

Date: 7-1-24

Multithreading Models & Thread States in OS.

Multithreaded Models:

- User threads: are managed by the kernel. If supported by the OS, they are managed without kernel support.
- Kernel threads: supported by the OS. A relationship must exist between user and kernel level threads.

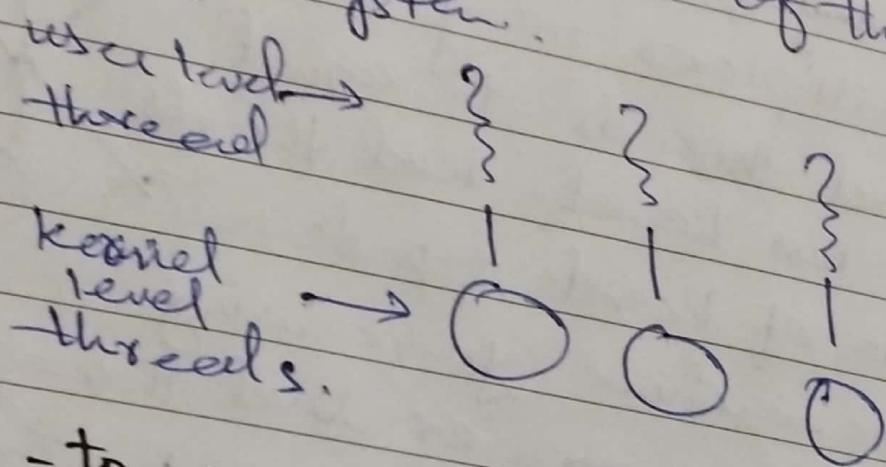
1. One-to-One Model:

- maps each user thread to kernel thread
- when a thread makes a blocking system call, another thread can be run
- Multiple threads can be run parallelly on multiprocessors. → Advantage
- Drawback: Every user thread requires to make a kernel thread. It can burden the performance of an application.

14TH
AUGUSTHAPPY
INDEPENDENCE
DAY

Date:

- there is restriction on number of threads supported by system.



2. Many - to - one Model:

- maps many user level threads to one kernel-level thread.
- Thread management is done ⁱⁿ user space using Thread Library; hence it is efficient → Advantage.
- Drawback: If a thread makes a blocking system call entire process will be blocked.

One thread can access kernel one at a time so no parallel execution on multi processors

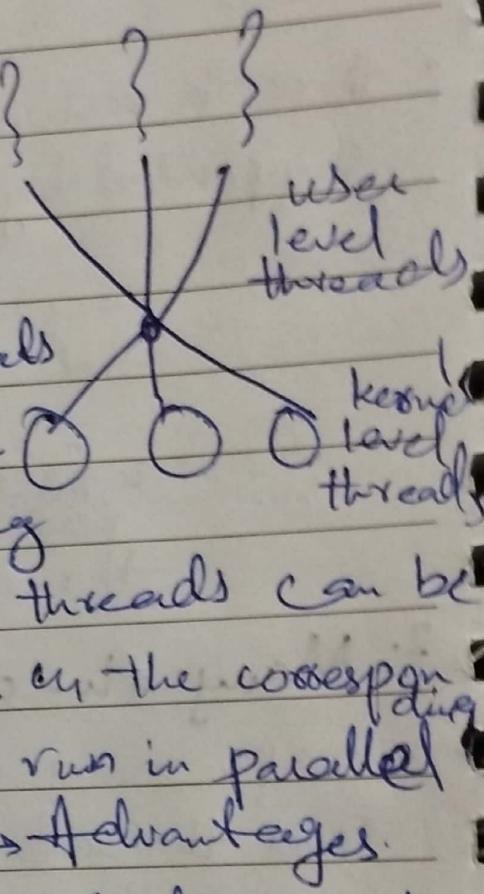
 14TH
AUGUST

 HAPPY
INDEPENDENCE
DAY

Date: _____

3. Many-to-Many Model:

- multiplexes many user-level threads to equal or less kernel-level threads
- No. of kernels depends on the machine you're using
- As many user level threads can be created as necessary. On the corresponding kernel level threads can run in parallel on a multiprocessor → Advantages.
- When a thread makes blocking system call, the kernel can schedule another thread to execute.


 14TH
AUGUST

 HAPPY
INDEPENDENCE
DAY

Date: _____

Relationship Blw Threads & Processes.

Threads : Processes

① 1 : 1

Each thread of execution is a unique process with its own address space & resources. Allows parallel execution.

Ex: UNIX. (Traditional OS).

② M : 1

One process contain multiple threads.

A process has an address space & resources, which every thread shares.

Ex: Modern OS, Linux, Windows.

③ 1 : M

A thread can migrate from one process environment to another through inter-process communication.

In distributed OS, such as in Clouds OS.

Ex Ra(Clouds), Emerald

14TH
AUGUST



HAPPY
INDEPENDENCE
DAY

Date: _____

(a) M:N

- Both 1:M & M:1 approaches.
- Ex: Experimental OS (TRIX),
- Multiple threads are mapped to multiple processes
- Threads in processes are managed independently.

Symmetric Multiprocessing, SMP.

A form of multiprocessor that allows the OS to execute on any available processor or on several available processors simultaneously.

- Kernel can execute on any processor.
- Each process does self-scheduling from the pool of available processes or threads, allowing parallel execution.

A multiprocessor OS must provide all multiprogramming ^{functions} & other features to accommodate multiprocessors.

14TH
AUGUSTHAPPY
INDEPENDENCE
DAY

Date: _____

Multiprocessing OS Design Considerations

① Simultaneous concurrent processes / threads

② Scheduling: Since any processor can do scheduling, so there must be some scheduling policy to avoid. and corruption

③ Synchronization: Synchronization is a facility to enforce mutual exclusion, event ordering for multiple active processes.

④ Memory Management:

Paging mechanism, when different processes share the same page or a segment to enforce data consistency.

⑤ Reliability & fault tolerance, effective process failure handling.

14TH
AUGUSTHAPPY
INDEPENDENCE
DAY

Date: _____

Microkernels:

- A small privileged OS core.
- Provides process scheduling, memory management & interprocess communication services.
- Relies on processes to perform some other functions ~~by OS kernel~~; user-level processes called servers.

Benefits of Microkernel Organization.

- Extensibility: New services can be added or existing ones can be modified without modifying core kernel.
- Flexibility: Components can be added or replaced without affecting core kernel.
- Portability: Change needed to adapt a new processor is done in microkernel.

14TH
AUGUSTHAPPY
INDEPENDENCE
DAY

Date: _____

- Reliability : Modular design can be tested rigorously.
- Distributed System Support . Messages are sent without knowing the targeted machine , maintains isolation.
- Object Oriented OS, Components are objects, interconnected by an interface.

Microkernel Design:

- Low-level memory management : mapping each virtual page to a physical page frame.
- Interprocess communication: Facilitates communication b/w processes.
- I/O & interrupt management.



Date: 7-Jan-24

CONCURRENCY: Mutual Exclusion & Synchronization

The central purpose of OS design is to manage processes & threads.

Multiprogramming: Management of multiple processes within a uniprocessor system.

Multiprocessing: Management of multiple processors within a multiprocessor

Distributed Processing: Management of multiple processes executing on multiple, distributed computer system.

Multiple processes or threads, being executed on distinct systems without knowing the essential details of processing architecture or targeted machines.

Concurrency is the fundamental of all these techniques in OS design.

14TH
AUGUSTHAPPY
INDEPENDENCE
DAY

Date: _____

Concurrency: Ability of an OS to execute multiple set of tasks or processes simultaneously.

Concurrency arises in three different context.

1. Multiple Applications: Multiprogramming, allow processing time to be dynamically shared between no. of active applications
2. Structured Applications: an extension modular design & structured programming. Some applications can be effectively programmed as a set of concurrent processes.
3. Operating System Structure: The same structuring advantages apply to system programs. OS itself is implemented as a set of process / threads.

