

Main function of operating system must be
• provide maximum throughput (more no. of process in limited time)
• provide convenience (user friendly)

DETAILED
Date / /
@ Pg No.

Operating Systems

Applications

OS

Hardware

What is OS and some types of OS:-

it is a software acting as an interface b/w computer, hardware and software (user application)

- manages resources
- manages incoming user commands
- provides services for successful execution of programs

(User Interface) Functions of OS (Device management)

Device management

process management memory management file system management

Some types of OS - • Windows • macOS • Linux • Unix - Android
• JOS

→ RTOS (Real time Operating System) - Used for devices requiring deterministic and real-time response

used in embedded

systems, control systems,
IoT devices

processing

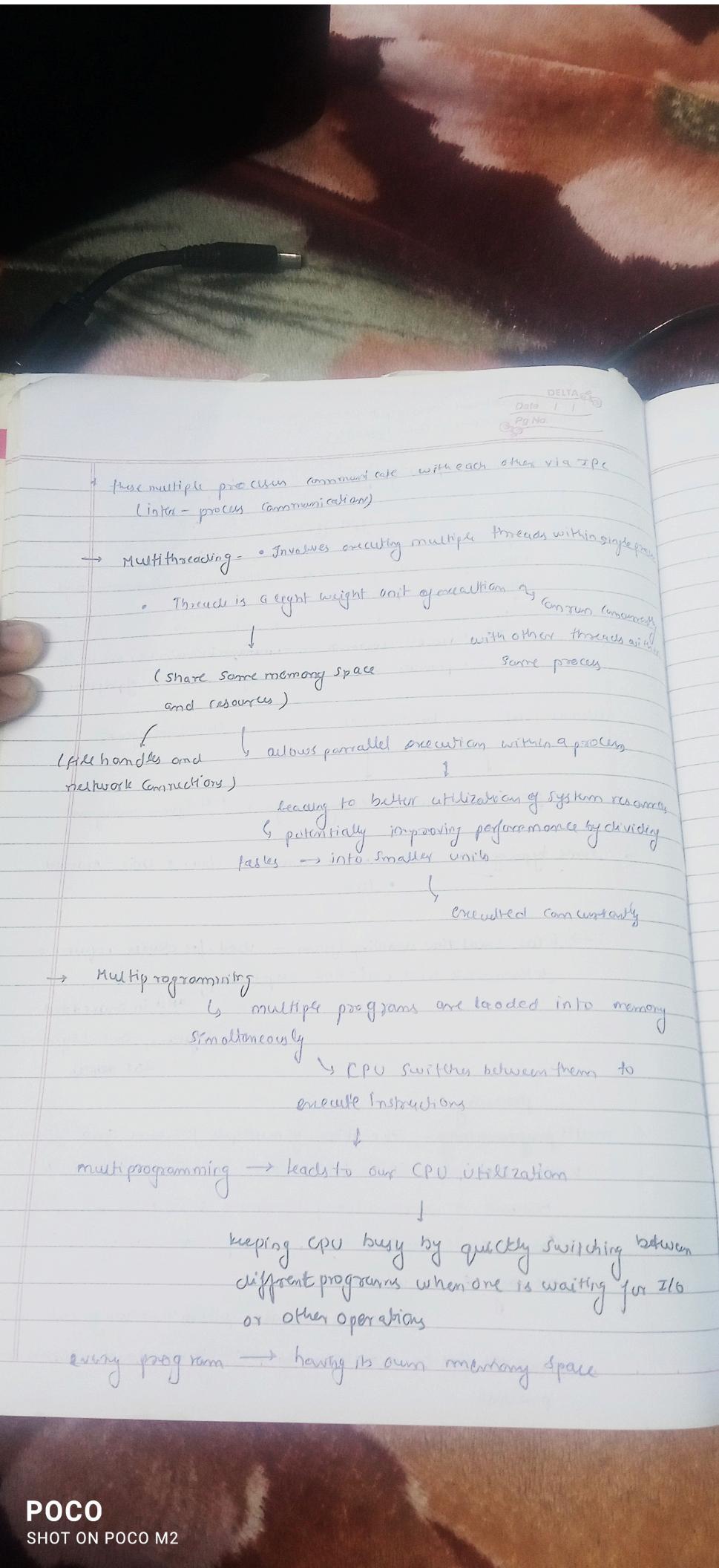
multiprogramming - Execution of multiple processes on a system with multiple CPU's or CPU cores, each process is an instance of a running programme and multiple processes can execute concurrently, each process having its own memory spaces and resources

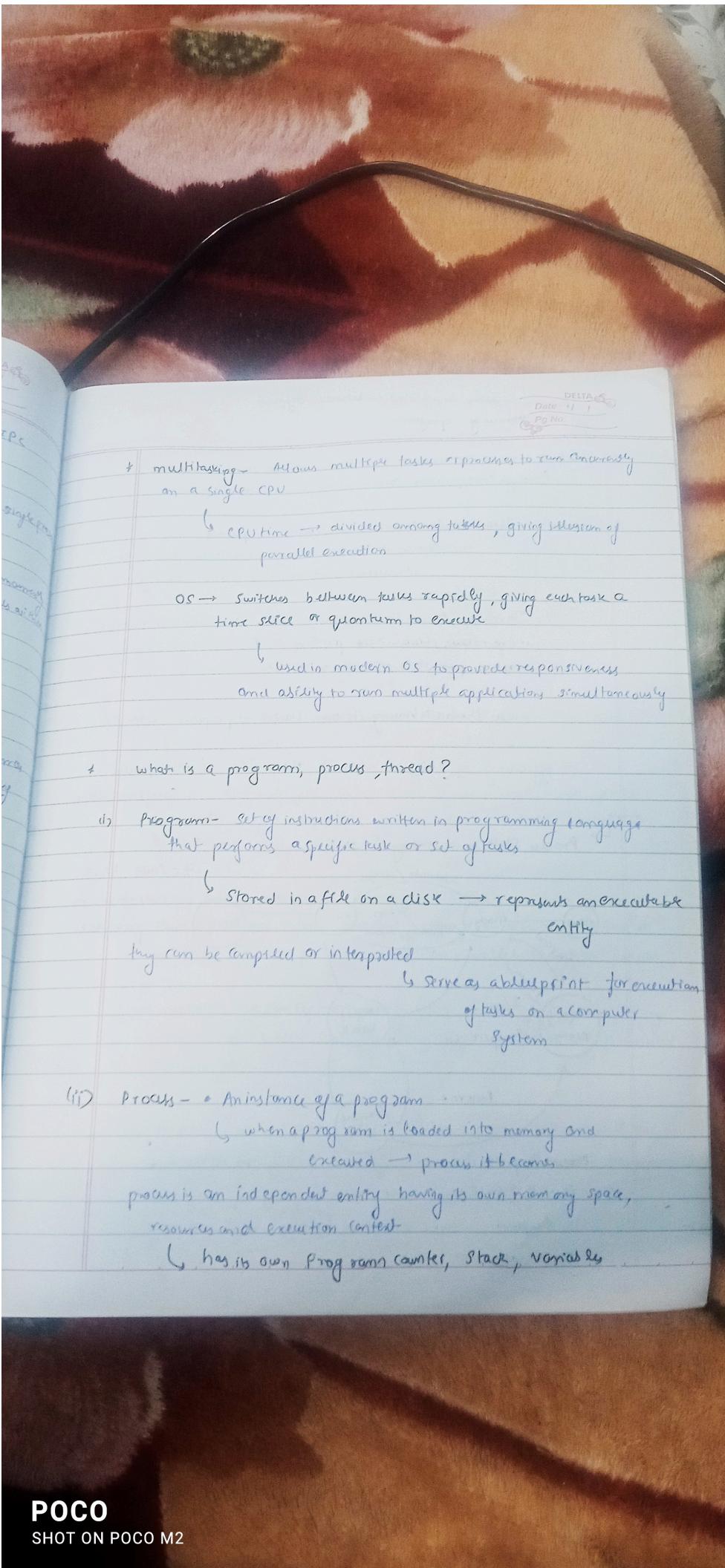
↳ Aim to increase system's throughput and

