

Operating System

Introduction

→ A program that acts as an intermediate b/w a user of a computer and the computer hardware. (user & computer hardware)

• Goals of OS

- Execute user programs and make solving user problem easier.
- Make computer system more convenient to use.
- Use the computer hardware in efficient manner.

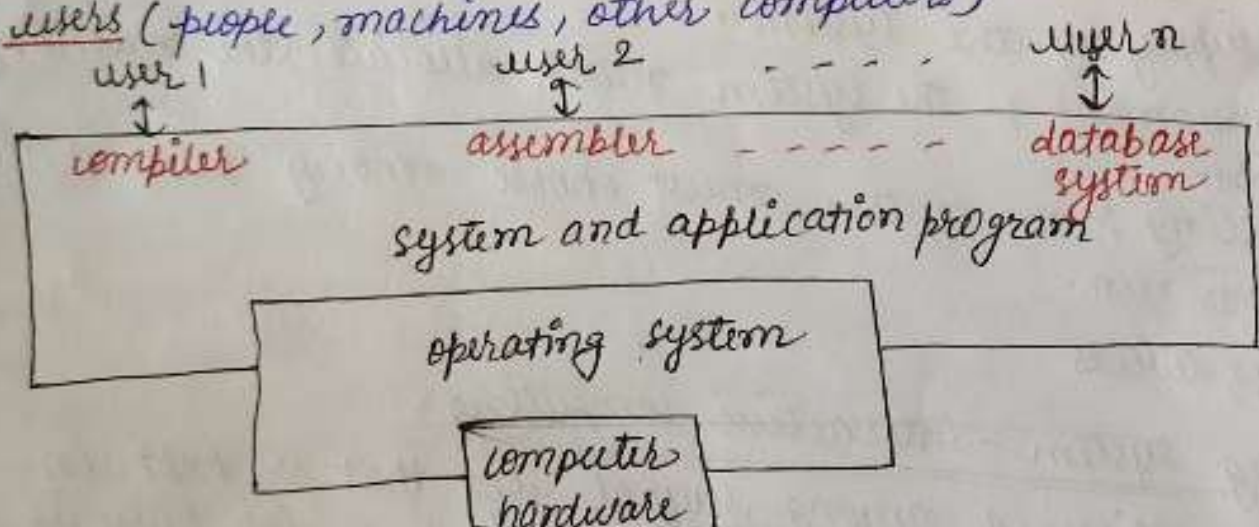
computer system components

1. Hardware: provides basic computing resources (cpu, memory, I/O devices).

2. operating system: controls and coordinates the use of the hardware among the various applications programs for the various users.

3. Application programs: defines the ways in which the system resources are used to solve the computing problems of the users.

4. users (people, machines, other computers)



⊛ operating system definitions

i) resource allocator: manages and allocates resources.

ii) control program: controls the execution of user program and operations of I/O devices

iii) kernel: the one program running at all times (all else being application programs)

⊛ mainframe systems

→ Reduce setup time by batching similar jobs.

→ Automatic job sequencing: automatically transfers control from one job to another. first rudimentary operating system.

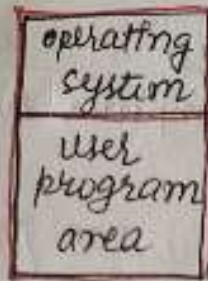
→ Resident monitor

→ initial control in monitor.

→ control transfer to job.

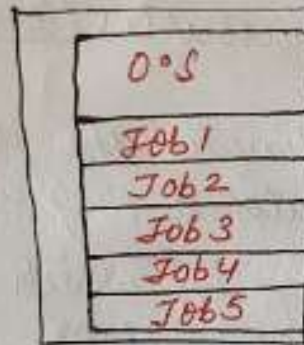
→ when job completes control transfer back to monitor.

memory layout for a simple Batch system.



★ Multiprogrammed Batch system.

several jobs are kept in main memory at the same time, and the CPU is multiplexed among them.



★ OS features needed for multiprogramming.

- i) I/O routine supply by the system.
- ii) memory management :- The system must allocate the memory to several jobs.
- iii) CPU scheduling :- the system must choose among several jobs ready to run.
- iv) Allocation of devices

★ Time sharing system - interactive computing.

- The CPU is multiplexed among several jobs that are kept in memory and on disk (the CPU is allocated to a job only if the job is in memory)
- A job swapped in and out of the memory to the disk
- online-communication between user and computer system when the operating system finishes the execution of one command, it seeks for next "control statement" from the user's keyboard.
- on-line system must be available for users to access data and code.

★ Desktop Systems

- i) personal computer:- computer system dedicated to a single user.
- ii) I/O devices:- keyboard, mice, display screen (monitor) etc.
- iii) user convenience and responsiveness.
 - can adopt technology developed for larger operating system.
 - often individuals have sole use of computer and do not need advanced CPU utilization or protection features.
 - may run several different types of operating systems (Windows, UNIX, Linux).

★ Parallel System

- multiprocessor system with more than one CPU in close communication.
- Tightly coupled system:- processors share memory and clock; communication usually takes place through the shared memory.
- Advantages of parallel system
 - Increased throughput
 - Economical
 - Increased reliability
 - i) graceful degradation
 - ii) fail-soft system.

★ Parallel system (cont.)

