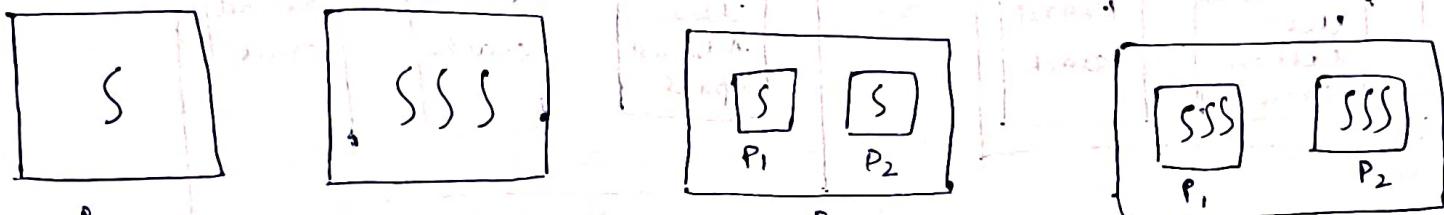


17/10/23

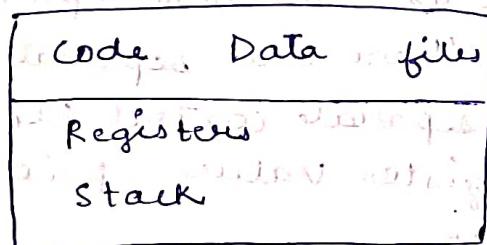
Threads

- Process → Program in execution
- Threads → Process is divided into multiple / single threads
→ basic unit of "CPU utilization"

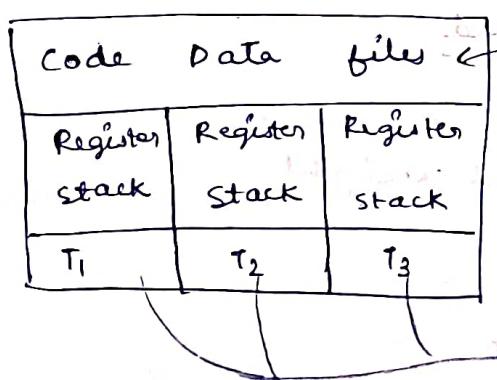


P_1 P_2 P_3 P_4
single process single process multiple processes multiple processes
single thread multiple threads single thread multiple threads

Process :



Threads :
(multiple)
threads

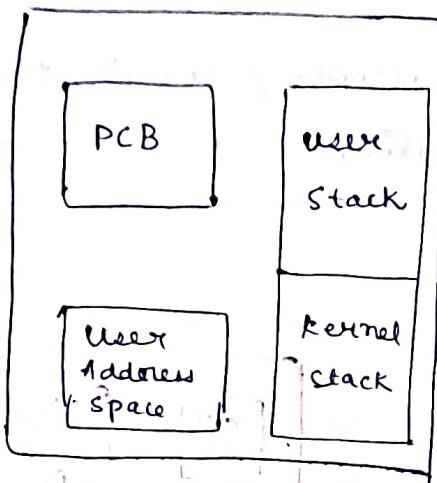


shared by all
threads
process is modified
by each thread
individual threads
with individual registers
& stack

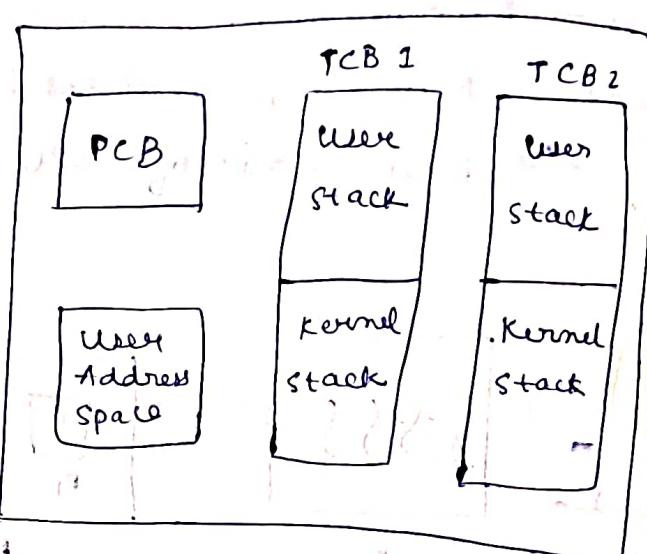
A thread is the basic unit of CPU utilization. We can have single/multiple threads in a system. Multithreading refers to the ability of OS to support multiple concurrent paths of execution within a single process. In a multithreaded environment, the thread shares code, data & files of the process & maintain their own registers & stacks.

