

Computer.

alibas
class 7

A computer is a general-purpose machine capable of executing programs.

A program is a sequence of instructions.

→ Instruction cycle (not changed)

Fetch, decode, execute

only universal thing in Computer Science

Calculator

Not a computer

Not fetch, decode, execute

perform operations but does not execute instruction.

→ A single program of computer can simulate a calculator.

SmartPhone

It is upto you, whether you see it a general purpose device or not.

CS
↓
Computer Science

is the study of computer, computing processes.

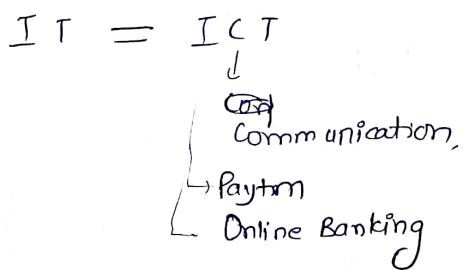
CE

↓

Computer Engineering.

There is no sharp boundary between these two.

Engineering → Application of real world problems.



* Is it important for computer to be electronic device?

Need not to be an electronic device.

Computer = Hardware + Software

(Non-physical)

(Non-tangible components)

→ programs

→ Source code + files,

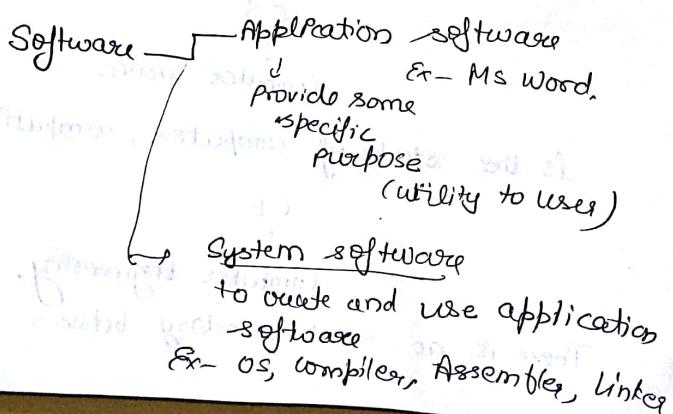
+ device drivers, user

User-manual, is meant to be read by the end user.

Documentation → Read by developers if developing diff. parts of computer modules

Software is a collective noun,

(Softwares) X



Runtime Environment (part of system software)
does not provide any specific functionality

Operating System

An operating system is a piece of system software that manages all resources of the computer, hardware or software, and provides a platform for easy and efficient execution of programs.

Software resource - files, database

Why we use?
OS provides the easiest way for interacting with hardware

→ How it is invented?
How it comes?

Charles Babbage → Analytical Engine

First design of computer,
(Never completed)

Ada Lovelace

She developed Bernoulli
Read about Analytical engine,

first working computer developed in mid-1940s

Konrad Zuse → K4. (computer) [fetch, decode, ...]

No OS → ENIAC

but there were

some control programs

OS evolved from control programs

University of Manchester

late-1950s

developed ATLAS

Computer Hardware

Major part \rightarrow CPU

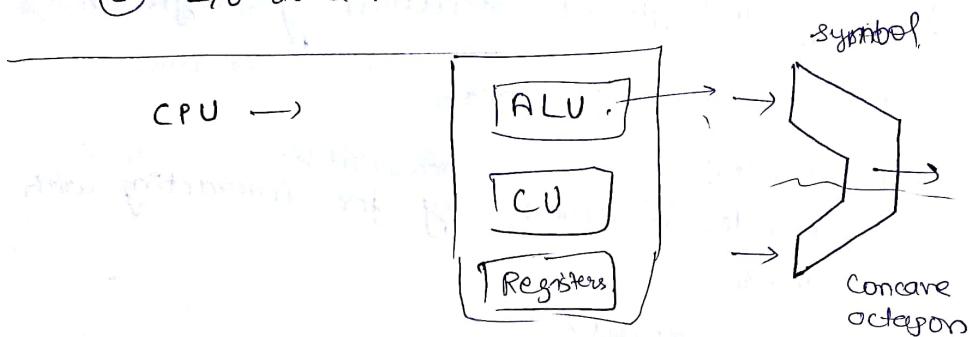
\hookrightarrow Fetch, decode, execute. is done by CPU

- To perform tasks

(1) Memory

store data, info & programs,

(2) I/O device :- interact with devices,



CU \rightarrow two ways.
hardwired. \rightarrow flip-flop, register
microprogrammed \rightarrow control programs,

CPU in Laptop \rightarrow Processor.

Microprocessor, used as CPU computers.

CPU become a single chip,
Faster, less energy, less space.

A computer that uses microprocessors as CPU
is called micro-computers.

Multi-core

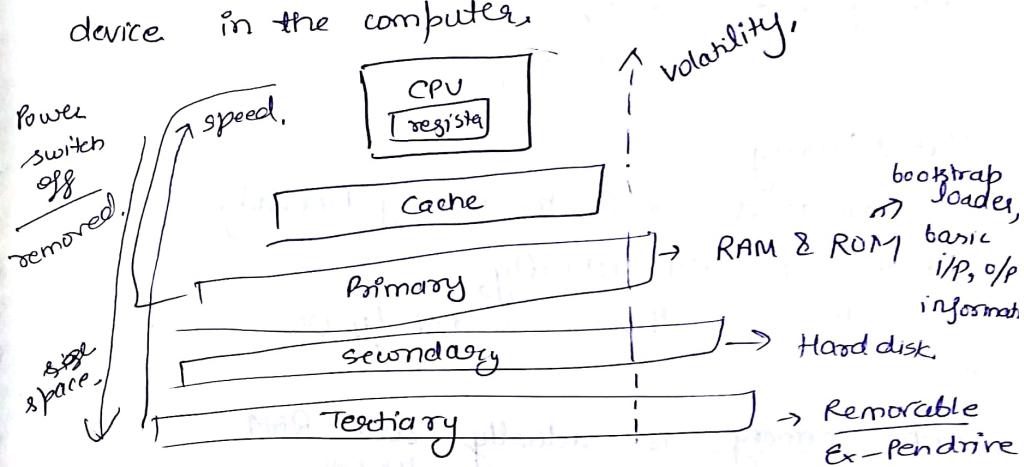
A processor with multiple sets of (ALUs, CUs)
on the same chip,



- Multiple core provides parallel-computing processing.
- Ensure concurrency

Memory hierarchy.

Arrangement of different types of memory device in the computer.



ROM

0th address

↓ start bootstrap loader

Random access can
also be performed
in ROM

RAM

Programs currently
executing
and data associated
with it

RAM access

Any location of RAM
can be accessed
at any time.

When computer switched off,

OS resides in Secondary storage (Hard disk)

When On, It execute bootstrap loader in ROM first address.

OS load into RAM.

till that time computer is On, major part of OS
resides into RAM.

Cost per byte,

Register > Cache > Tertiary

What is there at level i is a subset of level $i+1$.

Most persistent level,

Tertiary

Virtual Memory

You cannot see in the memory hierarchy

It does not exist actually,

It is an illusion created by OS.

Cache memory is actually static RAM,
flip flops.

dynamic RAM capacitors used.

(more capacity, more bits per unit area)

Cache is used for frequently requested data.

Cache is a small fast memory.

Cache is used for frequently requested data.

Cache is a small fast memory.

Cache is used for frequently requested data.

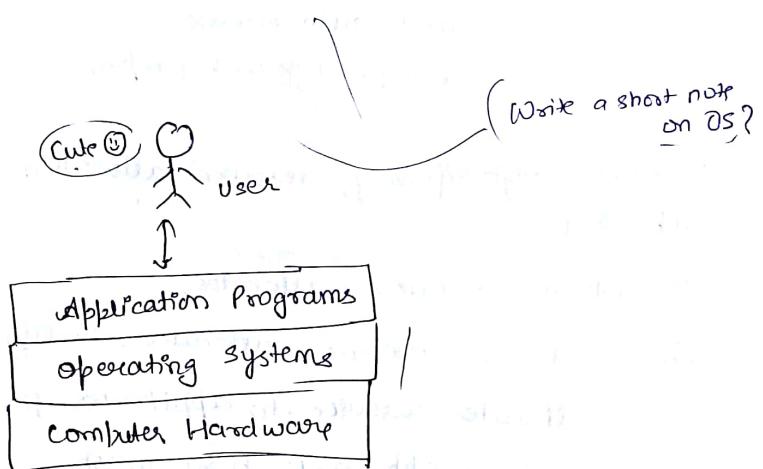
Cache is a small fast memory.

What OS do in your computer?
or your phones?

class-3
4/1/2018

- ↳ OS manages all the hardware
- ↳ OS facilitates all the execution of application programs
- ↳ OS act as an intermediate b/w user and computer hardware
- ↳ OS has to be designed to convenient to use and efficient.

Block diagram



- Does user directly interact with OS?
- ↳ Yes, but it is quite limited, (like shutdown, etc)

OS ensures that all the hardware and software are used properly.

OS Design Goals :-

what should be our goals?

- (1) Convenient (ease of use) → personal computers.
- (2) Efficiency (manage the resources in best possible way)
Multitasking is one way of efficiency → Super computers.
- (3) Energy conservation (for mobile specifically) → phones
- (4) Autonomous (it should have minimal user interface) → Embedded Systems.
Minimal user interface
(Ask in microprocessor class)
Convenience and autonomous are used for different purposes.

To achieve high efficiency, resource allocation is important.

OS acts as resource allocator.

- (a) OS acts as resource allocator. Justify,
OS allocates resource to application programs.
Ex - which app. gets how much resource,
for how much time is decided by OS.
Ex - Browser, Ms word, Paint allocate resources

- (b) OS acts as a control program.

OS evolved from control programs, send signals to h/w & s/w. It controls all h/w & s/w, when which h/w & s/w should work.

First OS :- ATLAS.

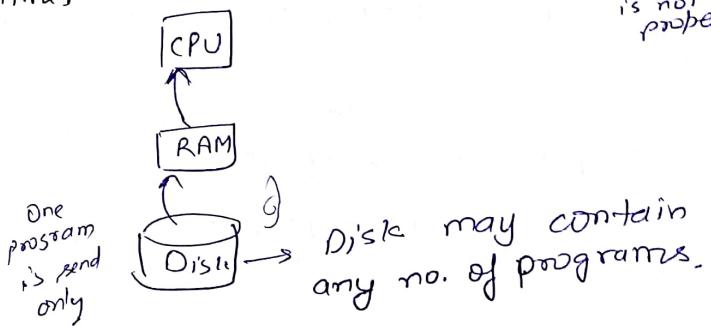
Our
{OS Name → XAAN}

Types of Operating System:- [eg. - MS DOS] (1983)

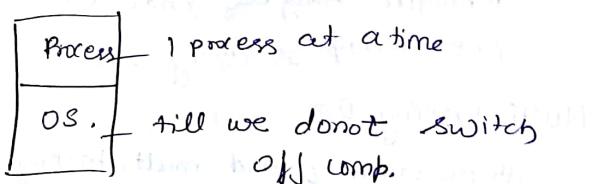
① Single process OS.

Process is a program execution.
→ Only one process is loaded into RAM at one time,

(resource utilization is not proper)



- They are inefficient.



② Batch Processing OS :- (eg. ATLAS)

Earliest OS developed (1957).

Inputs given using punch cards, primitive input devices.

Many single process clubbed together to form a batch, and computer executes the programs

one after other in a batch and output provided in a batch using primitive printer.

