

OS

1) What is OS?

- Interface b/w user & H/W
- Set of utilities

2) Services :-

- Program execution
- CPU scheduling
- Memory, file, resource management
- Error detection
- Accounting
- security
- Inter process communication.
- I/O operations

3) Goals:

- | | | |
|---------------|-------------|---------------------------------------|
| • convenience | • Efficient | • Reliable |
| • Scalable | • Portable | • Robust
handle
unseen
error |

→ Types of OS:

1) Uniprogramming OS: only 1 program can reside in MM at a time.

• Disadv: inefficient utilization of resources

2) Multiprogramming OS: allows more than 1 process in MM at a time.

• If we increase degree of multiprog., CPU utilization increases but upto a certain limit.

* Degree of Multiprogramming: No. of processes in MM at a time.

• TYPES :- Preemptive, Non-Preemptive

• Non-pre, process leaves CPU if it completes execution or goes for IO op.

NOTE: When system starts, bootstrapping process is done to bring OS from hard disk to MM.

3) Multitasking OS (Time sharing OS)
Multiprogramming + preemptive round robin scheduling

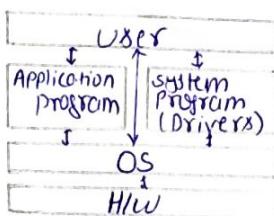
4) Multiprocessing OS: required on system with multiple processors.

• TYPES: Shared memory system (Tightly coupled) Distributed system (Loosely coupled)

5) Multiuser OS: All users can execute their separate processes together on a system having single CPU.

NOTE: Microsoft OS are not multiuser whereas Linux OS are.

6) Real Time OS:
OS gives deadline to processes.



• TYPES: Soft Realtime, Hard Realtime

7) Embedded System OS: used to communicate b/w mechanical machine & processor.

8) Handheld devices OS: of mobiles.

* System call: way for appl. programs or users to interact with OS. Invoked using slow interrupt.

1) PARTS OF OS:

- Kernel: responsible for all functionalities
- Shell: used to interact with kernel.

2) In windows we have 3 shells:-
CMD, PowerShell, GUI

3) DUAL MODE OF OPERATION: for protection

• User mode (mode bit 1): appl. process is running & application process can perform. Only non-privileged operation

• Kernel/System/Supervisor mode: (mode bit 0): to do privileged opr. Appl. process makes system call. (sw interrupt)

1) Program is binary file stored in disk.

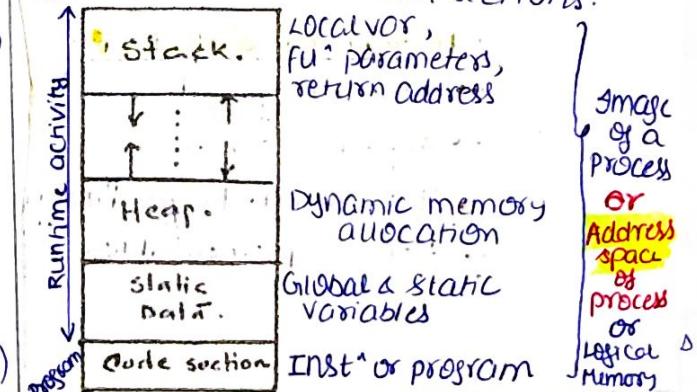
2) Running program is called Process

- Unit of execution (CPU)
- Locus of control (OS)
- Schedulable/Dispatchable unit

3) Process as a data structure:

- Definition: program
- Representation/Implementation
- Operations
- Attributes

4) Process Representation: Every process is stored in memory in 4 sections.

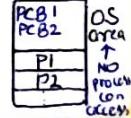


5) Operations of a process:-

- Create
- Schedule
- Wait
- Suspend, Resume
- Terminate

6) To control process, OS maintains

1 PCB for each process.
PCBs are stored under OS area in MM & processes cannot access it.

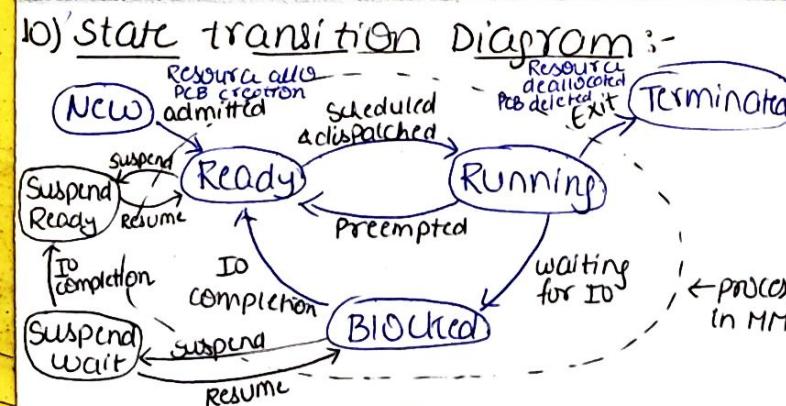


- 7) Each process has a unique process ID.
- 8) Process control block (PCB): collection of attributes of process (PID, PC, GPR's, list of files, devices, type, size, memory limit, priority, state, etc)
- Also called process descriptor.

9) The content of the PCB at the moment is called context of that process.

- 10) Context switching: saving CPU register values into PCB of current running process & loading other process context into CPU registers.
- Program of OS which does context switching is called **Dispatcher**.

- * For device queue of disk, scheduler is required, for all other IO devices scheduling is FCFS.
 - * If CPU utilization is less, long term sch. increases degree of multiprogramming.
 - * Swapped out process is still a process as PCB is present in MM.
 - * PCB is updated only by context switching
 - * If parent process is killed by some other process, its child process will also terminate.
 - * The state of a process is defined by the current activity of the process.
- NOTE: These formulas are with assumption that process does not require IO.



- 11) Running \rightarrow Terminated / **voluntary** transition done by process with
- Running \rightarrow waiting
- Rest all transitions are done by OS.

- 12) If a process is killed forcibly, transition may take from Ready / waiting to terminated.

NOTE: While running also process is in MM only context is with CPU.

- 13) Process \rightarrow CPU bound (CPU extensive) \rightarrow IO bound (IO extensive)
- For better resource utilization, OS maintains a good mixture of both and **IO bound process is scheduled before CPU bound process**.

Process scheduling

1) Queues:

- Job Queue • Ready Queue • Device Q (Newstate) (Ready state) (each IO, Queue of PCB device has)

2) Types of schedulers :-

- a) LTS: New \rightarrow Ready (Job scheduler)
- b) can increase degree of multiprogramming
- b) STS: Ready \rightarrow Running (CPU scheduler)
- c) does not affect degree of multiprogr.

- c) HTS: swaps from MM to swap space of disk & back to MM.

- PCB is still present in MM.
- can increase or decrease degree of multiprogramming.

NOTE: Swap space is under OS control

- 3) If swapping is done based on priority of process it's called **ROLING**

- 4) A process either in ready or waiting state can be swapped

- 5) Zombie process: if PCB of terminated process remains in MM.

CPU Scheduling: (short term sch.)

1) Goals:-

- Minimize WT & TAT
- Maximize CPU utilization & throughput
- No starvation

2) Scheduling Times:

- @ Arrival time: - time at which process is admitted by long term scheduler.

- ⑥ Burst / Service time: - time required for execution of a process.

- ⑦ Completion time

- ⑧ Turn around time: - total time from arrival to completion.

$$TAT = CT - AT = BT + WT$$

- ⑨ Waiting time

$$WT = TAT - BT$$

- ⑩ Response time: - time from arrival to first response by CPU

3) Scheduling length:

$$L = \max(AT) - \min(AT)$$

- 4) Throughput - NO of process executed per unit time. $T = \frac{n}{L}$

5) In non-preemptive, $WT = RT$ (g no op.)

Scheduling Algo :- (Gantt chart always starts from 0)

1) FCFS : Criteria: AT Tie breaker: PID Mode: Non-preem.

- Only FCFS suffers from Convo effect.
- No starvation

* Convo effect: Large process is scheduled first (among process having small processes) it will make other processes to wait lot.

2) SJF: Criteria: min BT Tie breaker: FCFS Mode: Non-preem.

- Min avg WT among all process.

* Longest job first may suffer from Convo effect.

* Starvation: If new processes keeps coming & executing and old process vation keeps waiting for indefinite time.

* Generally throughput other algo $\leq SJF \leq SRTF$

NOTE: If the process running is scheduled again, then context switch is not required. starvation (Required only in RR.)

NOTE: Optimal process algo $\rightarrow SRTF$
Optimal non-preem. algo $\rightarrow SJF$

3) SRTF: Criteria: min remaining BT Tie breaker: FCFS, Mode: Preemptive

- Min avg WT among all scheduling algo
- Better throughput, no convo effect in continuous run
- Starvation

NOTE: No practical implementation of SJF & SRTF as BT is not preknown

NOTE: Min no. of context switches in preempt execution = $(n-1)$

* If all process arrive at same time, SRTF reduces to SJF.

4) HRRN: favours shorter job also reduces WT of longer job.

Criteria: Response ratio ($\frac{WT+BT}{BT}$) $\propto \frac{WT}{BT}$

Tie breaker: SJF

Mode: Non-preemptive

- NO STARVATION

5) Priority based:

Criteria: Priority

Tie breaker: given in ques

Mode: Preemp / Non-prem.

- Better RT for real time situation.
- Starvation for low priority process.

Solution of starvation: Aging (possible only with dynamic priority)

Aging: After waiting for specific time, priority is increased.