TCB = thread control block

Single thread



Multi-thread



In a single threaded model, the process includes, ~~it~~ PCB & user address space as well as user & kernel stack to manage the behaviour of execution of processes. In a multithreaded environment, there is still a single process control flow & a user address space associated with the process but now there are separate stacks for each thread as well as separate control block for each thread that contains register values, priority & other thread related state info.

Advantages of Multithread

- 1) Creation is easy
- 2) Termination is light weight
- 3) Context switching
- 4) Resource Management
- 5) Communication

- It takes less time to create a new thread in an existing process rather than creating a brand new process.
- Less time to terminate a thread than a process.
- Less time to switch b/w 2 threads, than to switch b/w processes.
- Resource management is better.
- Communication in most OS, communication b/w independent processes require the intervention of the kernel to provide protection & the mechanisms needed for communications. As threads in a process share memory & file, they can communicate without invoking the kernel.

States of thread

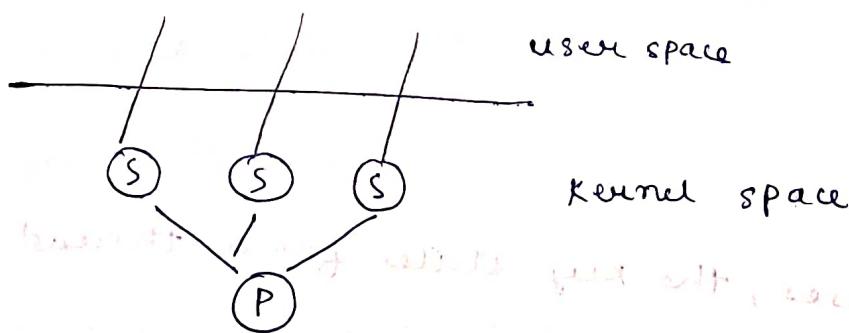
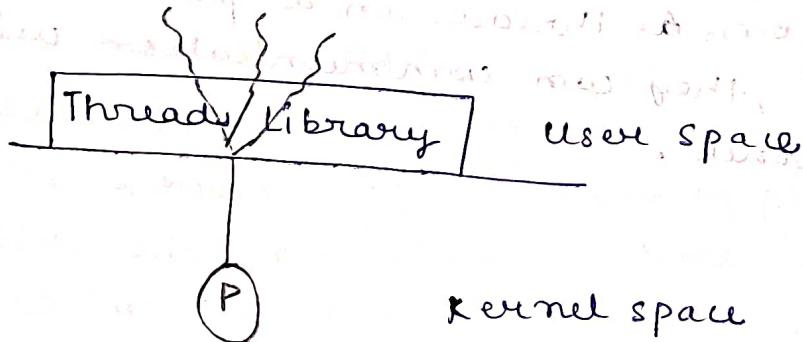
- Ready, Running
- Blocked
- Unblocked
- Finish
- Spawned

- As like processes, the key states for a thread are running, ready.
- When a thread needs to wait for an event it will block. And the processor will run the next ready thread.
- When the event for which a thread is blocked occurs, the thread is unblocked & sent to the ready state.
- When a thread completes, the register & stacks are deallocated.
- A thread within a process, may spawn another thread.

providing an "inst" pointer, arguments for the new thread.