provides faster execution by dividing workload across multiple processes

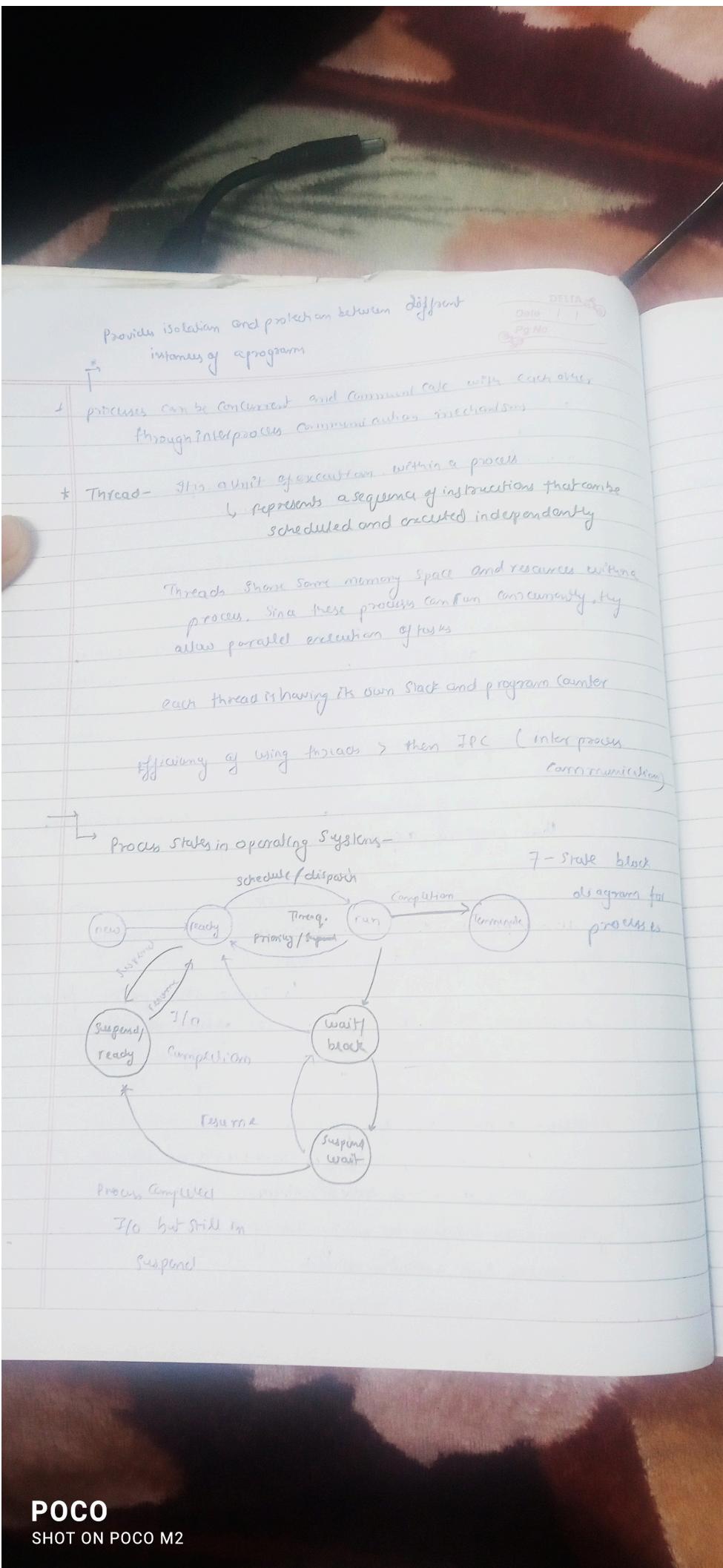
POCO

SHOT ON POCO M2





POCO
SHOT ON POCO M2



From our new state, process gets created and is been forwarded to ready state via long term scheduler

Running state being a queue format accepts the request that has been first input and removes as first output



This removal (dispatch) of processes is known as Scheduling

FIFO

* CPU allocated to the running stage → now there are 2 conditions

(i) If process gets completed, it goes to the termination stage and process gets over

(ii) If process demands for I/O operation, if you to wait / block, process is on hold, after completion of I/O request → this process is sent back to our ready state for execution

→ running state works on the theme of multitasking, process in running state is provided a fixed amount of time (time quantum) → Since we are having a unique processor and no parallel execution possible

* moreover, the concept of priority also matters for a process, if any other process is been currently executed, it would execute prioritized process first