→ Program in a batch were of same type (all arithmetic)

→ At one time, one program was executed.

→ High inefficient and difficult to use.

③ Multiprogramming OS.

With the advancement of technology, micro computers were made.)

Load a few processes in memory simultaneously

- Initially one process is executed, when that process halts, another process is executed in mean time,
- CPU not sitting idle
- Computer became efficient
- CPU always had one program to execute



Ex - THE

→ Technique ~~Hanging Snares~~ in 1960
first multiprogramming OS.

④ Multi-tasking OS.

Multiprogramming and multi-tasking are different.

Multi-tasking is enhancement of multiprogramming

→ Logical execution of multiprogramming OS.

→ Multiple processes load in memory

→ CPU executes little bit of one process, then little bit of another process.

and keep executing processes in some order

switching b/w different processes. (in microseconds)

is very fast, then which make illusion that

every program runs concurrently, simultaneously

To have multitasking OS, we need to have hands on computer,

Kernel

is a part of OS, that interacts with the hardware directly. Not whole OS is interact with h/w. Most imp. activity of OS performed by kernel.

✓
where user can directly interact with computer
(like our computer with keyboard/mouse)

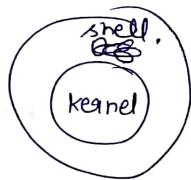
(In early 1960s., computers are not touched, there were punch cards)
so these computers can not have multitasking.

- Kernel is saved in RAM.
- From disk till the time comp. is shut down
- stays in RAM, till comp. is switched on
- Most important part of OS
- Kernel is small; so that it occupies less space in RAM, so that RAM contains more programs.
- If kernel is large, system become less secure
(have more keys due to large size)
- Most imp. activities are implemented within the kernel and less imp. activities outside kernel.
(here kernel can be small)
- When a kernel is very small so that it is minimal in size and provide base min facilities are called Microkernel.

Very small in size, but only 1 microkernel in any system.

Shell = Command Line Interpreter,
(MS DOS had it),

It is that part of computer which directly
interacts with user. User directly provides
commands to shell, then shell gets command
executed using other programs of OS.



Books :-

Operating System (Concepts / Principles).

By Silberschatz, Galvin, Gagne.



Dinosaur

Kernel

8/11/2018
Class-4

- * Kernel resides in RAM for entire duration.
- * Other utilities are implemented as separate modules outside the kernel.

Your program cannot interact with hardware; Kernel can. How?

Dual mode operation of computer system.

(1) User mode

(Everything outside, kernel works, shell of OS works in user mode).

(2) ~~Kernel~~ supervisor / system / privileged mode,

Anything that runs in user mode, does not have the permission to interact with hardware.

* Kernel works in kernel mode.

* Any part of code executing in kernel mode can interact with the hardware.

→ How will computer know in which mode it is working?

mode bit \Rightarrow stores bit that indicates in which mode we are working.

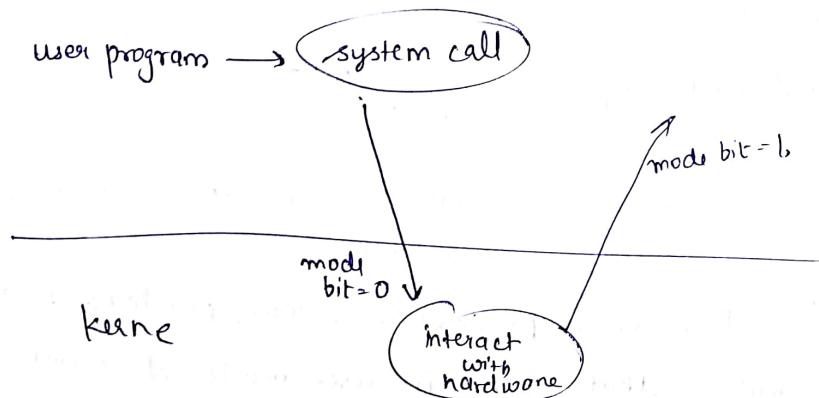
Values 0 or 1.

0 \Rightarrow kernel mode

1 \Rightarrow user mode

System call \rightarrow is a request, always invokes from user mode to execute some part of code in kernel mode.

- requesting kernel to do some work which program is not permitted to do by itself.
- If kernel has to do something, you need to switch to the kernel mode.

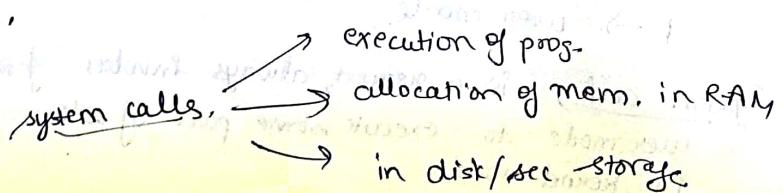


- Some instructions in microprocessor can be executed only in kernel mode
Privileged instruction if to be executed, you have to enter kernel mode
- Why can't you execute whole user program in kernel mode?

→ Massive security problem

↳ hardware can be damaged either genuine mistake or malicious

tasks to kernel have to be invoked by system calls.



System calls are implemented as C functions

system calls

↳ C functions

fork, getpid, getppid, wait, exec

related to execution of program.

What is dual mode operation?

2 or 3 marks
What is system call.

(PROCESS
(CREATION)
later)

Duties of an operating system

① Process Management.

Process is any program that is being executed,
only then process.

↳ includes every activity related to process

↳ starts with creation of process

↳ creation of user process

↳ creation of system process,
(generated by OS itself)

deleting of process
also a duty of OS

Suspension of process, } also a duty of OS.

Resumption of process }

* Interprocess communication

duty of OS.

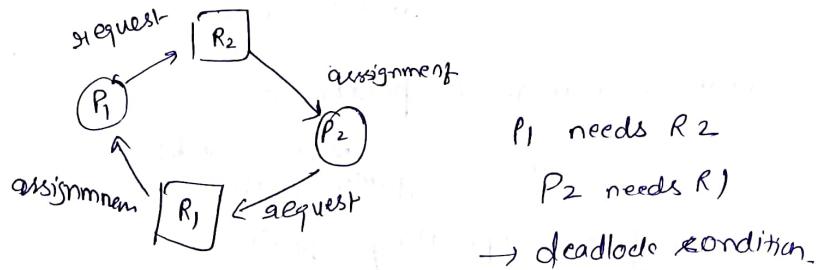
* Synchronise of process.

One program may have wait to other program.

~~OS Help~~:
Duty of OS help in process synchronization,

* In OS strange things happen.

Multiple process try to jump on one resource



P_1 needs R_2

P_2 needs R_1

→ deadlock condition

It is again duty of OS to

handle such situation.

(all these comprises process Management)

till midsem

we are discussing only process Management)

3 or 4 chapters

* Concept of process

↳ Threads

↳ Scheduling Process

↳ process synchronisation

↳ Deadlock

Q2

Memory Management

RAM only.

OS & one or more processes in RAM, which

computer on,

* Program & related data

* free space

It is duty of OS to keep track of free spaces in RAM.

If OS doesn't know where space, it cannot load programs.

↳ OS keeps track of which part of memory containing which program.

Memory management → 2 chapters.
After midsem.

#(3) Storage Management

Secondary memory.

or
hard disk.

↳ It is the job of OS to manage hard disk.

↳ logically stored files

file system → logical concept.

↳ It is the duty of OS to keep track of ~~file~~ file system.

blocks }
files → sectors } rotate at high speed
→ cylinders } schedule access to
→ plates/disk } diff. parts of disks,

Storage Management → file system mgt.

Storage Management → mass storage mgt.

Later study.

(4) Caching

(5) Input/Output management

It is the duty of OS to ensure that the user can use input/output devices properly.

Operating System Services :-

two parts

provide a large number of services.

① Helpful to the users

- (a) execution of our program
- (b) an user interface

CUI

Character
or
Command

(also known
as the
shell)

GUI

Graphical
(One of the biggest
success in
computers)

GUI was popularised in 1980s in USA and Europe.

By 1990s to India,

② C I/O operation.

- (a) OS helps in file system manipulation.
- (b) helps in interprocess communication
- (c) Error detection (bad alloc ~~free~~)

③ helpful to the system

do not help the user directly but
help the system to work properly,

(a) Resource allocation,

(b) facilitates protection and security.

(Tries to ensure diff. resource of computer are
being used properly.)

Operating System structures. -

Q how the whole source code of OS is organised in OS.

Windows code is billions of lines.

OS are large programs.

Q How to manage different modules.

Q How different modules are arranged.

(1) Monolithic (one stone),

(Entire source code is jumbled into one module)

Any part of OS can interact with any part of OS or hardware,

(Everything at same level.

Ex - MS DOS, Unix, Linux



① Firstly it is

insecure

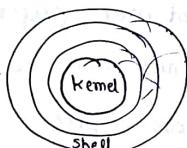
② If we would have to change some code, it would be difficult.

(2) Layered architecture

E.g. Tux

Any program at i^{th} layer

can access $i-1^{\text{th}}$ layer



Only kernel interacts with hardware.

* Improves software engineering aspect of OS

* Improves security aspect of OS.

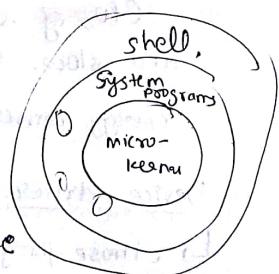
(3) Micro-kernel based architecture:

(layered microkernel) Here kernel becomes minimal.
Hence it is highly secure

Outside - shell

System programs implemented as module

Eg - Mach, MINIX, (Education purpose OS)



Real Time Operating System (RTOS)

9/1/2018
Class-5, 6

- ↳ special purpose OS
 - for specialized applications, Ex- Robots, moon rovers,
- rigid time constraints
 - (if not follow whole system fail).
- are generally on execution of program.
- two types
 - ↳ Hard RTOS
 - ↳ Soft RTOS

Hard → It is absolutely necessary to follow time constraints.

Soft → it is expected that time constraint will be met.

If not met, performance will degrade.

Ex- Linux can be considered as Soft RTOS.

Q) SHORT note on these topics :

Operating Systems for smartphones -

When CPU for smartphones are designed, they are slow.

CPU's of smart phones are intentionally kept simple and slow. So that heat does not generate.

Energy conservation,

Device drivers (Software)

- ↳ those programs that allow us to interact with different I/O devices,

Q Is ~~the~~ device driver a part of OS?

- controversial
- Installation from outside
not always needed
 - once you install, it becomes
part of OS.

→ Read in detail in chapter I/O management.

Case studies :- on OS

① CP/M

- * Control program or monitor.
- * used in late 70s, nearly 80s.
- * made for 8 bit microprocessor
- * Gary Kildall,
- * sold it to ~~so much~~
- * write about architecture, R

Read it

② MS DOS

Disk OS

Read about Architecture,

③ UNIX

④ Linux

⑤ Windows

Turing award
most prestigious award in computer science
crypto, theoretical CS, AI.

ACM

Association of Computer Management.
gives award for lifetime achievement award.

3 people got this award

→ K L Thompson

→ FP Brooks Jr.

wrote a book

Mythical Man Months.

Unit of effort.

is a unit

like km, m

1 man month = working on a particular job for one month.

1 man year. = can be two men working in some project for 6 months,

→ Barbara Liskov,

(Professor in MIT)

* (Open Ended Project)

develop a very similar minimal OS.

PC → easier

Smartphones → difficult

Unit-2

Process and Process Scheduling.

A process is a program that is in execution,

In OS,

process & job

used synonymously,

process has several components.