Types of Thread

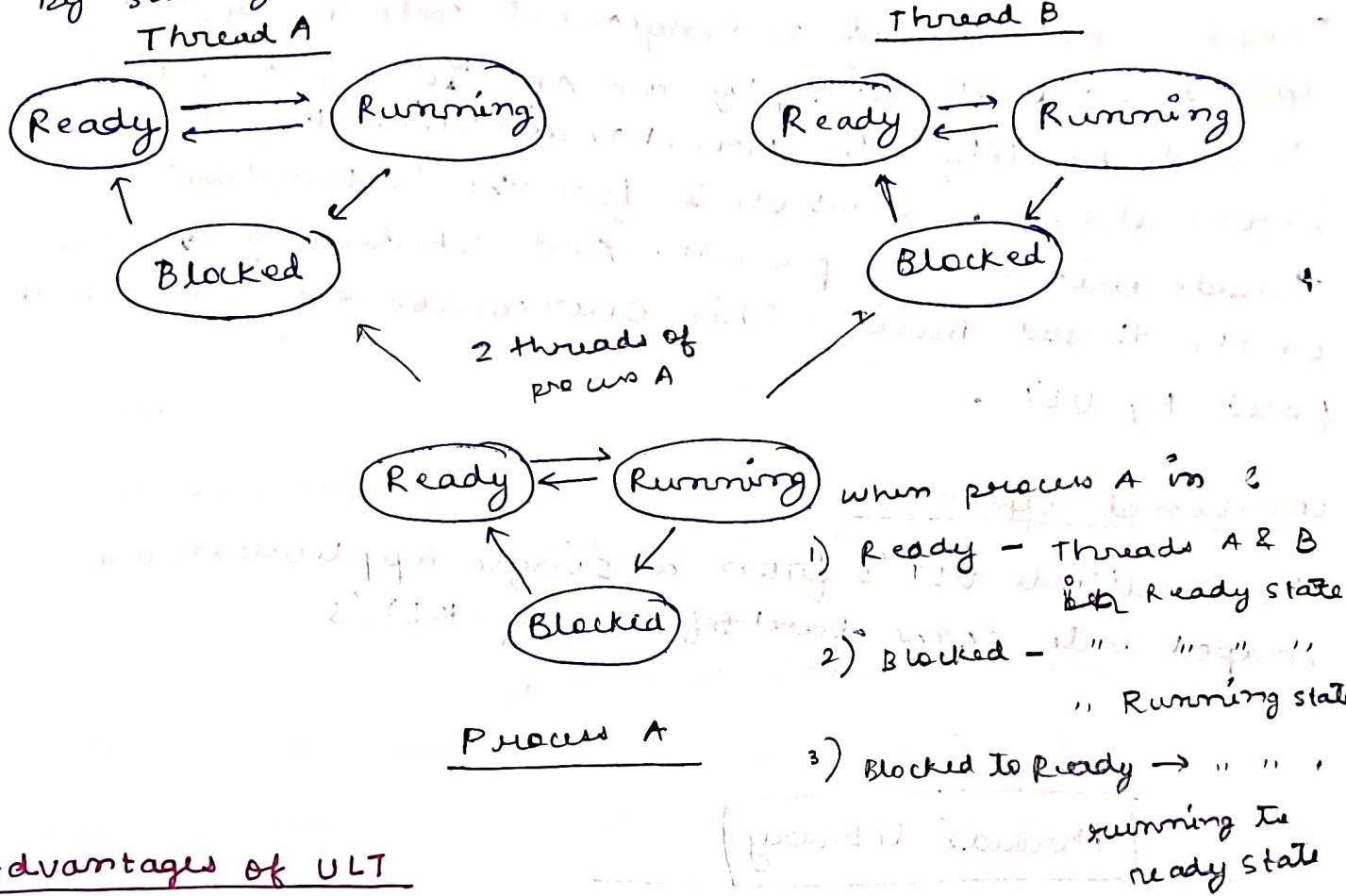
- 1) user-level thread
 - created by user
 - Not recognized by OS
- 2) Kernel level thread
 - created by OS
 - Recognized by OS



1. User-level thread

In user-level thread, all the thread management work is done by the application & kernel is not aware of the existence of threads. Any application can be programmed to be multithreaded using a thread library which is a package of routines used for ULT management. The thread library contains the code for creating & destroying threads. For passing msgs

& data b/w threads. For scheduling thread execution & by saving & restoring thread context.