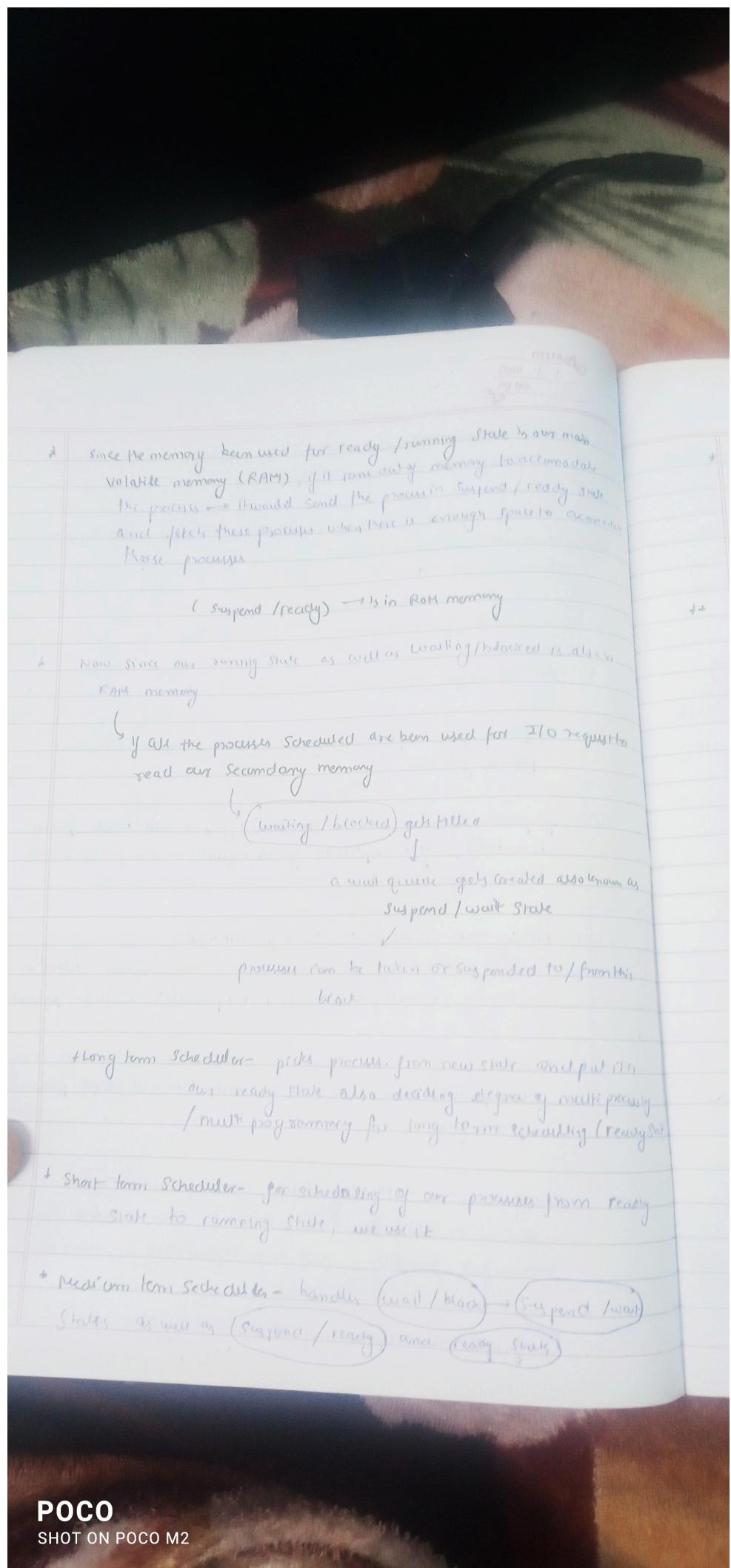
* Case of preemption in running state: When a process takes time T_i for execution and if this time quantum < T_i → execution stops, process gets forcefully removed and repeats itself for same time quantum after wait block

↳ if CPU runs a process without taking care of time

↳ non preemptive process

POCO

SHOT ON POCO M2



+ processes have two types of time

CPU time

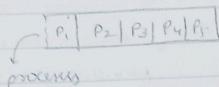
I/O time

Frequency of → short time > medium time > long time
Schedules schedules schedules schedules

++ CPU Scheduling Algorithms -

used in operating systems to determine the order in which processes are executed on the CPU

(i) First Come First Serve (FCFS) -



the process that arrives first is executed first, follows non-pre-emptive approach that means once a process starts running, it continues until it completes or voluntarily gives up the CPU, queue (FIFO) been used

Simple to understand but may lead to poor utilization of the CPU if long processes arrive before shorter ones

(ii) Shortest Job Next (SJN) or shortest-job first (SJF)

• It selects the process with the shortest total execution time next.

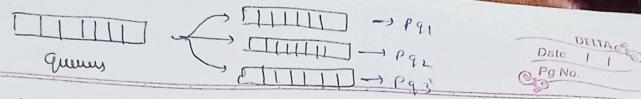
• It can be either non-preemptive or pre-emptive

process continues executing until it gets completed

if a new process with shorter burst time arrives, current running process may be pre-empted

POCO

SHOT ON POCO M2



(v) Multi-level Queue Scheduling - dividing ready queue into multiple priority queues, each with its own scheduling algorithm

processes are initially kept (placed) in highest priority queue and can move between queues based on pre-defined criteria

→ This approach allows for the differentiation of processes based on their priority & characteristics

every queue can individually use a different scheduling algorithm, such as FCFS, SJF or RR suitable for process within queue

(vi) Shortest remaining time (SRTF) - / SRTF

A pre-emptive version of the shortest job first which we have where the processor is allocated to the job closest to completion

↳ process with smallest amount of time remaining until completion is selected to execute.

(vii) Longest remaining time first (LRTF) / LRN

↳ a pre-emptive version of longest job first scheduling algorithm, used by operating system to program incoming processes for use in systematic way → schedules those processes first which have longest processing time remaining for completion

(viii) Highest response ratio next (HRRN)

↳ non-preemptive CPU scheduling algorithm, finding response ratio of all available processes and select the one with highest response ratio.

* * burst time → execution time

DELTA
Date 1 / 1
Pg No. 3

SJN/ SJF aims to minimize the average waiting time and is suitable when burst times are known in advance

(iii) Round Robin (RR) -

preemptive scheduling algorithm assigning a fixed time quantum (e.g. 10 milliseconds) to each process in circular manner

if execution time of process < time quantum → moves back to the ready queue allowing next process to execute

RR

providing fair execution to all processes but may suffer

high context switching overhead

(may not be efficient for long-running processes)