→ Text Section (program)

in ML,

→ current activity represented by values
of the registers.

→ program stack,

• subroutines & function calls,
whenever there is fn call,
return address is saved,

• used to evaluate multioperator
expression,

→ data section

stores the global variables in this.

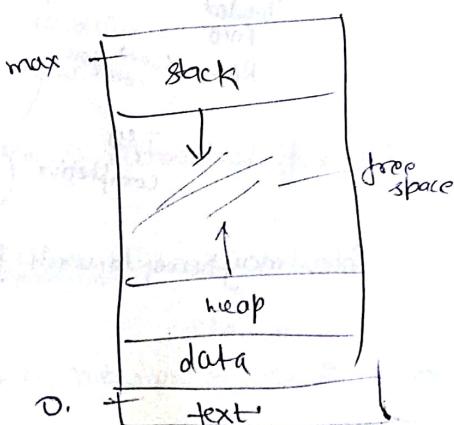
→ .heap

to store dynamically allocated variables.

Process is a bigger entity than a program.

memory allocated to
a program

logical address base
of a process,



→ Size of stack always varies.

→ Stack grows downward,

→ heap grows upward

They should not collide.

If collide, unpredictable things occur.

→ there are different techniques to check this,

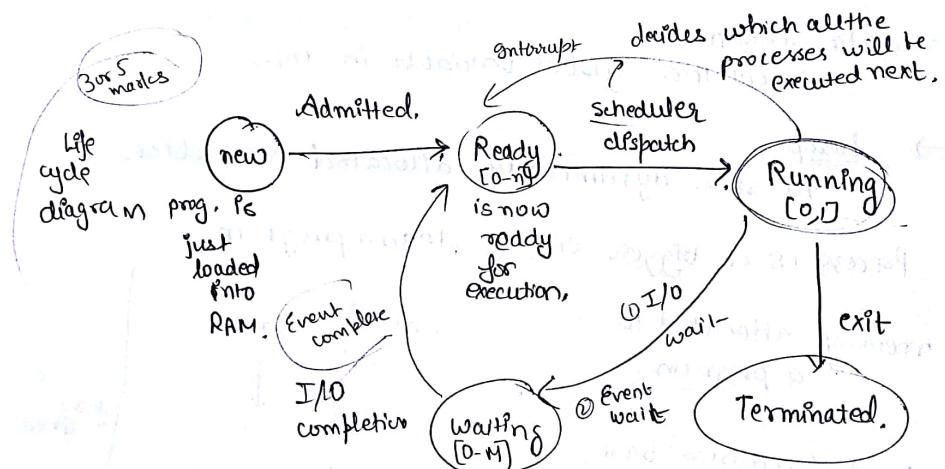
A program is a passive entity.

A process is an active entity.

Process state

→ tells what the process is actually doing right now.

→ State of a process changes quite a few time during its lifetime



→ Process may have to wait for an event to occur

Interrupts

- > the current instruction is always complete.
- in case of interrupt.
- > every process starts from new stage

OS uses

Process Control Block (PCB) ↴

- ↳ It is data structure
- ↳ stores all the imp facts.
- for each process OS has a PCB
- whenever process is in new state, its PCB is created.
- In RAM, there is PCB
- PCB is created by OS, OS has its own area & some free area.

what all stored in PCB

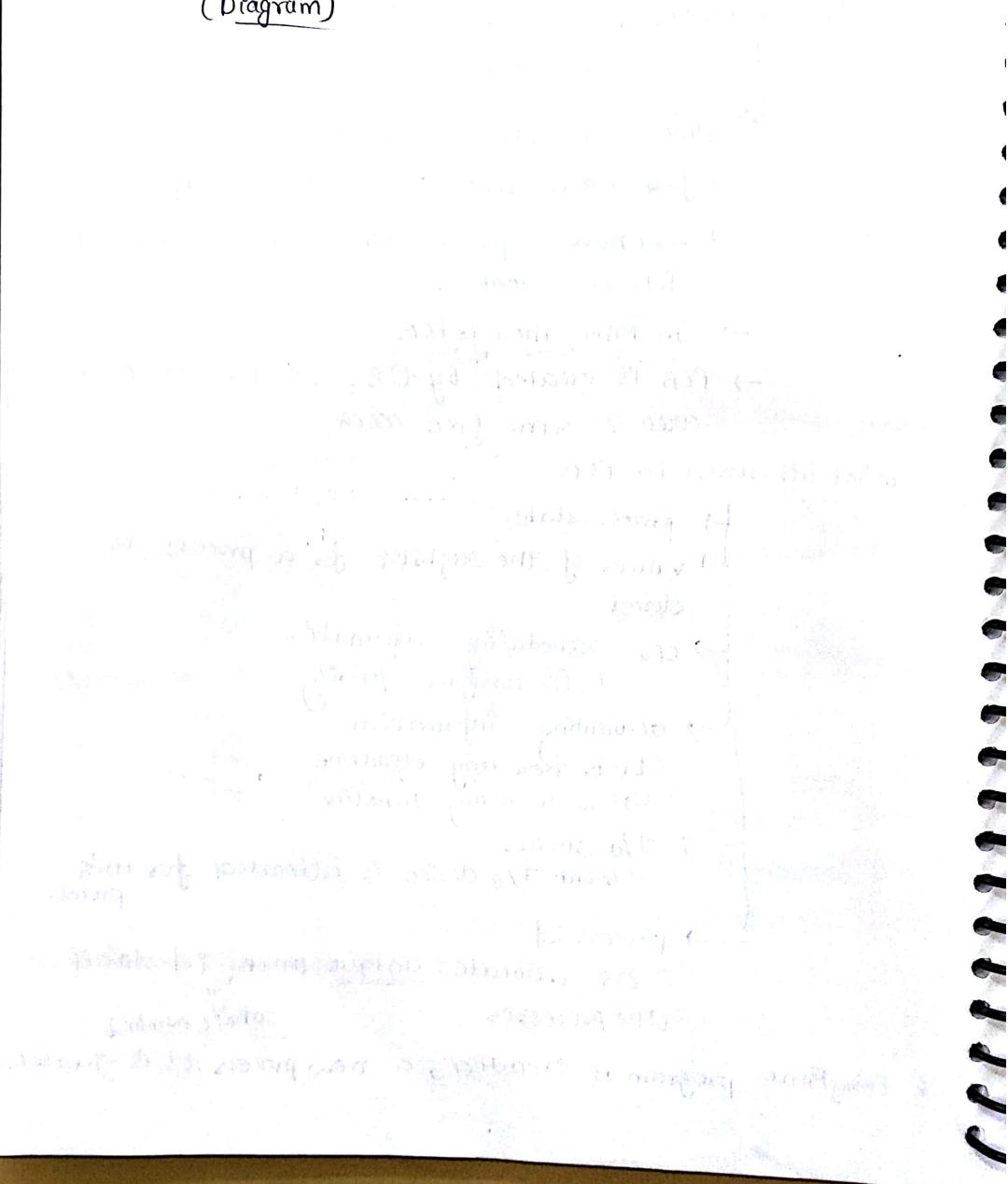
- ↳ process state
- ↳ values of the register for a process is stored
- ↳ CPU scheduling information
↳ OS assigns priority
- ↳ accounting information,
↳ is there any deadline
↳ is there any timeline
- ↳ I/O status
 which I/O device is allocated for this process.
- ↳ process id
 OS allocates unique process id to all the processes
 whole number

& Everytime program is executed, a new process id is generated.

Same program can create multiple process
and can create different ids for the
same program.

Process id of kernel is zero (in Linux) type

Page - 6
(Diagram)



Process Scheduling.

OS is resource allocated,

Most important resource : CPU.

Duty : To use that resource to max,

done by OS
↳ selects one process, to be next executed,

Multiprogramming

↳ objective
to use maximum use of CPU.

Objectives of multtasking :-

Switch CPU b/w processes.

→ To implement process scheduling, we have to use data structure,

Scheduling queues.
↳ itself is a ~~data~~ structure,
stored PCB as elements
(PCBs of processes).

Types of scheduling queues.

↳ Job queue

(there is only one job queue in system)
from birth to death, (process)

↳ ready queue

(also one ready queue)
store PCBs of all such processes which are in ready state
(loaded into RAM).

All these queues are implemented as
linked lists.

Q) short note on scheduling queues.

You should draw diagram as we're in the notes.

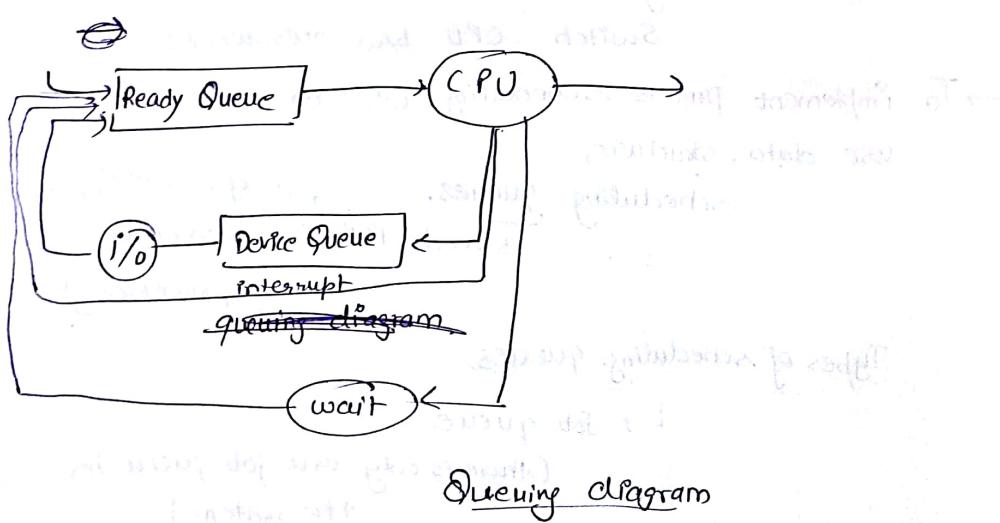
→ Device Queue

two or more processes may be use same I/O devices.

most I/O devices are not sharable,
so there is device queue for them.

Q) How many device queue?

Each device will have its own device queue.



→ Which will be given priority of I/O?

As in queue,

FIFO.

Priority Queue → Then queue will be manipulated.

→ Who will take this decision?

Scheduler.

→ is module of OS,

→ decides which process to be executed next.

3 types

① long-term scheduler ≡ job scheduler.

(decides which of the process to be loaded into memory at what time.)

→ controls the degree of multiprogramming
no. of processes in RAM:

Why long-term?

long term scheduler invoke infrequently.

(only after a long time)

(it use quite less)

Two types of
process

→ CPU-bound (more calculations)
use CPU for maximum time
for life time.

→ I/O-bound (Print, MS-Word, Music Player)

→ long term scheduler loads a proper mix of CPU-bound
and I/O-bound into the memory
ratio is maintained
duty of long-term scheduler.

② short-term scheduler

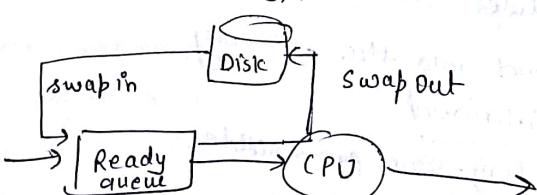
→ it will invoke quite frequently.

→ very fast.