14TH
AUGUSTHAPPY
INDEPENDENCE
DAY

Date: _____

The basic requirement for support of concurrent processes is the ability to enforce mutual exclusion.

Mutual exclusion is the ability to exclude all other processes from a course of action while one process is granted the ability.

Concurrency & Shared data:

- Concurrent processes may share data to support IPC.
- Threads in a process share same address space (global).
- Concurrently sharing data may cause problems e.g. lost updates.



Date: _____

Key terms related to concurrency:

① **Atomic Operations:** An ~~function~~ indivisible function or action with a sequence of instruction & is not interruptable. The instructions will execute as a group or will not execute. It won't effect system state. Atomicity ensures isolation from concurrent processes.

② **Critical Section:** A section of code in the process, that requires any shared resources. By that must not be executed while another process is in a corresponding section of code.

Simply, a resource is in one process use if it is in critical section. Second process also needs the resource but can't access it until process one completes.

14TH
AUGUSTHAPPY
INDEPENDENCE
DAY

Date:

③ Deadlock: deadlock arises when two processes are waiting for each other to do something, either releasing a resource or generating an output.

④ Livelock: livelock arises when processes changes their states for no useful reason. In response to other processes, without any useful work.

⑤ Mutual Exclusion: It is a requirement, it ensures that when a process is in a critical section using shared resources, no other process can be in a corresponding critical section to access the same resources.

⑥ Race Condition: A situation in which multiple threads / processes read or write a shared data item. The final result depends upon their relative timing of execution. Lastly, which threads executed? (Order of the execution)

14TH
AUGUSTHAPPY
INDEPENDENCE
DAY

Date: _____

- ④ Starvation: A situation in which a process is not being chosen by the dispatcher although it is runnable.

Principle of Concurrency:

In a multiple-processors system, it is possible to interleave the execution of multiple processes and to overlap them.

Both techniques can be viewed as examples of concurrent processing by both present the same problem.

In multiprogramming system, the relative speed of execution of processes can't be predicted.

It depends on other processes activities, OS interrupt handling techniques, by scheduling policies of OS.

14TH
AUGUSTHAPPY
INDEPENDENCE
DAY

Date: Sept 8, 2018

* Speak with honesty, think with sincerity and act with integrity.

Difficulties of Concurrency:

- ① Sharing of global resources. e.g. read or write operations on a global variable.
- ② Manage or allocate resources optimally e.g. a process has a resource or suspend in between. the resource is locked with the process leading to deadlock.
- ③ Difficulty to locate programmatical errors. Because are non-deterministic or unpredictable.

OS Concerns

Design & Management issues raised by the existence of concurrency.

- The OS must be able to keep track of processes (using process control block).
- OS must allocate & deallocate resources to processes.
- OS must protect data on physical resources concerned with each process against other processes.
- Functioning or output of process, must be independent of speed of execution relative to other concurrent processes.

Process Interaction:

They are classified as the degree of awareness upto which they are aware of each other existence.

- Processes unaware of each other:
 - independent process, not attended to work together
 - processes executing independently are in competition of resources. The OS must regulate the access of resources.
 - Timing of process may be affected.

- Potential Control Problems: Mutual exclusion, deadlock, starvation.
- Processes indirectly aware of each other.
 - These processes are not necessarily aware of each other by process IDs, but they share ^{access to} some object, e.g. I/O buffer.
 - They are related as cooperative.
 - Such processes exhibit cooperation in sharing the common object.
 - Results of one process may depend on others as well as time of process.
 - Potential Control Problem: Mutual exclusion, deadlock, starvation, data coherence (a relationship b/w data maintain).
- Processes directly aware of each other.
 - These processes are aware of each other by they can communicate by process ID.
 - They exhibit cooperation by communication.
 - Results of one process may depend on others as well as time of process.
 - Potential Control Problem: Deadlock, Starvation.

Requirements for mutual exclusion:

1. Mutual exclusion must be enforced: only one process at a time is allowed into its critical section, among all processes having similar critical section for the same resource.
2. A process that halts in it is not in its critical section must not interfere other processes.
3. No deadlock or starvation: Process should not be delayed indefinitely, requiring access to critical section.
4. If critical section is free, a process requiring access to it must be permitted to enter in it.
5. No predictions or assumptions are made about relative process speed or no. of processes.
6. A process remain in its critical section for a finite time.

The Critical Section Problem:

- When a process executes, that exploits resource or manipulates shared data, the process is in critical section for that.
- For execution of critical section, mutual exclusion must be enforced, to ensure only one process is

its critical section at a time.

- Mutual exclusion assures , data coherence .

Data coherence: there is a shared ~~data~~^{data} that is in a relationship with other data, they should maintain their coherence after the sequence of processes .

- Deadlock: Process P₁ has Resource R₂ as it is in its critical section . Process P₂ has resource R₁ as it is in its critical section . P₁ is waiting for R₁ as P₂ is waiting for R₂ , both are waiting for each other to leave their critical section , they are deadlocked .
- Starvation: A process is repeatedly denied access to a resource (critical section) . Process P₁, P₂, P₃. P₁ ~~is~~ is in critical section , P₂ or P₃ are waiting P₁ free , P₂ occupies critical section , P₂ free , before P₃ , due to priority P₁ access critical section it is repeating , P₃ is in starvation .

Mutual Exclusion: Hardware Support

1. Interrupt Disabling:

- It works in uniprocessor systems.
- Uniprocessor system: process can interleave but can't overlap (execute concurrently).
- One process at a time will be in critical section, disabling interrupt, no other process can demand critical section for a resource.
- To guarantee mutual exclusion, it is sufficient to disable interrupts.
- Disabling interrupts, interleaving will be reduced.
- It will not work in multiprocessor systems.

2. Simple-Machine Instructions

- In multiprocessor system, processors share same main memory.
- There isn't any interrupt mechanism b/w processors to achieve mutual exclusion.
- At hardware level, access to a memory location, will not allow any other access. The access is blocked for any other instruction.

Compare & Swap Instruction.

- There are two atomic operations, compare & swap.
- First, at a certain memory location, it'll compare with a test value, if it matches, it'll swap ~~update~~ the referenced location ^{with new value} otherwise not.
In both cases, it'll return old value.
- From the example code, where bolt is that memory location initialized to zero, test value is 0 & new value is 1. Then, memory will be updated in it will have access to critical section, if ~~old value~~ is 1, then it will do nothing & wait for its turn to enter in critical section.
- The wait mode is \rightarrow busy waiting or spin waiting.
- A process leaves critical section & resets bolt to 0. The other process which did not get access for CS, was checking continuously & will get access now. but only one at a time.

Exchange Instruction:

- This is another instruction to achieve mutual exclusion.
- bolt is in memory & key is an local variable or register.
- bolt in key value will be exchanged.
- when bolt is zero, process enters in critical section & after completion, it'll set bolt to 0.
- If a ~~process~~^{bolt} is 1, then it'll wait to enter in critical section.

$$\text{bolt} + \sum_i^{} \text{key}_i = n.$$
- bolt = 0 \Rightarrow no process is in critical section.
- bolt = 1 \rightarrow exactly one process is in critical section or the process, whose key value is 0.

Properties of Special Machine Instructions.

Advantages:

- applicable to any number of processes, bolt is used in multiprocessor system sharing main memory.
- simple & easy to verify.
- multiple critical sections can be handled each with their own bolt variable.

