end{array} \right. \end{array} \right. \end{array} \right.$

$P_1 \rightarrow P_2 \rightarrow P_1 \rightarrow P_2 \dots$  mutual exclusion as at a time only one process executes at a time -

$P_2$  can't proceed till  $P_1$  executes so progress is not satisfied so this fails. (if  $P_2$  doesn't want to execute)

- It achieves mutual exclusion
- Here processes are executed alternatively in critical section  $P_1 \rightarrow P_2 \rightarrow P_1 \rightarrow P_2 \dots$
- Busy waiting: waste of power and time  
(frequently  $P_1$  checks value of turn, and while statement - )

### Problems for attempt 1

- It had a common turn flag that was modified by both the processes.
- This required processes to alternate
- Due to this reason progress also not satisfied.

### Solution:

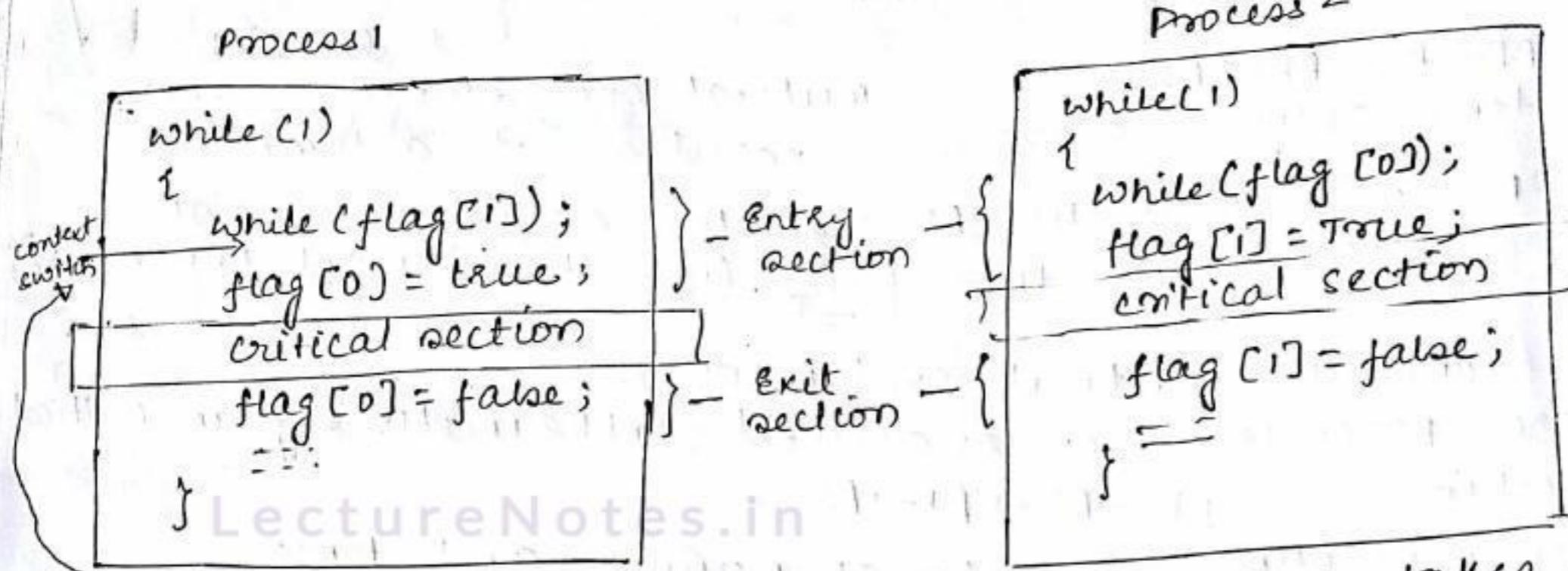
- It should have two flags: one for each process

### → Software solution (attempt 2)

10/8/17

- The problem with attempt 1 is that it does not retain sufficient information about the state of the process. It remembers only which process is allowed to enter into the critical section.
  - To avoid this problem we can replace 'turn' variable with the following array.
- Boolean flag[2];  
flag[0] = flag[1] = false

shared variable.  
 $\text{flag}[0] = \text{flag}[1] = \text{false}$



{ mutual exclusion fails when context switch takes place over here ... and both process enter into critical section }

- It does not guarantee the mutual exclusion
- It means process 1 and process 2 are in critical section at same time.

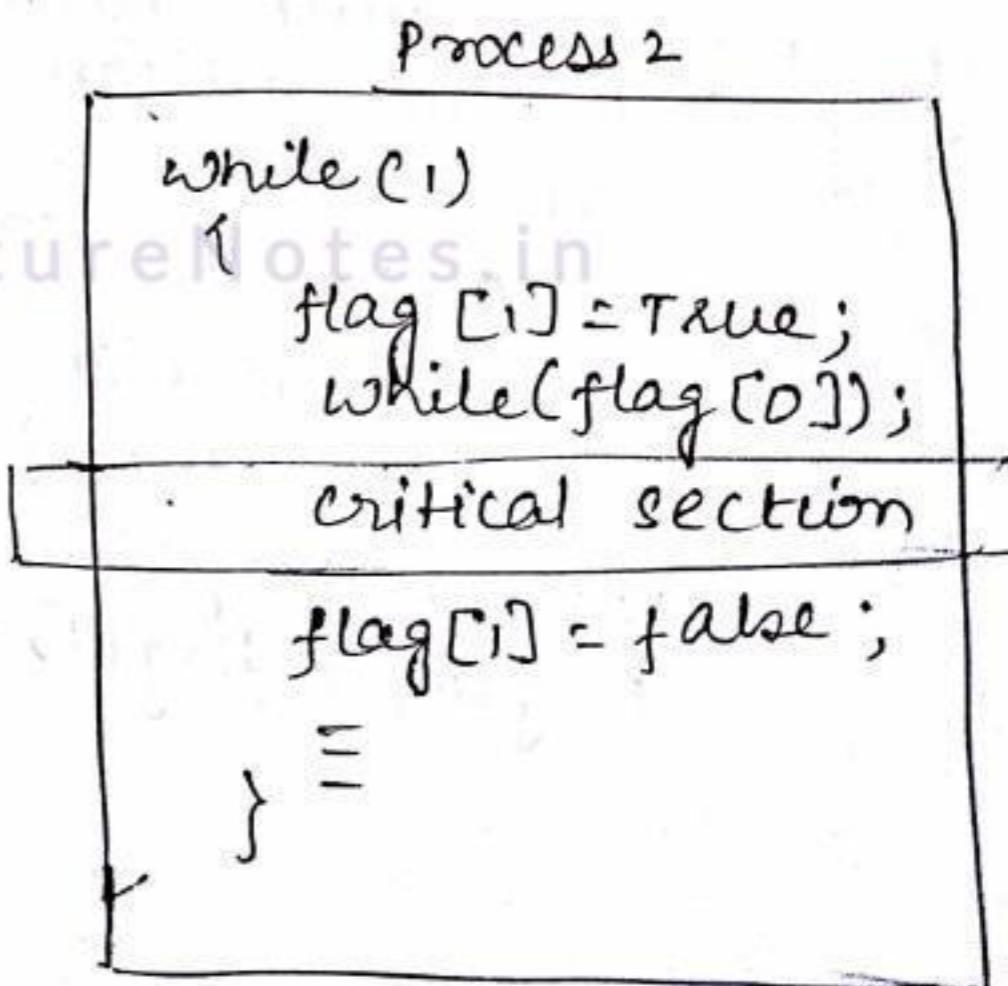
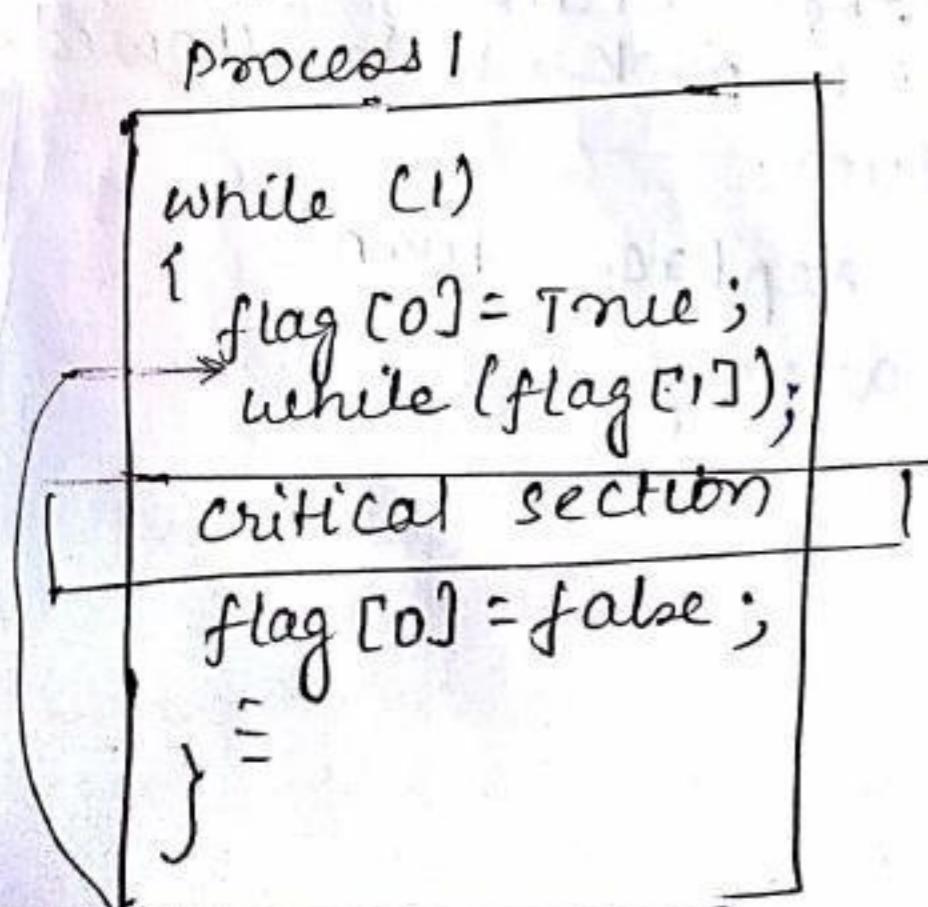
#### - problem with attempt 2

The flags are set after we break from the while loop

#### → software solution (attempt 3)

shared variable

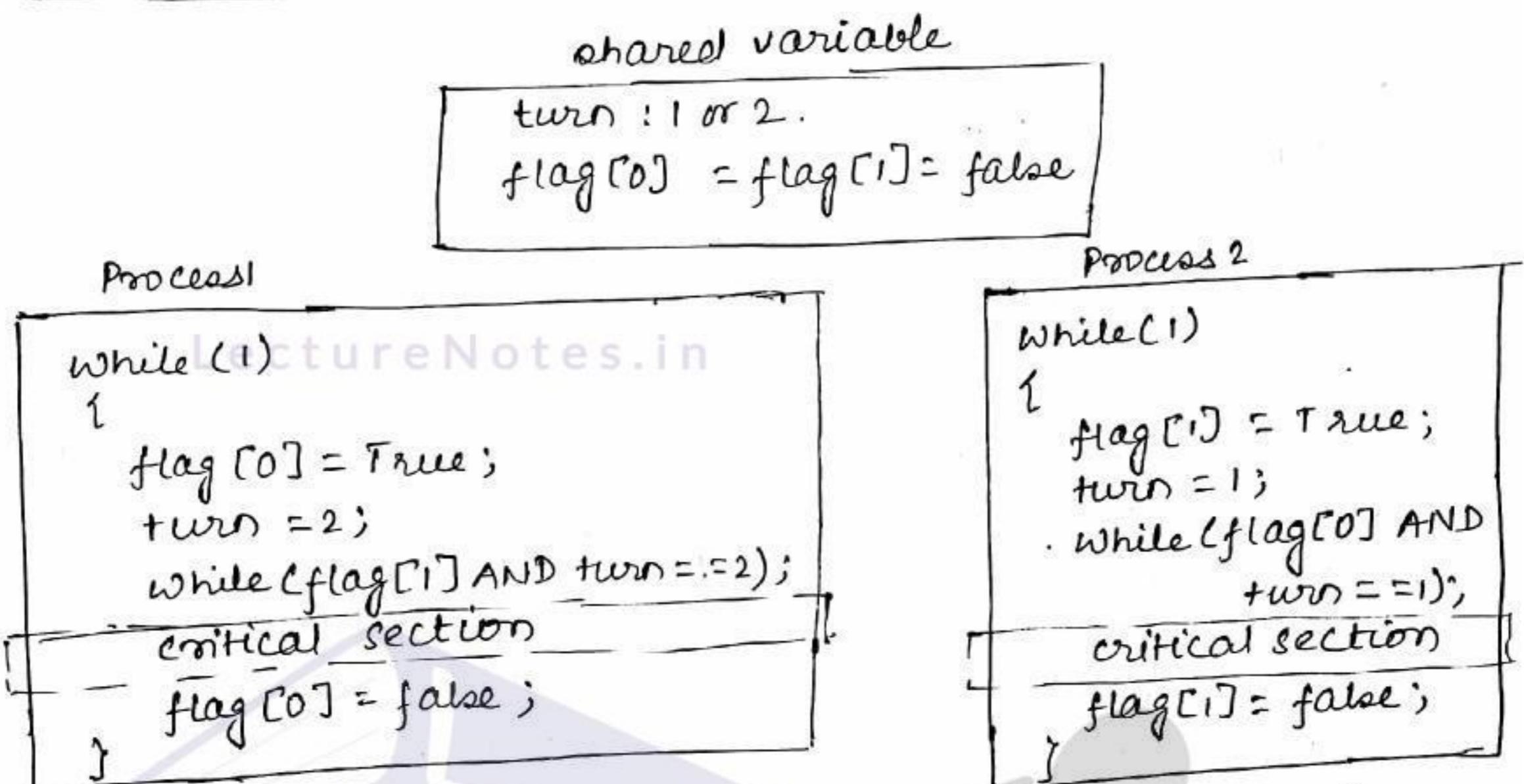
$\text{flag}[0] = \text{flag}[1] = \text{false}$



C.S. deadlock

- It achieves mutual exclusion
- It does not achieve progress as these two processes are deadlocked.

→ Peterson's solution



{  
 $\text{flag}[0] = \text{true}$ : means process 1 is interested to enter into critical section.  
 $\text{turn} = 2$ : means turn of process 2.  
 The deadlock is broken over here due to 'turn' variable because turn can be either 1 or 2 therefore the tie is broken. Only one process enters into the critical section at a time so it solves critical section problem for two processes.  
 It also satisfies all the three required cond's.

11/8/17

→ SEMAPHORE

It is a special kind of integer variable which can be initialized and can be accessed only by two atomic operations such as `PC()` or `wait()` operation and `ve()` or `signal()` operation. atomic means no context switch in b/w n.

```

semaphore s ;
P { wait(s)
  {
    while (s <= 0);
    s --;
  }
}
V [signal(s)
  {
    s++;
  }
]

```

$s$  is a variable of semaphore type  
it can do  $P()$ ,  $V()$ ,  $wait()$  or  
 $signal()$  opn.

- semaphore mutex;

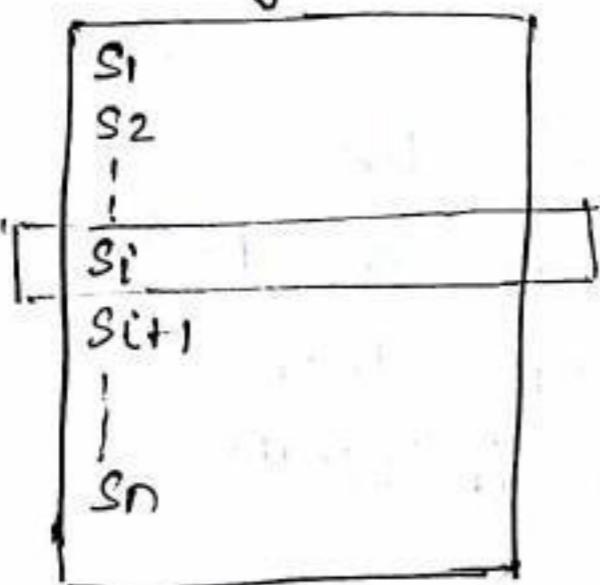
mutex = 1;

P<sub>i</sub> :  $P(mutex)$   
critical section  
 $V(mutex)$   
Remainder section

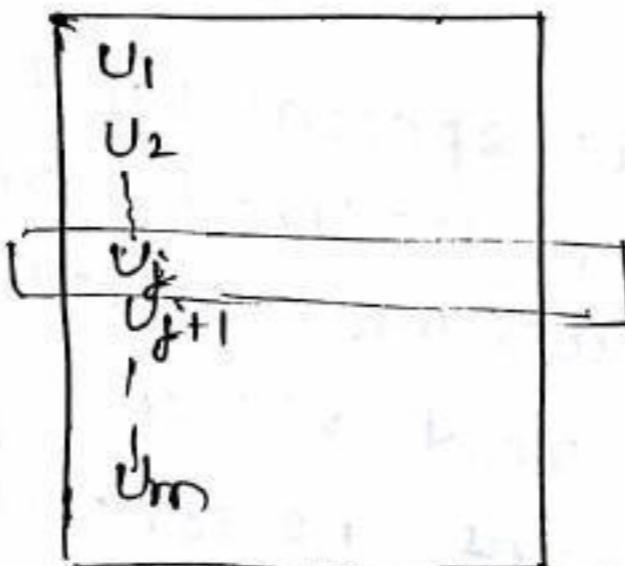
P<sub>j</sub> :  $P(mutex)$   
critical section  
 $V(mutex)$   
Remainder section

- using semaphore, we can implement mutual exclusion very easily. progress  $\rightarrow$  bounded waiting
- semaphore uses - process synchronization.
  - o There are two types of semaphore variable
    - i) counting semaphore : The value of counting semaphore can range over an unrestricted domain (can store any integer value).
    - ii) binary semaphore : The value of binary semaphore can range only between 0 and 1.

Program 1



Program 2

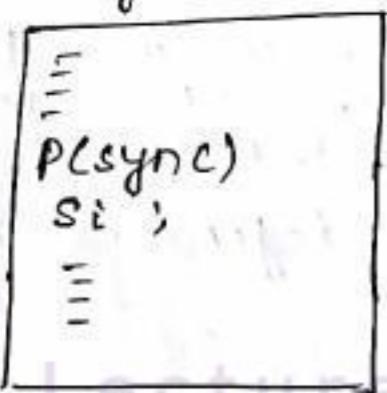


$U_j \rightarrow S_i$  so Restriction  
Uj must be executed before si for right output

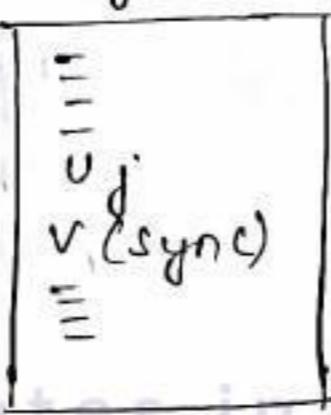
semaphore sync;

sync = 0;

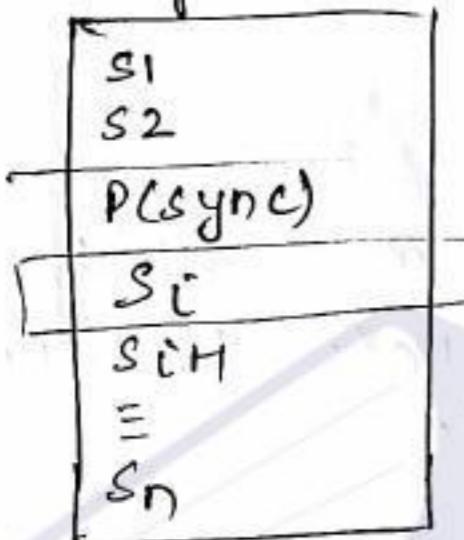
program - 1



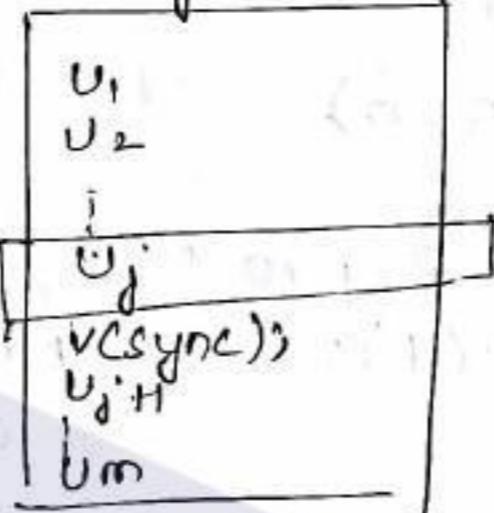
Prog - 2



Prog 1



Prog 2



- semaphore implementation (to save time & power of CPU due to busy waiting)

- Busy waiting : it is consuming CPU time without doing any useful work.
- This busy waiting can be eliminated by modifying the semaphore structure.

Now we define semaphore structure with two fields

i) integer                    ii) list

i.) integer : it contains the value of semaphore variable.

ii.) list : it contains no. of processes that are waiting.

C:-

typedef struct

```

    {
        int value;
        struct process *list;
    }Semaphore;
  
```

- wait() semaphore operation can be defined as :
 

```

wait(semaphore *s)
{
    s->value--;
    if (s->value < 0)
    {
        add this process to s->list;
        block();
    }
}
      
```

// Now process is not in running state but in block state.
- signal() semaphore operation can be defined as :
 

```

signal(semaphore *s)
{
    s->value++;
    if (s->value >= 0)
    {
        remove a process P from s->list;
        wakeup(P); // block state to running state.
    }
}
      
```

1618117

### → Synchronization hardware:

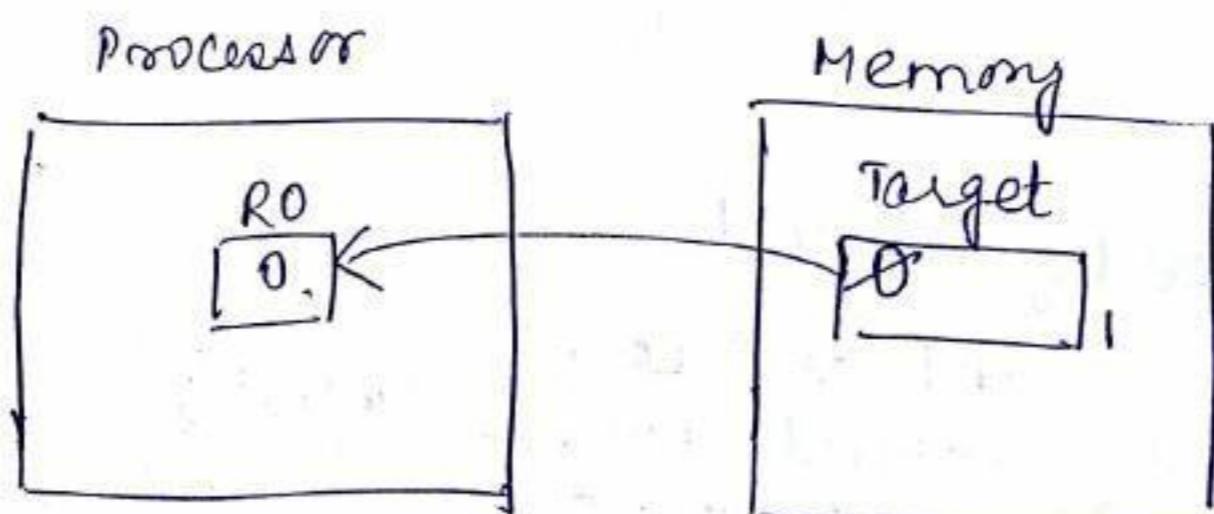
- Some hardware instructions can be used effectively in solving the critical section problem.

boolean test-and-set(boolean \*target)

```

{
    boolean rev = *target;
    *target = True;
    return rev;
}
      
```

# times	return
1	0
2	1
3	1



- first time it returns 0 rest time 1.
- Mutual exclusion implementation using test-and-set function.

Process 1	lock = 0
<pre> while (1) {     while (test-and-set (&amp;lock) == 1);     critical section;     lock = 0. } </pre>	achieves mutual exclusion but not others

Note\* Intel processor does not support test-and-set function.

atomic

```

int compare-and-swap (int *value, int expected,
                      int new-value).
{
    int temp = *value;
    if (*value == expected)
        *value = new-value;
    return temp;
}

```

memory

value
01.

1st value = 0      expected = 0  
                   temp = 0.  
                   new-value = 1.  
                   value = 1

2nd

- mutual exclusion implementation with compare-and-swap function

Process 1
<pre> while (1). {     while (compare-and-swap (&amp;lock, 0, 1) != 0);     critical section;     lock = 0; } </pre>

mutual exclusion satisfied  
not not satisfied

- Another algorithm using test-and-set function, instruction that satisfies all the three critical section requirements.

global variables:

boolean waiting[n] = false;
boolean lock = false

$P_i$ : while (1).

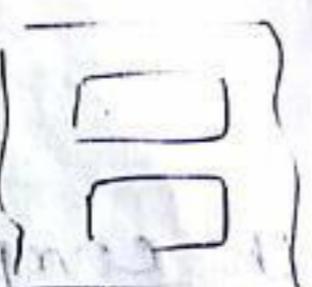
These two will execute when processes are interested

```
waiting[i] = true; //  $P_i$  is waiting to enter into the c.s.
key = true; // local variable
while(waiting[i] && key)
    key = test-and-set(&lock);
waiting[i] = false;
```

critical section:

```
j = i + 1 % n;
while((j != i) && !waiting[j]);
j = (j + 1) % n;
if(j == i)
    lock = false; // coming back to initial state(false)
else
    waiting[j] = false; //
```

$P_j$ :



allows next process to enter into critical section.

waiting[i] = false // initially

True false

key = false 0.

Now  $P_j$  wants to enter while we are in c.s of  $P_i$

waiting[j] = true key = true

It satisfies all three condns. for critical sections.

17/8/17

fork()

→ Mutex locks

Mutex  
Mutual exclusion.

- The hardware based sol' to all the critical section problem are generally inaccessible to application programmers.
- To avoid this problem, best designer build a software tools to solve critical section problem which is called mutex lock.
- Mutex locks is used to protect critical section and also prevents race condition.
- A process must acquire the lock before entering to the critical section and release the lock when it exits from the critical section.

while (1)

{

acquire lock

    critical section

release lock

}

acquire

{

    while (!available);  
    available = false;

release

{

    available = true;

}

due to busy waiting other process seem to be spinning so it is also called spin lock.

spinlock

- The main disadvantage of this implementation is busy waiting.
- while a process is in its critical section, any other process that tries to enter its critical section must loop continuously in the call to acquire().
- These types of mutex lock is called spin lock, because the process spins while waiting for the lock to become available.
- Deadlock and starvation due to semaphore.

- A deadlock problem will arise when two or more processes are waiting indefinitely for an event that can be caused only by one of the waiting process.

Semaphore S, Q;  
 $S=1; Q=1;$

P0.  
 wait(S);  
 wait(Q);  
 [-----]  
 critical section  
 signal(S);  
 signal(Q);

P1.  
 wait(Q);  
 wait(S);  
 [-----]  
 signal(Q);  
 signal(S);

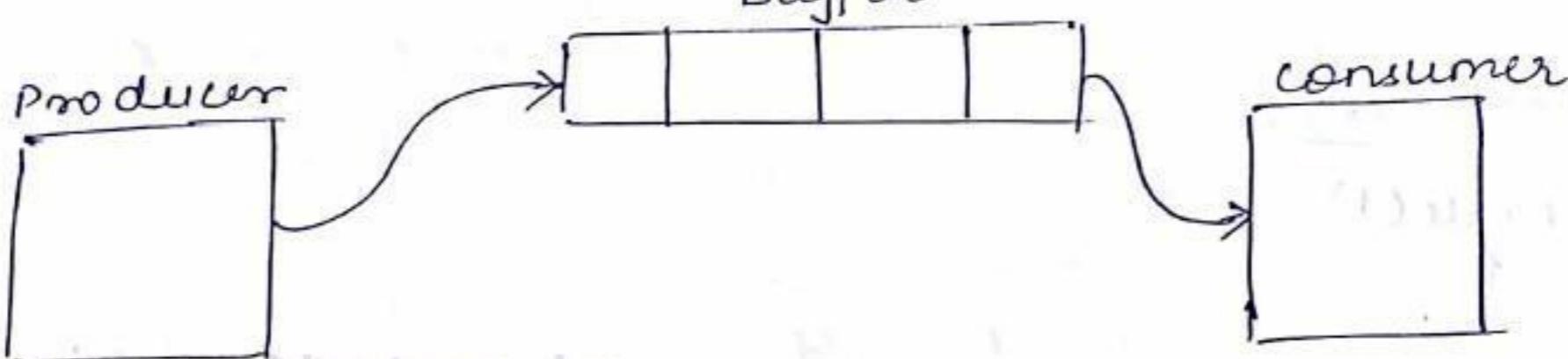
Rough. {  
 $S=1$   
 $Q=1$   
 wait(S)  $\rightarrow$   $S \neq 0$ .  
 (a)  $\neq 0$ .  
 c.s enters.  
 signal(S)  $S=1$ .  
 $Q=1$  then P1 can execute.

Rough {  
 wait(S).  
 context switch.  $S=0$ .  
 Now execute P0 or P1  
 wait(Q). busy waiting  
 deadlock  
 wait(S)  
 busy waiting }  
 wait(Q).  $Q=0$ .

{ if two signal ~~one~~<sup>count</sup> together then 2 processes may enter into c.s.

→ classic problems of synchronization (uv2)

i) producer-consumer problem / Bounded-buffer problem.



LectureNotes.in

- It is also known as Bounded-buffer problem as buffer size is bounded
- producer produces an item and stores in buffer, consumer consumes the item from buffer
- It has 2 problems
  - i) producer is interested to produce but the buffer is full
  - ii) consumer is interested to consume but the buffer is empty.
- Producer and consumer share the following data-structure.

LectureNotes.in

int n;

Semaphore mutex = 1; // mutual exclusion accessed to the buffer.

Semaphore empty = n; // initially all the n positions of buffer are empty.

Semaphore full = 0; // zero buffer is full means all are empty.

Producer

{ ... }

18/8/17

2218

(188) `int n;`  
semaphore mutex = 1;  
semaphore empty = n;  
semaphore full = 0;

x 0 x 0  
x n - - o  
x x 1 0

(2.)

### producer:

`while(1)`

{

`/* produce an item m in next-produced */`

`wait(empty);`

`wait(mutex);`

`/* add next-produced to the buffer */`

`signal(mutex);`

`signal(full);`

}

entry section  
c.s

E.S

### consumer:

`while(1)`

{

`wait(full);`  
`wait(mutex);`

`/* remove an item from buffer to new-consumed */`

`signal(mutex);`

`signal(empty);`

`/* consume the item m next-produced */`

--

}

22/8/17

### (2.) Reader-writer problem

1.  $R_1, R_2, R_3 \dots R_n$  ✓
2.  $R_1, W_1$  ✗
3.  $W_1, R_1$  ✗
4.  $W_1, W_2$  ✗

semaphore wrt = 1;

semaphore mutex = 1;

int readcount = 0;

#### writer

while(1)

{  
  | wait(wrt);  
  |    |  
  |    | /\* writing is performed \*/  
  |    |  
  |    | signal(wrt);  
  |    |  
  | } .

writer gets blocked  
(ii) case

#### Reader:

while(1)

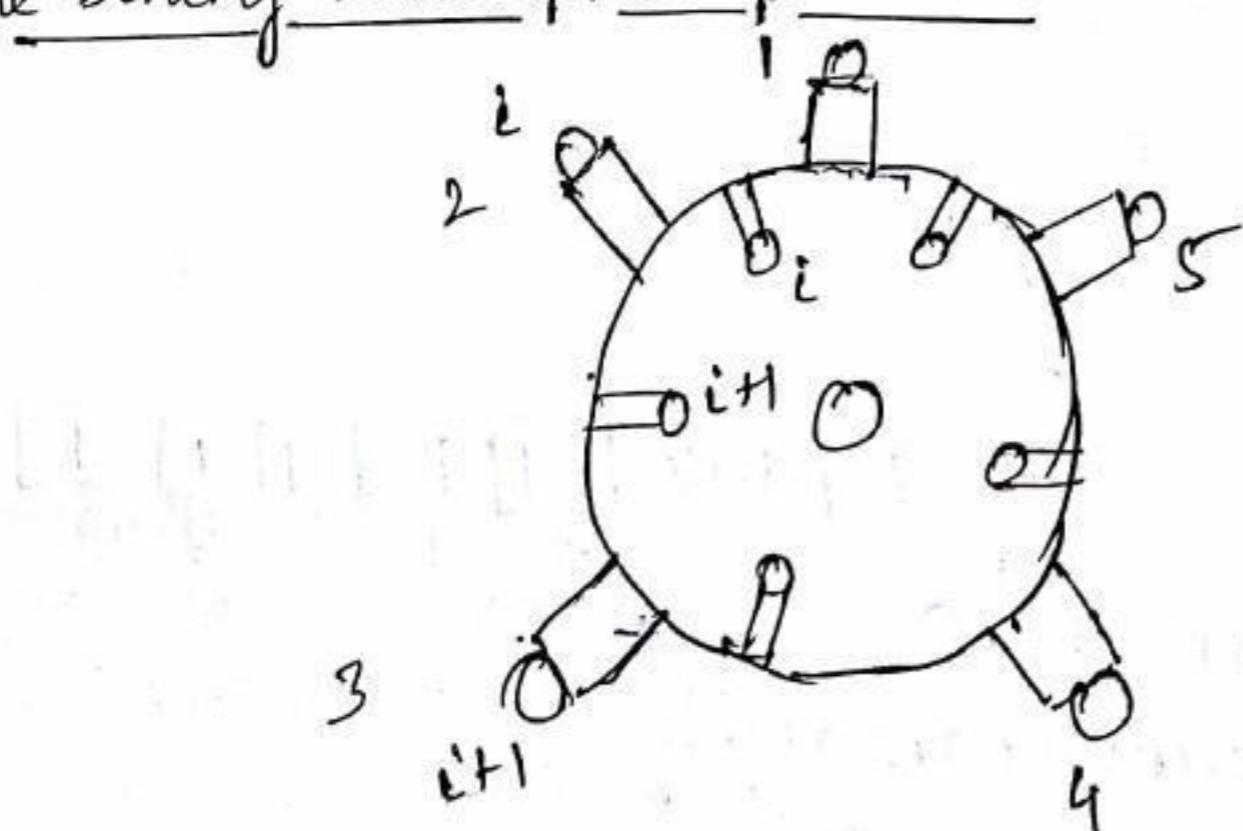
{  
  | wait(mutex);  
  |    |  
  |    | read-count++;  
  |    | if(read-count == 1)  
  |    |    | wait(wrt);  
  |    |    | wait(mutex);  
  |    |    | signal(mutex);  
  |    | } .

c.s.

/\* Reading is performed \*/

wait(mutex);  
read\_count--;  
if(readcount == 0)  
  signal(wrt);  
signal(mutex);  
} .

### (3.) The Dining Philosopher problem



- when a philosopher thinks, she does not interact with her colleagues.
- from time to time a philosopher gets hungry and tries to pick off the two chop sticks that are closest to her (left and right)
- A philosopher may pick up only one chopstick at a time
- she cannot pickup a chop stick that is already in the neighbour's hand.
- when a hungry philosopher has both chopsticks at the same time, she eats without releasing the chopstick
- when she finisheds eating, she puts down both the chopsticks and start thinking again.
- A philosopher pick up the chopstick by executing wait() operation and releases the chopstick using signal() operation.

Semaphore chopstick[5]; // Initialize to 1.

do

{

wait [chopstick[i]];  
wait [chopstick[(i+1) % 5]];

c.s.

[ /\* eat for a while \*/ ]

signal [chopstick[i]];

signal [chopstick[(i+1) % 5]];

[ ]

/\* think for a while \*/

} while(true)

Chopstick	1	1	1	1	1	1
	0	1	2	3	4	

- This guarantees that no two neighbours are eating simultaneously.

- Suppose that all five processes become hungry at the same time and each grabs for left chopstick and the right chopstick is not available for any philosopher. This leads to deadlock.
- Possible remedies <sup>to avoid</sup> of deadlock are
  - i) allow atmost four philosophers to sit simultaneously on the table
  - ii) allow a philosopher to pick up her left chopstick only if both chopsticks are available.
  - iii) an odd no. of philosopher picks up her left chopstick first and then right chopstick whereas even-numbered philosopher picks up her right chopstick first and then left chopstick.

## MONITORS.

### MONITORS.

Although semaphore provides a convenient and efficient mechanism for process synchronization but misusing the semaphore results in error.

- 1.)  $\left( \begin{array}{l} \text{wait(mutex)} \\ \text{c.s} \\ \downarrow \\ \text{signal(mutex)} \end{array} \right) \Rightarrow \begin{array}{l} \text{signal(mutex)} \\ \text{c.s} \\ \downarrow \\ \text{wait(mutex)} \end{array}$   
 problem: several process enter into the critical section simultaneously, violating mutual exclusion requirement.

- 2.)  $\text{wait(mutex)}$   
 $\text{c.s}$

$\text{wait(mutex)}$

problem: deadlock occurs.

- 3.) Suppose that a process omits  $\text{wait(mutex)}$  or  $\text{signal(mutex)}$  or both, <sup>then</sup> either mutual exclusion is violated or deadlock will occur.

4) To avoid these problems, monitor is introduced.

```
monitor monitor-name
{
    /* shared variable declaration */
    function p1()
    {
    }

    function p2()
    {
    }

    function p3()
    {
    }
}
```

initialization\_code(-)

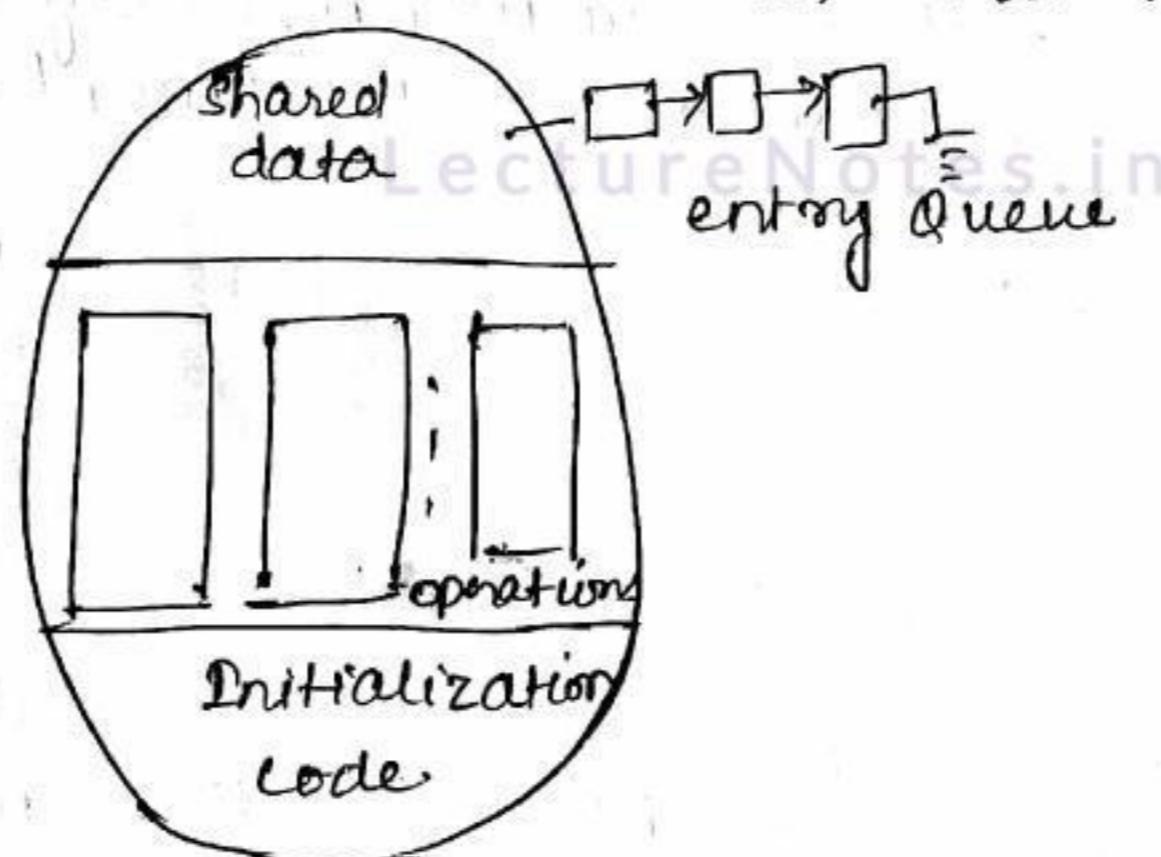
{ }

}

Syntax of monitor

#### → Monitor usage

A monitor type is an ADT that includes set of programmer defined operations that are provided with mutual exclusion within the monitor.



\* Schematic view of \* monitor

- A programmer who needs to write synchronization scheme can also define variables of type 'condition' condition x, y;
  - The operation `x.wait()` means that the process involved in this operation is suspended under `o` until another invokes i.e. `x.signal()`
- the fork() system call.

23/8/17

LectureNotes.in

```

Parent           child - returns 0.
{
    fork();      returns
    pid = ____;  CPID.
}

```

child processes  $\rightarrow 2^n$  processes  
 $\leq \{$  (fork) ~~(child)~~ child process  $2^n - 1$

system call `fork()` is used to create process. It takes no arguments and returns a process ID. The purpose of `fork()` is to create a new process which becomes the child process of the caller.

see pdf

LectureNotes.in

Scanned by CamScanner

After MidsemDeadlocks

A process request for resources. If the resources are not available at that time, the process enters a wait state. Waiting processes may never again change state, because the resources that have requested are held by other waiting processes. This situation is called deadlock.



Example of resources are -

CPU cycles, memory space, I/O devices etc.

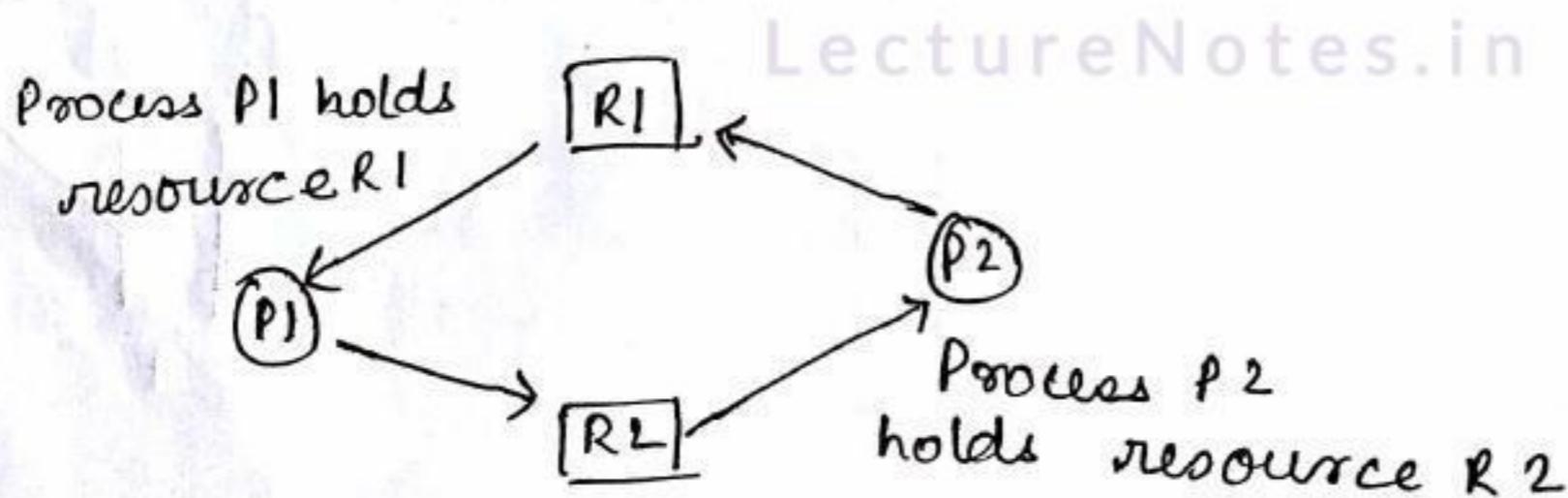
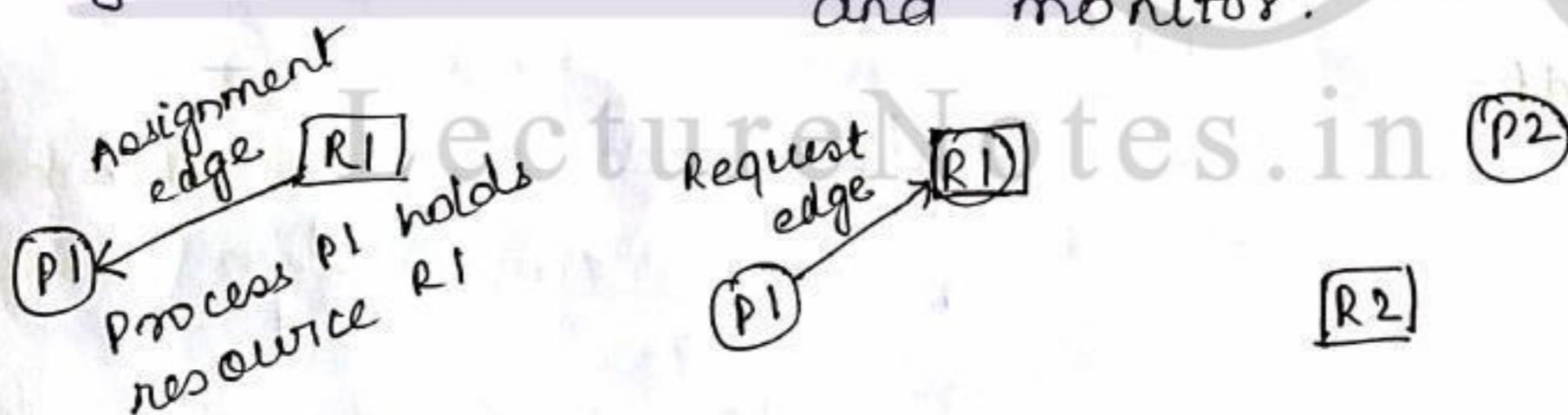
A process must request a resource before using it and must release the resource after using it.

A process may utilize a resource in following sequence -

1. Request (System call)
2. Use
3. Release (System call)

There are two types of resource -

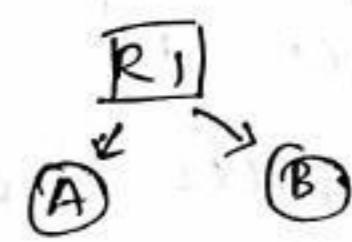
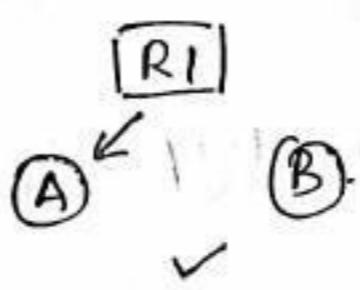
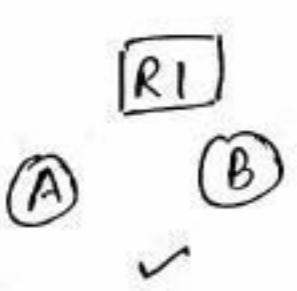
1. Physical Resource - CPU, Memory, I/O.
2. Logical resource - Read only files, semaphore and monitor.



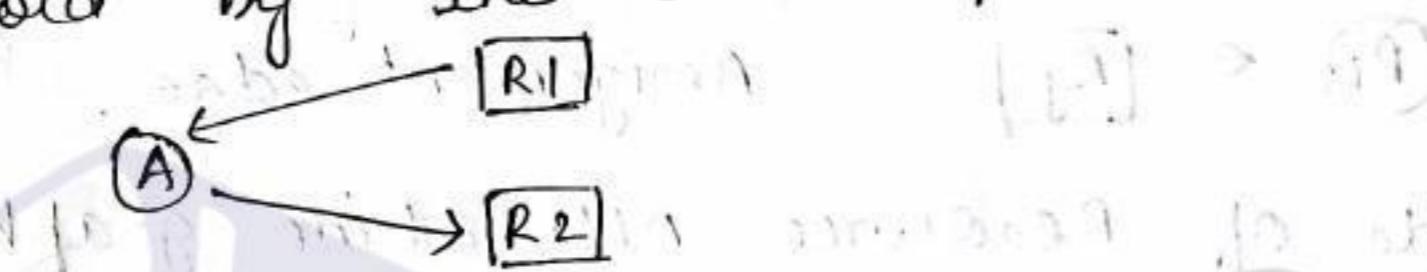
If there is a cycle in the graph, there may be chances of deadlock.

Necessary conditions for deadlock:-  
A deadlock situation can occur if the following 4 conditions hold simultaneously in the system -

1. Mutual exclusion - Only one process may use the resource at a time.

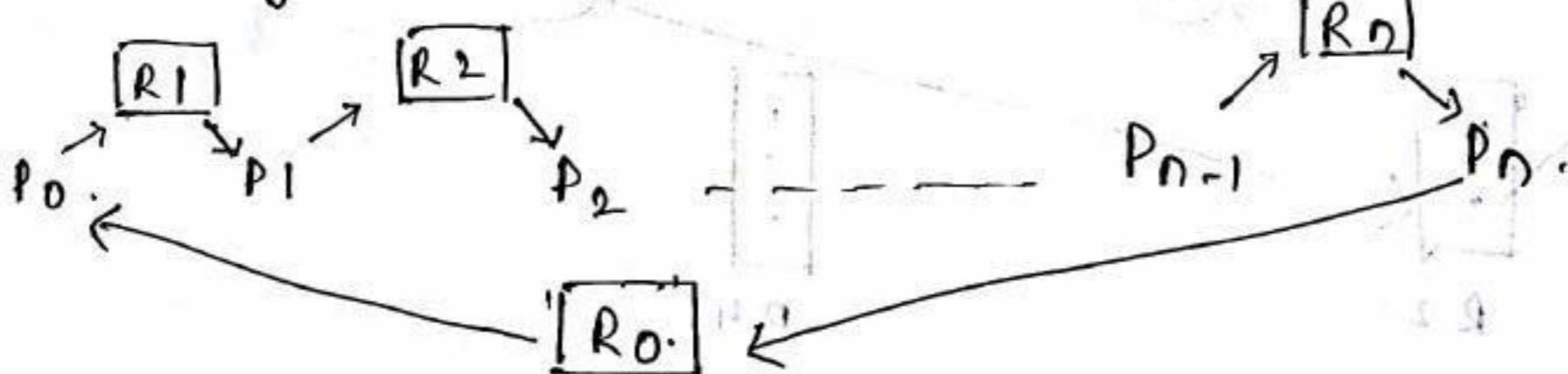


2. Hold and wait - A process must be holding at least one resource and waiting to acquire additional resources that are currently being held by the other processes.



3. No preemption : Resources cannot be preempted in the middle of the execution. A resource can only be released if the process has completed its task.

4. circular wait : A set  $\{P_0, P_1, \dots, P_{n-1}, P_n\}$  of waiting processes must exist such that  $P_0$  is waiting for a resource that is held by  $P_1$ ,  $P_1$  is waiting for the resource that is held by  $P_2$ ,  $\dots$ ,  $P_{n-1}$  is waiting for a resource that is held by  $P_n$  and  $P_n$  is waiting for a resource held by  $P_0$ .



## Resource Allocation graph.

deadlock can be described in terms of a directed graph called resource allocation graph. This graph consists of set of vertices edges  $E$ .

vertices are of two types -

1.  $\circ$  process

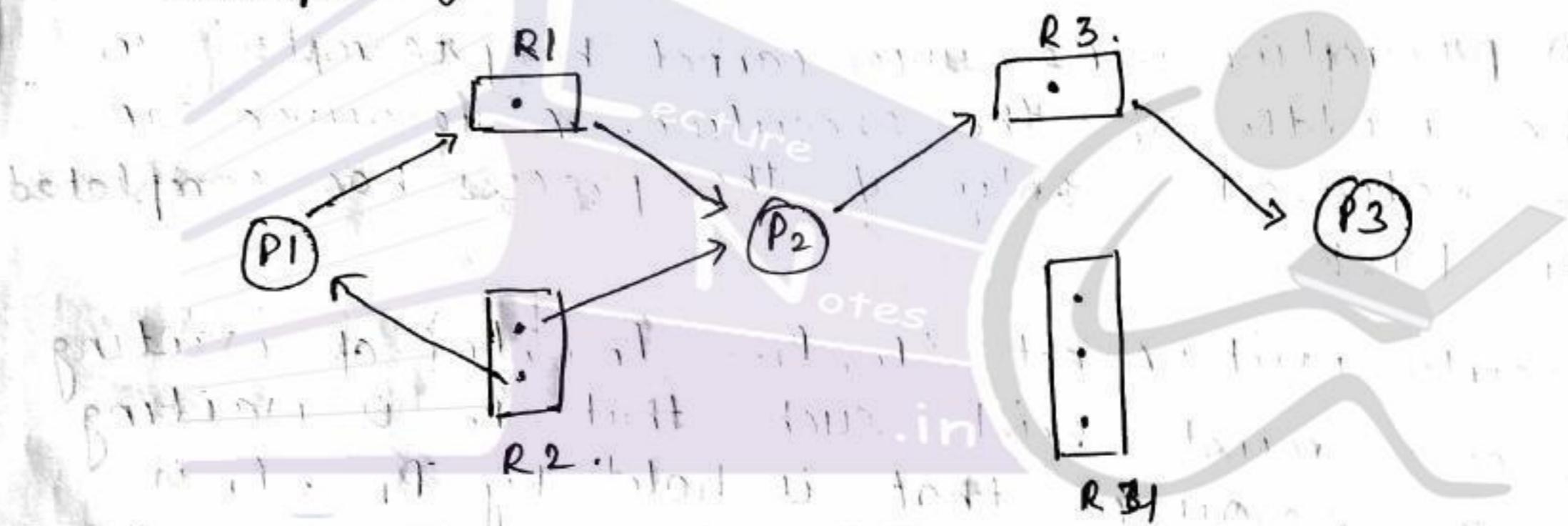
2.  $\square$  resource

edges of two types -

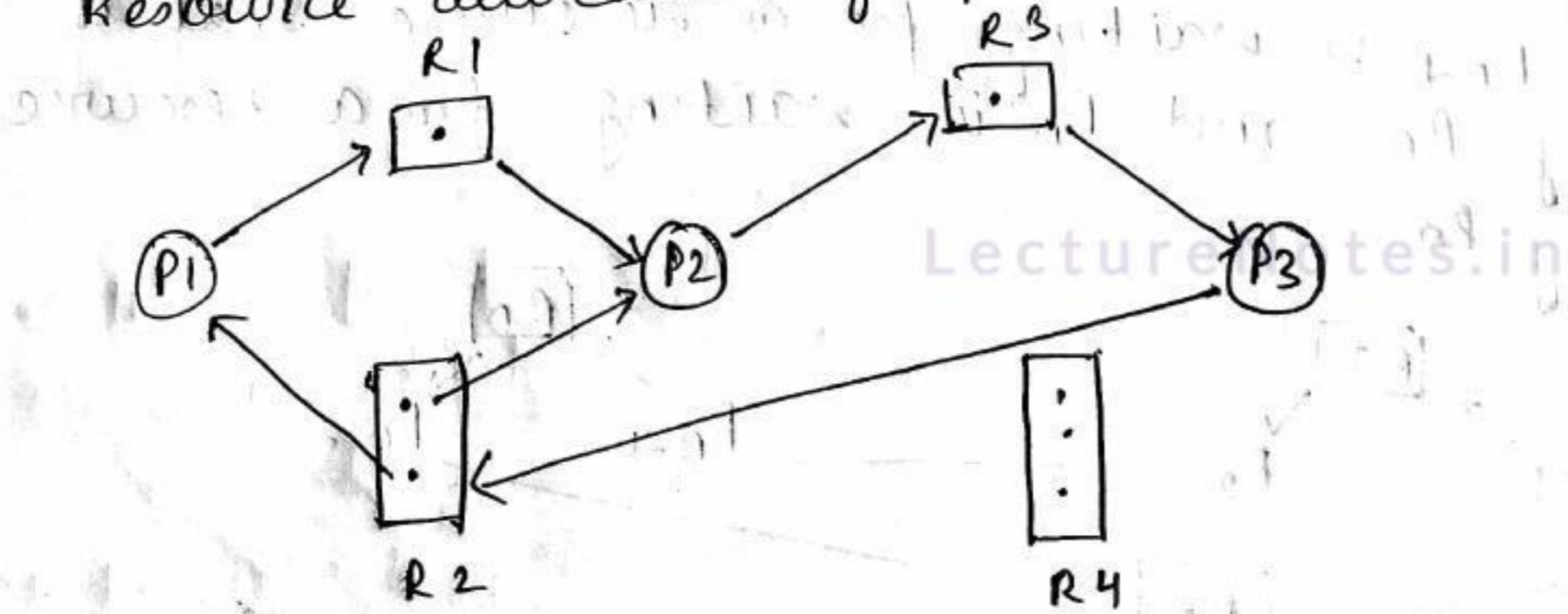
1.  $(P_i) \rightarrow [R_j]$ , Request edge

2.  $(P_i) \leftarrow [R_j]$ , Assignment edge.

example of Resource Allocation graph



Resource allocation graph without deadlock.

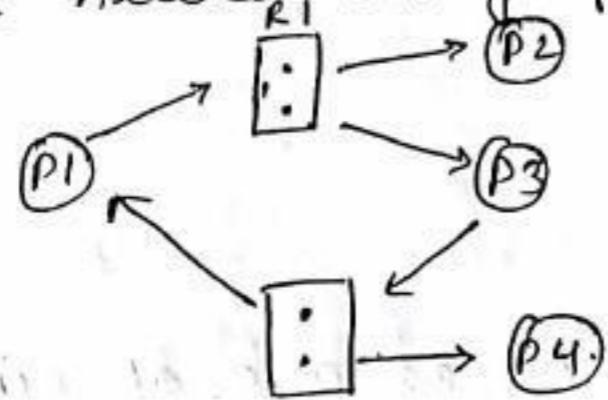


1.  $P_1 \rightarrow R_1 \rightarrow P_2 \rightarrow R_3 \rightarrow P_3 \rightarrow R_2 \rightarrow P_1$

2.  $P_2 \rightarrow R_3 \rightarrow P_3 \rightarrow R_2 \rightarrow P_2$ .

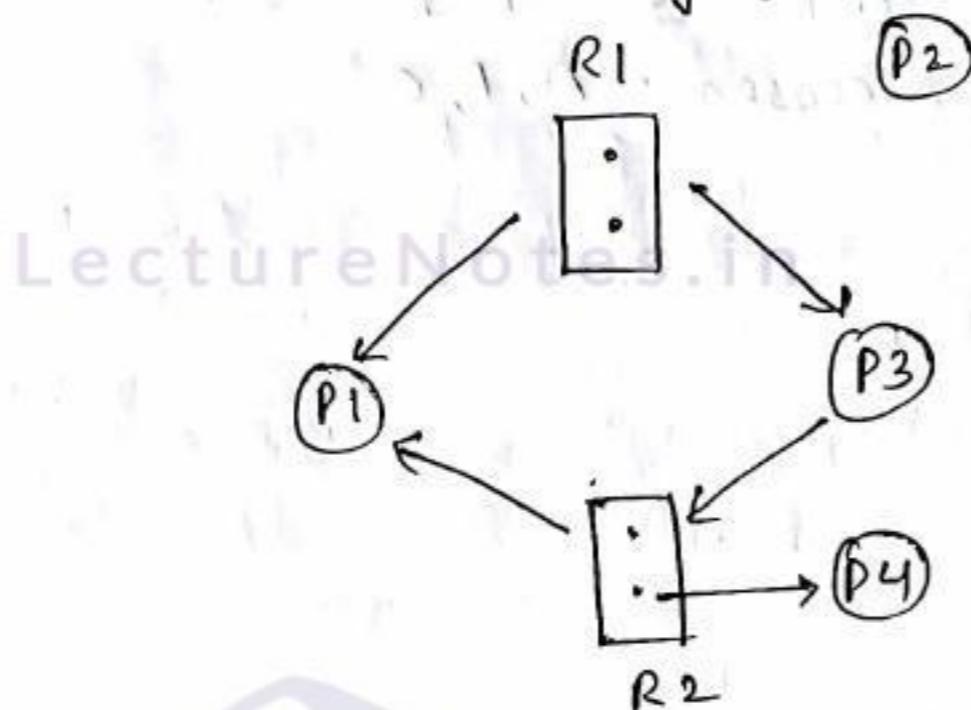
Deadlocked processes are  $P_1$ ,  $P_2$ ,  $P_3$ .

Resource Allocation graph with deadlock.



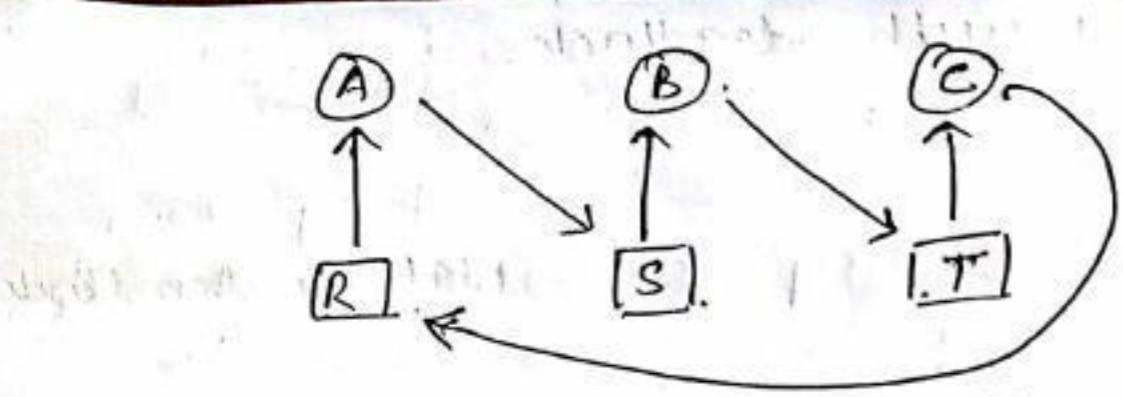
Not in deadlock

Resource allocation graph without deadlock



- If the graph contains no cycle then the processes are not in deadlock.
- If the graph contains a cycle then a deadlock may exist.
- If each resource type has exactly one instance and a cycle exist in resource allocation graph implies deadlock.
- If each resource type has several instances, then a cycle does not necessarily imply that a deadlock has occurred.  
so, here a cycle in the graph is necessary but not a sufficient condition for the existence of deadlock.

- Q.) 1. A requests R  
2. B " S  
3. C " T  
4. A " S  
5. B " T  
6. C " R.