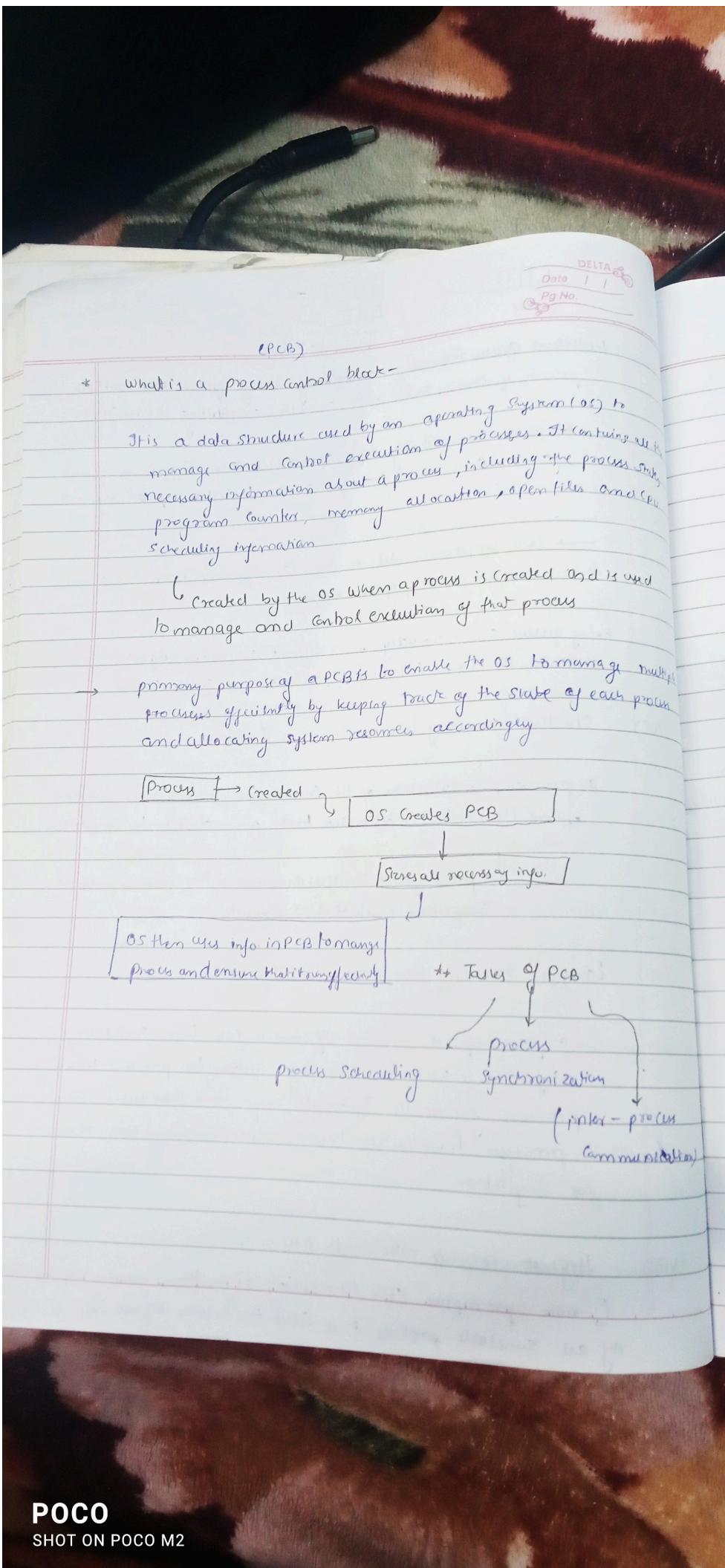
(iv) Priority Scheduling - priority → 4 3 1 2
$$\begin{bmatrix} P_1 & P_2 & P_3 & P_4 \end{bmatrix}$$

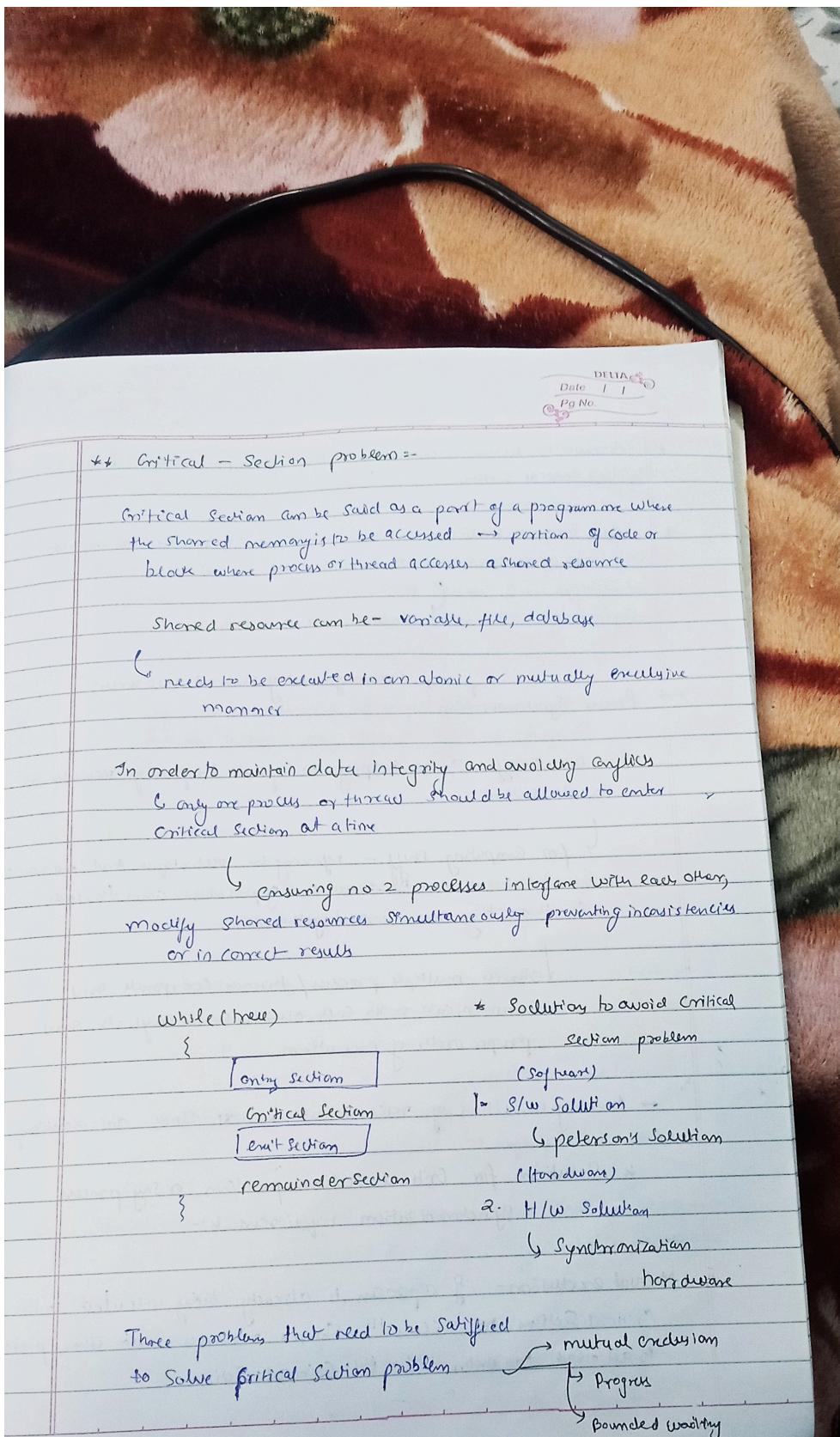
Assigns priority value to each process, CPU is allocated to the process with highest priority, can be pre-emptive and non-pre-emptive

Pre-emptive - if a higher priority process arrives, currently running process may be pre-empted

Non-preemptive - process continues executing until it completes or voluntarily gives up the CPU

problem of starvation for lower priority processes never get chance to execute





POCO
SHOT ON POCO M2

DELTA
Data 1 / 1
Pg No.