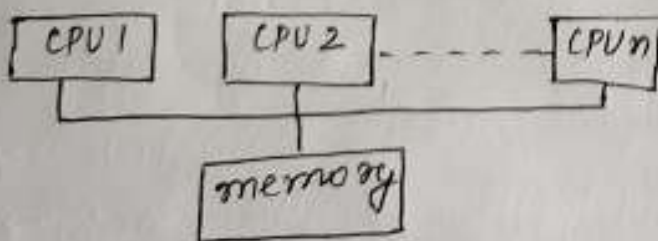
★ Symmetric multiprocessing (SMP)

- each processor runs an identical copy of the O.S.
- runs without performance deterioration.
- most modern O.S. supports SMP.

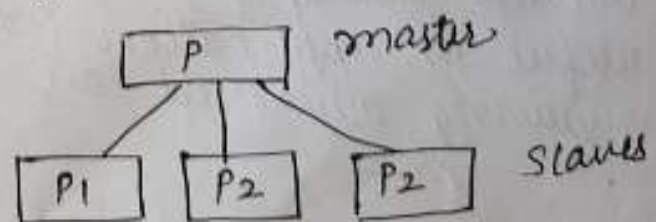
★ Asymmetric multiprocessing (AMP)

- Each processor is assigned a specific task; master processor schedules and allocates work to slave processor.
- more common in extremely large system.

Symmetric multiprocessing



Asymmetric multiprocessing



★ Distributed system

- Distributes the computations among several physical processors.
- loosely coupled system :- each processor has its own local memory; processors communicate with one another through various communication lines, such as high-speed buses or telephone lines.
- Advantages of distributed system
 - Resource sharing
 - Computation speed up - load sharing
 - Reliability
 - Communications.

④ Distributed systems (cont)

- Requires networking infrastructure
- local area network (LAN) or wide area network (WAN)
- may be either client-server or peer-to-peer system.

★ Real time system

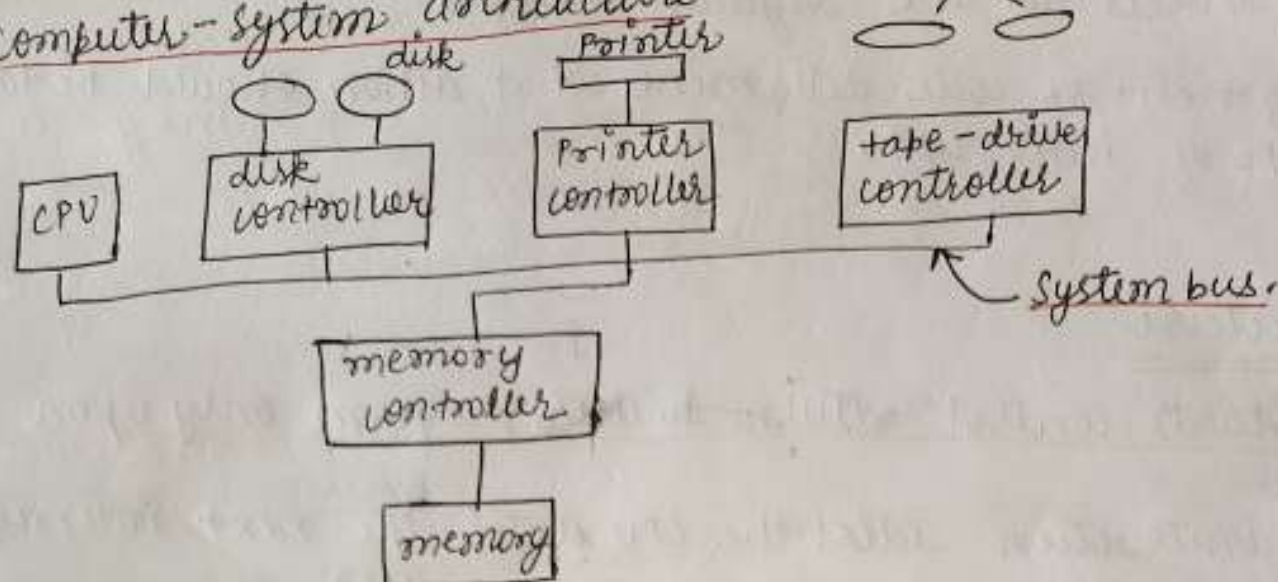
- often used as control device in a dedicated application such as controlling scientific experiment, medical imaging system, industrial control system, and some display system.
- well defined fixed-time constraints
- Real-time system may be either hard or soft real-time.

Real time system (cont.)

- ★ → Hard real-time :- secondary storage limited or absent, data stored in short term or ROM
 - conflict with time-sharing systems, not supported by general purpose O.S
- ★ → Soft real-time
 - limited utility in industrial control of robotics.
 - useful in application (multimedia, virtual reality) requiring advanced O.S features.

Computer-system structures

① Computer-system architecture



★ Computer-system operations.

- ① I/O devices and the CPU can execute concurrently.
- ② Each device controller is in charge of a particular device type.
- ③ Each device controller has a local buffer.
- ④ CPU moves data from/to main memory to/from local buffers.
- ⑤ I/O is from the device to local buffer of controller.
- ⑥ Device controller informs CPU that it has finished its operation by causing an interrupt.

→ functions of Interrupts

- ① Interrupt transfers control to the interrupt services routine generally, through the interrupt vector, which contains the addresses of all the services routines.
- ② Interrupt architecture must save the address of the interrupt instruction.
- ③ Incoming interrupt are disabled while another interrupt is being processed to prevent a lost interrupt.
- ④ A trap is a software-generated interrupt caused either by an error or a user request.

⑤ Interrupt Handling

- i) The O.S preserves the state of the CPU by storing registers and the program counter.
- ii) → Determines which type of interrupt has occurred.
 - i) Polling :- in polling CPU waste lots of CPU cycles by repeatedly checking the command ready bit of every command.
 - ii) vectored interrupt system :- an alternative to a polled interrupt, which requires that the interrupt handler to send a signal to each device in turn in order to find out which one send the interrupt requests.

interrupt can generate at any instant of the time, CPU keeps polling the devices at the regular interval.