6) Round Robin: (Time sharing)

Criteria: AT + Q ; Q \rightarrow time quantum
Mode: preemptive

No. of context switch = $\left(\sum \frac{BT}{Q}\right) - 1$

NOTE: In RR, it will be counted as 2 context switch.

\downarrow \downarrow \downarrow \downarrow
Very Small CPU eff ≈ 0 Smal very interactive system Large less interactive system \downarrow Very large degrada to FCFS if $Q > \text{max}(BT)$

No starvation

- Better interactivity
- Best response time. (min)
- Avg WT & TAT is more

NOTE: Throughput & idle period is same with all algorithm.

7) LRTF: Tie breaker: PID Preemptive

- If all process have arrived byon completion of any process then termination is in order of PID.

$\frac{1}{P_0} \frac{2}{P_1} \frac{3}{P_2} \frac{4}{P_3} \leftarrow \Sigma BT$

8) Multilevel Queue Scheduling:

- using different algo for different queues based on requirement.

\rightarrow Fixed Priority Preemptive Sch:-

- All queues have different sch algo priority decreased may be preemptive or non-preemptive.

- Process can't move from 1 queue to other.

- Starvation

- Preemption will happen with arrival of a process in higher priority queue.

\rightarrow Time Slicing: CPU time is divided among multiple queues.

9) Multilevel Feedback Queue Sch:

- Process may change queue, flexible

- Starvation for low priority queue

* If IO is also there, scheduling is like CPU, IO, CPU. $WT = TAT - (CPU + IO + CA)$

- * For n process, CPU utilization = $1 - p^n$
 $p \rightarrow$ fraction of time for IO
- * If there is limited IO, process will execute IO opr. in FCFS.
- All preemptive may behave as non-preemptive
- Only 1 process in MH
- all process arrived together.
- SRTF - later arriving process is longer
- Priority - " " has low prior.
- RR - if Δ is very large
- LRTF - $\frac{1}{\Delta}$ BT of all process is 1

Thread:

- 1) Lightweight process
- 2) Multiple threads of a process can run **parallelly** in multicore system.
- * Shared among thread | unique for each thread

Code Section	Slack
Data Section	Thread id
Heap	Register set
OS Resources (PTBR)	Program counter
Open files	
Address space	
- 3) Thread → user level thread (**not known to kernel**)
 → Kernel Level thread
- 4) In user level threading, process itself does context switching b/w threads.

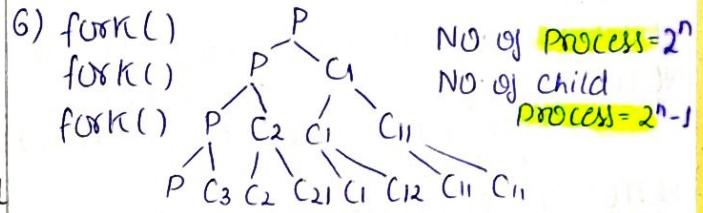
* User Thread

- Multithreading in user process.
- Created without Kernel intervention.
- **context switch faster.**
- If 1 thread is blocked, OS blocks entire process
- Generic & can run on any OS (implemented by libraries)
- **Faster** to create and manage.
- created using `fork()`
- 5) Thread context switching is faster than process context switching.

- 1) System call → requesting service from kernel via **Software Interrupt**.
- Slow interrupt generated when further inst can't be executed.

NOTE: Slow interrupt known as **Trap**. Generally interrupt is given to CPU to OS.

- 2) No system call is made to use stack or heap space.
- 3) Child process will start from next stmt of `fork()`.
- 4) Parent & child can run concurrently (not parallelly).
- 5) `fork` takes no parameter & returns an integer.
 - ve value : if creation of child process is unsuccessful.
 - zero : returned to child process.
 - +ve value : PID of child process is returned to parent process



- 7) `fork` creates child not thread. Child gets **copy of all things of parent**.
- 8) Unless all child processes are terminated, parent can't terminate.
- 9) By default consider child is created successfully.

- * Each kernel level thread has its own **PCB** pointing to same address space.

* Multithreading Model:

1:1, Many:1, Many:Many
 $\xrightarrow{\text{userlevel}}$ $\xrightarrow{\text{kernel level thread}}$

* System call categories:

- Process control
- File management
- Device
- Security

* **fork**: system call to create new process in Linux & UNIX OS.

* `if (fork())` Parent process will execute if, child will execute else.

* **concurrently**: one after other (time sharing)

parallelly, simultaneously: $\xrightarrow{\text{userlevel}}$

NOTE: `printf("Hi")` `printf("Hi\n")`
`fork() fork()` `fork() fork()`

Hi Hi Hi Hi

Hi

Print is buffered, buffer flushes only on using `in` or `flush()`, without flushing buffer is passed to child processes.

Process Synchronization

- 1) Process types: Independent process, Cooperating/Coordinating/Communicating
- 2) Synchronization needed b/w Cooperating process, not independent process.
- 3) Problems without synchronization: inconsistency, Data loss, deadlock.
- 4) Sync is required for only that part of process which is involved in communication i.e. critical section
- 5) Communication b/w process can be by:
 - Message passing
 - Sharing some common resource \rightarrow H/W
- 6) Critical section: code section where shared variable can be accessed.
- 7) OS has to provide sync to CS.
- 8) **Race condition:** when final result of concurrent processes depends on the sequence in which the processes complete their execution.
- 9) When there is mutual exclusion, race condition is not there.

- 10) If 2 processes P1 & P2 r executing concurrently, they can be executed as
 - P1 then P2
 - P2 then P1
 - Concurrently & P2 finishes 1st
 - " & P1 finishes 1st

Critical Section Problem:

- 1) If there is CS, OS needs to provide sync. This is CS problem.
- 2) A perfect solution to CS problem must satisfy all 3 requirements:
Mutual exclusion, progress, bounded waiting.

- NOTE: concurrent : executing 1 after other
parallel : executing simultaneously
- 3) Mutual exclusion: At a time Only 1 process can enter into CS.
 - Progress: If no one is in CS & some one want to go, let it go in CS.
 - Bounded waiting: A process should not repeatedly enter CS keeping other processes waiting.

- 4) How to check these 3 in question :-

@ Mutual exclusion :

only possible with single processor \rightarrow

- 1 process in CS, another tries to enter \rightarrow show that 2nd process will get blocked in entry code
- 2 or more processes are in entry code \rightarrow atmost 1 will enter CS.

(b) Progress : (no deadlock)

- No process in CS, P arrives \rightarrow P enters.
- 2 or more processes in entry code \rightarrow show that atleast 1 will enter CS (no deadlock)

(c) Bounded Waiting:

- 1 process in CS, another process waiting to enter \rightarrow show that first process after exiting CS cannot reenter CS; waiting process will get in.

5) Progress \rightarrow Bounded Waiting

Ex: Process P1 repeatedly enters CS.
Process is there but not BW.

Bounded waiting \rightarrow Progress

Ex: All processes are in deadlock in entry section. BW is there but not progress.

6) Progress + Bounded waiting \rightarrow NO starvation

7) Progress & Bounded waiting are independent of each other.

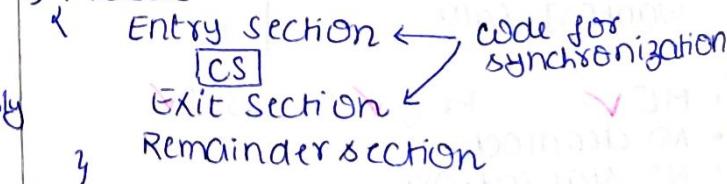
8) Bounded waiting & deadlock are not related to each other.

Deadlock \rightarrow may or may not be BW
BW \rightarrow may or may not be deadlock

Ex: deadlock but no BW

Deny entry of P1, always allow P2, on first entry of P3 deny entry of everyone.

9) Process



10) While(true) \rightarrow While(true); Infinite Loop

11) Preemption can happen anytime.

12) **Progress \rightarrow NO deadlock**
deadlock \rightarrow NO progress

NOTE: If interrupts are disabled, process can't have context switch ie preemption not possible. It's a h/w soln of sync.

⇒ 2-Process Sol^r: (S/W Solutions)

1) Sol^r: Lock variable Sol^r:

LOCK = False

Po

```
while (True)
  while (LOCK);
    LOCK = True;
    CS
  LOCK = False;
```

PI

```
while (True)
  while (LOCK);
    LOCK = TRUE;
    CS
  LOCK = False;
```

process can execute only no. of times.

ME X
Progress ✓
BW X

- Starvation possible: 1 process keeps running again & again, other starves.
- No deadlock
- Busy waiting

2) Sol^r: Turn variable Sol^r:

turn = 0;

Po

```
while (true)
  while (turn != 0);
    CS
  turn = 1;
```

PI

```
while (true)
  while (turn != 1);
    CS
  turn = 0;
```

ME ✓
Pr X
BW ✓

- Starvation
- No deadlock
- Busy waiting
- If initially turn=0, Po may starve

3) Sol^r: Peterson's Sol^r: Perfect sol^r for 2 process

want[2] = {false, false}

turn ;

Po

```
while (true)
  want[0] = True
  turn = 1
  while (want[1] & turn == 1);
    CS
  want[0] = False
```

PI

```
while (true)
  if want[1] = True
    turn = 0
    while (want[0] & turn == 0);
      CS
    want[1] = False
```

BW ✓

- ME ✓
- Pr ✓
- No deadlock
- No starvation
- Busy waiting

* Busy waiting :- A process gets a chance to run on CPU but keeps executing certain stmts only for waiting for CS thus wasting CPU time.