Advantages of ULT

1. Scheduling can be application specific
2. Thread switching does not require kernel mode privileges, bcoz all the thread management are within the user space of a single process
3. ULT can run of on any process, no changes are made to the kernel

Disadvantages of ULT

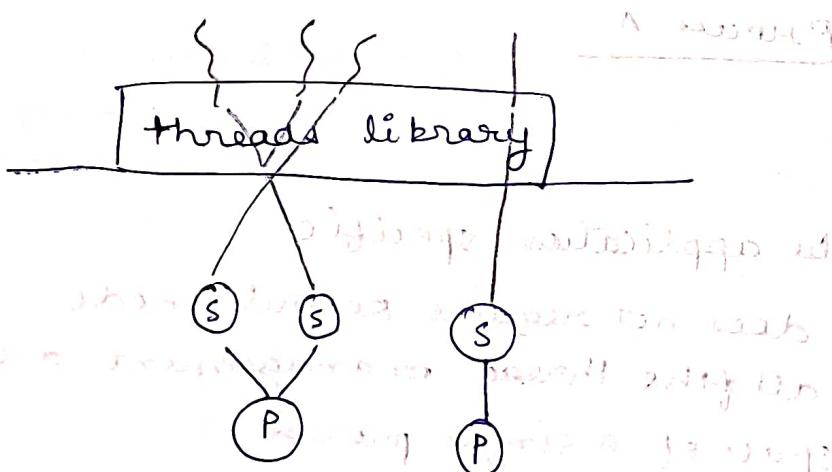
1. Many system calls are blocking in a typical process so blocking a process would block all its threads.
2. In a multithreaded application, it cannot take advantage of multiprocessor, as a kernel assigns 1 process to a processor at a time.

2. Kernel level threads

There is no thread management code in the application level, simply an API to the kernel. The kernel maintains kernel context thread facility. The kernel maintains individual information as a whole & for the individual threads within the process. And scheduling is done on the thread basis. This overcomes the 2 drawbacks faced by VLT.

Combined approach

The multiple VLT's from a single application are mapped into some ~~set of~~ ^{no. of} KLT's



18/10/23 4 types of relationship holds for the processes

1. 1:1 relationship (One thread one process)

Each thread of execution is a unique process with its own add. space & resources

2. M:1 relationship (Many threads one process)

A process defines an address space & of dynamic resource ownership. Multiple threads may be created & executed within that process

3. 1:M relationship (One -> To many)

A thread may migrate from one process environment to another. This allows a thread to be easily moved among distinct systems.

For eg. There is a main program that makes use of a I/O subprogram. The main program can spawn new process to handle the I/O, so there are 3 ways to implement this:

- The entire program can be implemented as a single process but we have memory constraints here.
- The main program & I/O subprogram can be implemented as 2 separate processes but it has a overhead of creating subordinate process & context switching.
- So the solution for these problems is, we treat the main program & I/O subprogram at a single activity & implement it as a single thread. One address space can be created for main program & one for the I/O subprogram, so that the threads can be moved b/w 2 address space as execution proceeds.

4. M:N relationship (Many threads & many processes)

It combines the attributes of many threads & many processes.

Multicore & Multi threading

T ₁	T ₂	T ₃	T ₁	T ₄	T ₂	T ₃
----------------	----------------	----------------	----------------	----------------	----------------	----------------

: single core

T ₁	T ₃	T ₁	T ₃
----------------	----------------	----------------	----------------

: Multicore

T ₂	T ₄	T ₂	T ₄
----------------	----------------	----------------	----------------

Challenges :

1. Data splitting (which thread to which core)
2. Balance (each thread gives equal weightage work)
3. Data Dependency (O/P of 1 is input to another & process is sequential)
4. Testing & Concurrency

Multithreading programming provides mechanisms for more efficient use of multicores. On a single core system, the execution of threads will be interleaved over time.