⑦ separate segment of code determine what action should be taken for each type of interrupt.

★ I/O structure

⑧ After I/O starts, control return to user program only upon I/O completion.

- wait instruction idles the CPU until the next interrupt.
- wait loop (contention for memory access)
- At most one I/O request is outstanding at a time, no simultaneous I/O processing.

⑨ After I/O starts, control returns to user program without waiting for I/O completion.

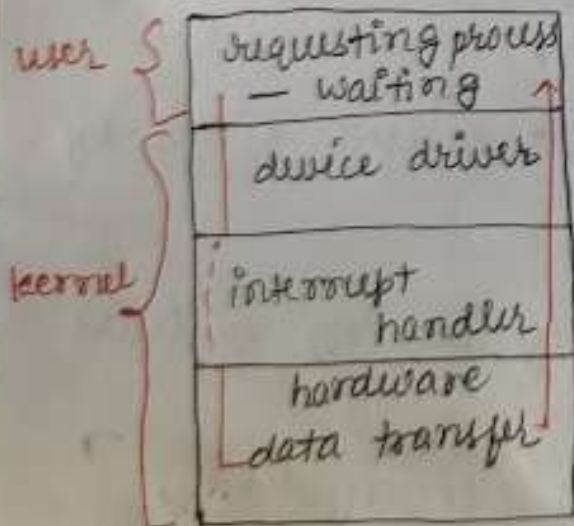
→ system call :- request to the O.S to allow user to wait for I/O completion.

→ device - status table contains entry for each I/O device indicating its type, address and state.

→ O.S indexes I/O devices table to determine device status and to modify table entry to include interrupt.

Two I/O methods

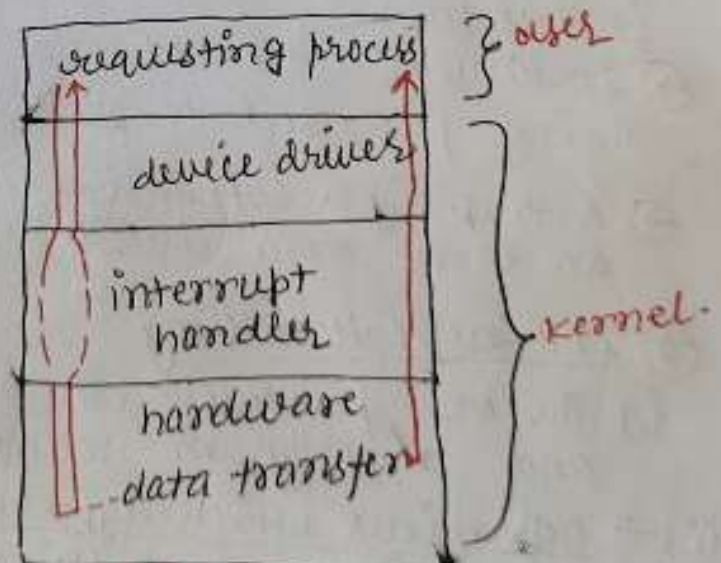
synchronous



→ fast transmission

→

Asynchronous



slower transmission due to extra bit gap.

Direct memory access structure

- ⊕ used for high-speed I/O devices able to transmit info. at close to memory speeds.
- ⊕ Device controller transfers blocks of data from buffer storage directly to main memory without CPU intervention.
- ⊕ only on interrupt is generated per block, rather than the one interrupt per byte.

⊕ Storage structure

- ⊕ main memory :- only large storage media that the CPU can access directly.
- ⊕ secondary memory :- extension of main memory that provides large non volatile storage capacity.
- ⊕ magnetic disk :- Rigid metal or glass platters covered with magnetic recording material
 - Disk surface is logically divided into tracks, which are subdivided into sectors.
 - The disk controller determines the logical interaction b/w the device and the computer.