Ex: while(...);

NOTE: Strict alternation means after Po, only PI can execute whether PI waits for CS or not.

⇒ Hardware sol^r for synchronization

• Using hardware atomic inst^r are implemented which is used to provide synchronization.

1) Test And Set (TAS) :-

- This inst^r takes a variable, return its current value & sets it to True. i.e. combining stmt 1 & 2 of Sol^r.

LOCK = False

Po

```
while (True)
  while (TestAndSet (LOCK));
    CS
  LOCK = False
```

PI

same code

- ME ✓
- Progress ✓
- BW X
- Starvation
- No deadlock
- Busy waiting

2) Swap() :-

- TAKES 2 var & swap their value.
- LOCK = False ← Global & shared var
- KEY ← Local var to each process

Po

```
while (True)
  Key = TRUE
  while (Key == TRUE)
    Swap (Key, LOCK);
    CS
  LOCK = False
```

PI

same code

- ME ✓
- Progress ✓
- BW X
- Starvation
- No deadlock
- Busy waiting

⇒ Sync. Tools provided by OS :-

- Semaphore
- Monitor
- Not in syllabus

* Semaphore :

- 1) It is an integer variable which can be accessed by following fun^r only:
 - Wait() / P() / Degrade() / down() : decrement by 1
 - Signal() / V() / Upgrade() / up() : increment by 1

These are atomic functions.

2) Types of semaphore :-

- Binary semaphore : can take value only 0 & 1
- Counting semaphore : can take any integer value (By default)

NOTE: If explicitly given in ques^r $x++$ is atomic then consider else consider it non atomic by default

3) Semaphore with busy waiting:

```

Wait(S)
{
    while(S <= 0);
        S--;
}
    } Signal(S)
    {
        S++;
}

```

NOTE: This implementation is perfect for binary & non-ve semaphore.

- If Binary semaphore $S=0$,
Wait(S) \rightarrow Busy waiting
- If $S=1$, Signal(S) $\rightarrow S$ remains 1

4) Critical Section soln:-

$S=1 \leftarrow$ Binary semaphore

P_0	P_1
while (True)	Same code
{ Wait(S)	
CS	
Signal(S)	

Priority inversion may be possible

- This code will work fine with n processes.
- This is a perfect soln of mutual exclusion

- 5) Binary semaphore used to provide soln to critical section problem.
- Counting semaphore used to control access to a resource with multiple instances.

- 6) Characteristics of semaphores:
- Provide mutual exclusion if implemented properly
 - Soln may have busy waiting, deadlock, starvation, priority inversion.
 - Semaphores are slow soln i.e. slow independent.

NOTE: Deadlock: processes get stuck, can't proceed anyway (infinite waiting)
Starvation: Indefinite waiting.
• Deadlock is a form of starvation

Deadlock \rightarrow Starvation

Q: A process wants to run $S2P() \& 23V()$. We need process to stuck at $25P()$. What should be initial value of S ? Stuck at $25P()$ means it executes $27P()$
 $\therefore S - 27 + 23 = 0 \Rightarrow S = 4$

- 7) To run 2 processes in strict alternation min no. of binary semaphore required = 2
 • For 3 processes to run in specific order $\Rightarrow 3$ binary semaphore required
 8) If semaphore $S=2$ in soln at point (4), max 2 processes can enter in CS together.

9) Priority inversion: high priority process stopped or preempted bcoz of low priority process.

Ex: P1 is in CS, P2 (higher priority) comes, P2 will wait until P1 leaves CS
 P_2 will be CPU, P1 cont. complete CS

10) Semaphore without busy waiting:

```

Wait(S)
{
    if (-S < 0)
        block();
}
    } Signal(S)
    {
        if (+S <= 0)
            unblock();
}

```

• Works perfectly for binary & non-ve semaphore.

• Soln to critical section probm same as in point 4.

• When a process is blocked, it's added to a queue. The order in which processes are blocked, they will execute CS in same order.

- ME ✓ Prop ✓ BW ✓
- No starvation; if we use stack or priority queue then starvation possible.
- No deadlock
- Priority inversion may be possible

⇒ Mutex lock (Spinlock):

- 1) Exactly same as binary semaphore, used to provide mutual exclusion to CS thus preventing race condition

2) A process must acquire lock before entering CS & release lock when it exits CS.

3) Soln to CS problem:

```

while (True)
{
    acquire (LOCK)
    if (LOCK);  $\leftarrow$  Busy
    CS
    release (LOCK)
    LOCK = TRUE;  $\leftarrow$  Busy waiting
}
    } release (LOCK)
    LOCK = FALSE;

```

4) It's called spinlock bcoz process spins while waiting for the lock to become available.

5) If CS takes very less time, in such case spinlock is preferred bcoz moving process to blocked state will take more time.

⇒ Classical problems:

- There are some classical problems which requires synchronization.
 - Producer-consumer problem
 - Reader-writer problem
 - Dining philosopher problem

A) Producer-consumer problem: (Bounded Buffer Problem)

Multiple producers → common buffer → Multiple consumers

- 1) Buffer size is fixed.
- 2) Producer must be blocked if buffer full
Consumer must be blocked if buffer empty
- 3) Variables: (shared variables)
 - Mutex: binary semaphore to take lock on buffer for mutual exclusion.
 - Full: counting semaphore to denote no. of occupied slots in buffer.
 - Empty: counting semaphore to denote no. of empty slots in buffer.
- 4) Initialization: Mutex=1, Full=0, Empty=N

5) Producer()

```

1. Wait(Empty)
  Produce an item
2. Wait(mutex)
  Add to buffer
3. Signal(mutex)
4. Signal(Full)
  
```

6) consumer()

```

1. Wait(Full)
2. Wait(mutex)
  Remove an item from buffer
3. Signal(mutex)
  consume item
4. Signal(Empty)
  
```

Ques asks, if these stmts are swapped what happens → check two conditions:
buffer is full & buffer is empty.

- 7) If stmt 1 & 2 are swapped in Producer, Deadlock possible when buffer is Full.
- 8) If stmt 1 & 2 are swapped in Consumer, Deadlock possible when buffer is Empty.
- 9) If stmt 1 & 2 are swapped in both producer & consumer, deadlock possible in both case when buffer is full or empty.
- 10) If we use 1 binary & 1 counting semaphore then also deadlock is possible

B) Reader-Writer Problem:-

Many readers → common file ← Many writers

- 2) If writer is active then all other readers & writers will be blocked.
- 3) If reader is accessing file, then readers are allowed but writers will be blocked
- 4) Variables: (shared variables)
 - Mutex: Binary semaphore to take lock on readcount variable for mutual exclusion
 - Wrt: Binary semaphore to take lock over file

• readcount: integer variable denotes no. of active readers.

5) Initialization: Mutex=1, Wrt=1, readcount=0

6) Implementation:

Reader()

```

1. Wait(Mutex)
2. readcount++
3. if(readcount == 1)
4.   wait(Wrt)
5. Signal(Mutex)
  
```

Read

```

6. Wait(Mutex)
7. readcount--
8. if(readcount == 0)
9.   Signal(Wrt)
10. Signal(Mutex)
  
```

```

writer()
{
  1. Wait(Wrt)
    write
  2. Signal(Wrt)
}
  
```

last reader unlock file

f) We can not compare semaphore value that's why taken readcount as int variable not semaphore.

g) In Reader(), if we remove line (S&G) or (3&8), only 1 reader will be allowed to read at a time.

h) In Reader(), if stmt 2 is placed before 1, deadlock is possible. (if last reader leaves without executing stmt 8&9)

i) In Reader(), if stmt 3&4 are placed after stmt 5 then it's possible writers can access file while readers are reading.

j) In Reader(), if stmt 8&9 are placed after stmt 10, then deadlock is possible if last reader leaves without unlocking wrt.

C) Dining Philosophers Problem:-

1) By default consider 5 philosophers & 5 forks.



2) Philosopher: process

Chopstick: shared resource

3) K philosophers seated around a circular table. There is 1 chopstick b/w each pair of philosopher.

A philosopher must pick up its two adjacent chopsticks in order to eat. 1 chopstick can be picked up by any 1 of its adjacent philosopher.

4) Variable: $ch_{\{k\}} = \{1, 1, 1, \dots, 1\}$

Binary semaphore array for chopsticks.

5) Solution:

P_i

while (true)

{ 1. wait ($ch_{\{i\}}$)

2. wait ($ch_{\{(i+1) \bmod k\}}$)

eating

3. Signal ($ch_{\{i\}}$)

4. Signal ($ch_{\{(i+1) \bmod k\}}$)

6) Deadlock possible : if all processes execute Stmt 1 One by one.

Deadlock \rightarrow starvation

7) Starvation possible : some philosopher may repeatedly take forks (chopstick) leading other to starve.

8) Ways to avoid deadlock :

• Have no. of philosophers $\leq (K-1)$ chopsticks

• A philosopher should only be allowed to pick their chopstick if both are available at same time. (Prevents hold & wait)

• Break the symmetry :

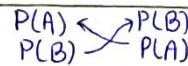
1 philosopher picks first left then right fork, while all others picks first right then left. (For 1 process swap Stmt 1 & 2 in solution.)

OR

Odd philosophers pick first left then right, while even philosophers pick first right then left.

9) Min no. of forks required to avoid deadlock with 5 philosophers = 6

\Rightarrow Deadlock :-



1) Deadlock : If 2 or more processes are waiting for such an event which is not going to occur ever.