- decides which of the process currently in the memory will be executed next (given ^{control of} CPU).
- they are softwares
- also known as CPU scheduler
activity known as CPU scheduling.
↳ has to be very fast.
- Most of the OS ex - Ubuntu, Mac, Windows used only short term scheduler
 As in computers, suppose to be used by one user
use will be automatically do long-term scheduler

③ Medium term scheduler

- ↳ is invoke less frequently than short term and more frequently than long term
- ↳ also known as Swapper,
- ↳ If long term realizes that program is using a lot of memory then medium term will take such program from RAM to disk.



∴ That process is assigned
 if x is executed half, swapped out,
 (then resumed from where it left once
 swapped in.)

Modern day OS also don't use Medium term.

\xrightarrow{MTS} improves CPU & I/O bound processes ratio.

Context switching when

① if one process is terminated.

② if current process is sent to wait stat.

what to be done next?

① when we have to switch, state of current process.

should be stored to PCB & state of next should be loaded/reloaded, ~~pure overhead~~.

→ Context switching is pure overhead:

↳ happens very frequently

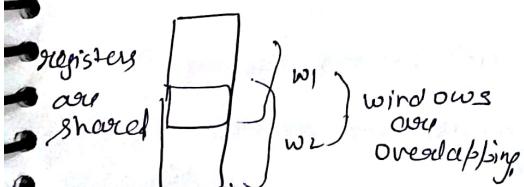
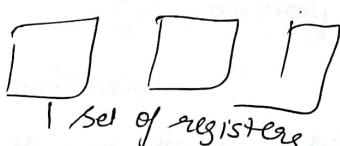
we should minimize the time of context switching.

Reduced Instruction Set Computer (RISC)

Used concept of overlapped register windows,

(Read about) IT

(from Morris Mano)
2-3 page.



CWP

Current Window Pointer

→ Overlapped will be necessary in case of interprocessors
Communication

↳ Assign each process one register window. Only
one window can be accessed at a time.

→ There are certain registers common to 2 window
& some are for a specific windows.

→ There is a register called CWP whose value
changes when window is to be changed.



Process Creation

15/11/18
class 8

Process

↳ user processes (doing something)
↳ system processes. (for user directly)

↳ are parts of the operating system.

→ Connecting program to a process is process creation.

.1 Who can create a process?

↳ One process can create another process.

parent of
those other processes

children processes

→ Kernel is the first process to be loaded into system.

→ Kernel is the ancestor of all other processes.

→ Kernel creates more processes using shell, shell creates more processes and so on.

→ If we draw a tree, root node is represented by kernel.

→ Once a process is created, it requires resource in

h/w, s/w, CPU, memory.
file database

who will give these resources?

→ A parent process can give resources to children.

→ A process can directly request resource from OS.

It depends on OS to determine:

- ① Both child & parent will execute concurrently and will compete with each other for control of CPU.
- ② If the parent process is waiting for children to return value, then child process executed first.

The child process can be a clone of parent process or it may have a new program loaded.

It has same code and same data set,

Which system call is used to create a new process?

Ans- fork, what does it return?

```
int main()
{
    int value;
    value = fork();
    if (value)
        printf("Parent");
    else
        printf("Child");
    return 0;
}
```

if it is used to create two processes.
→ fork will return to zero to the child process
→ return process ID of child it is creating to parent process.

This program is used in 2 processes, for parent process, non-zero value returned,

Output :- Parent or Child
Child. or Parent - wait

④ How will process know its own id?

```
int main()
{
    printf("%d", getpid());
    printf("%d", getpid());
    return 0;
}
```

child process may not be clone of parent.

```
int main()
{
    int value;
    value = fork();
    if (!value)
        exec("c:/calculator.exe");
    return 0;
}
```

Once a process has been executed, till the last instruction, it goes to terminated state.

Process termination

Deletion of PCB is needed from memory.

All resources allocated have to returned.

All files needed, needed to closed.

Sometimes, a parent process may also terminate a child process.

g) may have to kill
→ parent process may require to kill its own
child process kill)

Reasons ?

- it is taking too many time,
- Parent asked child to do some calculation, that calculation not required any more than parent. Kill child.
- A process whose parent process is terminated is called Orphan process.
- Some OS do not allow orphan process to execute. or some allow?
- Some OS require parent process to be alive for orphan to execute otherwise tree of process degenerate into forest.)
- If parent process is terminating, it has to kill its children process and can lead to grand children, being terminated.
This is cascading termination of processes.

Zombie process

A process which is terminated but its PCB is still in the system. because that process is supposed to return a value to its parent process.

Interprocess communication

↳ Two process may require to communicate with each other.

→ Information Sharing,

→ Communication speed up

→ same software tools as multiple processes
(↑ modularity)

Q ways → ① shared memory

(multiple processes loaded into RAM)

P1 can write in shared memory which can be read by P2.



② Message passing.

(No shared memory)

→ Processes can send messages to each other.

→ Process can send data using process format.

Q types of system calls (pseudo calls, abstract system call that has to be implemented)

① send (P, message);

 message to be send
 an integer

 (id of recipient process)

② receive (id, message);

 receive message, send to a particular process

GD of process
which has
sent the
process

Threads

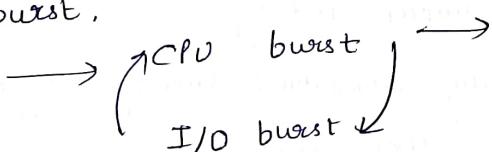
- > Earlier,
 - one process has sequence of instruction,
executed one after other,
- > Now,
 - process contains one or more thread.
- > Each thread is a sequence of instructions.
 - /
- > Multiple threads can be executed - parallel and scheduled independently, this program paradigm, called multi-threaded programming.
- > A thread is also known as light weight process.
- > Different threads of same process share the same resources, memory, files, resources.
- > A thread has its own state
- > thread can be synchronized
- > thread create more threads.
- > process have its own state of memory, resources
so, it is different from thread
- > Thread can wait for other threads,
- > Context switching become more fast,
 - what note

When a process starts exec, it goes to the CPU & performs calculations for some time.

16/11/2018

Class-9-10

- * When process starts perform CPU burst.
- * After that process, does some I/O for some time I/O burst.



Execution of process is a sequence of alternate CPU burst & I/O burst.

CPU Burst → Within a short duration several CPU processes are being executed.

When will CPU scheduler invoke?

Under four particular conditions,

- C1 : when a process goes from running → waiting.
CPU scheduler invoke.
- C2 : running $\xrightarrow{\text{low priority}}$ ready,
so that high priority execute.
then, CPU scheduler invokes.
- C3 : waiting \rightarrow ready,
let current process does CPU burst, then I/O burst -
CPU control is relinquished, so low priority is executed
- C4 : running \rightarrow terminated
running process is terminated.

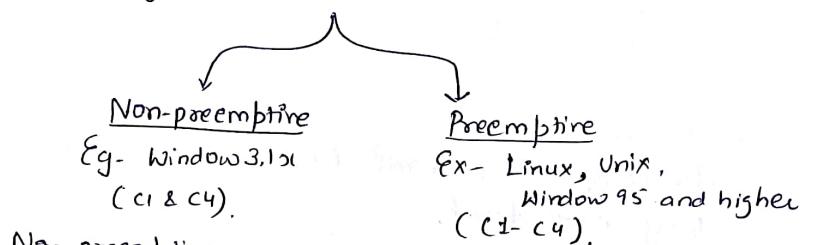
In case of C1 and C4, CPU scheduler compulsory
C2 and C3 are optional

Interrupt \Rightarrow C2

\rightarrow On what basis, these CPU scheduler will work?

that logic is implemented by :-

C CPU scheduling algorithms. (must be very fast,
(simple as possible))
takes logic that select one of the many
ready processes.



Non-preemptive
Once a process gets control of CPU, is allowed to execute as long as it wants.
i.e. till termination or need I/O

No other process of OS can evict it out of the processor (CPU)

- \rightarrow also known as cooperative scheduling
- \rightarrow it is simple to implement
- \rightarrow they are not efficient.

Preemptive

A running process can be evicted out of CPU \rightarrow priority based
 \rightarrow another logic based

→ requires extra hardware

In case of C2 and C3, CPU scheduler may or may not be invoked.

MS-DOS

↓
single process OS

so there is no concept of CPU scheduling

Dispatcher

↳ Model in OS called dispatcher

↳ part of software

↳ will give the control of the CPU to the chosen process ~~chosen~~ by CPU scheduler.
selected

steps used by dispatcher

① switch context.

② switching to user mode

(change the value of mode bits)

③ jump to the proper location in the user program

→ Time taken by the dispatcher to do all the functions is called Dispatcher Latency.

termination or suspension of process to start or resuming another process

→ Hardware dependant
→ pure Overhead

→ value is constant, must be minimum

→ no useful work can be done in it.

There are some parameters, according to which, we compare different CPU schedulers called

~~CPU Scheduling~~

Scheduling Criteria

① CPU Utilization (should be highest)

For what % of time, CPU is busy (80-90%)

② Throughput (as high as possible)

No. of processes completed/unit time.

③ Turnaround time (should be low)

Time taken from its submission to its termination,

↳ time spent in CPU + I/O + ready queue + devices

① and ② is calculated for the entire system.

③ can be calculated for individual processes.

Average turnaround time for all the processes

④ Waiting time (minimize)

Time spent by the process in the ready

queue (waiting only for CPU)
Only (not for I/O.)

Can be calculated for individual process

If for whole system, take average

⑤ Response time (minimum)

time spent ^{from} to the submission of the process to the arrival of first response.

↳ can be for individual

↳ take average for system.

min. variance

* Variance of response time of all processes
should be minimum.

Gantt Chart used ^{study} for CPU scheduling also.

↳ type of notation when "CPU scheduling"
↳ used in management.

(* CPU scheduling numerical
(some) in exam.)

CPU scheduling algorithms

① First-Come First Service (FCFS). Algorithm:

↳ Non preemptive algorithm.

↳ fair

↳ simple to implement.

↳ performance will not be good.

throughput & waiting time & turnaround time will be high.

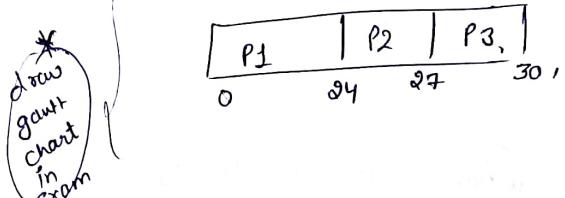
'Conroy effect' occurs

↳ stuck up in ready queue
for a long time.

Q _i	<u>Process</u>	Arrived time (ms)	(CPU burst-time (ms))
	P ₁	0	24
	P ₂	1	3
	P ₃	2	3

Ready queue : P₁, P₂, P₃.

Gantt chart:



Calculate different parameters:

① throughput

$$3/30 = 0.1 \text{ process/ms.}$$

Show calculations

② turn around time

$$P_1 = 24 \text{ ms.}$$

$$P_2 = 27 - 1 = 26 \text{ ms.}$$

$$P_3 = 30 - 2 = 28 \text{ ms.}$$

$$\text{Average} = \frac{24 + 26 + 28}{3} = 26 \text{ ms.}$$

③ waiting time

$$P_1: \text{when it get control} - \text{when it arriving} = 0 - 0 = 0 \text{ ms.}$$

$$P_2: 24 - 1 = 23 \text{ ms.}$$

$$P_3: 27 - 2 = 25 \text{ ms.}$$

$$\text{Average} = \frac{23 + 25}{3} = \frac{48}{3} = 16 \text{ ms.}$$

CPU utilization X

if we calculate, it will become 100%
as we did not consider dispatcher,

Response time X

it is difficult to tell when the first
response is coming?