+ Remainder section - remaining portions of the program excluding critical section

+ Race condition - When our final output depends on the order in which variables are accessed
↳ race condition

→ Process synchronization → (management of shared resources)

This is like a traffic signal helping to regulate the flow of vehicles at an intersection

↳ for computing stuff - referring to techniques and mechanisms used to coordinate the execution of processes and threads so that they can work together

↳ ensures multiple processes / threads cooperate and communicate with each other to avoid conflicts and proper order of execution

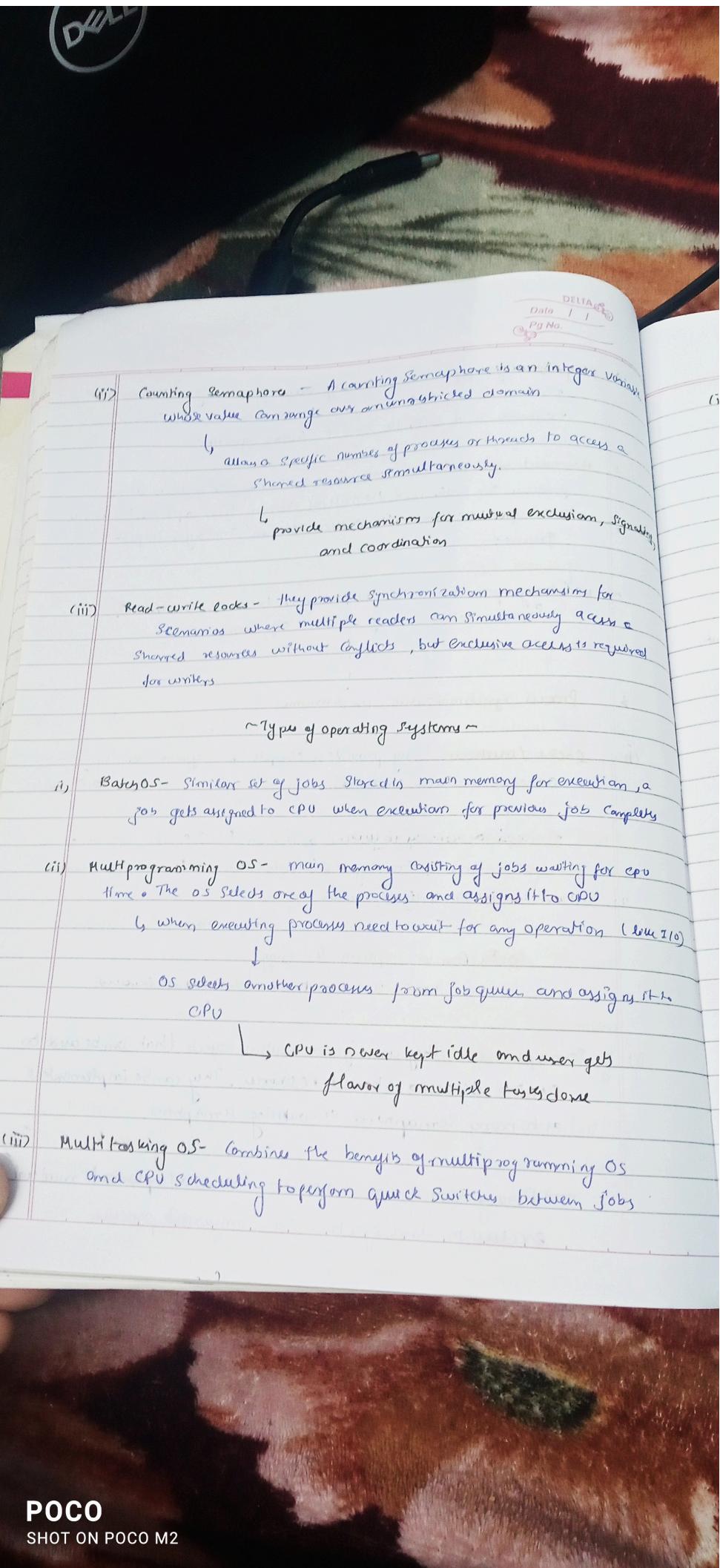
→ Avoiding conditions of race condition, deadlock, data inconsistency

% Solution for Critical section problem or say process synchronization requirements:-

④ Mutual exclusion - If a process is already being executed in its critical section to access its shared resources, no other process is allowed to enter into the critical section