2) Deadlock occurs if we do not have good synchronization solution.

3) Necessary conditions for deadlock :-

• Mutual Exclusion : Resource can be used by only 1 process at a time.

• Hold & Wait : holding some resource & waiting for some other resource.

• Circular Wait : processes waiting for resources in circle.

• NO preemption of resource ;

process can't forcibly take resource from some other process.

4) These are necessary conditions for deadlock, not sufficient.

• Deadlock \rightarrow All 4 condition satisfied

• NOT having any condition \rightarrow deadlock will never occur

• Having all 4 conditions \rightarrow deadlock may or may not occur.

5) If all processes needs at most 1 resource each, deadlock cannot occur as hold & wait can't occur.

6) Resource allocation graph :-

• Process P_i \rightarrow Resource R_j

• Resource with multiple instances R_j

• From process to resource :

• From resource to process :

NOTE: Process requests for resource, OS will grant instance if available.

Strategies for dealing with deadlock :

1) Just ignore the problem (Ostrich)

• Pretend that deadlock never occurs in the system.

• Restart the system manually, if system seems to be deadlocked.

• Used by most OS including Unix.

2) Deadlock prevention :

• Prevent any 1 of 4 necessary conditions so deadlock will never occur.

① Preventing Mutual Exclusion :

• All processes should be independent ie no shared resource.

• All process should have their own resources

Both are practically not possible.

② Preventing Hold & Wait :

• Process will either hold or wait but not together.

• If all resources are available then acquire or just wait for all.

• If process needs some more resource then it will release all acquired resources first then wait.

• Problem : Starvation ; may not know required resources at start of run.

- ② Preventing NO preemption: starvation
- Take resource from process forcibly.
 - It is not a viable sol' as process will have to rollback, all progress is lost.

③ Preventing circular wait:

- Give all resources unique number.
- If a process holds resource R_i then it can request R_j only if $j > i$.
- If $j < i$ then process should release R_i & try to acquire R_j first.
- Its practical sol' but may lead to starvation.

3) Deadlock Avoidance:

- requires knowledge of resource requirement option
- We manage resource allocation so that deadlock can never occur.
 - Safe State:** System can allocate resources to process & avoid deadlock
 - Unsafe state:** deadlock may or may not occur.
 - Before granting resource, os checks if resulting state is safe then grant.
 - Algo to avoid deadlock: Banker's Algo.
 - Deadlock avoidance will never leave system in unsafe state.



* **Banker's Safety algo:** To check if the system is in safe state or not.

Process	Allocation			Max			Available			Need		
	A	B	C	A	B	C	A	B	C	A	B	C
P ₀	0	1	0	7	5	3	3	3	2	7	4	3
P ₁	2	0	0	3	2	2				1	2	2
P ₂	3	0	2	9	0	2				6	0	0
P ₃	2	1	1	2	2	2				0	1	1
P ₄	0	0	2	4	3	3				4	3	1

Allocation: present allocation of resources

Max: Max requirement of resources A,B & C

Available: available resources in system

• **Need = Max - Allocation**

- Now check with available resources, anyone's need can be fulfilled or not.

Available

A B C

3 3 2 → + Allocation [P₁]

5 3 2 ← + Allocation [P₃]

7 4 3 → + Allocation [P₃]

P₀ 7 5 3
P₂ 10 5 5
P₄ 10 5 7

Safe Sequence:

<P₁, P₃, P₀, P₂, P₄>

- ④ All process can run in any ~~order~~ sequence then its in safe state
- Total no. of resources = $\sum \text{Allocation} + \text{Available}$
 - If $\text{Available} \geq \text{Max}$ Need then processes can run in any order. no. of process NO. OF POSSIBLE SAFE SEQUENCES = $N!$

* Banker's Resource Request algo:

- To check if the request can be granted.

- If P_i requests 2 instance of A & 3 instance of C ie $R_i = \langle 2, 0, 3 \rangle$

- i) If $R_i \leq \text{Need}$: goto step(ii)
else raise error

- ii) If $R_i \leq \text{Available}$ goto step(iii)
else wait

$$\begin{aligned}\text{Available} &= \text{Available} - \text{Request}_i \\ \text{Allocation}_i &= \text{Allocation}_i + \text{Request}_i \\ \text{Need}_i &= \text{Need}_i - \text{Request}_i\end{aligned}$$

- iii) Run safety algo, if safe then request granted else rejected.

- If 2 requests R_i & R_j are made, possible options :-

- i) both can't be granted

- ii) both can be granted

- iii) only i

- iv) only j

- v) Either i or j, not both

NOTE: If R_i is granted then R_j is granted then both can be granted in any sequence.

NOTE: Deadlock prevention is more strict than avoidance & allows less concurrency.

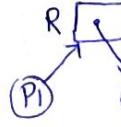
As, prevention may reject request even if resulting state is safe.

4) Deadlock Detection & Recovery

Detection:

- Deadlock detection is done using wait for graph.

- Wait for graph is made from resource allocation graph.



P₁ is waiting for resource held by P₂ ie P₁ is waiting for P₂

(P₁) → (P₂) wait for graph

1) In wait for graph:

- if **no cycle** \Rightarrow **NO deadlock**
- if cycle
 - if only 1 instance per resource type.
 \Rightarrow **Deadlock**
 - if resources have **multiple instances**
 \Rightarrow **Possibility of deadlock**

2) When resources have multiple instances deadlock detection is done using specific form of bankers algo.

In detection we just check given **requests** can be fulfilled or not.

Process	Allocation			Request			Available		
	A	B	C	A	B	C	A	B	C
P ₀	0	1	0	0	0	0	0	0	0
P ₁	2	0	0	2	0	2			
P ₂	3	0	3	0	0	0			
P ₃	2	1	1	1	0	0			
P ₄	0	0	2	0	0	2			

After	Available			If all process can complete	\Rightarrow No deadlock
	A	B	C		
P ₀	0	0	0		
P ₂	3	1	3		
P ₁	5	1	3		
P ₃	7	2	4		
P ₄	7	2	6		

NOTE: When allocation & Max given
 \rightarrow Deadlock avoidance
When allocation & Request given
 \rightarrow Deadlock detection

4) When should detection be done :-
• Do after every resource allocation
• Or do when there is some clue like performance degradation.

Recovery from deadlock:

1) Inform to user & let him handle manually. (Restart system).

2) Process termination :-

- Terminate all processes involved in deadlock.
- Terminate processes one by one until deadlock is broken.
Factors deciding which process to terminate : process priorities, how close it is to finish, how many resources it holds & requires, etc.

3) Resource preemption :-

• Take care of process to rollback & process should not starve.

To find min no. of resources system should have so that deadlock can never occur :-

- Allocate (max) no. of resources to all processes
- Have 1 more available resource so that all process can run.
- Min no. of resource = $(\sum(\text{max-i})) + 1$ so that deadlock can never occur
- Max no. of resources so that deadlock can occur = $\sum(\text{max-i})$

Q: 3 process, 4 instance of a resource, each process request max K instance. Max value of K that will always avoid deadlock?

$$(K-1)3 + 1 = \text{Min no. of instance to avoid deadlock}$$
$$\Rightarrow (K-1)3 + 1 \leq 4 \Rightarrow 3K \leq 6 \Rightarrow K_{\text{max}} = 2$$

\Rightarrow **Types of locks :-**

- 1) Semaphore
- 2) Mutex (Spinlock)
- 3) Deadlock : processes will be blocked (in waiting state). Will not waste CPU time.
- 4) Livelock : system is in livelock when the processes do repeated work without any progress i.e. doing useless work. wastes CPU time.

Ex: P_i has to acquire R, if R is with some other process, kill it & P_i will take R.

Any process can kill another process & acquire resource, this can go in cyclic fashion ensuring livelock.

- **Deadlock & Livelock are mutually exclusive**, at any point of time only one can happen in a system.
- **Deadlock / Livelock \rightarrow No progress**
- **Deadlock / Livelock \rightarrow starvation**

5) **Reentrant lock**: It is a mutual exclusion mechanism (type of mutex) that may be locked multiple times by the same process/thread who releases previous lock & w/o calling while taking lock, count is deadlock incremented by 1 every time. A resource is locked until count returns to zero.

6) Non Reentrant lock: process/thread can't own same lock multiple times.

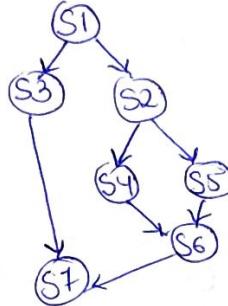
- If a process/thread tries to acquire already owned lock, it will get blocked. It's a **deadlock**.

⇒ PARBEGIN / PARENDD used to specify concurrent statements.

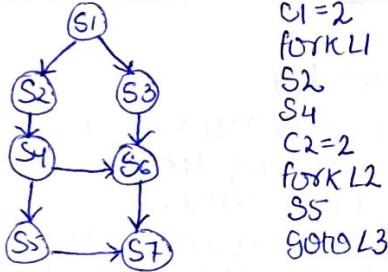
- Also called as Fork/Join, cobegin/coend

Ex: begin

S1
 {
 parbegin
 { S3
 begin
 S2
 {
 begin
 end
 end at
 1 simt
 }
 end
 parbegin
 { S4
 S5
 parend
 S6
 end
 parend
 end
 ctrl



⇒ fork/Join:



S1
C1=2
fork L1
S2
S4
C2=2
fork L2
S5
Goto L3
L1: S3
L2: join C1
S6
L3: join C2
S7
C2 initial value
is 2 means 2
process which
will execute it
will merge at S7

fork L1 means a child process is created from label L1.