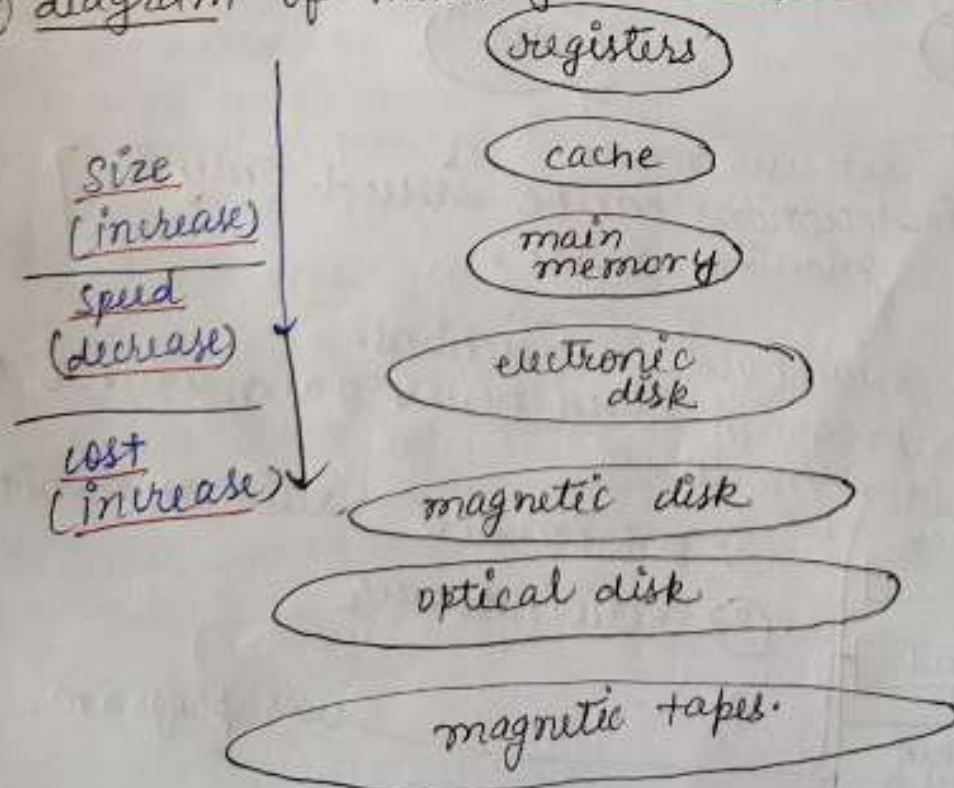
⊕ Storage Hierarchy

- ⊕ Storage system organized in hierarchy.

→ speed
→ cost
→ volatility

- ⊕ Caching :- Coping into faster storage system main memory can be viewed as a last cache for secondary storage.

⊕ Diagram of memory or storage-device Hierarchy



★ Caching

- use of high speed memory to hold recently - accessed data.
- Requires a cache management policy.
- Caching introduces another level in storage hierarchy. This requires data that is simultaneously stored in more than one level to be consistent.

④ Hardware Protection

→ i) Dual-mode operation

- I/O protection
- memory protection
- CPU protection.

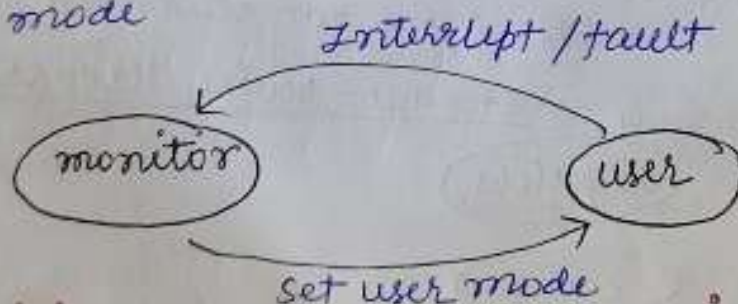
★★ Dual-mode operation.

- Sharing system resources requires O.S to ensure that a incorrect program cannot cause other programs to execute incorrectly.
- Provide hardware support to differentiate b/w at least two modes of operation.

1. User mode :- execution done on behalf of user.
2. monitor mode :- (also kernel mode or system mode) - execution done on behalf of O.S.

★ Dual-mode operation (cont.)

- ⑤ mode bit added to computer hardware to indicate the current mode :- monitor (0) or user (1)
- ⑥ when an interrupt or fault occurs hardware switches to monitor mode



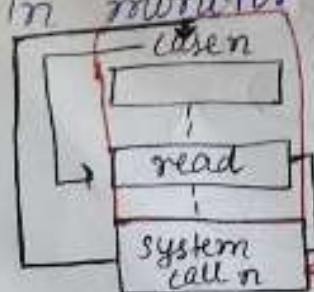
[Privileged instructions can be issued only in monitor mode.]

★ I/O Protection

- All I/O instructions are privileged instructions.
- must ensure that user program could never gain control of the computer in monitor mode.