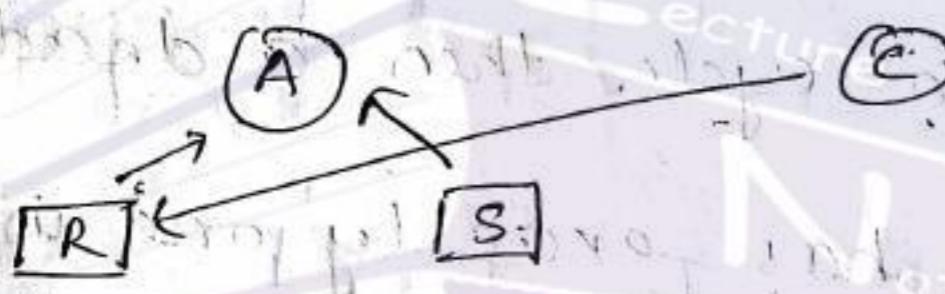


There is cycle and all resources are of single instance hence deadlock.

cycle :  $A \rightarrow [S] \rightarrow [T] \rightarrow [C] \rightarrow [R] \rightarrow A$

Deadlocked processes : A, B, C.

- Q. 1. (A) requests [R]
2. (C) " [T]
3. (A) " [S]
4. (C) " [R]
- 5 (A) releases [R]
6. (A) releases [S]



NO cycle hence no deadlock.

Methods for handling deadlock - we can deal with the deadlock problem in one of the three ways

1. Deadlock prevention
2. Deadlock avoidance
3. Deadlock detection and recovery.

1.) Deadlock prevention - One can prevent deadlock by avoiding a policy that eliminates one of the conditions such as mutual exclusion, hold and wait, no preemption, and circular wait.

2.) Deadlock avoidance - Deadlock avoidance requires the advanced additional information of resource allocation of the processes. By additional knowledge it can be determined whether there is a deadlock or not.

3) Deadlock detection and recovery - we can allow the system to enter into the deadlock state by holding 4 conditions- mutual exclusion, hold or wait, no preemption and circular wait and take action to recover it.

#### → Deadlock prevention

Deadlock prevention is a set of methods for ensuring that ~~at least~~ atleast one of the four conditions ME, hold and wait, no preemption and circular wait cannot hold.

• Mutual exclusion - Mutual exclusion means only one process can use the resource at a time. It means resources are not shared by no. of processes at a time. so the resources are non-shareable in nature. To avoid this condition use protocol which convert non-shareable resource to shareable resource.

ex : Non shareable - Printer

shareable - Read only files.

Printer is not shared by no. of processes at a time, so we cannot convert non-shareable to shareable node.

Disadvantage - we cannot avoid deadlock by converting non-shareable to shareable for all resources

- Hold and wait - to avoid hold and wait we must ensure that it does not hold any other resource. Two protocols are used -
  - 1.) It allows a process to request and allocate all its resource before the execution begins.
  - 2.) Allows a process to request only when the process has none.

It means, a process can request any additional resource by releasing all the resources that is currently allocated. ex - printer, pendrive etc.

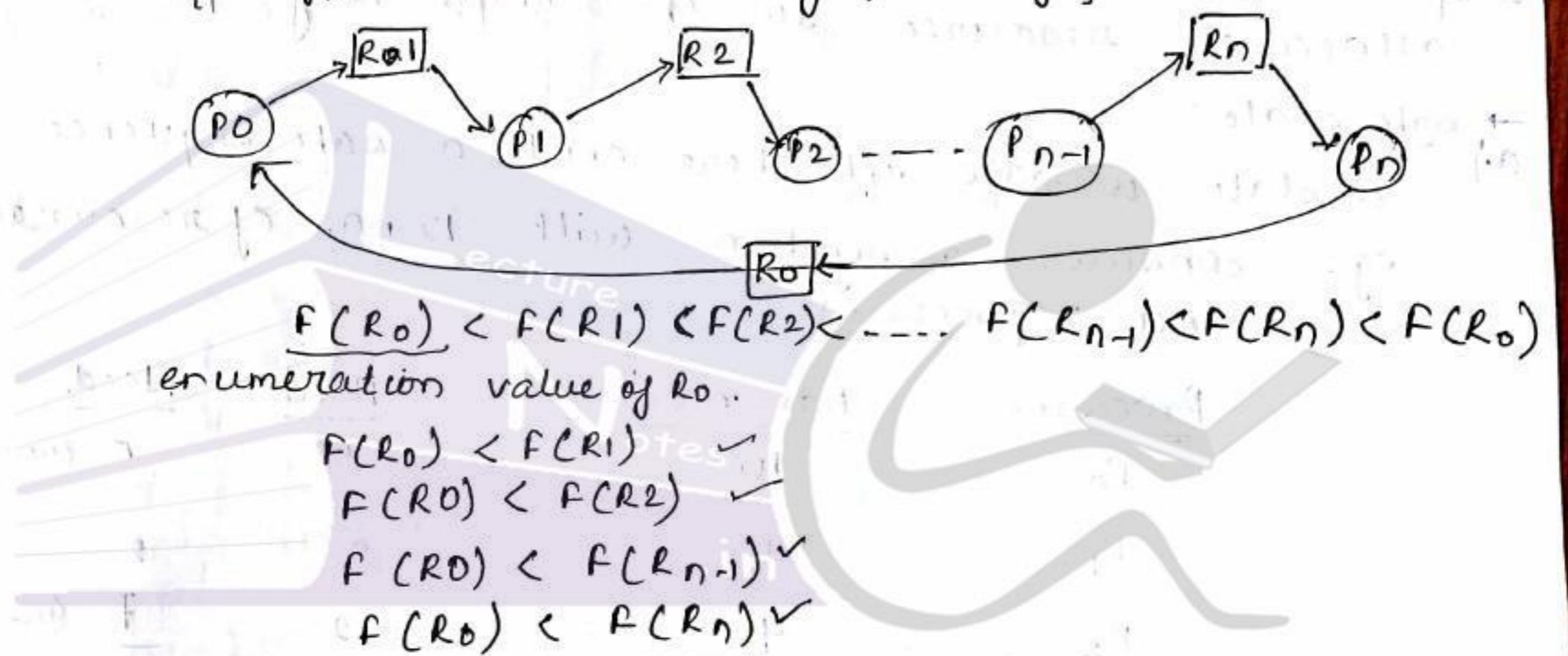
Disadvantage :- Resource utilization is low and starvation may occur.
- No preemption - No preemption means resources are not released in the middle of the execution. Two protocols required are -
  - 1.) If a process is holding same resource and requests other resources that cannot be immediately allocated to it, then all the resources currently being held are preempted.
  - 2.) If a process requests some resources, we first check whether they are available or not. If they are available, we allocate them. If the resources are not available we check whether they are allocated to some other resources that is waiting for additional resource. If so we preempt the desired resource from the waiting process and allocate them to the request process.

Disadvantage : It cannot be applied to resources like printer, cd & drivers, tape drivers etc.

12/9/17

- circular wait condition never holds ordering of all the resources in an increasing order of enumeration
- to prevent deadlock following protocols are used.
  - each process can request resources only in an increasing order of enumeration.
  - a process can request resources of type  $R_j$  if and only if  $F(R_j) > F(R_i)$
  - whenever a process request an instance of resource type  $R_j$ , it has released only resources arrived  $R_i$  such that  $F(R_i) \geq F(R_j)$

$\{F(R_j)\}$ : Resources held by process  $j\}$



enumeration value of  $R_0$ .

$$\begin{aligned} F(R_0) &< F(R_1) \\ F(R_0) &< F(R_2) \\ F(R_0) &< F(R_{n-1}) \\ F(R_0) &< F(R_n) \end{aligned}$$

Here circular wait condition breaks and becomes false.

By transitivity we get

$$F(R_0) < F(R_0)$$
 which is impossible so

there can be no circular wait

disadvantage: circular wait prevention may be inefficient, slowing down processes and denying resource access unnecessarily.

## 2.) Deadlock avoidance

- Deadlock avoidance algorithm dynamically examines the resource allocation state to ensure that a circular wait condition can never exist.
- Resource allocation state is defined by no. of available and allocated resources and the maximum demand of the processes.
- If a process request for a resource the avoidance algorithm checks before the allocation of resources about the state of the system.
- If the state is safe, system allocates resources to the requesting process.
- If the state is unsafe, then the system do not allocate resource to the requesting process.

→ safe state

- a.) A state is safe iff there exist a safe sequence
- eg:- consider a system with 12 no. of resources
- (1) 1 and 3 processes

<u>Processes</u>	<u>Max. need</u>	<u>Hold</u>	<u>Need</u>
P <sub>0</sub>	10	5	5. (unsafe)
P <sub>1</sub>	4	2	2
P <sub>2</sub>	9	2	7 (unsafe)

currently available resources = Total no. of resources - Total no. of currently holding resources

$$12 - 9 = 3$$

At time t<sub>0</sub>, granting 3 resources to P<sub>0</sub> and P<sub>2</sub> will leave the system in unsafe state.

At time t<sub>1</sub>, 2 resources are allocated to process P<sub>1</sub>, the currently available =  $\underbrace{3 + 2}_{3 - 2} = 1$ .

$$\left\{ \begin{array}{l} 3 - 2 = 1 \\ 1 + 4 = 5 \end{array} \right.$$

When all resources are released by P<sub>1</sub>.

At time  $t_2$ , 5 resources are allocated to  $P_0$ . It becomes safe and get executed.

After execution currently available =  $5+5 = 10$ .

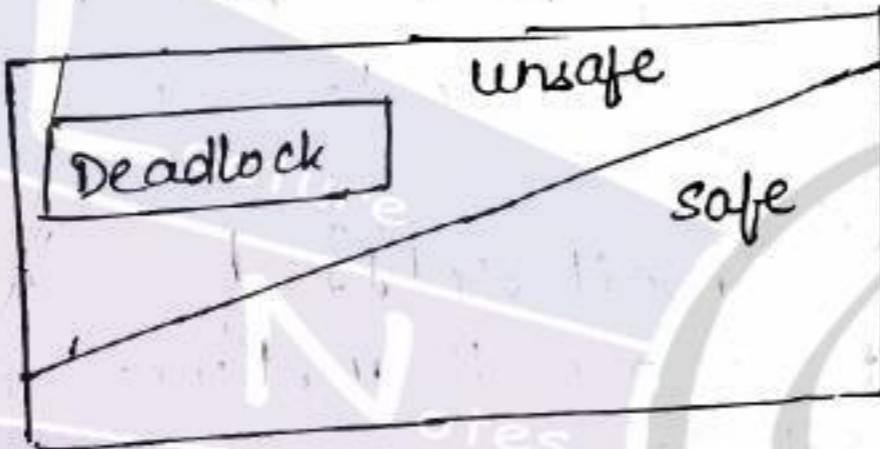
At time  $t_3$ , 7 resources are allocated to  $P_2$ , it becomes safe and get executed.

After execution  $9+3 = 12$ .

currently available

The safe sequence is  $\langle P_1, P_0, P_2 \rangle$

- If there is safe sequence in your system, your system is safe so no deadlock.
- A safe state is not a deadlock state.
- A deadlock state is an unsafe state but not all unsafe states are deadlock.
- An unsafe state may lead to a deadlock.

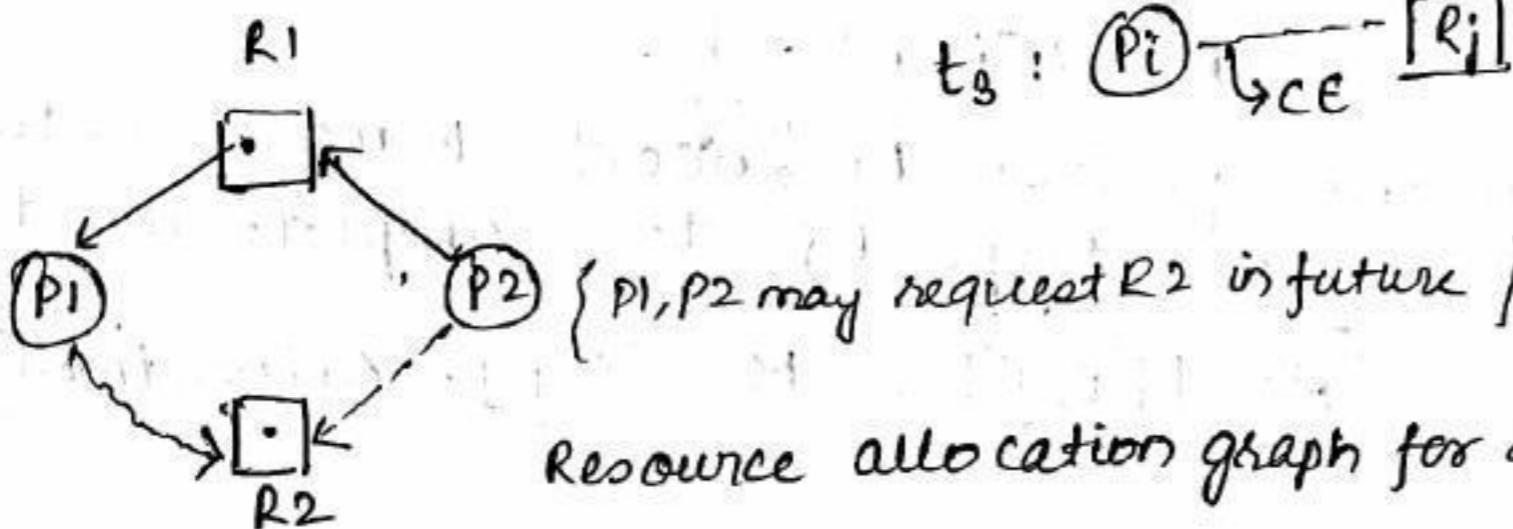
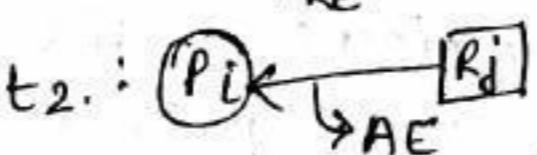
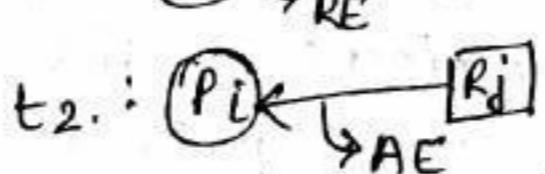
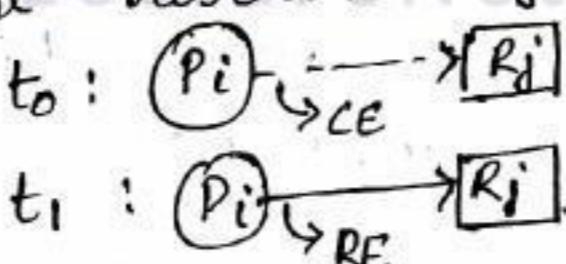


13/9/17

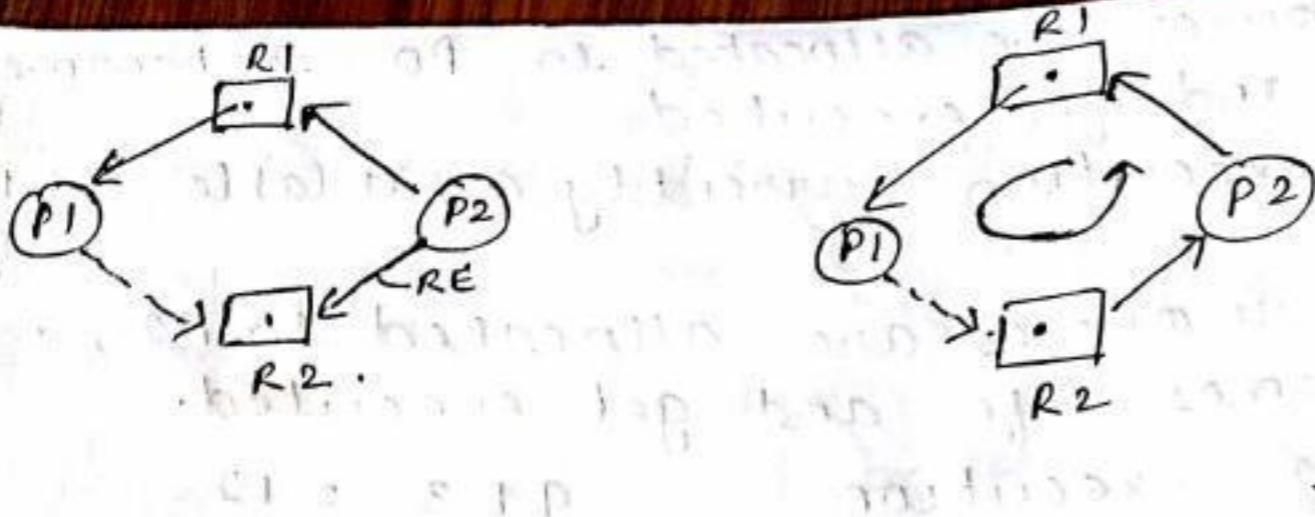
b) Resource allocation graph algorithm.

for deadlock avoidance, resource allocation graph contains a resource allocation system with only one instance of each resource type.

claim edge: A claim edge  $P_i \rightarrow R_j$  indicates that process  $P_i$  may request resource  $R_j$  at some time in future to:



Resource allocation graph for deadlock avoidance.



If  $P_1$  requests  $R_2$ , then there will be cycle. This happens due to this is unsafe state.

### → Banker's Algorithm.

- Resource allocation graph algorithm for resource allocation system with each resource type.
- Let  $n$  be the no. of processes and  $m$  be the no. of resource type.
- Available : it is a vector of length  $m$  indicates no. of available resources of each type.  
 $\text{Available}[j] = k$ .  
 It means there are  $k$  instances of resource type  $R_j$  are available.
- Max :  $\text{max}[i, j] = k$ .  
 It means process  $P_i$  request atmost  $k$  instances of resource type  $R_j$ .
- Allocation : no. of resources of each type currently allocated to each process.  
 $\text{Allocation}[i, j] = k$ .  
 It means process  $P_i$  is currently allocated with  $k$  instances of resource type  $R_j$ .
- Need : It indicates the remaining resources needed by each process  
 $\text{Need}[i, j] = k$ .  
 It means process  $P_i$  <sup>may</sup> need  $k$  more instances of resource type  $R_j$  to complete its task.

$$\boxed{\text{Need}[i, j] = \text{Max}[i, j] - \text{Allocation}[i, j]}$$

### i) safety algorithm

This algorithm finds whether the system is safe or unsafe.

#### • Step 1:

Let work and finish be vectors of length m and n.  
work := Available

Finish[i] = false for  $i = 1, 2, \dots, n$ .

#### • Step 2:

Find an i such that

- a) Finish[i] = False  $\{/\!\!/\text{ process must not be completed}\}$
- b) Need $_i \leq \text{work}$ .  $\{/\!\!/\text{ need should be} \leq \text{available}\}$

If not such i is found, go to step 4.

#### • Step 3:

work := work + Allocation;

Finish[i] = True.

$\{/\!\!/\text{ we need to release when process executes completely.}\}$

goto step 2  $\{/\!\!/\text{ for other processes here only 1 process gets completed.}\}$

#### • Step 4:

If Finish[i] = true for all i, then the system is in safe state.

### ii) resource request algorithm

When a request for a resource is made by process  $P_i$  then following actions are taken:-

#### • Step 1:

If  $\text{Request}_i \leq \text{Need}_i$ , goto step 2 otherwise an error will generate.

#### • Step 2:

If  $\text{Request}_i \leq \text{Available}$ , goto step 3. Otherwise  $P_i$  must wait till the resources are available.

#### • Step 3:

The system allocates the requested resources to process  $P_i$  by modifying the following:

Available := Available - Request $_i$

Allocation $_i$  := Allocation $_i$  + Request $_i$

Need $_i$  := Need $_i$  - Request $_i$

If resource allocation state is safe  
is completed and process  $P_i$  is allocated with resources  
If the state is unsafe then process  $P_i$  must wait  
for request:

example.

consider a system with 5 processes  $P_0$  to  $P_4$  and 3  
resource type A, B and C. Resource type A has 10  
instances, B has 5 instances and C has 7 instances.

- i.) check the system is safe or unsafe
- ii.) Suppose process  $P_1$  request one more instance of  
resource type A and 2 more instances of resource  
type C. Can the request be granted?

14/9/17.

At time  $T_0$

Processes	Allocation (Hold)			Max			Available			Need.			
	A	B	C	A	B	C	A	B	C	A	B	C	
$P_0$ (false)	0	1	0	7	5	3	3	3	2	7	4	3	X✓
$P_1$ ( $\cancel{A}$ ) $(\cancel{C})$	2	0	0	3	2	2				1	2	2	✓
$P_2$ ( $F$ )	3	0	2	9	0	2				0	1	1	✓
$P_3$ ( $\cancel{F}$ ) $(\cancel{T})$	2	1	1	2	2	2				4	3	1	✓
$P_4$ ( $\cancel{F}$ ) $(\cancel{T})$	0	0	2	4	3	3							

Ans:-

$$\text{i) Available} = \text{Total-allocation} \\ (10 \ 5 \ 7) - (7 \ 2 \ 5) \\ = (3 \ 3 \ 2)$$

ii) apply safety algorithm

step1 work = available =  $(3 \ 3 \ 2)$   
finish[i] = false

step2 find i such that  $\text{Need}_i \leq \text{work} \& \text{finish}[i] \neq \text{POX}$

Process -  $P_1 \checkmark \quad \text{Need}_i \leq \text{work}$

if  $\text{Need}_i \leq \text{work}$   
 $\text{Need of } P_1 (1 \ 2 \ 2) \leq (3 \ 3 \ 2)$

$$\begin{aligned}
 \text{work} &= \text{work} + \text{Allocation} \\
 &= (332) + (200) \\
 \boxed{\text{work} = (532)} &\quad \text{finish}[i] = \text{true}
 \end{aligned}$$

Need  $P_3(011) \leq \text{work}(532)$

$\text{work} = \text{work} + \text{Allocation}$

$$(532) + (211)$$

$$\boxed{\text{work} = (743)}$$

Need  $P_4(431) \leq \text{work}(743)$

$$\begin{aligned}
 \text{work} &= \text{work} + \text{Allocation} \\
 &= (743) + (002)
 \end{aligned}$$

$$\boxed{\text{work} = (745)}$$

Next loop.

Need  $P_0(743) \leq \text{work}(745)$

$$\begin{aligned}
 \text{work} &= \text{work} + \text{allocation} \\
 &= (745) + (010)
 \end{aligned}$$

$$\boxed{\text{work} = (755)}$$

Need  $P_2(600) \leq \text{work}(755)$

$$\begin{aligned}
 \text{work} &= \text{work} + \text{allocation} \\
 &= (755) + (302)
 \end{aligned}$$

$$= (1057).$$

safe sequence:  $\langle P_1, P_3, P_4, P_0, P_2 \rangle$ .

Hence the system is in safe state.

	A	B	C
P1	1	0	2

ii) Request : Resource, request algorithm.

1.)  $\text{Request}_{P_1} \leq \text{Need}_{P_1}$ .

$$(102) \leq (122) \checkmark$$

2.)  $\text{Request}_{P_1} \leq \text{Available}$

$$(102) \leq (332) \checkmark$$

If these cond's are false then request will not be granted here granted.

$$\text{available} = \text{available} - \text{Request}_{P_1}$$

$$(3 \ 3 \ 2) - (1 \ 0 \ 2)$$

$$= (2 \ 3 \ 0)$$

15/9/17

3.) Dec

$$\text{allocation}_{P_1} = \text{allocation}_{P_1} + \text{request}_{P_1}$$

$$(2 \ 0 \ 0) + (1 \ 0 \ 2)$$

$$= (3 \ 0 \ 2)$$

$$\text{Need}_{P_1} = \text{Need}_{P_1} - \text{Request}_{P_1}$$

$$= (1 \ 2 \ 2) - (1 \ 0 \ 2)$$

$$= (0 \ 2 \ 0)$$

The new state of the system will be:-

Process.	Allocation			Max			Available			Need		
	A	B	C	A	B	C	A	B	C	A	B	C
P <sub>0</sub>	0	1	0	7	5	3	2	3	0	7	4	3
P <sub>1</sub>	3	0	2	3	2	2	0	2	0	0	2	0
P <sub>2</sub>	3	0	2	9	0	0	2	0	0	6	0	0
P <sub>3</sub>	2	1	1	2	2	2	0	1	1	0	1	1
P <sub>4</sub>	0	0	2	4	3	3	3	4	3	4	3	1

Now you must determine whether this <sup>new</sup> system is safe or not. To do so we execute safety algorithm again and find the safe sequence as.

This : safe  $\langle P_1 \ P_3 \ P_4 \ P_0 \ P_2 \rangle$  which satisfies our safety requirement. Hence we can grant the request for process  $P_1$ .

15/9/17

### 3.) Deadlock detection

If a system does not employ either deadlock prevention or deadlock avoidance algorithm then a deadlock may occur so a system must provide.

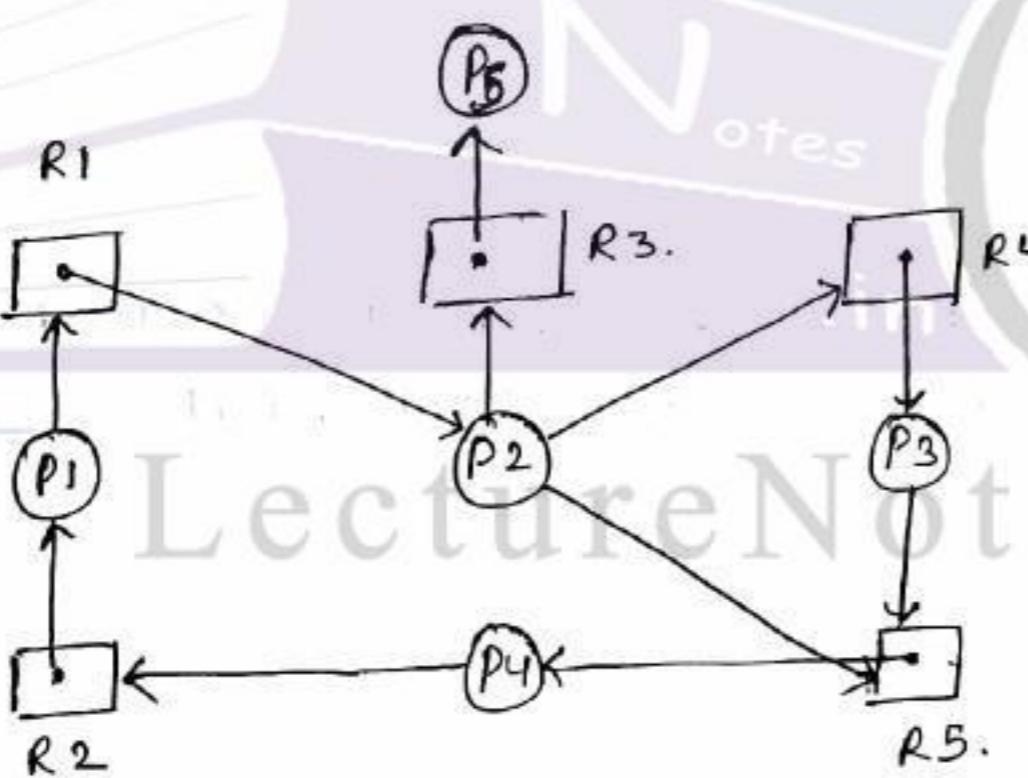
- i) an algorithm that examines the deadlock has occurred or not (deadlock detection algorithm).
- ii) an algorithm to recover from deadlock.

#### a.) single instance of each resource type:

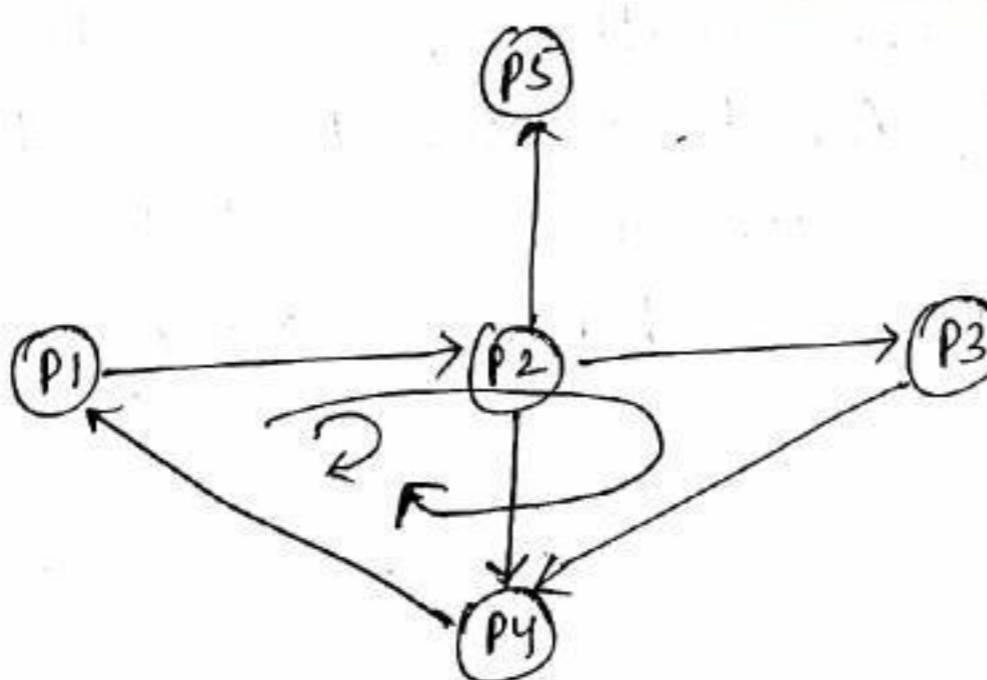
If all the resources have only a single instance then wait-for graph is used as deadlock detection algorithm.

We obtain wait-for graph <sup>from</sup> for resource allocation graph by removing resource nodes and collapsing the appropriate edges.

A deadlock exist in the system iff the wait for graph consists a cycle.



(Resource allocation graph).



$$P_1 \rightarrow R_1 \quad R_1 \rightarrow P_2 \\ \therefore P_1 \rightarrow P_2$$

wait-for-graph.

b.) multiple instances of a resource type.  
The deadlock detection algorithm is applicable to a system with multiple instances of each resource type.

### deadlock detection algorithm

Step 1: let work & finish be vectors of length  $m$  and,

respectively.

work := Available for  $i=1, 2, \dots, n$ .

if Allocation $_i \neq 0$ , then Finish $[i] = \text{false}$ .

otherwise, Finish $[i] = \text{true}$ .

Step 2: for an index  $i$  such that

a) Finish $[i] = \text{false}$

b) Request $_i \leq \text{work}$

if no such  $i$  exist, go to step 4.

Step 3:

work := work + Allocation $_i$

Finish $[i] = \text{true}$

goto step 2.

Step 4:

if Finish $[i] = \text{false}$  for some  $i$ ,  $1 \leq i \leq n$ , then  
the system is in deadlock state.

example

example

we consider a system with 5 processes P<sub>0</sub> to P<sub>4</sub> and these resources are A, B and C. Resource type A has 7 instances, B has 2 instances and C has 6 instances. Suppose at time T<sub>0</sub> we have following resource allocation state.

Process

Process	Allocation			Request			Available		
	A	B	C	A	B	C	A	B	C
P <sub>0</sub>	0	1	0	0	0	0	0	0	0
P <sub>1</sub>	2	0	0	2	0	2	1	1	1
P <sub>2</sub>	3	0	3	0	0	0	1	1	1
P <sub>3</sub>	2	1	1	1	0	0	1	1	1
P <sub>4</sub>	0	0	2	0	0	2	1	1	1

LectureNotes.in

$$\text{Step 1: } \text{work} = 000; \text{ available} = 000.$$

$$\text{Step 2: } \text{Request}_0(000) \leq \text{work}(000) \checkmark$$

$$\text{Step 3: } \text{work}_0 = \text{work}_0 + \text{allocation}_0 \\ 000 + 010.$$

$$\text{Request}_1(202) \leq \text{work}(010) \times$$

$$\text{Request}_2(000) \leq \text{work}(010) \checkmark$$

$$010 + 303 = 313.$$

$$\text{Request}_3(100) \leq \text{work}(313) \checkmark$$

$$\text{Request}_4(002) \leq \text{work}(524) \checkmark$$

$$\text{work} = 313 + 211 = 524.$$

$$\text{Request}_1(202) \leq \text{work}(524) \checkmark$$

$$\text{work} = 524 + 200 = 724.$$

safe sequence :  $\langle P_0, P_2, P_3, P_4, P_1 \rangle$

so the system is not in the deadlock state.

→ suppose the process P<sub>2</sub> makes an additional request for an instance of resource type C then find out which processes are in deadlock state.

- P<sub>2</sub> Request  
001

$$\text{work} = \text{available} = 000.$$

$$\text{Request}_0(000) \leq \text{work}(000) \checkmark$$

$$\text{work} = 000 + 010 = 010.$$

$$\text{Request}_1(2002) \leq \text{work}(010) \times$$

$$\text{Request}_2(001) \leq \text{work}(010) \times$$

$$\text{Request}_3(100) \leq \text{work}(010) \times$$

$$\text{Request}_4(002) \leq \text{work}(010) \times$$

P<sub>0</sub> is true  
rest  
all false.

Now the system is in deadlock because after execution of P<sub>0</sub>, the no. of available resources is not sufficient to fulfill the requests of other processes. Thus the processes P<sub>1</sub>, P<sub>2</sub>, P<sub>3</sub>, P<sub>4</sub> are deadlocked.

## → RECOVERING FROM DEADLOCK.

- When detection algorithm exists, then we can deal with deadlock manually or automatically.
- To break the deadlock either we apply:
  - i) process termination
  - ii) Resource preemption
- iii) process termination:

To eliminate deadlock by aborting a process and releasing all its resources. We basically use two methods for process termination.

- a.) abort all deadlock processes
- b.) abort one process at a time until the deadlock cycle is eliminated.  
Criteria for termination  
(-priority      -duration of computation  
-no. and type of resources allocated.      -<sup>type</sup> no. of processes.

### ii) Resource preemption:

We preempt some resources from process and give this to other processes until the deadlock cycle is broken.

Preemption of resources should follow these issues.

- a.) selecting a victim: we must determine which resources and which processes are to be preempted so that it will minimize the cost.

- b.) rollback: since the resources are preempted from a process, it cannot continue its normal execution thus we must rollback the process to some safe state ~~or~~ and restart it from that state.

Generally it is difficult to determine a safe state so we go for total roll-back.

c) starvation: to avoid starvation we should take care that resources will not always be preempted from same process. we must ensure that a process can be picked as a victim only for a finite no. of times.

### Deadlock complete

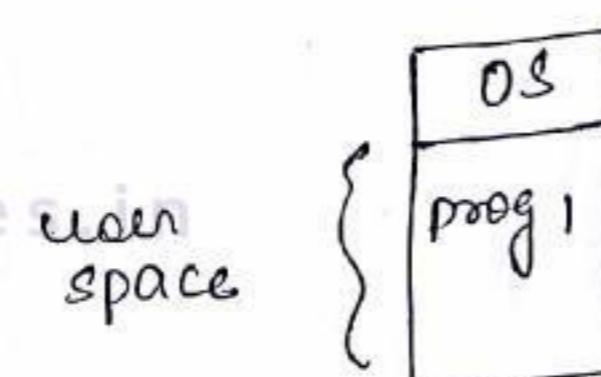
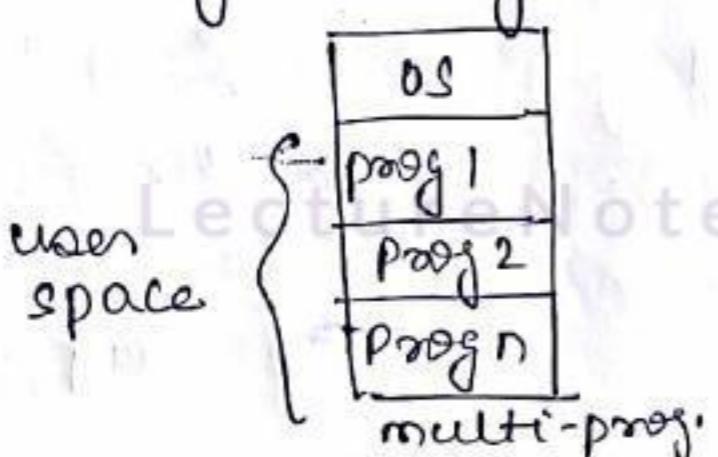
LectureNotes.in



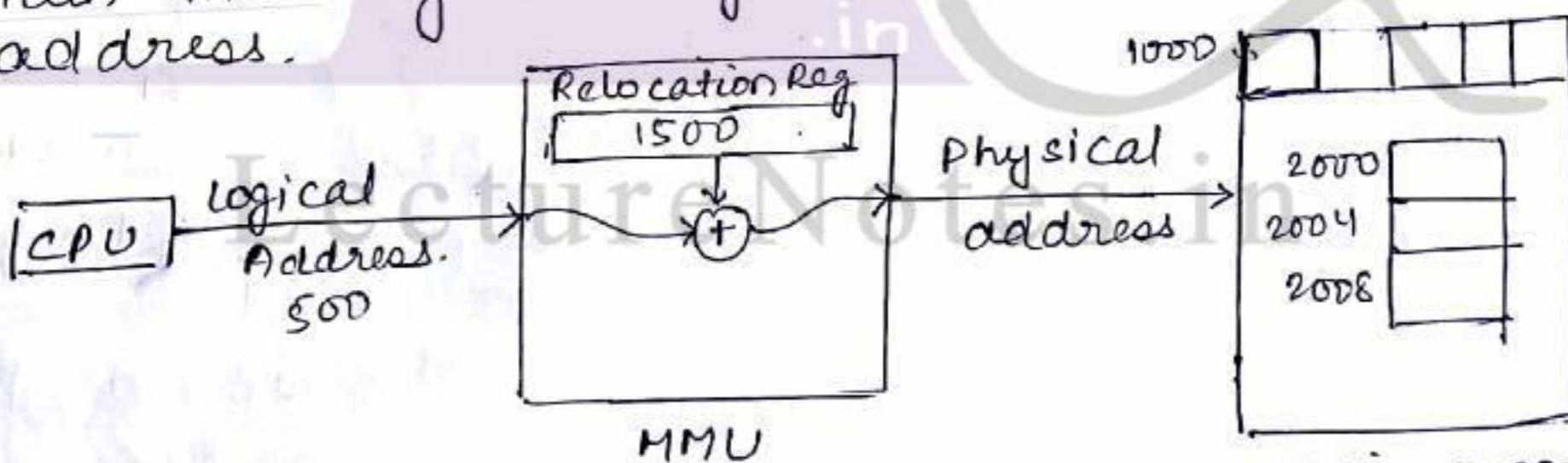
LectureNotes.in

LectureNotes.in

- In uniprogramming main memory contains OS and a single program.
- In multi-programming main m/m contains OS and multiple programs.
- The dynamic partition of memory by the OS is called memory management.



- memory is a large collection of semi-conductor cells which is capable of storing 1 bit of information either 0 or 1.
- Address binding:  
 1. Symbolic address      2. Absolute address      3. Relocatable address  
address found in Prog. at the time of compilation at the time of execution
- logical address space vs critical physical address space  
 - CPU generates one address called logical or virtual address.  
 - Main memory having some address called physical address.



$$\text{Physical Add} = \text{logical add} + \text{content of Relocation Register}$$

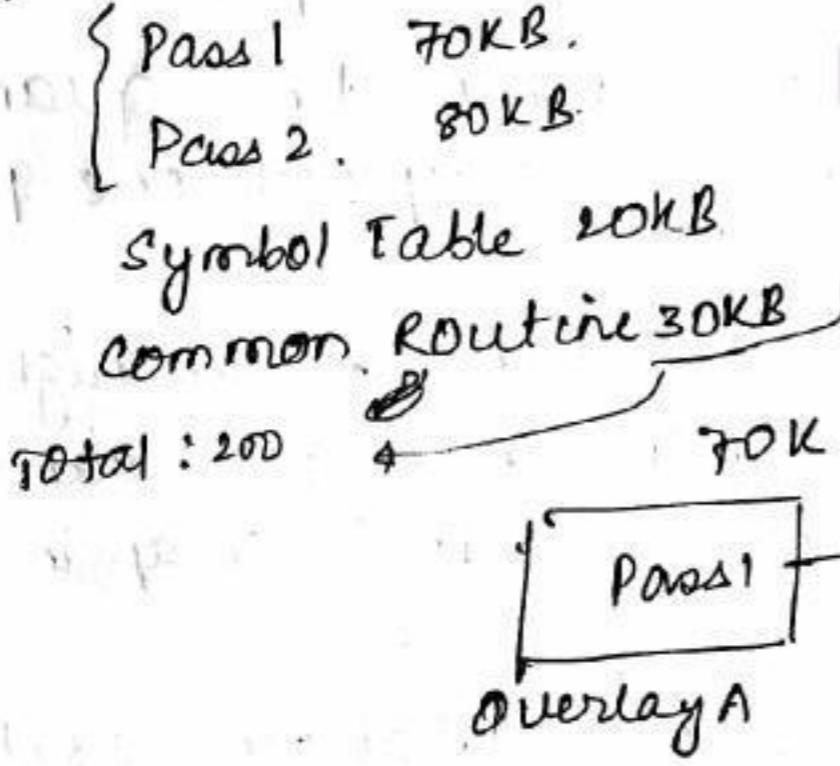
- dynamic loading
  - loading means load the program from secondary memory to main memory.
  - Types
    1. static or compile time loading
    2. dynamic or run time loading.
  - If all the loading routines (functions) are loaded in the main memory at the time of compilation

- is called static loading.
- If all the load routines are loaded in the main memory at the time of execution then it is called dynamic loading.
- In dynamic loading, routines are not loaded until it is called.
- Advantage (dynamic loading)
  - i) an unused routine is never loaded so we can obtain better memory utilization.
- Dynamic linking
  - Linking of library files before execution time is called static linking.
  - Linking library files at the time of execution is called dynamic linking.

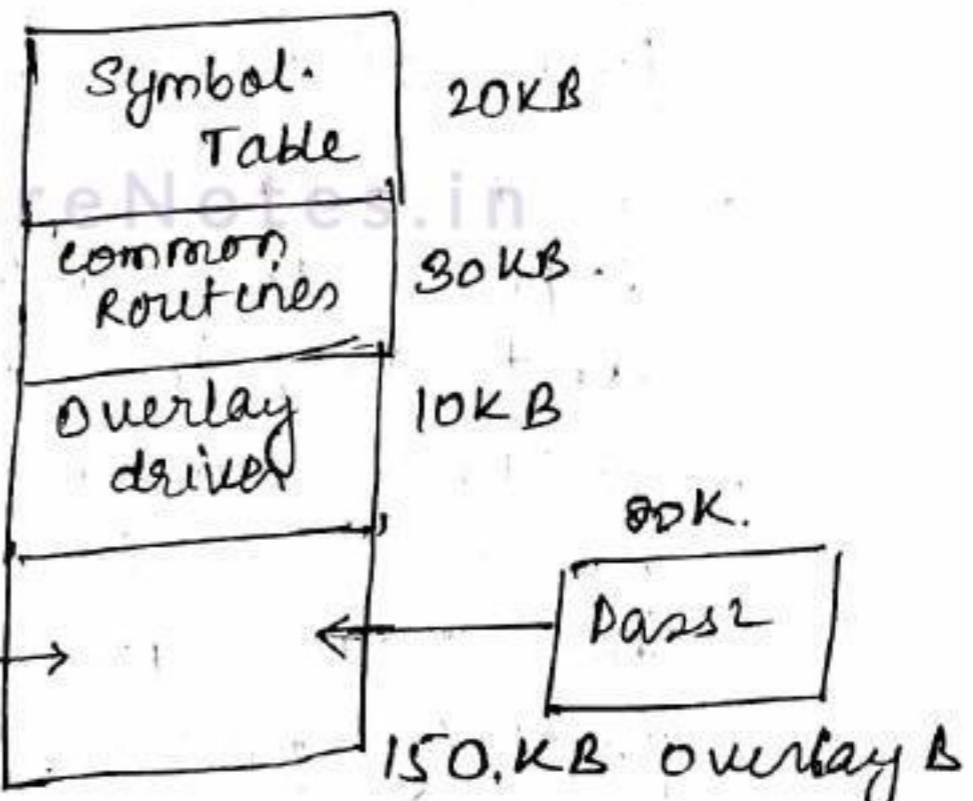
### Overlays

- To enable a process to be larger than the amount of memory allocated to it then we use overlays.
- The meaning of overlays is to keep in memory only those instructions or data that are needed. and other instructions loaded after the other other.
- when other instructions are needed, they are loaded into m/m space occupied previously by the instructions which are no longer needed.

12/10/17



Overlay for 2-pass assembler



To load everything at once we require  $200\text{ KB}$  of memory but if we have only  $150\text{ KB}$  available, we cannot run our process so we define two overlays.

Overlay A contains system table, common routines and Pass 1 similarly overlay B contains system table, common routine and Pass 2.

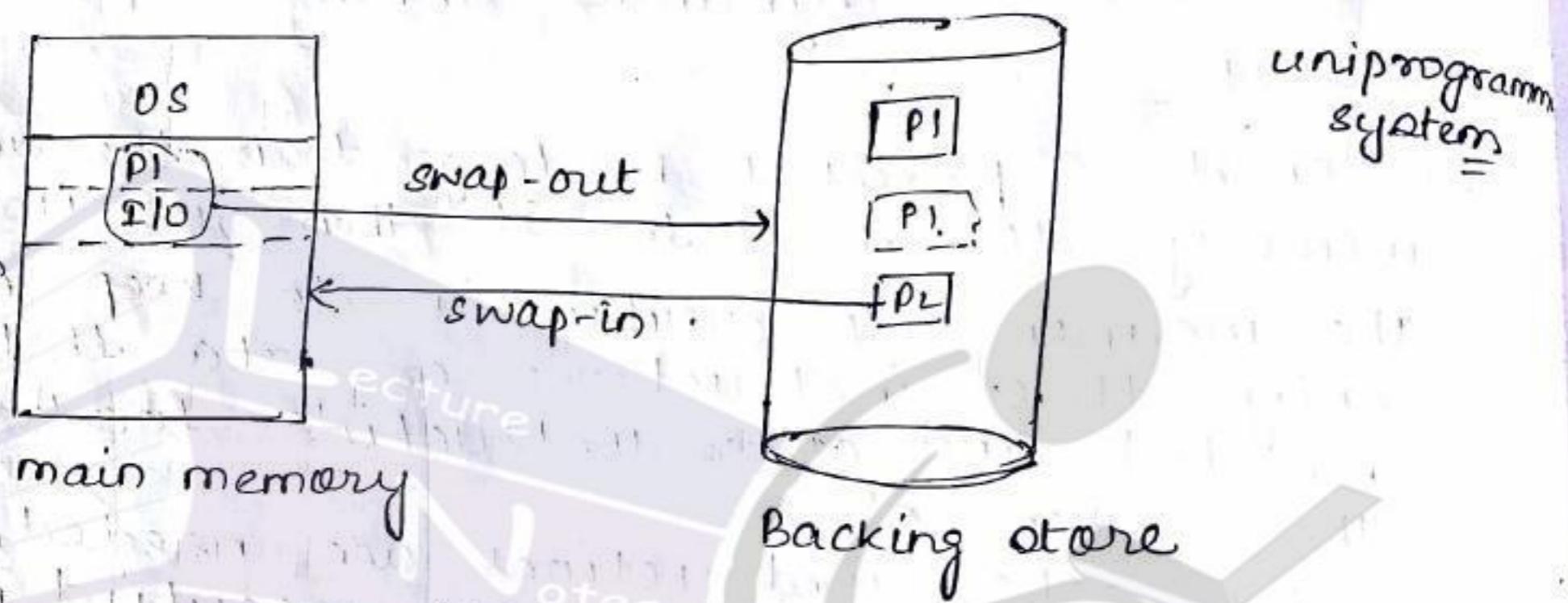
size of overlay A is  $20 + 30 + 70 = 120\text{ KB}$

size of overlay B is  $20 + 30 + 80 = 130\text{ KB}$

{Disadvantage : slower due to division}.

## • Swapping —

while executing P1 we need some I/O  
so less CPU utilization  
move it to (P1)  
backing store  
(swap-out)  
and execute P2.



- swapping is a method to improve the main memory utilization
- switching a process from main memory to disc is called swap-out and switching a process from disc to main m/m is swap-in.
- In Round Robin algorithm when time quantum expires the m/m manager swap out one process and swap-in another process
- Similarly in priority scheduling, if high-priority process arrives the m/m manager can swap-out the lower priority process and swap in high priority process.
- It is also sometimes called roll-in-roll-out.

## contiguous memory allocation

- Main memory must accommodate OS and user-process
- Memory protection:> protecting the OS from the user process and protecting one user process from another user-process is called memory protection
- > we can provide protection using 2 registers : relocation register and limit register
- > relocation register is base register and
- > limit register contains range of logical address.

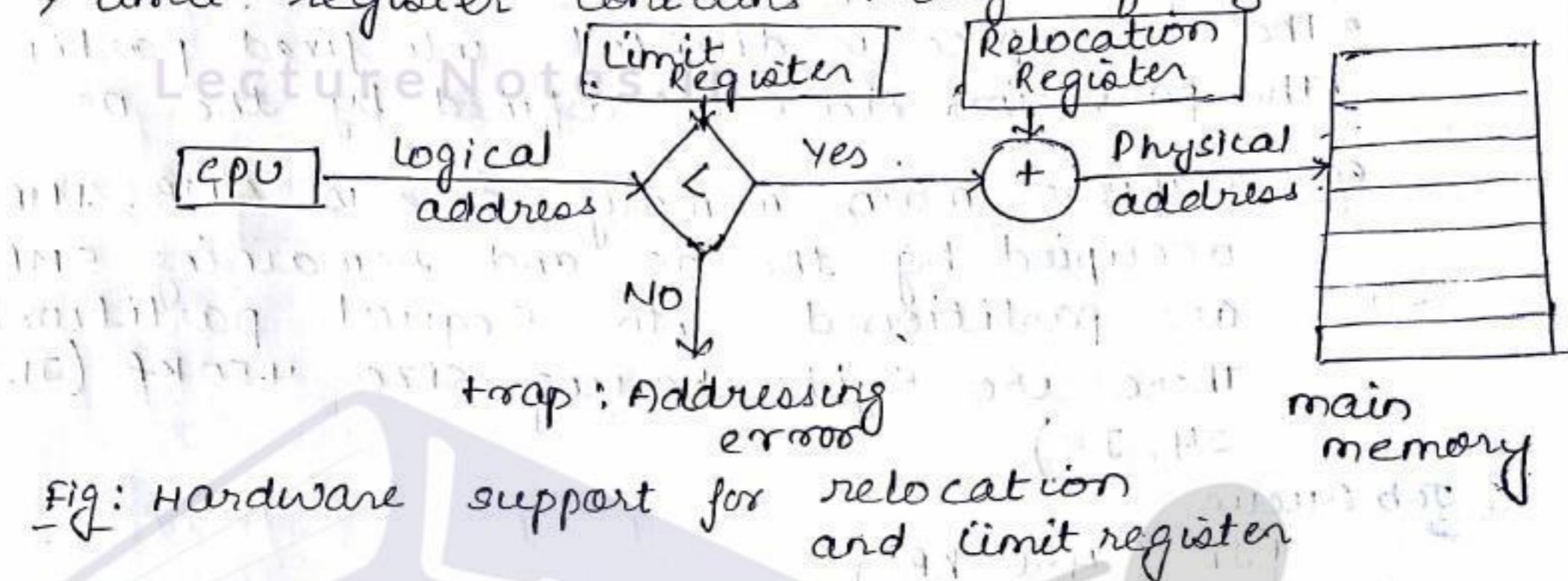
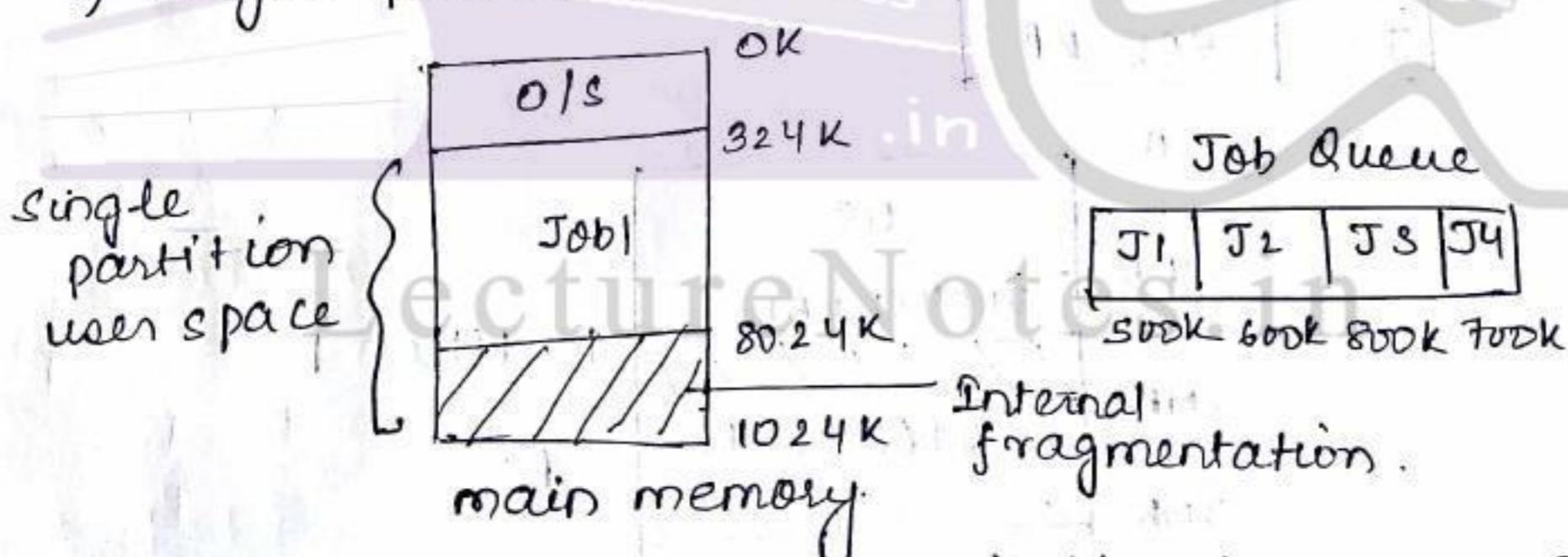


Fig: Hardware support for relocation and limit register

## memory allocation

- > There are many methods for memory allocation.
- i) single partition allocation:



Here single partition is available for user space. Only one job can be loaded in this user space. This main memory consist of only one process at a time because user space is treated as a single partition.

advantage: it is simple

disadvantage: i) poor utilization of processor.

ii) poor utilization of memory.

iii) user's job must be limited to size of available main m/m

## ii) multiple partitions allocation.

13/10/17

This method is implemented in 3 ways!

- fixed, equal multiple partitions
- fixed variable multiple partition
- dynamic, multiple partition

### a) fixed equal multiple partitions

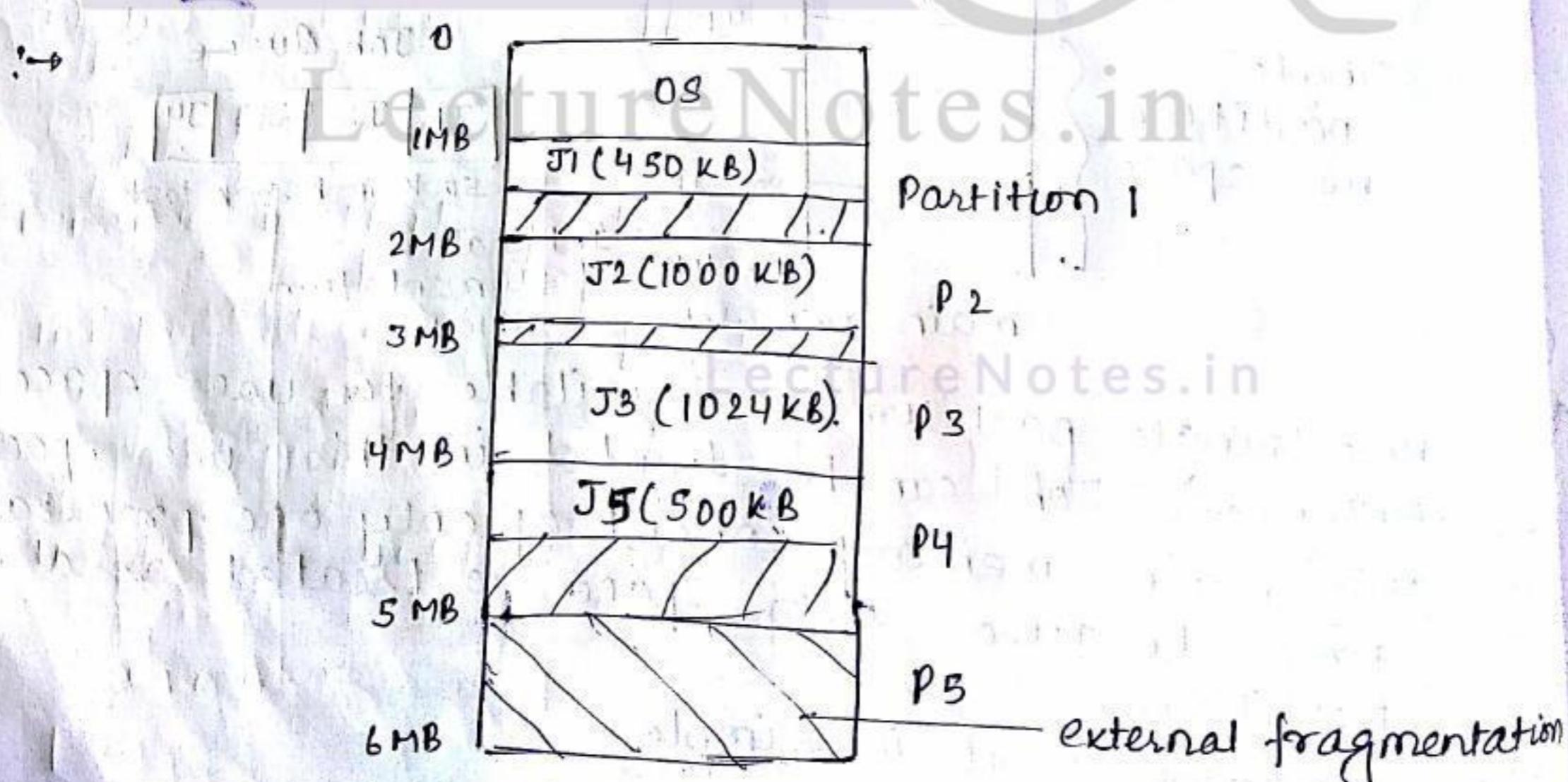
- The user space is divided into fixed partitions.
- The partition size are defined by the OS.

Eg: suppose main memory size is 6MB, 1MB occupied by the OS, and remaining 5MB are partitioned into 5 equal partitions. There are 5 jobs having size 450KB ( $J_1, J_2, J_3, J_4, J_5$ )

Job Queue

$J_1$	450 KB
$J_2$	1000 KB
$J_3$	1024 KB
$J_4$	1500 KB
$J_5$	500 KB

Find the internal and external fragmentation.



$$\text{Total internal fragmentation} = (1024 - 450) + (1024 - 1000) + (1024 - 1024) + (1024 - 500) = 1122 \text{ KB}$$

{ total space is  $(1122 + 1024) > 1500$  but still cannot be allocated as contiguous space is needed.

### advantage

- i.) supports multiprogramming.
- ii.) efficient utilization of CPU.
- iii.) simple and easy to implement

### disadvantage

- i.) supports internal and external fragmentation.
- ii.) The single free area may not be large enough for partitioning a single partition.

### b) fixed variable multiple partition

In this technique the user space of main memory is divided into no. of partitions but the partition size are different length.

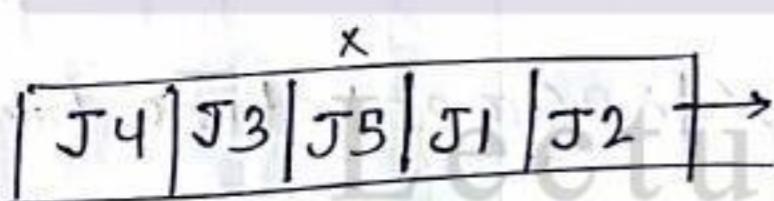
eg :-

Job Queue

Job	Size (KB)	Arrival Time
J1	825	10
J2	600	5
J3	1200	20
J4	450	30
J5	650	15

Partition	Size (KB)
P1	700
P2	400
P3	525
P4	900
P5	350
P6	625

Find internal and external fragmentation using this technique.



Job Queue  
According to arrival time

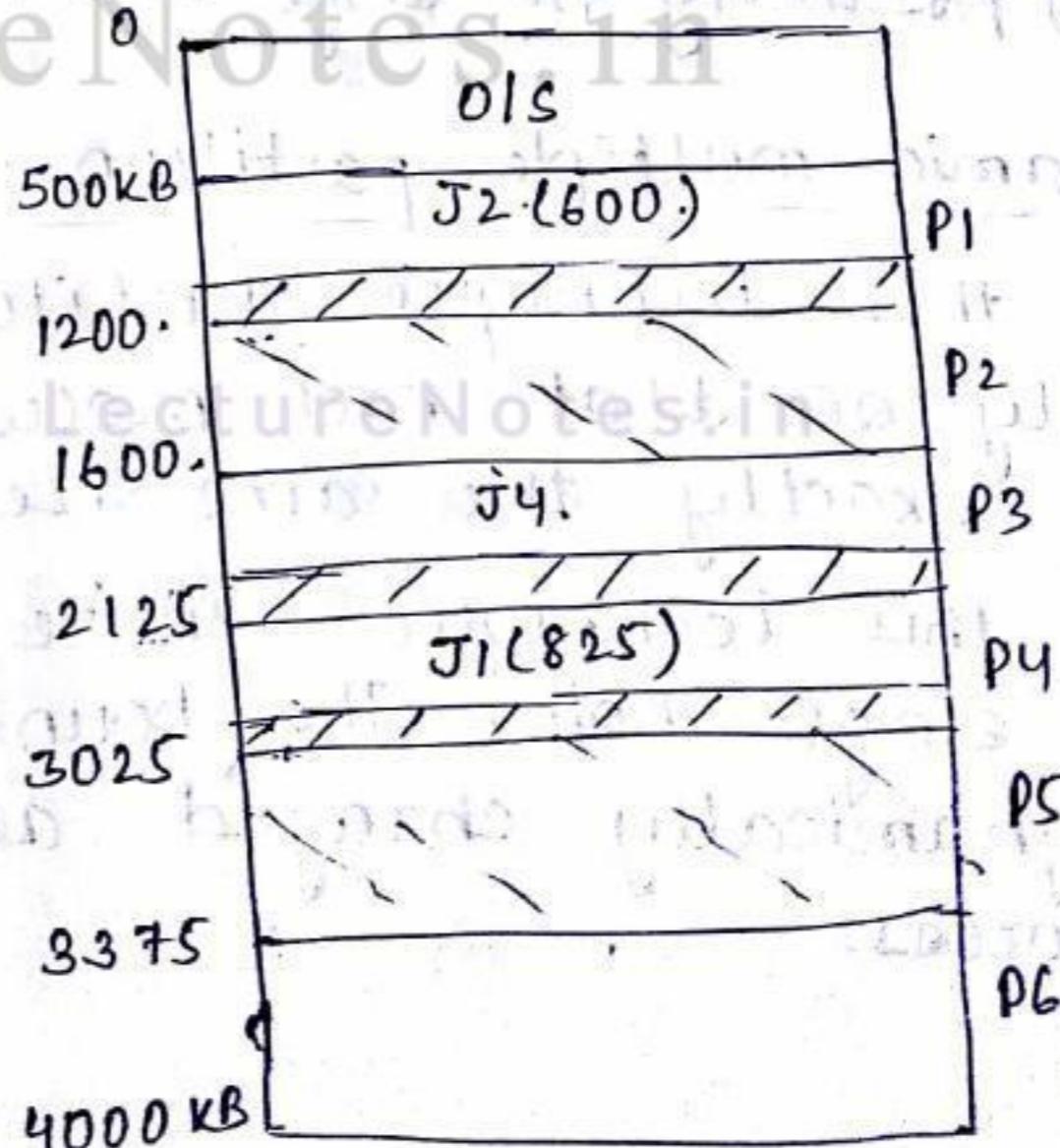
Total internal fragmentation

$$= (700 - 600) + (900 - 825)$$

$$+ (525 - 450) \\ = \underline{250 \text{ KB}}$$

Total external fragmentation

$$= 400 + 350 + 625 \\ = \underline{1375 \text{ KB}}$$



problem: the size of  $J_2$  is 600 KB, it is loaded in partition 1. so internal fragmentation is 100 KB. But if job 2 is loaded in partition 6 then internal fragmentation is 25 KB.