* Let n processes & k instances of a resource. Each Pi requires max Si instances of resource. To ensure deadlock does not occur:

$$\sum_{i=1}^n (S_i - 1) + 1 \leq k \Rightarrow \sum_{i=1}^n S_i \leq (k+n)$$

* I/O redirection: can use file for ip, op.
Command <file>: take ip from file
Command >file: put op into file

* Reference Count is used by garbage collector.
Memory location whose reference count becomes 0 is deallocated.

NOTE Monitor is a syn tool, it ensures mutual exclusion but **not deadlock**.

Memory Management

- 1) Functions - Memory allocation, deallocation and protection.
- 2) Goal: min wastage of space, ability to run larger program with limited space (virtual memory).
- 3) Memory management techniques:
 - Contiguous
 - Fixed partition
 - Variable partition
 - internal fragmentation
 - External fragmentation
 - Non contiguous
 - Paging
 - Internal fragmentation
 - External fragmentation
 - Segmentation

NOTE: virtual memory can't be implemented using contig. M. alloc.

- Fixed partition Cont. M. alloc :-
- MM divided into **fixed no. of positions**. Each position can accommodate **exactly one process**.
 - Partition allocation policy:-
First fit, Best fit, Worst fit, Next fit
min internal frag. **max internal fit frag.**
 - Degree of multiprogramming = NO. of partitions in MM.

- Variable partition Cont. M. alloc :-
- Partition is created on arrival of process, of size equal to process.
 - Sol. of external frag.: - compaction.
 - **Worst fit is better for it.** (time consuming)

Non cont. M. alloc :-

Process is divided into partitions & stored at different locations.

@ Paging :-

- Process is divided into equal size partitions called pages. MM is also divided into same size frames.
- Each page can be stored at any frame in MM. This info is maintained in Page table for each process.
- Each Page Table entry = frame no + extra bits
- NO. of entries in PT = NO. of Pages in process
- PT size = NO. of Pages * 1 entry size
- NO. of Pages in process = $\frac{\text{Process size}}{\text{Page size}}$
- NO. of frames in MM = $\frac{\text{MM size}}{\text{Page size}}$

• NO. of bits for offset = $\log_2(\text{Page size})$

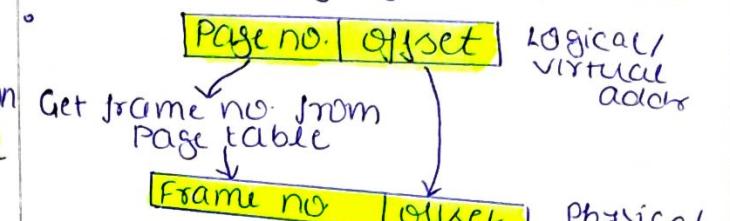
By default Byte addressable → Byte in page

• NO. of bits for page no. = $\log_2(\text{No. of pages})$

• NO. of bits for frame no. = $\log_2(\text{No. of frames})$

• Logical address: Byte no. of process. CPU knows only this not physical address.

• Physical Address: MM addr at which corresponding byte is stored.



• CPU generates logical address. OS converts it into physical.

• To find Page no. entry in page table, skip p entries.

• OS maintains a page table for each process. Page table is also stored in MM in OS area.

• Page table base register (PTBR) stores base address of page table.

• Its present in PCB of process.

• Collection of logical addresses is called **logical address space** or **virtual space**.

• Collection of physical addresses is called **physical address space** or **MM size**.

NOTE: If page table entry is 17 bits & page table size = $2^{20} \times 17$ bits and it uses approx page table size then consider Ans $2^{20} \times 3B$ i.e 3MB NOT 2MB we may have extra bits, not less.

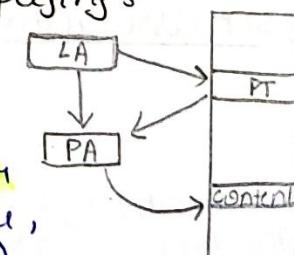
1) Performance of paging :-

- Time required to access content or Effective memory access time = **2 * T_{MM}**

(One for page table, One for content)

- If page table is very small & stored in a register
 $EAT = E T_{reg} + T_{MM} = T_{MM}$

* Translation lookaside buffer (TLB) is used to reduce EAT, by storing frequently used page table entries.



2) CPU demands LA

↓
Search in TLB

Hit

↓
get PA

Miss

↓
Access PT in
MM to get PA

Access MM for content

• TLB is a hardware.

3) EMAT with TLB = $H(T_{LB} + T_{MM}) + (1-H)(T_{LB} + 2 \cdot T_{MM})$
(NOT considering miss penalty) = $T_{LB} + T_{MM} + (1-H)T_{MM}$

4) Hit ratio = $\frac{\text{No. of hit}}{\text{Total no. of reference}}$

5) TLB always have hierarchical access

6) Few entries of current executing process page table is in TLB.

7) TLB Mapping:

a) Direct Mappings :- uses addressable TLB (1 way set associative mapping)

• TLB entry no. = (page no.) % (No. of entries in TLB)

• Tag bits is used to distinguish b/w entries fighting at one slot of TLB.

• LA:

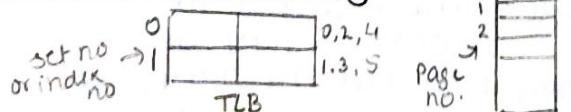
Page Number	Offset
-------------	--------

Tag	TLB entry no.
-----	---------------

• No. of entries in TLB = $\frac{\text{TLB size}}{1 \text{ page table entry size}}$

NOTE: Tag bits is not stored in TLB, its stored in a separate metadata storage.

b) Set Associative mapping :-



• Associativity: no. of PT entries in 1 set

• Tag bits is stored for each entry in TLB.

• TLB set = $\frac{\text{Page no.}}{\text{No. of sets in TLB}}$

• LA:

Page no	Offset
---------	--------

Tag	Set no.
-----	---------

• 4 way set assoc. means 4 entries in 1 set

c) Fully associative mapping:

• All entries in a single set

Page no	Offset
---------	--------

Tag

• Implemented using content addressable memory (Associative memory).

• Its very fast & costly.

• Page no. in LA is parallelly compared to all cells & corresponding page table entry is taken.

accessing same address again ↗

NOTE: Based on temporal locality of reference, PT entry is brought to TLB.

NOTE: No. of bits in logical addr may be more than no. of bits in physical addr in case of demand paging.

1) Paging suffers from internal fragmentation (no external frag)

Int. frag. = $\text{Page size} - (\text{Process size} \% \text{ page size})$

2) Reducing page size may reduce internal fragmentation.

3) Reducing page size increases page table size. (bcz no. of pages increases).

4) In paging, it's possible that unrelated contents (like stack, heap) may be present in same page. So we went for segmentation.

5) Paging incurs extra memory overhead as it takes extra pages in memory for page table.

* By default, address space means logical address space.

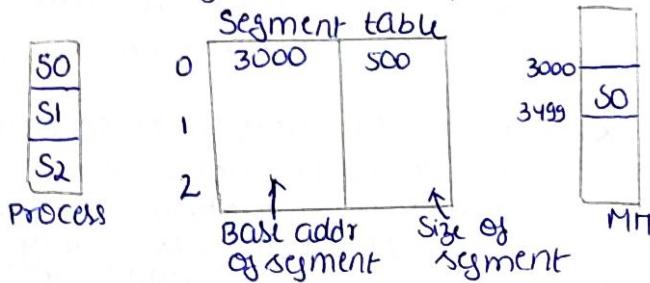
NOTE: If there is no virtual memory concept in system, hardware support for memory management is no longer needed.

In virtual memory system, also (Memory management unit) is required to translate logical address into physical address.

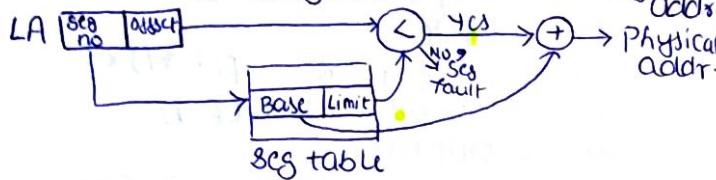
b) Segmentation :-

- Process is divided into logically related partitions called segments.
- Segments may be of variable size.
- MM is not divided, where space is available, segment is kept.

4)



5) Converting logical add. into Physical add.



6) Seg. fault : Trying to access memory beyond limit. Its trap to OS (slow interrupt) \rightarrow CS

7) If limit is 500, it means segment has byte/word 0 to 499.

8) Offset bit size in logical add. is decided based on max segment size.

9) Segmentation suffers from External fragmentation. Use compaction to avoid it.

⇒ virtual Memory &

* Size of virtual memory (VAS) is limited by size of disk, MAR

* In virtual memory system, practically

Physical addr. space \leq Virtual addr. space \leq Disk Addr. space
(MM size) \leq (Process size) \leq (Disk size), int

◦ Conceptually may have any relation

NOTE: 2 level virtual memory means 2 level paging.

Valid bit is 1 for pages that are in MM for rest pages it's 0.

5) CPU requests logical add.
Check valid bit of requested page in page table

Get frame no. & access content from MM

o>
Page not present in MM (Page fault)

6) Page fault is a software interrupt. OS takes control, brings faulted page from secondary memory to MM, updates page table entry. Then the inst. is restarted. It's called page fault service.

7) If empty frame is available, OS keeps new page into it otherwise replaces an existing page.