Assumption = 3 ms.

CPU scheduling Algorithms.

class-11

18/11/18

Not Attended

⑧ Shortest Job First (SJF) algorithm.

Process with smallest ~~burst~~ next CPU burst is selected for execution. (in ready queue).

Two or more process with same next CPU burst, FCFS breaks Tie.

⇒ Shortest-next CPU burst scheduling

Optimal ~~burst~~ performance: avg. waiting/turnaround/
response time will be minimum.

→ Best possible CPU scheduling also.

How do you find next CPU burst?

~~Theoretical~~ practical

" cannot be implemented.

Impossible to predict next CPU burst without executing the program.

We have 24ms, 3ms type in yesterday, tues but that wasn't affecting the logic of scheduling.

→ we can try to predict next CPU burst time.

Concept of Exponential avg.

Predict next CPU burst time of particular process:

t_{n+1} = trying to predict time of $(n+1)$ CPU bursts.

'n' CPU burst have already occurred.

$$\text{Predicted time } t_{n+1} = \alpha t_n + (1-\alpha) \bar{t}_n$$

\swarrow actual time
 \nwarrow taken in n^{th} CPU burst.

$$0 \leq \alpha \leq 1$$

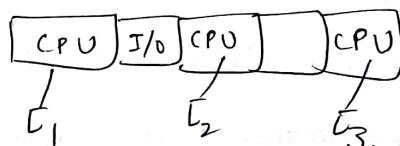
↓
constant base case

| For recursion

\bar{t}_0 = constant (use any number)

or
system avg

Predict this value by average of t of all system processes.



If $\alpha \rightarrow 0$, $t_{n+1} \approx t_n$, i.e. static situation.
Predicted time will be constant.

if $\alpha \rightarrow 1$, $t_{n+1} \approx t_n$, i.e. last CPU burst time

Recommended,

$$\alpha = 0.5$$

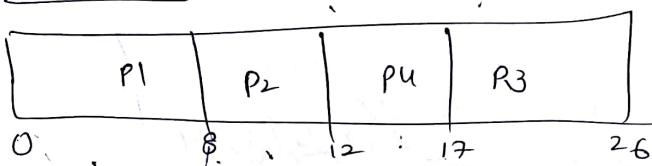
Q5 :- 4 processes.

	Arrival time(ms)	CPU burst time(ms)
P1	0	8
P2	1	4
P3	2	9
P4	3	5

Non-preemptive algorithm is used.

Calculate throughput, avg. turn around time, avg. waiting time.

Gantt chart



ready queue: ~~P1, P2, P3, P4~~

continues till 8 ms.

throughput : $\frac{4}{26} = 0.15 \text{ proc/ms}$

turn around

$$P1 = 8 - 0 = 8$$

$$P2 = 12 - 1 = 11$$

$$P3 = 26 - 2 = 24$$

$$P4 = 17 - 3 = 14$$

$$\text{Avg} = \frac{8 + 11 + 24 + 14}{4} = 14.25$$

waiting time

$$P1 = 0 - 0 = 0$$

$$P2 = 8 - 1 = 7$$

$$P3 = 17 - 2 = 15$$

$$P4 = 12 - 3 = 9$$

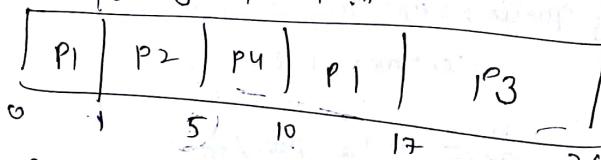
$$\text{avg. waiting time} = \frac{31}{4} = 7.75 \text{ ms.}$$

- Preemptive version of shorted FCA is available.
- If next CPU burst time of new process is less than current processes, then new process can preempt other processes.
- ⇒ also called shortest-remaining-time-first algorithm.

Q) Recalculate:

Ready queue: ~~P1, P2, P3, P4~~

{ only ± preemption}



P1 will be stopped and P2 starts as $7 > 4$

Throughput \propto same

≈ 0.15 processes/ms.

Turnaround

$$P_1 = 17 - 0 = 17$$

$$P_2 = 5 - 1 = 4$$

$$P_3 = 26 - 2 = 24$$

$$P_4 = 10 - 3 = 7$$

$$Avg = \frac{52}{4} = 13$$

short process

gaining q lot

and long process losing.

losing.

⇒ Turnaround time reducing

If not mentioned
assume
 $DL = 0$

Large process, waiting process increases
short process, waiting process decreases

Waiting process:

$$P_1 : (0-0) + (10-1) = 9$$

$$P_2 : (1-1) = 0$$

$$P_3 : 17 - 2 = 15 \text{ ms}$$

$$P_4 : 5 - 3 = 2 \text{ ms}$$

$$\text{Avg} = 6.5$$

	Throughput	Avg. turnaround time	Waiting time
FIFS.	0.15	15.25	8.75
SJF (Non preemptive)	0.15	14.25	7.75
SJF (Preemptive)	0.15	13.00	6.50

$DL = 1 \text{ ms}$ (considering dispatch latency,
CPU utilisation reduces)

NP 0.13 16.00 9.5

P 0.13 16.00 9.50

These values become worse in general

Preemptive more context switching
more latency \Rightarrow less CPU utilisation.

If $DL = 0$ ms.

If using non-preemptive algo, & girding all processes together,

Response time \approx waiting time

Else if using preemptive algo,

Response time may be diff waiting time.

Process 1: $W_1 = 10$, $T_1 = 10$, $R_1 = 10$
Process 2: $W_2 = 10$, $T_2 = 10$, $R_2 = 10$
Process 3: $W_3 = 10$, $T_3 = 10$, $R_3 = 10$
Process 4: $W_4 = 10$, $T_4 = 10$, $R_4 = 10$
Process 5: $W_5 = 10$, $T_5 = 10$, $R_5 = 10$
Process 6: $W_6 = 10$, $T_6 = 10$, $R_6 = 10$
Process 7: $W_7 = 10$, $T_7 = 10$, $R_7 = 10$
Process 8: $W_8 = 10$, $T_8 = 10$, $R_8 = 10$
Process 9: $W_9 = 10$, $T_9 = 10$, $R_9 = 10$
Process 10: $W_{10} = 10$, $T_{10} = 10$, $R_{10} = 10$
Process 11: $W_{11} = 10$, $T_{11} = 10$, $R_{11} = 10$
Process 12: $W_{12} = 10$, $T_{12} = 10$, $R_{12} = 10$
Process 13: $W_{13} = 10$, $T_{13} = 10$, $R_{13} = 10$
Process 14: $W_{14} = 10$, $T_{14} = 10$, $R_{14} = 10$
Process 15: $W_{15} = 10$, $T_{15} = 10$, $R_{15} = 10$
Process 16: $W_{16} = 10$, $T_{16} = 10$, $R_{16} = 10$
Process 17: $W_{17} = 10$, $T_{17} = 10$, $R_{17} = 10$
Process 18: $W_{18} = 10$, $T_{18} = 10$, $R_{18} = 10$
Process 19: $W_{19} = 10$, $T_{19} = 10$, $R_{19} = 10$
Process 20: $W_{20} = 10$, $T_{20} = 10$, $R_{20} = 10$
Process 21: $W_{21} = 10$, $T_{21} = 10$, $R_{21} = 10$
Process 22: $W_{22} = 10$, $T_{22} = 10$, $R_{22} = 10$
Process 23: $W_{23} = 10$, $T_{23} = 10$, $R_{23} = 10$
Process 24: $W_{24} = 10$, $T_{24} = 10$, $R_{24} = 10$
Process 25: $W_{25} = 10$, $T_{25} = 10$, $R_{25} = 10$
Process 26: $W_{26} = 10$, $T_{26} = 10$, $R_{26} = 10$
Process 27: $W_{27} = 10$, $T_{27} = 10$, $R_{27} = 10$
Process 28: $W_{28} = 10$, $T_{28} = 10$, $R_{28} = 10$
Process 29: $W_{29} = 10$, $T_{29} = 10$, $R_{29} = 10$
Process 30: $W_{30} = 10$, $T_{30} = 10$, $R_{30} = 10$

(3) Priority Scheduling

Class 12
Experiments

Each process is assigned a priority and it should be whole number.

↳ whether system or user process

→ higher priority process will be executed first.

→ In many OS, priority can be 0-7
and in some OS,

0-1023.

→ In some OS, 0 is highest priority
and in some 7 is highest priority.

All depends on OS developer.

factors on which priority depends,

1.) Internal / technical factors.

2.) External / organisational factors.

→ It depends on whose process it is
whether student, faculty, admin.

→ no depends on non-computers

→ non-technical factors,

Internal factors → Tim limit, memory requirement,
no. of files it requires, no. of I/O devices it
requires.

Priority is assigned when PCB is created and it's
stored in PCB.

FIFO queue \rightarrow priority queue
names

\rightarrow Priority queues can be both preemptive and non-preemptive.

\rightarrow SJF is a special case of Priority scheduling.

SJF priority of $\frac{1}{\text{next CPU burst}}$.

Starvation

Indicates wait for processor.

Method to solve

① Aging

\rightarrow gradually increase the priority of process waiting for a long time.

\rightarrow set limit to increase in priority (say upto 255).

② Priority Inversion

Priority goes high

Assign the priority of high to low.

<u>Process</u>	<u>Arrival Time (ms)</u>	<u>CPU burst time (ms)</u>	<u>Priority</u>
P1	0	8	7
P2	2	4	5
P3	4	6	0 (max)
P4	6	4	1

① Non-Preemptive priority scheduling.

Ready Queue

P1, P2, P3, P4

Gantt Chart

	P1	P3	P4	P2
0	8	14	18	22

a) Throughput = $\frac{4}{22} = 0.18 \text{ proc/ms.}$

b) Turnaround time

$$P_1 = 8 - 0 \\ = 8 \text{ ms.}$$

$$P_2 = 22 - 2 = 20 \text{ ms.}$$

$$P_3 = 14 - 4 = 10 \text{ ms.}$$

$$P_4 = 18 - 6 = 12 \text{ ms.}$$

$$\text{Avg} = \frac{8+20+10+12}{4} = \frac{50}{4} = 12.5 \text{ ms.}$$

c) Waiting Time

$$P_1 = 0 - 0 \text{ ms.}$$

$$P_2 = 18 - 2 = 16 \text{ ms.}$$

$$P_3 = 8 - 4 = 4 \text{ ms.}$$

$$P_4 = 14 - 6 = 8 \text{ ms.}$$

$$\text{Avg} = \frac{28}{4} = 7 \text{ ms.}$$

② Preemptive Priority Scheduling

Ready Queue: P1, P2, P1, P3, P2, P4

Gantt Chart:

P1	P2	P3	P4	P2	P1
0	2	4	10	14	16

22.

a) Throughput

$$= \frac{4}{22} = 0.18 \text{ proc/ms}$$

b) Turnaround time

$$P_1: 22 - 0 = 22 \text{ ms},$$

$$P_2: 16 - 2 = 14 \text{ ms},$$

$$P_3: 10 - 4 = 6 \text{ ms},$$

$$P_4: 14 - 6 = 8 \text{ ms},$$

$$\text{Avg} = \frac{22+22+6+8}{4} = 12.5 \text{ ms},$$

c) Waiting Time

$$P_1: (0-0) + (16-2) = 14 \text{ ms},$$

$$P_2: (2-2) + (14-4) = 10 \text{ ms},$$

$$P_3: (4-4) = 0 \text{ ms}$$

$$P_4: (10-6) = 4 \text{ ms}$$

$$\text{Avg} = \frac{14+10+4}{4} = \frac{28}{4} = 7 \text{ ms},$$



④ Round Robin (RR) Scheduling

23/11/18
Class 13/4

OS developer define some time quantum

Time quantum \approx time slice
 $\rightarrow 15\text{ ms. (let's say)}$

Processor will be allowed to execute ~~a~~ process
 at max one time quantum.