Diseadvantages :

- Busy waiting is employed for processes consuming processor's time
 - Starvation is possible, if more than one process are waiting for CS, any process can be denied access.
 - Deadlock is possible: P1 is in CS, P2 with higher priority interrupts P1. P2 asks for resource using by P1, it will not be granted due to mutual exclusion.
- 27

SEMAPHORES:

- Semaphore is an integer variable initialized with a non-negative value.
- It is used for signaling a process for cooperation b/w them
- Two operations.
 - ① decrement semaphore, SemWait operation, if value becomes -ve, then the process will be blocked.
 - ② increment semaphore, semSignal operation,



Date _____

after incrementing, if $\text{value} \leq 0$, then ^{the} blocked process due to `semWait` will be unblocked.

- To begin, it can initialize with ≥ 0 value.
- value \rightarrow +ve integer, shows that no. of processes that can issue a wait & immediately continue to execute.
- If the semaphore value = 0, (by initializing to 0 or by. no. of processes = initial semaphore values issued wait), then the next process to issue wait will be blocked. as Semaphore \Rightarrow -ve.
- -ve value of semaphore = no. of processes waiting to be unblocked.
- When semaphore is -ve, each signal (`semSignal`) will unblock one process.
- We can't say it before, that the decrement will block a process.
- After ~~an increment~~, a process is woken up in ~~processes~~ they both run concurrently.

Date _____

- Signaling Semaphore doesn't ensure a process is waiting, unblocked processes can be one / zero.

Binary Semaphore:

- initialized by 0 or 1.
- SemWait operation checks if $s = 0$, then the process will be blocked. If $s = 1$, the value is changed to 0 by process executes.
- SemSignal operation checks if there are blocked processes ($s \leq 0$). If yes, then process is unblocked in $s = 1$.

Non-binary Semaphore is also called Counting Semaphore or General Semaphore

Mutex:

- Like binary Semaphore
- difference is, a process that locked the mutex (set value to 0) must be the one to unlock it (set value to 1).

- A queue is used to maintain blocked processes for both counting by binary semaphores.
- It uses FIFO, process blocked first will be unblock first.
- Strong Semaphore : is the one that uses FIFO policy.
- Weak Semaphore , is the one that does not specify any order, by which processes are unblocked .
- Strong semaphores guarantee freedom from starvation as it is convenient while weak doesn't (isn't).

Producer / Consumer Problem:

semWait (Semaphore s) {

 s = s - 1;

 if (s < 0)

 add process in blocked list

 else

 return;

semSignal (Semaphore s) {

 s = s + 1;

 if (s ≤ 0)

 wake up process in ready queue

 else

 return;

Date _____

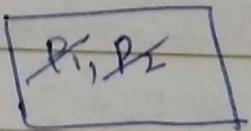
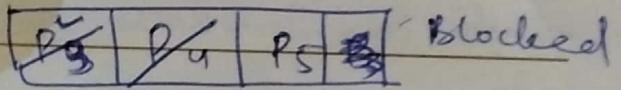
If s is a +ve integer > 0 , then we can say
~~that~~ if s is the no. of processes that can enter into CS.

If s is 0, then no next process can execute in Critical section as will be blocked.

If s is -ve integer < 0 , then it shows the total number of blocked processes.

Let $s = 2 \times X \otimes + 2 \times S$

semWait $\rightarrow 2 \times S - 2 - 1$



P₃

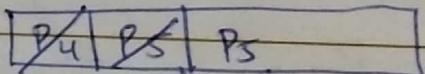
P₃, P₄

semSignal.

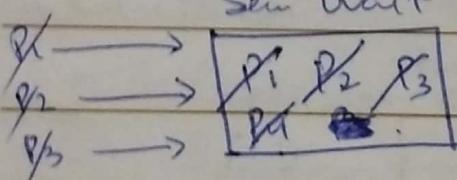
$$B - 3 + \frac{3}{2} \quad \frac{3 \cdot W}{2S} \quad 3u - 2d$$

Let $s = 3 + 2 \times X \otimes + 1 \times S$

semWait $\rightarrow X \otimes + X \otimes - 1$

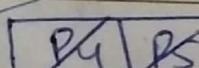


B



P₄ ←

P₅ ←



R.

semSignal.

Pos Bin any Semaphore.

SemWait (Semaphore S) {

if ($s == 1$) .

$s = 0$;

else {

Block this process

} :

}

SemSignal (Semaphore S) {

if & is blocked list
is empty

$s = 1$;

else {

wake up a process
from blocked list

}

$s = X 0$

↑
Success

enter Critical
Section

$s = 0$

unsuccessful
Block the
process

$s = \emptyset 1 \leftarrow$ (if block list)
is empty

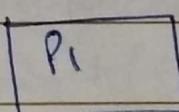
$s = X 1 \leftarrow$ (no change)

if $s = 0$ is initiali
zed
to 0, all process
will be blocked,
deadlocked

P₁,

SemWait

cs



Sem Signal

Let $s = \emptyset \leftarrow$

Let $s = X \emptyset X 0$

P₂,

SemWait

cs

Sem Signal.

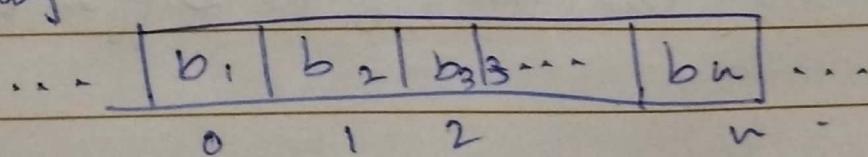
Blocked

Ready

Date 15-1-24.

The Producer / Consumer Problem.

- A common problem faced in concurrent processes.
- Statement: "There are one or more producers generating some type of data (records, characters) and placing these in a buffer. There is a single consumer that is taking items out of the buffer, one at a time. The system is constrained to ~~be~~ prevent the overlap of buffer operations. That is, only one agent (producer or consumer) is allowed to access buffer at any one time".
- Problem: A producer won't try to add data if it is full and consumer won't try to remove item from an empty buffer.
- Producer process to add data into a buffer
- Consumer process to remove data from a buffer
- Buffer is shared between both processes.
- They are cooperative.
- Let's assume a buffer is an infinite linear array





Date _____

Producer Function

produce() {

 while (true) {

 // produce item v

 b[in] = v;

 in++;

}

Consumer function

consume() {

 while (true) {

 while (in <= out)

 // do nothing

 w = b[out];

 out++;

 // consume item w.

 ↑
 consume first make
 sure, it does not read
 from an empty buffer.
 → It will consume
 only if $in > out$

↓
producer generates an
item v, store it into a
buffer, and increment an
index → in

Implementing the system using Binary Semaphore.

- Keep track the number of items in the buffer using variable n ($n = in - out$)
- Semaphore s is used to force mutual exclusion
- Semaphore delay is used to make consumer wait using semWait, if the buffer is empty.

semWait B(Semaphores) {
 if ($S == 1$)
 $S = 0$;
 else
 Block
 Date: 3
 int n; $n = 0$.

Sem Signal B (Semaphore S) }
 if (block list is empty)
 $S = 1$
 else wake up
 }.



binary_semaphore & $S = 1$, delay = 0;

producer() {

while (true) {

produce();

semWait B(S);

append();

$n++$;

if ($n \geq 1$)

semSignal B(delay);

semSignal (S);

consumer() {

semWait B(delay);

while (true) {

semWait B(S);

take();

$n--$;

semSignal B(S);

consume();

if ($n = 0$)

semWait B(delay);

Let's run this.

S. n delay. Semaphore

Blocked

X 0 0 1 SemWait B(S).

[C]

semSignal B(delay)

S delay n Semaphore processes.

X 0 1 0 1 wait(S), signal(D), signal(S) P

X 0 1 0 0 0 wait(D), wait(S), signal(S), wait(D) C

(X) 0 1 0 1 0 wait(S), signal(D), signal(S) P

X 0 0 X 0 -1 ^{Prev C/last} wait(D), wait(S), signal(S), wait(D), wait(S), n-, signal(S)

These is no data, producer is generating, consumer is trying to get from empty



```
* SemWaitB(Semaphore S) {  
    if (S == 1)  
        S = 0  
    else  
        block
```

```
SemSignalB(Semaphore S) {  
    if block list is empty  
        S = 1  
    else  
        wake up.
```

value of n has reached -1, means, no element in the buffer, consumer is trying to get. It leads to deadlock condition.