8) Demand Paging: \rightarrow demand segmentation also possible.

◦ Pure demand paging: initially when process is created, no page is loaded in MM.

◦ Demand paging: initially few pages are loaded into MM.

9) Valid bit is used to ensure page hit/miss.

NOTE: Physical memory is MM (RAM).

10) Page swap time saving: have a extra bit "Dirty or Modified bit" in each page table entry.

◦ Dirty bit 0 means only read is done on page, at time of replacement, overwrite it.

◦ Dirty bit 1 means write is done on page, at time of replacement write back page to secondary memory.

11) A victim page (to be replaced) is selected based on page replacement policy.

NOTE: Max process size = Max virtual memory size = Max secondary memory size.

12) Eff. memory access time

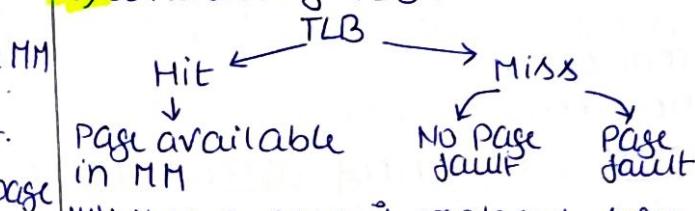
$$= P * \text{page fault service time} + (1-P) * \text{time w/o page fault}$$

$$P: \text{page fault rate} = \frac{\text{No. of page fault}}{\text{No. of memory access}}$$

◦ Time w/o page fault = $2T_{MM}$

$$= T_{MM} \text{ (Page table in registers)}$$

13) Considering TLB :-



14) When a page is replaced from MM & if its entry is present in TLB, its made invalid in page table.

NOTE: If a process requires more heap than allocated, then system call is made & process is reallocated memory.

⇒ Page replacement policies

Used to find victim page.

① First in First Out (FIFO)

- Belady Anomaly: sometimes for some special page reference sequences, increasing no. of frames increases no. of page fault. → Page which is replaced is next requested one
- It happens in FIFO policies which are not based on stack.
- Poor performance i.e. more page faults.

② Optimal policy:

- Page JD jutne me jaldi usenhi hona hai, replace it.
- Use FIFO to break tie.
- not practically possible.
- Min. no. of page fault.
- # Page fault of other policies \geq # Page fault of optimal policy
- Requires future knowledge of program and is time consuming.

③ Least Recently Used (LRU): Practically its best

- Replace page which has not been used since longest period of time.
- Use FIFO to break tie.
- Generally, LRU is better than FIFO.

NOTE: EMAT is for 1 page access.

- Total time for access of x pages

$$= \text{No. of hit} * (2T_{MH}) + (\text{No. of hit}) * \text{Page fault service time}$$
- $\text{EMAT} = \frac{\text{Total time}}{\text{No. of page access}}$

* Counting algo: LFU, MFU

④ Least frequently used (LFU):

- Jiski frequency kam hai past me, replace it.

• Tie breaker: FIFO

⑤ Most frequently used (MFU):

- Jiski freq. zyada in past, replace it.

⑥ Last in first out (LIFO):

- Only last page is replaced always.

⑦ Most recently used (MRU):

- Replace page referred most recently.

NOTE: Locality of reference: page reference made by a process is likely to be one of the page used in last few references.

- Temporal LOR: accessing same address again
- Spatial LOR: accessing nearby address again.

* Let $p \rightarrow$ no. of pages in a sequence
 $f \rightarrow$ no. of frames
 $r \rightarrow$ no. of repetition of sequence
If initially all frames are empty & $f < p$ and all pages are distinct in seq.
Then, no. of page fault in

- FIFO = P_{fif}
- Optimal (if $f \geq r$) = $P_{opt}(x-1) * f$
- LIFO = $(x-1)p_f - (x-1)(f-1)$
- MRU = Optimal
- LIFO - Optimal = $(x-1)$ if $f \geq r$

NOTE: Virtual memory can't be implemented using contiguous memory allocation.

⇒ Frame Allocation: depends on Addressing mode, instruction set architecture

- 1) Every process must have enough pages in MM to complete an inst

2) Types: Min no. of frames

- Equal allocation: Every process is allotted equal no. of frames.
- Proportional allocation: Frames are allotted proportional to the process size.

⇒ Page Replacement: local global

1) Local Replacement:

- Page being replaced be in a frame belonging to same process.
- No. of frames belonging to a process will not change.
- Processors control their own page fault rate, do not affect page fault of other processes.

2) Global Replacement:

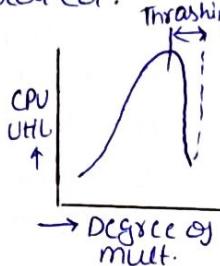
- Can replace a page of itself or other process from a set of all frames allocated to user processes.

- High priority processes can increase their allocation at the expense of lower priority process.
- More efficient use of frames & better throughput.

3) In local replacement, only 1 page table is updated & in global replacement 1 or 2 page table is updated.

⇒ Thrashing:

- 1) With increase in degree of multiprogramming, CPU utilization increases & after a point of time it decreases rapidly.



- 2) Thrashing: when CPU spends more time on page fault service rather than process execution.

- 3) On seeing CPU utilization low, long term scheduler schedules more & more no. of process. CPU is thrashed.

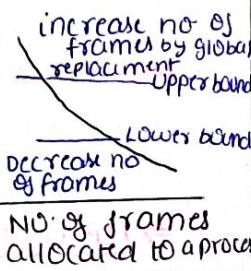
- 4) How to handle thrashing?

@ Working set model:-

- Each process is provided with its working set in MM.
- At present, set of pages process requires frequently is its working set.

⑥ Page fault frequency

- Upper & lower bound of page fault rate is decided & OS tries to keep PFrte of each process within range.



$$* \text{Hit ratio} = \frac{\text{No. of hit}}{\text{Total no. of page reference}}$$

$$* \text{Page no. of process} = \left\lfloor \frac{\text{Address}}{\text{Page size}} \right\rfloor$$

- * Random page replacement policy (may behave like FIFO): selecting a page randomly for replacement. It may suffer from belady anomaly.

⇒ Multilevel Paging:

- 1) If page table can't fit in 1 page, it's also divided into pages & stored randomly on frames in MM.
- We have other page table to store its mapping.

* Helps to reduce page table size. In demand paging, pages that are used PT entries of only those pages that are present in MM.

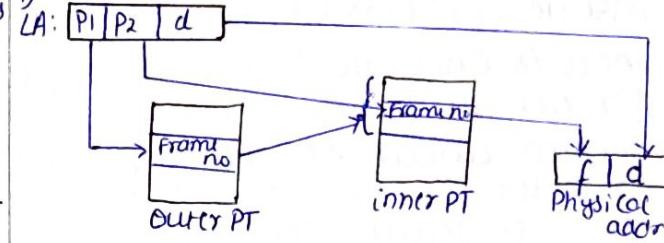
2) LA:

Page no. of Process	Offset (d)
---------------------	------------

Page no. $\xrightarrow{\text{inner PT}}$ $\xrightarrow{\text{which entry of page table}}$
inner

3) This mapping is done until the outer most page table is stored on a single page.

4)



5) PTBR stores address of outer most page table.

6) More page table levels → translation overhead increases.

7) Every PT is stored in memory.

8) Total pages required for page table
 $= \sum \text{all levels}$.

9)

P1 P2 d

 NO. of pages for Outer page table = 1
(always)

NO. of pages in inner page table = 2^{P_1} = 8

10) Page table entry (frame no. + extra bits) will be same across all levels.

11) Total size process takes = process size + (no. of pages * page size)
PT takes

Q: Page size = 2KB, PT entry = 4B
Logical add size = 8GB

P1 P2 P3 d(11)

 :: 3 levels needed
33 bits

NO. of PT entries in 1 page = $\frac{2^P}{2^2} = 2^9$

NO. of pages in outermost PT = 1

" " " " middle PT = $2^{P_1} = 2^4$

" " " " inner PT = $2^4 * 2^9$

Q: PT entry size = 4B, LA size = 16GB
3 level paging, outer PT occupies complete page. Page size = ?
Let page size = 2^P ; NO. of PT entries = $\frac{2^P}{2^2} = 2^{P-2}$

P1 P2 P3 d

P-2 P-2 P-2 P

in 1 page

12)

P1 P2 P3 d

 This is not possible.
4 6 8 All inner PT's will have max entries in each page.
Outer may have \leq max no. of entries.

P1 P2 P3 d

2 8 8

is correct.

13) 3 level page table means 3 times MM access required to get physical addr.

14) Problem with virtual memory: If PT is too big, so many entries is invalid.

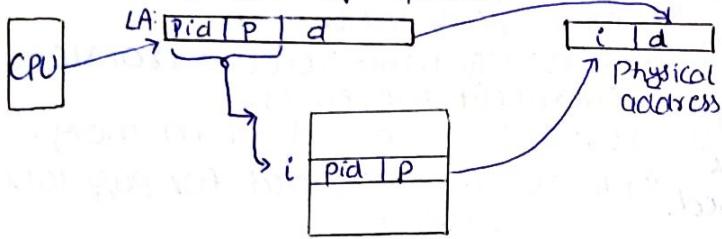
Inverted Page Table:-

1) It has entry for each frame of MM which contains:

Page no., Process Id, Extra bits

2) There is only one inverted pagetable for all processes.

3) We do **Linear searching**, if found its index (frame no.) is used else its page fault. Its time consuming but saves lot of space.



4) If we skip 'i' no. of entries in inverted PT, we are at entry of frame no. i.