On the system with multicore, concurrency means the thread can run. However this proposes few challenges :

→ Dividing activities : The application must be examined to find areas that can be divided into separate concurrent tasks

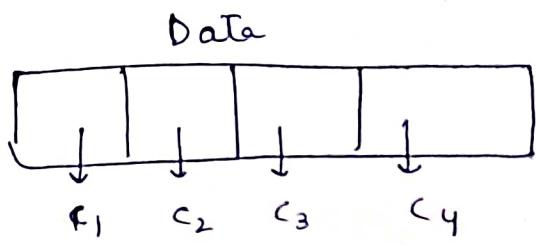
→ Balance : It must be ensured that the tasks perform equal work of equal value

→ Data splitting : As the applications are divided into tasks, the data accessed & manipulated must be divided to run on separate cores.

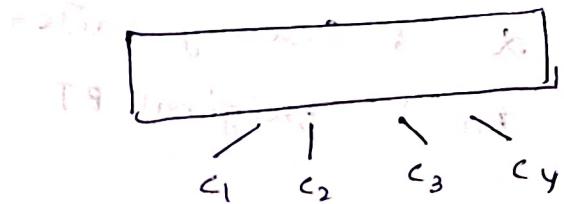
- Data dependency : The execution of tasks must be synchronized to accommodate data dependency
- Testing, debugging & dependency : When a program is running parallelly on multicore, there are different execution paths. Testing & debugging such programs is difficult.

(1)

(1) Data level parallelism



(2) Task level parallelism



There are 2 ways to parallelise the workload of a process —

- (1) Data level parallelism / Fine-grained threading
- (2) Coarse grained threading / Task level parallelism

(1) Data level parallelism

It focuses on distributing subsets of same data across multiple competing cores & performing the same operation on each core. It is also known as fine-grained threading.

(2) Task level parallelism

It involves distributing not data but tasks across multiple competing cores. An eg would be different threads performing a unique operation on the same array of elements. Aka coarse grained threading.

In SJF

To predict time (we can use the below methods) of next process :

(1) Simple Avg (of before processes)

(2) Exponential Avg.

$$T_{(n+1)} = \alpha t_n + (1-\alpha) T_n$$

where $T_{(n+1)}$: predicted time of next process

α : scaling factor

t_n : original BT

$$\textcircled{1} \quad \alpha = 0.5 \rightarrow P_1 \quad P_2 \quad P_3 \quad P_4$$

$$\therefore T_1 = 10 \quad 6 \quad 4 \quad 6 \quad 4$$

$$T(5) = ?$$

$$T(5) = \alpha t_4 + (1-\alpha) T_4 \rightarrow \text{(i)}$$

$$T_2 = \frac{1}{2} \times 6 + \frac{1}{2} \times 10 = 8$$

$$T_3 = \frac{1}{2} \times 4 + \frac{1}{2} \times 8 = 6$$

$$T_4 = \frac{1}{2} \times 6 + \frac{1}{2} \times 6 = 6$$

$$T_5 = \frac{1}{2} \times 6 + \frac{1}{2} \times 6 = 2 + 3 = 5$$

$$\textcircled{2} \quad P_1 \quad P_2 \quad P_3 \quad P_4 \quad \left\{ T(5) = ? \right.$$

$$5 \quad 2 \quad 6 \quad 3 \quad \left\{ \alpha = 0.5 \right.$$

$$T(1) = 5$$

Solve using simple averaging & exponential averaging

Sol" (1) Simple Avg

$$\bar{x}(3) = \frac{5 + 2 + 6 + 3}{4} = \frac{16}{4} = 4$$

(2) Exponential Avg.

$$\bar{x}(2) = \frac{1}{2} \times 5 + \frac{1}{2} \times 5 = 5$$

$$\bar{x}(3) = \frac{1}{2} \times 2 + \frac{1}{2} \times 5 = \frac{7}{2} = 3.5$$

$$\bar{x}(4) = \frac{1}{2} \times 6 + \frac{1}{2} \times 3.5 = \frac{9.5}{2} = 4.75$$

$$\bar{x}(5) = \frac{1}{2} \times 3 + \frac{1}{2} \times 4.75 = \frac{7.75}{2} = 3.875$$

25/10/23

Process Synchronization

- (1) Serial \rightarrow one after another } Process running
 (2) Parallel \rightarrow running at the same time }

Eg - P_1 ~~assuming two~~ P_2 ~~if it's~~ B shared = 5

```

    { int sharedx; { int j;
    x++; sleep(1); if parallelly run them
    w(x); } } } O/P = 6
    Race cond           but output should be 7
    
```

A race condition occurs when multiple processes/thread read & write data items so that the final result depends on the order of execution of instructions in the multiple processes.