④ Memory

- ① trap to monitor

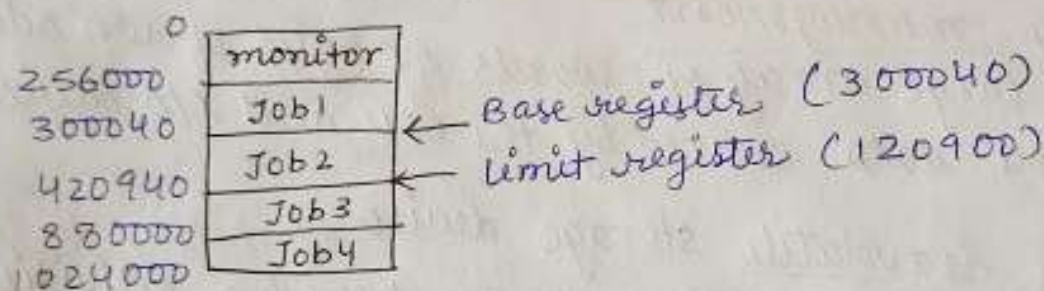


- ② perform I/O resident monitor
- ③ return to user

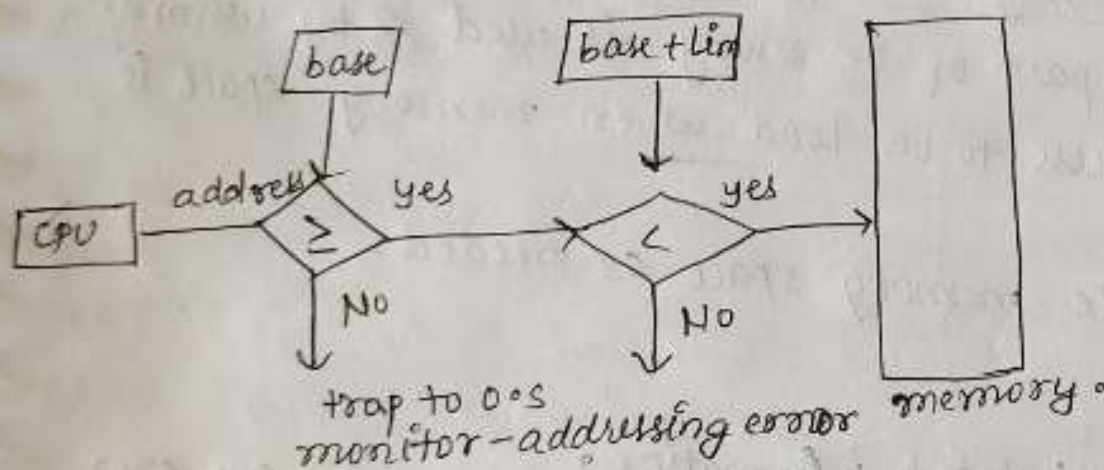
★ memory Protection

- must provide memory protection at least for the interrupt vector and the interrupt services routines.
- In order to have memory protection, add two registers that determines the range of legal addresses a program may access
 - Base register:- holds the smallest legal physical memory address.
 - limit register:- contains the size of the range.
- memory outside the defined range is protected.

use of a Base & limit register



★ Hardware address Protection



★ Hardware protection.

- when executing in monitor mode, the O.S. has unrestricted access to both monitor and user's memory.
- The load instructions for the base and limit registers are privileged instructions.

★ CPU protection

- Timer - interrupt computer after specified period to ensure O.S. maintain control.
 - timer is decremented every clock tick.
 - when timer reaches the value 0, an interrupt occurs.
- timer commonly used to implement time sharing.
- timer also used to complete the current time.
- load-timer is a privileged instruction.

★ Process management

- A process is an program in execution.
- process need resources, including CPU time, memory, files & I/O devices, to accomplish its task.

★ Operating system's work in process management.

- i) Process creation and deletion.
- ii) process suspension and resumption.
- iii) Provision of mechanisms for
 - a) process synchronization.
 - b) process communication.

★ main-memory management

- memory is a large array of words or bytes of own address.
- quickly accessible data shared by the CPU and I/O devices.

④ main memory is a volatile storage device.

⑤ operating system is responsible for the following activities in connections with memory management.

- i) Keep track which part of memory is used & by whom.
- ii) Decide which process to be load when memory space is available in.
- iii) Allocate & deallocate memory space as needed.

★ file management

- ⑤ collection of relative/related information.
- ⑤ commonly, file represent program (both source & object file) and data.

⑤ operating system's work.

- file creation and deletion.
- Directory creation and deletion.
- Support of primitive for manipulating files and directories
- mapping file onto secondary storage.
- file backup on stable (non-volatile) storage media.

* I/O system management

→ It consists of :-

- a) A buffer-caching system
- b) A general device interface. b/w device-driver.
- c) Drivers for specific hardware devices.

⊛ Secondary-storage management

→ Non-volatile

→ Large size

→ modern computer system use disks.

⊕ work of ~~operation~~ operating system :-

- i) Free space management.
- ii) Storage allocation
- iii) Disk scheduling

⊕ Networking (Distributed system)

- i) collection of processors that do not share memory or clock. each processor has its own memory.
- ii) Processors are connected through a communication network.
- iii) communication takes place using a protocol.
- iv) provides access to various system resources.

↓
benefits of shared resources.

- a) Computational speed up
- b) Increased data availability
- c) Enhanced reliability

⊛ Protection System

→ Protection refers to a mechanism for controlling access by programs, processors, or users to both system and user resources.

works of it :-

- a) distinguish b/w authorized and unauthorized usage
- b) specify the controls to be imposed.
- c) provide the means of enforcement.

★ Command - Interpreter System

→ many command give to the OS by control statements which deal with:

- a) process creation and management.
- b) I/O handling
- c) secondary - storage management.
- d) main - memory management.
- e) file - system access.
- f) protection
- g) networking.

⑦ Operating System Services

- Program execution :- system capability to load a program into memory and to run it.
- I/O operations :- provide ways for I/O.
- file system manipulation :- programs capability to create, read, write and delete files.
- communications :- exchange of info. b/w processors on same computer or on different system tied together by network. Implemented via shared memory or message passing.
- Error detection :- ensure correct computing by detection of errors in memory hardware, in I/O system etc.

★ Addition works of OS to ensure efficient system operations.

- a) Resource allocations :- allocating resources to multiple users or multiple jobs running at the same time.
- b) Accounting :- keep records of which user used how much and what resources. (for usage statistics).
- c) Protection :- ensuring that all access to system resources is controlled.

★ System calls

⑦ Provides the interface b/w running program & the operating system.

→ available as assembly - language instructions.

EX C, C++ to replace assembly - language instr.

★ Three general methods for system calls.

- Pass parameters in register.
- Store the parameter in a table in memory, and the table address is passed as a parameter in registers.
- Push (store) the parameters onto the stack by the programs, and pop off the stack by OS.

★ Types of system calls

- Process control.
- file management.
- Device management.
- Information maintenance.
- communication.

★ System Programs

→ it provides convenient environment for program development and execution

divided into :-

- file manipulation
- status information.
- file modification.
- Programming language support.
- Program loading and execution.
- communication.
- Application programs.