POCO
SHOT ON POCO M2

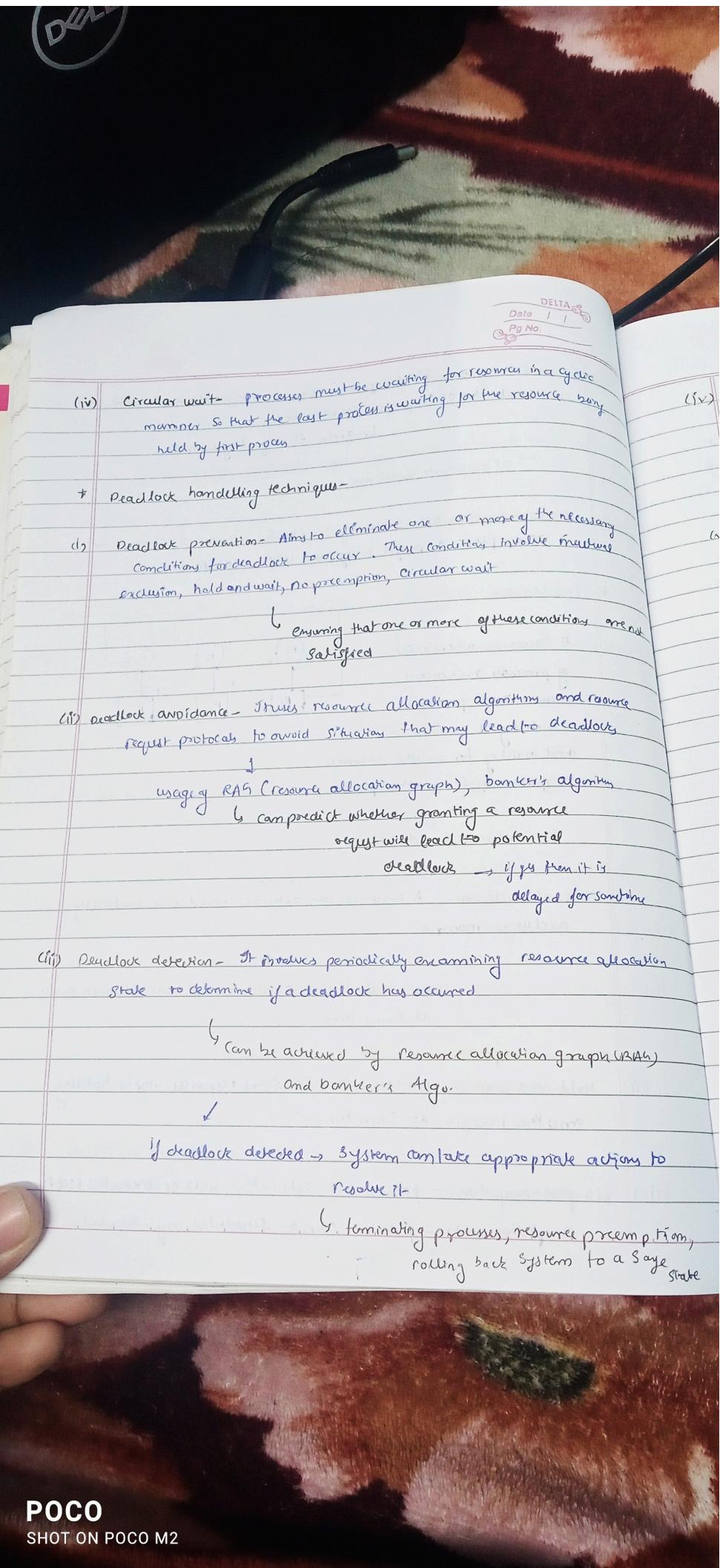
- + **Priority Pre-emption** - If one process is not being executing presently in its Critical Section, it can only wait until the remainder releases the Shared Resource and enter the Critical Section again.
- + **Bounded waiting** - There is a bound on the number of times other processes can enter into the Critical Section after a process has made a request to access the Critical Section and before request is granted.
 - ↳ Every process would get a chance to enter the Critical Section without any chance of starvation.
- **Process Synchronization mechanisms**-
 - (i) **Locks / mutexes** - They provide a simple and effective way to achieve mutual exclusion.
 - ↳ Allow only 1 process or thread to acquire lock at a time, ensuring exclusive access to shared resource or critical section.
 - ↳ *hardware implementation*
 - ↳ Locks can be implemented using *Software Constructs*
 - (ii) **Semaphores** - They are synchronization objects that can be used to control access to shared resources. They can be implemented as binary semaphores or counting semaphores.
 - **Binary Semaphores (0 or 1)** - Used to implement mutual exclusion and synchronize concurrent processes.