\rightarrow Ready queue is implemented as a circular queue
 \rightarrow It is always preemptive.

If Time Quanta is very large,
 RR will be same as FCFS.

If time quanta is very small
 there is so much preemption.

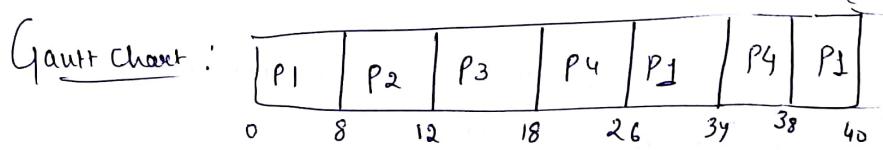
Choose time quanta such 80% of CPU burst is less than
 time quanta.

\rightarrow Once it is decided, it is hard coded in OS,
 if we have to change, we make new version and
 change programming.

Ex \rightarrow	Process	Arrival time (ms)	CPU burst time (ms)
	P1	0	18
	P2	2	4
	P3	4	6
	P4	6	12

$$TQ = 8\text{ ms.}$$

Ready Queue : P₁, P₂, P₃, P₄, P₁, P₄, P₁



(a) Turnaround Time

$$P_1 : - 40 - 0 = 40 \text{ ms}$$

$$P_2 : 12 - 2 = 10 \text{ ms}$$

$$P_3 : 18 - 4 = 14 \text{ ms}$$

$$P_4 : 38 - 6 = 32 \text{ ms}$$

Dispatch latency = 1

Completion

$$= 68$$

85.11%

$$\text{Avg} = \frac{40+10+14+32}{4} = \frac{96}{4} = 24 \text{ ms}$$

(b) Throughput

$$: \frac{4}{40} = 0.1 \text{ proc/ms}$$

(c) Waiting Time

$$P_1 : (0-0) + (26-8) + (38-34) = 4 + 18 = 22 \text{ ms}$$

$$P_2 : (8-2) = 6 \text{ ms}$$

$$P_3 : (12-4) = 8 \text{ ms}$$

$$P_4 : (18-6) + (34-26) = 12 + 8 = 20 \text{ ms}$$

$$\text{Avg} : \frac{22+6+8+20}{4} = \frac{56}{4} = 14 \text{ ms}$$

$$T_\theta = \sin \theta$$

Ready Queue : $\theta_1, \theta_2, \rho_1, \rho_2, \rho_3, P_2, P_4$

Count check : 01 11 02 01 3 2 4 | 1 3 4 | 1 3 4 | 1 4 | 1 4 | 1 0

→ let us assume, that the process goes from running to ready
put first,

a) Turnaround time

$$P_1 = h - ah_{\text{loss}}.$$

$$P_{22} : 12 - 2 = 90 \text{ ms.}$$

$$f_3 : g_4 - 4 = 20 \text{ ms},$$

$$P_{14} : 38 - 6 = 32 \text{ m/s.}$$

$$Av\delta = \frac{50+52}{4} = \frac{102}{4} = 25.5 \text{ m/s.}$$

$$5) \frac{\text{Waiting Time}}{P_1} = \frac{(6-4)}{(14-8)} + \frac{(80-14)}{(86-32)} + \frac{(60-86)}{(30-51)} + \frac{(38-34)}{(38-34)}$$

$$9 + 6 + 4 + 4 + 6 = 30 \text{ ms.}$$

$$\begin{array}{l} \text{P2: } (4-2) + (10-6) \oplus = 6 \text{ ms}, \\ \text{P3: } -(8-4) + (16-10) + (22-18) \ominus = 4+4+6 = 14 \text{ ms}, \\ \text{P4: } (11-1) + (21-9) + (28-36) + (32-30) + (36-34) \end{array}$$

$$= 84 - (18 - 6) + (18 - 6) + (18 - 6) + (18 - 6) = 84.$$

grows, $g = g_0 + g_1 t + g_2 t^2 + \dots$

$$\Delta t = \frac{60}{80+80} = \frac{60}{160} = 15.5 \text{ ms.}$$

CPU utilization = 68.97%

Dispatch latency = 1.

⑤ Multilevel Queue Scheduling

→ Ready queue is partitioned into multiple queues,

Each queue will be there for ~~one~~^{each} type of process.

for ex - separate background process, user process,
text processing.

Each of these queue, will be having its own
scheduling algorithms.

two types of scheduling

① ↳ inter queue

② ↳ intra queue,

① we can use

(a) Preemptive priority,

Each queue will be having its priority.

(b) RR (80% time for ~~one~~ foreground process

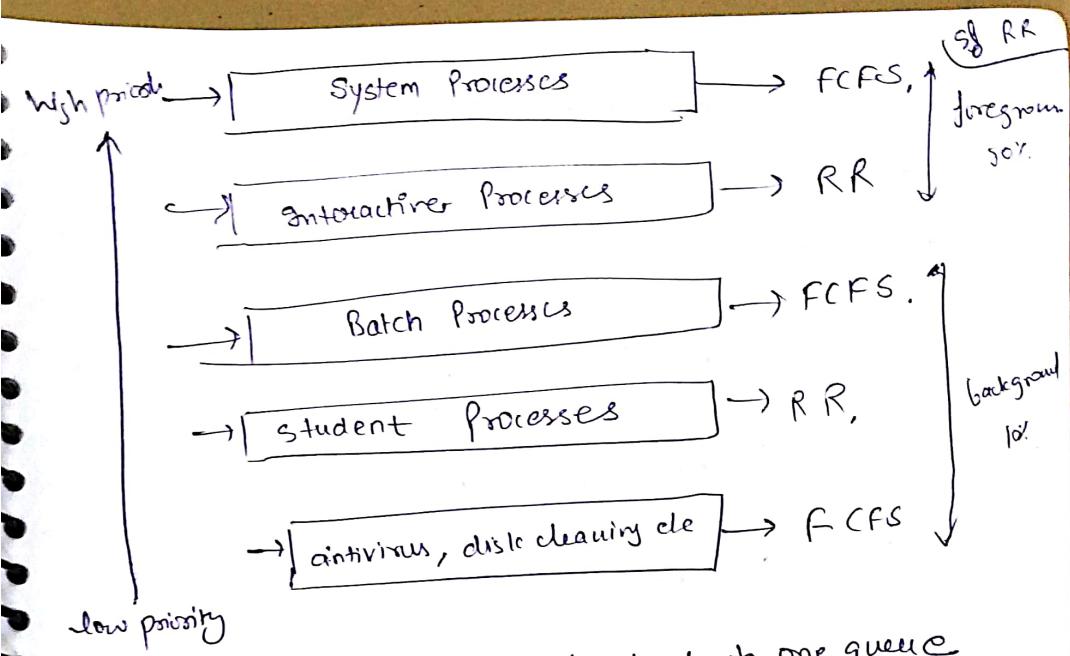
and 20% time for background process)

ex: Antivirus

② we can use

it is upto OS developer,

book says it is FCFS



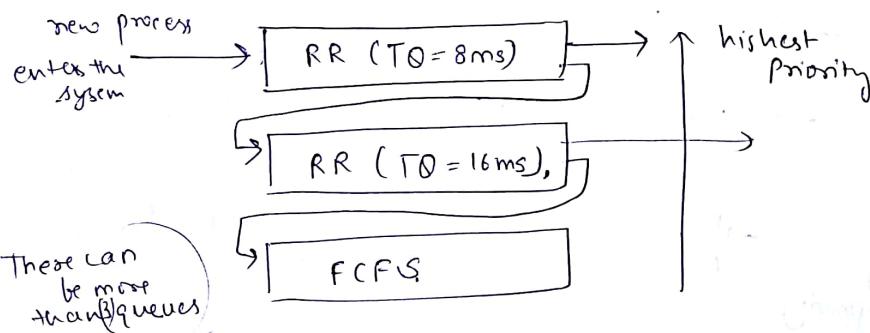
- * One process is assigned/located to one queue permanently,
 - * It ensures low waiting time for high priority processes.
 - * Other processes too have satisfactory performance.
- Problem
- It is bit complex, and CPU scheduler becomes slow.

Q) Ex-5 on Pg 12

⑥ Multilevel feedback queue scheduling.

→ It is decomposed into multiple queues.

→ Each queue has its own priority.



cpc burst time

if a process has less than 8 ms, then it
terminated.

else

it go to next priority

→ Each queue is using its own algorithm, called intra queue scheduling.

→ interqueue scheduling also,

If two processes, which queue to give,
use preemptive priority.

In ⑤ algo.,
process is assigned to queue for eternity,
process here can move in queue.

→ A process waiting too long in a low priority queue may be allowed to move higher priority queue.

One numerical given in notes,
do it manually.

Generally numerical didn't come from algo ⑤ and ⑥

but → sometimes bad luck. ⑦ |

Parallel Computing :-

→ If you can execute more than one instruction at a particular time.

→ Anything more than ~~one~~ one (1) can be (1,5)

→ Multicore processor

A single processor has more than one core, which can have ALU, AC, PC, registers, cache

support.

→ Parallel computer
Distributed computer

have multiple microprocessor on the same motherboard

Parallel computer = highly coupled computer

Distributed computer = loosely coupled computer

- There is only one OS.
Ex- for parallel computing.
 - VLIW
 - Cluster Computing
 - Grid Computer
 - Cloud Computing
 - GPU
- Read from Wikipedia

Multiple Processor Scheduling

- We have multiple CPU.
- We have to schedule them together.

Multiple processor system

- homogeneous
 - All processors are same
 - Symmetry in connection
 - RAM, peripheral devices are same

heterogeneous

(otherwise)

CPU scheduling

- Asymmetric
 - there is one processor
 - one ready queue

Symmetric

- Each processor is self scheduling
- Each processor has its own ready queue
 n processors \rightarrow n queues

Symmetric \Rightarrow

Two interesting issues

① Processor affinity

affinity of a process

or

necessity of a process to

be executed,

→ we migrate a process, only when it is strongly needed,
or it's worth.

Processor affinity

hard → we must execute process on same processor.
Ex - if process has opened lot of files

soft

soft → if possible, we execute on same processor
Ex - processes% complete

Load Balancing

→ All the processors should be equally loading.

→ for optimal use.

① push migration

if diff. in no. of process

then one processor can push processes,

② pull migration

lightly loaded can pull processes.

But there is some threshold, limit to every push or pull

③ Combination of both

load-balancing,

balance b/w ~~migration~~ & processor affinity.

Also consider affinity.

if a process wants to continue execution on
a processor, then that process should not
pull/push.

chapter 3, 4, 5 done

Let,

initially: $i = 10$.

$P_0 : R_1 = 10$

$P_1 : R_2 = 10$

$P_0 : R_1 = 11$

$P_1 : R_2 = 9$

$P_0 : i = 11$

$P_1 : \boxed{i = 9}$

if we order is
different,

we execute P_1 before P_0

then,
~~at~~ the end,

result will be different.

This problem is known
as a race condition.