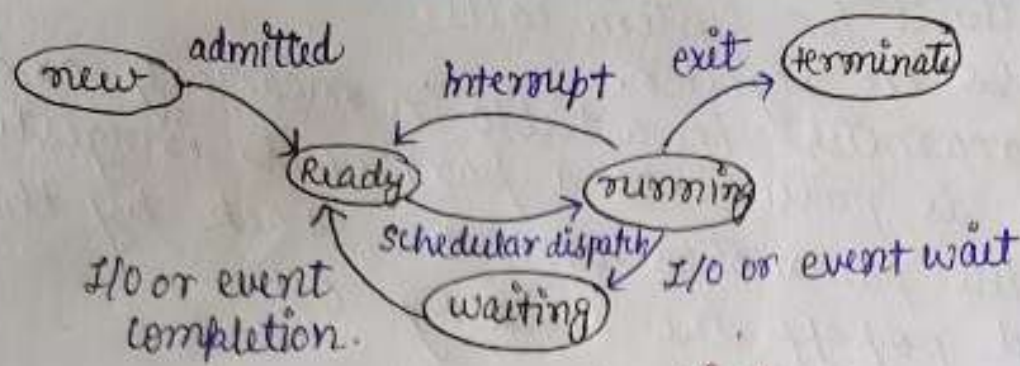
Chapter 4 : Processes

★ Process concept.

Process :- a program in execution
• sequential execution

Process status

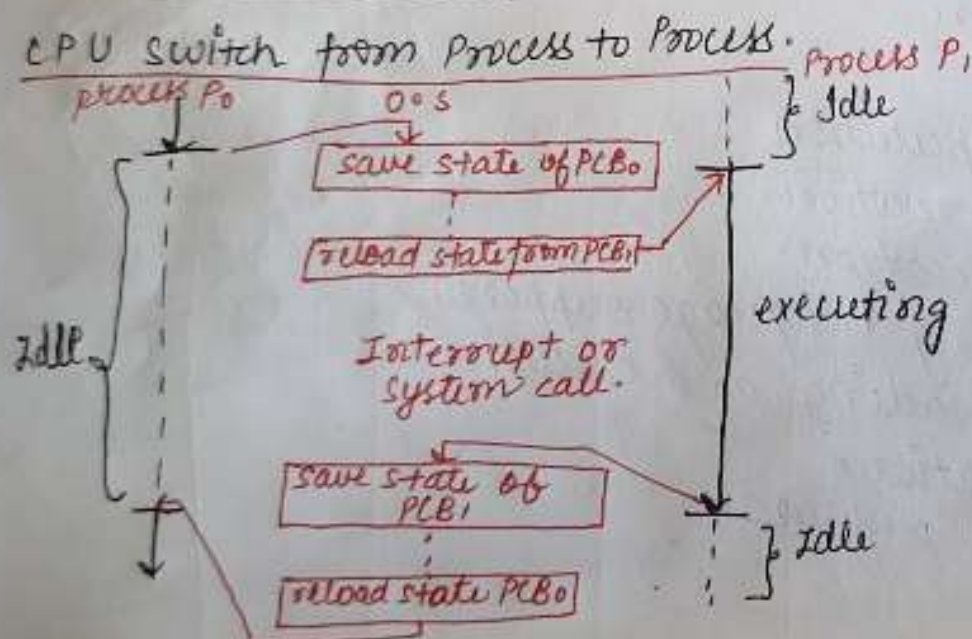
- New :- The process is being created.
- running :- Instruction is being executed.
- waiting :- The process is waiting for event to occur.
- Ready :- The process is waiting to be assigned to a process.
- terminated :- The process has finished execution.



Process state diagram

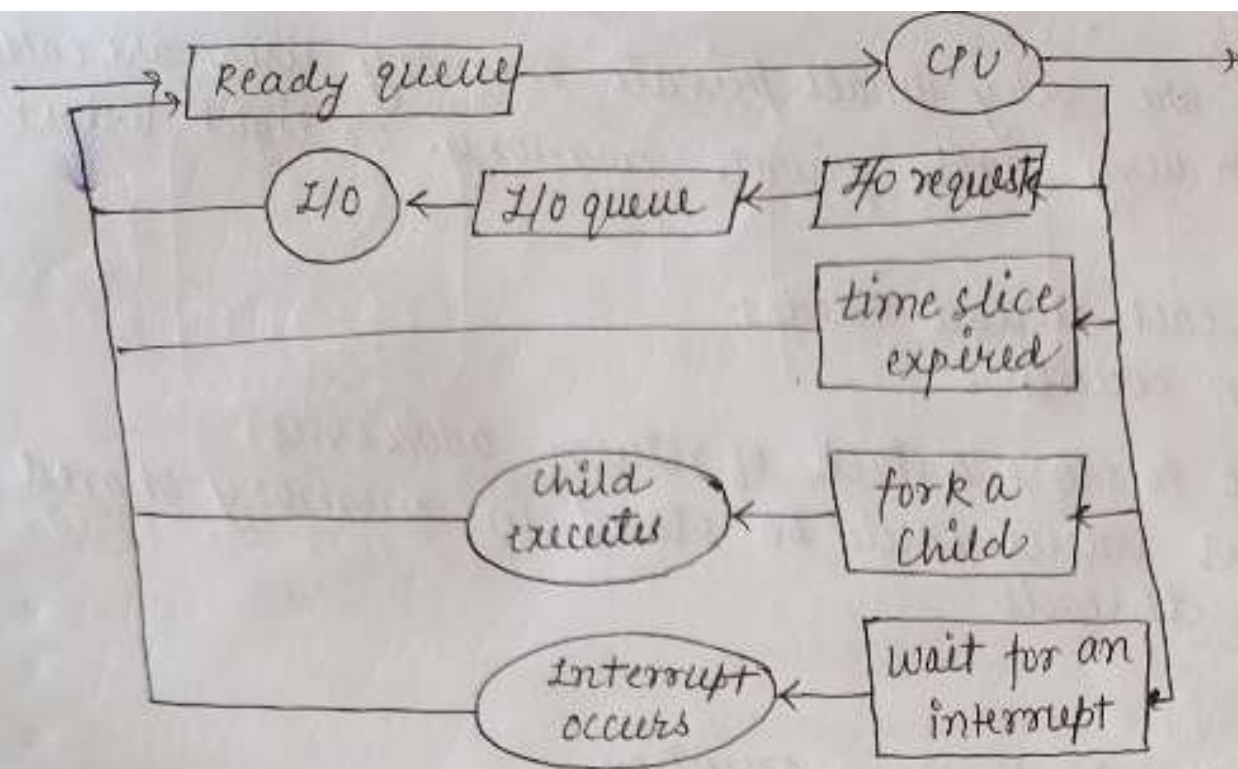
★ Process control Block (PCB)