To avoid this problem we use

i.) first-fit algorithm: allocate the partition that is big enough starting from the first till the end. we can stop searching as soon as we find a free hole (partition) that is large enough.

ii.) best-fit algorithm: allocate the smallest hole or partition that is big enough.

iii.) worst-fit algorithm: allocate the job to the largest hole.

advantage of fixed variable multiple partition

- i.) It supports multiprogramming.
- ii.) efficient utilization of CPU and memory. (best & first fit)
- iii.) simple and easy to implement

disadvantage:

- i.) This technique supports internal and external fragmentation.
- ii.) possibility for large external fragmentation (worst-fit)

### c) Dynamic multiple partition

- In this technique partitions are created dynamically so that each process is loaded into partition of exactly the same size as that of process size.
- In this technique entire user space is treated as a single hole. The boundaries of partition are dynamically changed and depend on size of the process.

eg:

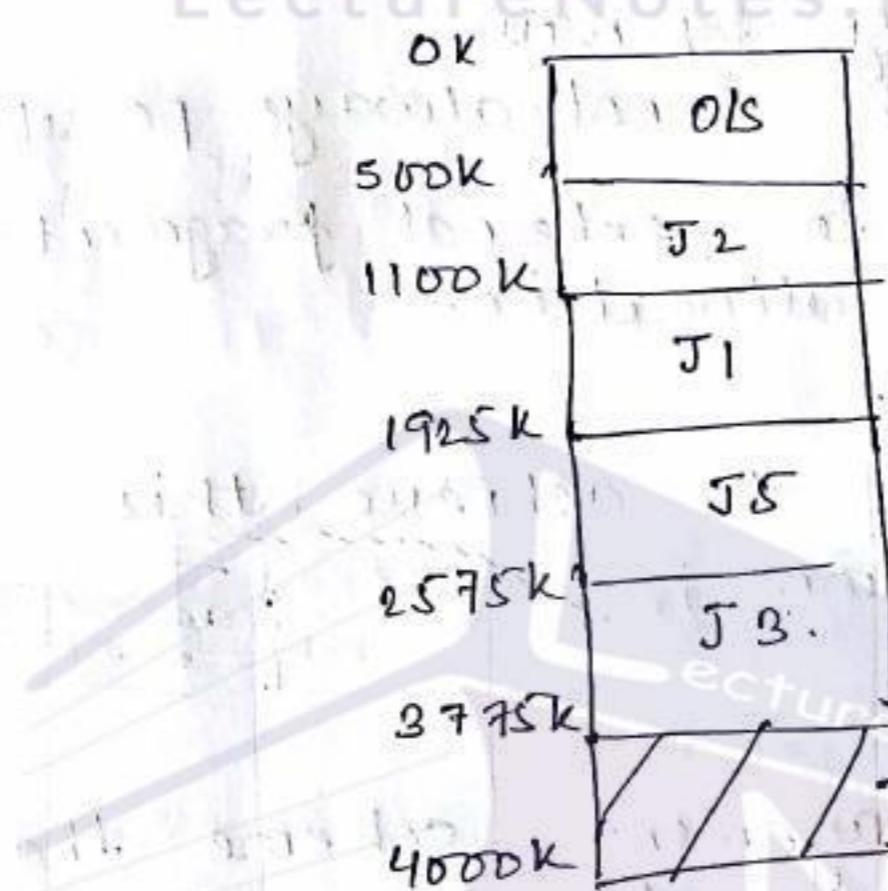
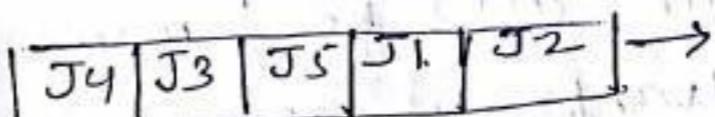
partition  
Job 2  
station  
an  
ree

19/10/17

### Job Queues

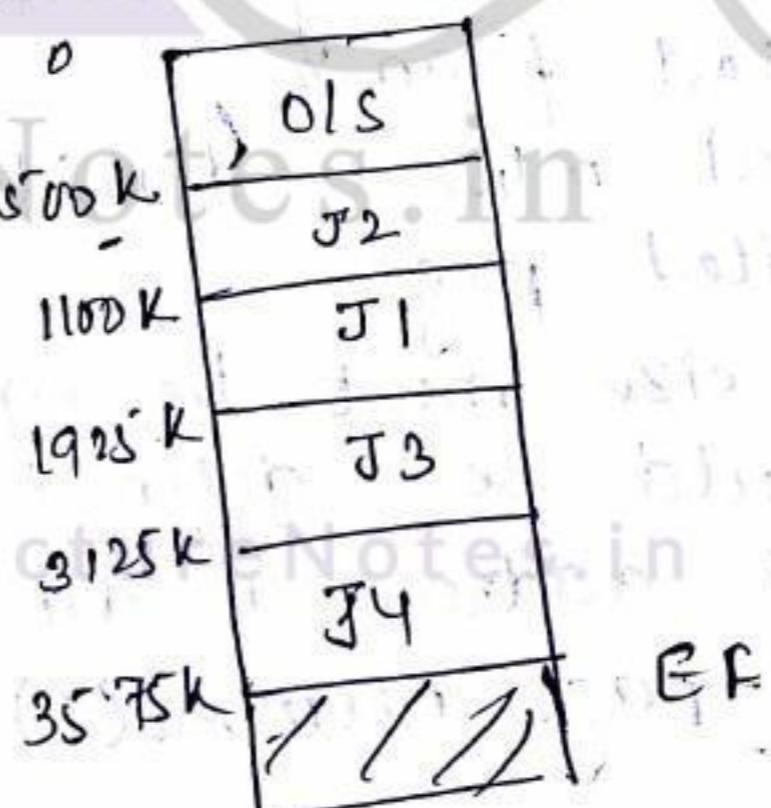
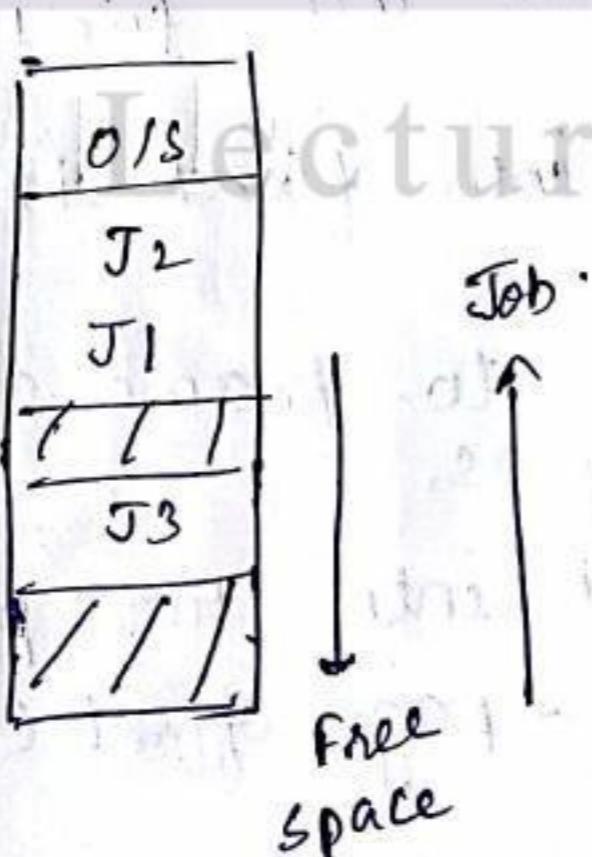
Job	size	Arrival Time
J1	825	10
J2	600	5
J3	1200	20
J4	450	30
J5	650	15

### Job Queue



After all job allocation, only 225KB of memory space is available so we cannot allocate J4 into it.

Suppose after sometime J5 completes its operation and release its memory. so total free space =  $225 + 650 = 875K$



### Advantage:

- It does not support internal fragmentation because partitions are dynamically changed.
- efficient utilization of mem and CPU.

### Disadvantage

- It supports external fragmentation.

### → compaction .

solution to the problem of external fragmentation is called compaction. The goal of compaction is to place all free memory space in one large block.

advantage: It improves memory utilization.

disadvantage: It improves performance.

expensive.

is not always possible.

Note\*1: Another solution to external fragmentation is non-contiguous allocation.

### → Non-contiguous allocation :-

There are two techniques to achieve this

- i) paging
- ii) segmentation

{ address of ram - Phys  
" " process - logic

### • Paging .

Paging is a mm management scheme that permits physical address space of a process to be non-contiguous.

Physical memory is divided into fixed size blocks called frames.

logical mm is divided into fixed size blocks called pages.

Page size must be equal to frame size which should be defined by OS.

logical address is divided into two parts

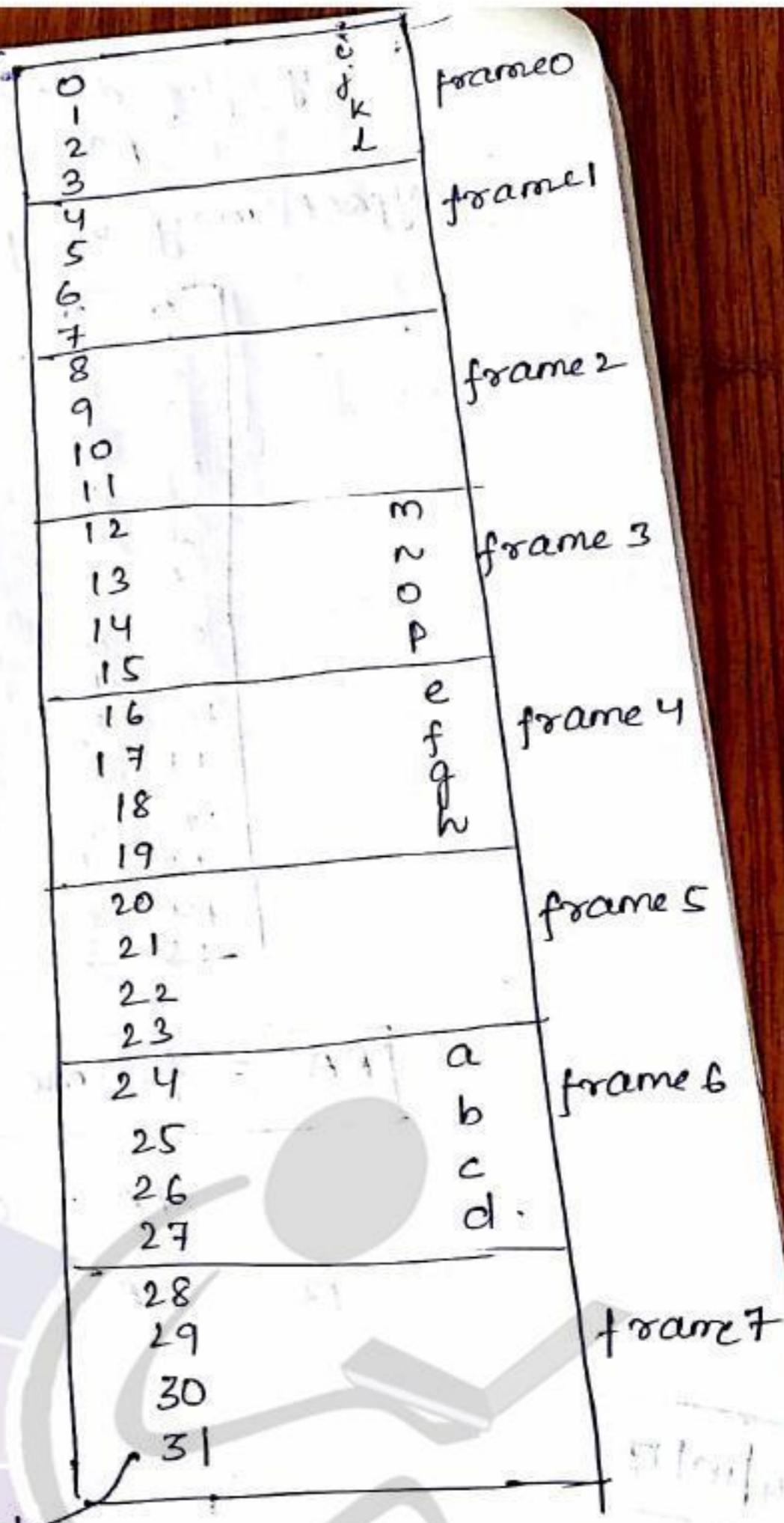
- page number( $P$ )

- page offset( $d$ )

0	a
1	b
2	c
3	d
4	e
5	f
6	g
7	h
8	i
9	j
10	k
11	l
12	m
13	n
14	o
15	p

logical address

logical memory.



$$\text{Page no.} = \frac{\text{Logical address}}{\text{page size}}$$

$$\text{offset} = \frac{\text{logical address} \% \text{Page size}}{\text{Page size}}$$

→ Page 0 stored in frame 6.

0	6
1	4
2	0
3	3

Page Table

- Page number is used as an index into the page table
- Page table contains the base address of each page of the physical memory

if  $\log_2 \text{add} = 9$   
 $\text{page no.} = \frac{9}{4} = 2$   
 offset may vary from  $0 \text{ to } 3$

0	0	a
1	1	b
2	2	c
3	3	d
4	0	e
5	1	f
6	2	g
7	3	h
8	0	i
9	1	j
10	2	k
11	3	l
12	0	m
13	1	n
14	2	o
15	3	p

$$\text{offset} = 9 \% 4 = 1$$

eg. 15  
 $\text{log. page no.} = \frac{15}{4} = 3$   
 $\text{offset} = \frac{15 \% 4}{4} = 3$

$\text{PA} = \text{frame no.} \times \text{page size} + \text{offset}$

$$\text{PA} = 0 \times 4 + 1 = 1$$

$$\text{PA} = \cancel{3 \times 4 + 1} = 13$$

$$3 \times 4 + 13 \% 4$$

$$= 12 + 1 = 13$$

24/10/17

→ Paging hardware :-

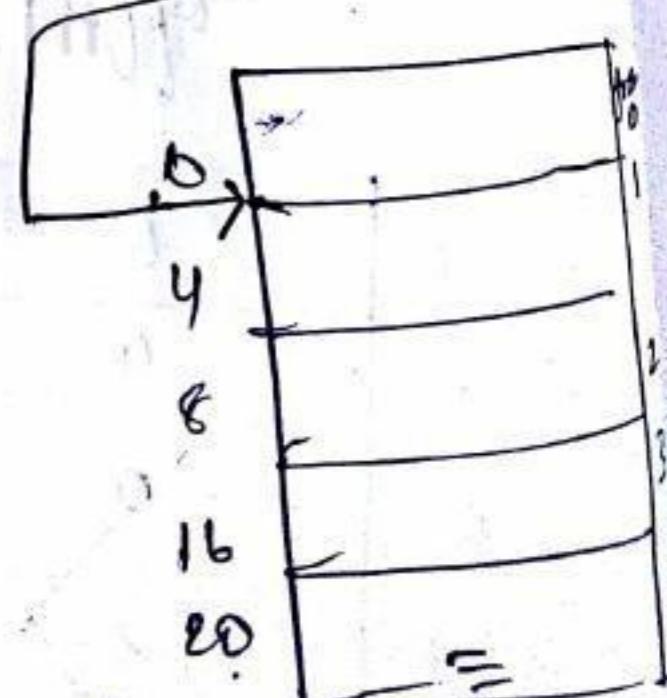
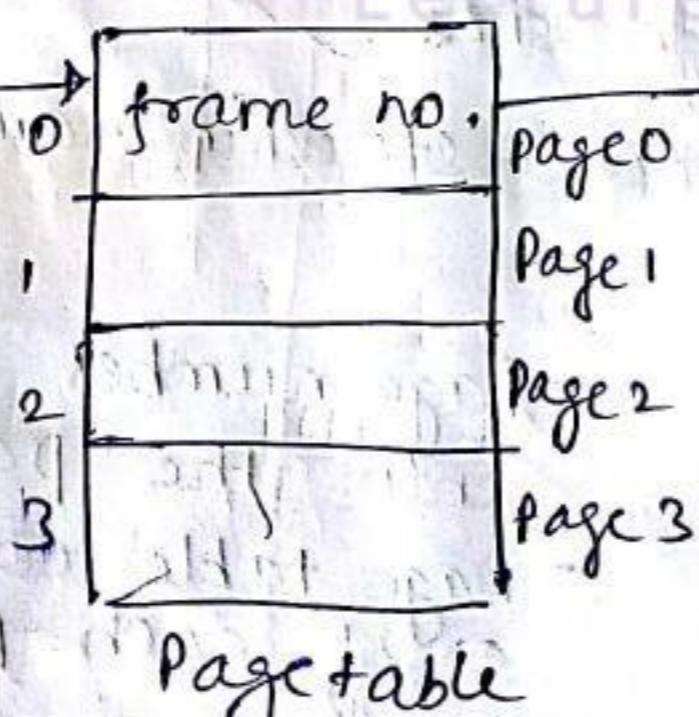
Logical Add.

Page NO.(P) | offset(d)

frame NO. | offset(d)

Physical Add.

CPU



Paging Paging hardware

Physical memory

## Advantage of Paging :

- i) no external fragmentation but we may have some internal fragmentation.
- ii) It supports <sup>time</sup> sharing system.
- iii) It supports virtual memory.
- iv) sharing of common code is possible.
- v) To avoid internal fragmentation small page size are desirable.

## Disadvantage:

- i) It supports some amount of internal fragmentation.
- ii) If the no. of pages are high it is difficult to maintain the page table.

## Hardware support of page table :-

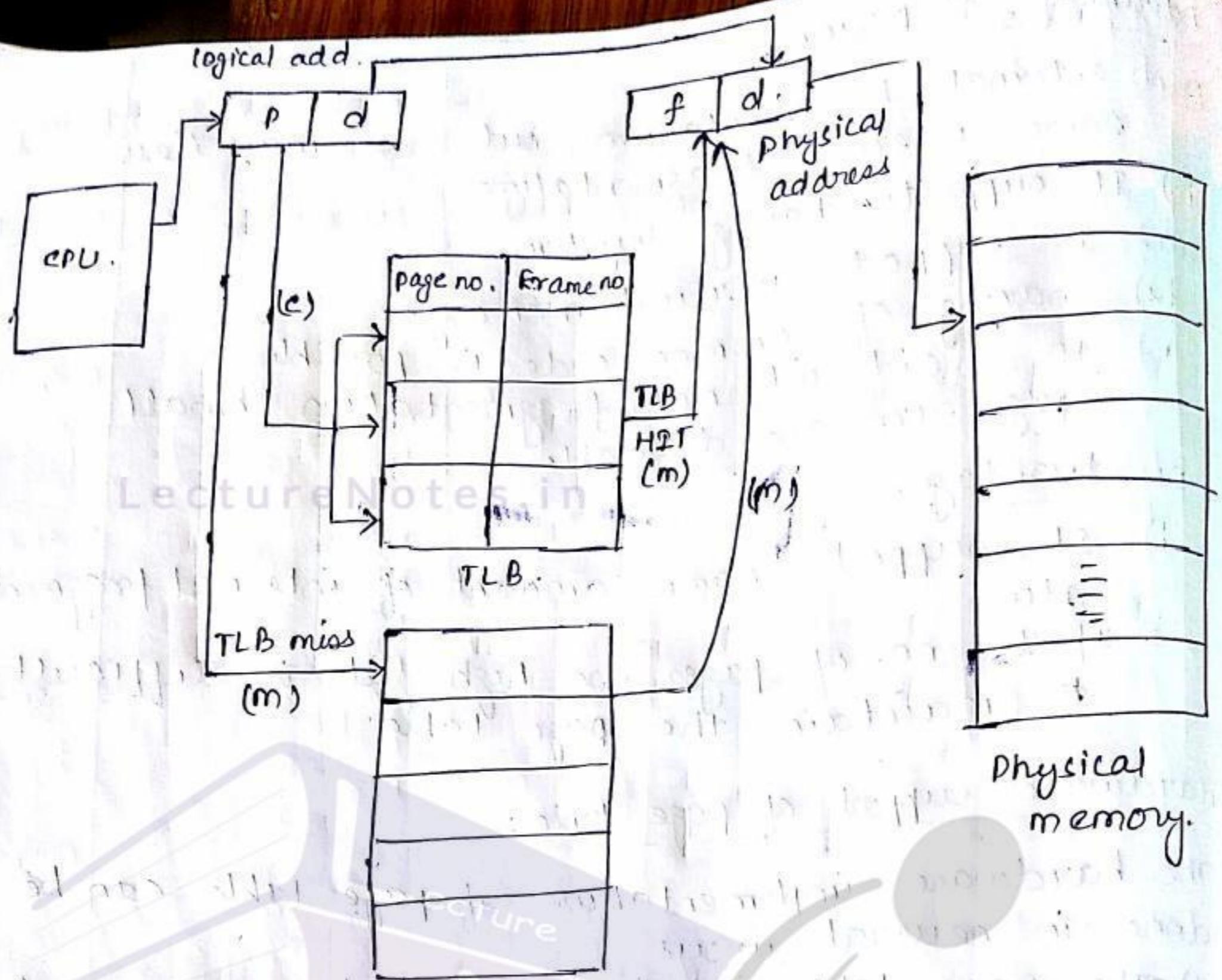
The hardware implementation of page table can be done in several ways

- i) The page table is implemented by using a set of a dedicated register.  
advantage : accessing the page table will be faster.  
disadvantage : design of large page table using registers are not feasible.
- ii) keep the page table in main memory.  
advantage : we can create a large size page table.  
disadvantage : it takes more time to access the memory location (by a factor of  $^2$ ).

- iii) The solution to this problem is translation look aside buffer (TLB) : It is a special small, fast hardware called TLB.

This TLB is a associative memory or high speed memory. (content addressable mem (CAM) principle)

eg: when you search some part of name of a file and file comes )



let  $h$  be the TLB hit ratio,  $c$  is the cache access time and  $m$  is memory access time then effective memory access time =

$$\text{EMAT} = h * (c + m) + (1 - h) * (c + 2m)$$

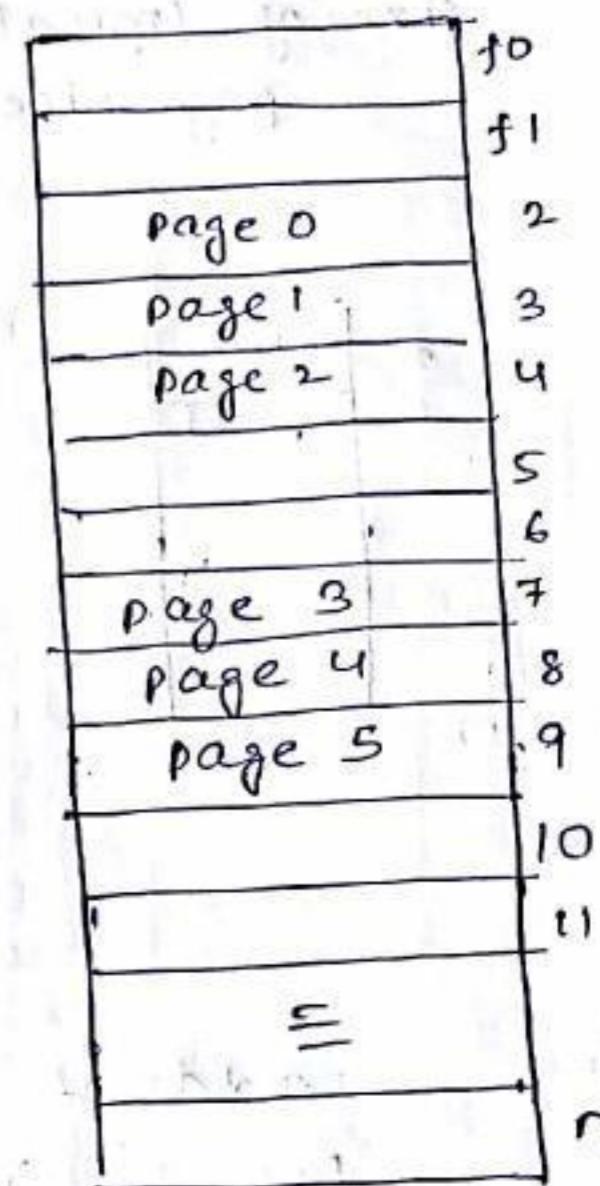
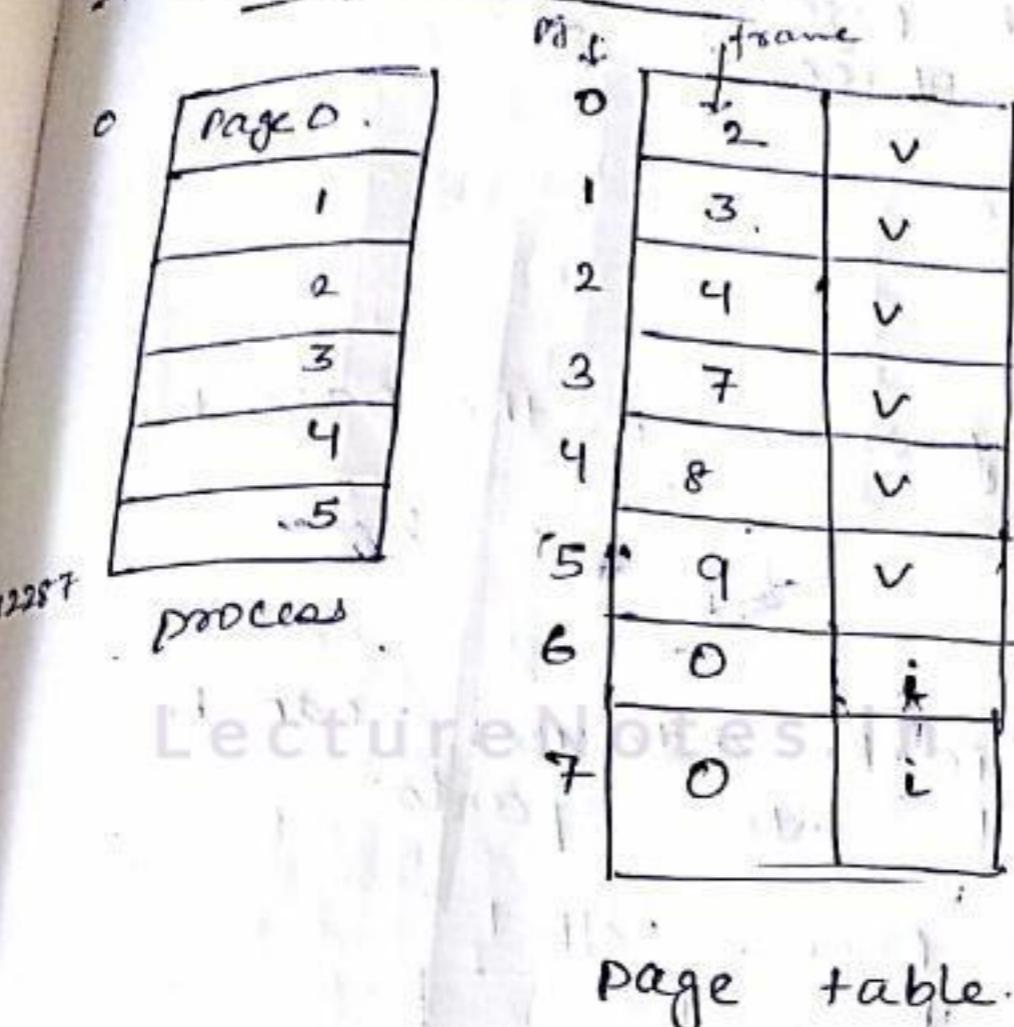
Q: A TLB having 80% hit ratio 20 nanosec to search TLB, 100 ns to access the memory. Find the effective memory access time.

$$h = 0.8 \quad c = 20 \text{ ns} \quad m = 100 \text{ ns}$$

$$\text{EMAT} = 0.8 \times (120) + 0.2 \times (20 + 200)$$

$$\begin{aligned} &= 96 + 44 \\ &= 140 \text{ ns} \end{aligned}$$

→ memory protection.

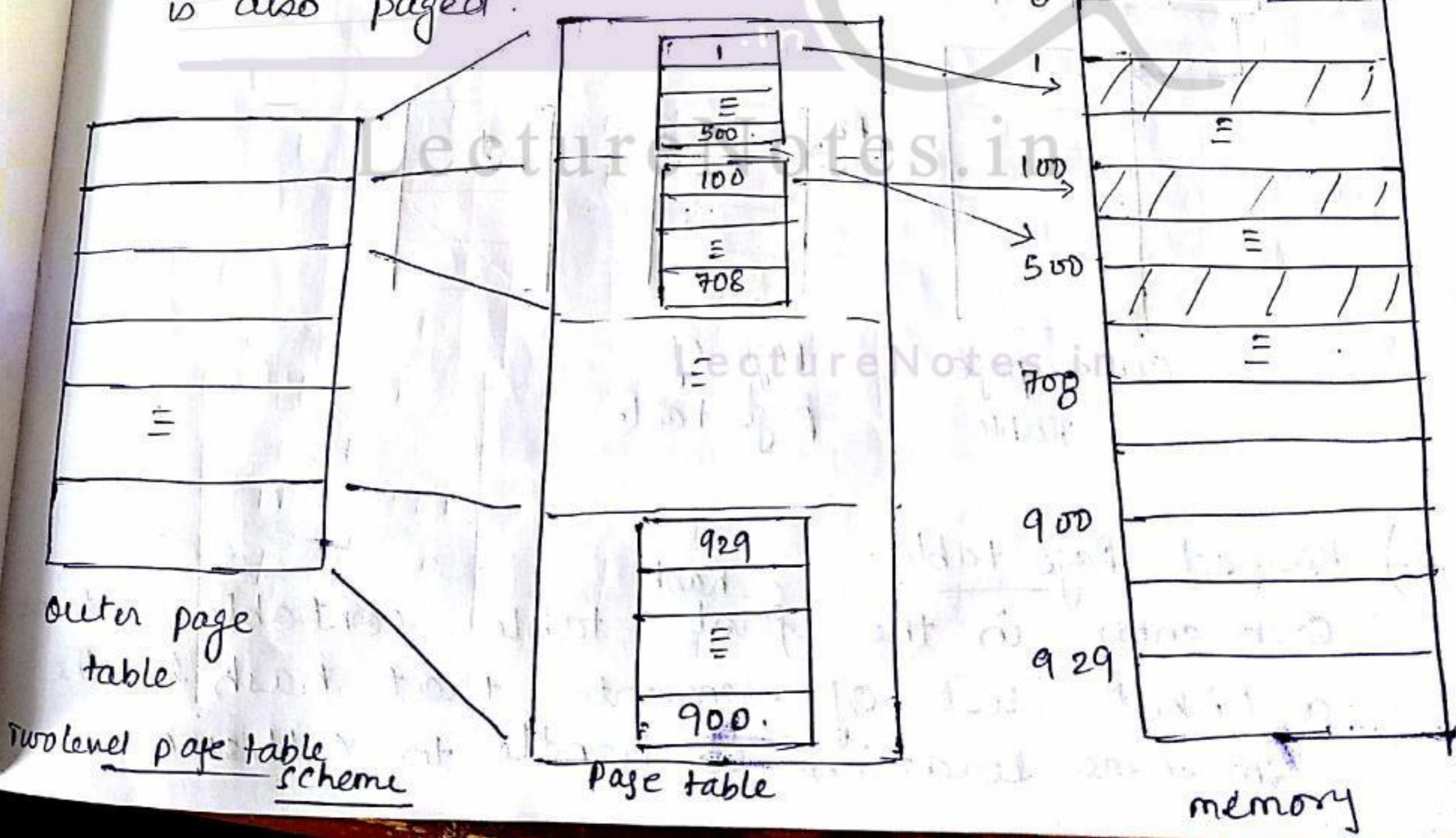


physical memory.

→ structure of the page table :-

i) Hierarchical paging.

In two level paging system, the page table itself is also paged.



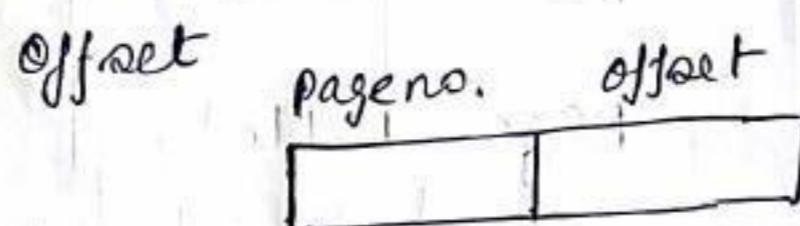
e.g.: A 32 bit machine with page size of logical address space is  $2^{32}$  size of 4 KB.

$$\begin{aligned}\text{page size} &= 4 \text{ KB} \\ &= 2^2 \times 2^{10} \\ &= 2^{12}\end{aligned}$$



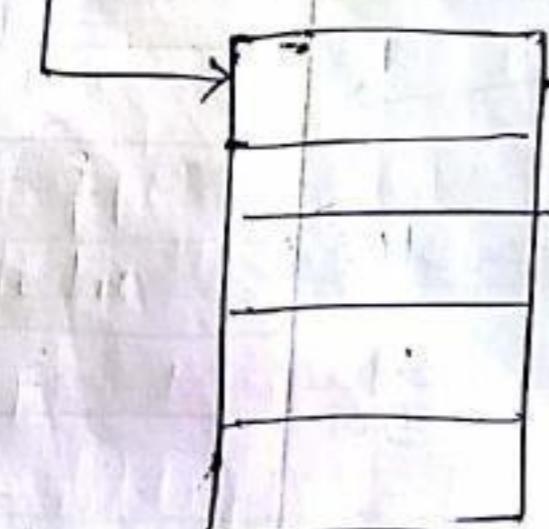
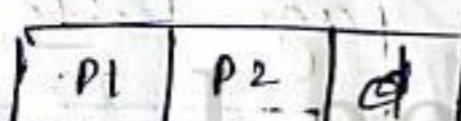
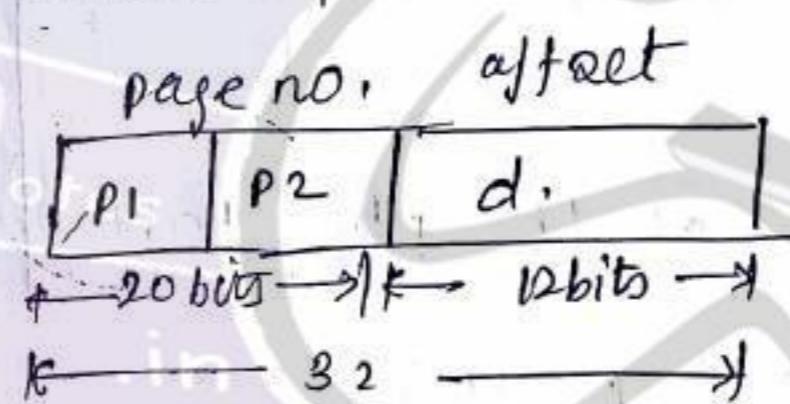
$$\begin{aligned}\text{No. of entries in the page table} \\ = \frac{2^{32}}{2^{12}} = 2^{20}\end{aligned}$$

This logical address can be divided into 2 parts



to represent 2 addresses - 21 bit

for  $2^{12}$  add → 12 bits reqd.

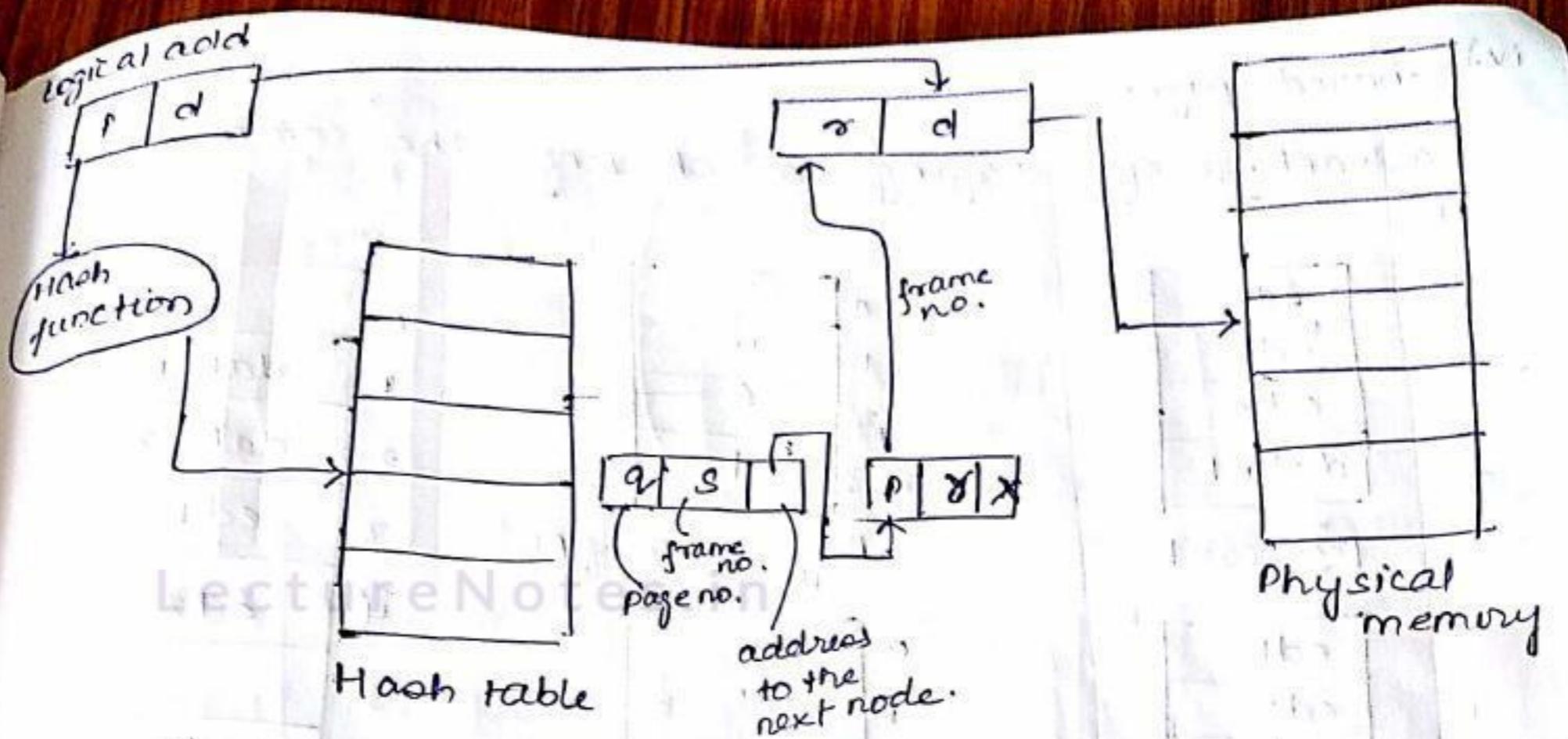


Outer page  
Table

Page of  
page table

## ii.) Hashed Page table.

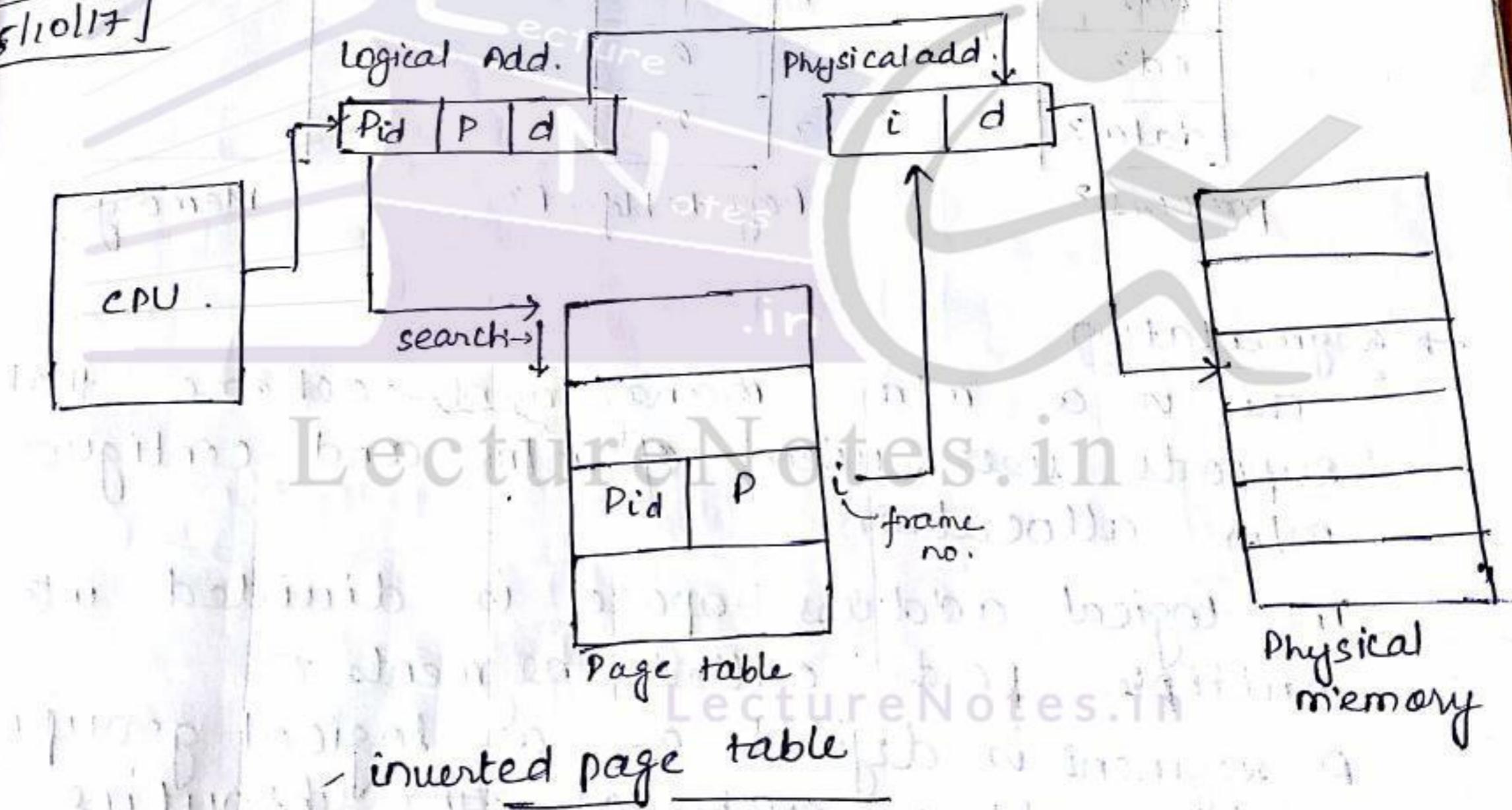
Each entry in the first table contains a linked list of elements that hash to the same location to handle collision.



### iii) Inverted page table-

each process has a page table associated with it.  
An inverted page table has one entry for each page or frame

25/10/17]



In this case only 1 page table is available in the system. Each virtual address consists of a triplet  $\langle$  Process id, page no., offset  $\rangle$ . In each inverted page table entry contains  $\langle$  process id, page no.  $\rangle$

#### iv) shared pages

advantage of paging is sharing the common code

ed1
ed2
ed3
data1

process 1

ed1
ed2
ed3
data2

process 2

ed1
ed2
ed3
data3

process 3

0	3
1	4
2	6
3	1

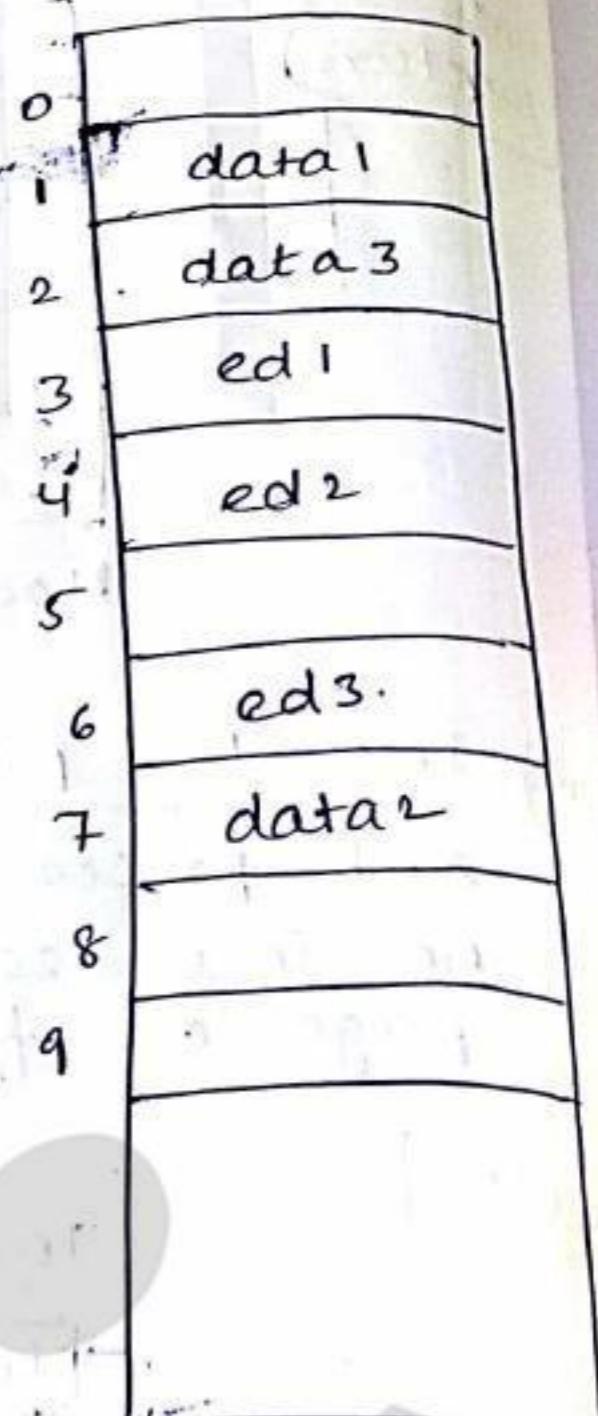
Page table of P1

0	3
1	4
2	6
3	7

Page table for P2

0	3
1	4
2	6
3	2

Page table for P3



Memory

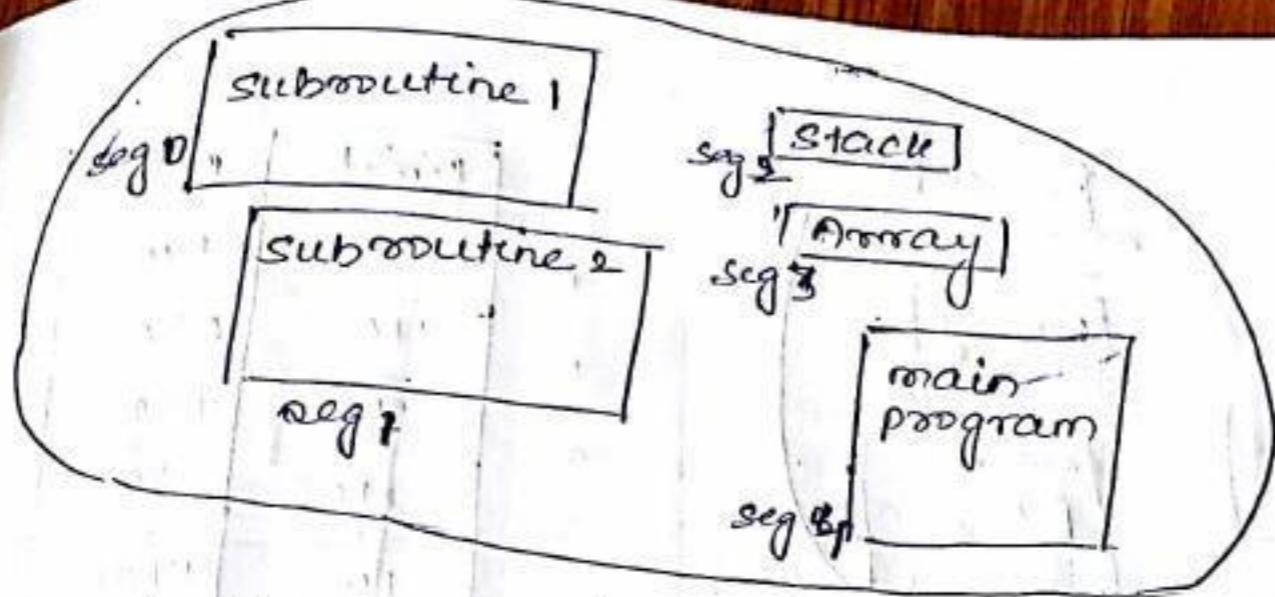
#### → segmentation

This is a m/m management scheme that supports user view of m/m and contiguous m/m allocation.

The logical address space is divided into multiple parts called segments.

A segment is defined as a logical grouping of instructions such as ~~subg~~ subroutine, data area, stack, array, etc.

Each segment has a name and a length. Every program is a collection of segments. Segmentation is a technique of managing these segments.



A logical address consists of two tuples <segment no., offset>.

### Segmentation hardware

A segment table is used to map two-dimensional user defined address into one-dimensional physical address.

each entry in the segment table has segment base and segment limit. Index of the segment table is the segment number.

logical address

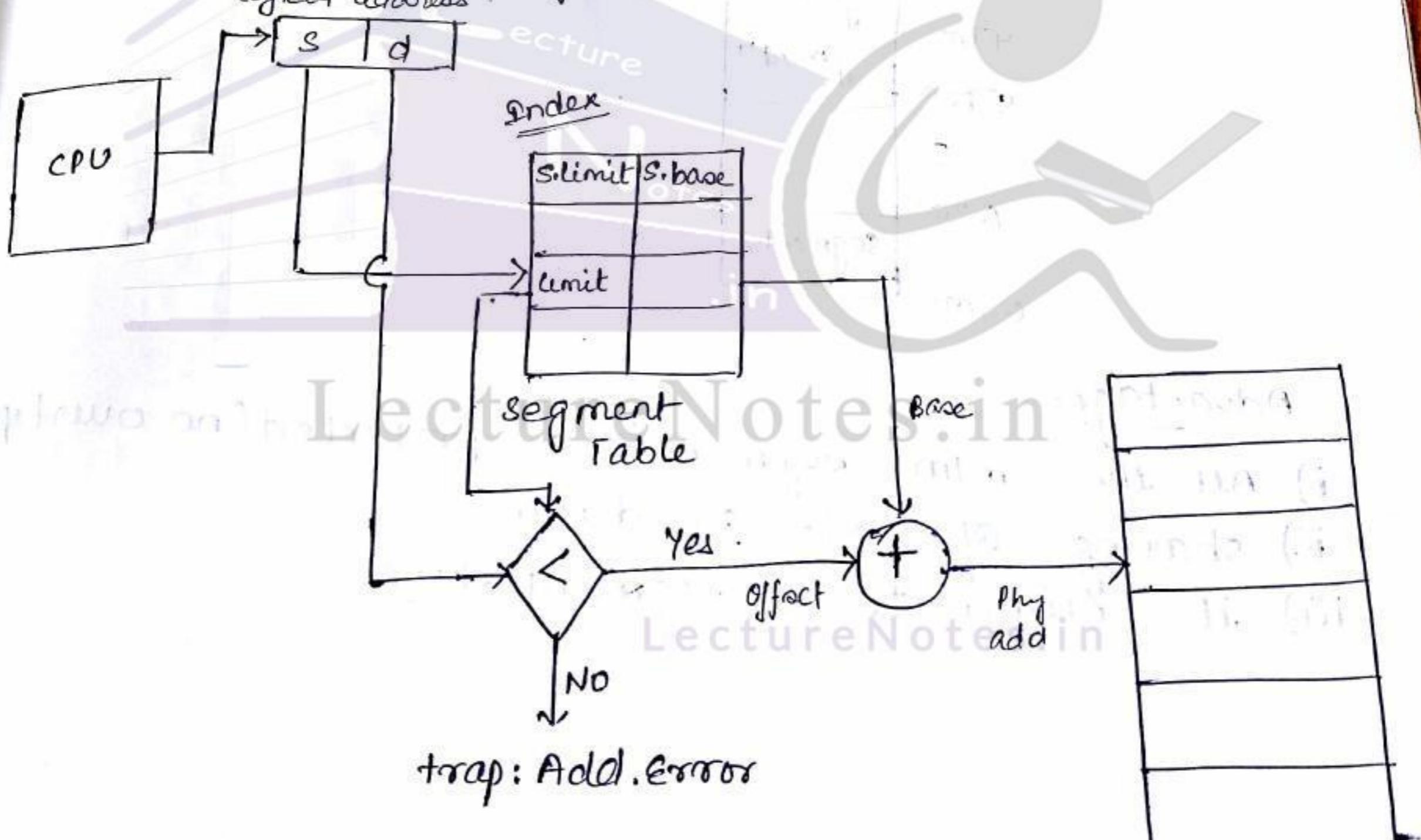
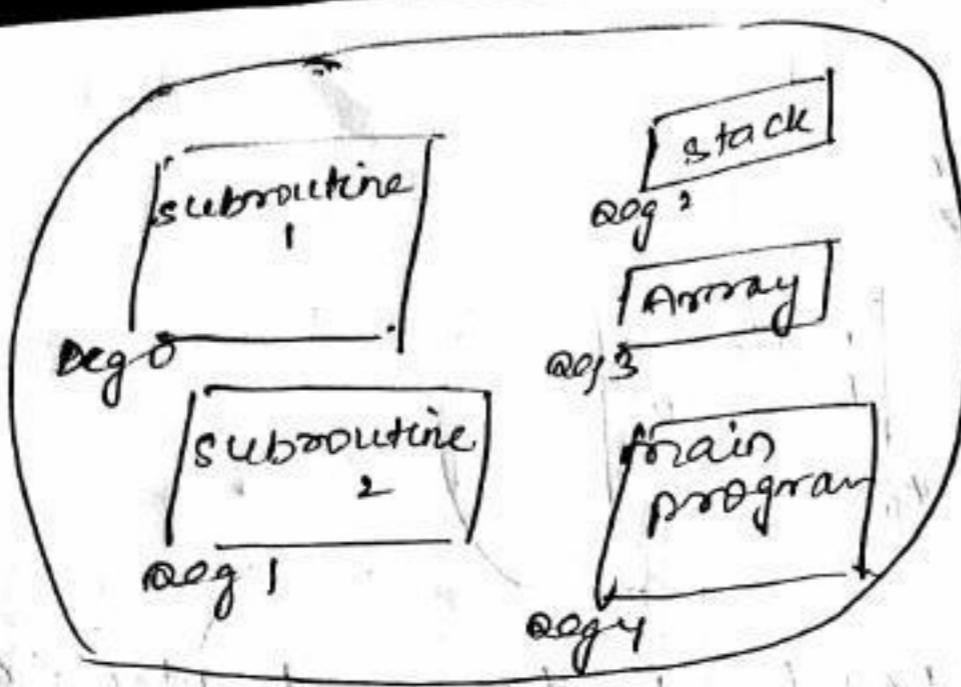


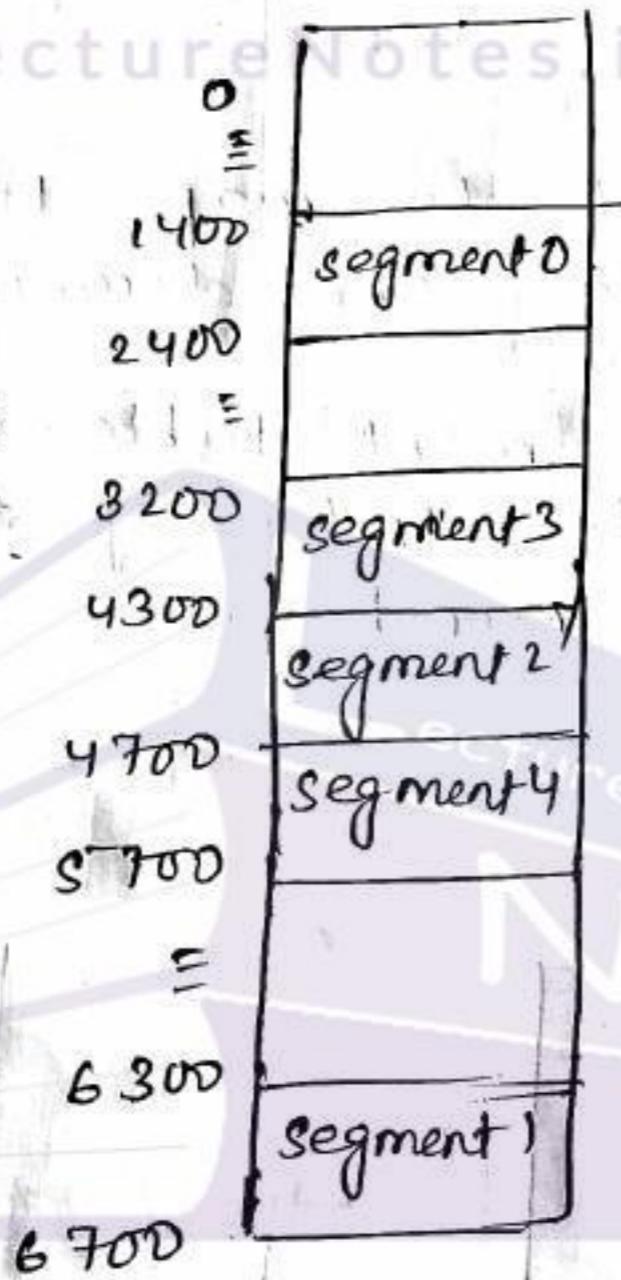
Fig: Segmentation hardware

physical  
memory



	limit	base
0	1000	1400
1	400	6300
2	400	4300
3	1100	3200
4	1000	4700

Segment table.



{ 3500 valid address  
 { 2500 invalid address

### Advantages.

- i) All the mm segments are protected (no overlapping)
- ii) sharing of code or data.
- iii) it eliminate fragmentation.

26/10/12

## VIRTUAL MEMORY.

Virtual m/m is implemented by using demand paging. It is similar to paging system with swapping. In demand paging, a page is not loaded into the main m/m from the secondary m/m until it is needed. So a page is loaded into the main memory by demand. This technique is called demand paging. Here a lazy swapper is used to swap the pages b/w main m/m and secondary m/m.

A lazy swapper never swaps a page into the memory unless that page will be needed.

### Demand Paging hardware:-

Page no.	A
0	B
1	C
2	D
3	E
4	F
5	G
6	H

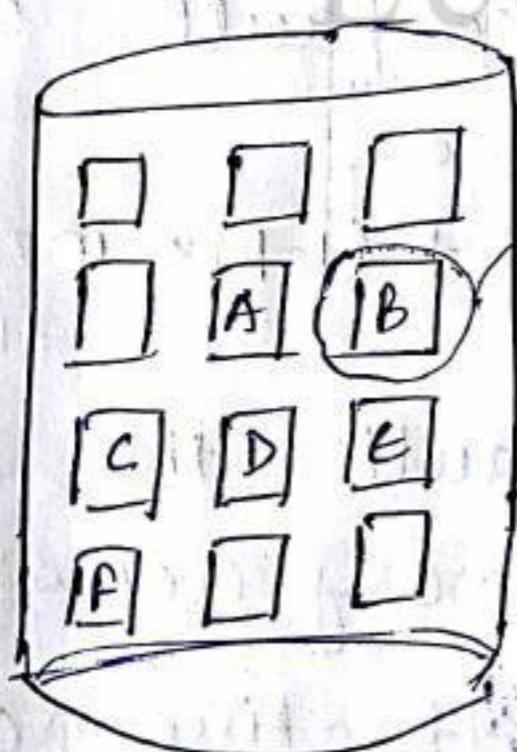
logical memory

Page no.	frame no.	valid, invalid bit
0	2	v
1	3	iv
2	4	v
3		i
4		i
5	7	v
6		i
7		i

Page Table

0	
1	
2	A
3	B -
4	C
5	
6	
7	F
8	
9	
10	
=	=

main memory.



page fault

time taken  
page fault service time

### Page Fault.

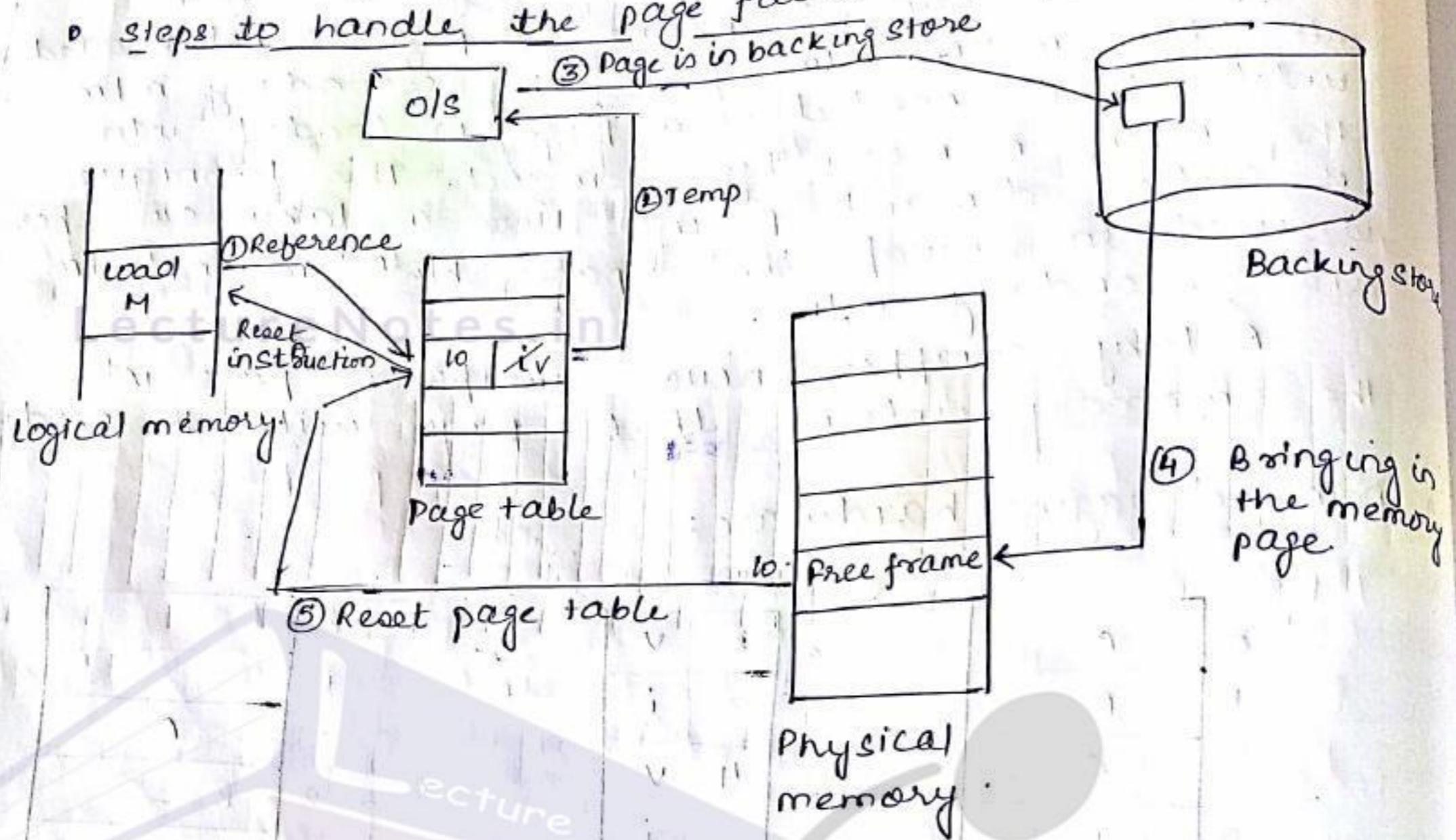
When the process needs a page, that page is not available in memory this

to execute a particular

available in memory this

situation is called page fault. It means the OS fails to bring the required page into memory.