→ we have to do
avoid this

because it will unpredictable
results.

we can avoid
race condition

by Process Synchronization

P_0 and P_1 should
consult ~~to~~ each other
before changing any value

When P_0 executed, P_1 should not interrupt or vice-versa

↳ To ensure,
that only one process can manipulate data at one
time

shared variable ↗ both can
read,

→ Process synchronization is the task of a programmer

→ OS provides some mechanism to ensure process
synchronization.

Critical Section Problem

is a way to model cooperating processes

We have, $P_1, P_2, P_3, \dots, P_n$. (processes)

all of them are using some shared variable,
whole code is not accessing

but only part of code is accessing shared variable, called critical section,

The critical section problem is to ensure that only one process is executing its critical section.

Program.

```
do {  
    Remainder section  
    Entry section.  
    [critical Section]
```

a process requests permission, or seeks permission, to enter its own critical section

only enter when other n-1 processes agrees otherwise wait here

```
Exit section  
Remainder section  
} while(TRUE);
```

section that says everyone that its own critical section is complete, off and other process can not execute its critical section.

Three properties

→ necessary condition for a solution of critical section problems —

(Short question can come - ?

Please use terms as mentioned.

1/2/2018
class-18

① Mutual Exclusion

→ It is a locking mechanism.

→ At one time, one of the process can access shared variables

Both read & write

② Progress

$k \leq n$

Which process to come in critical section next will be decided by the all processes which are in entry section ~~and~~ or in the exit section.

- * Decision can not be delayed indefinitely.
- * Which process gets to progress next.

③ Bounded Waiting

There should be upper limit of the no. of times other processes ~~can allow~~ enter in its critical section after one process request to enter into its critical section, otherwise it will lead to deadlock before that request is granted. Otherwise, it will lead to starvation.

Peterson's Solution

Peterson proposed a solution for critical section problem. (only for 2 processes),

p_0 and p_1

In this solution,

there are two shared variable,

```
int turn;
```

```
Boolean flag [2];
```

if $\text{flag}[0] = \text{True}$

P_0 wants to
to enter

if P_0 exists,

into its
critical
section

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

same with P_1 ,

$\text{flag}[1] = \text{True}$ False.

Code for process i , P_i —

```
do {  
    flag[i] = true;  
    turn = i; // something.  
    while (flag[j] && turn == j);  
    // do nothing.  
    critical section.  
    if flag[i] = false;  
    remainder section  
} while (true);
```

P_i give chance
to other process

when P_i will
enter
deciding factor

Q Prove that peterson's solution satisfy the 3 properties?

bounding waiting, (threshold value)

in this case it is one (1)

Semaphores.

5/2/2018

Class-19

(Not Attended)

- Low level construct which can be used to solve critical section problem.
- better than Peterson's solution.
- A semaphore is an integer variable, but it is not an ordinary integer.
- It can be accessed by 2 atomic operations.
Atomic → "indivisible"
either the whole operation occurs or it does not occur at all,

What are these 2 atomic operations?

↳ wait()
↳ signal()

- Semaphores can be used to ensure process synchronization among multiple processes.
- can solve multiple process critical section problem.
- only 1 process can access a semaphore at a time.

Assuming s is any semaphore :-

```
wait(s)
{ while (s <= 0)
    ; // busy waiting           " Busy doing
      ;" )                      " nothing "
  s
signal(s) { s++; }
```

Semaphores are of 2 types

① Binary Semaphore

- can have only 2 values (0 or 1)
- also called mutex/mutex lock (can be used to lock a process)
- always initialized to 1

Let S be a binary semaphore:-

```
do {  
    [wait(S)]  
    CS  
    [Signal(S)]  
    RS;  
} while (true);
```

② Counting Semaphore

When the process are accessing multiple clone copies of the same resource,

- initialized to N (no. of copies of that shared res)
- control access to resources that have multiple instances.

Semaphores can create 2 problems -

- Starvation
- Deadlock.

$$\begin{array}{c} P_0 \\ \hline \text{wait}(S) \\ \text{wait}(Q) \\ \equiv \\ \text{Signal}(S) \\ \& \text{Signal}(Q) \end{array}$$
$$\begin{array}{c} P_1 \\ \hline \text{wait}(Q) \\ \text{wait}(S) \\ \equiv \\ \text{Signal}(Q) \\ \& \text{Signal}(S) \end{array}$$

Deadlock Sample

Producer Consumer Problem (Boundary-buffer problem)

- Producer processes are producing some abstract items.
- There are some consumer processes consuming those abstract items.

Semaphores

denotes the no. of buffers which are currently empty

mutex = 1; // access the buffer
empty = N;
full = 0;

Producers

```
do { // produce an item } RS
    wait (empty); // entry section
    wait (mutex);
    // add item to the buffer → (S;
    signal (mutex); // Exit section
    signal (full); // shared resources
                    // i.e. buffer is being accessed)
```

Consumers

```
do { wait (full);
      wait (mutex);
      // remove an item from buffers
      signal (mutex);
      signal (empty);
      // consume the item }
```

3 while (true);

Readers - Writers Problem

Shared data -

int readcount = 0;

Semaphores -

(mutex = 1;) binary semaphore
wsl = 1; binary semaphore

Writers

```

do {
    entry crit;
    wait(wsl);
    // write
    signal(wsl);
    gwhile(true);
}
    
```

Readers

```

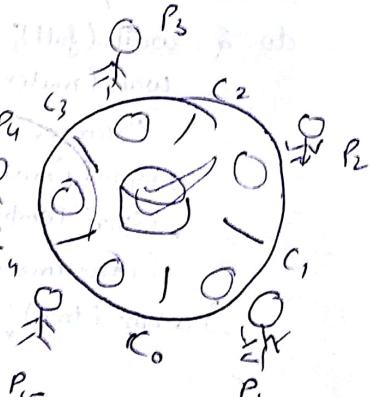
do {
    entry crit;
    wait(mutex);
    readcount++;
    if (readcount == 1)
        wait(wsl);
    signed(mutex);
    // read
    wait(mutex);
    if (readcount == 0)
        signal(wsl);
    signal(mutex);
}
    
```

class 20
6/2/2018

Next Critical Section

Dining Philosophers Problem

You have to write processes for philosopher, we have only 5 chopsticks and there should be synchronization.



5 (chopsticks aren't identical)

* In case of buffer location was irrelevant,

i. Here order matters!

5 semaphores needed.

Semaphore

Semaphores

```

chopsticks[5];
P;
do {
    wait (chopsticks[i]);
    wait (chopsticks[(i+1)%5]);
    // Eat.
    signal (chopstick[i]);
    signal (chopstick[(i+1)%5]);
    // Think
} while (true);

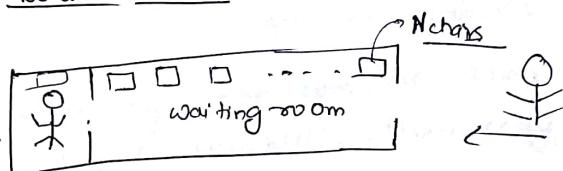
```

what will be happen if all the philosophers get hungry at the same time?

Deadlock

No will able to eat, all will die (deadlock)

Q) Sleeping Barber Problem



Case 1:- No empty chair, goes back.

Case 2:- Some vacant chairs, goes & sit.

Case 3:- All vacant, goes into barber's room

barber is busy with a customer

when barber gets free, catches one of the customer and cut his hair.

barber is free and is sleeping

↓
customer wakes up the barber.

Barber

wait - seat[n];

do for barber
use counting sem

Deadlocks.

Multiple process trying to access same h/w or s/w.

Class-2

8/2/2018

(Not Attended)

Request
Permission granted } = waiting state
use resource
release that resource.

request → use → release

- Processes are waiting indefinitely
- OS handles deadlock → Job of OS.

Resources : h/w or s/w.

Resource types:

1. Resources which have some functionality and can be used in exchange of each other → instances of that particular resource

type

e.g. - B, C, D, E → general purpose resource type

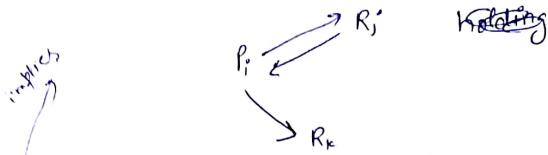
- if necessary not complete, then will not happen but if complete, then may or may not occur.
- if sufficient complete, then will happen, but if not complete, then may or may not occur.

↗ (Read monitor on your own)

Necessary conditions for deadlocks:

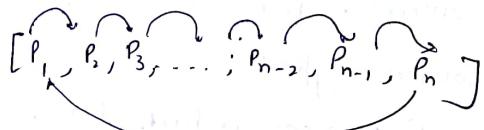
- 1) Mutual exclusion: There must be one or more non sharable resource in system. Ex: printer, register.

2.) Hold & wait : Now process P_i is holding R_j .



3.) No preemption : There are some resources which can be preempted like HL register, suppose one process is using HL & some other process needs it, so it can snatch, at the same time, there are certain resources like printer, speakers which can't be preempted,

4.) Circular wait :



a process waits for

a resource held by the next process.

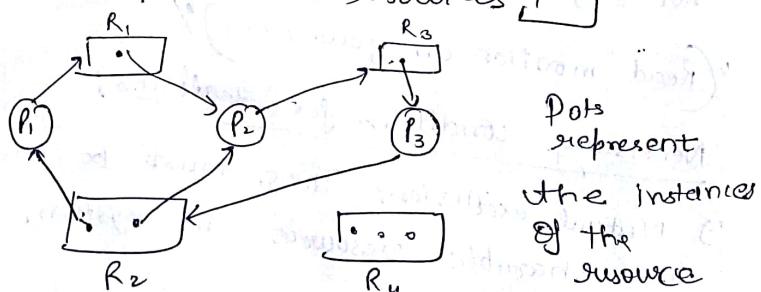
If there is hold & wait, only then we can have Circular wait.

If we can ensure any one of the four is false, then no deadlock.

Resource allocation Graph.

Circular nodes to represent the processes (P_1)

Rectangles to represent the resources (R_1, R_2, R_3, R_4)



Request Edge: when a process requests for a resource.

Assignment Edge: when a process is assigned a resource.

If there is a cycle, there can be deadlock. (necessary condition).

there are 2 cycles: $P_1 R_3 P_3 R_2 P_2 \& P_1 R_1 P_2 R_3 P_3 R_2 P_1$

if each resource has only one instance, then occurrence of cycle is sufficient condition for deadlock.

If in a graph, some resources have one instance & some have multiple instances, then if in a cycle if each resource has one instance then sufficient condition for deadlock.

Methods for deadlock handling :-

- Deadlock prevention: (prevention is better than cure).
To ensure atleast 1 cycle detection $O(n^2)$
 $O(\epsilon)$,
of the 4 condition can't occur.
- Deadlock avoidance: Each process tells the OS its maximum requirement of resources as soon as it enters. Now it becomes the job of OS to schedule processes such that no deadlock occurs.
- Deadlock detection & recovery from deadlock.

- Only one of the above three is required.
- If deadlock occurs, some processes hang.
- Deadlocks are very rare, once in a year.
- Modern OS don't handle deadlocks.
- The user has to kill such processes on their own (lot ④),

(4)

Deadlock Prevention

Class-22

12/2/18

(Not attended)

① Mutual Exclusion

② Hold & Wait:

① → Protocol: Request & get all the resources at the beginning (no hold & wait).

Hold & wait condition made false, now resource utilization and starvation of processes.