A helping variable can be made in consumer in critical section.

Producer

Consumer

```
int n = 0;
```

binSemaphore S = 1, delay = 0

producer() {

consumer() {

while(true) {

int m;

produce();

SemWaitB(delay);

SemWaitB(S);

while(true) {

append();

SemWait(S);

$n++$;

take();

if ($n == 1$)

$n = -1$;

SemSignalB(delay);

$m = n$;

SemSignalB(S);

SemSignalB(S);

consume();

if ($m == 0$)

SemWaitB(

delay);

delay);

Date _____

Blocked

~~AF~~

S d n

Semaphore processes.

X₁ Ø₁ Ø₁

wait(S), n++, signal(D), signal(S) P

X₁ X₀ X₀

m=0 wait(D), wait(S), signal(S) Block C.

X₁ Ø₁ Ø₁

wait(S), n++, signal(D), signal(S) P

X₁ X₀ X₀m=0 wait(D), wait(S), n--, signal(S),
m=20 wait(D), wait(S), n--, m=-1, sig(S)X₁ X₀ X₀to wait(D), wait(S), n--, m=0, signal(S)
, wait(D) Block CX₁ Ø Ø₁

wait(S), n++, signal(delay), signal(S)

X₁ Ø X₀ X₀wait(D), wait(S), n--, m=0,
signal(S), wait(D),

now, consumer won't get any data from buffer, until the producer has produced any data.

Solution using General (Counting) Semaphore:

- Here, n is a semaphore now, but its value is the number of data items in buffer.
- no delay semaphore here.



Semaphore $n=0, s=1$; Date _____

producer() {

while(true) {

produce();

semWait(s);

append(); \exists CS.

semSignal(s);

semSignal(n);

}

consumer() {

~~semWait(n);~~

semWait(s);

take(); \exists CS.

semSignal(s);

consume();

}

① Suppose, semSignal(s) or semSignal(n) are interchanged.

- This makes semSignal(s) to be done in producer's critical section, with no interface by consumer or other producer.
- This would not affect program, because consumer must wait on both n & s semaphores.

② Suppose, semWait(n) & semWait(s) operations are interchanged.

- This is not acceptable, since semWait(n) ensures semaphore = 1, in this case, there is a data in buffer, then

Date _____

consumer should after its critical section.

- But if semWait(s) implements before consumer allow the critical section to run, then if $n=0$, it will be the condition that consumer is trying to get data from an empty buffer.
- The system is deadlocked.
- So, the operations should be implemented in the correct sequence to avoid deadlock by enforce mutual exclusion, also avoiding race condition (?).

Solution using infinite buffer, buffer is circular.

Producer:

```
while(true) {
```

// produce an item v.

```
    while ((in+1)%n == out)
```

// do nothing

```
b[in] = v
```

```
    in = (in+1) % n
```

Consumer:

```
while(true) {
```

```
    while(in == out)
```

// do nothing

```
w = b[out]
```

```
    out = (out+1) % n
```

// consume item w.



SemWait(semaphore s) {

s--;
if (s.count < 0) {
 block
}

SemSignal (semaphore s) {

s.count++;
if (s.count <= 0) {

} wake up

const int sizeOfBuffer = _____ Date - let e = 10

semaphore n=0, s=1, e = size of Buffer

producer () {

while (true) {

 producer();

 SemWait (s);

 append();

 SemSignal (s);

 SemSignal (n);

}

consumer () {

 while (true) {

 SemWait (n);

 SemWait (s);

 take();

 SemSignal (s);

 SemSignal (e);

 consume();

}

- This implementation uses counting semaphores

- The process will block, if producer inserts in full buffer or consumer tries to get data from an empty buffer.

- It will unblock on consumer \rightarrow item inserted in producer \rightarrow item removed.

S e n

X _{P1} B ₂ D ₁	wait(e), wait(s), signal(s), signal(n)	P
X _{P1} B ₃ X _D	wait(n), wait(s), signal(s), signal(e)	C
1 3 D-1	wait(n), ^{block} ,	C
X _{P1} B ₂ -X _D	wait(e), wait(s), signal(s), signal(n)	P
X _{P1} B ₃ D	wait(s), signal(s), signal(e) wake up	C

Date _____



Total = 9. in out
n = Filled = 3. 8 0
n 4 1

0	a
1	b
2	c
3	item 1
4	
5	
6	
7	
8	

S	n	e	
1	03	6	
01	4	8	Producer
01	3	6	Consumer

no presumption
example.

Total = 8. in out

n = 3 8 4 0 1

0	a
1	b
2	c
3	d
4	
5	
6	
7	

S	n	e	
1	3	8	
0		4	Producer
01	2	5	Consumer
01	3	5	Producer



Date _____

Monitors:

- Monitor is a programming language construct that provides equivalent functionality to that of semaphores by that is easier to control.
- The construct have been implemented in multiple programming languages e.g. Java by have also been implemented as a program library.
- Monitor lock, can be placed on any object.
- E.g. all linked lists can be locked with a lock or every linked list can have multiple locks for each or every node in the list can have a lock.

Monitor & with Signal:

- A software module with one or more procedures, an initialization sequence & local data.

Characteristics:

- local variables are accessible only by monitor's procedure. they're private.
- process enters monitor by invoking one of its procedures.

Date _____

- Only one process may be executing in monitors at a time. Others trying to invoke monitors are blocked.
- The first two characteristics resembles object oriented design. Indeed an OO. OS / PL can implement monitors as an object with special characteristics.
- Enforcing ~~suspending~~ one process at a time makes monitors to ~~enforce~~ mutual exclusion.
- For concurrent processing, monitors must ensure synchronization.

Condition Variables:

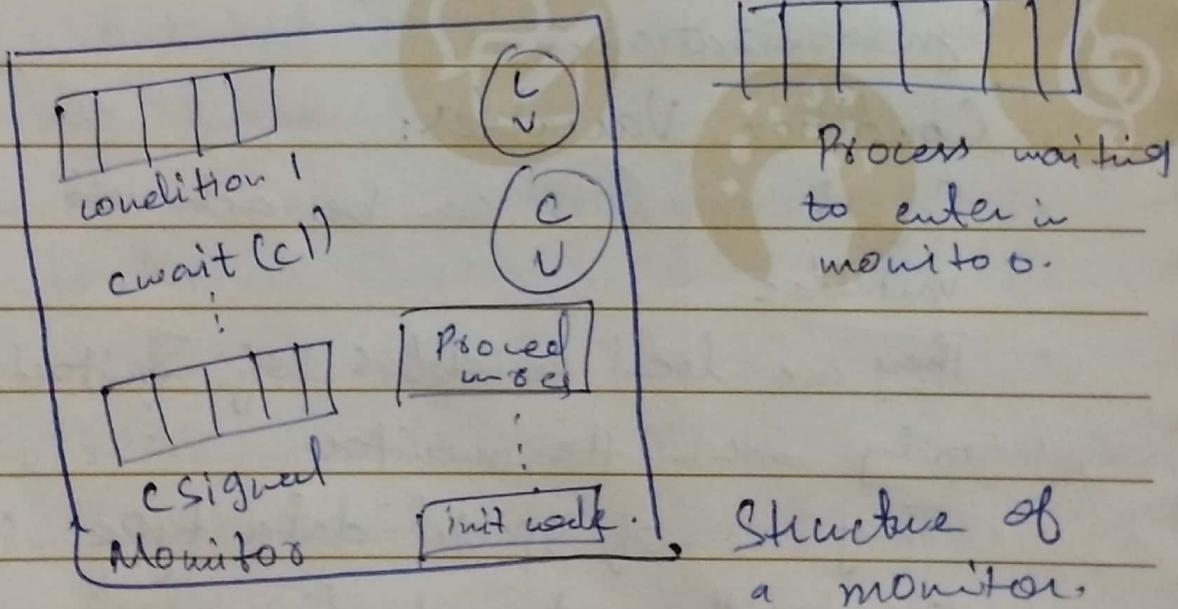
- Synchronization can be achieved using condition variables.
- They are local variables of monitor. Is accessible only within the monitor.
- They are of special data type & are operated on by these two functions.