5) Difficult to implement page sharing.

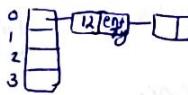
* In simple paging, page sharing is implemented by having some page table entries pointing to same frame in MM.

Hashed PT :-

1) Entry of pages that are present in MM are hashed into PT as in chaining.

Ex: 16 pages in a process

4 hashed location in PT
Slot no = Page no % 4



2) Each process has its own hashed PT.

3) Each node contains:

Page no., PT entry, pointer to next node

* TLB always have hierarchical access.

TLB and Multilevel Paging :-

1) Single level Paging :-

• Default EMAT = $2 \times T_{MM}$

• $EMAT = H(T_{TLB} + T_{MM}) + (1-H)(T_{TLB} + 2 \times T_{MM})$
 $= T_{TLB} + T_{MM} + (1-H)T_{MM}$

H: Hit in TLB TMM: MM access time

T_{TLB}: TLB access time

2) 2 Level Paging :-

• Default EMAT = $3 T_{MM}$

$$\begin{aligned} \bullet EMAT &= H(T_{TLB} + T_{MM}) + (1-H)(T_{TLB} + 2T_{MM}) \\ &= T_{TLB} + T_{MM} + (1-H)2T_{MM} \end{aligned}$$

K-Level Paging :-

• Default EMAT = $(K+1)T_{MM}$

• $EMAT = T_{TLB} + T_{MM} + (1-H)K \times T_{MM}$

$\Rightarrow TLB$ with virtual memory: including page fault



1) Single level :

$$EMAT = H(T_{TLB} + T_{MM}) + (1-H) * [T_{TLB} + T_{MM} + (1-P)T_{MM} + P * \frac{PF \text{ service time}}{PT \text{ lookup}}]$$

$$2) 2\text{-Level Paging :}$$

$$EMAT = H(T_{TLB} + T_{MM}) + (1-H) * [T_{TLB} + 2T_{MM} + (1-P)T_{MM} + P * \frac{\text{service time}}{\text{page fault rate}}]$$

For K level
put $K \times T_{MM}$

Including dirty page :-

1) Single Level :-

$$EMAT = H(T_{TLB} + T_{MM}) + (1-H) * [T_{TLB} + T_{MM} + (1-P)T_{MM} + P \{ d * \frac{\text{service time for dirty page}}{\text{dirty page}} + (1-d) * \frac{\text{service time for non dirty page}}{\text{non dirty page}} \}]$$

d: fraction of dirty page
P: Page fault rate

2) In case of K-level Paging, just change lookup time (real one) to $K \times T_{MM}$.

NOTE: Look for TLB hit \rightarrow Page fault

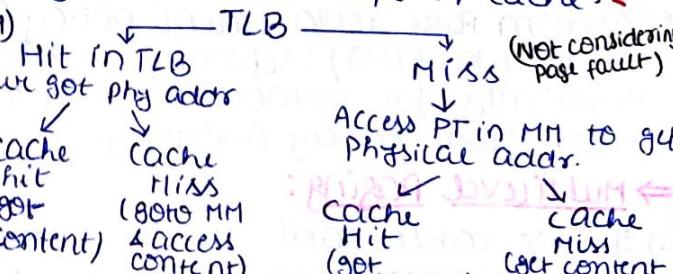
NOTE: If page transfer time b/w MM & disk for dirty page take it 2x. For non dirty page take it x.

TLB with cache :-

1) TLB stores page table entries
Cache stores process content

2) We can access cache from physical address only (not logical address)

3) Page table is not kept in cache.



$$4) EMAT = H[T_{TLB} + H_{cache} * T_{MM} + (1-H_{cache})(T_{MM} + T_{disk})]$$

$$+ (1-H_{TLB})[T_{TLB} + T_{MM} + H_{cache} * T_{MM} + (1-H_{cache})(T_{MM} + T_{disk})]$$

$$= T_{TLB} + T_{MM} + (1-H_{TLB})T_{MM} + (1-H_{cache})T_{MM}$$

* change lookup time (real one) to $K \times T_{MM}$ in case of K-level Paging.

File System Implementation

- 1) We're going to see how we store info/data i.e. files on magnetic disk.
- 2) Addressing unit of disk = unit of atomicity = **BLOCK OR SECTOR**.
Here we will term sector as block.
- 3) Disk reads/writes in terms of block, not bytes.
- 4) File: Data in some format (.c, .exe, .pdf)
- 5) File system: module of OS to manage files

6) File Operations: filename, read, write, etc
 7) Objective of FS: efficient fetching, should not take much time to know where file is stored on disk.

8) Directory: A file with special bit set in metadata to indicate directory.

NOTE: In UNIX everything is file.

NOTE: We are talking with reference to UNIX OS

9) File Metadata: Name, inode no., location, size, protection, time, date, user.

10) Metadata is stored in data structure called **inode** (in Unix) i.e. index node.

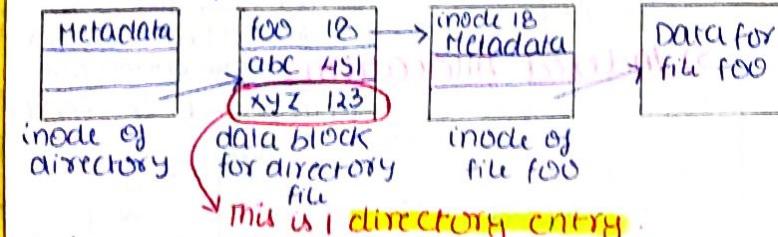
11) For every file or directory \rightarrow 1 inode.
Every inode has unique inode number.

12) File (or directory) = **Metadata + data (blocks)**

13) Inode for a file :-



14) Inode for a directory :-



15) Inode of directory can be implemented like this also :

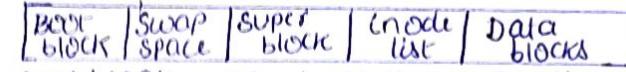
16) File system job is to translate (filename, offset) into block no. i.e. disk address.

17) Inode is used to keep track of sectors allocated to a given file.

File directory structure:

- 1) Single level directory
 - Two files can't have same name
- 2) Two level directory
 - All user names are visible but one user can't access other user data
- 3) Tree structured directory
 - used in modern system
 - If anyone has access to root folder can access everything inside.

→ Disk is divided into some parts:



- Boot-block: OS, bootstrap program
- Super block & # files, free blocks, etc

1) Disk has **magnetic medium** of storage.

2) Disk partitioning / Formatting :-

- Physical (low level): creating tracks & sectors; done by manufacturer
- Logical (high level): creating logical partitions (drives); done by user

3) Types of logical partition:

- Primary (C drive): contains OS + user files
- Extended: contains user files, can't install OS.

4) We may have different OS installed on different primary partition.

NOTE: n blocks are numbered 0 to n-1.

1) Types of file system: FAT32, NTFS, HFS+, EXT2/EXT3/EXT4 (Mac)

2) To open a file, CPU calls OS & OS calls file system to open.

NOTE: Any block can't contain content of 2 different files.

* A directory can't have 2 files of same name & type.

⇒ Disk blocks :-

- 1) Total disk size = No. of blocks * 1 block size
- 2) Disk block address = Disk block no.
- 3) Disk block add. size = $\log_2(\text{no. of blocks})$

⇒ Free space management:

- Which blocks are free, data is maintained in file system.
- 1) **Free List:** maintain linked list of free block no (address).
- 2) **Bitmap:** For each block 1 bit is stored (1: Empty / 0: Occupied)
For n blocks, n bit info
Block $\begin{matrix} 0 \\ 1 \\ 0 \\ 0 \\ 1 \\ \dots \\ 0 \\ 1 \end{matrix}$
- Free list does not require searching of free block, bitmap requires searching.
- Free list is faster
- Free list size is variable.

⇒ File allocation Methods :-

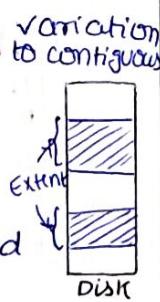
1) Contiguous Allocation:

- File is stored on consecutive blocks.
- Metadata (inode) contains starting block address & no. of blocks.
- Fast sequential & random access.
- Low storage overhead.
- Fragmentation: Internal & External
Here, compaction is known as defragmentation.
It's time consuming process.
- Difficult to grow file

NOTE: All allocation may have internal fragmentation.

2) Extent-based allocation :

- Multiple contiguous regions per file (like segmentation)
- Each region is an extent.
- Metadata: For each extent, start & size of extent is stored in inode.

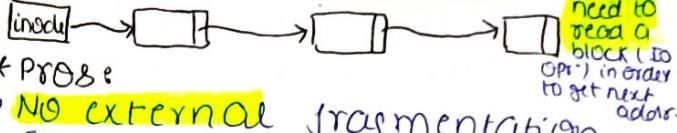


* Pros:

- Low storage overhead
 - File can grow overtime.
 - Sequential & random access.
- * Cons:
- It helps with external fragmentation but still a problem.
 - Internal fragmentation.

3) Linked Allocation:

- All blocks of a file are in linked list. Each block has a pointer to next block.
- Metadata: inode contains pointer to first node & last node.



* Pros:

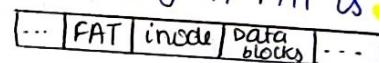
- No external fragmentation.
- Files can easily grow.

* Cons:

- Internal frag.
- Slow sequential access, difficult random access NOT possible.
- Large storage overhead (One pointer per block).