R

(2)

② Protocol :- should release its current resources before requesting for more resources.

Releasing the current resource & then asking for more resources. Not result into low resource utilisation

→ If a process requires several resources at same time, then they (processes) may end up starving.

③ No Preemption

Some resources like files, registers, memory location can be preempted,

De

If we can make all resources preemptable,

④ Circular wait:

Circular wait implies hold & wait. If we can ensure that there is no circular wait, then no deadlocks (hold & wait can be still there). Arrange the resources types as $R_1, R_2, R_3, \dots, R_n$ (fixed order).

Protocol: a process can request resources in an increasing order only when process has obtained needed instances of R_i , it can ask for R_2 & so on. & it has got resources for R_2 , it cannot ask for R_1 again. No circular wait happen.



Requires multiple resources of R_i , it should ask the same resources at same time.

② Protocol: if a process requests for R_i , then it must release.

R_j for $\forall j \geq i$

Eg. $i = 3$

R_1, R_2, R_3

process requesting for R_3 resource, it must release all R_j 's where $j \geq 3$, "release (age walk)" resources if acquired by mistake.

Disadvantage:

↳ low resource utilization

↳ low throughput

- all instances of a resource must be acquired at the same time

- Difficult to implement.

Deadlock avoidance

requires matrix addition & subtraction.

Problem: matrices can be very large & their size $n \times n$ is static, but varies.

Safe state :- State of the system which denotes the maximum no. of resources that can be allocated w/o creating resources.

Safe sequence :- Sequence of processes

$\langle P_1, P_2, P_3, \dots, P_n \rangle$ if for each process P_i , request for resources which should be less than or equal to ~~some~~ sum of resources of P_j ($j < i$) & the free resources.

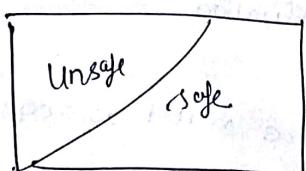
Sequence $\langle P_1, P_2, \dots, P_n \rangle$ is safe if for each P_i , request for resources which should be less than or equal to the resources that P_i may still request can be satisfied by the free resources plus all resources held by P_j with $j < i$.

A process must execute with its own resources, plus free resources & resources of processes on left.

if safe, then no deadlock

if unsafe, there may be a deadlock

Unsafe state is a necessary condition for deadlock



Banker's Algorithm

Must (High probability)

(Numericals)

class-23,

24

② 13/2/18

Whenever a process, asked for more resources,
then OS calculate if it is possible to allocate,
then it will granted resource.

→ whether allocating resource leave in a safe state
or not.

We have,

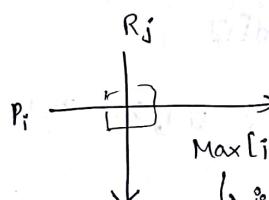
n processes

m resource types

Data structures →

→ Available [m] // gives no. of free instances of each resource

→ Max [$n \times m$]



b_{ij} is the max. no. of instances of resources R_j required by process i .

→ Allocation [$n \times m$]

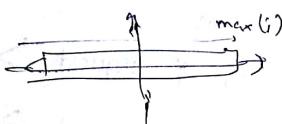
Now, many instances of resource R_j is allocated to process p_i .

→ Need [$n \times m$]

Max - Allocation

denotes how many instances of R_j can be further requested by p_i

→ Max [i] → denotes



Algorithm has two parts

- 1) ① Safety
- 2) ② Request

Resource Request

① Safety Algorithm checks whether a system is in a safe state or not.

Steps
1) Let $Work[m]$ and $Finish[n]$ be vectors.
 m integer
 n boolean.

$Work = Available$

$Finish[i] = False, \forall i = 0, 1, 2, \dots, n-1$

2.) find a process P_i such that $(Finish[i] == False \& Need[i] \leq Work)$

3.) if no such P_i can be found then
 goto 4

3). $Finish[i] = true$

$Work = Work + Allocation[i]$

Goto 2

4.) if $Finish[i] == true, \forall i = 0, 1, 2, \dots, n-1$, then

 state is safe
else unsafe

Time complexity is $O(mn^2)$.

② Resource Request Algorithm

i) is requesting for resources and $Request[m]$
gives the request

Steps
1) if $request > Need[i]$, then error

2) if $request > Available$, then wait

3) Pretend to allocate the resources.

$$\text{Available} = \text{Available} - \text{Request}$$

$$\text{Allocation}[i] = \text{Allocation}[i] + \text{Request}$$

$$\text{Need}[i] = \text{Need}[i] - \text{Request}$$

If the resultant state is safe, then allocate resources.
else undo changes Allocation, need, and available.

Q) A system has 4 processes, namely P₁ to P₄ and three resource types. A system in a safe state, with available = [1 1 1], Max =

$$\begin{bmatrix} 0 & 1 & 1 \\ 6 & 2 & 4 \\ 4 & 3 & 5 \\ 1 & 2 & 1 \end{bmatrix}$$

and

$$\text{Allocation} = \begin{bmatrix} 0 & 0 & 1 \\ 4 & 2 & 3 \\ 3 & 3 & 3 \\ 0 & 2 & 1 \end{bmatrix}$$

→ Which of the following request should be granted?

a) Request = [1 0 1] from P₃

b) Request = [1 0 1] from P₂

Soln

RR algo —

$$\text{Need} = \begin{bmatrix} 0 & 1 & 0 \\ 2 & 0 & 1 \\ 1 & 0 & 2 \\ 1 & 0 & 0 \end{bmatrix}$$

Step 1 X

Step 2 X

We go to step 3

$$\text{Available} = [0 \ 1 \ 0]$$

$$\text{Allocation} = \begin{bmatrix} 0 & 0 & 1 \\ 4 & 2 & 3 \\ 4 & 3 & 4 \\ 0 & 2 & 1 \end{bmatrix} \quad \text{Need} = \begin{bmatrix} 0 & 1 & 0 \\ 2 & 0 & 1 \\ 0 & 0 & 1 \\ 1 & 0 & 0 \end{bmatrix}$$

D) It is Resultant state.

Now we have to check, whether it is safe state or not.

We check it by Safety Algorithm.

Initially	Work			Finish
	0	1	0	F F F F
Select P1	0	1	1	T F F F
Select P3	4	4	5	T F T F
Select P2	8	6	8	T T T F
Select P4	8	8	9	T T T T

this is a sequence
(need not be unique)

Yes, request can be granted.

Next week SE presentation
Syllabus will be completed on Thursday

Your favorite topic
(for 5 minutes)
(no ppt)

from monday
(1 to 10 roll number)
(If somebody is absent)
(He/she will get bengali result)

Deadlock Detection

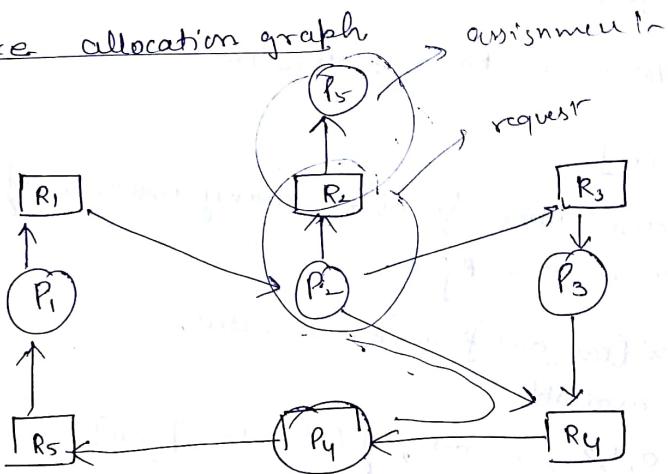
Class 25

15/2/18

Handling deadlocks

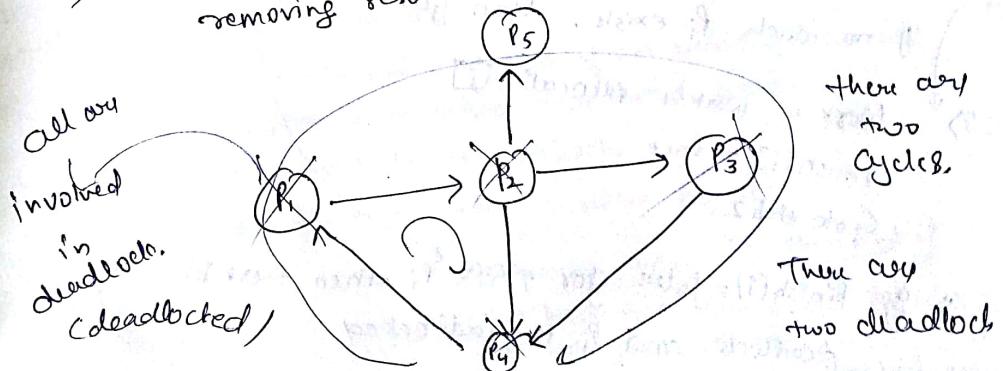
Let us assume, that each resource type has only one instance,
then we can use graphical methods for detection.

Resource allocation graph



If there is a cycle in a graph, then there will be a deadlock.

wait-for graph
removing resource nodes



Cycle detection algorithm

↳ Time complexity $O(n^2)$

①

use this providing all resource type has single instance

(To detect deadlock)

~~be exception with~~

in which resource type has multiple instance

we use

Modified Version of Banker's Algorithm

Request [$n \times m$]

↳ which process asking for how many instance of a particular resource R_j

1) Let work[m] and finish[n] be vectors.

work = available

for $i = 0, 1, 2, \dots, n-1$, if ($\text{Allocation}[i] = 0$)

then $\text{finish}[i] = \text{false}$

else

$\text{finish}[i] = \text{true}$,

2) Find a P_i such that $(\text{finish}[i] == \text{false}) \& (\text{Request}[i] \leq \text{work})$

if no such P_i exists, then goto step 4.

3) $\text{Work} = \text{Work} + \text{Allocation}[i]$

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

Goto step 2

4) if $\text{Finish}[i] = \text{false}$ for some P_i then there is a deadlock and P_i is deadlocked

Time Complexity : $O(m \times n^2)$

(In this topic, no numerical sol.)

Q7, How often to use this algo?

→ at fixed time interval,

→ when you are suspicious,

When you think there may be deadlock.

Ex → CPU utilization falls, beyond
a threshold,

(CB)

Recovery from deadlocks.

→ Process termination, or kill processes,

(terminate all the deadlocked process)

but you are decreasing overall system performance

(not efficient)

→ Kill one process, then check again that deadlock removed or not,

Kill processes one by one,
(bit efficient),

→ (3), Resource Preemption.

take away off resources from waiting process
and give it to other deadlocked process

→ we have to use rollback,

→ we have to select a victim process.

↳ a process may be starved if
we have not used good logic
or not good selection.

• Same process is getting selected again and again

3 types → Deadlock avoidance
Deadlock prevention
Deadlock detection }
Completed
(till midsem)

- Coding → process synchronization
↳ Peterson algorithm or changes in them
↳ pseudo code of classical problem in
process synchronization
↳ code for monitors.

→ Theoretical question

(in depth) (not simple)

- Read Threads (question can come), (short note)
↳ 4.1, 4.2 (6th edition)
→ numericals (main source,
↳ CPU scheduling
↳ banker
↳ use of banker in detection

1.1, 1.4 - 1.9, 1.11
(2.1 - 2.5), 2.7, 2.10
3.1, 3.2, 3.3, 3.4
4.1, 4.2
5.1, 5.2, 5.3, 5.4
6.1 - 6.7
7.1 - 7.7