→ cwait(c): Suspend execution of calling process on condition 'c'. The monitor ^{is} now available to use by another process.

`csignal(c)`: resume the execution of some process blocked after a cwait on the same condition.

If there are several such processes, choose one of them; if there is no such process, do nothing.

- The wait or signal operations are different than semaphore operations.
- The signal is lost, if a process in monitor signals & no task is waiting on the condition.



- A process enters in monitor by invoking one of its procedures
- One process enters in monitor at a time in next

will be blocked or placed in a queue

- A process in the monitor may be temporarily blocked ^{enter in condition queue} on some condition 'c' by issuing `cwait(c)`. Then it waits for 'c' to change to ready.
- `csignal(c)` will alert that condition has changed

Message Passing:

- Synchronization and Communication are two basic requirements when processes interact.
- Synchronization to enforce mutual exclusion.
- Communication to exchange information as the processes are cooperative.
- Message passing approach can do both.
- This can be implemented in distributed systems, shared-memory multiprocessor & uniprocessor systems.
- The function of message passing includes two primitives

- ① `Send (destination, message)`
- ② `receive (source, message)`.

- A process sends information using send primitive, message is information, ^{destination} ~~destination~~ of another process.

- A process receives information using receive primitive message → information, source of of the process.

General Message Format:

Header {	message type	→ Type to disseminate b/w different messages
	destination id	→ the targeted processes
	source id	→ the current processes
	message length	→ number of characters
	control information	→ e.g. pointer of a linked list of messages; a priority field
Body {	Message contents	→ actual message. "Hello world".

Date 17 - 1 - 24

DEADLOCKS

- A set of blocked processes, each holding a resource, and waiting for other processes to acquire another resource's. The processes are deadlocked.
- It is a permanent blocking of a set of processes that either compete for system resources or communicate with each other.
- Example # 1
Two processes P_1 and P_2 . A system has two disk drives. P_1 and P_2 respectively holding one of disks. Each need other drive. Deadlocked :)
- Example # 2

Semaphore 'a' and 'b' \rightarrow initialized to 1.

P_0

P_1

semWait(a) \rightarrow process after CS $a=0$ semWait(b) \rightarrow Process enters CS for b $b=0$
 semWait(b) \rightarrow $b=0$, blocked semWait(a) \rightarrow $a=0$, blocked

Deadlocked :)

System Model:

- A system consists of finite no. of resources to be distributed among no. of competing processes.
- Resources include Memory space, CPU cycles, files & I/O devices
- If a system has two processors, then ~~the~~ each resource type has two instances
- A process must request a resource before using it & must release the resource after using it.
- A process can request multiple resources.
- → The sequence
 - Request: Process requests resource. If not available, the process must wait until grants access.
 - Use: Process can operate / utilize / use the resource.
 - Release: Process releases the resource to be used by another process requesting for it.
- The request & release are system calls.

Deadlock Characterization:

There are some conditions which ensures deadlocks.

Necessary Conditions:

1. Mutual Exclusion: Only one process can use a resource at a time by other processes requesting it are blocked.
 2. Hold & wait: A process ^{must} hold at least one resource and waiting for the assignment of another one.
 3. No preemption: A process can't be preempted. A resource can't be removed from a process due to priority of another process.
 4. Circular waits: A closed chain of processes exist, such that each process holds at least one resource needed by the next process in chain. (necessary & sufficient condition)

Set of processes $\{P_0, P_1, \dots, P_n\}$.

 $P_0 \xrightarrow{\text{wait}} P_1, P_1 \xrightarrow{\text{wait}} P_2, \dots, P_n \xrightarrow{\text{wait}} P_0$.

This condition implies hold & wait condition.
- Mutual Exclusion, Hold & wait and no preemption tells possibility of a deadlock.
 - Adding circular wait tells the existence of deadlock.

Resource Allocation Graph:

- Directed Graph, deadlock can be described using graphs.
 - Also called System Resource-Allocation Graphs
 - Set of vertices called nodes divided in two Process nodes and Request nodes.
- $\{P_1, P_2, \dots, P_n\}$ $\{R_1, R_2, \dots, R_n\}$
- A directed edge from process to a resource $P_i \rightarrow R_j$.

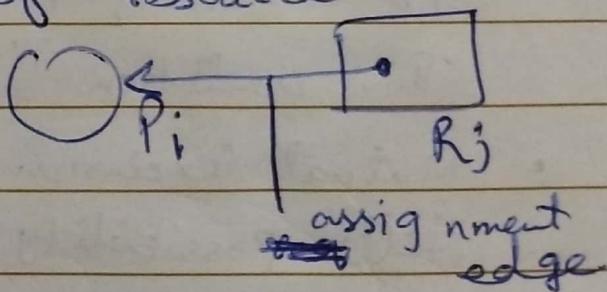
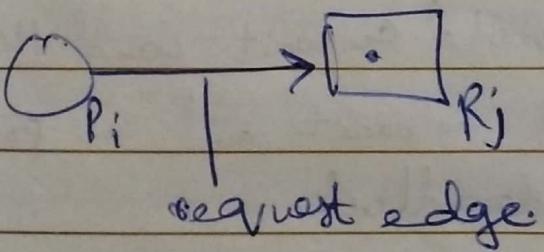


It signifies, a process P_i has requested an instance of resource R_j . i is currently waiting

- A directed edge from resource to process $R_j \rightarrow P_i$ or $P_i \leftarrow R_j$

It signifies that an instance of resource R_j is assigned/allocated to process P_i

- represents instance of resource

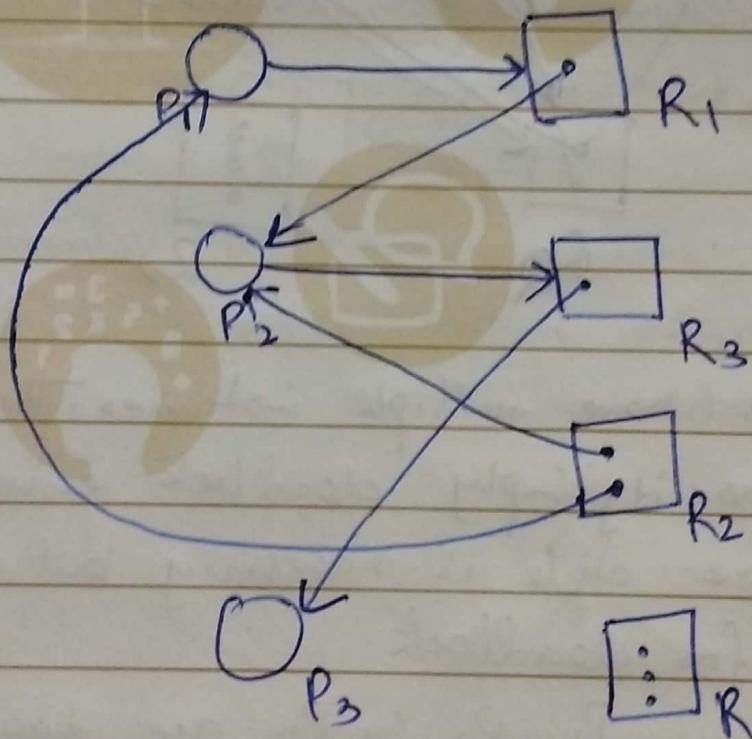


Date _____

- $P = \{P_1, P_2, P_3\}$.
- $R = \{R_1, R_2, R_3, R_4\}$.
- $E = \{P_1 \rightarrow R_1, P_2 \rightarrow R_3, R_1 \rightarrow P_2, R_2 \rightarrow P_2, R_2 \rightarrow P_3, R_3 \rightarrow P_3\}$.