Pointer	process state
process number	
program counter	
registers	
memory limits	
list of open files	



⊛ Process scheduling Queues

- Job queue :- Set of all the process in the system.
 - Ready queue :- Set of all processes residing in main memory, Ready, waiting to execute.
 - Device queue :- Set of process waiting for I/O devices.
- processes migrate from one queue to another.

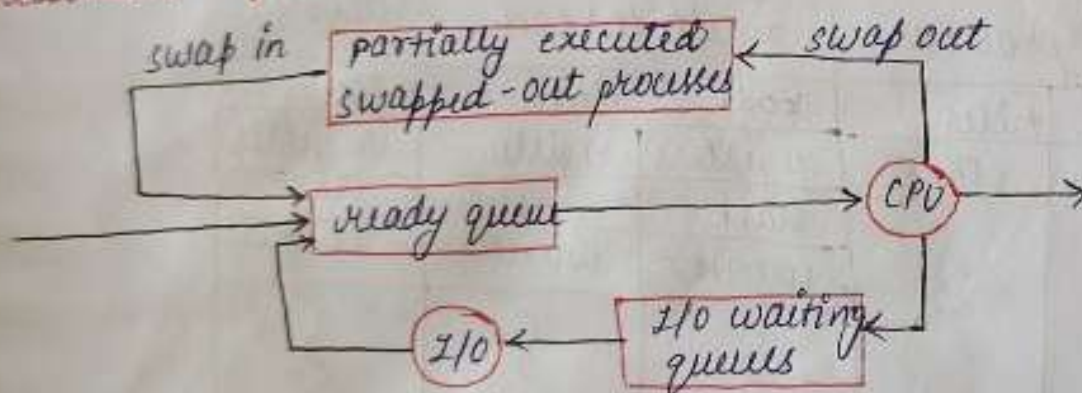


★ Schedulers

★ long-term scheduler :- selects which process should be brought to Ready queue.

★ short-term scheduler :- selects which process should be executed next and allocates CPU.

Addition of medium term scheduling



schedulers (cont.)

→ short-term scheduler is invoked very frequently (ms) (fast)

→ long-term scheduler is invoked very infrequently (seconds, min.) (may be slow)

→ long term scheduler controls the degree of multiprogramming.

Processes can be described as either

a) I/O bound process :- spends more time doing I/O than computations, many short CPU bursts.

b) CPU-bound process :- spend more time on computation, few long CPU burst

② context-switch :- when process switch it store its current state.
→ time dependent hardware.

Chapter 5 :- Threads

→ Processes which are very small private memory user are called thread or with very little private memory. or light weight processes.

→ At a minimum each thread needs:-

→ a) program counter.

b) a place to store a stack of return addresses.

c) all other values could be stored in memory shared by all threads.

⊛ Threads

→ alternative model of program execution

→ A process creates threads through a system call.

→ threads operate within process contents.

Process split into two parts

a) Resource state remains with process.

b) CPU state is associated with thread.

Switching b/w threads is less overhead than processes.

single thread processes			multithread processes.		
code	data	files	code	data	files
registers	counters	stack	registers	register	register
			stack	stack	stack
			counter	counter	counter

Process

a) It has heavy weight or resource intensive

b) process switching need interaction b/w OS & process

c) has its own memory & file resources

d) one process blocked then no other execution

e) use more resources

~~Process~~ Threads

a) It is light weight taking lesser resources than a process.

b) no need interaction with OS.

c) share same set of open files child processes.

d) if one blocked then other execute.

e) use less resources.

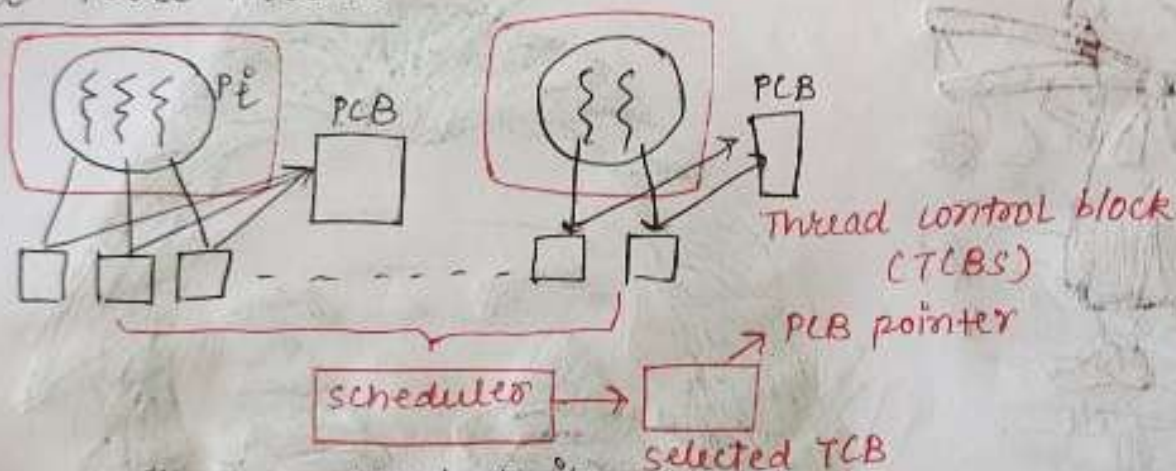
Benefits

- i) more efficient communication
- ii) simplification of design
- iii) Responsiveness
- iv) Resource sharing
- v) Economy
- vi) utilization of MP architecture