- steps to handle the page fault



- Pure demand paging

In Pure demand paging never bring page into memory until it is required

- performance of demand paging

can be measured by computing effective access time.

let  $P$  be the probability of page fault  $0 \leq P \leq 1$   
then, the effective access time =

$$P * \text{page fault service time} + (1-P) * \text{memory access time}$$

If we do not have any page fault the effective access time = memory access time ( $\therefore P = 0$ )

If page fault increases then effective access time also increases.

- page replacement Algorithm  
when page fault occurs, the page replacement is needed. Page replacement means select a victim page in main memory and replace that page with required page from the backing store.

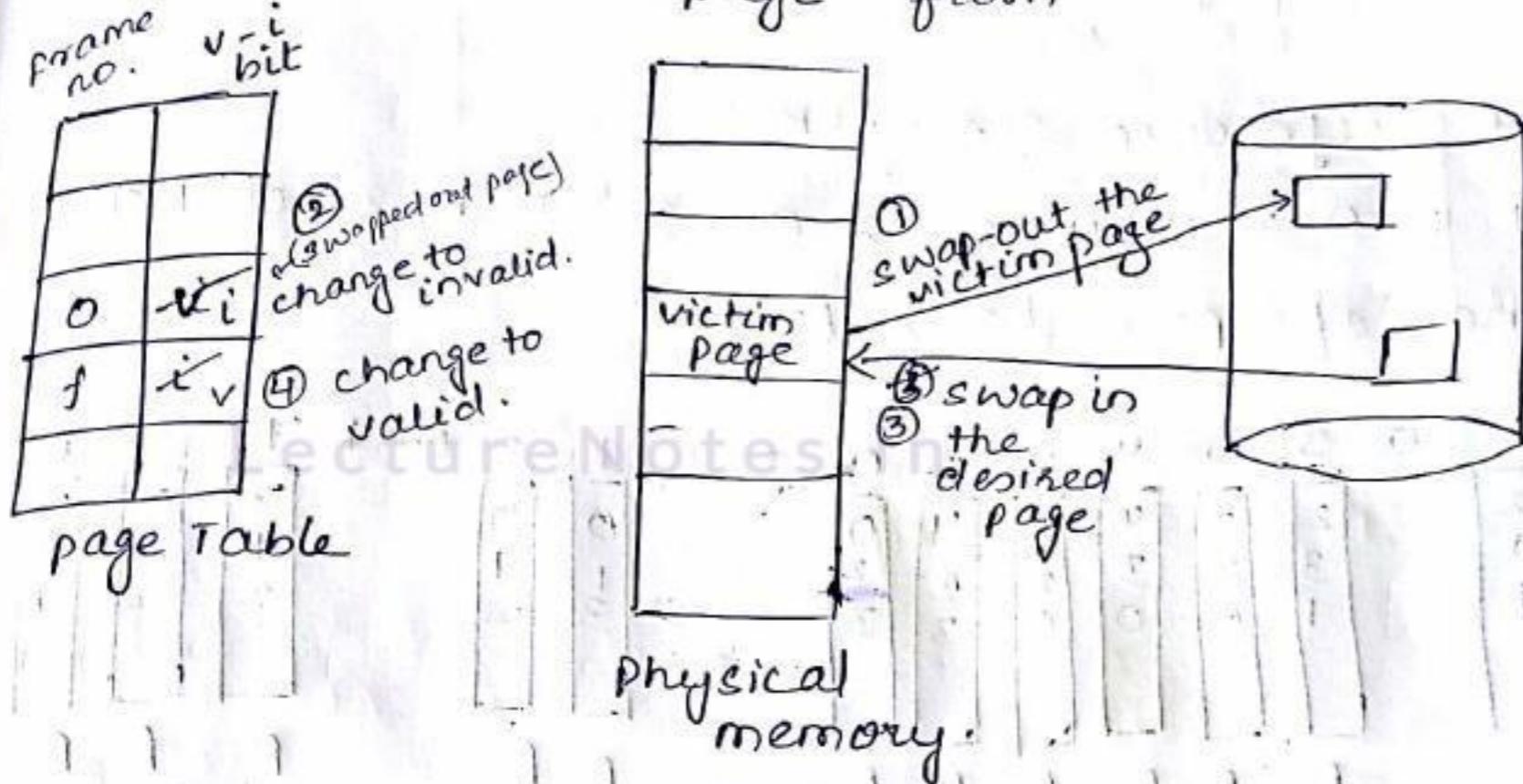
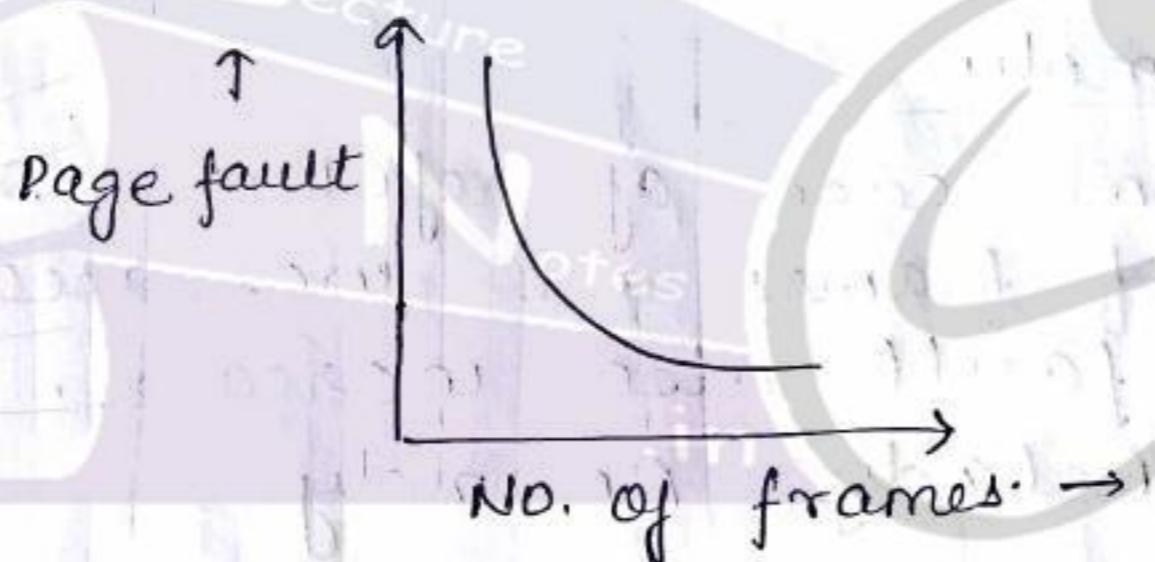


Fig: Page replacement working principle.

If the no. of frames increases the number of page fault decreases.



LectureNotes.in

### FIFO Replacement Algorithm

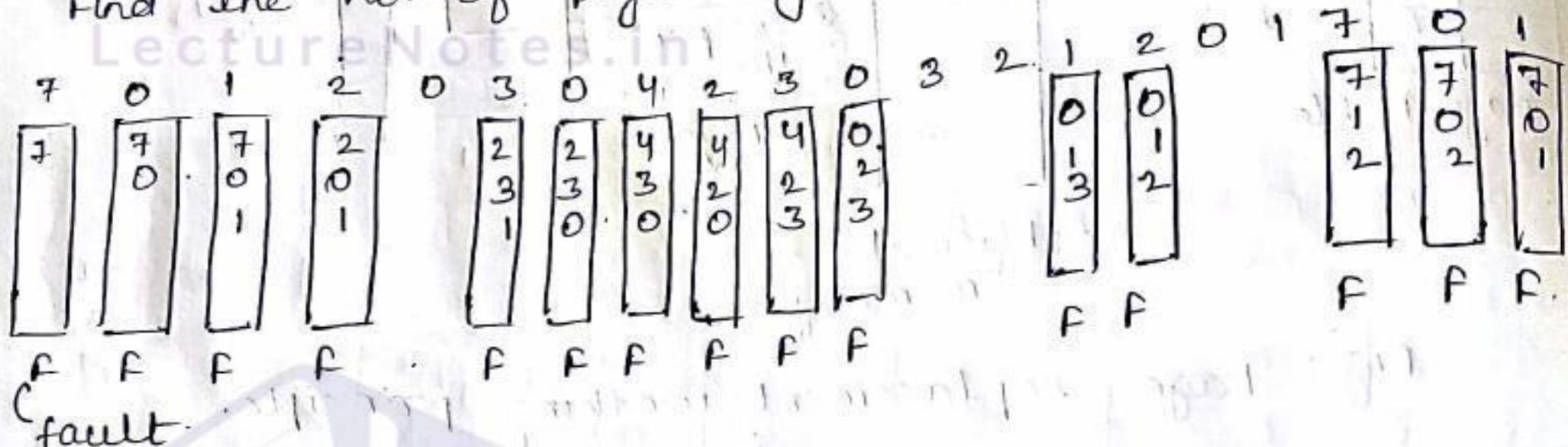
Replace the page of main memory on the principle of FIFO. It means replace the page which comes first.

Example -

Suppose reference string is

7 0 1 2 0 3 0 4 2 3 0 3 2 1 2 0 1 7 0 1

Find the no. of page layouts -



(fault)

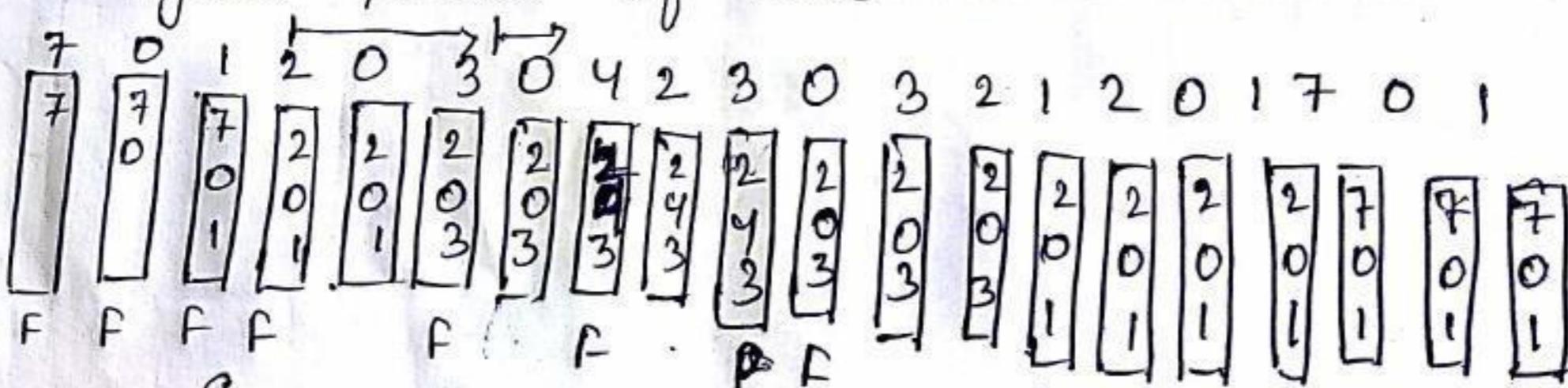
so, no. of page fault is 15.

### Belady's Anomaly

This is some case of reference string in which if the no. of frames are increased then the no. of page fault also increases. This situation is called Belady's Anomaly.

### Optimal page representation Algorithm

Replace the page that will not be used for longest period of time



so, the no. of page fault is 9.

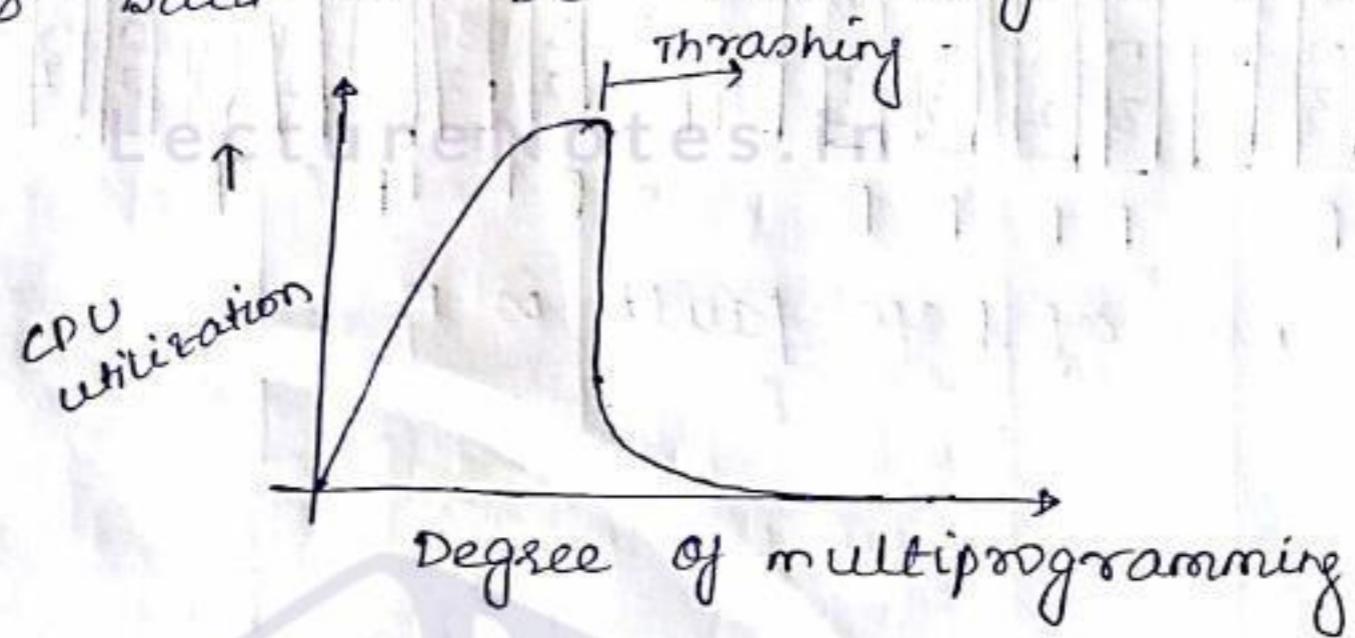
least Recently used (LRU) Page replacement Algorithm  
Replace the page that has not been used for  
the longest period of time.

7	0	1	2	0	3	0	4	2	3	0	3	2	1	2	0	1	7	0	1
7	0	7	2	2	0	1	3	3	3	0	3	2	3	2	3	2	0	7	0
F	*	F	F	F	F	F	F	F	F	F	F	F	F	F	F	F	F	F	

so the no. of page fault is 12.

LectureNotes.in

<sup>VV2</sup> Thrashing :- If the no. of processes submitted to the CPU for execution are increased, the CPU utilization also increases but increasing the process continuously at certain time the CPU utilization falls sharply and sometimes it reaches to 0. This situation is said to be 'thrashing'.



If the process does not have required no. of frames, it will cause quickly page fault. At this point, it must replace some page. since all its pages are in active use, it must replace a page that will not be needed again. As a result, it will give fault again and again.

If this process continues to fault, replacing the pages for which it hit the faults and bring back the so in right away. This high paging activity is called thrashing.

A page is thrashing if it is spending more time in paging than executing.

memory & virtual m/m  
complete

## File system Interface

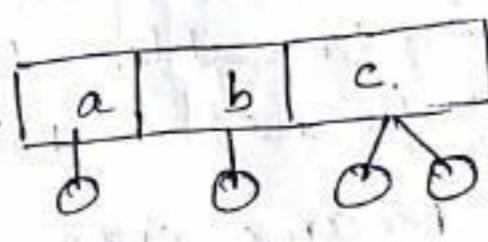
file access method (LQ).

- i) sequential access    ii) Direct access    iii.) Index sequential access.

directory structure (L8)

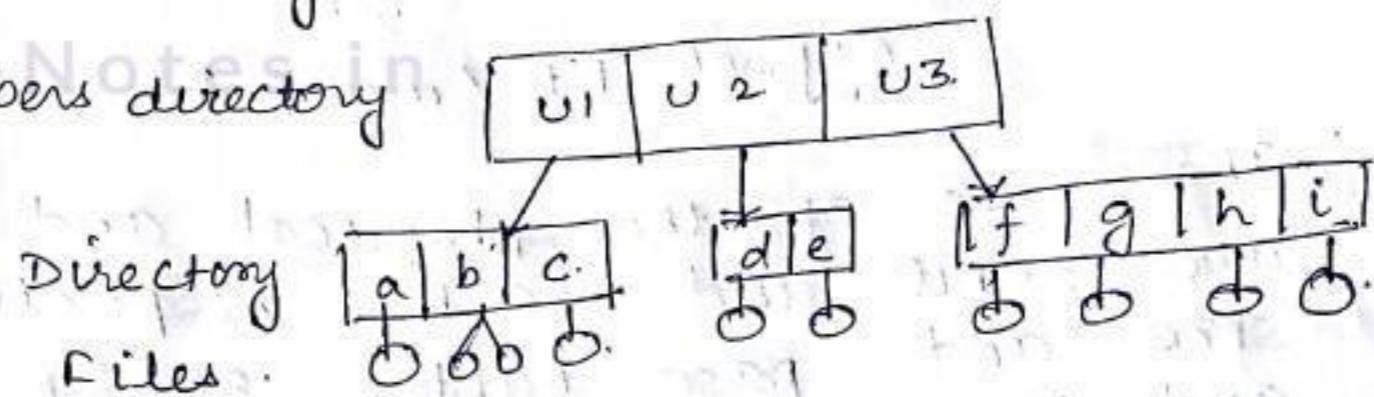
- i) single level directory.

directories  
files.

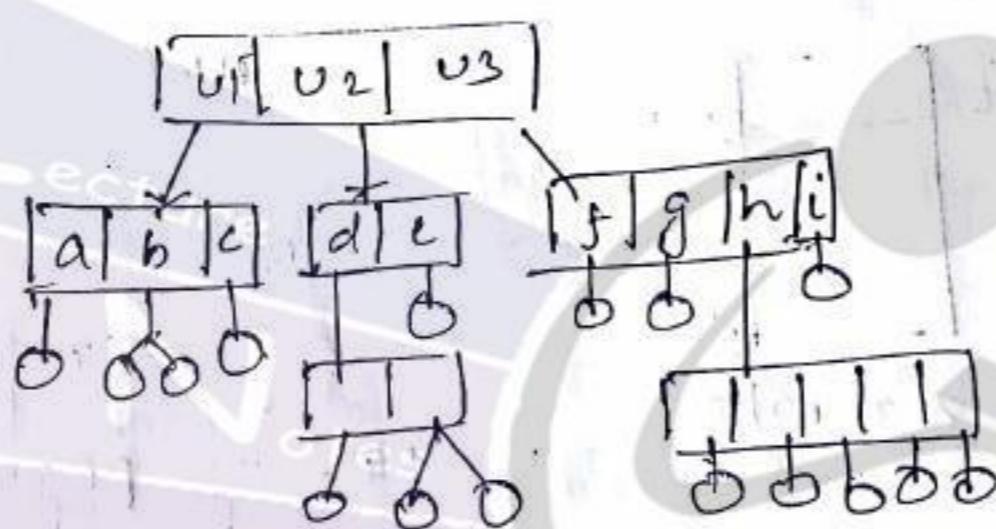


- ii) two level directory

users directory



- iii.) Tree structured directory!



→ problems (memory).

Q: calculate the page offset if the physical address is 24 bits and main memory having 128 frames.

$$\rightarrow 0 - 2^{24} \quad \text{size} = 2^{24}$$

Physical address = 24 bits

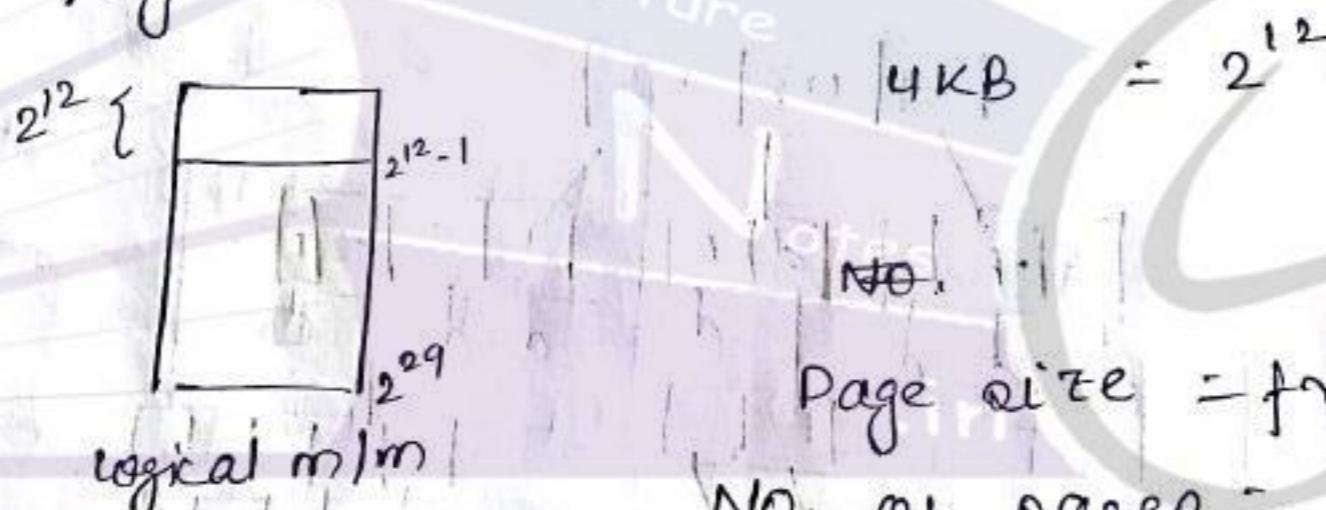
∴ Addressable locations =  $2^{24}$

No. of frames = 128 =  $2^7$

$$\text{Page size} = \text{frame size} = \frac{2^{24}}{2^7} = 2^{17} = 17 \text{ bits.}$$

Offset = 17 (max).

Q: The size of the physical and logical addresses are 25 bits and 29 bits respectively. If the frame size and page table entry size are 4KB and 4 byte respectively, what is the size of the page table.



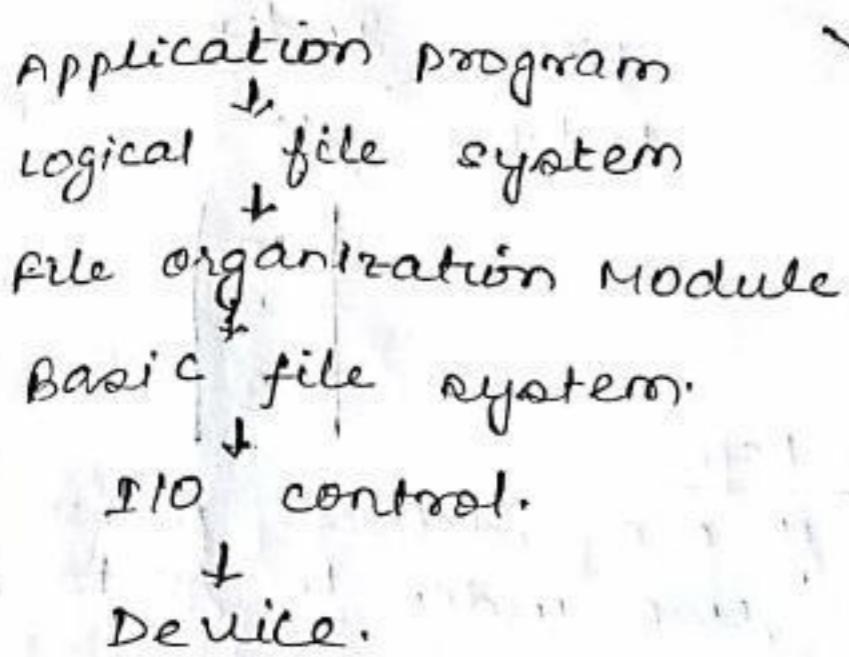
Page size = frame size = 4KB

$$\text{No. of pages} = \frac{2^{29}}{2^{12}} = 2^{17}$$

There is an entry for every page in page table  
∴ No. of entries in page table =  $2^{17}$ .

$$\text{size of page table} = \frac{2^{17} \times 4 \text{ B}}{\cancel{2^{17}}} = 512 \text{ KB}$$

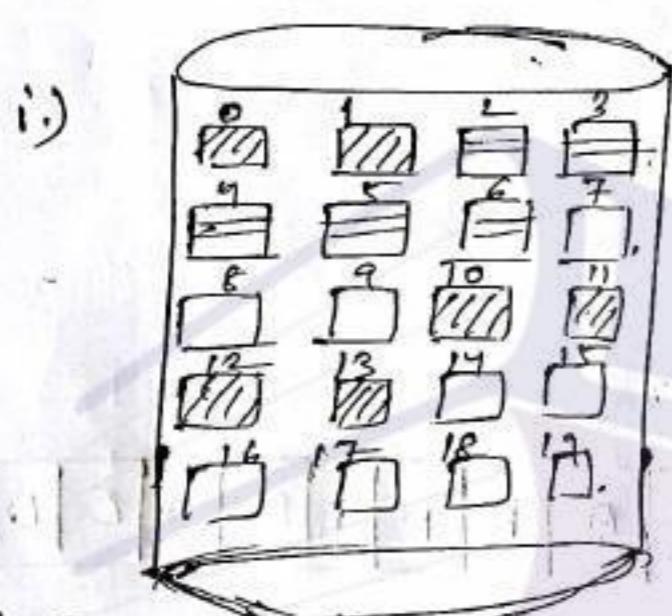
## Ch - File system implementation



## file system organization (LQ)

### Allocation methods (LQ)

- i) contiguous allocation
- ii) linked allocation
- iii) indexed allocation

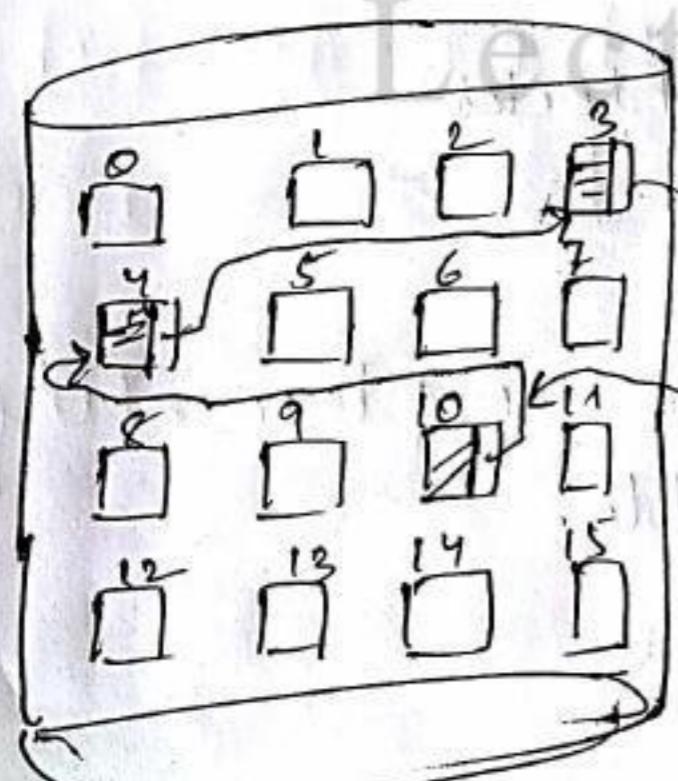


File	start	length
test.c	0	2
add.cpp	10	4
sub.java	2	5

### Disadvantage

- eliminated
- i) size of bit space should be available from before
  - ii) when size of file increases or decreases then fragmentation.

### linked allocation



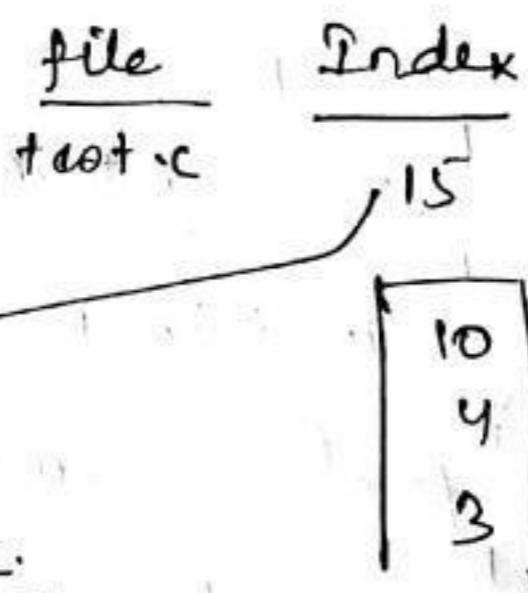
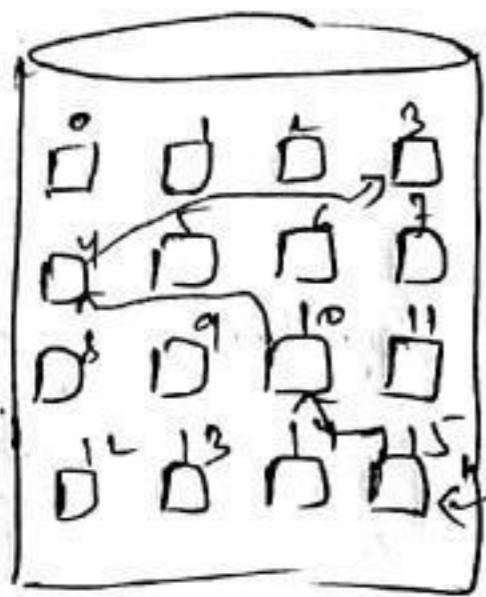
file	start	end
Test.c	10	3

### Disadvantage

- i) for eg:- 4th block can't be accessed directly. u have to go through 10.
- ii) If block 4 misses then u can't move forward.
- iii) mem occupied due to storng of address.

### iii) indexed allocation :-

Object  
3/1/17

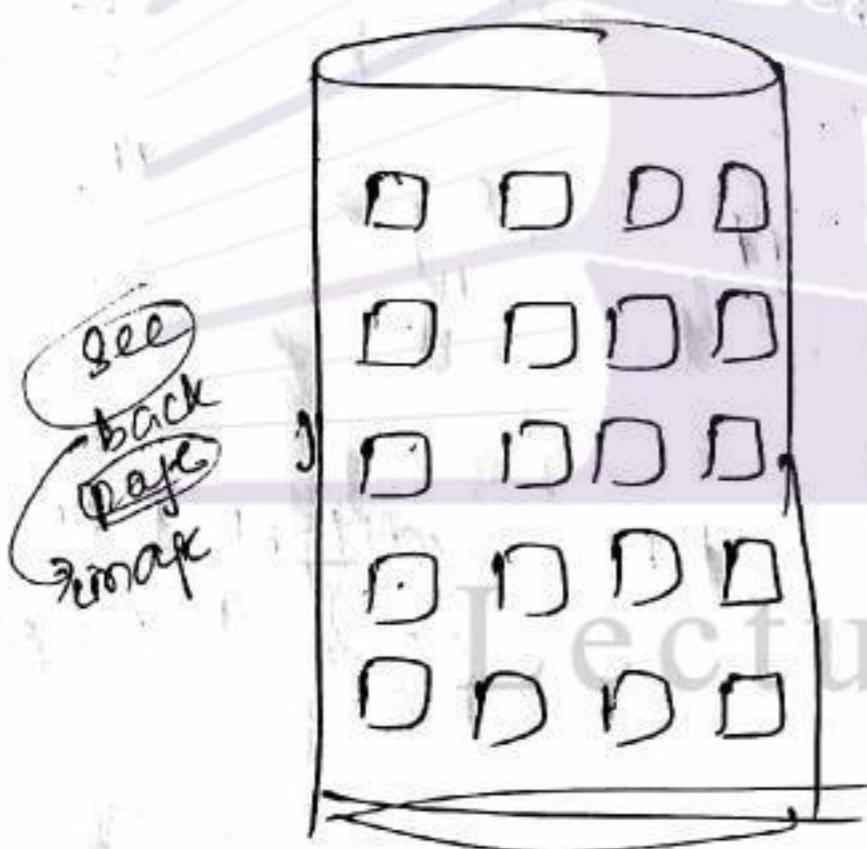


#### Advantage.

Even if one address misses then we can use index to go to next block

### → Performance consideration.

- i.) free space management
- ii.) bit vector (mainly)



Bit vector



- iii.) linked list

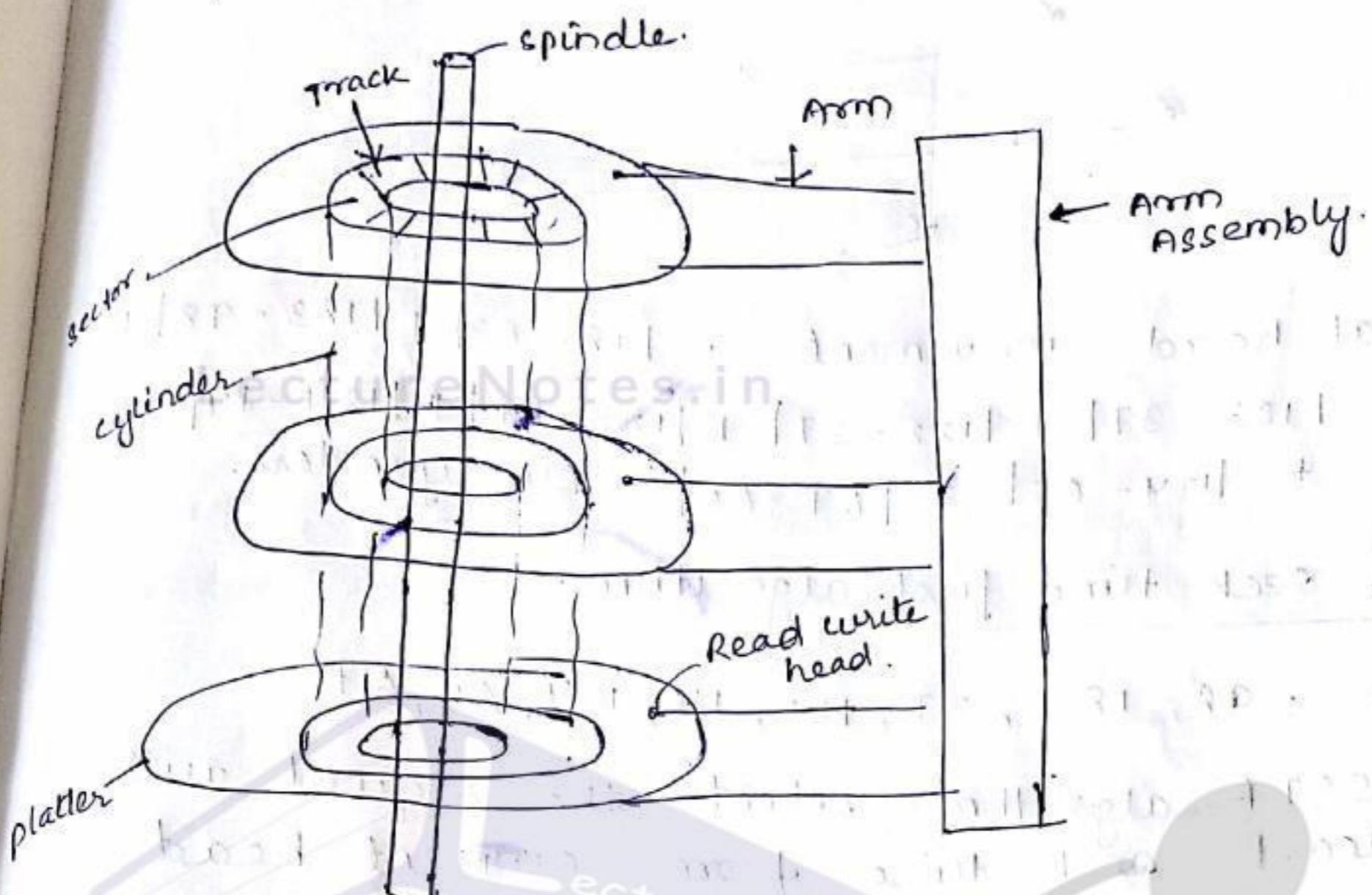
- iv.) counting (count the free blocks)

### → Directory implementation.

There are two methods for implementing directory linear list and hash table.

Mass storage structure

1. Magnetic Disc :-



1) Transfer rate :- It is the rate at which data flow b/w drive and computer.

2) Access time consist of two parts :-

a) seek time :- The time required to move the disc arm to the desired cylinder.

b) rotational latency :- The time required for the desired sector to rotate the disc head.

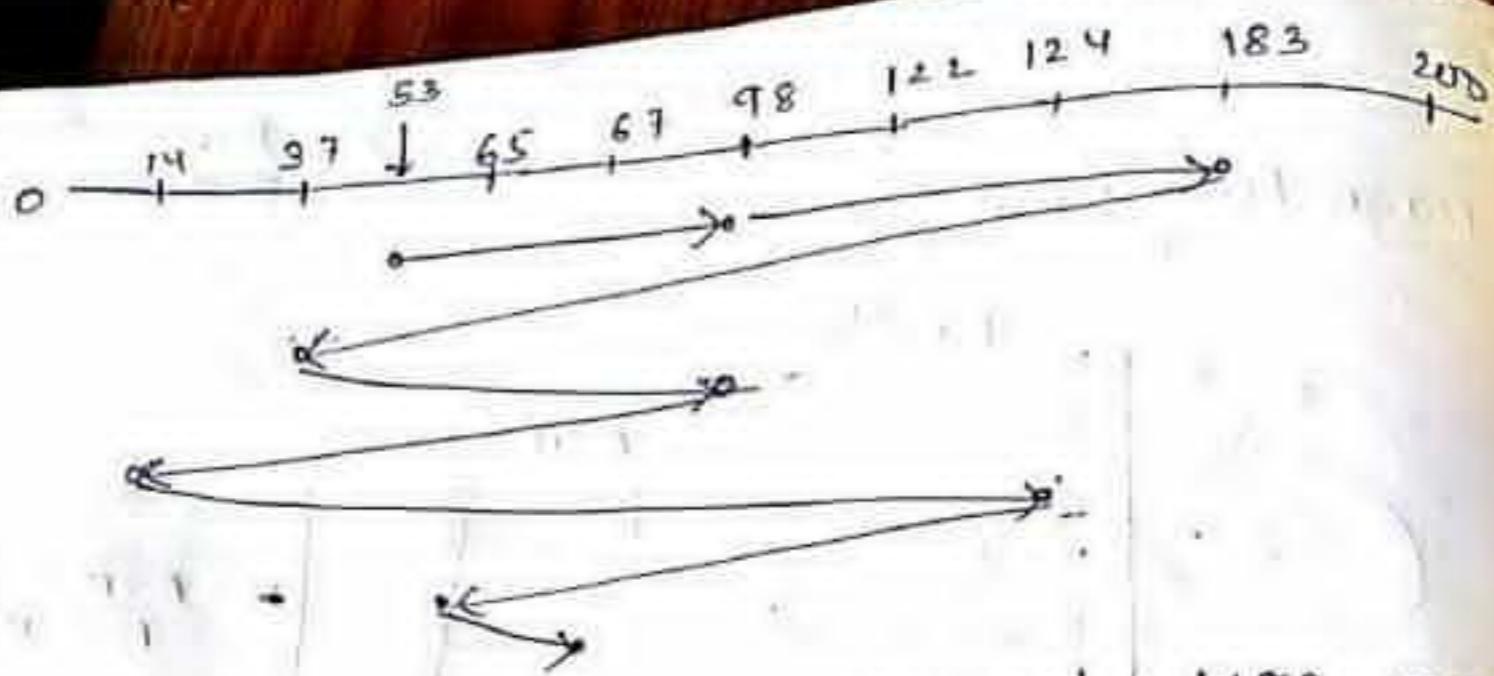
- Disc scheduling :-

- FCFS scheduling algorithm :-

A disc queue request the blocks on cylinders of following order.

Queue : 98, 183, 37, 122, 14, 124, 65, 67

Head start at = 53. Find the total head movement.

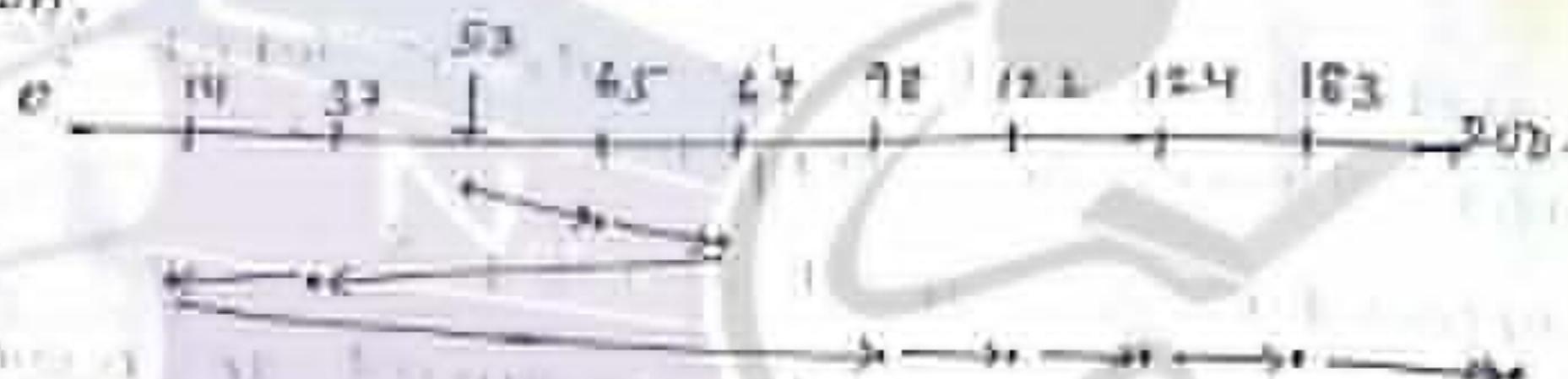


$$\therefore \text{Total head movement} = |198 - 53| + |183 - 98| + \\ |183 - 37| + |122 - 37| + |122 - 14| + |124 - 14| \\ + |124 - 65| + |67 - 65| = 640 \text{ cylinders.}$$

### - shortest seek time first algorithm

queue : 98, 183, 37, 122, 14, 124, 65, 67

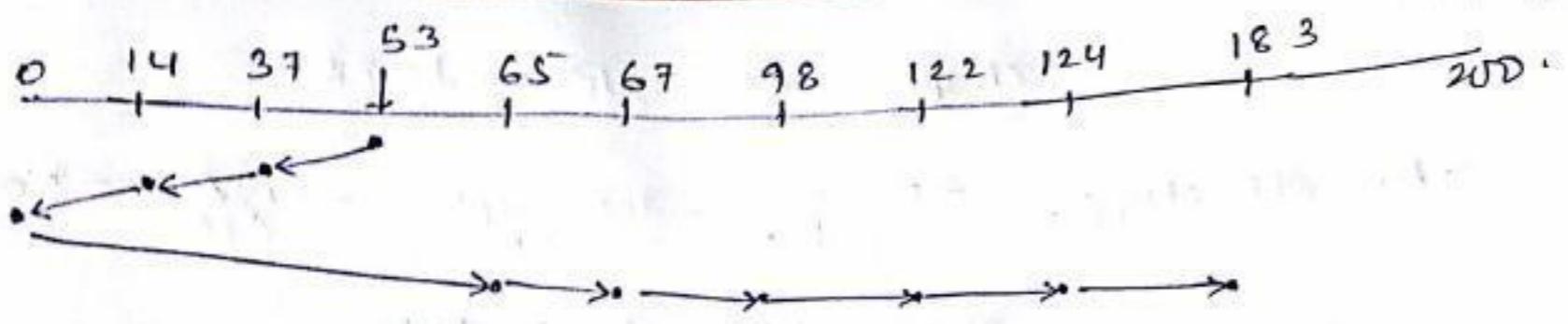
The SSTF algorithm selects the request with the least seek time from current head position.



$$\text{Total head movement} = |65 - 53| + |67 - 65| + |67 - 37| \\ + |37 - 14| + |698 - 14| + |122 - 14| \\ + |124 - 122| + |183 - 124| + |67 - 183| \\ = 236$$

### - SCAN / Elevator Algorithm

queue : 98, 183, 37, 122, 14, 124, 65, 67  
Head starts at 53

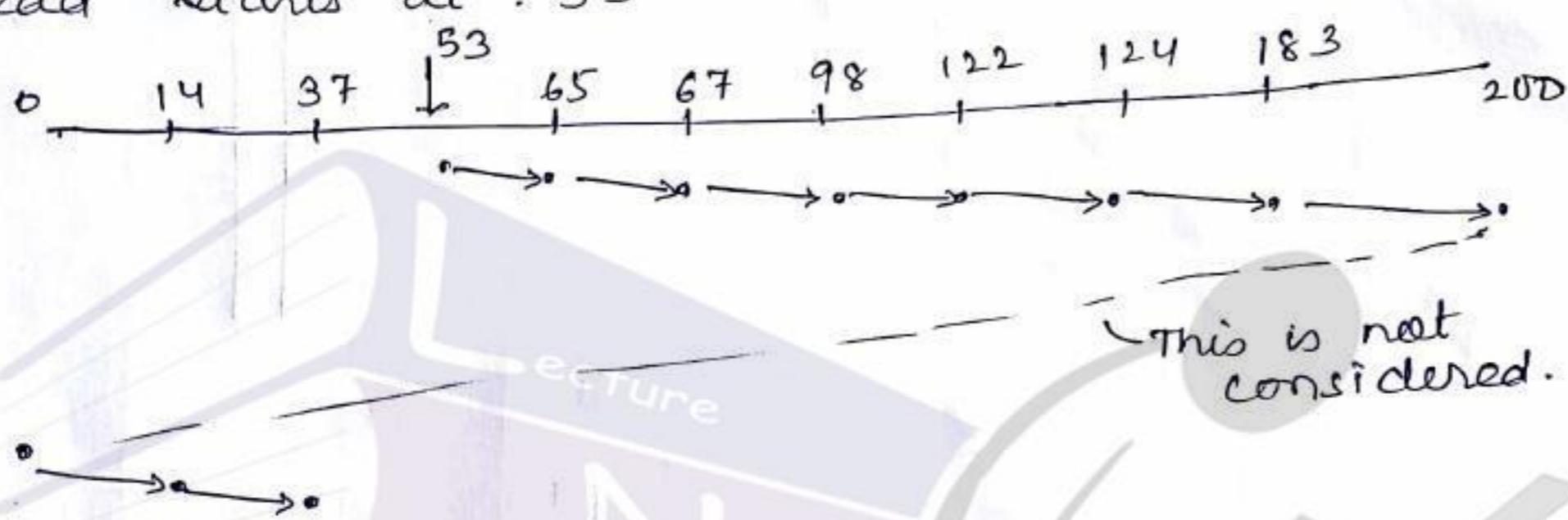


$$\text{Total head movement} = |53 - 37| + |14 - 37| + |14 - 0| + \\ + |65 - 0| + |67 - 65| + |98 - 67| + \\ + |122 - 98| + |124 - 122| + |183 - 124| = 236$$

### circular - Scan (C-SAN) Algorithm

Queue : 98, 183, 37, 122, 14, 124, 65, 67.

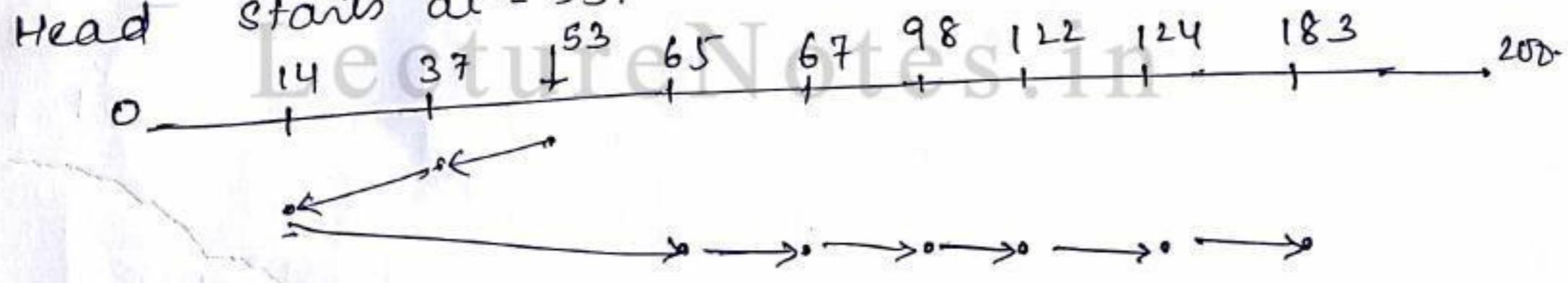
Head starts at : 53



### LOOK scheduling

Queue : 98, 183, 37, 122, 14, 124, 65, 67

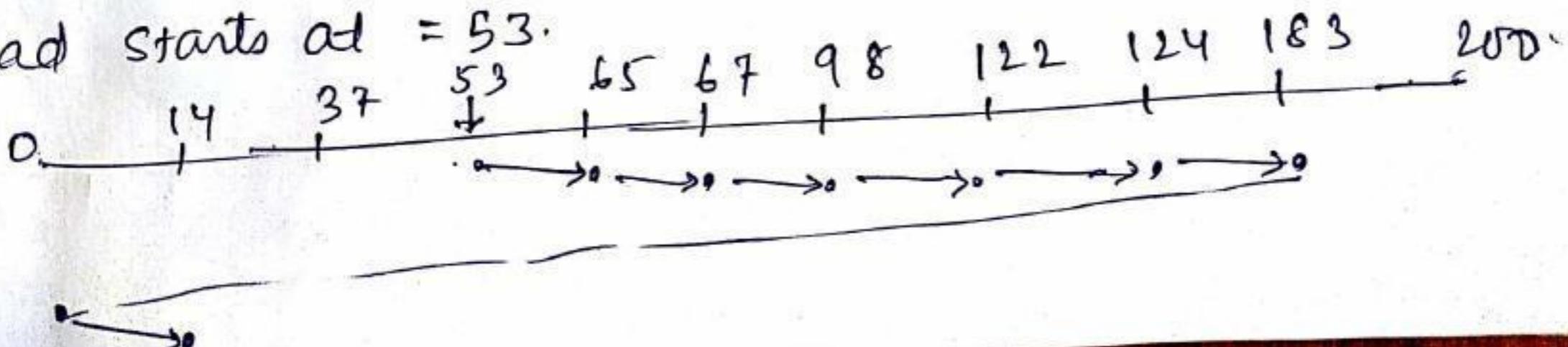
Head starts at = 53.



### circular look algorithm (C-LOOK)

Queue : 98, 183, 37, 122, 14, 124, 65, 67

Head starts at = 53.



7/11/17

chapter : 13 I/O systems

- bus structures
- Polling (S&H)
- Interrupt (L&Q)
- DMA (L&Q)
- I/O request hardware to hardware to (L&Q) open

chapter : 14 protection

- Domain structure:
  - Access matrix of figure A uses objects.
  - Access matrix with copy rights.
  - Access matrix with owner rights.

LectureNotes.in

ch15: security

virus worm

Friday  
Test mode

# File System

**Dr. Ajay Kumar Jena  
School of Computer Engineering  
KIIT University, Bhubaneswar**

LectureNotes.in

## ACKNOWLEDGEMENT

I am very much thankful to say that  
some contents of these slides are taken  
from the book

Operating System Concepts  
by  
Silberschatz, Galvin and Gagne.

# File Concept

- It is a collection of related information that is recorded in the secondary storage
- Contiguous logical address space (programs i.e. source and object forms)
- File control bloc: storage structure consisting of information about a file
- Types:
  - Data
    - Numeric, alphanumeric
    - Character, alphabetic
    - binary
  - Program
- Contents defined by file's creator
  - Many types
    - Consider **text file, source file, executable file**
- A file is a sequence of bits, bytes, lines, records.



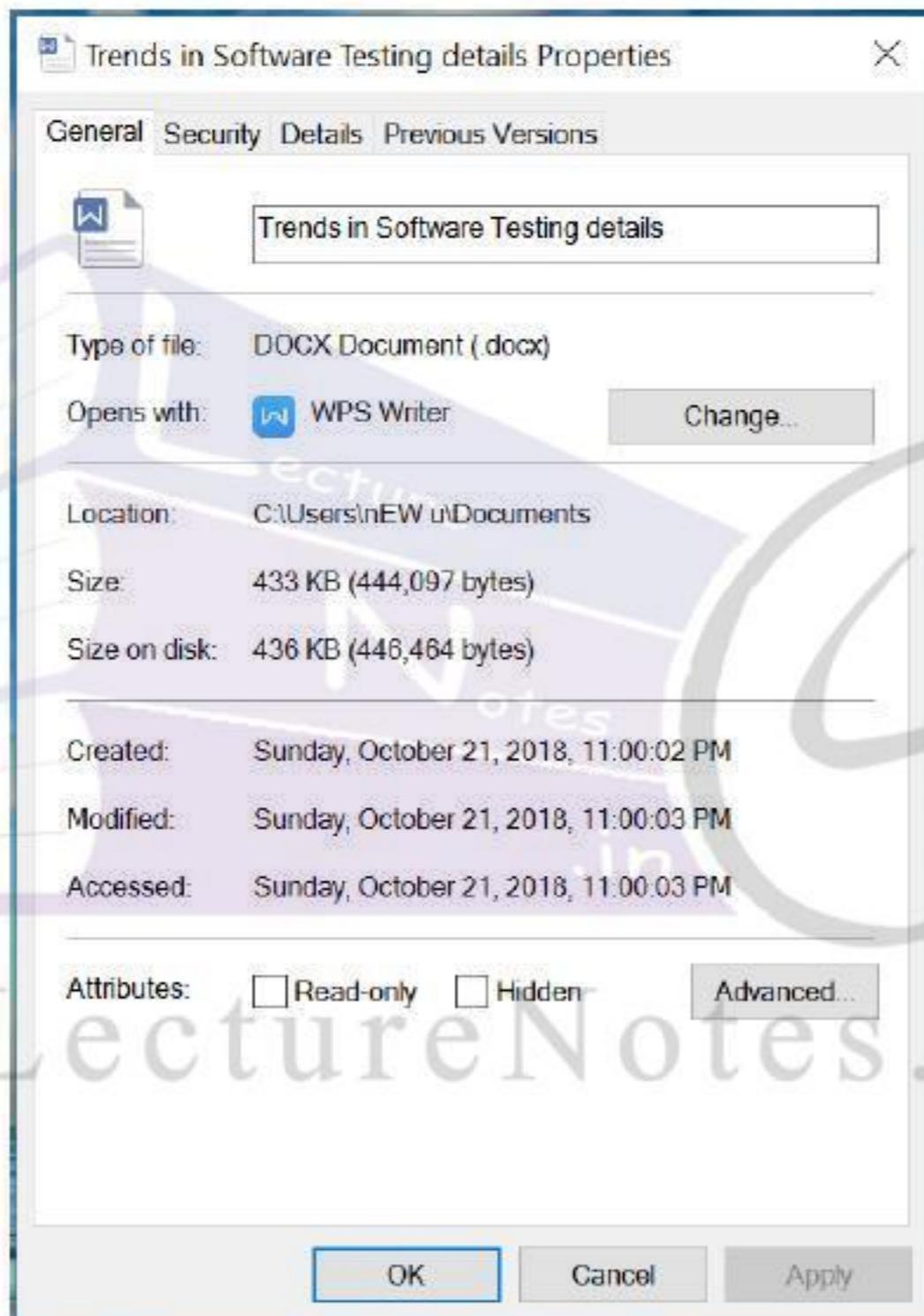
LectureNotes.in

LectureNotes.in

# File Attributes

- A file is named, for the convenience of its users and is referred to you by name
- **Name** – symbolic file name is the only information kept in human-readable form
- **Identifier** – unique tag (number) identifies file within file system
- **Type** – needed for systems that support different types
- **Location** – pointer to file location on device
- **Size** – current file size
- **Protection** – controls who can do reading, writing, executing
- **Time, date, and user identification** – data for protection, security, and usage monitoring
- Information about files are kept in the directory structure, which is maintained on the disk
- Many variations, including extended file attributes such as file checksum
- Information kept in the directory structure

# File info Window on Mac OS X



# File Operations

- File is an **abstract data type**
- **Create** – Space in file system must be found, an entry for the new file will be made in the directory
- **Write** – To make a system call to write at **write pointer** location
- **Read** – at **read pointer** location
- **Reposition within file - seek**
- **Delete**
- **Truncate**: User may erase the content of the file but wants to hold the attributes
- **Open( $F_i$ )** – search the directory structure on disk for entry  $F_i$ , and move the content of entry to memory
- **Close ( $F_i$ )** – move the content of entry  $F_i$  in memory to directory structure on disk

# File Types – Name, Extension

file type	usual extension	function
executable	exe, com, bin or none	ready-to-run machine-language program
object	obj, o	compiled, machine language, not linked
source code	c, cc, java, pas, asm, a	source code in various languages
batch	bat, sh	commands to the command interpreter
text	txt, doc	textual data, documents
word processor	wp, tex, rtf, doc	various word-processor formats
library	lib, a, so, dll	libraries of routines for programmers
print or view	ps, pdf, jpg	ASCII or binary file in a format for printing or viewing
archive	arc, zip, tar	related files grouped into one file, sometimes compressed, for archiving or storage
multimedia	mpeg, mov, rm, mp3, avi	binary file containing audio or A/V information

# File Structure

- None - sequence of words, bytes
- Simple record structure
  - Lines
  - Fixed length
  - Variable length
- Complex Structures
  - Formatted document
  - Relocatable load file
- Can simulate last two with first method by inserting appropriate control characters
- Who decides:
  - Operating system
  - Program



# LectureNotes.in Access Methods

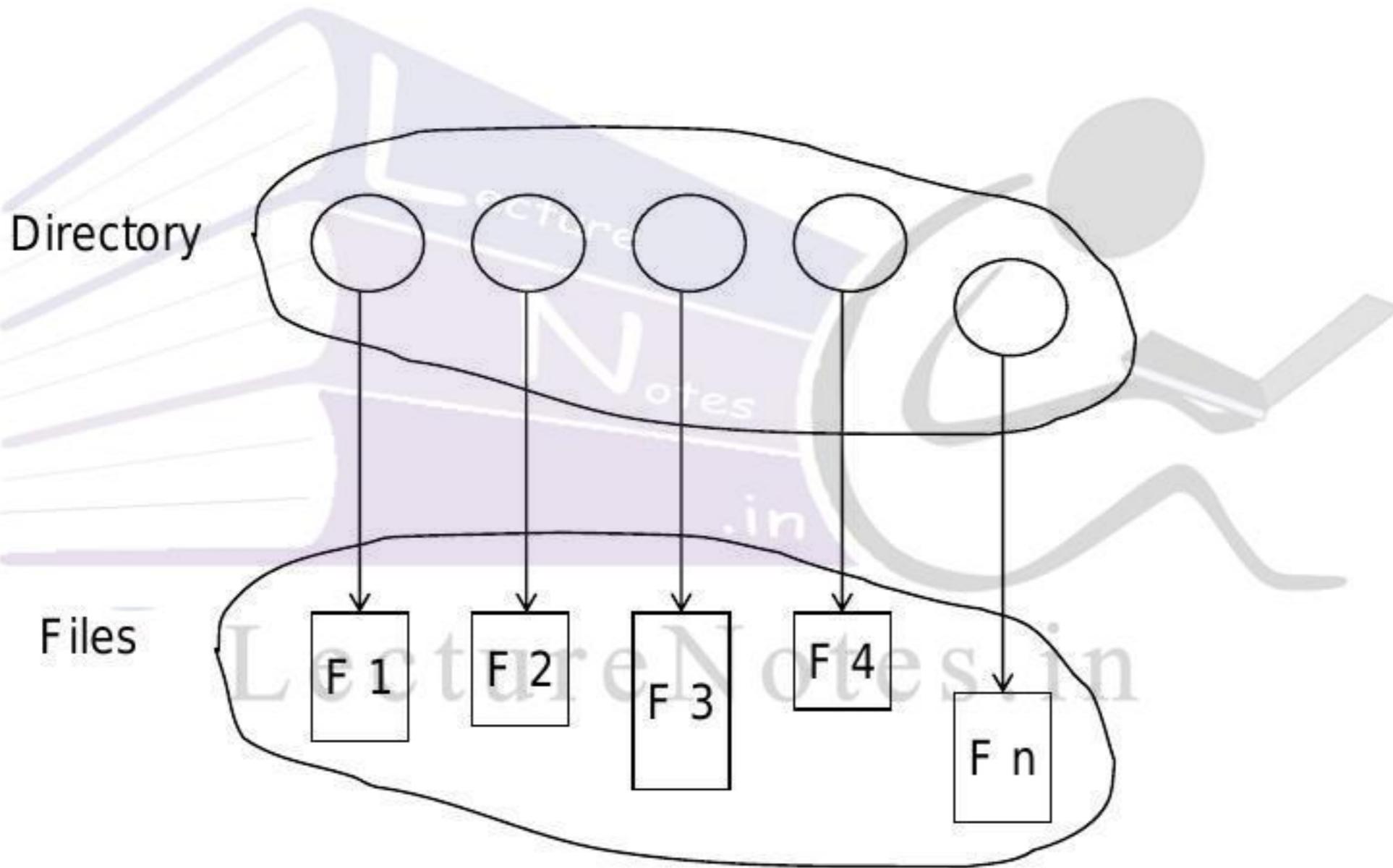
- The information in the file can be accessed in several ways.
- **Sequential Access.** The most simplest method. Information accessed sequentially one record after another
  - read next**
  - write next**
  - no read after last write (rewrite)
- **Direct Access** – When particular record/s are needed to be accessed out of many record. Randomly access any file block from HDD.
  - read n**
  - write n**
  - position to n**
  - read next**
  - write next**
  - rewrite n**

*n = relative block number*
- **Index Sequential Access:** Combination of sequential access and direct access. This access method involves maintaining an index. Te index is a pointer to the block.