R_1 (1 instance), R_2 (2 instances)

R_3 (1 instance), R_4 (3 instances)

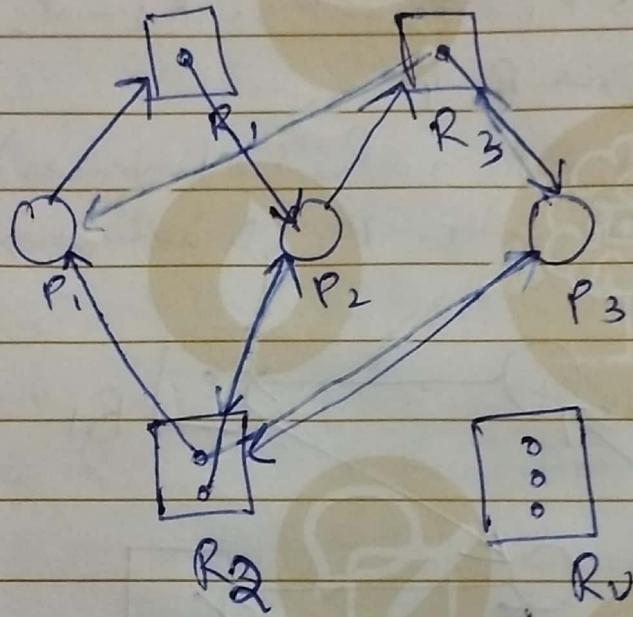


- If the graph contains a cycle, deadlock may exist.
- If there is a cycle and each resource has single instance, then deadlock exists.

cycle is a necessary & sufficient condition for deadlock occurrence

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

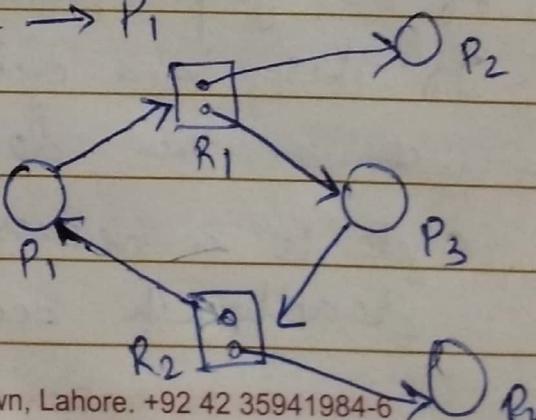
$P_2 \rightarrow R_3 \rightarrow P_3 \rightarrow R_2 \rightarrow P_2$.



- If resources have multiple instances, then cycle doesn't necessarily imply deadlock existence.
- In this case cycle is necessary but not sufficient condition for deadlock.
- Here processes P_1, P_2 & P_3 are deadlocked

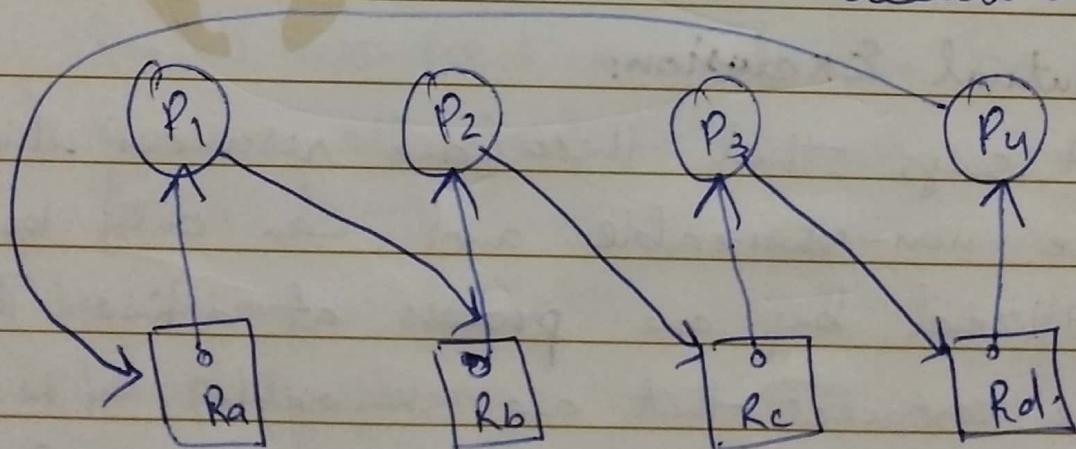
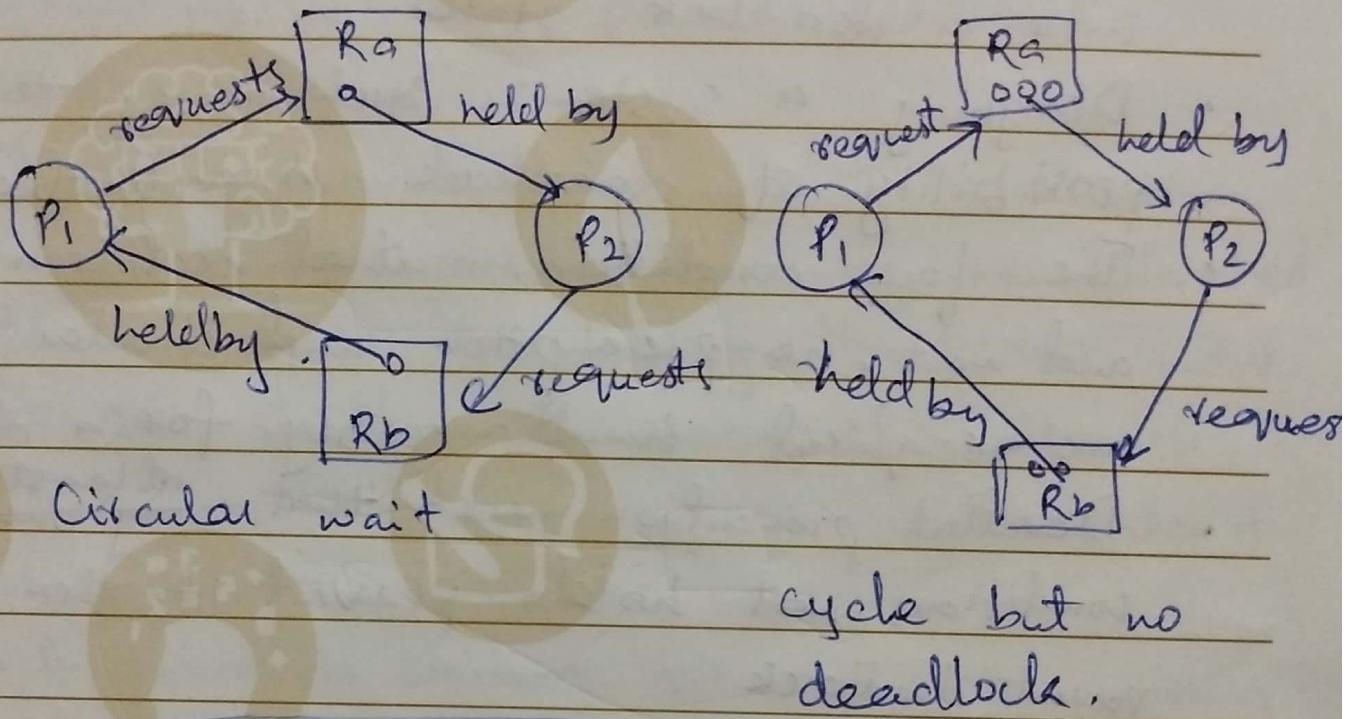
$P_1 \rightarrow R_1 \rightarrow P_3 \rightarrow R_2 \rightarrow P_1$

- No deadlock, as soon as P_1 & P_2 release resource R_2 & R_1 , they be used by P_1 & P_3 .