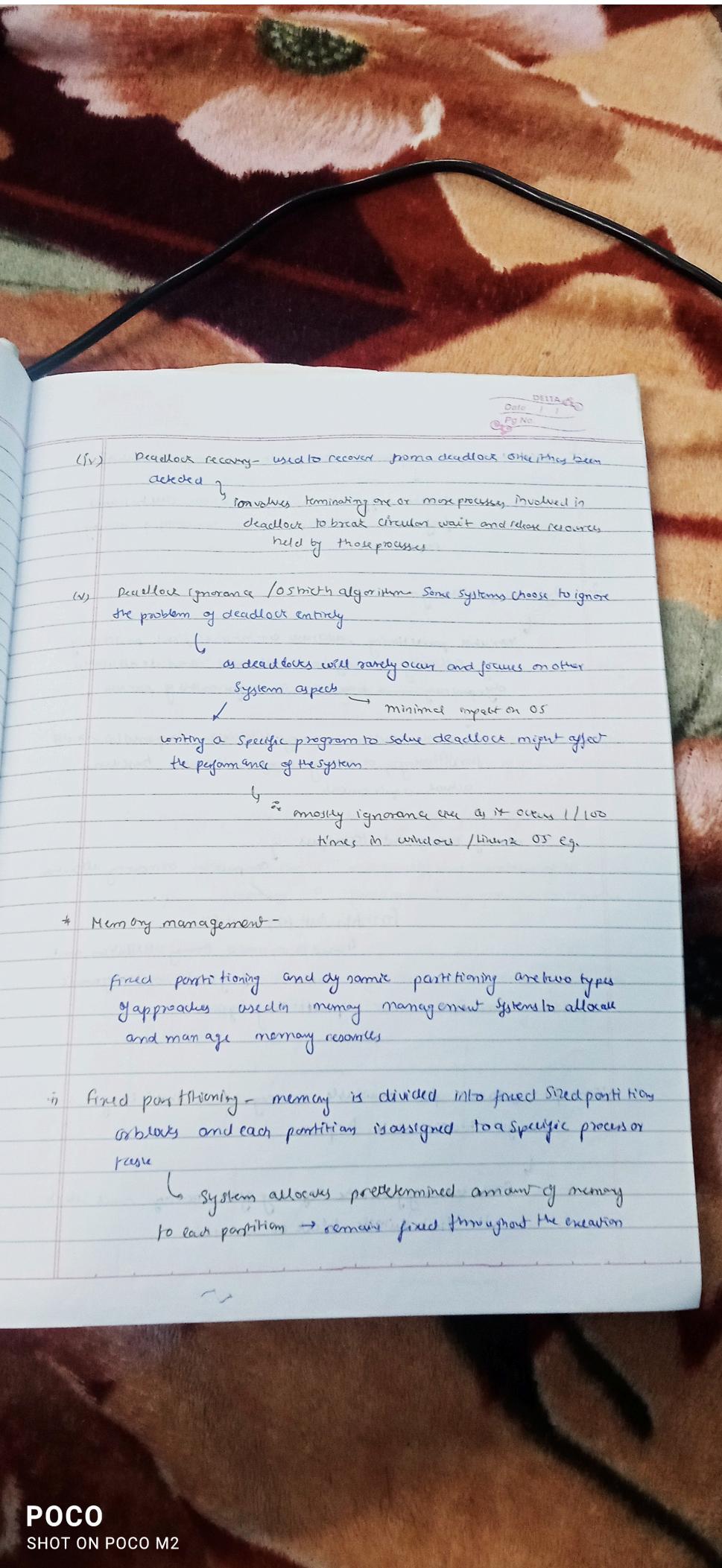


- DELLA
Date 11
Pg No.
- (iii) Time sharing OS - these systems require interaction with the user to interact the OS to perform various tasks
 ↳ given through keyboard
- OS responds with an off.
- (iv) RTOS (Realtime OS) - usually built for dedicated systems to accomplish a specific set of tasks within deadlines
- * Deadlocks
- A situation where a set of processes are blocked because each process is holding a resource and waiting for another resource required by some other process
- Diagram illustrating Deadlocks:
- ```

 graph LR
 P1((P1)) --> R1[R1]
 P2((P2)) --> R2[R2]
 P3((P3)) --> R3[R3]
 R1 --> P2
 R2 --> P3
 R3 --> P1

```
- \* Condition for deadlock -
- (i) Mutual exclusion - A resource can only be shared in a mutually exclusive manner  
 ↳ implying that two processes cannot use same resource at same time
- (ii) Hold and wait → process waits for some resources while holding another resource at same time.
- (iii) No pre-emption - process once scheduled will be executed till the completion. No 2 processes can be scheduled by scheduler.





→ relatively simple to implement and provides fast memory allocation

↓  
→ leads to inefficient memory utilization due to internal fragmentation → when processes have varying memory requirements

## 2. Dynamic partitioning -

↳ variable partitioning, addresses limitation of fixed partitioning by allowing memory to be allocated and deallocated dynamically → based on size requirements of processes

↓  
→ provides better memory utilization compared to fixed partitioning, as memory may be allocated based on actual requirement

→ memory management techniques ↳  
compaction or memory allocation

↙ first fit, best fit, worst fit

↳ used to optimize memory utilization and minimize fragmentation in dynamic partitioning systems

→ first fit - allocates first available memory block that is large enough to accommodate the process

↳  
→ starts searching from the beginning of the memory and selects the first suitable block encountered

(iii) Best-fit - Searches for the smallest available memory block that is large enough to accommodate the process.

(iv) Worst-fit - Arranging process is allocated the block of memory in which it leaves maximum gap.

#### \* Paging

It is a storage mechanism used to remove processes from the secondary storage into the main memory in the form of pages.

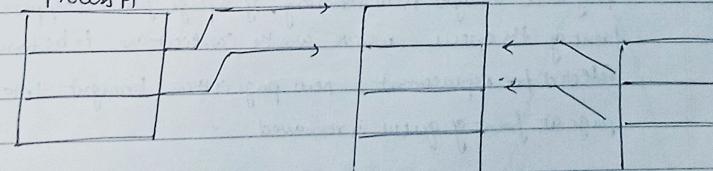
Main idea behind Paging → divide each process in form of pages  
the main memory will also be divided in form of frames (pages)

One page of process is stored in one of the frames of memory

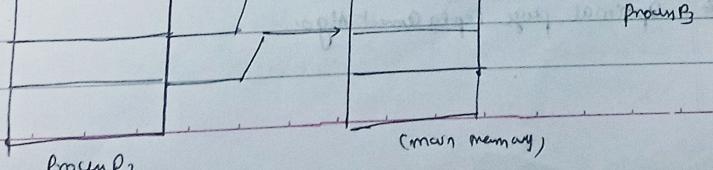
The pages can be stored at different locations of the memory but priority is always given to contiguous holes or frames

Pages of the process are brought into main memory only when they are required or else they remain in Secondary Storage

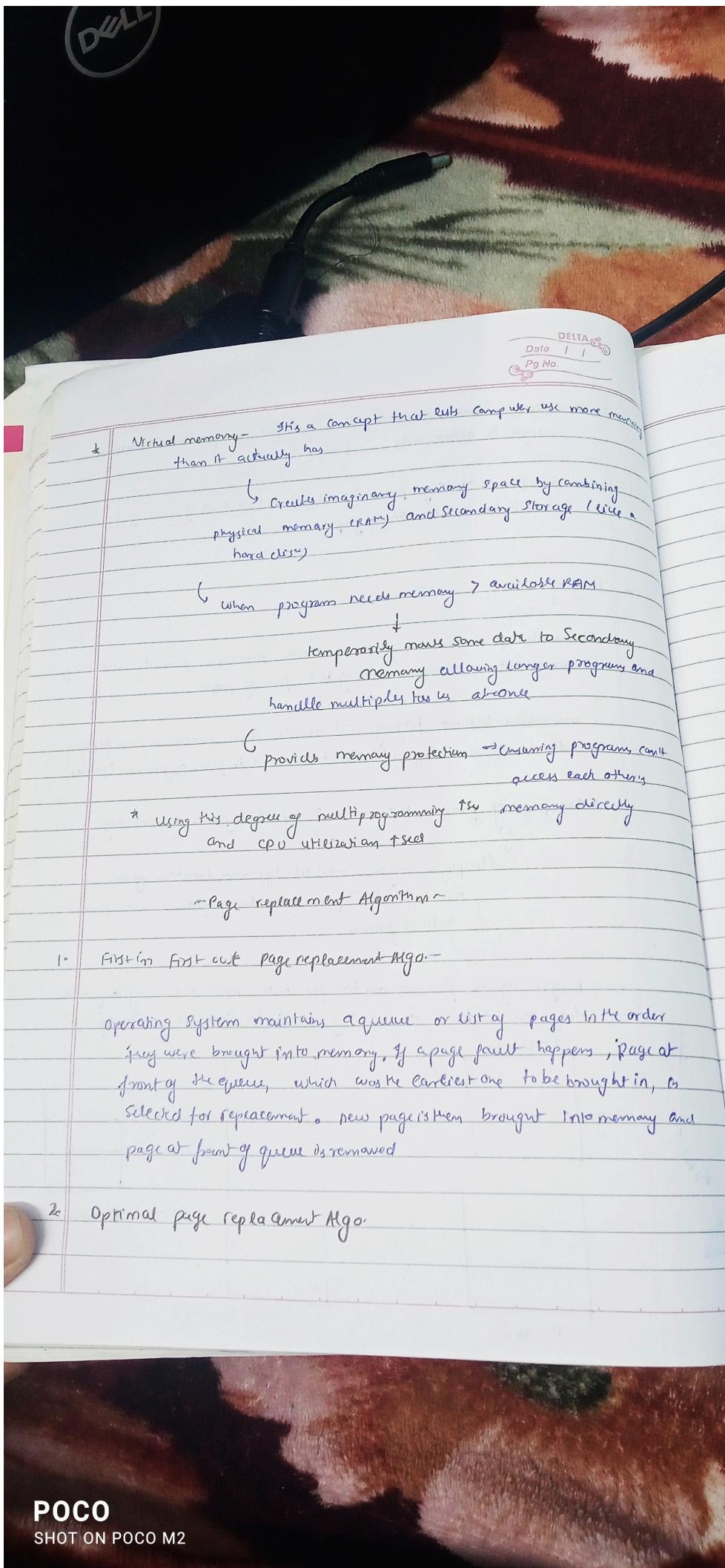
Process P<sub>1</sub>



Process P<sub>2</sub>



(main memory)