# Directory Structure

- Files are to be organised. How to store the file.
- A collection of nodes containing information about all files



Both the directory structure and the files reside on disk

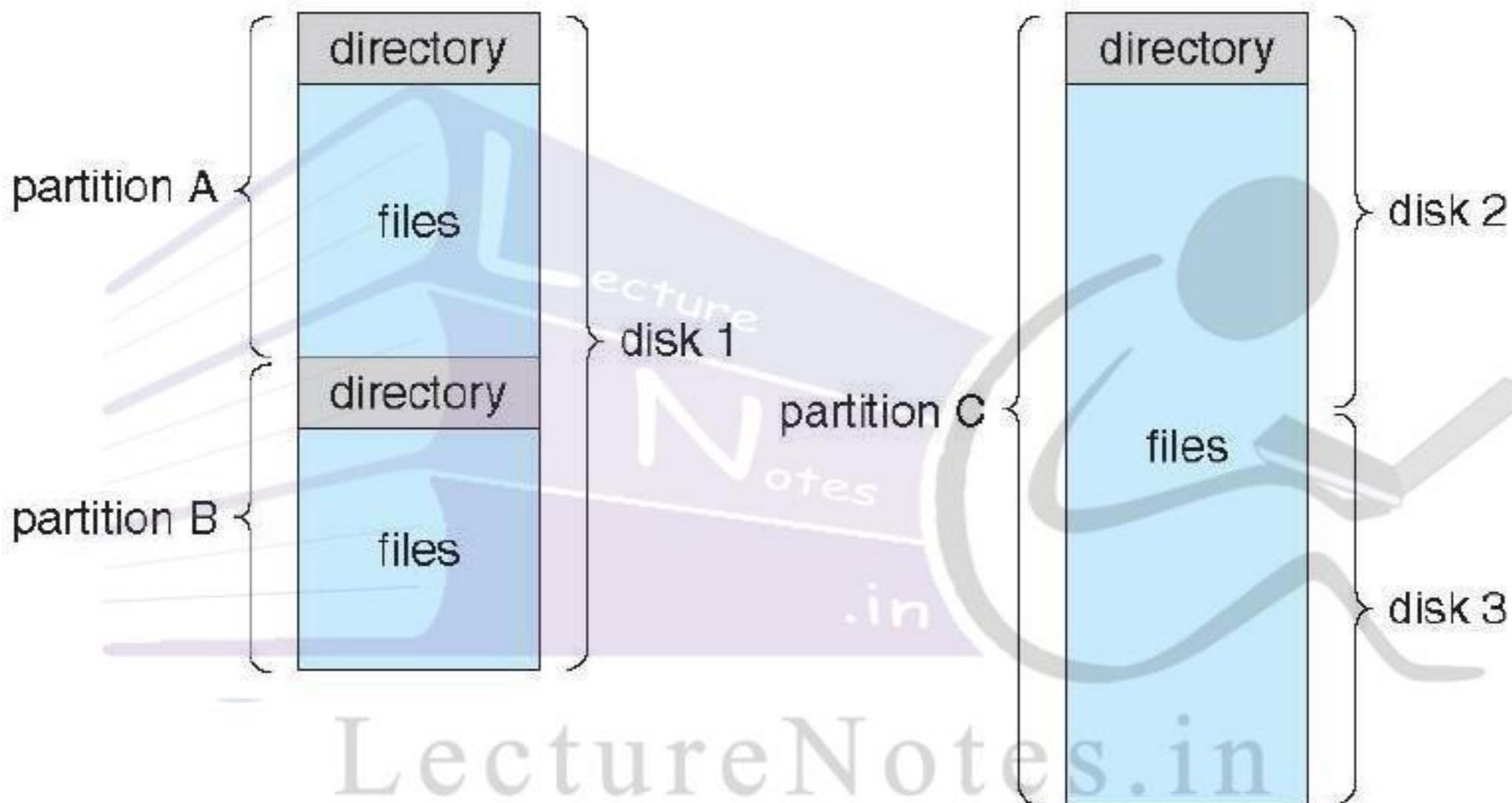
# Disk Structure

- Disk can be subdivided into **partitions**
- Disks or partitions can be **RAID** protected against failure
- Disk or partition can be used **raw** – without a file system, or **formatted** with a file system
- Partitions also known as minidisks, slices
- Entity containing file system known as a **volume**
- Each volume containing file system also tracks that file system's info in **device directory** or **volume table of contents**
- As well as **general-purpose file systems** there are many **special-purpose file systems**, frequently all within the same operating system or computer

LectureNotes.in

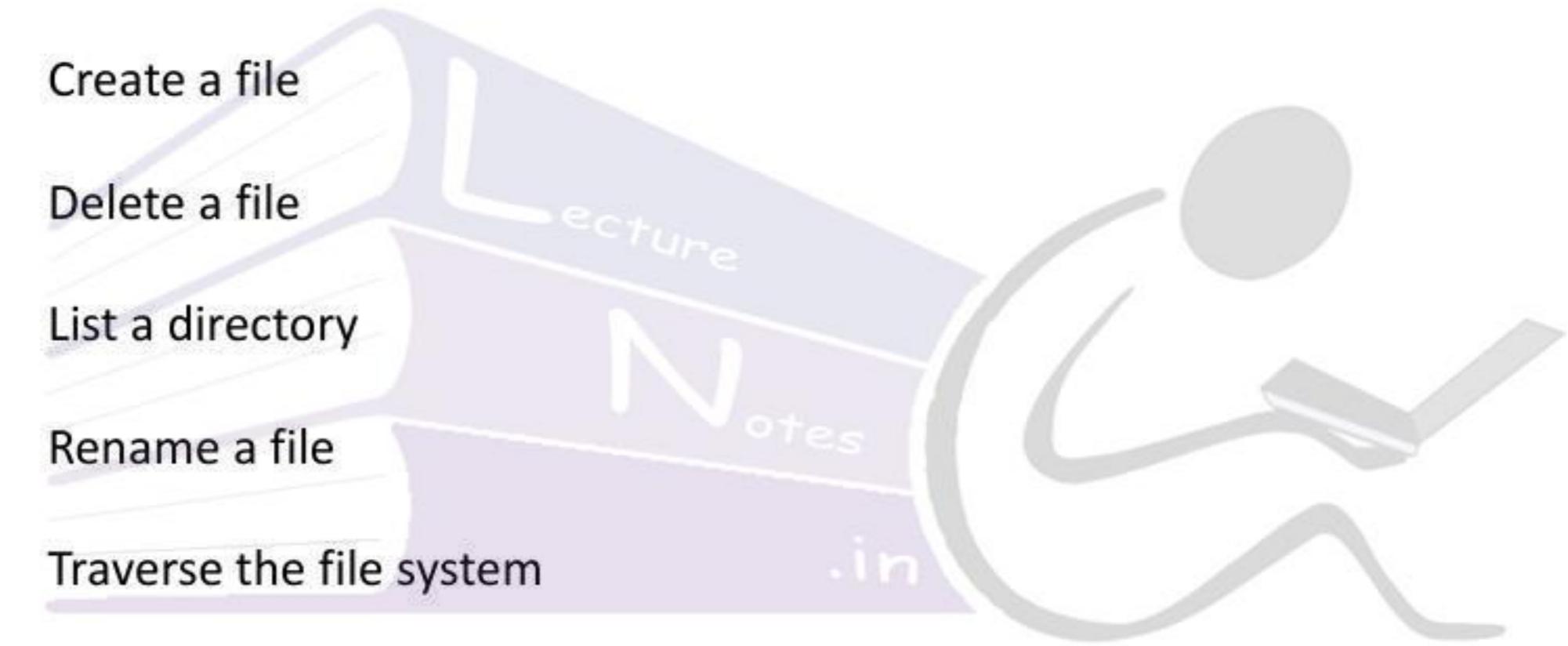
LectureNotes.in

# A Typical File-system Organization



# Operations Performed on Directory

- Search for a file
- Create a file
- Delete a file
- List a directory
- Rename a file
- Traverse the file system



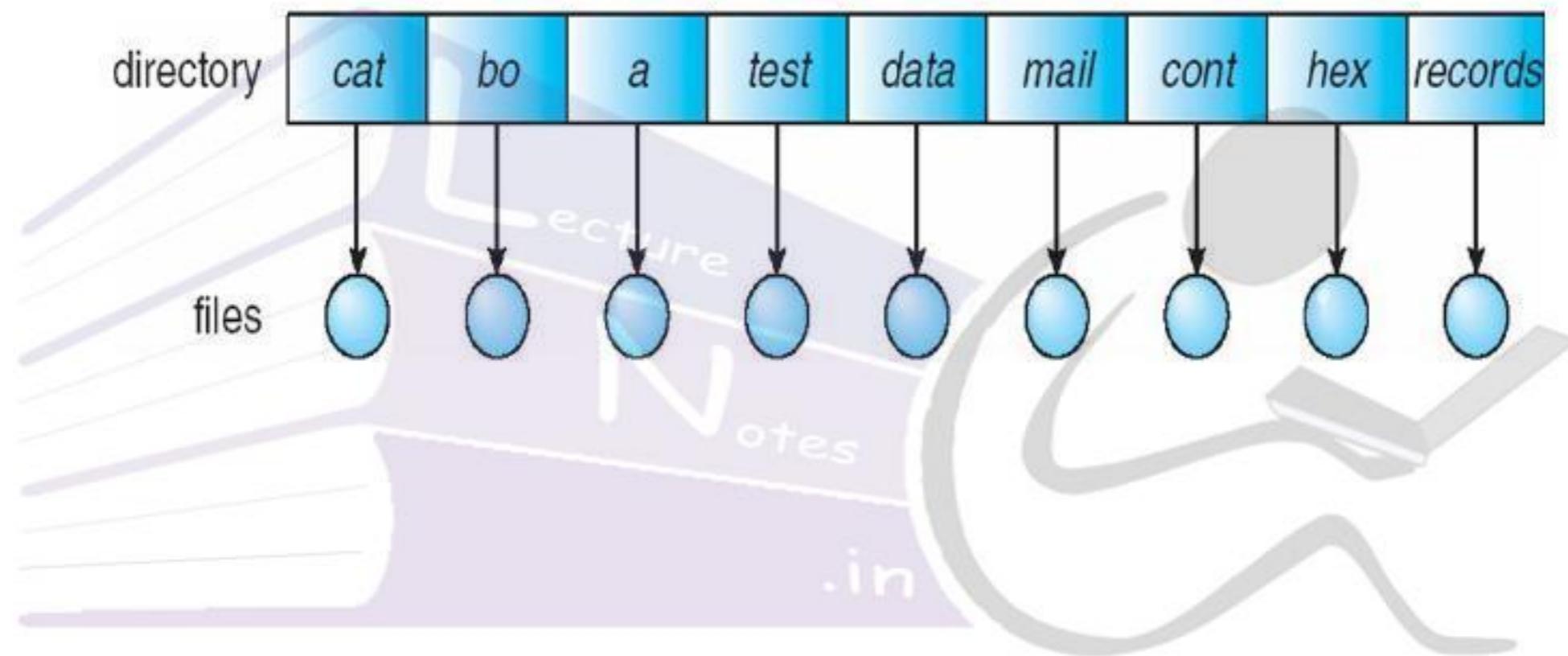
LectureNotes.in

LectureNotes.in

By Ajay Kumar Jena, KIIT Deemed to be  
University, Bhubaneswar

# Single-Level Directory

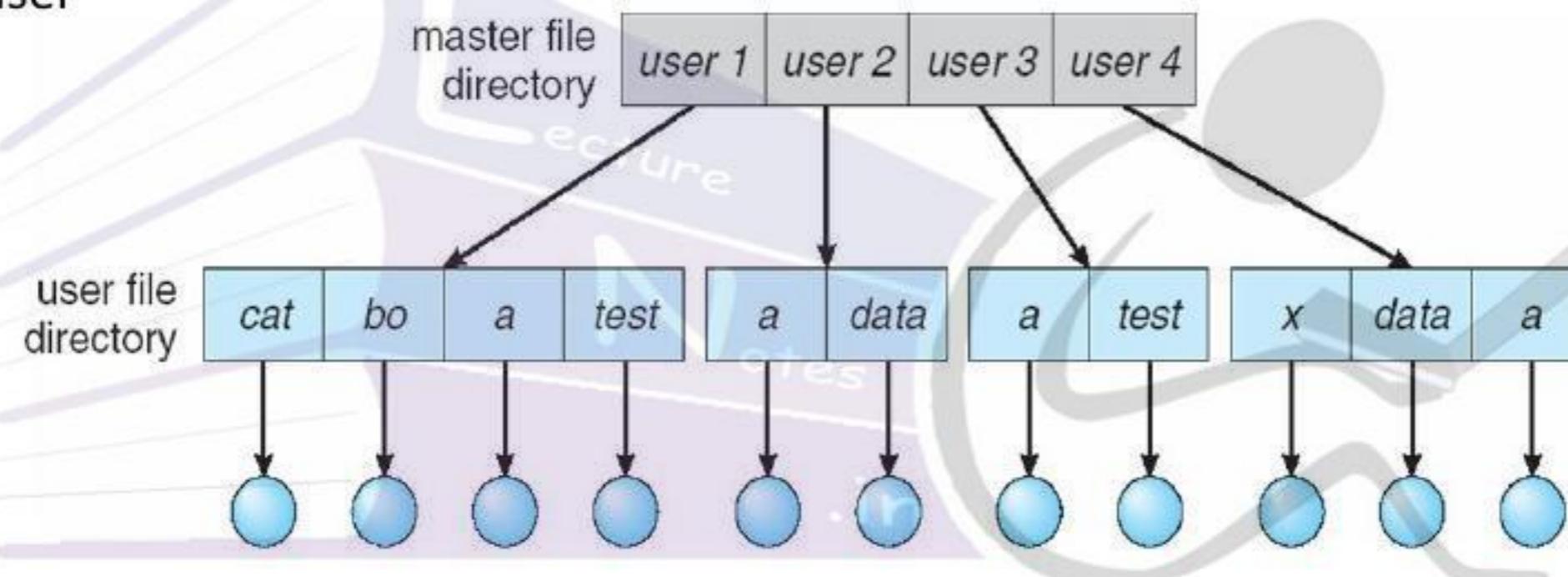
- The simplest directory structure is a single directory for all users.
- All files reside in one and the same directory.



- Naming problem (No two files having same name)
- Grouping problem

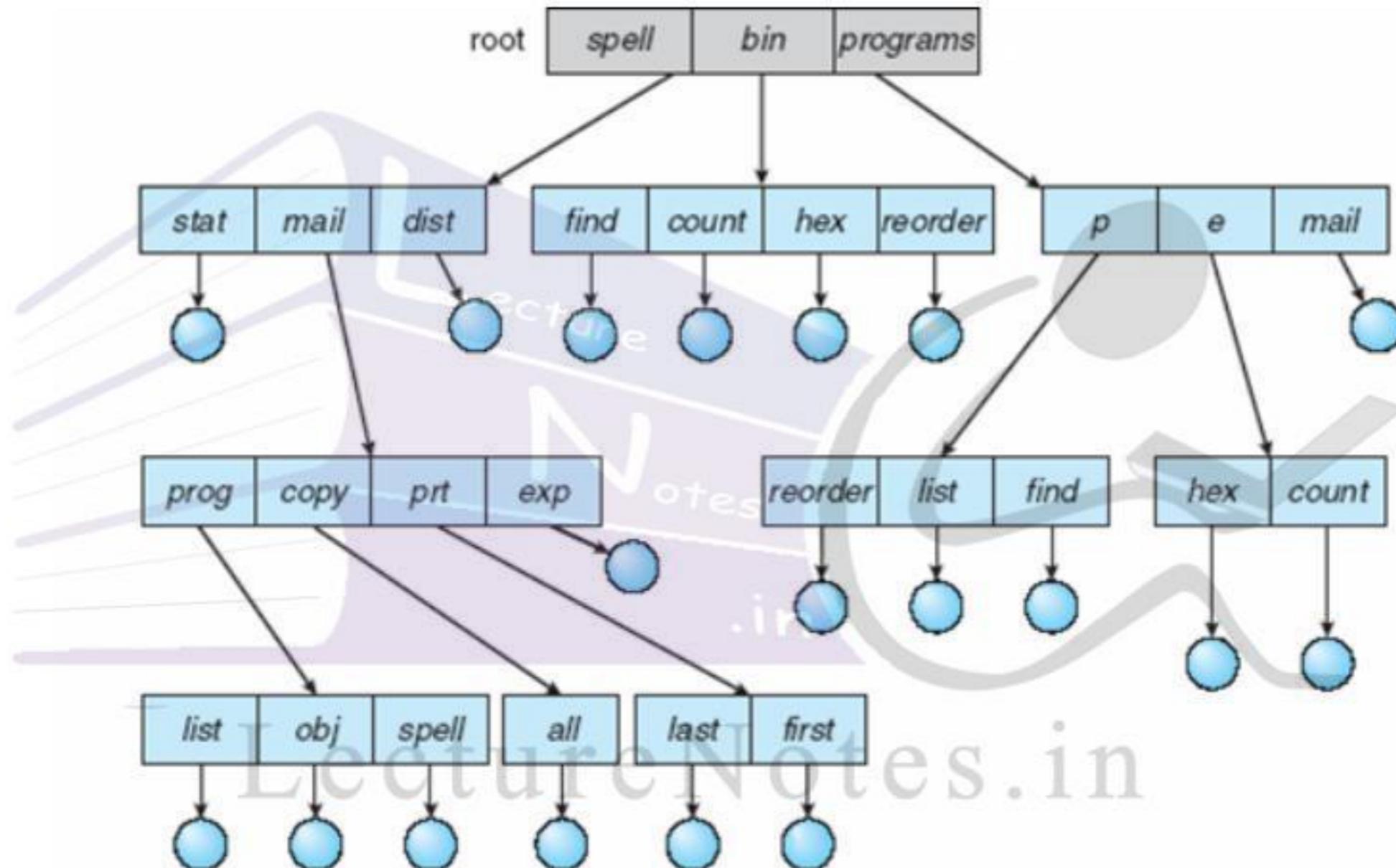
# Two-Level Directory

- The main drawback of single-level is to have unique file names by different users. It can be solved by creating separate directories.
- Each user has its own user file directory (UFD). Separate directory for each user



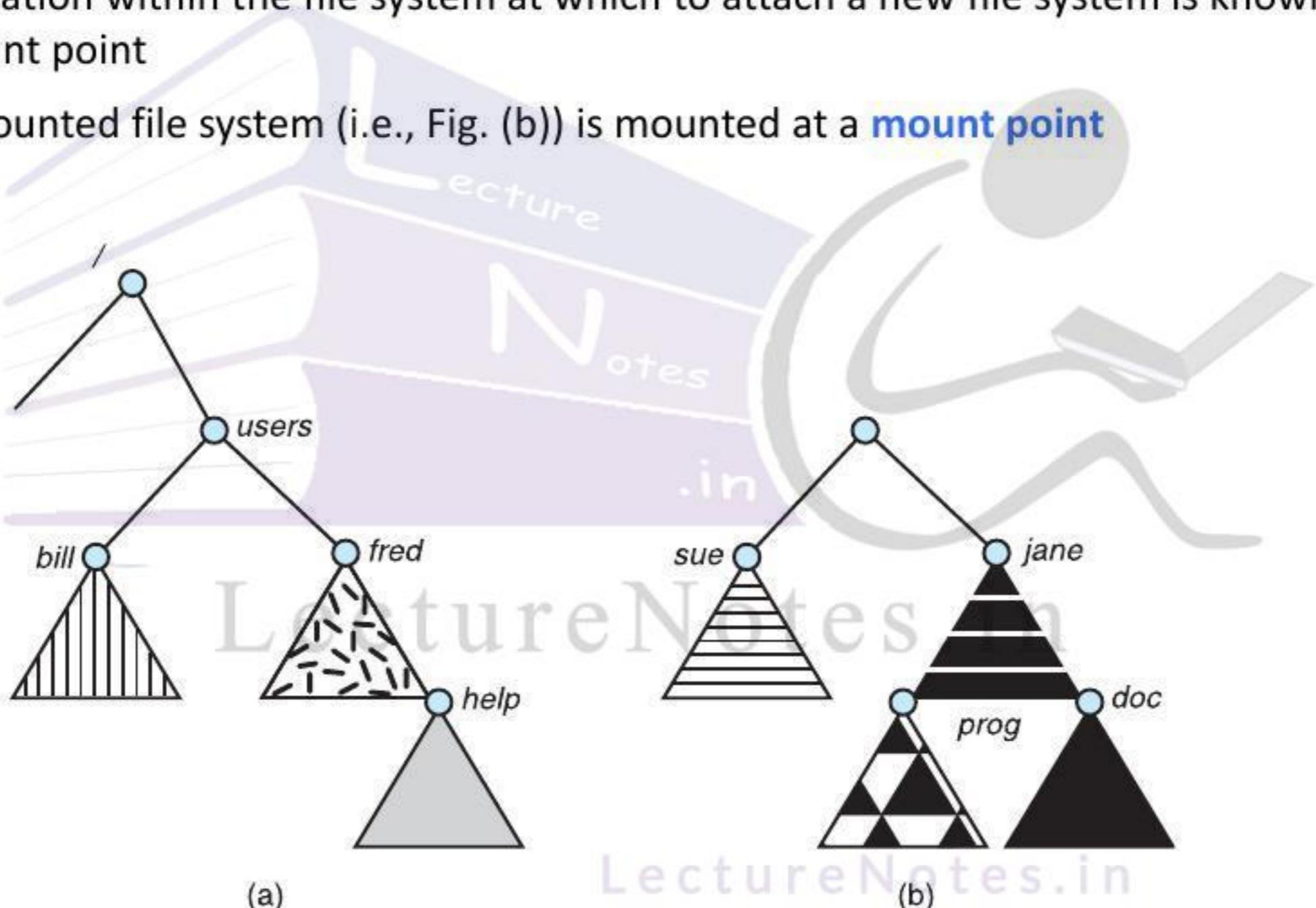
- One directory exclusive for each user. Path name
- Can have the same file name for different user
- Efficient searching
- No grouping capability

# Tree-Structured Directories

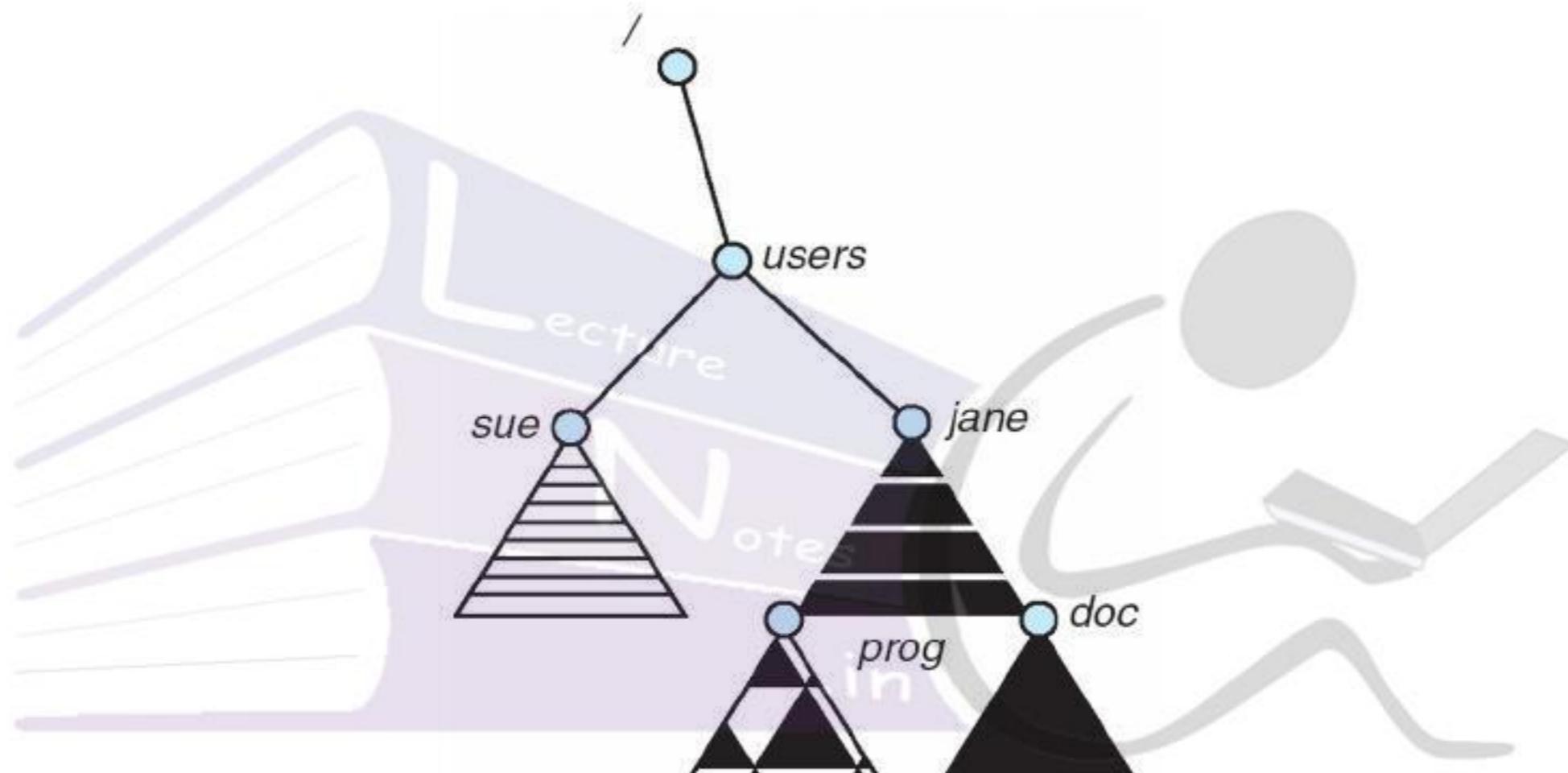


# File System Mounting

- As a file must open before it is used, a file system must be **mounted** before it can be accessed.
- The location within the file system at which to attach a new file system is known as mount point
- A unmounted file system (i.e., Fig. (b)) is mounted at a **mount point**



# Lecture Notes in Mount Point



LectureNotes.in

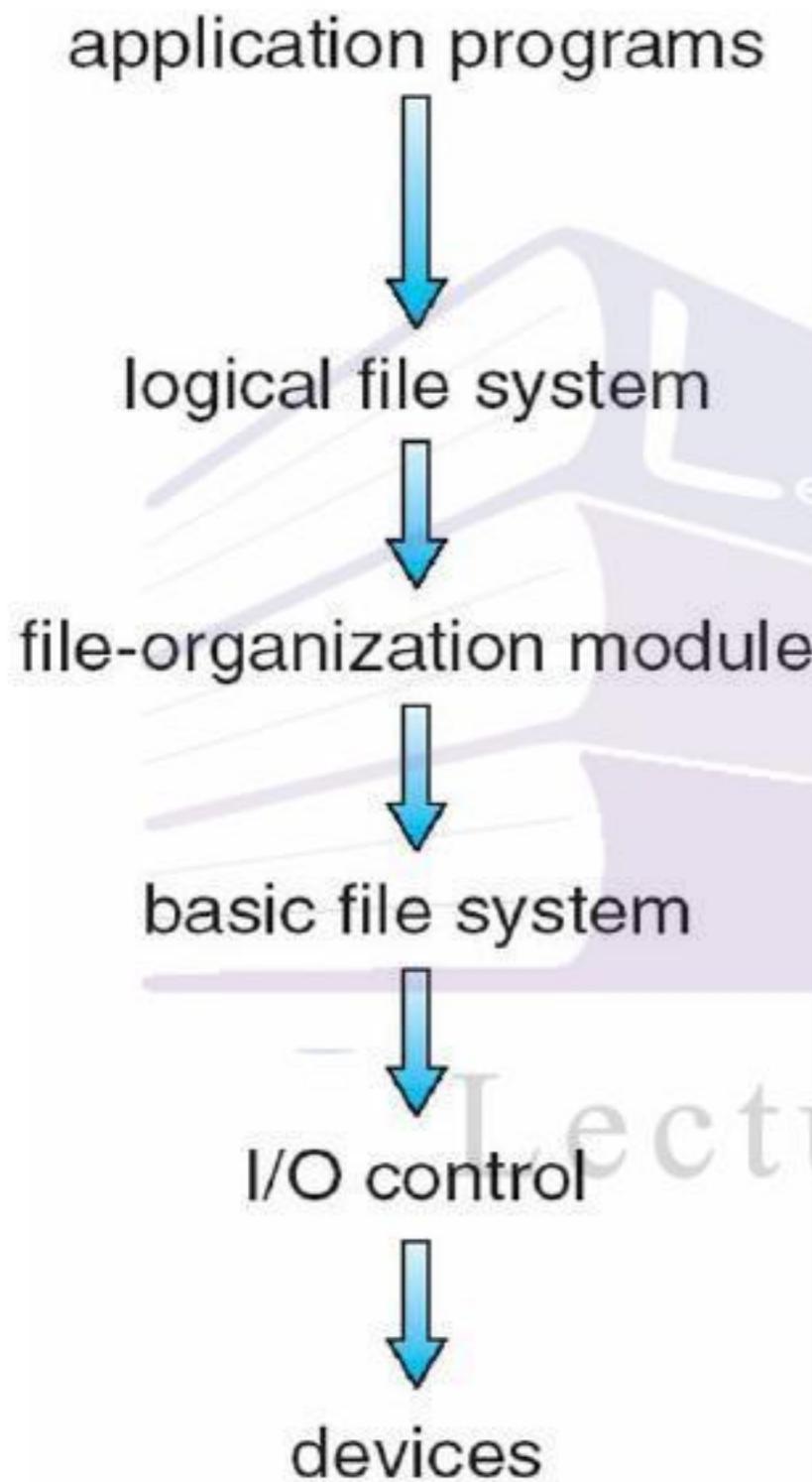
LectureNotes.in

By Ajay Kumar Jena, KIIT Deemed to be  
University, Bhubaneswar

# File-System Structure

- File is the most important part in computer. Bcoz all the data stores in files. It is the part of operating system. IO, Process, everything is a file. It is a sw which is managing all these files. It stores data and program.
- File structure
  - Logical storage unit, Collection of related information
  - Allocating disk space
  - Recovering the free space
- **File system** resides on secondary storage (disks)
  - Provided user interface to storage, mapping logical to physical
  - Provides efficient and convenient access to disk by allowing data to be stored, located retrieved easily
- Disk provides in-place rewrite and random access
  - I/O transfers performed in **blocks** of **sectors** (usually 512 bytes)
- **File control block** – storage structure consisting of information about a file
- **Device driver** controls the physical device
- File system organized into layers

# LectureNotes.in Layered File System



By designing any SW we use layering approach. We are going to divide all the task in to groups.

When one try to modify one layer it will affect the other layer.

Application program desires for a file i.e create, modify, append and read it will ask the request to the logical file system . LFS has meta data (who is the user, permission, when created with directory structure)

FOM: File is divided into logical blocks.  
Mapping of LB=Phy. Block.) also free space mgt.

BFS: Read/Write. Read one block. issue the command to the IO control to do it.

IOC: consists of device drivers, interrupt handler

Devices: IO control will access the device.

# File System Layers

- **Device drivers** manage I/O devices at the I/O control layer
  - Given commands like “read drive1, cylinder 72, track 2, sector 10, into memory location 1060” outputs low-level hardware specific commands to hardware controller
- **Basic file system** given command like “retrieve block 123” translates to device driver
- Also manages memory buffers and caches (allocation, freeing, replacement)
  - Buffers hold data in transit
  - Caches hold frequently used data
- **File organization module** understands files, logical address, and physical blocks
  - Translates logical block # to physical block #
  - Manages free space, disk allocation

# File System Layers (Cont.)

- **Logical file system** manages metadata information
  - Translates file name into file number, file handle, location by maintaining file control blocks (**inodes** in UNIX)
  - Directory management
  - Protection
- Layering useful for reducing complexity and redundancy, but adds overhead and can decrease performanceTranslates file name into file number, file handle, location by maintaining file control blocks (**inodes** in UNIX)
  - Logical layers can be implemented by any coding method according to OS designer

LectureNotes.in

LectureNotes.in

# File System Layers (Cont.)

- Many file systems, sometimes many within an operating system
  - Each with its own format (CD-ROM is ISO 9660; Unix has **UFS**, **FFS**; Windows has **FAT**, **FAT32**, **NTFS** as well as floppy, CD, DVD Blu-ray, Linux has more than 40 types, with **extended file system** **ext2** and **ext3** leading; plus distributed file systems, etc.)
  - New ones still arriving – **ZFS**, **GoogleFS**, **Oracle ASM**, **FUSE**

# File-System Implementation

- We have system calls at the API level, but how do we implement their functions?
  - On-disk and in-memory structures
- **Boot control block** contains info needed by system to boot OS from that volume
  - Needed if volume contains OS, usually first block of volume
- **Volume control block (superblock, master file table)** contains volume details
  - Total # of blocks, # of free blocks, block size, free block pointers or array
- Directory structure organizes the files
  - Names and inode numbers, master file table

# File-System Implementation (Cont.)

- Per-file **File Control Block (FCB)** contains many details about the file
  - inode number, permissions, size, dates
  - NFTS stores into in master file table using relational DB structures

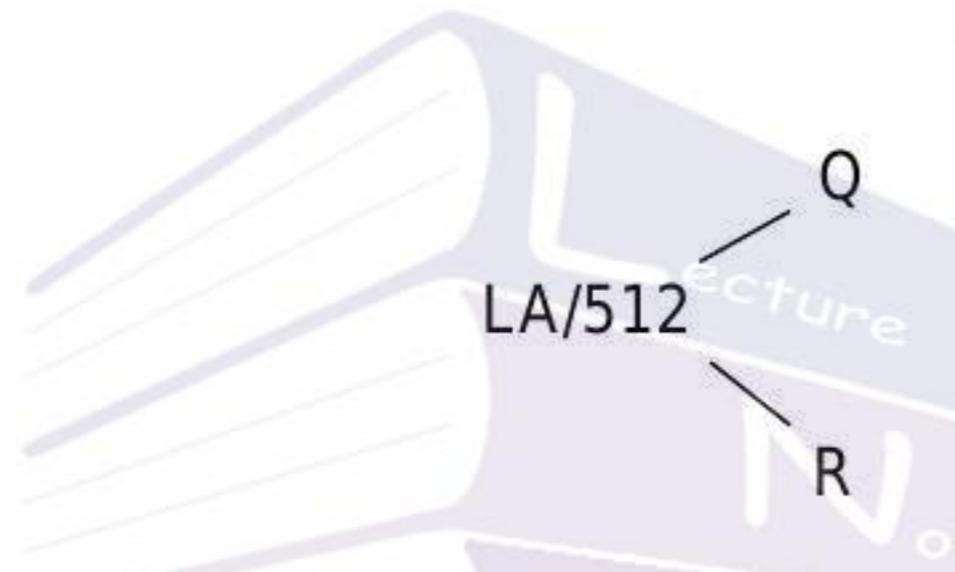
file permissions
file dates (create, access, write)
file owner, group, ACL
file size
file data blocks or pointers to file data blocks.

# Allocation Methods - Contiguous

- An allocation method refers to how disk blocks are allocated for files:
- The disk space should be utilized efficiently and files can access quickly
- **Contiguous allocation** – each file occupies set of contiguous blocks
  - Best performance in most cases
  - Simple – only starting location (block #) and length (number of blocks) are required
  - Problems include finding space for file, knowing file size, external fragmentation (suppose a file is deleted and in its place a new file to be inserted, due to lack of space it may not be possible), need for **compaction off-line (downtime)** or **on-line**

# Contiguous Allocation

- Mapping from logical to physical



Block to be accessed = Q +  
starting address

Displacement into block = R



directory		
file	start	length
count	0	2
tr	14	3
mail	19	6
list	28	4
f	6	2

# Extent-Based Systems

- Many newer file systems (i.e., Veritas File System) use a modified contiguous allocation scheme
- Extent-based file systems allocate disk blocks in extents
- An **extent** is a contiguous block of disks
  - Extents are allocated for file allocation
  - A file consists of one or more extents

LectureNotes.in

LectureNotes.in

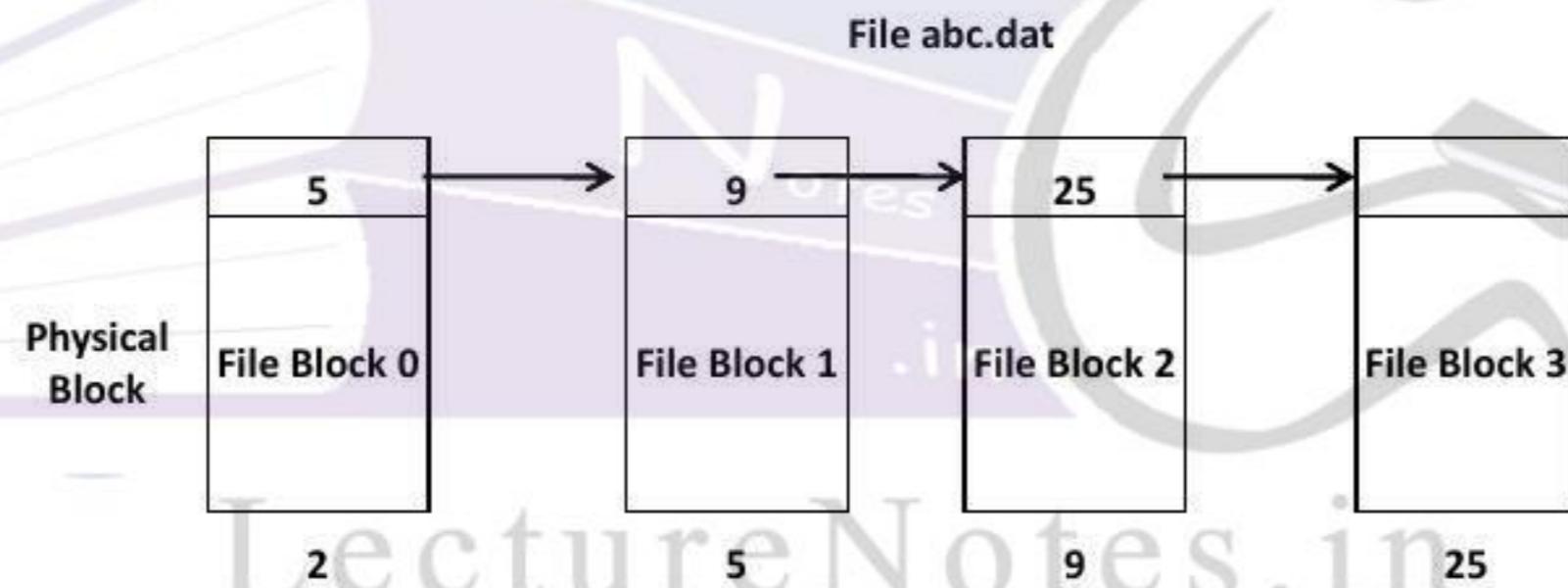
# Allocation Methods – Linked Non-contiguous

- **Linked allocation** – each file a linked list of blocks
  - File ends at nil pointer
  - No external fragmentation
  - Each block contains pointer to next block
  - No compaction, external fragmentation
  - Free space management system called when new block needed
  - Improve efficiency by clustering blocks into groups but increases internal fragmentation
  - Reliability can be a problem
  - Locating a block can take many I/Os and disk seeks

# Allocation Methods – Linked (Cont.)

## FAT (File Allocation Table) variation

- Beginning of volume has table, indexed by block number
- Much like a linked list, but faster on disk and cacheable
- New block allocation simple



Unlike contiguous allocation every disk block can be used

No space is lost due to fragmentation

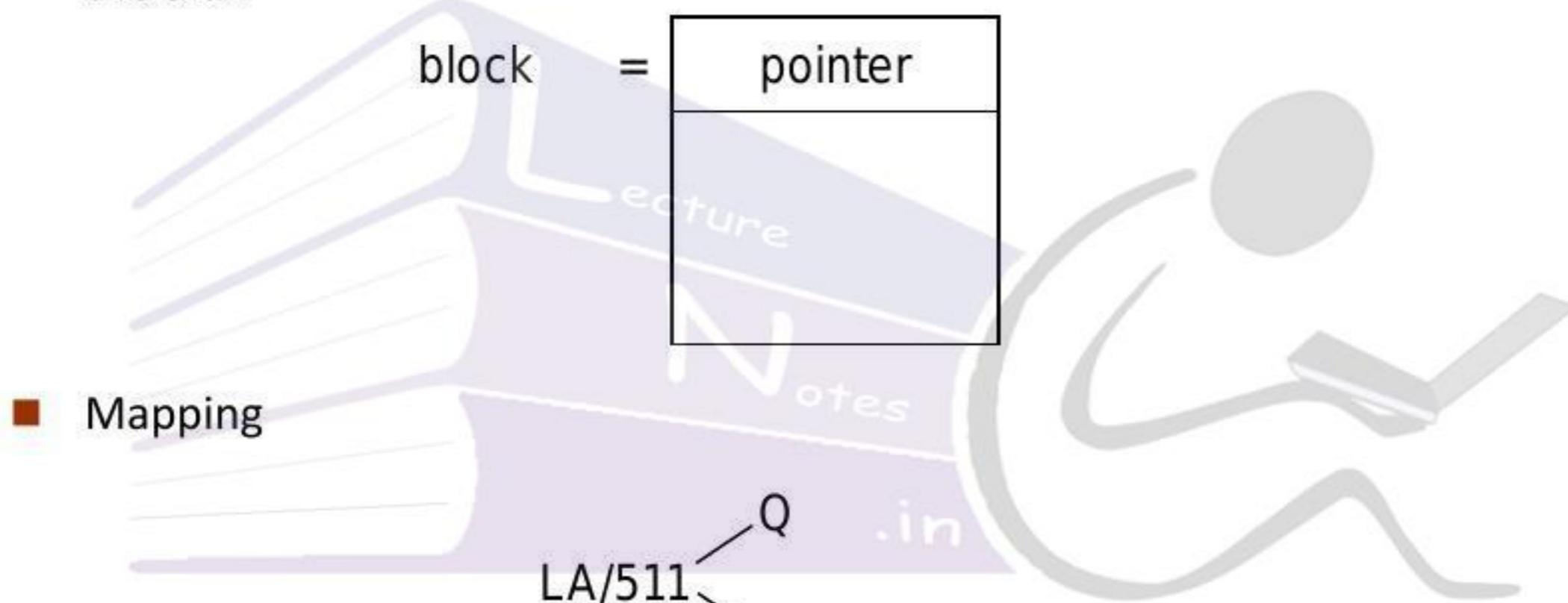
External fragmentation avoided

Disadvt: Random access is very slow

Pointer also takes some space

# Linked Allocation

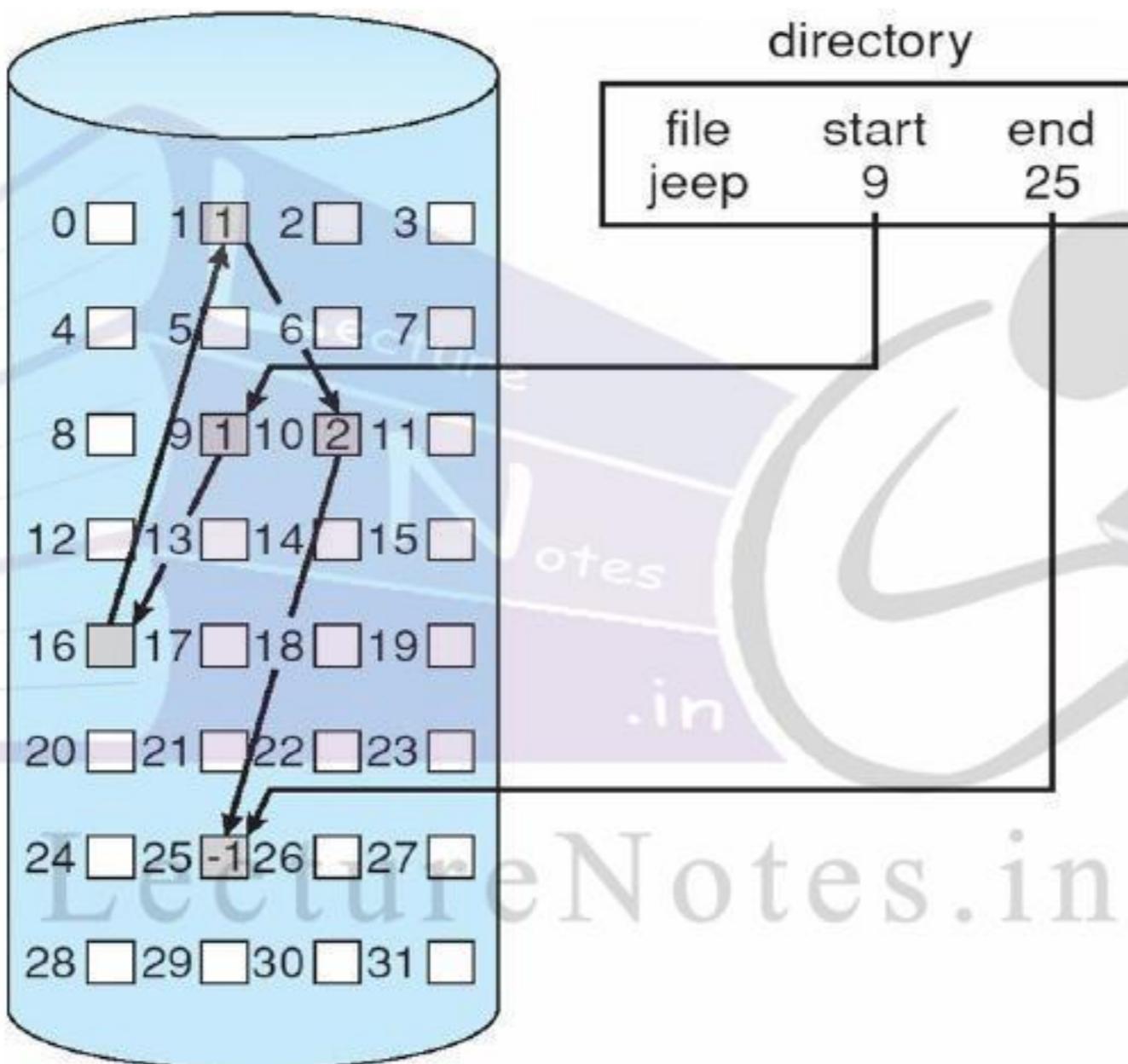
- Each file is a linked list of disk blocks: blocks may be scattered anywhere on the disk



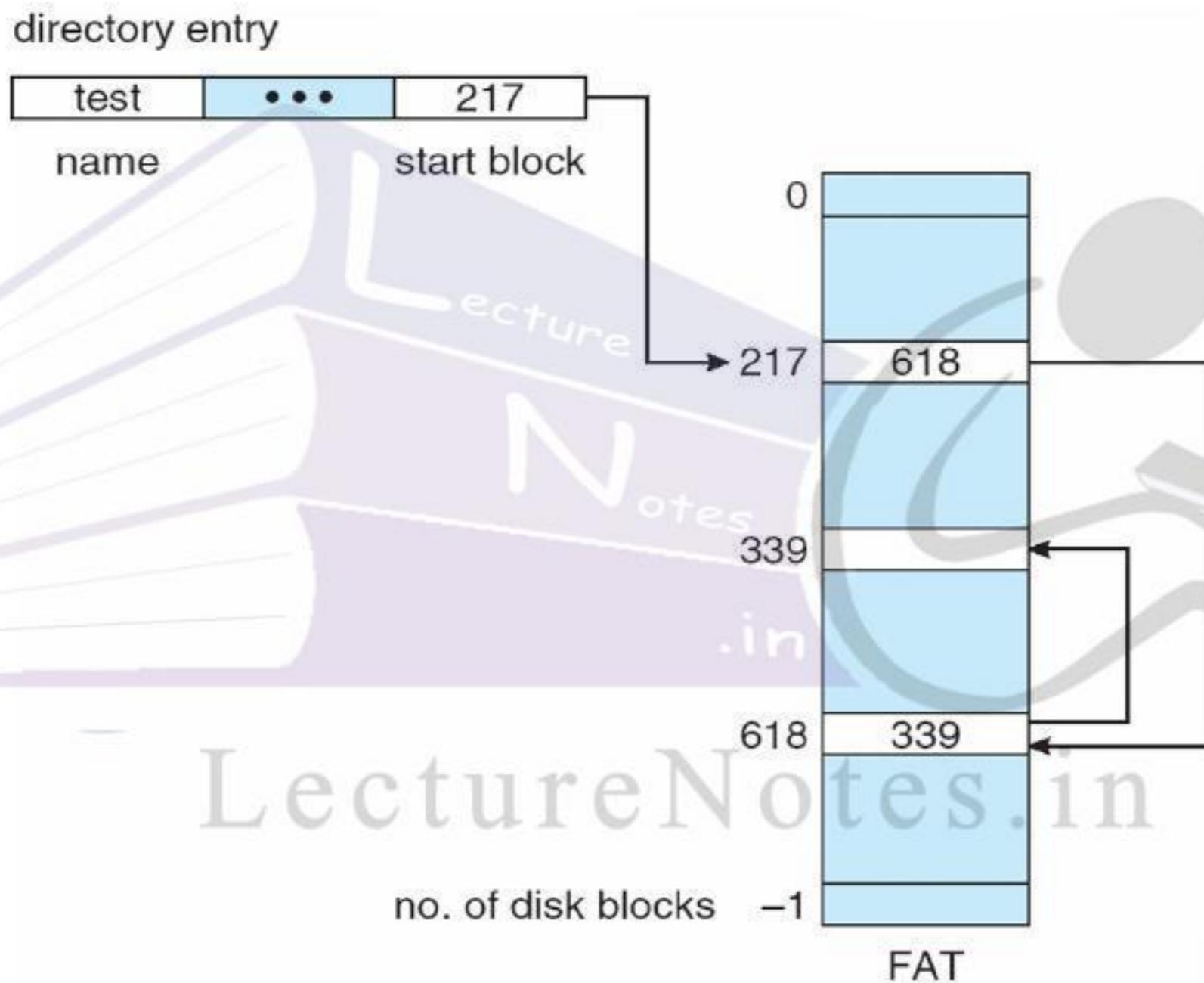
Block to be accessed is the Qth block in the linked chain of blocks representing the file.

$$\text{Displacement into block} = R + 1$$

# Linked Allocation



# File-Allocation Table

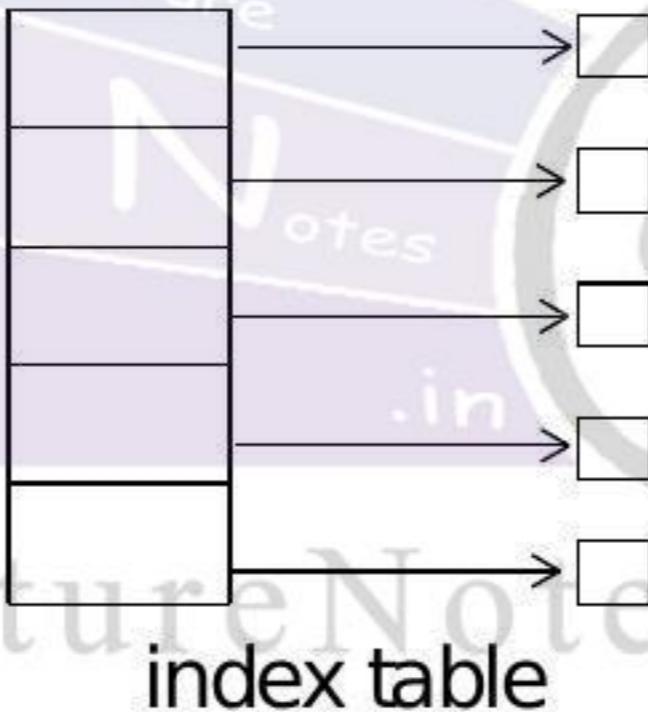


# Allocation Methods - Indexed

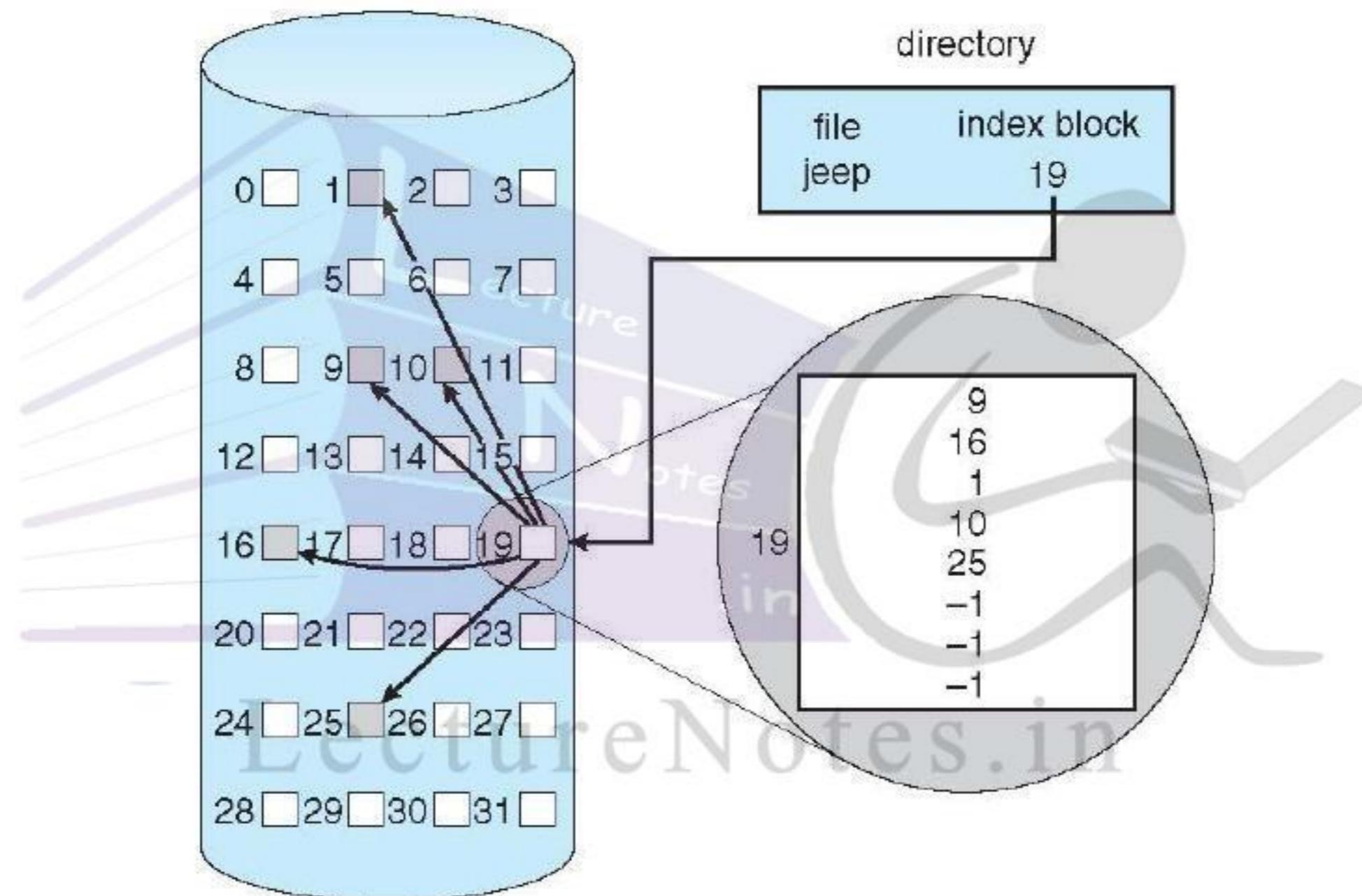
- **Indexed allocation**

- Each file has its own **index block(s)** of pointers to its data blocks

- Logical view



# Example of Indexed Allocation



# Indexed Allocation (Cont.)

- Need index table
- Random access
- Dynamic access without external fragmentation, but have overhead of index block
- Mapping from logical to physical in a file of maximum size of 256K bytes and block size of 512 bytes. We need only 1 block for index table

LA/512 → Q  
                  R

LectureNotes.in

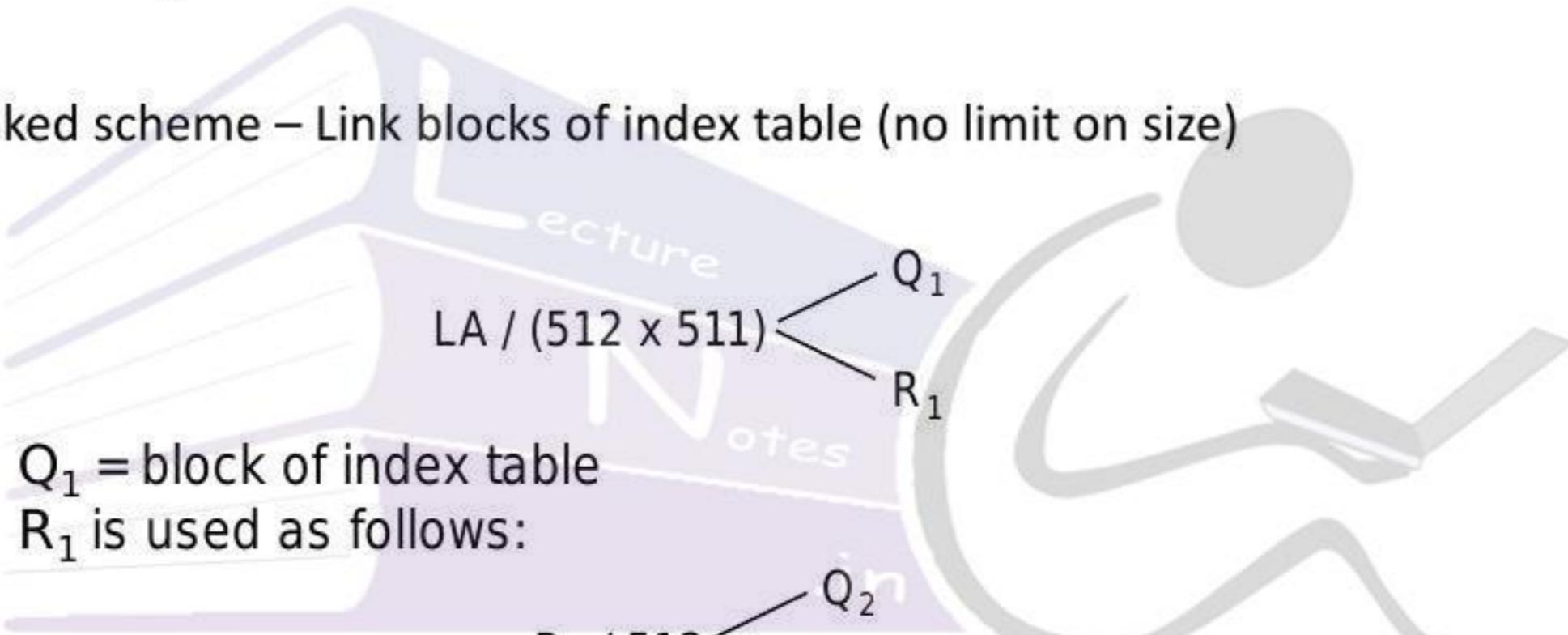
Q = displacement into index table

R = displacement into block

LectureNotes.in

# Indexed Allocation – Mapping (Cont.)

- Mapping from logical to physical in a file of unbounded length (block size of 512 words)
- Linked scheme – Link blocks of index table (no limit on size)



$Q_1$  = block of index table

$R_1$  is used as follows:

$R_1 / 512$        $Q_2$

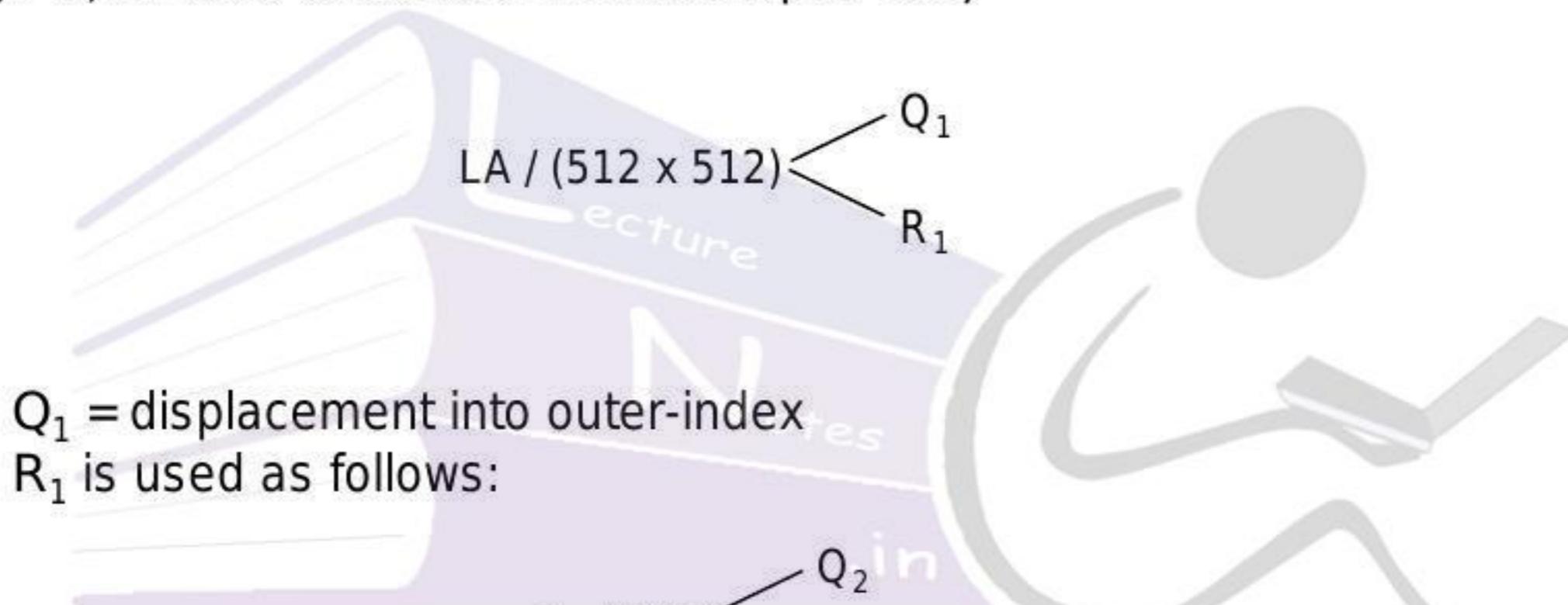
$R_2$

$Q_2$  = displacement into block of index table

$R_2$  displacement into block of file:

# Indexed Allocation – Mapping (Cont.)

- Two-level index (4K blocks could store 1,024 four-byte pointers in outer index -> 1,048,567 data blocks and file size of up to 4GB)



$Q_1$  = displacement into outer-index

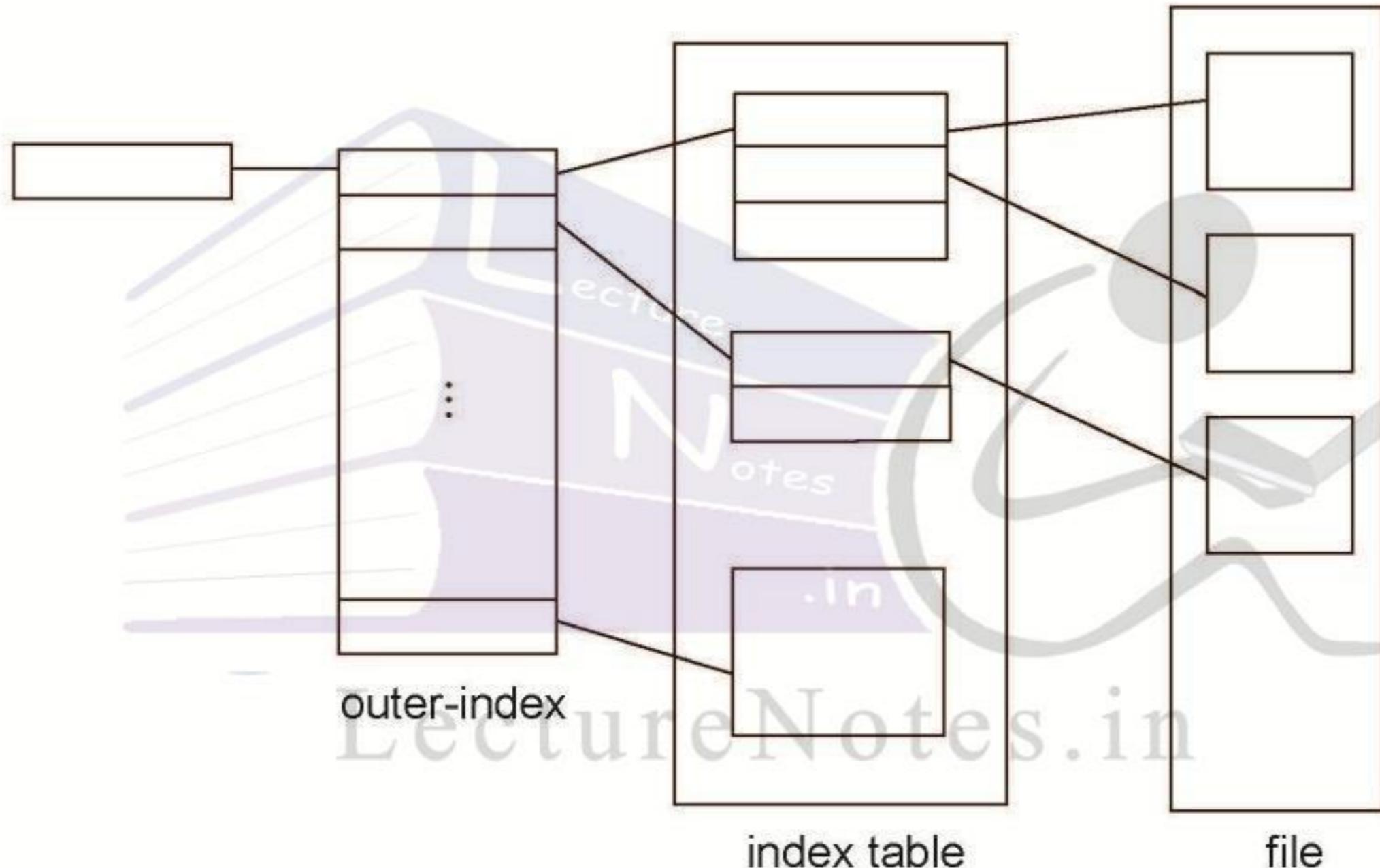
$R_1$  is used as follows:

A pointer labeled  $R_1$  points to a point on the stack. Two lines extend from this point to labels  $Q_1$  and  $R_1/512$ . Below the stack, the text "LectureNotes.in" is displayed.