Date _____

- If the graph has no cycle, no deadlock.
- If the graph contains a cycle, may or may not be deadlock





Date _____

Methods for Handling Deadlocks.

① Deadlock Prevention:

- It ensures that the system will never enter into a deadlock state.
- Designing a system in such a way, that the possibility of deadlock is excluded.
- The four conditions, mutual exclusion, hold and wait, no-preemption and circular wait are required simultaneously for a deadlock.
- Deadlock prevention work that ~~any one~~ ^{at least} of the conditions not hold, prevent it, then we can prevent deadlock.

1. Mutual Exclusion:

It says that there are resources which are non-shareable and can only be utilized by one process at a time. There, are resources that are shareable & hence don't require mutual exclusive access & thus can't be involved in a deadlock.

Pointer is a non-shareable resource.

Date _____

Read-only files are a shareable resource.

→ Conclusively, we can say that it is hard to prevent mutual exclusiveness due to non-shareable resources, hence can't prevent deadlocks

2. Hold & Wait:

- It can be prevented by starting the execution after every required resource has been allocated to the process or blocking the process until it gets simultaneous access for all.
- There are two disadvantages which makes it inefficient

- A process is waiting for a long time, waiting for all resource's requests to be fulfilled
- Resources allocated to a process, may remain unused for a long time.

Low resource utilization, starvation is possible.

3. No preemption:

- A process P_i holding some resources by request some other resources, if other resources can't be accessed, process P_i preempt all its current

resources and request them again together with additional resources.

- If a process P_1 request a resource, which is held by process P_2 , then OS may preempt resources from P_2 .
- This approach is practical only if the resources state can be saved & restored easily.

4. Circular Wait:

- It can be prevented by defining total ordering / linear ordering of resources types.

Let tape drive = 1

disk drive = 5

pointer = 12

a process can only request next resources in ascending order.

If process P_1 has asked for disk drive which is 5, then it can't access tape drive because.

$1 \not\leq 5$ but can access pointer because, $it \not\leq 5 \leq 12$

② Deadlock Avoidance:

- With deadlock avoidance, a decision is dynamically made whether the current resource allocation request, if granted, will lead to a deadlock.
- It requires future process resource requests.
- The system will know the complete sequence of requests; it can decide whether or not a process should wait in order to avoid deadlock.
- To implement this, the system will have information of available resources, resources allocated to each process by future requests.
- The simplest & useful model requires that each process declare the maximum no. of resources of each type, it may need.
- Deadlock avoidance algorithm dynamically examines resource-allocation states, to ensure that a circular wait condition can never exist.
- Resource allocation state: no. of available & allocated resources by max. no. of required resources.

Safe state:

- System can allocate resources in some order in no deadlock is encountered.
- Safe sequence \rightarrow System is in safe state
- P_i requesting resources & gets them, as P_j is not affected. It is safe sequence.
- If P_i requires a resource which is not available, then it can wait, until P_j has finished, then it will take them & continues its execution & terminates, then P_{i+1} continues & so on.
- If there is no such sequence, then it is unsafe.
- An unsafe state may lead to a deadlock.
- We have to avoid the system to enter in an unsafe state.

Date _____

③ Deadlock Detection:

- It follows the approach, that deadlock has occurred, now check the deadlock and apply an algorithm to recover back from it.
- Detection & recovery has to deal with various overheads.
 - ① Runtime cost to manage the information by executing detection algorithm
 - ② Potential losses inherent in recovering from a deadlock.

Recovery from Deadlock:

- After detecting deadlock, operator can be informed about deadlock & it will be recover manually.
- Another approach is to recover automatically.
 - ① Abort / Process termination.
 - ② Resources Preemption.

Process Termination:

- ② Terminate all deadlocked process at once
 - May be, the process have computed for a long time.
- ③ Terminate one process at a time until the deadlock cycle is eliminated.
After one termination, deadlock detection algorithm will check if deadlock is still or not. This reduces the overhead.

Resource Preemption:

- Preempt the resources from a process until the deadlock cycle has broken.
- Preempting resources requires following 3 issues to be addressed.
 - ① Selecting a victim:
which processes & which resources are to be preempted?
 - ② Rollback: after preemption, what should be done with that process? Roll it back

Date _____

to a safe state; totally, Restart it, ready que

③ Starvation: How to ensure that resources will not always be preempted from same process

cost factors. Only a finite no of times, a process should be chosen as victim to avoid starvation.

Add no. of rollbacks in cost factors.

MEMORY MANAGEMENT

- In uniprogramming system, main memory is divided into two parts.
 - ① one part for OS
 - ② one part for program currently executing.
- In multiprogramming systems, the user part of memory is further subdivided to accommodate multiple processes.
- This subdivision is dynamically carried out by OS, it is called memory management.

Memory Management terms

Frame: A fixed-length block of main memory.

Page: A fixed-length block of data that resides in secondary memory (disk). A page of data may temporarily be copied into a frame of main memory.

Date _____

Segment: A variable-length block of data that resides in secondary memory. An entire segment may temporarily be copied into an available region of main memory (segmentation) or the segment may be divided into pages which can be individually copied into main memory (combined segmentation & paging).

Memory:

- Collection of data in a specific format
- Consists of a large array of words or bytes, each with its own address.
- Internal storage area of the computer.
- It stores all kind of data & information

Types of Memory:

1. Physical Memory:

- Primary or main memory, hardware component inside a computer
- tangible & volatile (loses its contents when)

the power is turned off.

- directly accessible by CPU or other hardware components
- Ex: RAM, ROM

2. Logical Memory:

- a memory space that a program can use.
 - an abstraction of physical memory that the OS ~~can~~ provides to applications & processes
 - Programs can access to memory addresses
 - Provides a logical address space for each process, independent of actual physical memory
- Ex: CPU memory \rightarrow cache, embedded memory

3. Virtual Memory:

- An imaginary memory area supported by some OS (eg Windows).
- An extension of logical memory.
- VM is a memory management technique that provides an illusion to applications that they have access to large, contiguous block of memory even if the physical memory is fragmented or insufficient.

Date _____

- allow running processes to use more memory than physically available.
- using disk storage as an extension of physical memory
 - If the size of program is $>$ available memory (physical one) then virtual memory is used.

Available Memory:

- Amount of Physical memory not in use.
- Loading the OS takes up memory as well.
- When system boots up, available memory amount drops.
- Each time you open a program, it is loaded into your computer's memory, therefore reduces your available memory.

Memory Management:

- Functionality of an OS to handle & manage primary memory (RAM eg)
- Subdividing available memory among different processes is called memory management



Date _____

- keep track of all available memory locations that either it is allocated to some process or not.
- decides which process will get which part of memory at any time.
- OS need to achieve to broad tasks

Memory-Management-Unit (MMU)

- Each process must have enough memory space to execute in a processes can not run into memory space of another process and not be run into by another process.
- There are different types of memory in the system in all must be used properly to achieve effectiveness.

Memory Management Unit (MMU)

How OS handles: RAM ?

- It keeps track of parts of memory in use
- Allocate & deallocate memory to from processes when needed
- Swapping & paging b/w main memory & disk

Date _____

- The principal operation of memory management is to bring processes into main memory for execution by the processor.

Static & Dynamic Memory:

- ① Static Memory: SRAM - Static RAM
faster, less volatile.
requires more power & it is expensive.
- ② Dynamic Memory: DRAM - Dynamic RAM.
RAM must be constantly refreshed,
reenergized or it'll loose its contents.