Pages are replaced which would not be used for longest duration of time → providing best possible performance by reducing no. of page faults

↳ practically not possible in real world scenarios

knowledge of prior references  
also necessary to impossible

(LRU)

3. Least recently used -

Page will be replaced which is least recently used

LRU works on principle that pages that have been recently accessed are more likely to be accessed again and again in near future

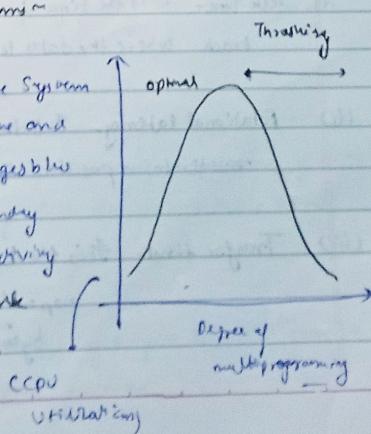
↳ helps to retain frequently accessed pages in memory.

↳ Reducing no. of page faults and thus overall system performance

### ~ Thrashing in Operating Systems ~

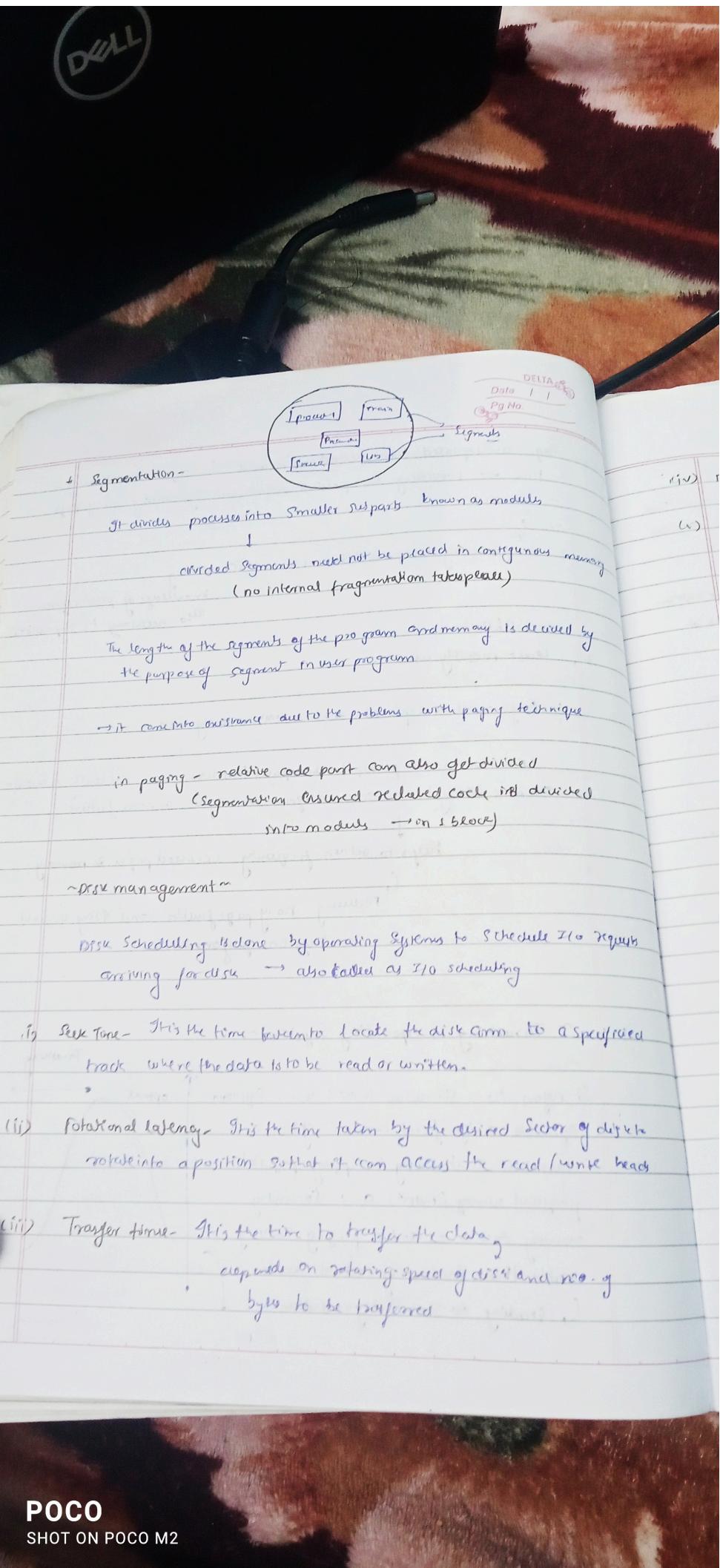
It refers to a situation in OS where system spends a significant amount of time and resources continuously swapping pages between physical memory (RAM) and secondary memory due to excessive paging activity rather than executing actual work

↳ leading to severe degradation



**POCO**

SHOT ON POCO M2



- DELTABOOK  
Date 1/1  
Pg No. 1
- (iv) Disk access time = seek time + rotational latency + transfer time
- (v) Disk response time = average time spent by a request waiting to perform its I/O operation  
↳ it is the response time of all requests
- Disk scheduling algorithms
- It is a key component of disk management that determines the order in which I/O (input/output) requests are processed and serviced by the disk controller  
↳ minimizes the seek time and rotational latency and maximizes overall disk performance
- (vi) FCFS (First Come First Served) → Simplest disk scheduling algorithm that processes requests in order they arrive  
↳ does not consider position of disk head or proximity of requests  
→ can result in potential delays if a long seek time between requests.
- (vii) SSTF (Shortest Seek Time First) → this algorithm selects the request with shortest seek time from current position of disk head → minimizes average seek time and reduces overall disk access time  
↳ may lead to starvation of requests located farther away from current position