$Q_2$  = displacement into block of index table

$R_2$  displacement into block of file:

# Indexed Allocation – Mapping (Cont.)



# Free-Space Management

- The main responsibility of the file system is to allocate block to each file and to keep track. Other responsibility is to know how many free space is there.
- File system maintains **free-space list** to track available blocks/clusters
  - (Using term “block” for simplicity)
- Bit vector** or **bit map** ( $n$  blocks) (for every block a bit is maintained)



$\text{bit}[i] = \begin{cases} 1 & \Rightarrow \text{block}[i] \text{ free} \\ 0 & \Rightarrow \text{block}[i] \text{ occupied} \end{cases}$

Block number calculation

(number of bits per word) \*  
(number of 0-value words) +  
offset of first 1 bit

CPUs have instructions to return offset within word of first “1” bit

# LectureNotes.in Free-Space Management (Cont.)

- Bit map requires extra space

- Example:

block size = 4KB =  $2^{12}$  bytes

disk size =  $2^{40}$  bytes (1 terabyte)

$n = 2^{40}/2^{12} = 2^{28}$  bits (or 32MB)

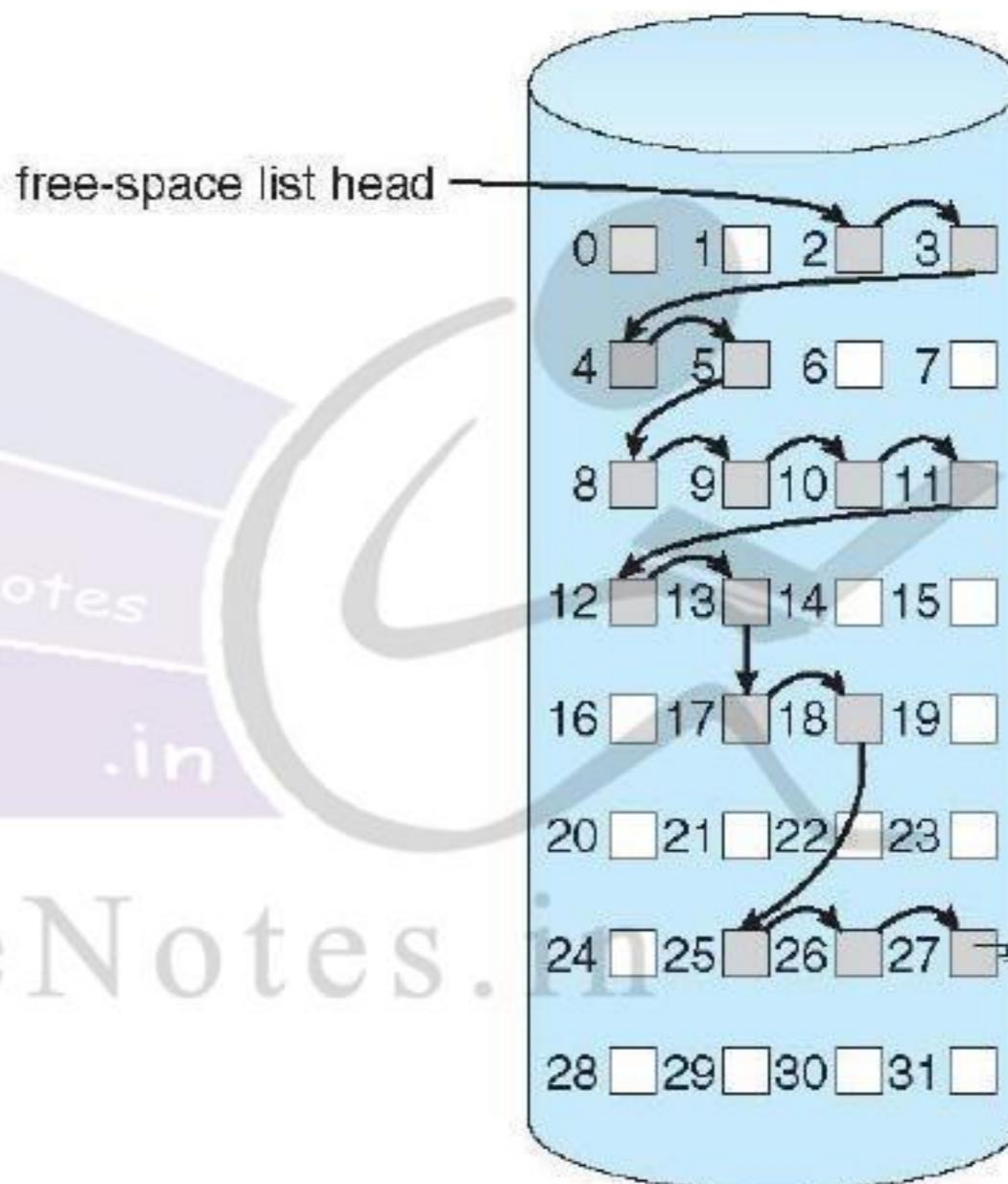
if clusters of 4 blocks -> 8MB of memory

- Easy to get contiguous files

# Linked Free Space List on Disk

- Linked list (free list)

- Cannot get contiguous space easily
- No waste of space
- No need to traverse the entire list (if # free blocks recorded)



# Free-Space Management (Cont.)

## ■ Grouping

- Modify linked list to store address of next  $n-1$  free blocks in first free block, plus a pointer to next block that contains free-block-pointers (like this one)

## ■ Counting

- Because space is frequently contiguously used and freed, with contiguous-allocation allocation, extents, or clustering
  - ▶ Keep address of first free block and count of following free blocks
  - ▶ Free space list then has entries containing addresses and counts

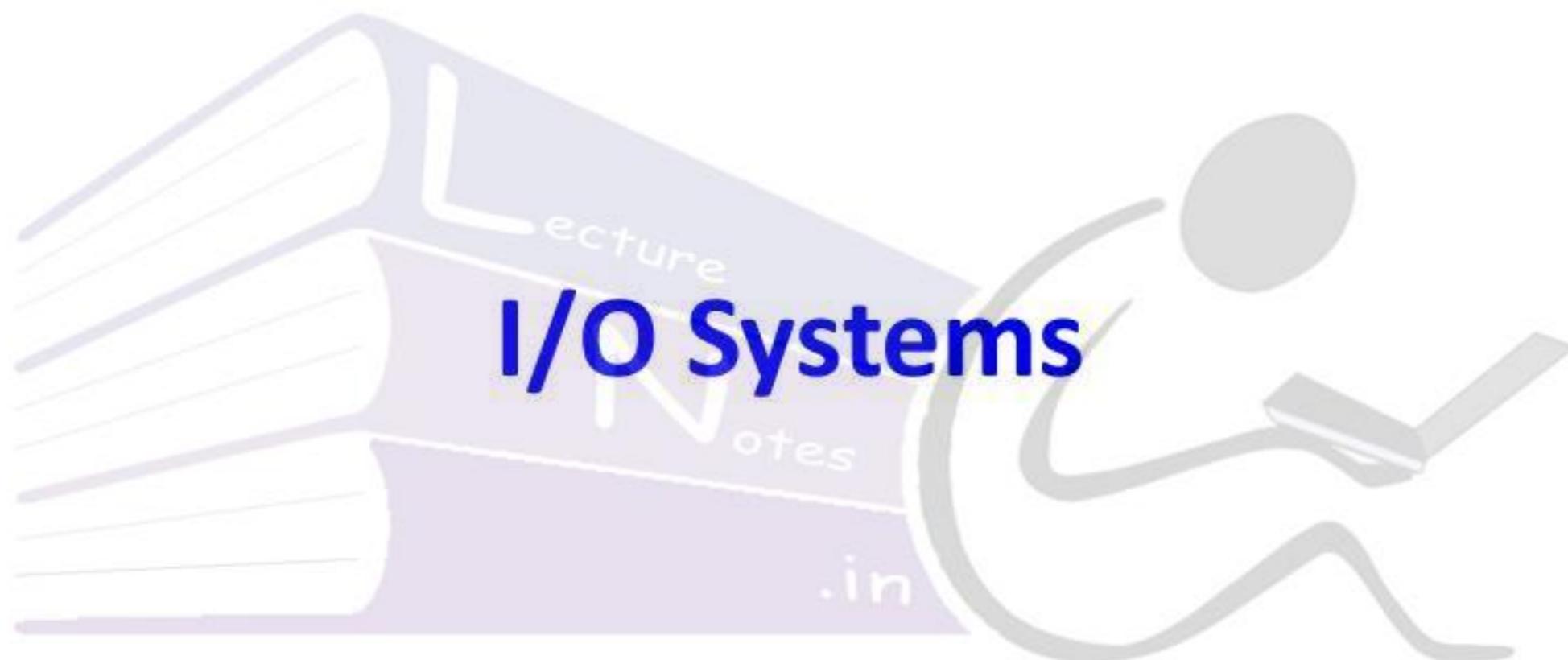
LectureNotes.in

LectureNotes.in

# Free-Space Management (Cont.)

## ■ Space Maps

- Used in **ZFS**
- Consider meta-data I/O on very large file systems
  - ▶ Full data structures like bit maps couldn't fit in memory -> thousands of I/Os
- Divides device space into **metaslab** units and manages metaslabs
  - ▶ Given volume can contain hundreds of metaslabs
- Each metaslab has associated space map
  - ▶ Uses counting algorithm
- But records to log file rather than file system
  - ▶ Log of all block activity, in time order, in counting format
- Metaslab activity -> load space map into memory in balanced-tree structure, indexed by offset
  - ▶ Replay log into that structure
  - ▶ Combine contiguous free blocks into single entry



LectureNotes.in

LectureNotes.in

By Ajay Kumar Jena, KIIT Deemed to be  
University, Bhubaneswar

# LectureNotes.in

## Overview

- I/O management is a major component of operating system design and operation
- It transfers the information between the internal storage (CPU, MM) and the External storage (HDD)
  - Important aspect of computer operation
  - I/O devices vary greatly
  - Various methods to control them
  - Performance management
  - New types of devices frequent
- Ports, busses, device controllers connect to various devices
- **Device drivers** encapsulate device details
  - Present uniform device-access interface to I/O subsystem



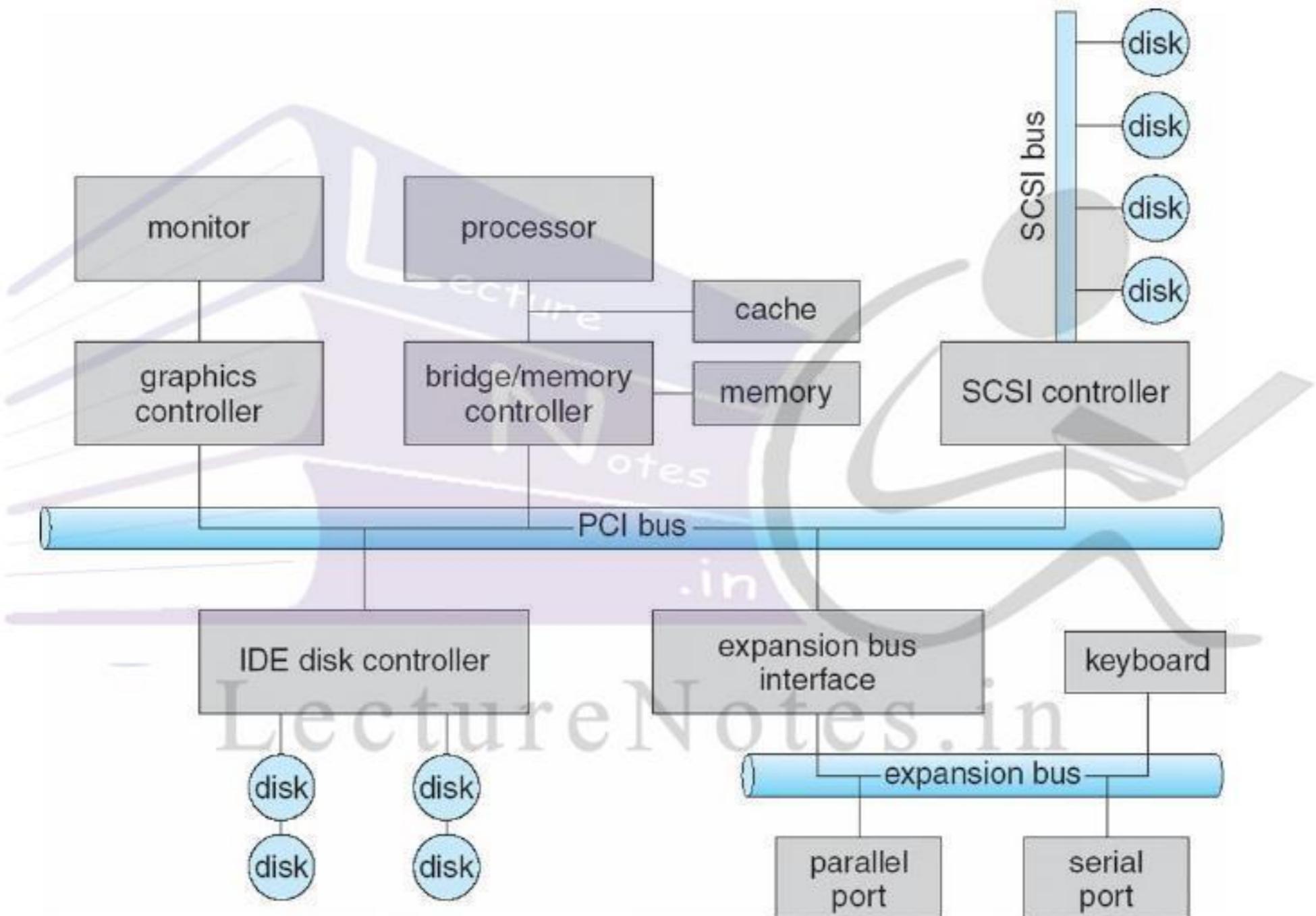
- Computer operates on a great many kind of devices
  - General category storage device: HDD, tapes
  - Transmission devices: N/W card, modem, router etc.
  - Human interface devices: KB, Mouse, screen etc.
- A device communicates with the system by sending signals over a cable or even through air
- The device communicates with the machine via a connection point (port). Serial, Parallel, USB
- A bus is a set of wires and a rigidly defined protocol that specifies a set of messages that can be sent on wires. Buses are used widely in computer architecture.
- **Device drivers** encapsulate device details
  - Present uniform device-access interface to I/O subsystem

# Lecture Notes in I/O Hardware

- Incredible variety of I/O devices
  - Storage      Transmission      Human-interface
- Common concepts – signals from I/O devices interface with computer for which special communication tracks are necessary.
  - **Port** – connection point for device

The device communicates with the machine via a communication point (port).
  - **Bus** – It is a set of wires and a rigidly defined protocol that specifies a set of messages that can be sent on wires. Buses are used widely in computer architecture.
    - ▶ PCI bus common in PCs and servers, PCI Express (**PCIe**)
    - ▶ **expansion bus** connects relatively slow devices
  - **Controller (host adapter)** – electronics that operate port, bus, device
    - ▶ Sometimes integrated
    - ▶ Sometimes separate circuit board (host adapter)
    - ▶ Contains processor, microcode, private memory, bus controller, etc
      - Some talk to per-device controller with bus controller, microcode, memory, etc

# A Typical PC Bus Structure



## A typical Peripheral Component Interconnection Structure

By Ajay Kumar Jena, KIIT Deemed to be

University, Bhubaneswar

# I/O Hardware (Cont.)

- I/O instructions control devices
- Devices usually have registers where device driver places commands, addresses, and data to write, or read data from registers after command execution
  - Data-in register, data-out register, status register, control register
  - Typically 1-4 bytes, or FIFO buffer
- Devices have addresses, used by
  - Direct I/O instructions
  - **Memory-mapped I/O**
    - ▶ Device data and command registers mapped to processor address space
    - ▶ Especially for large address spaces (graphics)

- I/O devices have two components.
  - Mechanical component (the device itself i.e. mouse, KB)
  - Electronic component of devices which controls the device is called device controller. The circuit inside the KB and mouse is the controller)
  - Tasks of controller
    - ▶ It converts serial bit stream to block of bytes
    - ▶ Performs error correction if necessary
- Ports, busses, device controllers connect to various devices
- **Device drivers** encapsulate device details
  - Present uniform device-access interface to I/O subsystem

# Device I/O Port Locations on PCs (partial)

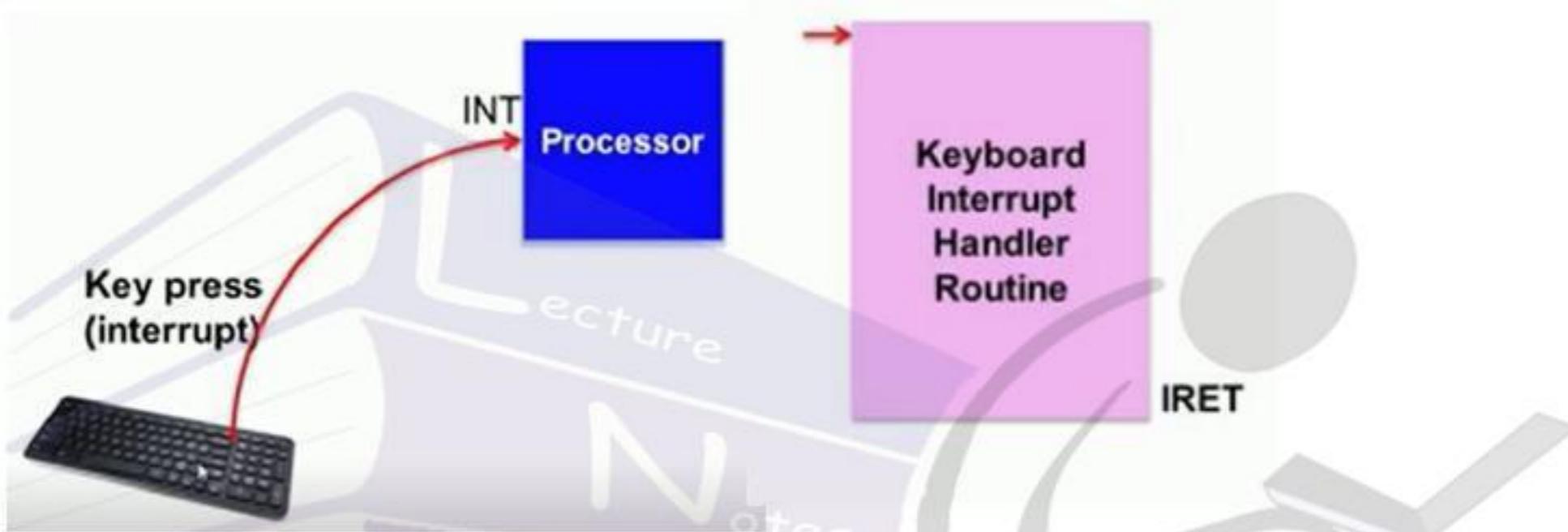
I/O address range (hexadecimal)	device
000–00F	DMA controller
020–021	interrupt controller
040–043	timer
200–20F	game controller
2F8–2FF	serial port (secondary)
320–32F	hard-disk controller
378–37F	parallel port
3D0–3DF	graphics controller
3F0–3F7	diskette-drive controller
3F8–3FF	serial port (primary)

# Lecture Notes in Interrupts

- Data transfer between the I/O and the system can be done
  - Programmed I/O
    - ▶ CPU has to wait for the I/O module which is a disadvantage
  - Interrupt controlled I/O
    - ▶ CPU does not wait for I/O. When I/O ready it interrupts CPU
  - DMA
- The H/w mechanism that enables a device to notify the CPU is called as interrupt.
- CPU **Interrupt-request line** triggered by I/O device
  - Checked by processor after each instruction
- **Interrupt handler** receives interrupts
  - **Maskable** to ignore or delay some interrupts
- **Interrupt vector** to dispatch interrupt to correct handler
  - Context switch at start and end
  - Based on priority
  - Some **nonmaskable**
  - Interrupt chaining if more than one device at same interrupt number

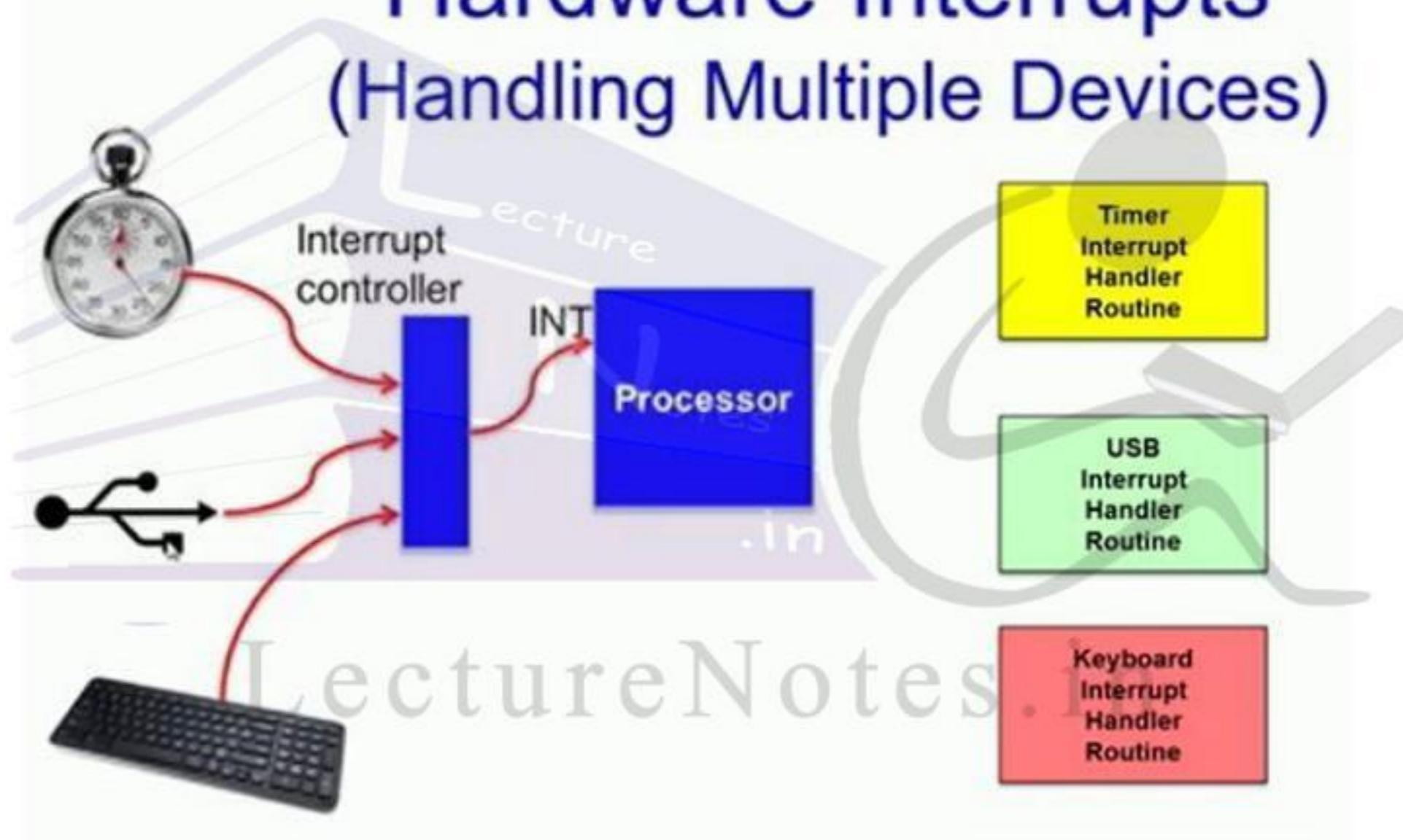
# LectureNotes.in

# HARDWARE INTERRUPTS

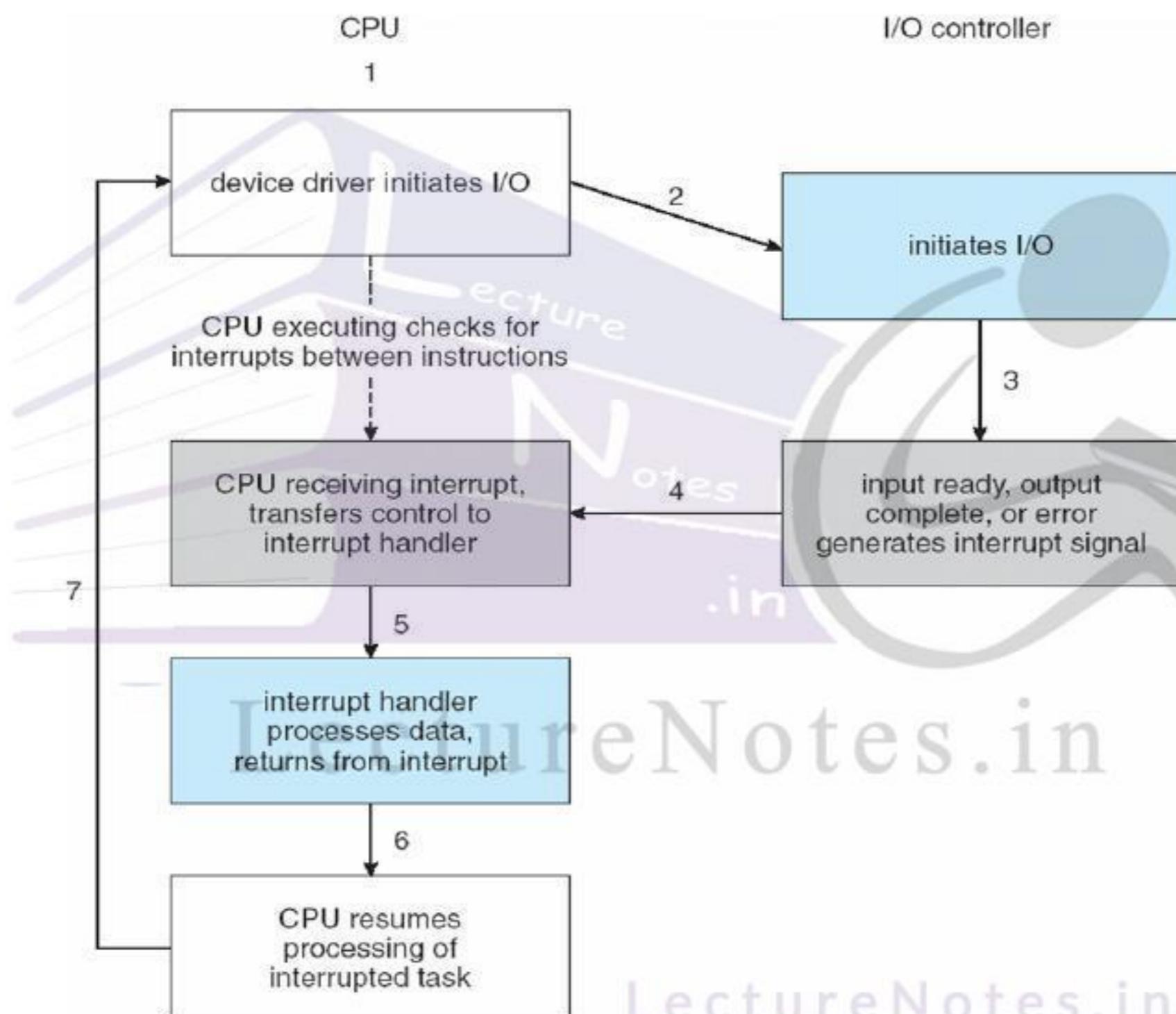


Now-a-days processors have dedicated pins known as INT pin or INTR pin. Devices like KB connects to the processor by INT pin. When a key is pressed in the KB an interrupt is generated in the CPU. The KB interrupt handler routine will be invoked.

# Hardware Interrupts (Handling Multiple Devices)

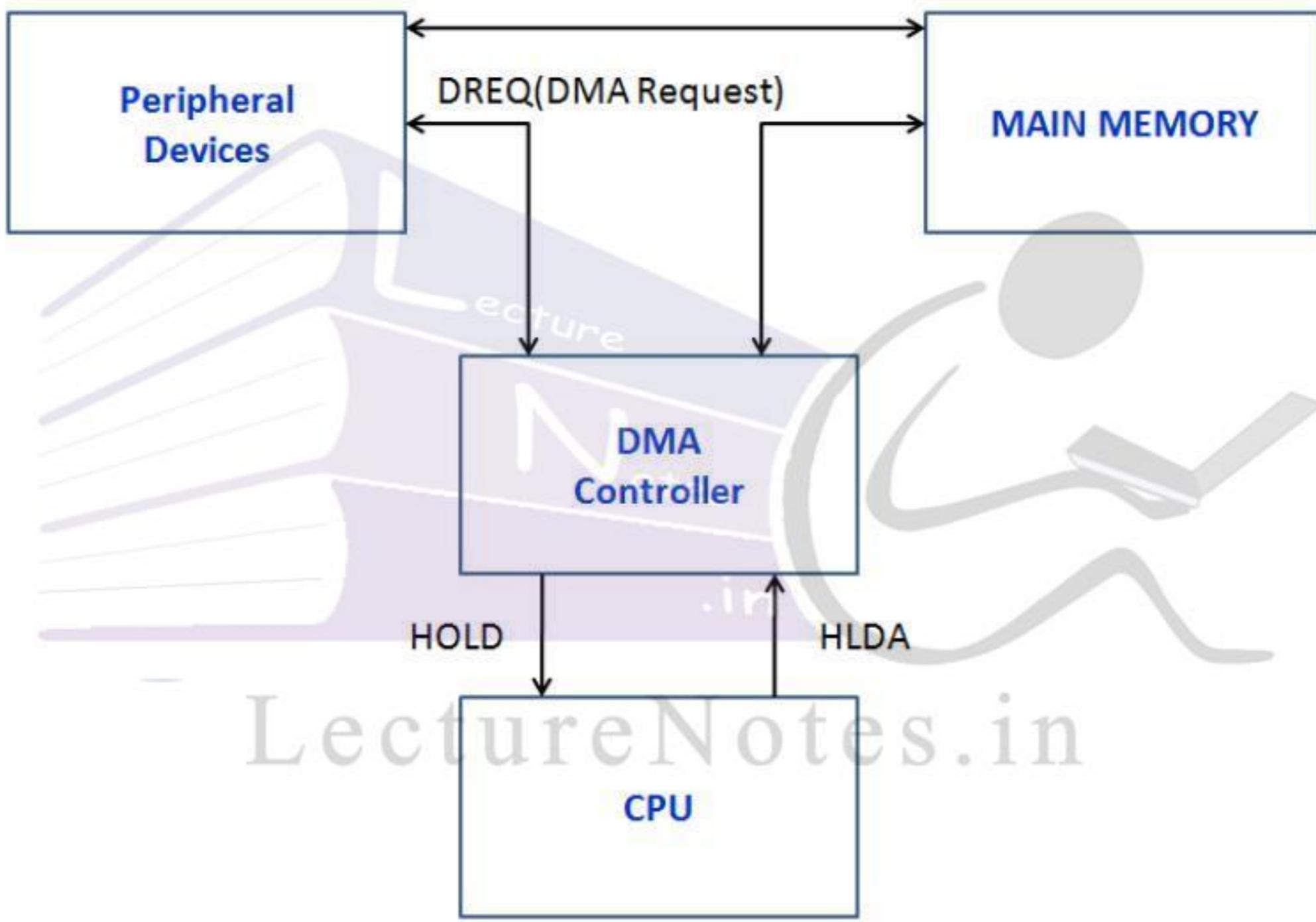


# Interrupt-Driven I/O Cycle

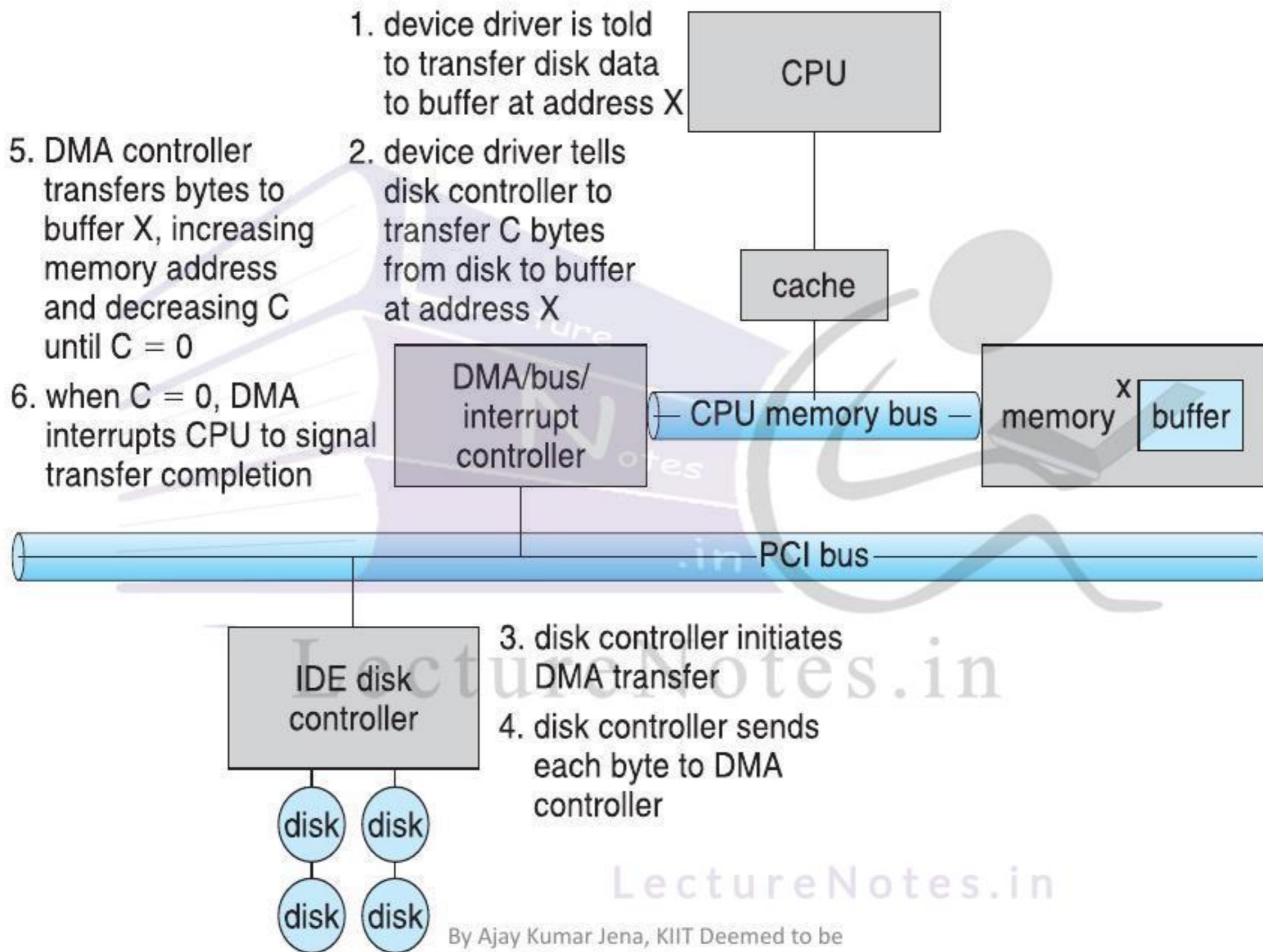


# Direct Memory Access

- Used to avoid **programmed I/O** (one byte at a time) for large data movement
- The IO devices will not be synchronized with CPU for transferring the data as IO devices are very slow.
- Requires **DMA** controller is having four parameters i.e. source address, destination address, no of blocks, Read/Write
- Bypasses CPU to transfer data directly between I/O device and memory
- OS writes DMA command block into memory
  - Source and destination addresses
  - Read or write mode
  - Count of bytes
  - Writes location of command block to DMA controller
  - Bus mastering of DMA controller – grabs bus from CPU
    - ▶ **Cycle stealing** from CPU but still much more efficient
  - When done, interrupts to signal completion



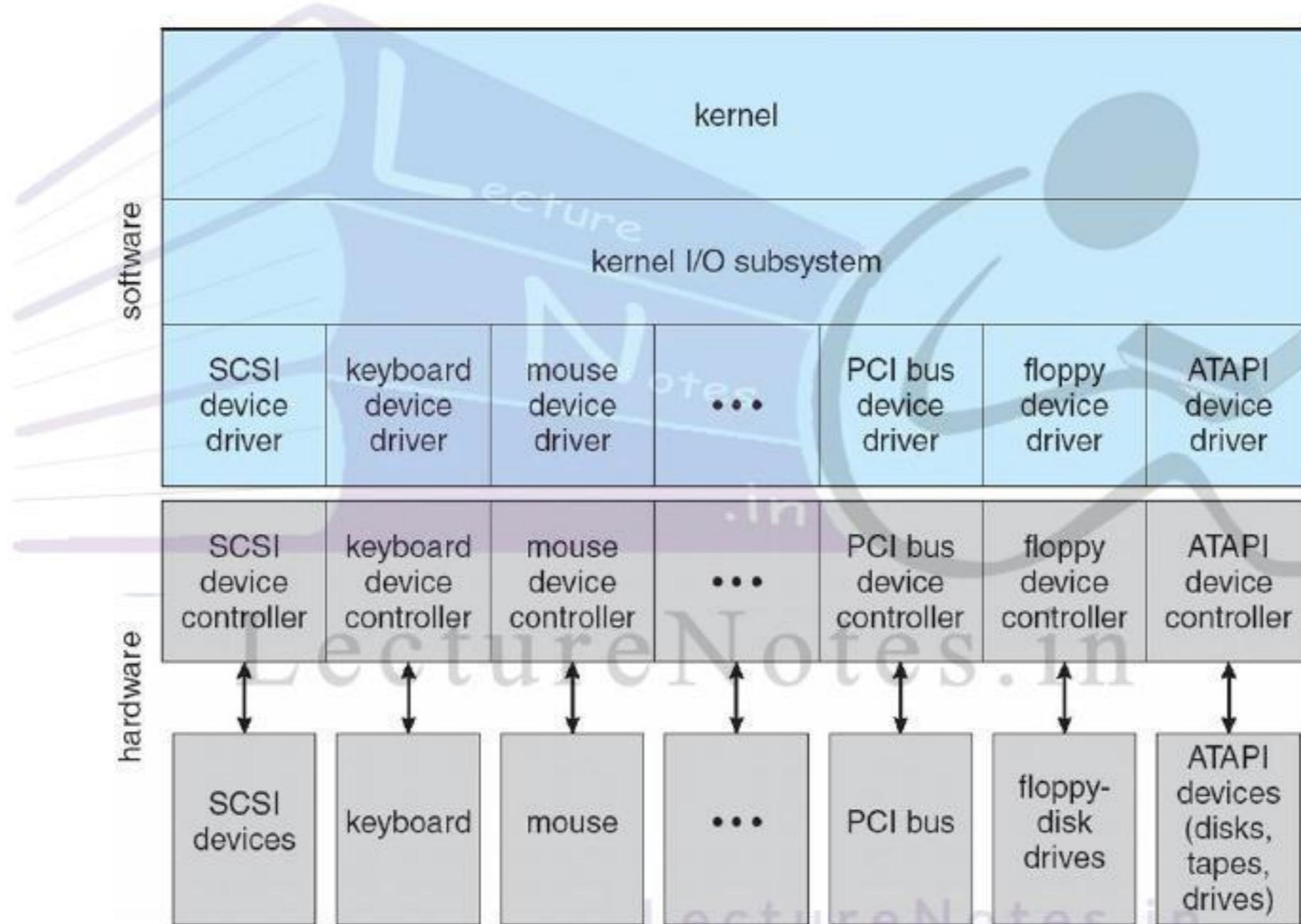
# Six Step Process to Perform DMA Transfer



# LectureNotes.in Application I/O Interface

- I/O system calls encapsulate device behaviors in generic classes
- Device-driver layer hides differences among I/O controllers from kernel
- New devices talking already-implemented protocols need no extra work
- Each OS has its own I/O subsystem structures and device driver frameworks
- Devices vary in many dimensions
  - **Character-stream or block**
  - **Sequential or random-access**
  - **Synchronous or asynchronous (or both)**
  - **Sharable or dedicated**
  - **Speed of operation**
  - **read-write, read only, or write only**

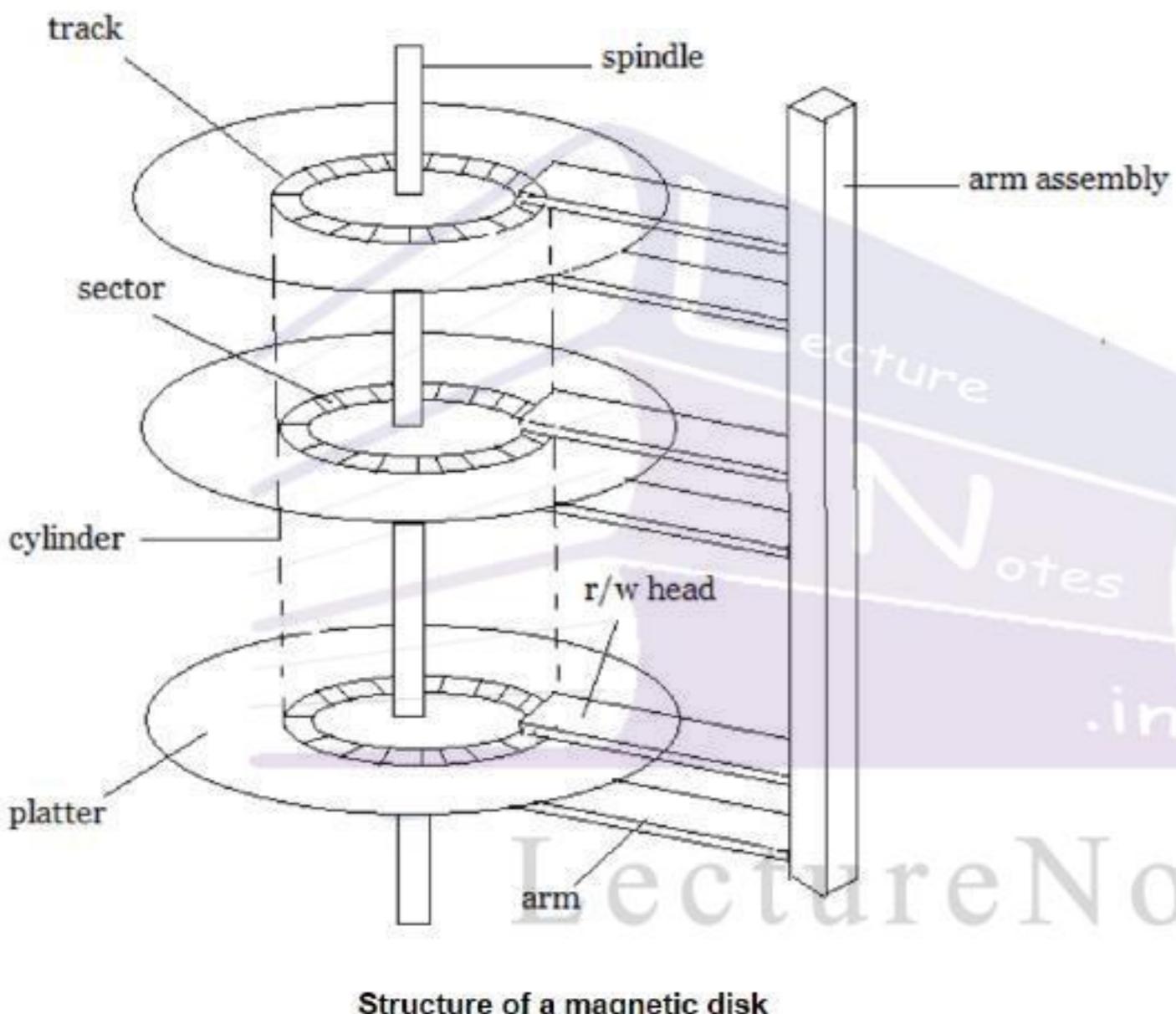
# A Kernel I/O Structure



# Characteristics of I/O Devices

aspect	variation	example
data-transfer mode	character block	terminal disk
access method	sequential random	modem CD-ROM
transfer schedule	synchronous asynchronous	tape keyboard
sharing	dedicated sharable	tape keyboard
device speed	latency seek time transfer rate delay between operations	
I/O direction	read only write only read-write	CD-ROM graphics controller disk

# Disk Structure



- Secondary storage devices are those devices whose memory is non volatile
- Secondary storage is also called auxiliary storage.
- Secondary storage is less expensive when compared to primary memory like RAMs.
- The speed of the secondary storage is also lesser than that of primary storage.
- Hence, the data which is less frequently accessed is kept in the secondary storage.
- A few examples are magnetic disks, magnetic tapes, removable thumb drives etc.

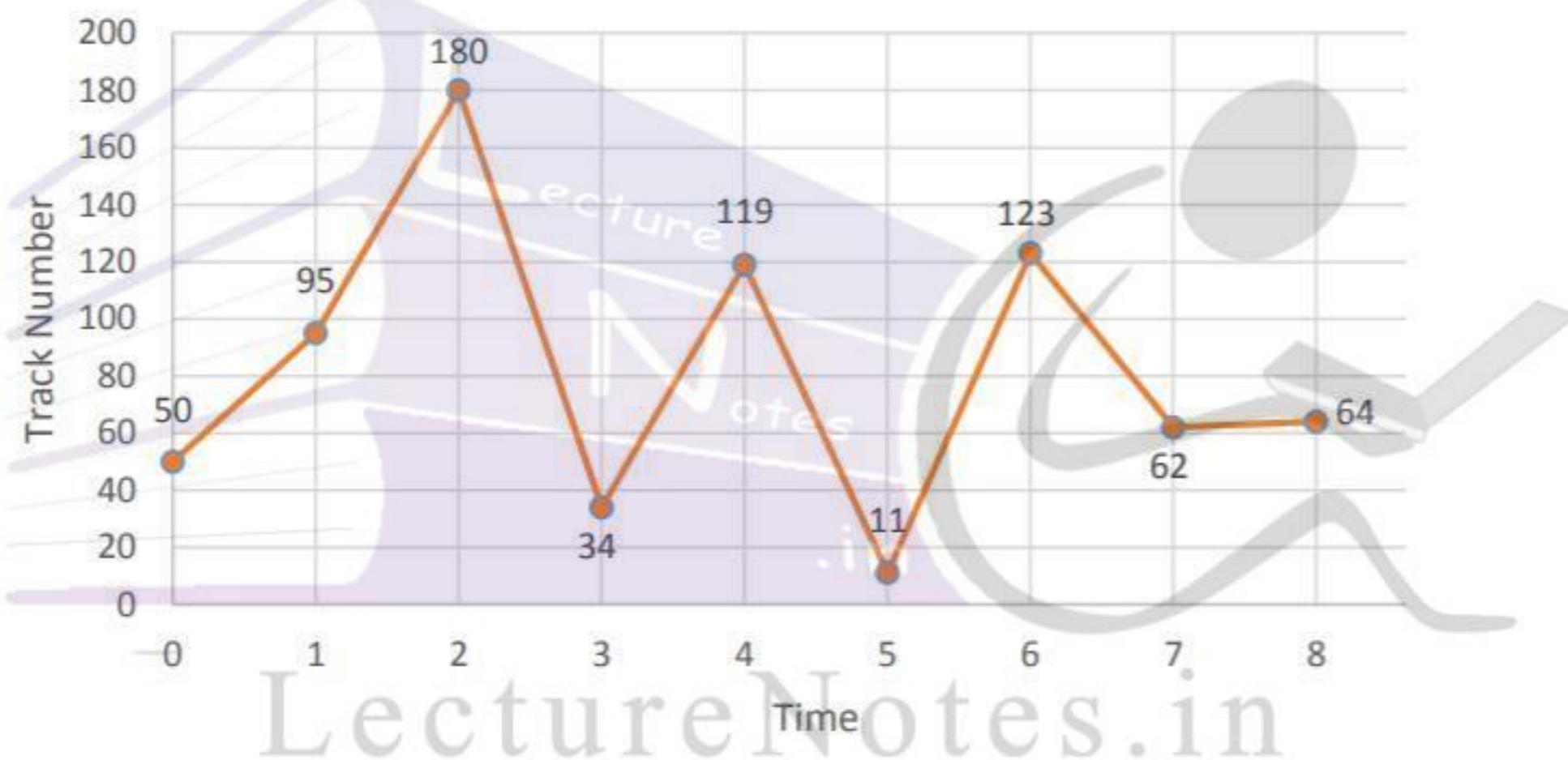
# Disk Scheduling

- A magnetic disk contains several **platters**. Each platter is divided into circular shaped **tracks**.
- The length of the tracks near the centre is less than the length of the tracks farther from the centre.
- Each track is further divided into **sectors**.
- Tracks of the same distance from centre form a cylinder.
- A read-write head is used to read data from a sector of the magnetic disk.
- The speed of the disk is measured as two parts:
- **Transfer rate:** This is the rate at which the data moves from disk to the computer.
- **Random access time:** It is the sum of the seek time and rotational latency.
- **Seek time** is the time taken by the arm to move to the required track.
- **Rotational latency** is defined as the time taken by the arm to reach the required sector in the track.

## FCFS Disk scheduling algorithm:

Given the following queue -- 95, 180, 34, 119, 11, 123, 62, 64 with the Read-write head initially at the track 50 and the tail track being at 199.

**Disk Scheduling Algorithms : FCFS**

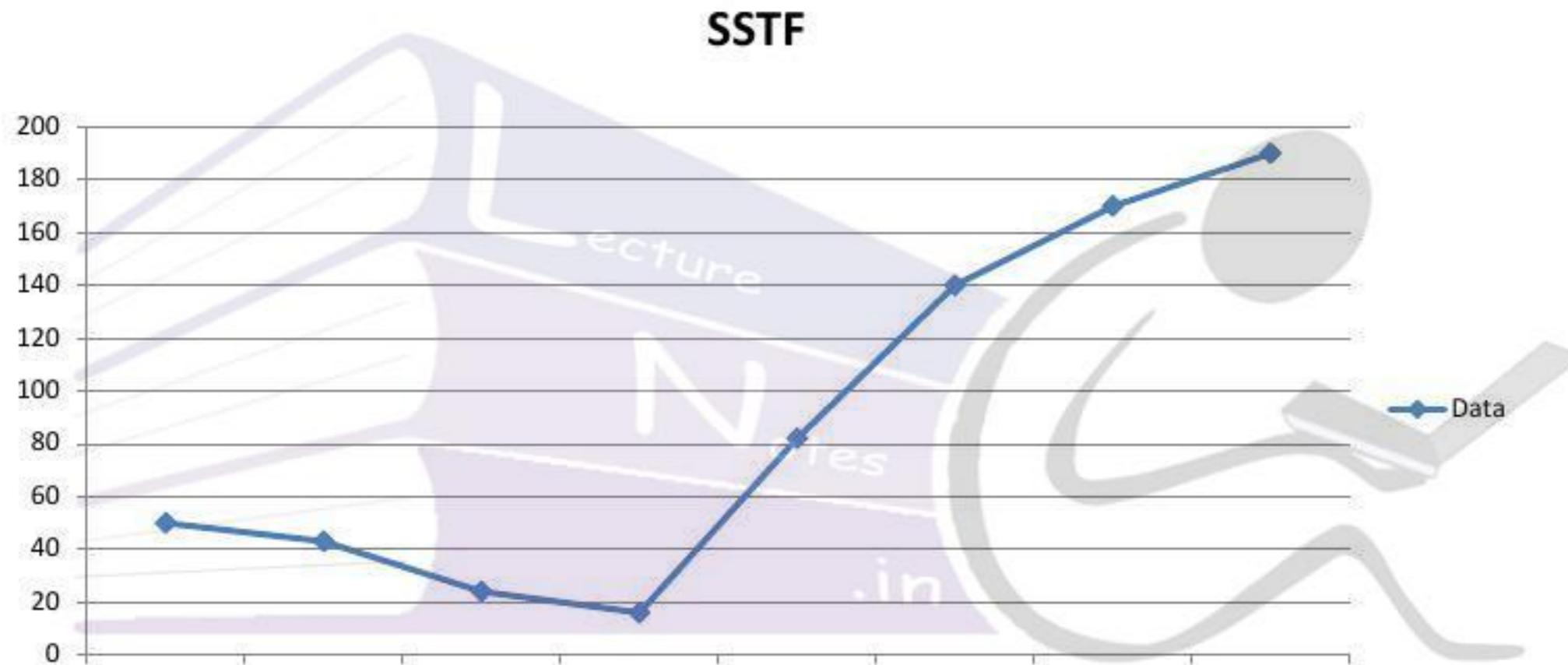


For this case it went from 50 to 95 to 180 and so on. From 50 to 95 it moved 45 tracks. If you tally up the total number of tracks you will find how many tracks it had to go through before finishing the entire request. The total head movement of 644 tracks.

$$\begin{aligned}
 \text{THM} &= (180 - 50) + (180 - 34) + (119 - 34) + (119 - 11) + (123 - 11) + (123 - 62) + (64 - 62) \\
 &= 130 + 146 + 85 + 108 + 112 + 61 + 2 = 644 \text{ tracks}
 \end{aligned}$$

## SSTF (Shortest Seek Time First ) Disk scheduling algorithm:

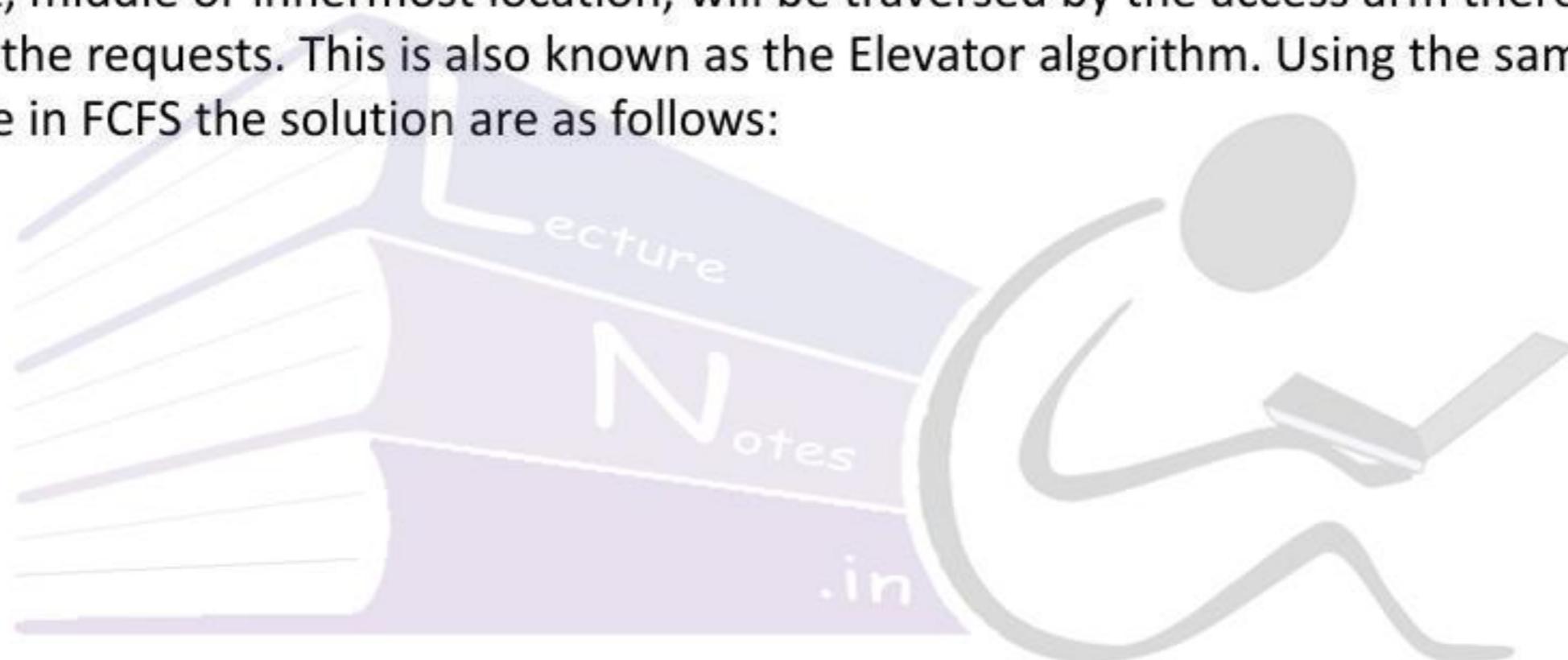
Given the following Request queue – 82, 170, 43, 140, 24, 16, 190 with the Read-write head initially at the track 50 and the tail track being at 199.