- Q ① $P_1()$ $P_2()$ initially $B = 2$
- { beginning of function { beginning of function
 - (1) $C = B - 1$ (3) $D = 2 * B$ for parallel execution
 - (2) $B = 2 * C$ (4) $B = D - 1$ to take care of
 - }

Ans Case 1 : {1 2 3 4}

$$C = 1, \boxed{B = \frac{2}{3}}, D = 4$$

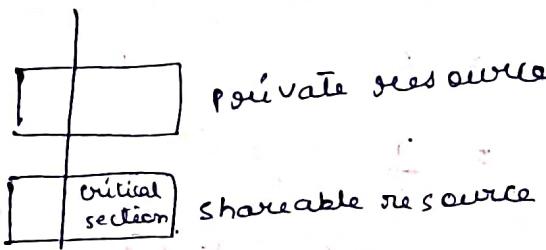
case 2 : {3 4 1 2}

$$D = 4, \boxed{B = \frac{8}{9}}, C = 2$$

case 3: $\{1 \ 3 \ 4 \ 2\}$ $C = 1$, $D = 4$, $B = 2$

case 4: $\{3 \ 1 \ 2 \ 4\}$ $D = 4$, $C = 2$, $B = 1$

** Process P



Critical Section

Shared resources are accessed by different processes at a time. The space where or the region where the processes access this shareable resource, that region is known as critical section.

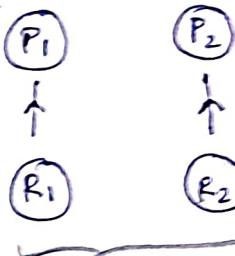
Concerns of OS

1. The OS should keep track of different processes which is done with the help of PCB.
2. The OS must allocate & deallocate various resources for each active process.
Resources can be:- Memory, I/O devices, files & processor time.
3. The OS must protect the data & physical resources of each process against unintended interference by other processes.
4. The functioning of a process must be independent of the speed at which its execution is carried out relative to the speed of another concurrent process.

Process Interaction

- (1) Process unaware of each other
- (2) Process indirectly aware of each other
- (3) Process directly aware of each other processes

Faci'ond
cooperation by sharing
cooperation by communication



Deadlock condition: $P_1 \rightarrow R_1 \rightarrow P_2 \rightarrow R_2 \rightarrow P_1$. R_1 is given to P_1 , P_1 wants R_2 and R_2 is given to P_2 . But now both P_1 " " P_2 requires the other resource. When both want to complete, neither leave each other resource so that they can complete. Therefore "deadlock condition" arises.

Criteria's to enter critical section

- Mandatory
- (1) Mutual Exclusion - 1 resource to one process at a time
 - (2) Progress - Only wanted processes shall enter critical section
 - (3) Bounded wait
 - (4) Architecture Neutral - work on all systems
- Optional

The criteria to access critical section is :

(1) Mutual Exclusion -

- Only 1 process at a time is allowed in critical section. Among all processes that have critical sections for the same resource / shared object.
- A process that waits in non-critical section must do so without interfering with other processes.

- It must not be possible for a process requiring access to critical section to be delayed indefinitely

- When no process is in critical section, any process that requests entry to its critical section must be permitted to enter without delay

(2) Progress

Only those processes should be allowed to enter the critical region who willingly want to enter critical region. No assumption are made about relative process speed / no. of processes.