4) FAT Table : (variation to linked)

- Store linked list pointers outside block in file allocation table.
- A section of disk is used to contain this table.
- FAT has 1 entry for each block of disk.
- Used in MSDOS & Windows.
- Still we do pointer chasing but can bring entire FAT in MM so can be cheap compared to disk access.
- Generally, entry in FAT is block no.

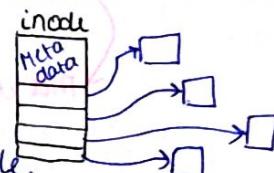


- Slow sequential access, random access possible.

- Large storage overhead for FAT table.
- Assuming all available storage space used to store one file, Max file size = Total disk size - overhead (FAT table)

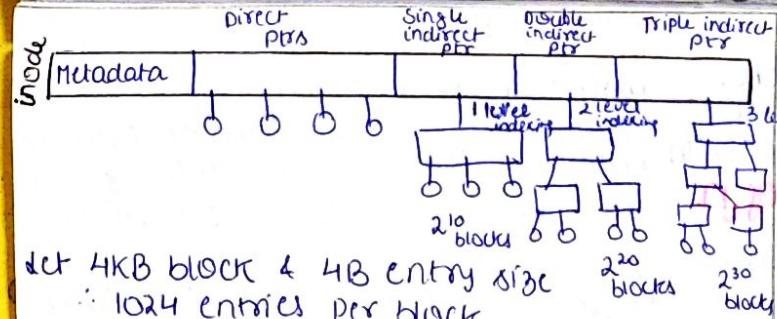
5) Indexed Allocation: → No external frag

- Index table is maintained in inode.
- If inode has space for 4 pointers, it can have atmost 4 sectors.
- Both sequential & random access possible.
- It requires large inode size so that file can grow as required.
- Large inode will have lot of unused entries initially.



6) Multilevel indexed : (variation to indexed)

Unit inode (index node) structure



NOTE: Min file size: atleast 1 block should be given.
Max file size = Disk size - overhead

Config: direct indexed
Block size = 2^{12} B. Each file has 13 direct
ptrs, 4 singly indirect, a doubly indirect
& a triply indirect ptr. No. of inodes
that can fit into a block. Pts size 4B
Pointers size = $(13+4+1+1)4B = 76B$
Assuming other metadata, consider it 128B
(closest \geq power of 2)
No. of inodes that can fit in a block = $\frac{2^{12}}{2^7} = 2^5$

* Performance of sequential access:
contiguous > linked > indexed

* Performance of random access:
contiguous > indexed > linked

→ Consider file consisting of 100 blocks.
Assume index block is in MM. In contiguous allocation, there is no free block to grow at beginning but can grow at end. Also block info to be added is present in MM. calculate no. of disk IO operations required
disk access

	contiguous	linked	Indexed
block is added at beginning	201	1	1
add in middle	101	52	1
add at end	1	3 (assuming last node add present in inode)	1
Remove from beg.	0	1	0
Remove from middle (15th block)	98	52	0
Remove from last	0	100	0

⇒ Disk Scheduling:

1) Requests made by CPU to disk are kept in disk (I/O) queue.

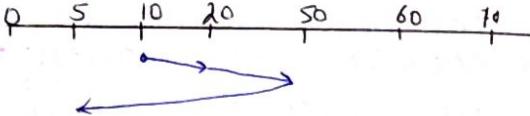
OS need to schedule these requests in order to improve performance of system.

2) To reduce seek time, we store file in cylinder.

Considering request is made for a cylinder. (track no.)

1) FIRST COME FIRST SERVE: (FCFS)

- Requests are served in same order in which it comes.



$$\text{No. of head movements} = (50-10)+(50-5) = 85$$

If for 1 cyl movement, disk head takes 2 ms, Seek time total = $85 \times 2 = 170$ ms

No. of times head changes direction = 1

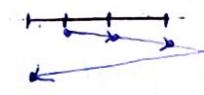
• No starvation

2) SHORTEST SEEK TIME FIRST: (SSTF)

- Fulfill the request which is closest to the current position of the head.
- Throughput increases (Best throughput)
- Average response time decreases.
- can cause starvation.
- High variance of response time

3) SCAN: (Elevator)

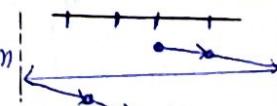
- Go in one direction to end fulfilling all requests, then go in other direction fulfilling requests until all requests are served.
- Always 1 time head direction changes at corner (very last cylinder).
- High throughput
- Low variance of response time.
- Long waiting time for requests just visited by disk arm.



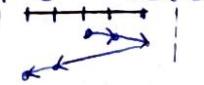
NOTE: By default, arm moves towards larger value.

4) C-SCAN: (Circular scan)

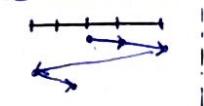
- without fulfilling any request go from 1 end to other ie requests are fulfilled in one direction only.
- Provides more uniform waiting time compared to SCAN.



5) LOOK: Similar to SCAN, don't go to corner, go only upto a valid request.



6) C-LOOK: Similar to CSCAN, we don't go upto corner.

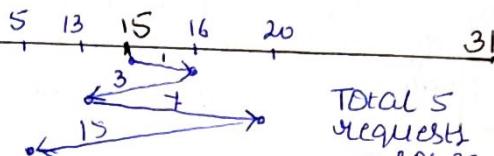


NOTE: In case of tieup continue going in same direction (in SSTF)

* If the total no. of tracks are n & head can start from any track. Max cardinality of request set so that the head changes its direction after servicing every request follows **SSTF** = $\log_2 n$

For max cardinality, start from middle.

Ex: Let $n=32$



$$\text{TOTAL } S \text{ requests} = \log_2 32$$

Use this formula when n is power of 2. Otherwise do manually.

A process may be multithreaded & some threads may generate disk request parallelly.

If OS is capable of handling **single sequential process** at a time, then only 1 request will be live at a time & hence there is no need of disk scheduling algo. We may use any sched. There will be no improvement in I/O performance of process.

Speed of CPU > disk > I/O

Priority of interrupt by CPU > disk > mouse > keyboard

A processor needs slow interrupt to obtain system services which need execution of privileged instructions.

Generally, devices with **high speed** are given higher interrupt priority & slow devices get low priority.

A multiuser / multiprocessing OS can not be implemented on HW that does not support address translation 2 modes of OS.

Overlay: Is used to run a program that is bigger than the size of physical memory (MM) on systems where virtual memory is absent by only keeping required data in memory at given time.

Overlaying: Is used to run a program which is longer than the address space of computer. False

Virtual memory: Is used to accommodate program which is longer than the address space of computer. False because by default address space means logical address space.

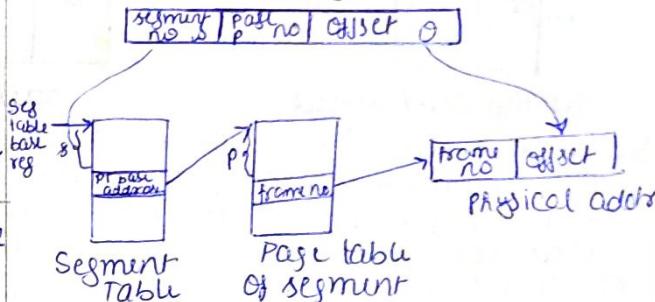
* Interrupt service routine can be written in high level language, also but its performance will be bad.

* Macro definition cannot appear within other macro definition in assembly language program.

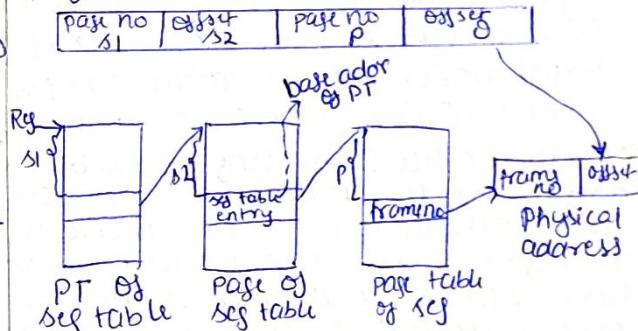
⇒ combination of Paging & Segmentation

Segments are divided into pages, there is page table for each segment.

1) Segmented Paging:

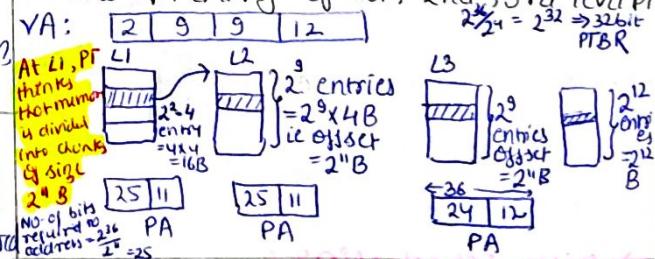


2) Paged Segmentation: Sometimes seg table is too big to keep in 1 block. So, seg table is also divided into pages.



GOB: PA [frame no | offset] → 36 bits VA: [L1 | L2 | L3 | offset] → 32 bits

Page: 4KB ; PT entry: 4B ; 3 level PT
Find no. of bits required for addressing the next level page PT or page frame in the PT entry of 1st, 2nd, 3rd level PT



* Overlay tree of a program:

What will be the size of partition in MM required to load & run this program?

In Overlay, we only need to have the part in MM which is required at that instance of time, either we need root-A-D or root-A-E or root-B-F or root-C-G.

$$\max(12, 14, 10, 14) = 14 \text{ KB}$$

If we have 14KB partition, we can run any of them.