Here the position which is closest to the current head position is chosen first. The request is serviced according to next shortest distance. Find the shortest seek time. To move the desired track. From 50, the next shortest distance would be 43. From 43 to 24 track. By moving accordingly we have 43,24,16,82,140,170,190.  
 $THM = (50-16) + (190-16) = 34+174=208$  tracks

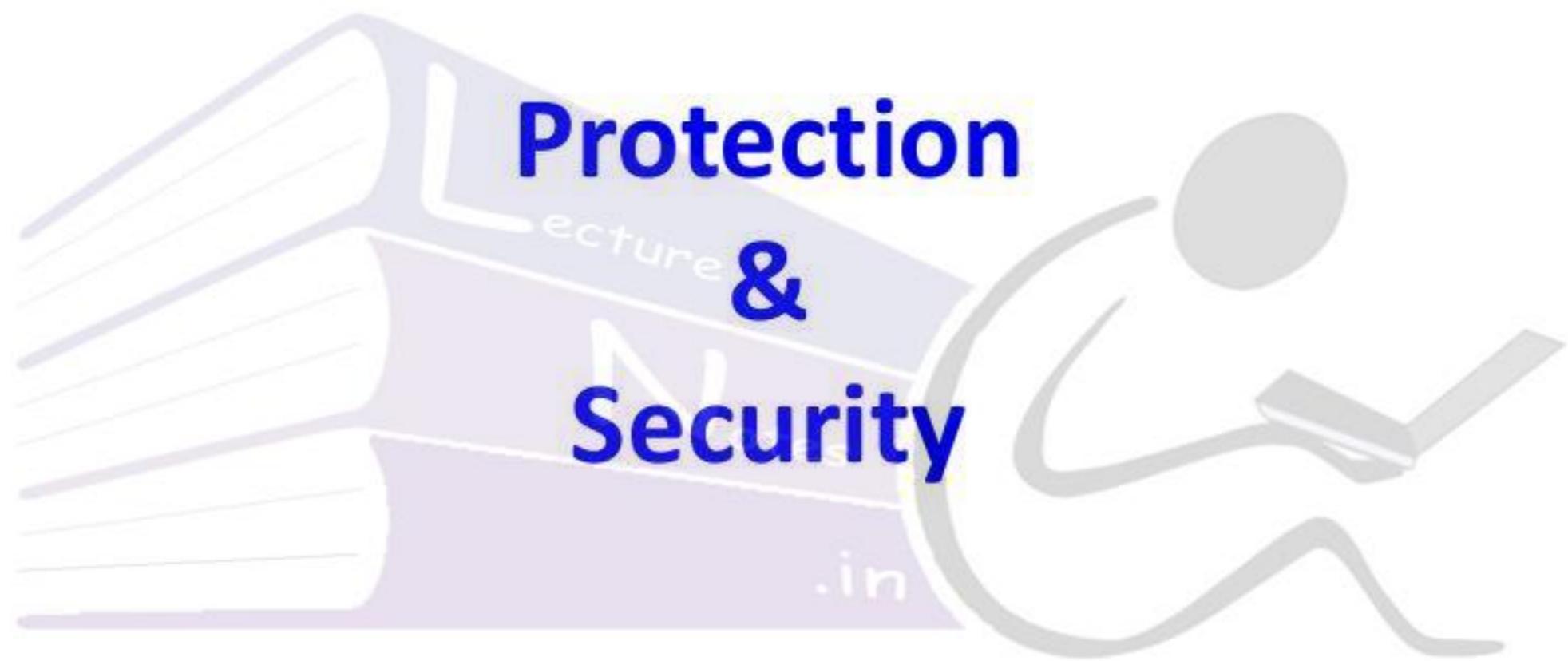
## SCAN Scheduling Algorithm

This algorithm is performed by moving the R/W head back-and-forth to the innermost and outermost track. As it scans the tracks from end to end, it processes all the requests found in the direction it is headed. This will ensure that all track requests, whether in the outermost, middle or innermost location, will be traversed by the access arm thereby finding all the requests. This is also known as the Elevator algorithm. Using the same sets of example in FCFS the solution are as follows:



LectureNotes.in

LectureNotes.in



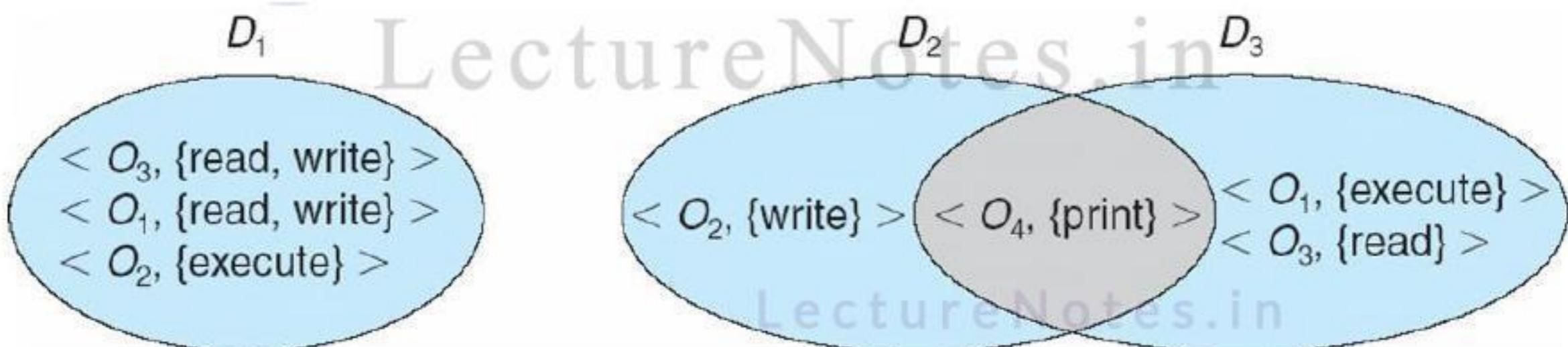
LectureNotes.in

# LectureNotes.in Protection

- A computer system is a collection of processes and objects.
- HW objects such as CPU, memory segment, disk, tape drives.
- SW objects like files, programs, semaphores
- Different operations are performed on different objects.
- CPU can only execute, Memory segment read, write tape drives can have read, write, rewound
- Data files can be created, opened, read, written, closed, deleted
- Program files can read, written, executed and deleted.
- A process should be allowed to access only those resources for which it has authorization.
- So the process operates within a protection domain which specifies that the resources the process may access.
- Each domain is defined a set of objects and the type of operations that may be invoked on each objects.
- The ability to execute an operation on the object is an access right.

# Lecture Notes in Domain Structure

- Protection Domain
  - User (user id)
  - Process (process id)
  - Procedures
- Access-right =  $\langle \text{object-name}, \text{rights-set} \rangle$  It is a order pair where *rights-set* is a subset of all valid operations that can be performed on the object
- Domain = collection of set of access-rights



# LectureNotes.in Access Matrix

- View protection as a matrix (**access matrix**)
- Rows represent domains (Users/groups/subjects)
- Columns represent objects (Resources)
- **Access (i, j)** is the set of operations that a process executing in Domain  $D_i$  can invoke on Object  $O_j$ . **Access control matrix** defines what access rights user u has for object O.

		object	$F_1$	$F_2$	$F_3$	printer
		domain				
User1	$D_1$	read		read		
User2	$D_2$				print	
User3	$D_3$		read	execute		
User4	$D_4$	read write			read write	

# Use of Access Matrix

- If a process in Domain  $D_i$ , tries to do “op” on object  $O_j$ , then “op” must be in the access matrix
- User who creates object can define access column for that object
- Can be expanded to dynamic protection
  - Operations to add, delete access rights
  - Special access rights:
    - ▶ *owner of  $O_i$*
    - ▶ *copy op from  $O_i$  to  $O_j$  (denoted by “\*”)*
    - ▶ *control –  $D_i$  can modify  $D_j$  access rights*
    - ▶ *transfer – switch from domain  $D_i$  to  $D_j$*
  - *Copy* and *Owner* applicable to an object
  - *Control* applicable to domain object

# Use of Access Matrix (Cont.)

- **Access matrix** design separates mechanism from policy
  - Mechanism
    - ▶ Operating system provides access-matrix + rules
    - ▶ If ensures that the matrix is only manipulated by authorized agents and that rules are strictly enforced
  - Policy
    - ▶ User dictates policy
    - ▶ Who can access what object and in what mode
- But doesn't solve the general confinement problem

# Access Matrix of Figure A with Domains as Objects

object domain \ object domain	$F_1$	$F_2$	$F_3$	laser printer	$D_1$	$D_2$	$D_3$	$D_4$
$D_1$	read		read			switch		
$D_2$				print			switch	switch
$D_3$		read	execute					
$D_4$	read write		read write			switch		

# Access Matrix with *Copy Rights*

object domain	$F_1$	$F_2$	$F_3$
$D_1$	execute		write*
$D_2$	execute	read*	execute
$D_3$	execute		

(a)

object domain	$F_1$	$F_2$	$F_3$
$D_1$	execute		write*
$D_2$	execute	read*	execute
$D_3$	execute	read	

LectureNotes.in

# Access Matrix With *Owner Rights*

object domain	$F_1$	$F_2$	$F_3$
$D_1$	owner execute		write
$D_2$		read* owner	read* owner write
$D_3$	execute		

(a)

object domain	$F_1$	$F_2$	$F_3$
$D_1$	owner execute		write
$D_2$		owner read* write*	read* owner write
$D_3$		write	write

(b)

# Modified Access Matrix of Figure B

object domain \ object domain	$F_1$	$F_2$	$F_3$	laser printer	$D_1$	$D_2$	$D_3$	$D_4$
$D_1$	read		read			switch		
$D_2$				print			switch	switch control
$D_3$		read	execute					
$D_4$	write		write		switch			

- The Security Problem
- Program Threats
- System and Network Threats
- Cryptography as a Security Tool
- User Authentication
- Implementing Security Defenses
- Firewalling to Protect Systems and Networks
- Computer-Security Classifications
- An Example: Windows 7



# The Security Problem

- **Threats:** It is the potential security violation
- When a program created by an user is used by another user there may be the chances of misuse.
- Examples: **Trojan horse:** It sits idle and transmits all user information to the attacker.
- Trap doors, stack/buffer overflow
- **System threats:** The OS files or resources are misused.
- Examples: WORMS: Small programs that replicates themselves. It contains malicious codes that cause major damage to the files of OS. They consume maximum resources thus denying the service to the user.
- **VIRUS:** Small programs that replicates and modification to the files.
- **Intruders (crackers)** attempt to breach security
- **Threat** is potential security violation
- **Attack** is attempt to breach security
- Attack can be accidental or malicious
- Easier to protect against accidental than malicious misuse

# Security Violation Categories

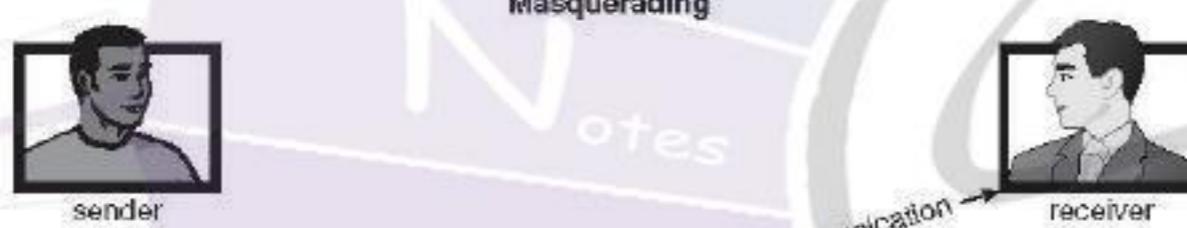
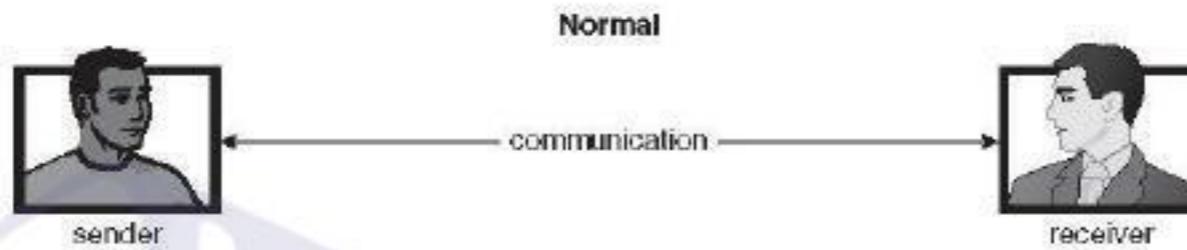
- **Breach of confidentiality**
  - Unauthorized reading of data
- **Breach of integrity**
  - Unauthorized modification of data
- **Breach of availability**
  - Unauthorized destruction of data
- **Theft of service**
  - Unauthorized use of resources
- **Denial of service (DoS)**
  - Prevention of legitimate use



# Security Violation Methods

- **Masquerading (breach authentication)**
  - Pretending to be an authorized user to escalate privileges
- **Replay attack**
  - As is or with **message modification**
- **Man-in-the-middle attack**
  - Intruder sits in data flow, masquerading as sender to receiver and vice versa
- **Session hijacking**
  - Intercept an already-established session to bypass authentication

# Standard Security Attacks



# Security Measure Levels

- Impossible to have absolute security, but make cost to perpetrator sufficiently high to deter most intruders
- Security must occur at four levels to be effective:
  - **Physical**
    - Data centers, servers, connected terminals
  - **Human**
    - Avoid **social engineering, phishing, dumpster diving**
  - **Operating System**
    - Protection mechanisms, debugging
  - **Network**
    - Intercepted communications, interruption, DOS
- Security is as weak as the weakest link in the chain
- But can too much security be a problem?

# Program Threats

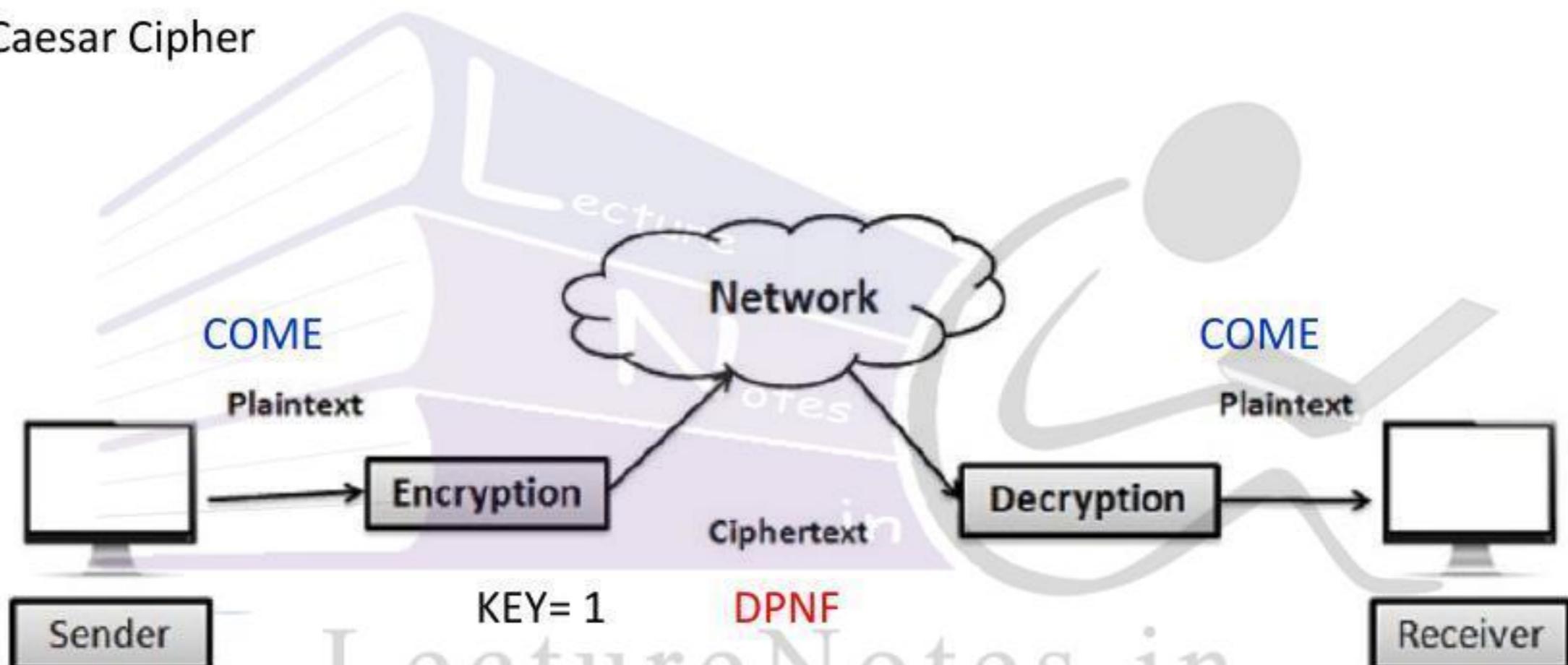
- Many variations, many names
- **Trojan Horse**
  - Code segment that misuses its environment
  - Exploits mechanisms for allowing programs written by users to be executed by other users
  - **Spyware, pop-up browser windows, covert channels**
  - Up to 80% of spam delivered by spyware-infected systems
- **Trap Door**
  - Specific user identifier or password that circumvents normal security procedures
  - Could be included in a compiler
  - How to detect them?

# Program Threats (Cont.)

- **Logic Bomb**
  - Program that initiates a security incident under certain circumstances
- **Stack and Buffer Overflow**
  - Exploits a bug in a program (overflow either the stack or memory buffers)
  - Failure to check bounds on inputs, arguments
  - Write past arguments on the stack into the return address on stack
  - When routine returns from call, returns to hacked address
    - Pointed to code loaded onto stack that executes malicious code
  - Unauthorized user or privilege escalation

# Encryption

- It is a technique that enables the sender of a message to encrypt the message, so that only a legitimate receiver can decrypt and read the message.
- Caesar Cipher



## ACKNOWLEDGEMENT

I am very much thankful to say that  
some contents of these slides are taken  
from the book

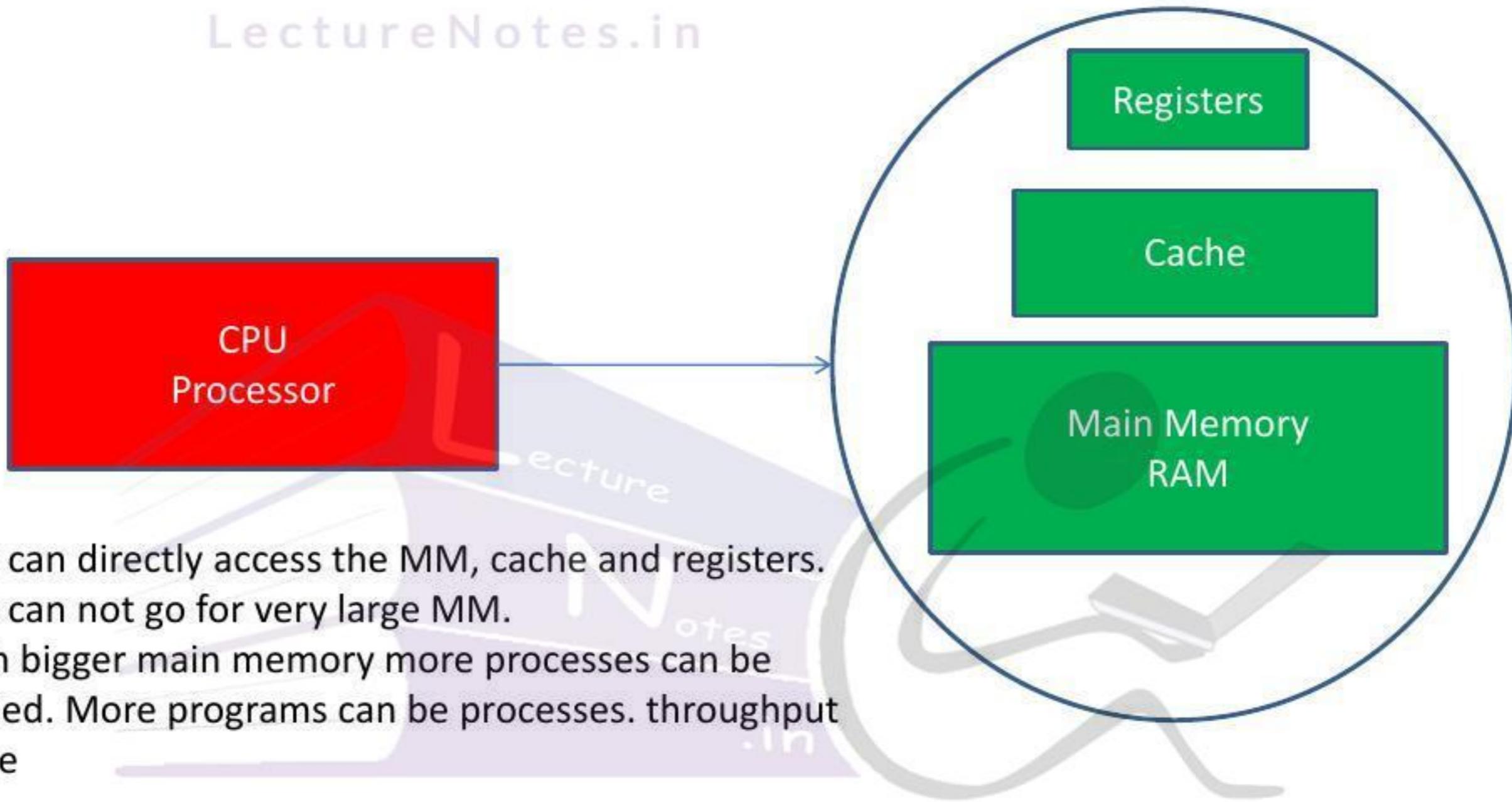
Operating System Concepts  
Silberschatz, Galvin and Gagne.

LectureNotes.in

# **MEMORY MANAGEMENT**

**Dr. Ajay Kumar Jena**

**School of Computer Engineering  
KIIT University, Bhubaneswar**



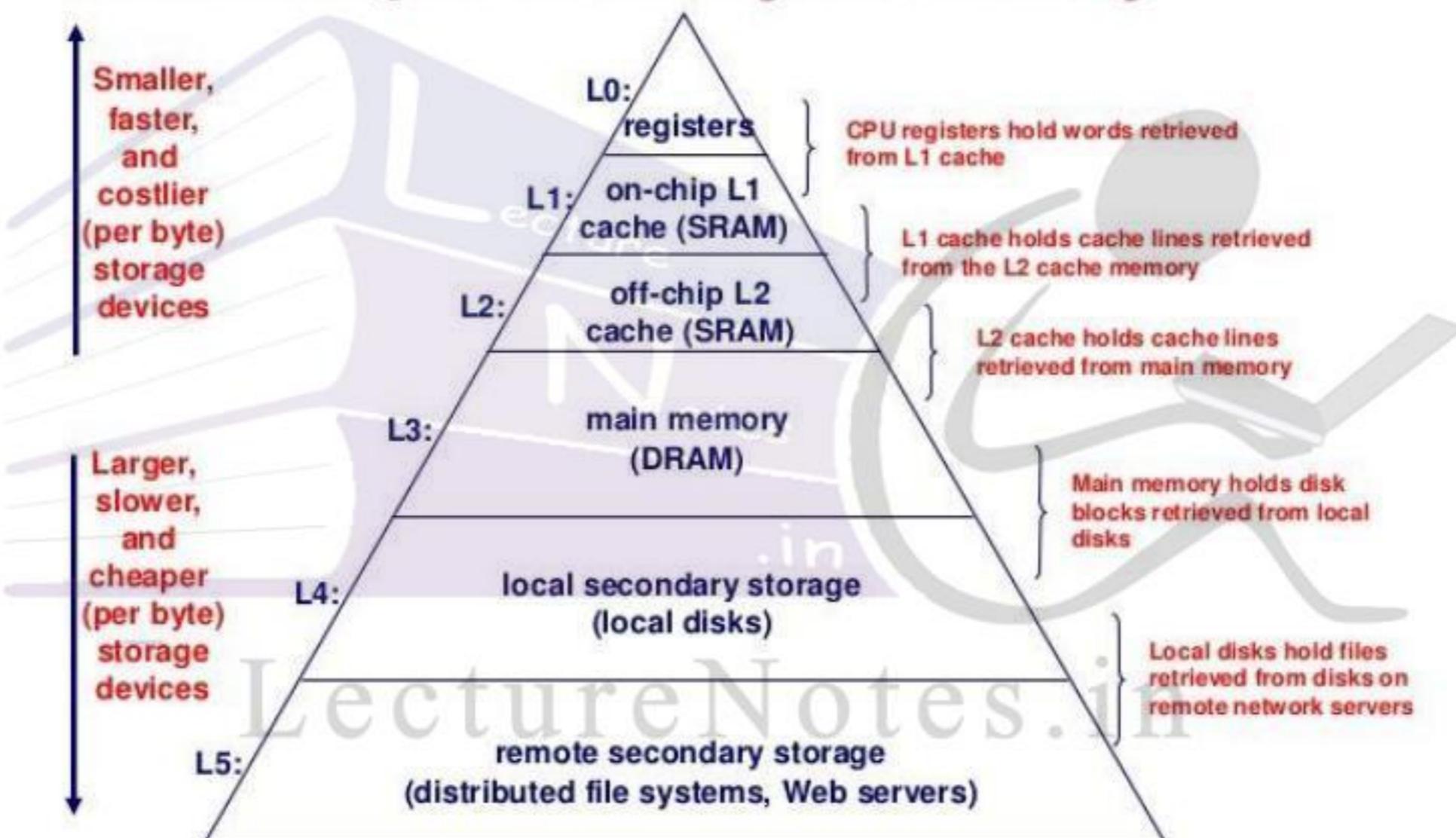
CPU can directly access the MM, cache and registers.  
One can not go for very large MM.

With bigger main memory more processes can be loaded. More programs can be processes. throughput more

### Criteria

1. Size (More)
2. Access Time (less)
3. Per unit cost (less)

## An Example Memory Hierarchy

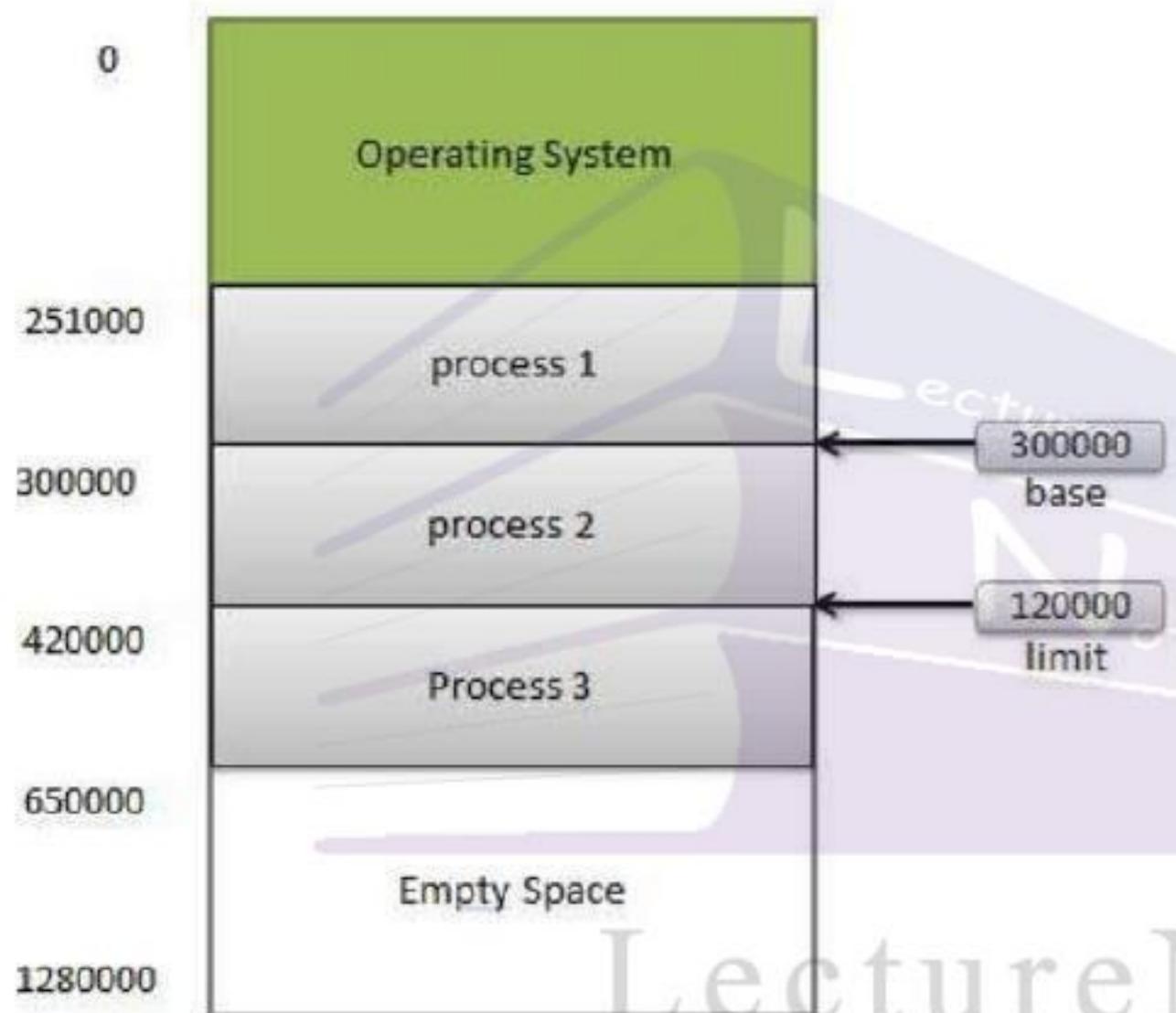


- In a uni-programming system, main memory is divided into two parts. One part for the OS(kernel) and another part for the program currently being executed.
- In a multiprogramming system, the “user” part of memory must be subdivided further to accommodate multiple processes.
- The task of subdivision is carried out dynamically by the OS is known as memory management.
- To improve the utilization of the CPU and speed of computer, we must keep several processes in memory i.e we must share memory

# Overview of Memory

- **Memory consists of large array of words or bytes, each with its own address.**
- **The CPU fetches instructions from memory according to the value of program counter**
- **Instruction execution cycle consists of, first fetch an instruction from memory. The instruction is then decoded and operands fetched from memory. The instruction has been executed on the operands, results may be stored back in memory.**

# Memory Protection

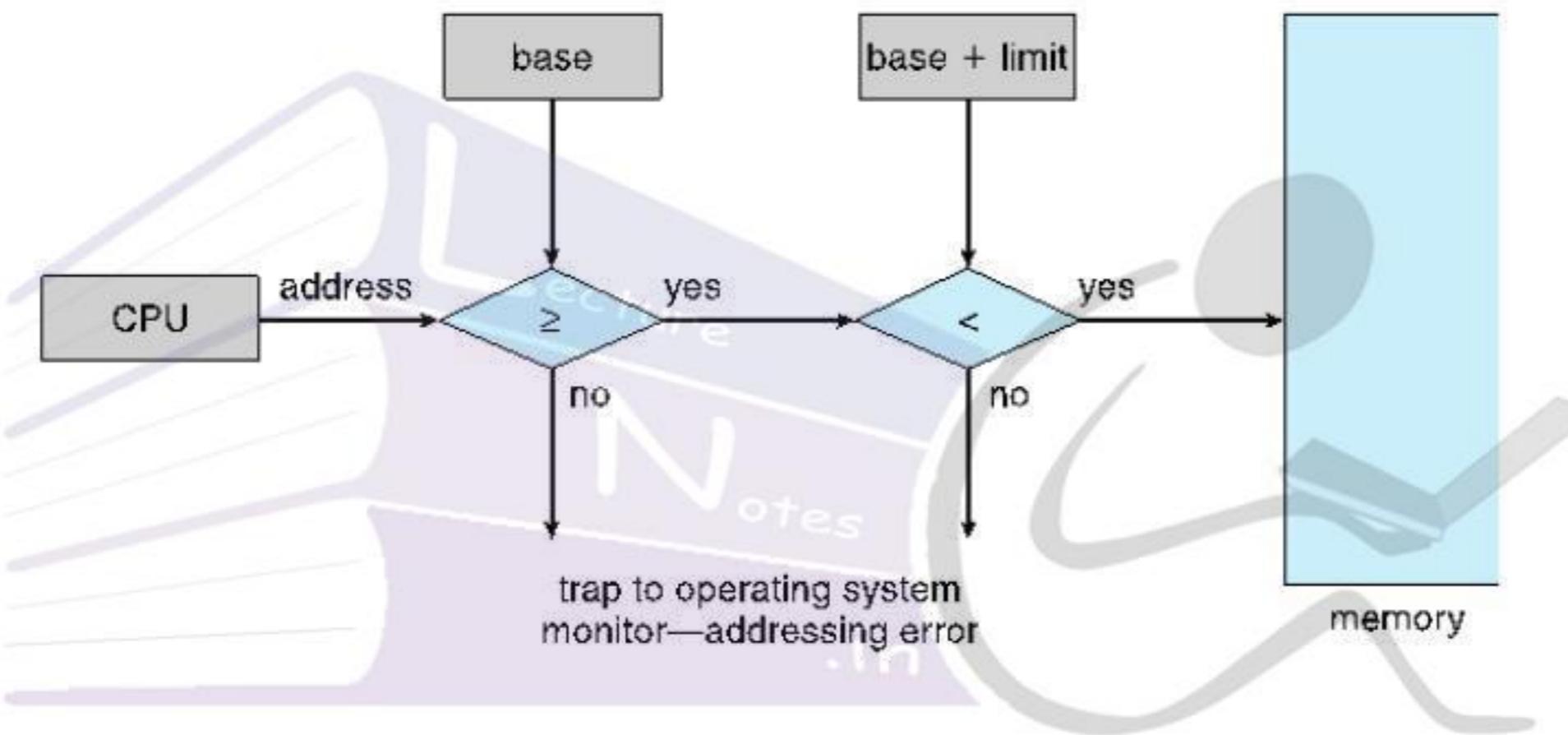


Base and limit register defines a logical address space

- Each process has to store in a separate memory space.
- For concurrent execution multiple processes must be loaded to memory.
- To separate memory spaces and to ensure that a process can access its legal address space base and limit registers are used.
- Here the program can legally access all address from 300000 through 420000.

# Memory Protection

- Protecting the O/S from user processes and protecting the user process from one another the base and limit registers are used.
- The base register contains the value of the smallest physical address and the limit register contains the range of the logical address.
- In this ex. base register holds 300000 and the limit register holds 120000, then the program can legally access all addresses from 300000 through 420000.
- Each logical address must be less than the limit register. The memory management unit (MMU) maps the logical address dynamically by adding the value in the base register. This mapped address is sent to the memory.

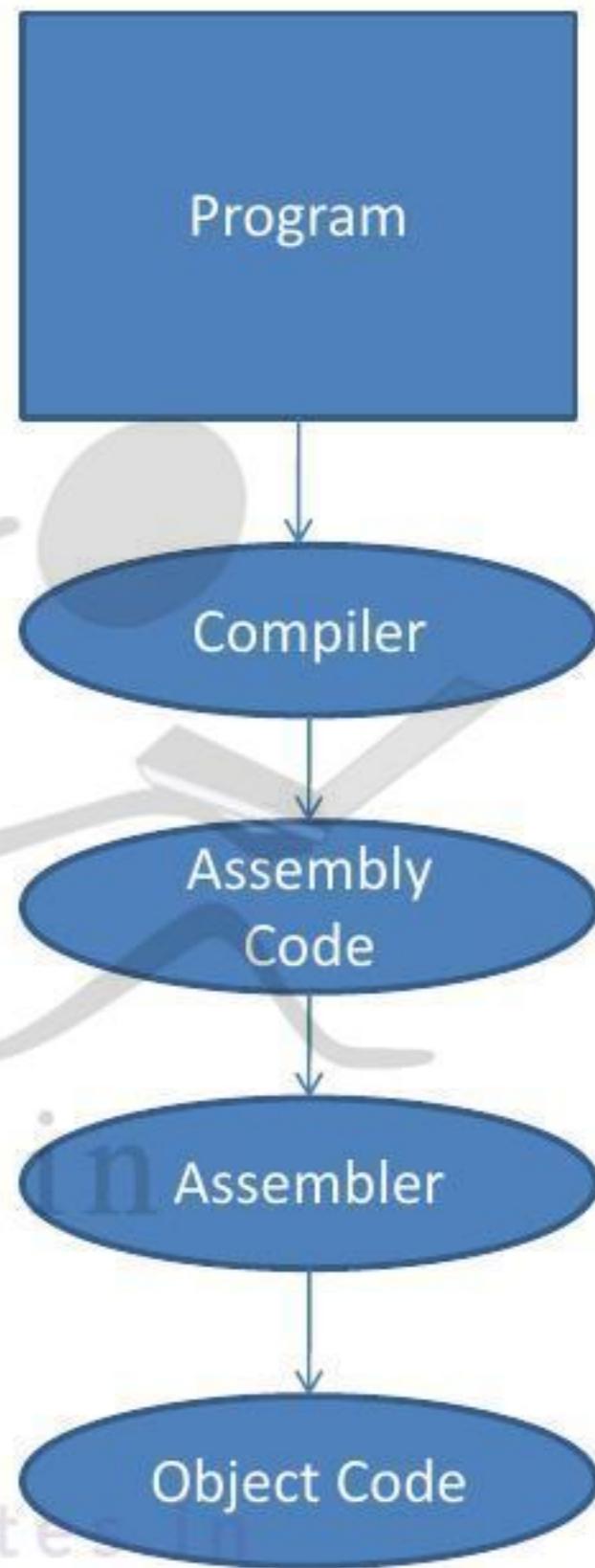
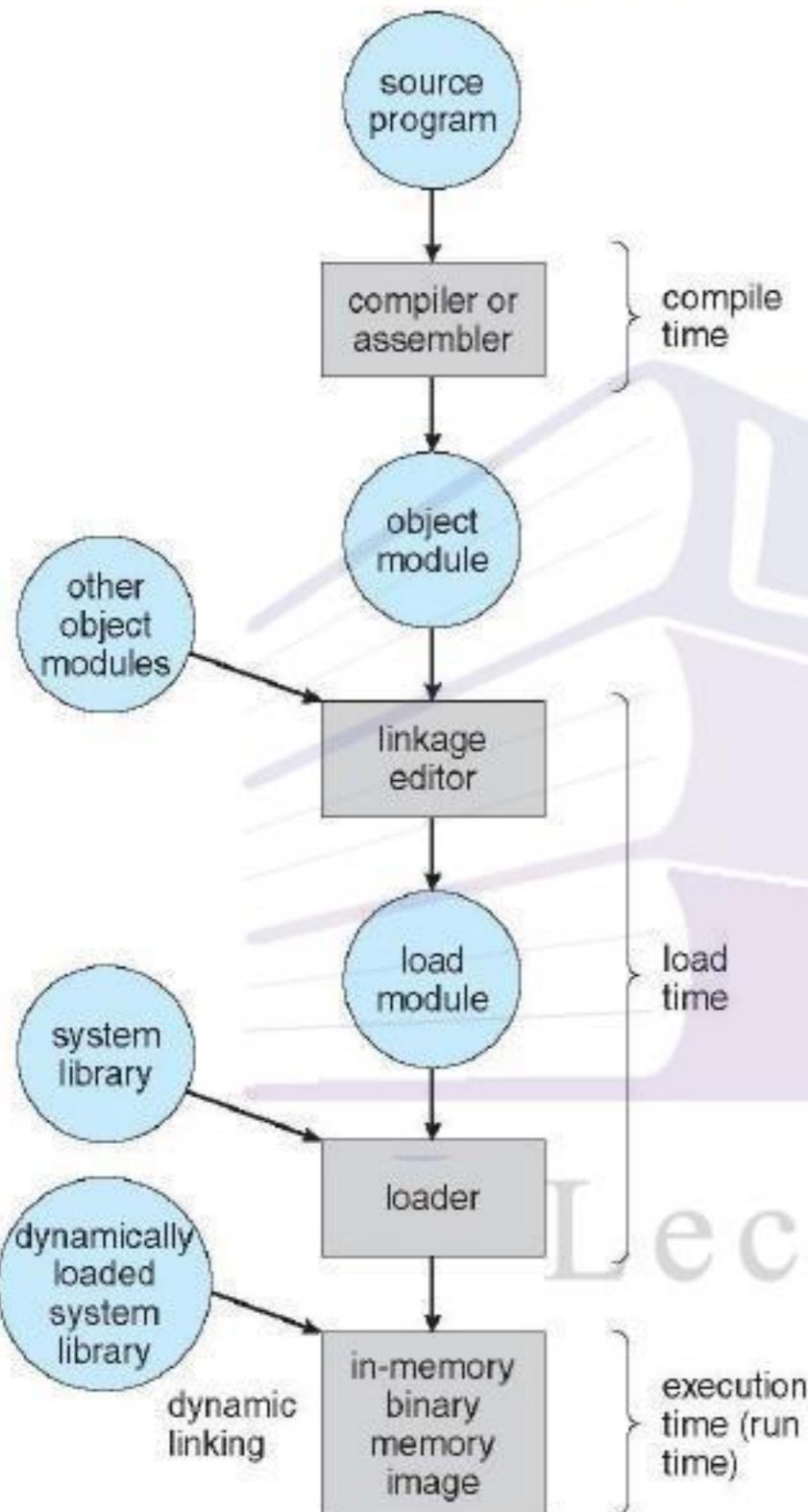


## Hardware address protection with base and limit register

# Address Binding

- A program resides on a disk as a binary executable file. The program must be brought into memory and placed within a process for it to be executed.
- The collection of processes on the disk that is waiting to be brought into memory for execution forms the input queue.
- Select one of the processes from the input queue and to load that process into memory. When the process is executed, it accesses instructions and data from memory. After execution, the process terminates and its memory space is free.
- Address in the source program are symbolic. The compiler will bind these symbolic addresses to relocatable addresses.

- The linkage editor or loader will bind these relocatable addresses to absolute addresses.
- Binding of instructions and data to memory addresses can be done:
  - Compile time: If memory location known a priori (known at the compile time where the process can reside), absolute code can be generated; must recompile code if starting location changes
  - Load time: If address is not known at compiler time, it must generate relocatable code
  - Execution time: Binding delayed until run time if the process can be moved during its execution from one memory segment to another



## Multistep processing of a user program

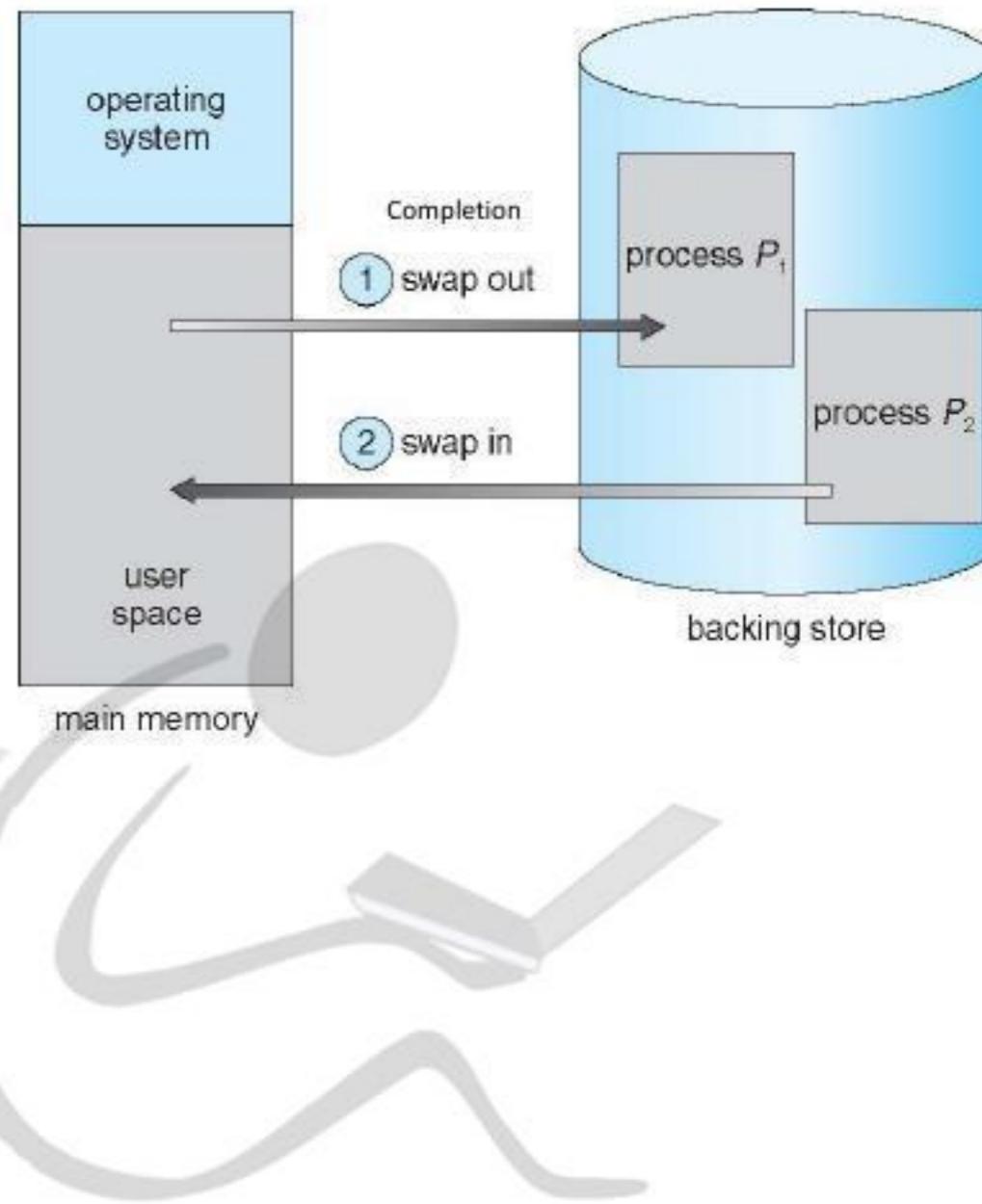
# Logical and Physical Address Space

- An address generated by the CPU is called as logical address.
- An address seen by the memory unit i.e. the one loaded into the memory address register of the memory is referred as a physical address.

LectureNotes.in

LectureNotes.in

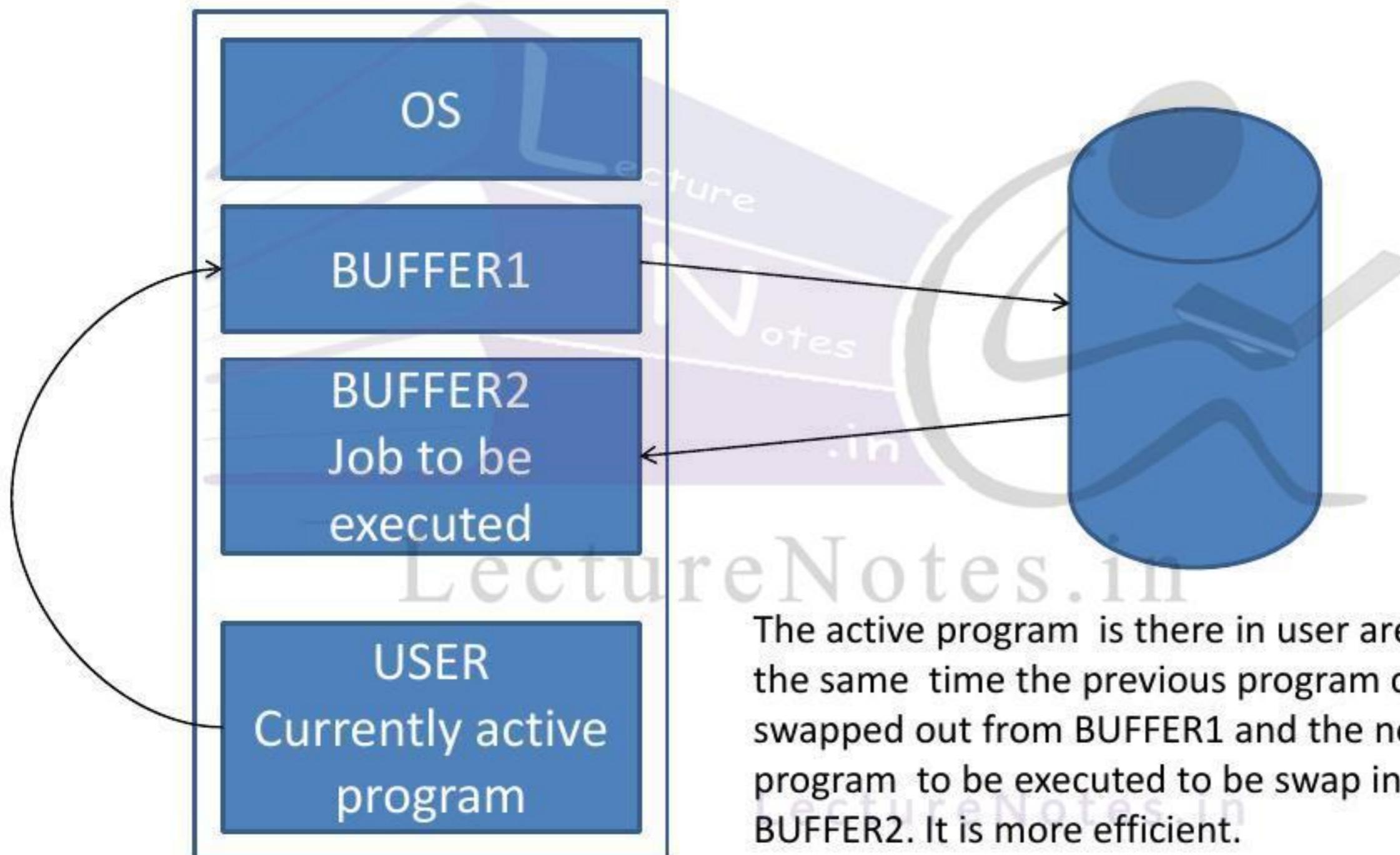
- Swapping is a method to improve the memory utilization.
- A process needs to be in memory for execution. A process can be swapped temporarily out of memory to a backing store and then brought back into memory for continued execution.
- Suspended due to I/O, terminated itself after completed all of its operation and abnormal termination at the time of trap.
- Switching a process from the main memory to disk is said to be “swap out” and switching from disk to main memory is said to be “swap in”.



In RR scheduling, when a quantum expires, the memory manager will start to swap out the process that just finished and to swap in another process.

- Suppose the main memory composed of 10 processes, assume this is the maximum capacity the CPU and currently process 9 is executing. In the middle suppose P9 needs some I/O, then CPU switch to another process and P9 is moved to the disk and another process is loaded in main memory in place of P9. After completion of the I/O of P9 it will be moved to main memory from the disk.

- In every swap out and swap in performs some I/O operation. Which takes a huge amount of time. The CPU utilization is getting very poor if more number swap will be there. CPU cannot start until the swap in and swap out operation is completed.
- To improve the CPU efficiency



The active program is there in user area and at the same time the previous program can be swapped out from BUFFER1 and the next program to be executed to be swap in into BUFFER2. It is more efficient.

# Dynamic Loading

- **Loading means load program or module from secondary storage device (disks) to main memory.**
- **Mainly divided into two types**
  - **Compile time loading / static loading**
  - **Runtime loading / dynamic loading**
- **All the routines are loaded in the main memory at the time of compilation is said to be static loading or compile time loading.**
- **The routines are loaded in the main memory at the time of execution or running is said to be dynamic loading.**
- **Only the main program will be loaded and the functions will be there in the sec. storage.**

# Dynamic Linking

- **Static linking** – system libraries and program code combined by the loader into the binary program image
- Dynamic linking –linking postponed until execution time
- Small piece of code, **stub**, used to locate the appropriate memory-resident library routine
- Stub replaces itself with the address of the routine, and executes the routine
- Operating system checks if routine is in processes' memory address
  - If not in address space, add to address space
- Dynamic linking is particularly useful for libraries
- System also known as **shared libraries**
- Consider applicability to patching system libraries
  - Versioning may be needed

# Contiguous Memory Allocation

- The main memory must accommodate both the OS and the user processes.
- The memory is divided into two partitions i.e. Resident of OS and User processes
- The OS may place either in low or high memory. Since the interrupt vector is often in low memory, programmers usually place the OS in low memory.
- Generally, we want several user processes to reside in main memory at the same time. So, in contiguous memory allocation, each process is contained in a single contiguous part of memory.

## Memory Management Technique

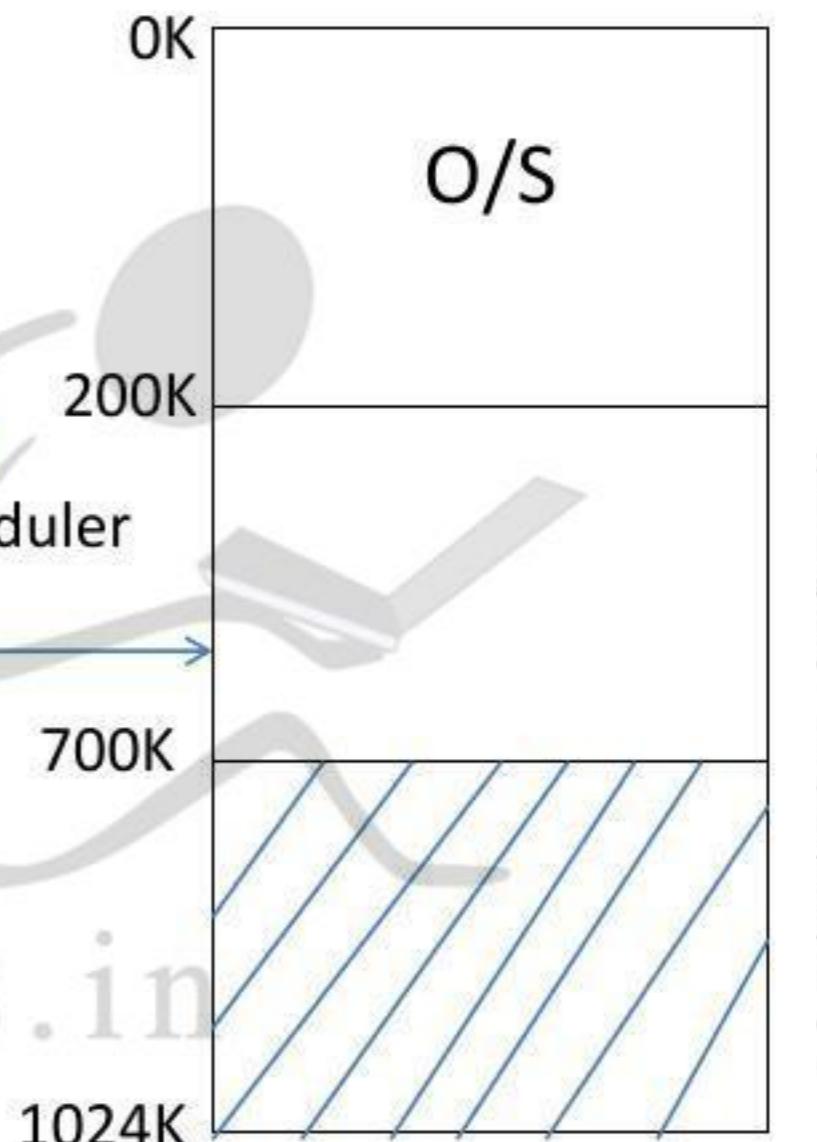
- Contiguous (ex. array)
  - Fixed Partition (Static)
  - Variable Partition (Dynamic)
- Non-contiguous (linked list)
  - Paging
  - Multilevel Paging
  - Inverted Paging
  - Segmentation
  - Segmented Paging



LectureNotes.in

# Single Partition Allocation

- The OS resides in the low memory and the remaining memory is treated as a single partition. This is available for user space. Only one job can be loaded in this user space.
- The main memory consists of only one process at a time because the user space is treated as a single partition.



- ADVT: It is simple and doesn't require great expertise.
- DISADVT: Poor utilization of memory
- Poor utilization of process waiting for I/O
- User's job is limited to the size of available main memory.

# Multiple Partitions memory Management

- This method can be implemented in three ways.
- Fixed equal multiple partitions memory mgt.
- Fixed variable multiple partitions memory mgt.
- Dynamic multiple partitions memory mgt.

# Fixed equal multiple partitions/ static partition

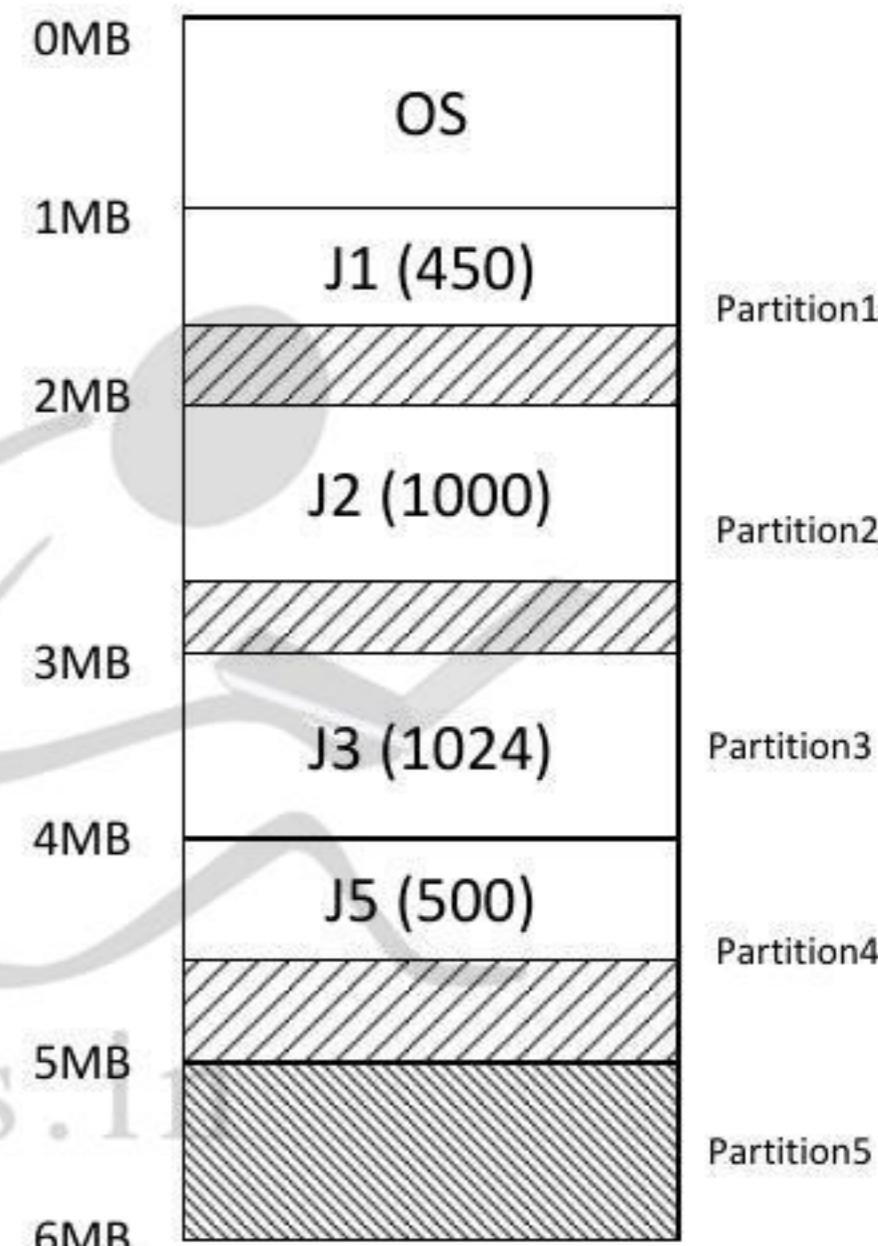
- The user space is divided into fixed partitions.  
The partition sizes are depending upon OS.

Total MM is of size 6MB. 1MB occupied by OS.

Remaining is partitioned into 5 equal fixed partitions.

J1,J2,J3,J4,J5 are 5 jobs to be loaded in to MM.

Job	Sizes
J1	450KB
J2	1000KB
J3	1024KB
J4	1500KB
J5	500KB



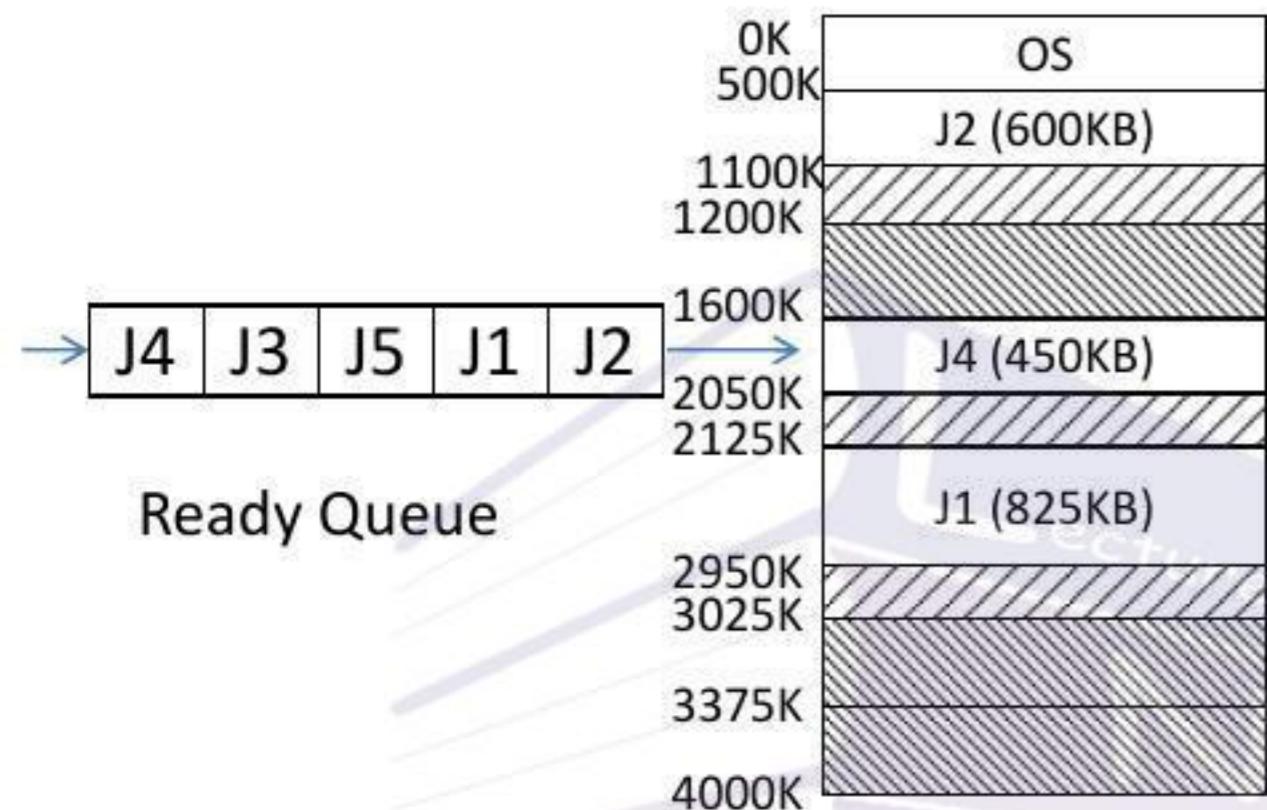
- Job1 is loaded into partition1. The size of par1 is 1024KB. Size of J1 is 450KB.
- Waste space=  $1024 - 450 = 574$ KB
- This known as **internal fragmentation**.

There is no enough space to load J4 bcoz the size of J4 is GT all partitions. The memory which is not allotted to any process is external fragmentation i.e. 1024KB.

Total internal fragmentation is =  $(1024 - 450) + (1024 - 1000) + (1024 - 500) = 1172$ KB

A partition of MM wasted within a partition is internal frag. And the wastage of an entire partition is external fragmentation.