Partitioning:

- ① Fixed Partitioning: The available memory, for use by processes. The simplest scheme is to partition available memory into fixed size regions with fixed boundaries.
It is static. A process may be loaded into equal or greater size of a partition.
It is simple w/ little OS overhead.

Internal fragmentation due to fixed sizes, inefficient use of memory.

② Dyna.

③ Dynamic Partitioning:

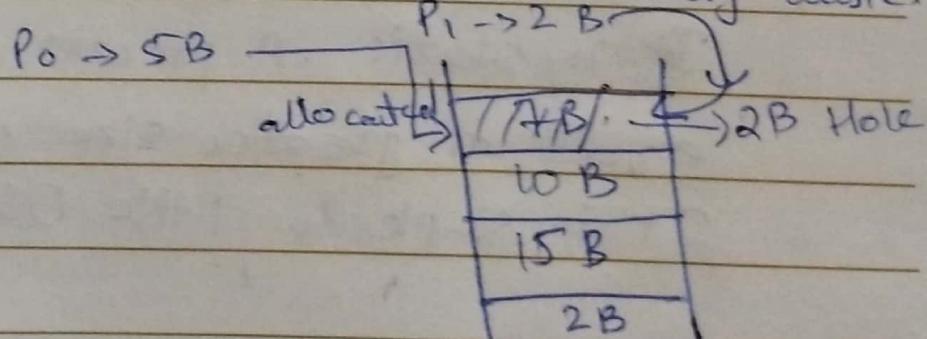
- Partitions are created dynamically; a process loaded into partition will be of same size.
- No internal fragmentation.
- Inefficient use of processor due to compaction to control external fragmentation.

3 Ways in which data is put in RAM

Placement algorithms:

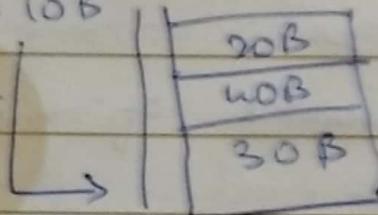
① First Fit:

- Begins to scan memory from the beginning and chooses the first available block that is large enough.
- Faster allocation but leads to memory waste.



2. Worst Fit:

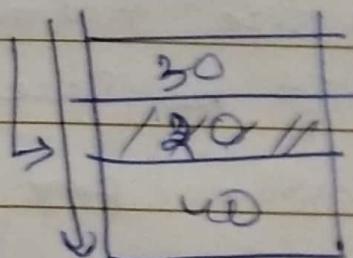
- It scans the entire available memory in place to allocate process in the largest block.
- Least efficient, waste of Memory & time
- It makes big wholes $P_1 \rightarrow 10B$ in memory. as slow as it traverse whole list.



3. Best Fit:

- Chooses the block that is closest in size to the requested process
- Best use of memory.
- Slow because traverse whole list then allocate.
- It'll choose closest blocks, but it can make small wholes in each block, - - .
- Less space but more time (time space tradeoff)

$$P_1 = 48B$$



Paging:

- Another memory management technique.
- A logical concept.
- A process is divided into small chunks called pages.
- Chunks of memory are called frames.

Process	Page no		Bytes.	frame no	
	0	1		0	1
P ₁	0	1	2 Bytes.	1	2 3
	1	2		2	4 5
		3		3	6 7
				4	8 9
				5	10 11
				6	12 13
				7	14 15

Process size = 4B
 Page size = 2B
 No. of pages = $\frac{4B}{2B} = 2B$
 2 Bytes in each page

Main memory

Page size & frame size
should be same

MM. Size = 16B

Frame size = 2B

Let assume, page 0 is in frame No. of frames = $16/2 = 8$
 In pge 1 is in frame 5 8B
 When CPU need any data of process, how do it know that where it is in main memory?



Date _____

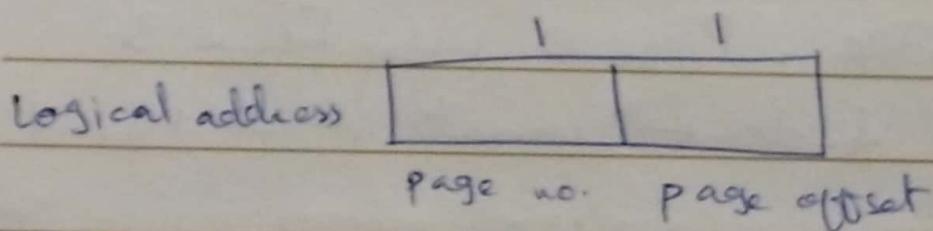
- A mapping function is used to convert the address of that required data, that is actually page's data into frame's address of main memory.^{by CPU}
- MMU does the mapping function. It uses page table.
- Page table contains frame numbers that where that page is actually present in the main memory.

Page table of a Process P₁

Page no	Frame no
P 0	F 3
P 1	F 5

Each process will have its own page table.

CPU works with logical address, logical address is made up of page number & page offset



no. of pages ~~=~~ 2^{\log_2} number of bits can be rep. in 1 bit.

Page no \Rightarrow 0, 1 \Rightarrow 1 Bit

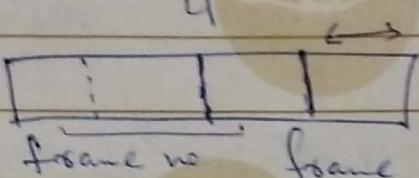
Page size contain 2 bytes over 0, 1.

SCHOOL
MY

Physical address is for main memory.

Logical addresses are converted into physical addresses.

Physical address



Our MM is of 16 Bytes, 16 can be represented in binary in 4 bits.

Physical address is of 4 bits.

Frame size is 2 Bytes, 2 can be represented in binary in 1 bit.

frame offset 1 bit

3 bit for frame no. 1 bit for offset.

as in an example CPU asks for ~~Byte~~ 1 Byte
its logical address \Rightarrow

0	1
PN	PO

In Page table, Page 0 has Frame 3.

The mapping function \Rightarrow converts it into physical address \rightarrow

0	1	1	1
PN	FO	return 7 Byte from MM.	

Date _____

Frames

- Frames are physical concept.
- In RAM.
- Small chunks are called frames in RAM.
- We can't see pages but we can see frames.

Mapping

- Converting virtual or logical address of pages by CPU into physical addresses of frames in MM.
- Performed by MMU.
- If an error occurs at any physical address, it retrieves back to CPU.

Paging Protection:

- Paging process is protected by adding valid / invalid bit
- Consider \rightarrow 14bit address space.
 $2^{14} = 16384$ bytes.
- address limit is set to 10468.

- If five processes P_0, P_1, P_2, P_3, P_4 are defined within this address space (i.e. their address will be less than 10468) then it is valid bit.
- P_5 has started before 10468, it alone is considered
- Remaining processes are invalid.
- Pages are internally fragmented in this way.
- The bit is in Page table.

Advantage & Disadvantages of Paging.

- ★ ① No external fragmentation.
- ☺ ② Simple memory management algo.
- ③ Easy swapping (equal sized pages of fixed size).
- ☹ ① Internal fragmentation.
- ② Page table consumes memory.

The wasted space in memory for each process is due to internal fragmentation consisting of only a fraction of the last page of a process.

Date _____

Fragmentation:

As processes are loaded & removed from memory, the available memory space is broken into little pieces. hence, it is not utilized or wasted. This is fragmentation in memory.

External fragmentation:

When there is enough memory space to satisfy a request but the available spaces are not contiguous; storage is fragmented into a large number of small holes.

Internal fragmentation:

If memory is divided into fixed partitions by processes all loaded in the partitions, if the process sizes are smaller than the partition, there is wasted space internal to the partition due to the fact that the block of data loaded is smaller than the partition.