② level of threads

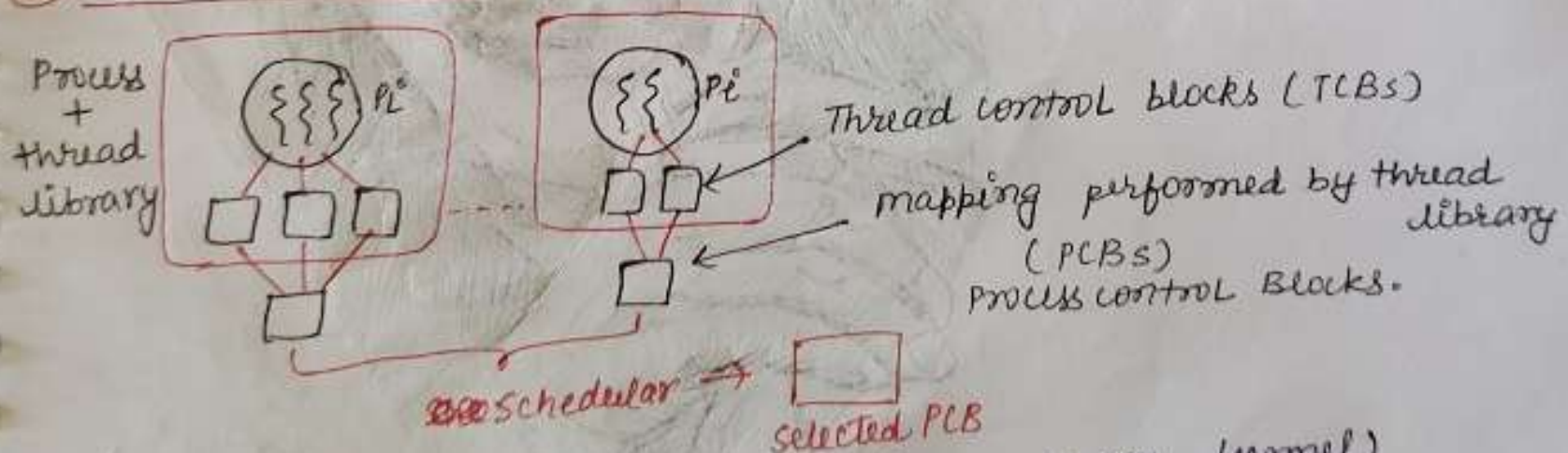
- i) kernel-level Threads:- managed by kernel.
- ii) user-level thread:- managed by thread library.
- iii) Hybrid Threads:- combine of both kernel & user-level threads

★ kernel-level thread



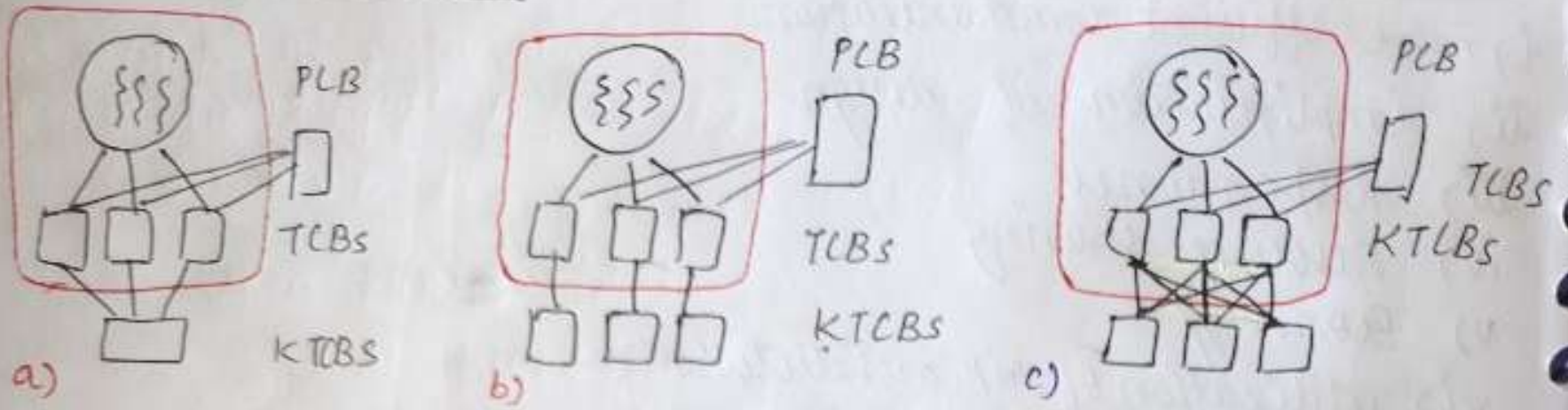
they are like processes but it has smaller amount of state information.

★ user-level threads



- fast thread switch b/w threads (because no kernel)
- one block all threads block
- No concurrency or parallelism.

Hybrid-Thread models



can provide a combination of parallelism and low overhead.