# Fixed variable multiple partitions



Partitions	
Partition	Sizes
P1	700KB
P2	400KB
P3	5254KB
P4	900KB
P5	350KB
P6	625KB

Job Queue		
Job	Sizes	Arival Time
J1	825KB	10ms
J2	600KB	5ms
J3	1200KB	20ms
J4	450KB	30ms
J5	650KB	15ms

- In this the user space is divided into a number of partitions but the partition sizes are different.
- The OS keeps track which partitions are occupied and which are free.
- J2 arrives first and after searching from low to high loaded to Partition1.
- J1 arrives next. P4 is large enough to accommodate. So J1 into P4 leaves 75KB free.
- J5 arrives no enough space to load that job. J5 has to wait till gets free memory.
- J3 arrives next, there is no large enough partition to load that.
- J4 arrives last the size is 450KB, partition P3 is large enough to load that one, loaded to P3.
- Total external Fragmentation=  $400+350+625=1375$
- Total internal fragmentation=  $100+75+75=250$  (J2 can be loaded P6 for efficient use)

To avoid this problem, we use (partition allocation policy)

1. First Fit algorithm: It allocates the process in the first free large enough partition.
2. Best Fit: It allocates the process in the smallest possible memory partition
3. Worst Fit: It searches for the largest possible memory partition

**Example:** There are five different size memory partitions of 100K, 500K, 200K, 300K, 600K (in sequence). By using first, best and worst fit algo place the processes of 250, 455, 189, 465K (in order).

First Fit
100KB
P1
500KB
P3
200KB
300KB
P2
600KB

Best Fit
100KB
P2
500KB
P3
200KB
P1
300KB
P4
600KB

Worst Fit
100KB
P2
500KB
200KB
300KB
P3
600KB
P1

A table called Partition Description Table (PDT) keeps information about the partitions.

Let us consider we have 1 MB of primary memory (RAM). Four partitions are marked as 200K, 250K, 150K and 500K respectively. To keep track of all partitions, a table is created as follows

P. No	P. Base	P. Size	P. Status
1	0K	200	FREE
2	200K	250	FREE
3	450K	150	FREE
4	600K	500	FREE

# Dynamic Partitions Memory Management

- In this method the partitions are created dynamically, so that each process will be loaded in the partition exactly the same size of the process.
- In this scheme the entire use space is treated as a “big hole”. The boundaries of the partitions are dynamically changed and the boundaries are depend on the size of the process.



Ready Queue

Job Queue		
Job	Sizes	Arival Time
J1	825KB	10ms
J2	600KB	5ms
J3	1200KB	20ms
J4	450KB	30ms
J5	650KB	15ms

- The jobs J2, J1, J5, J3 are loaded. The last job J4, the size is 450K, but the available memory is 225K ( $4000 - (500 + 600 + 825 + 650 + 1200)$ ). So the job J4 has to wait until the memory is available.

**ADVT:** This scheme support multiprogramming.

Efficient utilization of the processor and I/O devices.

Partitions are changes dynamically, so it does not suffer from internal fragmentation.

**DISADVT:**

This scheme suffers from internal as well as external fragmentation.

# Lecture Notes in Compaction

The solution to external fragmentation is compaction.

The goal of compaction is to place all free memory together in one large block.

- Compaction is not possible if relocation is static.
- Compaction is possible if relocation is dynamic done at execution time.
- The simplest compaction algorithm is to move all processes towards one end of memory so that all holes will be moved to the other end producing one large hole of available memory.

**ADVANTAGES:** This improves memory utilization and performance of system.

**DISADVANTAGES:**

This scheme is very expensive.

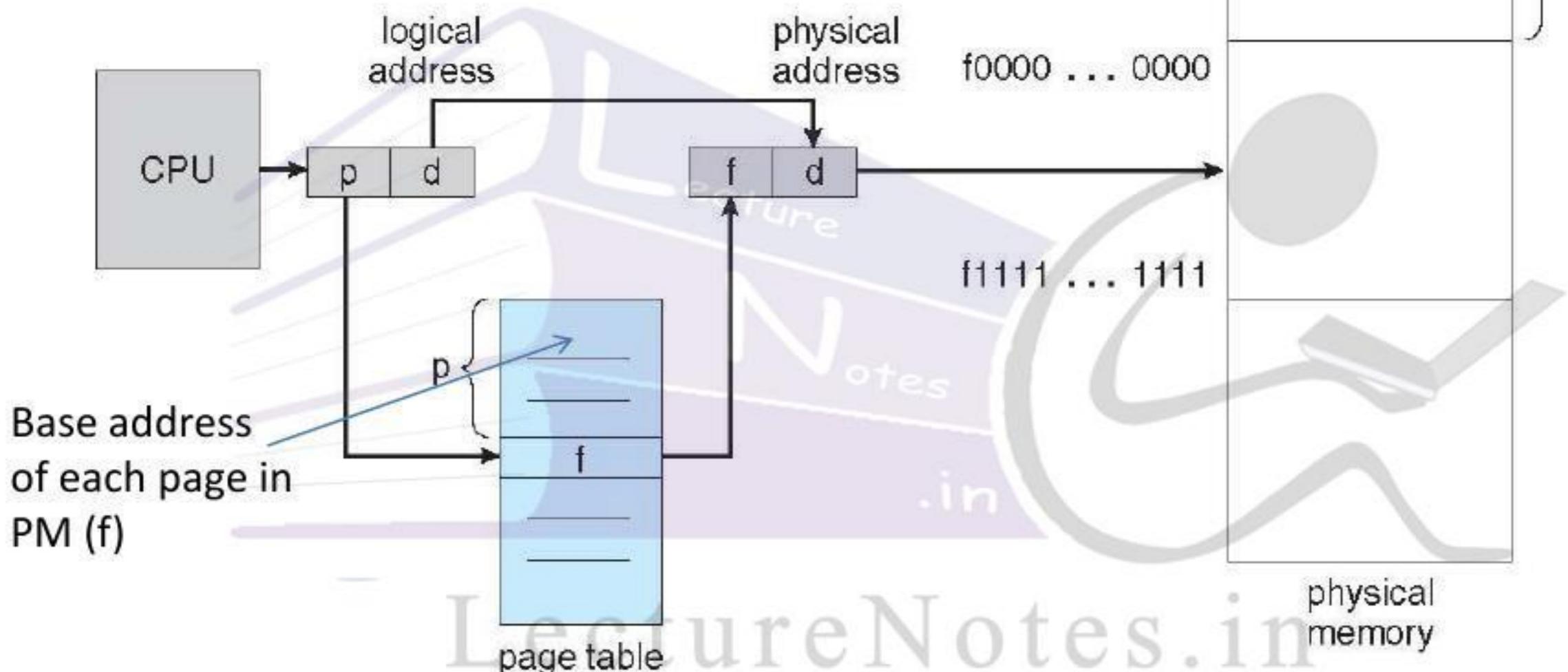
It is not always possible.

# Non-contiguous Memory Allocation

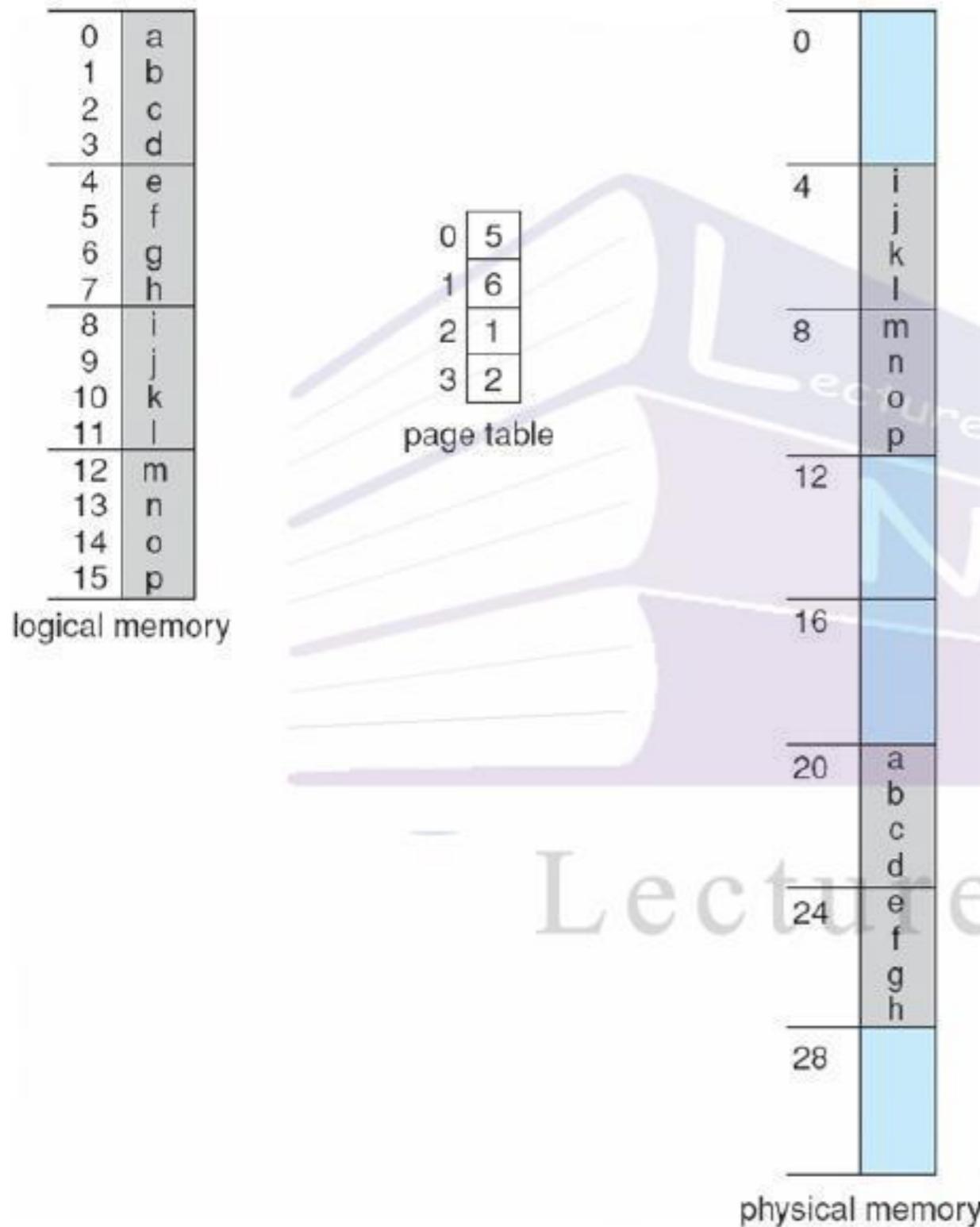
## Paging

- A contiguous memory may not be available always.
- For execution of a process no consecutive memory is required.
- A process is divided to a number of parts, which may be loaded into different area of RAM / primary memory.
- A process is divided to number of equal size blocks called pages.
- Primary memory or physical memory into same sized blocks called **frames**.
- **Size of the page = Size of the frame**
- A page of a process is loaded into any of the available frame.
- A Page Table keeps track of which page of a process is loaded into which frame of primary memory.
- For every process there is a page table.

# Paging Hardware



# Paging Example



Mapping of logical address to physical address.  
A page size =4 bytes and physical memory = 32 bytes( 8 pages)

Logical address of 0 is page 0, offset 0

Logical address of 1 is page 0, offset 1

Logical address of 2 is page 0, offset 2

Logical address of 3 is page 0, offset 3

Logical address of 4 is page 0, offset 0

Logical address of 5 is page 0, offset 1

Logical address of 0 is page=0, offset=0

Physical Address= Frame no X page size +offset

Physical addr of logical addr 0 is 20= 5X4+0

Physical addr of logical addr 3 is 23= 5X4+3

Physical addr of logical addr 4 is 24= 6X4+0

Physical addr of logical addr 13 is 9= 2X4+1

# Advt. and dis-advantages of Paging

- No external fragmentation
- Internal fragmentation bcoz only the last page of the process is not filled completely. On an average the last page of the process is half full, so it is called  $\frac{1}{2}$  of the page.
- Sharing of pages can be done.
- It supports virtual memory

## Dis-advantages:

- Extra space is required for sharing the page table
- Takes more time for address translation
- It suffers from page break. Suppose logical address space is 17KB, page size is 4KB, So this job requires 5 frames, the 5th frame consists of only 1KB, so remaining 3KB wasted. Therefore, it is called page break.
- If the number of pages are high, it is difficult to maintain page table.

# Lecture Notes in **Segmentation**

- Segmentation is a method in which a process will be divided into parts and put it in the MM looking to the related data i.e function, arrays.....
- Paging also same but in fixed size. (Example of C program, Inst1....Inst99 )
- Segmentation divided the program in variable sizes. It is not diving the process without looking into it.
- The user view is not same as the actual physical memory.
- Memory-management scheme that supports user view of memory contiguous memory allocation.
- A logical address space is a collection of segments. It is defined as logical grouping of instructions such as :

main program	procedure
function	method
object	local variables, global variables
common block	stack
symbol table	arrays

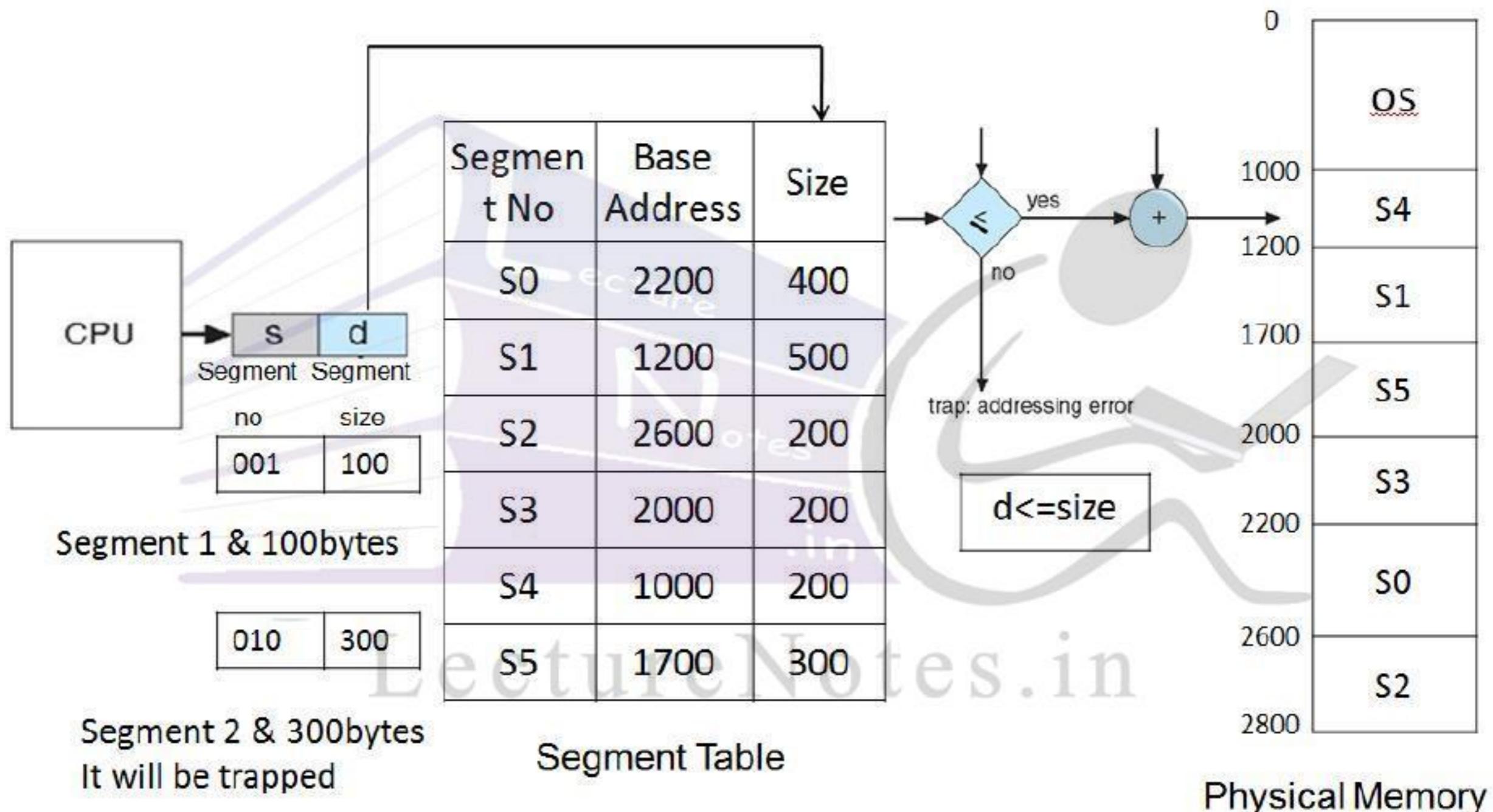
- Each segment has a name and length. The addresses specify both the segment name and the offset within the segment.

# Lecture Notes in Segmentation Architecture

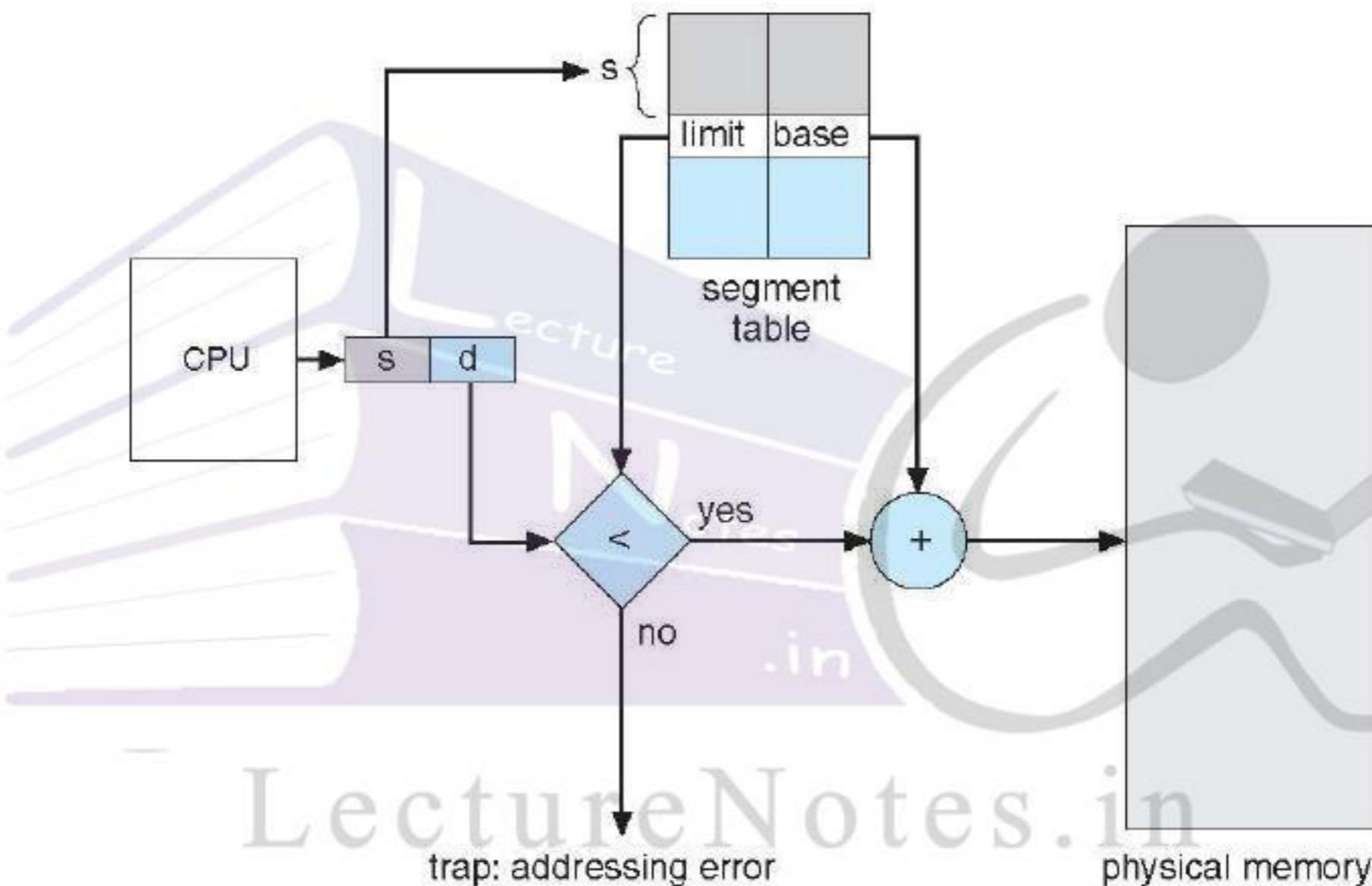
- CPU generates logical address
- Logical address consists of a two tuple:
  - $\langle \text{segment-number}, \text{offset} \rangle$ ,
- **Segment table** – maps two-dimensional physical addresses; each table entry has:
  - **base** – contains the starting physical address where the segments reside in memory
  - **limit** – specifies the length of the segment
- **Segment-table base register (STBR)** points to the segment table's location in memory
- **Segment-table length register (STLR)** indicates number of segments used by a program;
- segment number **s** is legal if  **$s < STLR$**

# LectureNotes.in

# Segmentation Hardware



# LectureNotes.in Segmentation Hardware



LectureNotes.in

# Comparison between Paging and Segmentation

## Paging

- Logical memory is partitioned into pages
- The logical address space is divided into pages by MMU
- Suffering from internal fragmentation
- O/S maintains a page table to map between frames and page.
- It doesn't support the user's view of memory.
- Page no. and displacement used to calculate absolute address(p,d).
- Multilevel paging is possible.
- O/S maintains a free frame list, so no need to search for free frame.

## Segmentation

- Logical memory is partitioned into segments.
- The logical address space is divided into segments specified by programmers.
- Suffering from external fragmentation
- O/S maintains a segment table for mapping.
- It supports the user's view of memory.
- Segment no. and offset is used to calculate absolute address(s,d).
- Multilevel segmentation is possible but no use
- O/S maintains the particulars of available memory.

# Paging using Translation Look-aside Buffer

- Page table is data structure implemented as a set of dedicated registers. The use of registers for small page table is efficient. But large page table these are not feasible.
- The page tables are kept in main memory. In the main memory where the program stores the page table also stores there.
- By accessing a logical address by the CPU is to access the page table in the main memory is same as accessing the page table first and the physical address next.
- So the main memory will be accessed twice or by factor 2.
- It is a penalty of time. To reduce this time.
- **The memory access is slow. (trade off speed and time)**

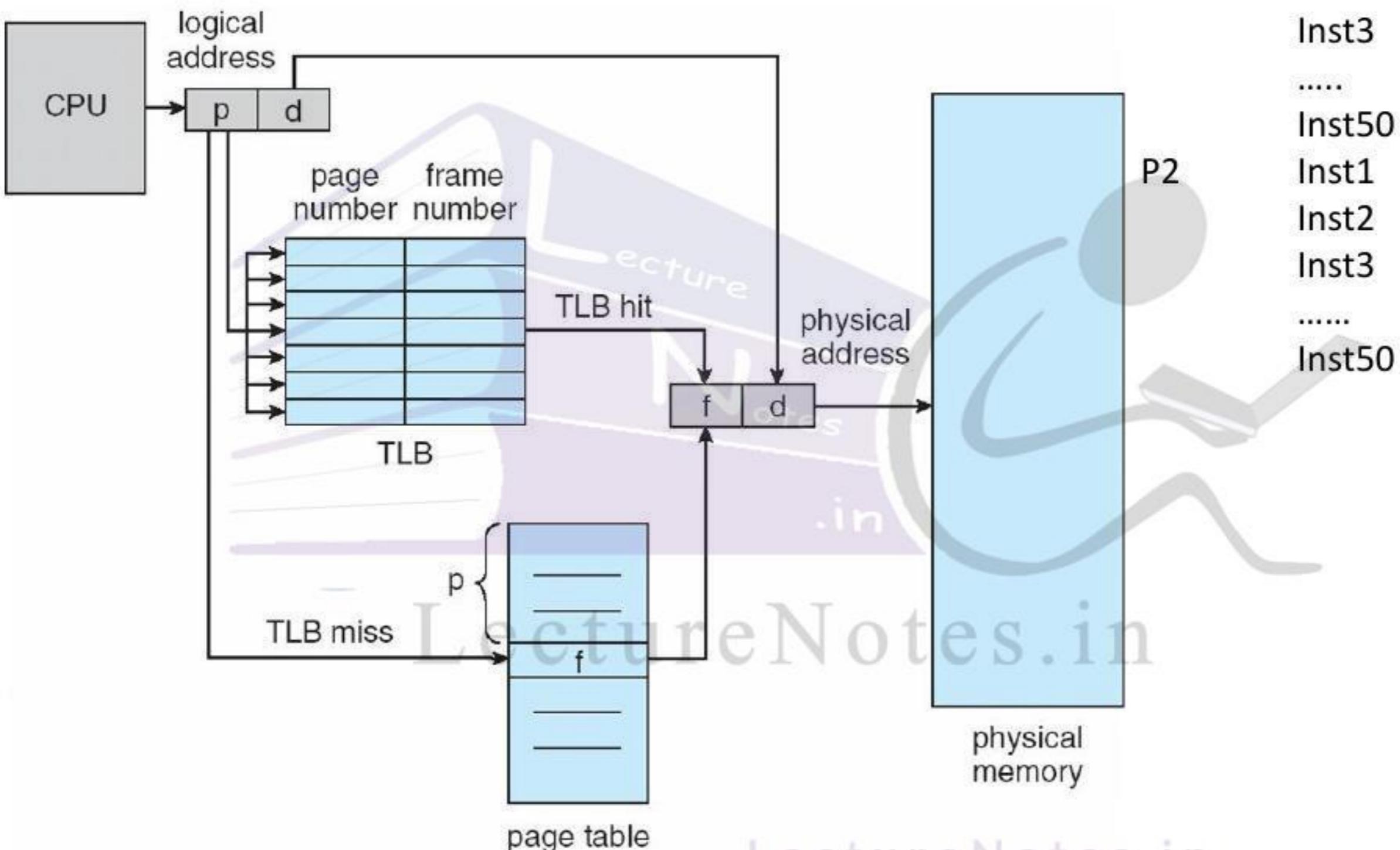
The solution is

LectureNotes.in

- A special small fast look up h/w cache called Translation Look-aside Buffer (TLB).
- While context switch one has clear the TLB so that next process can start fresh.

LectureNotes.in

# Paging Hardware With TLB



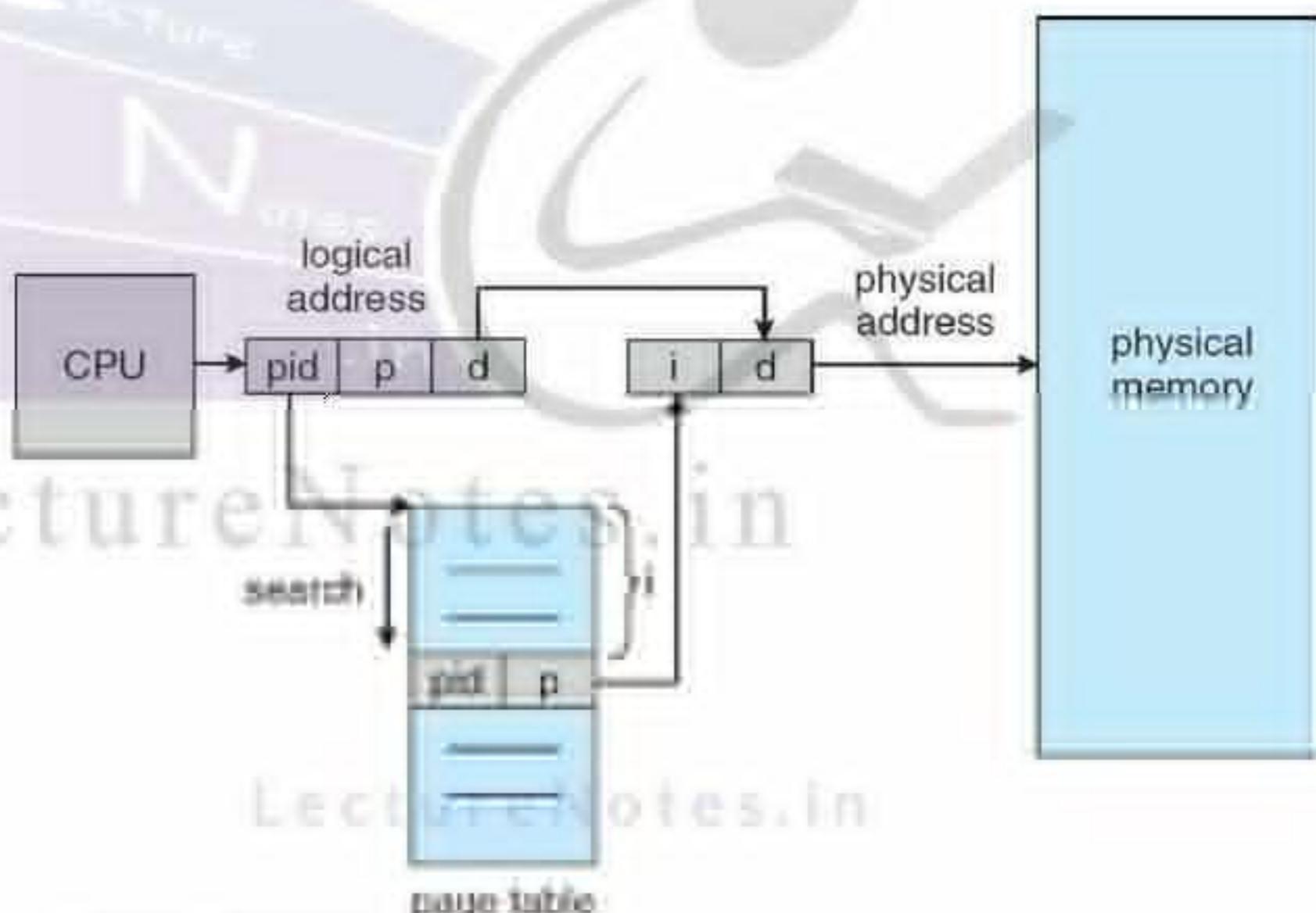
# Inverted Page Table

- Each process is having a page table which keep track of all possible logical pages, track all physical pages
- Each page table may consists of millions of entries. These tables may consume large amount of physical memory because for each process it creates a page table. As the MM is limited so frames are limited. So we have to think something
- To reduce the memory one can use inverted page table.
- Instead of keeping all the page table for all the processes, one can create a global page table which will handle all the page table.
- Here the searching time is more.

# Inverted Page Table Architecture

Global Page Table

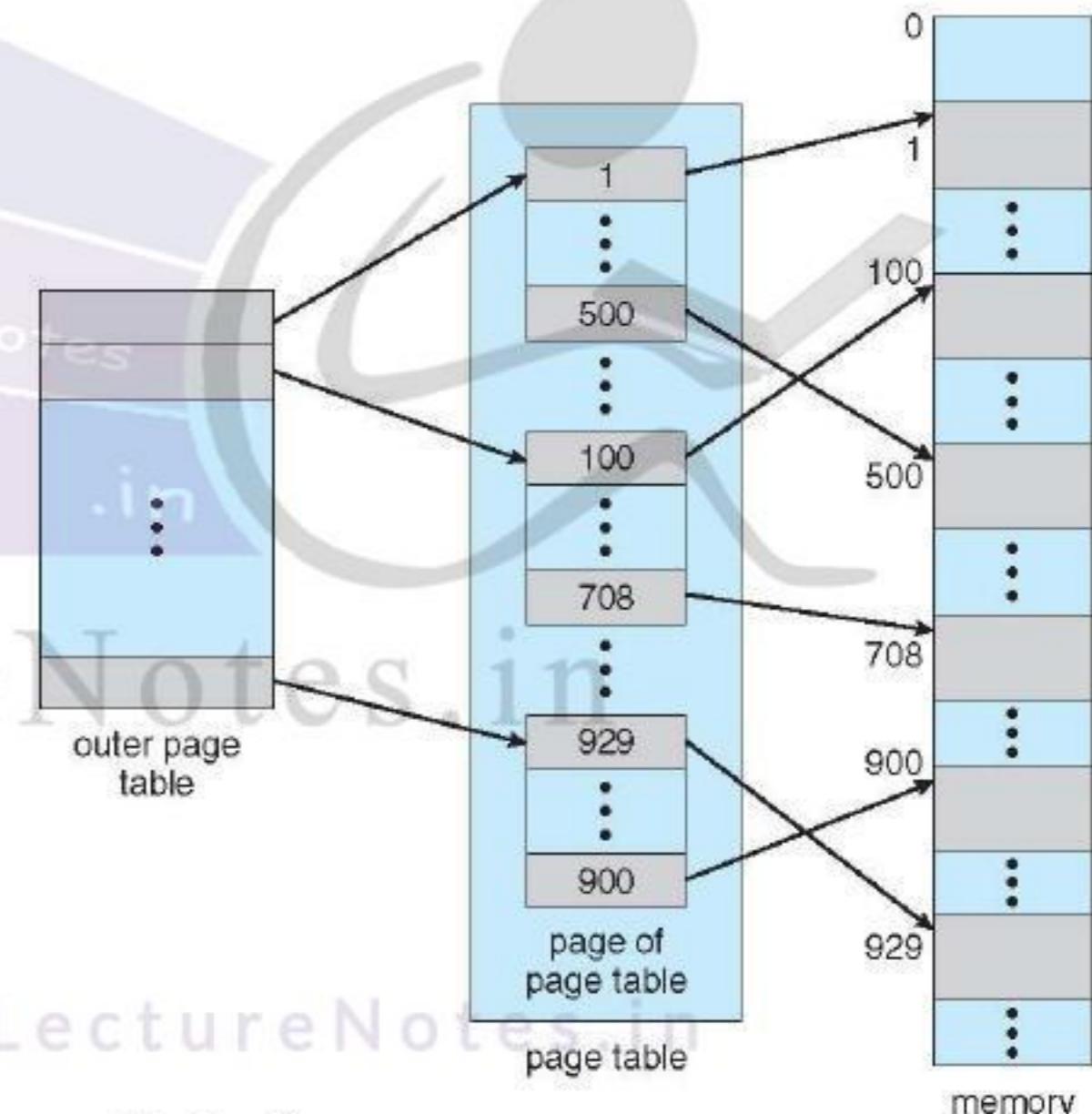
Frame No	Page No	Process ID
0	Page0	P1
1	Page3	P2
2	Page2	P3
3	Page2	P1
4	Page1	P3
5	Page0	P2



# Lecture Notes

## HIERARCHICAL PAGING

- Most of the modern computer system supports a large logical-address space( $2^{32}$  to  $2^{64}$ )
- So the page table becomes exclusively large.
- Example: Consider a system of 32-bit logical address space. If the page size of the system is 4KB( $2^{12}$ ), then a page table will consists of 1 million entries ( $2^{32}/2^{12}$ ).



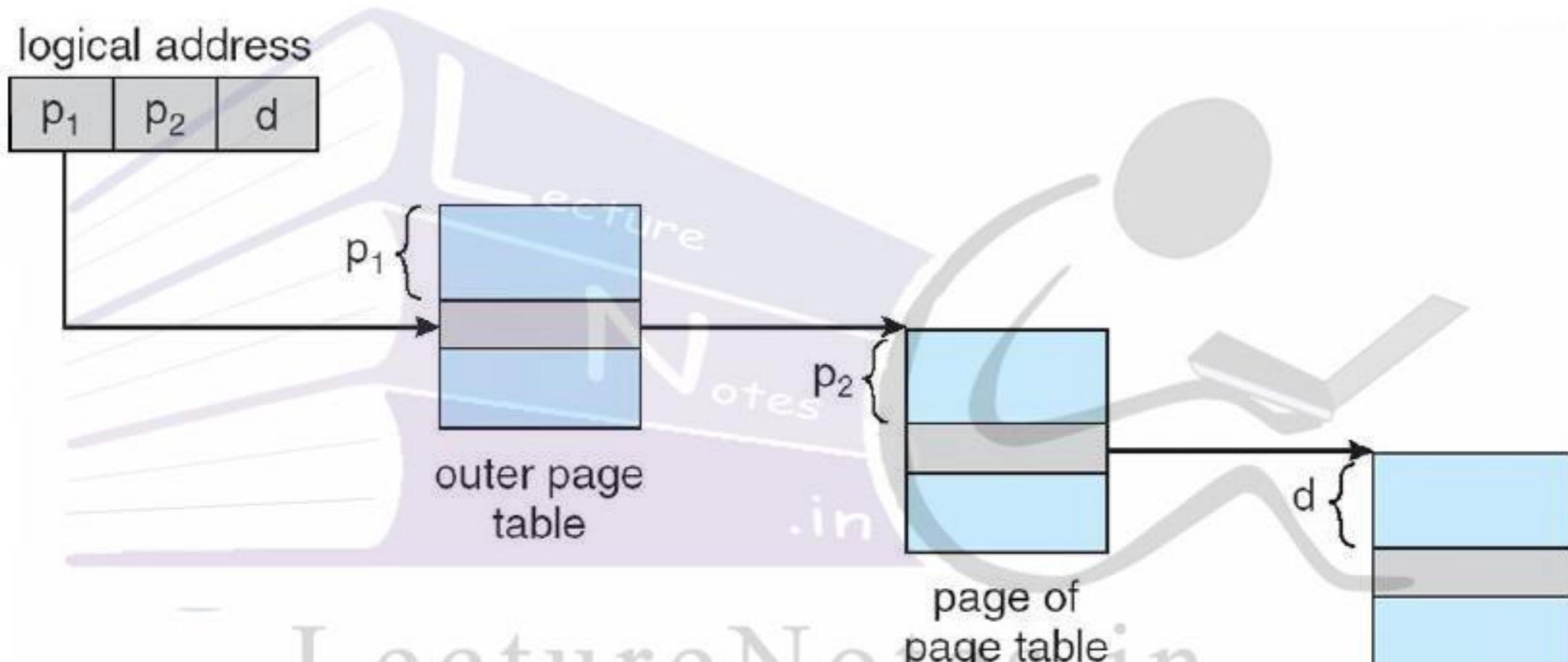
# Two-Level Paging Example

- A logical address (on 32-bit machine with 1K page size) is divided into:
  - a page number consisting of 22 bits
  - a page offset consisting of 10 bits
- Since the page table is paged, the page number is further divided into:
  - a 12-bit page number
  - a 10-bit page offset
- Thus, a logical address is as follows:

page number	page offset
$p_1$	$p_2$
12	10
	10

- where  $p_1$  is an index into the outer page table, and  $p_2$  is the displacement within the page of the inner page table
- Known as **forward-mapped page table**

# Address-Translation Scheme



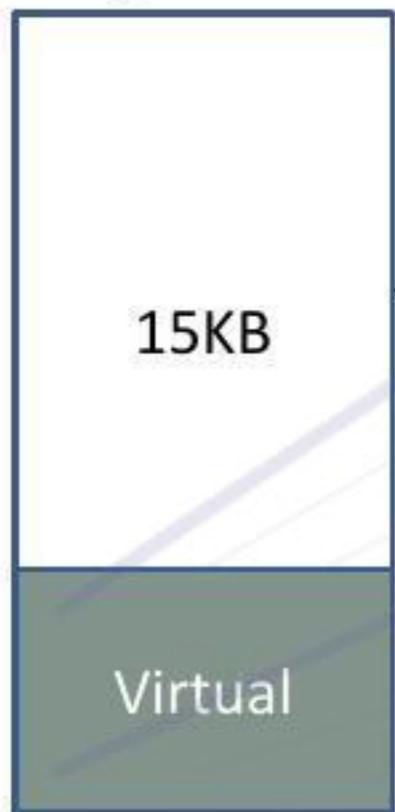
# Lecture Notes in Virtual Memory

- Whenever the process size is more than the MM the OS helps us in executing the process.
- Virtual memory is a technique that allows the execution of the processes that may not be completely in the memory.
- Programs can be executed even program size (logical address space) is greater than the physical available memory.
- The main advantage that the more larger programs can be executed in the physical memory which size is less.
- Ex. Suppose a program size is 15MB but available memory is 12MB. So,, 12MB is loaded in main memory, remaining 3MB loaded in secondary memory. When this 3MB will be needed fro the execution then swap out 3MB swap out 3MB from main memory and swap in 3MB from secondary memory to main memory.
- Virtual memory is a separation of user logical memory from physical memory. This allows an extremely large VM to be provided for programmer when only a small physical memory is available.

# Advantages of Virtual Memory

- Efficient main memory utilization
- Ability to execute a program that is partially available in the main memory.
- More programs can be loaded to the main memory so that degree of multiprogramming will be achieved and with more CPU utilization.
- Less I/O will be needed to load or swap each user program into memory. So, each user program would run faster.
- It allows processes to easily share files address spaces(logical memory) through page sharing. Due to this, it provides efficient mechanism for process creation.
- Virtual memory is implemented by demand paging.
- Here both the system and user will be benefited.

Program Size



RAM

Loaded

10KB

15KB

Virtual

Physical Execution  
Logical Execution

Seems to be  
extended by CPU

RAM

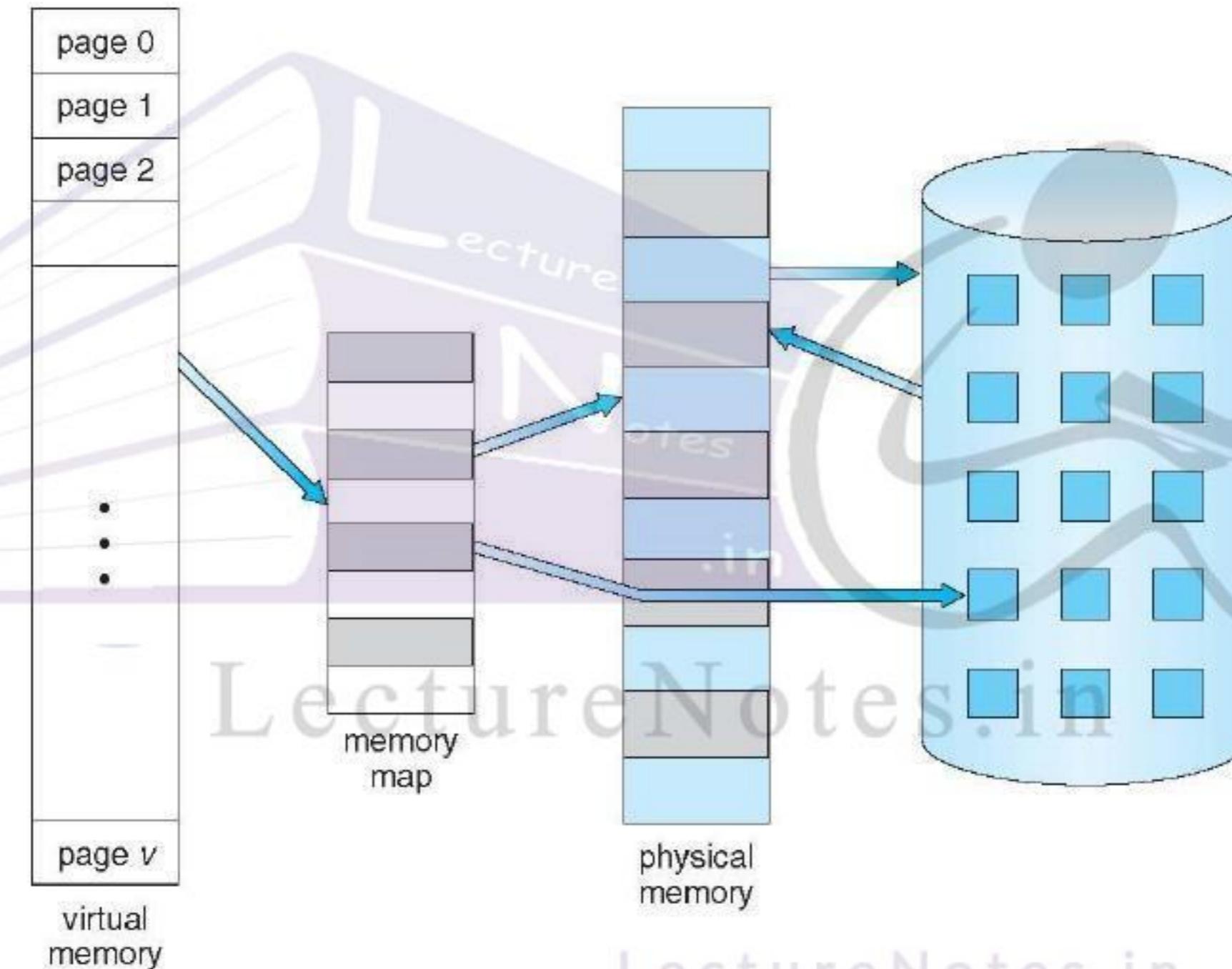
15KB

Virtual

10KB

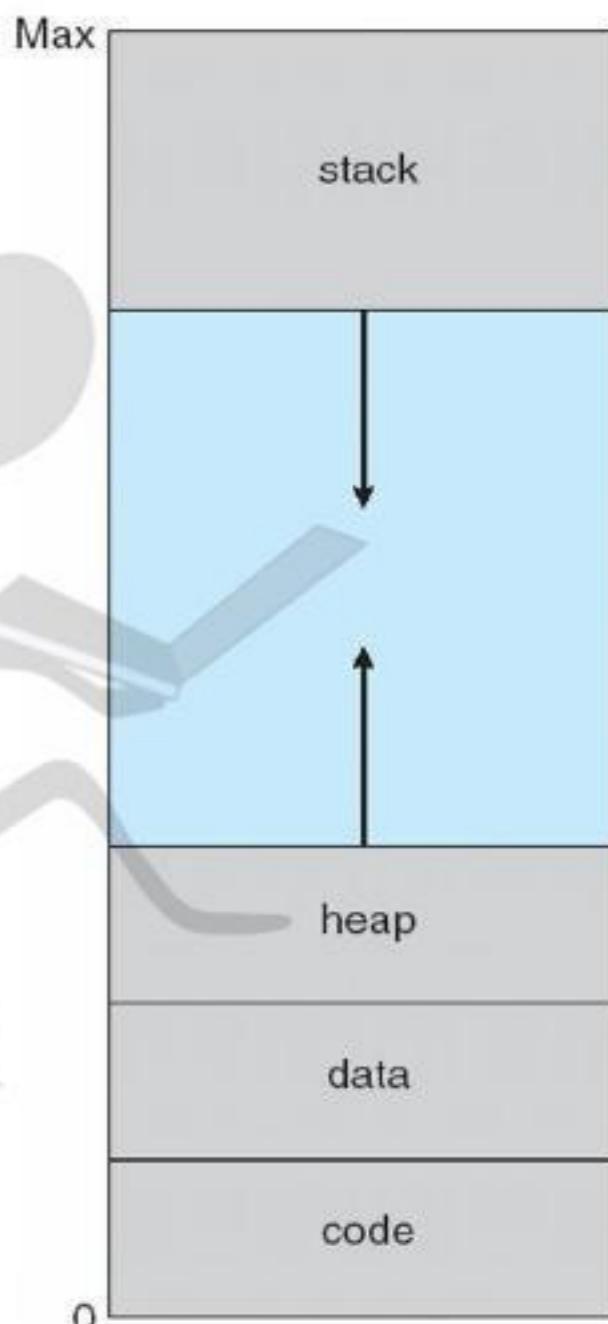
Load

# Virtual Memory That is Larger Than Physical Memory



# Lecture Notes in Virtual-address Space

- The virtual address space of a process refers to the logical (or virtual) view of how a process is stored in memory.
- Usually design logical address space for stack to start at Max logical address and grow “down” while heap grows “up”
  - Maxizes address space use
  - Unused address space between the two is hole
    - No physical memory needed until heap or stack grows to a given new page
- Enables **sparse** address spaces with holes left for growth, dynamically linked libraries, etc
- System libraries shared via mapping into virtual address space
- Shared memory by mapping pages read-write into virtual address space
- Pages can be shared during `fork()`, speeding process creation

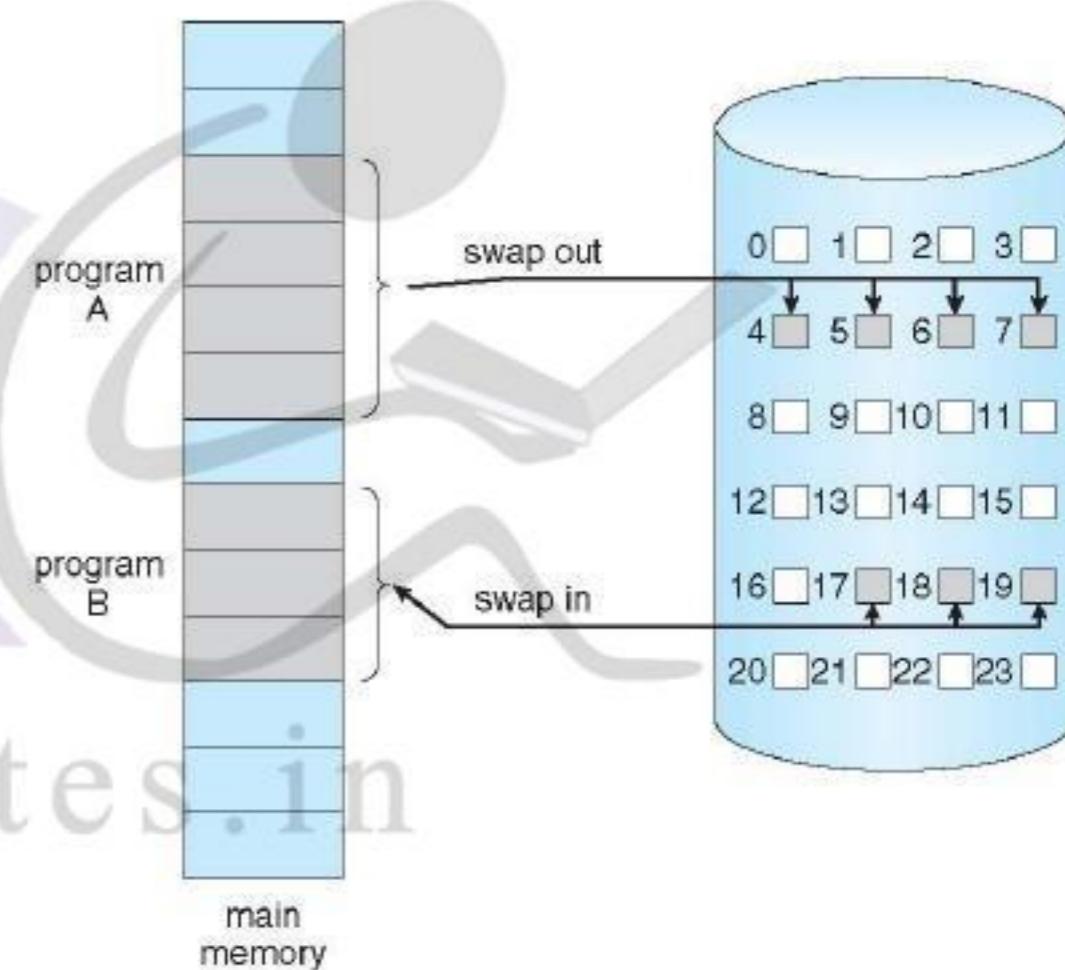


# Demand Paging

- The virtual memory is implemented by demand paging.
- The criteria of demand paging is “A page will not be loaded into the main memory from the SSD until it is needed”.
- The page will be loaded only when they are demanded during execution.
- A demand paging system is similar to a paging system with swapping where processes reside in secondary memory.
- Generally by using a lazy swapper a page loaded only when it is needed.

# Demand Paging

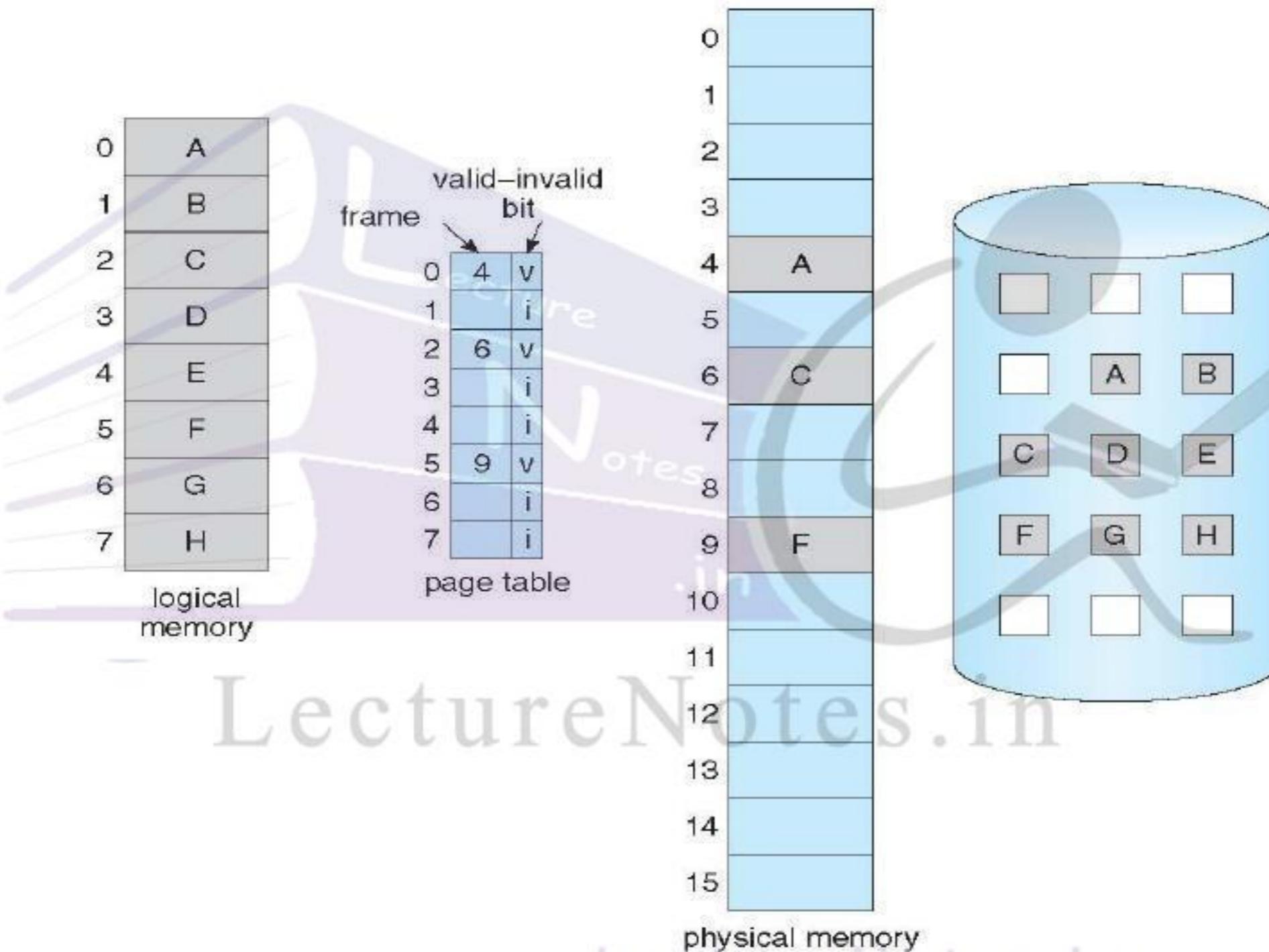
- Bring a page into memory only when it is needed
  - Less I/O needed, no unnecessary I/O
  - Less memory needed
  - Faster response time
  - More users
- Similar to paging system with swapping
- Page is needed  $\Rightarrow$  reference to it
  - invalid reference  $\Rightarrow$  abort
  - not-in-memory  $\Rightarrow$  bring to memory
- **Lazy swapper** – never swaps a page into memory unless page will be needed
  - Swapper that deals with pages is a **pager**



# Basic Concepts

- With swapping, pager guesses which pages will be used before swapping out again
- Instead, pager brings in only those pages into memory
- We need hardware support to distinguish between those pages that are in main memory and those pages that are in the disk.
- This can be done by using a “valid” and “invalid” bit.
- When this bit is “valid” this value indicates that the associated page is both legal and in memory.
- If the bit is “invalid” this value indicates that the page either not valid (i.e. not in the logical address space of the process) or is valid but is currently on the disk.

## Page Table When Some Pages Are Not in Main Memory

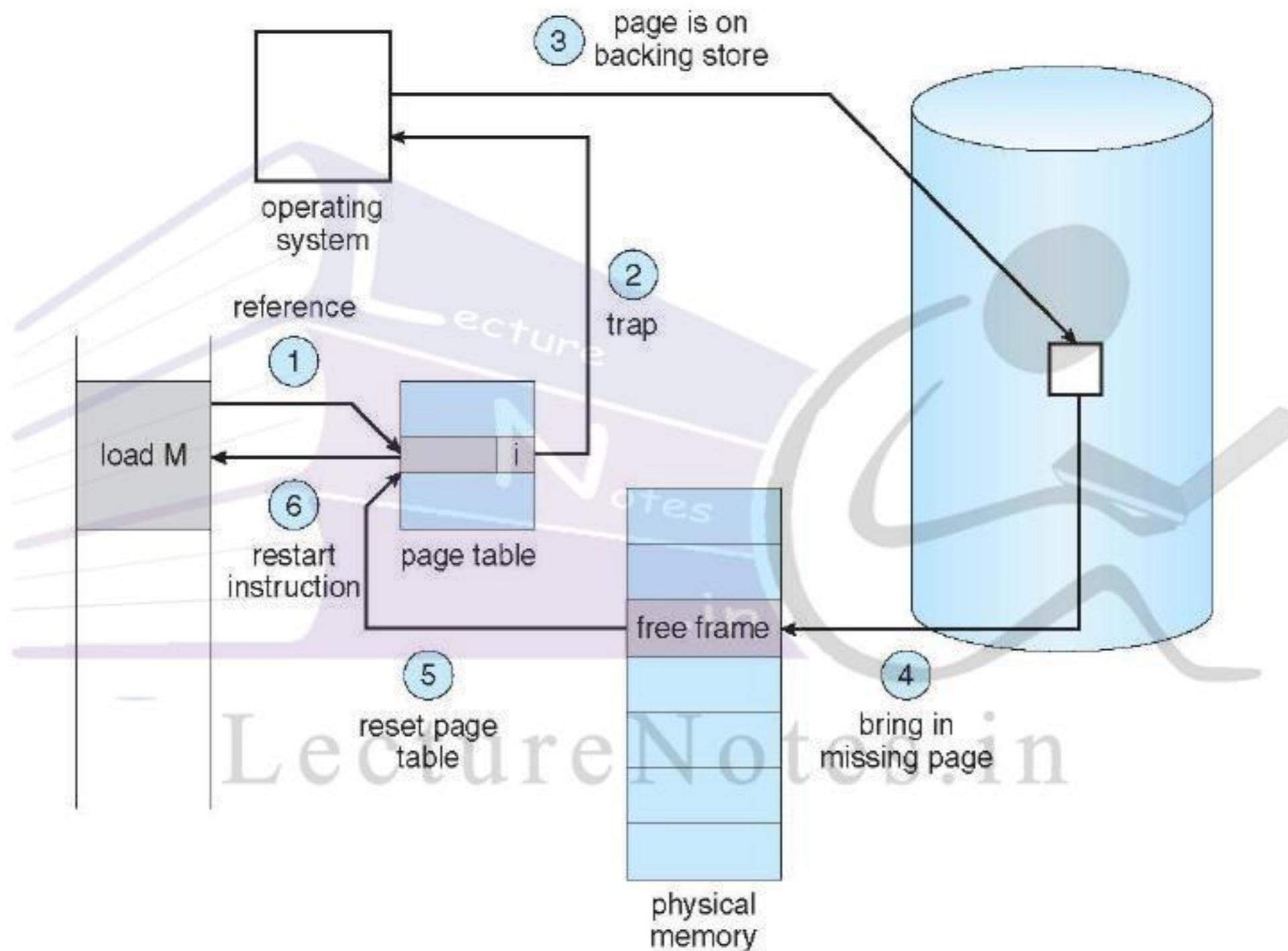


# LectureNotes.in

## Page Fault

- When a process needs to execute a particular page, that page is not available in the main memory then there will be page fault trap or interrupt.
- It means OS fails to bring the desired page into the memory.
- The procedure for handling the page fault is as follows:
  - The memory address requested is first checked, to make sure it was a valid memory request.
  - If the reference was invalid, the process is terminated. Otherwise, the page must be paged in.
  - A free frame is located, possibly from a free-frame list.
  - A disk operation is scheduled to bring in the necessary page from disk. ( This will usually block the process on an I/O wait, allowing some other process to use the CPU in the meantime. )
  - When the I/O operation is complete, the process's page table is updated with the new frame number, and the invalid bit is changed to indicate that this is now a valid page reference.
  - The instruction that caused the page fault must now be restarted from the beginning, ( as soon as this process gets another turn on the CPU. )

# Steps in Handling a Page Fault



# Pure Demand Paging

- When OS sets the instruction pointer to the first instruction of the process, which is on a non-memory resident page, the process immediately faults for the page.
- After this page is brought into memory the process continues to execute, faulting as necessary until every page that it needs is in memory. At this point it can execute with no more faults. This is known as pure demand paging.

LectureNotes.in

LectureNotes.in

# Performance of Demand Paging

## ■ Stages in Demand Paging (worse case)

1. Trap to the operating system
2. Save the user registers and process state
3. Determine that the interrupt was a page fault
4. Check that the page reference was legal and determine the location of the page on the disk
5. Issue a read from the disk to a free frame:
  1. Wait in a queue for this device until the read request is serviced
  2. Wait for the device seek and/or latency time
  3. Begin the transfer of the page to a free frame
6. While waiting, allocate the CPU to some other user
7. Receive an interrupt from the disk I/O subsystem (I/O completed)
8. Save the registers and process state for the other user
9. Determine that the interrupt was from the disk
10. Correct the page table and other tables to show page is now in memory
11. Wait for the CPU to be allocated to this process again
12. Restore the user registers, process state, and new page table, and then resume the interrupted instruction

# LectureNotes.in

## Performance of Demand Paging (Cont.)

- Three major activities
  - Service the interrupt – careful coding means just several hundred instructions needed
  - Read the page – lots of time
  - Restart the process – again just a small amount of time
- Page Fault Rate  $0 \leq p \leq 1$  (probability of page fault occur = $p$ )
  - if  $p = 0$  no page faults
  - if  $p = 1$ , every reference is a fault
- Effective Access Time (EAT)

$$\text{EAT} = (1 - p) \times \text{main memory access time (in nano sec.)} + p (\text{page fault overhead} + \text{swap page out} + \text{swap page in}) \text{ (In mili second)}$$

1 mili seceond = 10,00000 nano second

# Copy-on-Write

- Copy-on-Write (COW) allows both parent and child processes to initially **share** the same pages in memory
  - If either process modifies a shared page, only then is the page copied
- COW allows more efficient process creation as only modified pages are copied
- In general, free pages are allocated from a **pool** of **zero-fill-on-demand** pages
  - Pool should always have free frames for fast demand page execution
    - ▶ Don't want to have to free a frame as well as other processing on page fault
  - Why zero-out a page before allocating it?
- vfork () variation on fork () system call has parent suspend and child using copy-on-write address space of parent
  - Designed to have child call exec ()
  - Very efficient

## Before Process 1 Modifies Page C