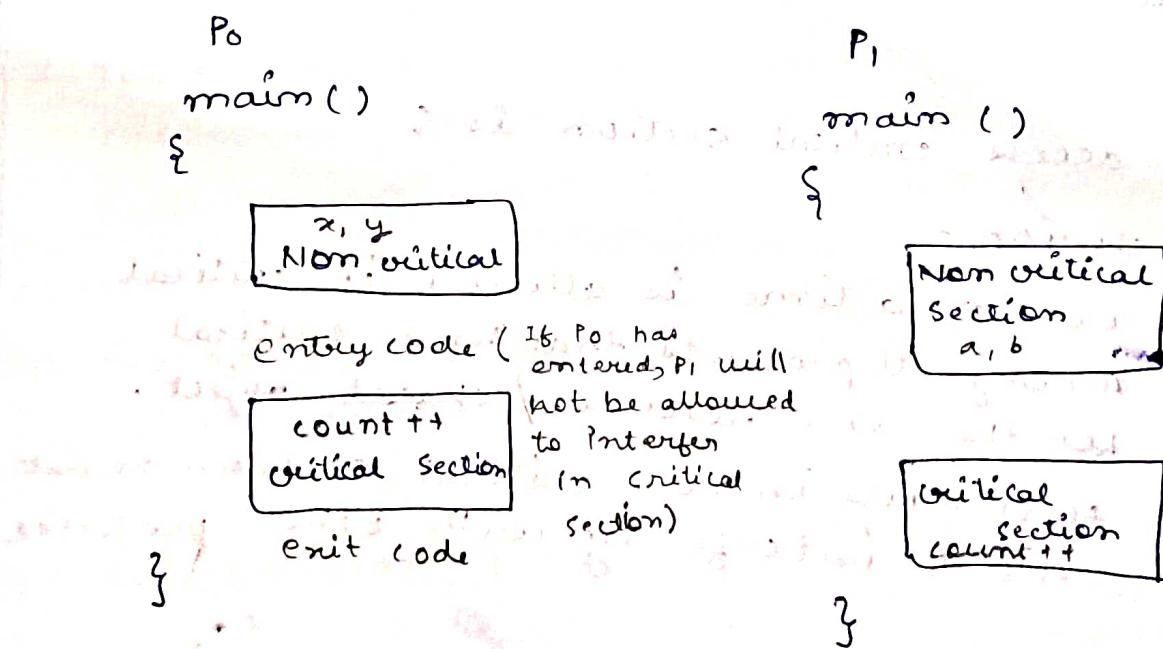
(3) Bounded wait

It is the max. time a process can wait. A process can remain inside the critical section for a finite time only.

Mutual exclusion & progress are mandatory criteria while bounded wait & architecture neutral are optional criterias.

27/10/23

Hardware Solutions



- Independent variables need not be added in the critical section. If any program wants to enter critical section, it will not access directly, it will execute some entry section

Therefore P_0 can't enter even though the CS is empty.

3. Flag variable

P_0

{

$$\text{flag}(0) = T$$

while ($\text{flag}(1) = \bar{T}$) ; ; (1 - 320)

CS

$$\text{flag}(0) = F$$

}

P_1

{

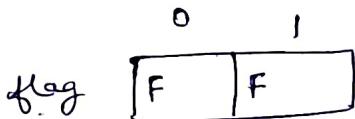
$$\text{flag}(1) = T$$

while ($\text{flag}(0) = \bar{T}$),

CS

$$(\text{since } \text{flag}(1) = F)$$

}



→ Here we have no progress

- This solⁿ guarantees mutual exclusion but creates another problem if both processes set their flags to true before either has executed the while statement then each process will think that the other has entered the CS causing deadlock.

4. using both variable (Peterson's solution) only for 2 process

P_0

{

$$\text{flag}(0) = T$$

$$\text{turn} = 1$$

while ($\text{turn} = 1 \& \text{flag}(1) = \bar{T}$);

CS

$$\text{flag}[0] = F$$

P_1

{

$$\text{flag}(1) = T$$

$$\text{turn} = 0$$

while ($\text{turn} = 0 \& \text{flag}[0] = \bar{T}$);

CS

$$\text{flag}[1] = F$$

}

- satisfies both mutual E & progress
- The mutual exclusion is preserved in this solution. This solves the problem of progress. ~~This~~ This also follows ~~progress as problem of progress as~~
 - 1) If P_1 has no interest in CS, as it will set flag [1] = F
 - 2) P_1 can use CS repeatedly, this case is not possible as P_1 gives opportunity to P_0 by setting the turn to 0 before it enters CS. Hence we can say Peterson solution satisfies all mandatory criteria i.e. Mutual exclusion & progress.

31/10/23 Hardware Solutions

(1) ~~Test & Set~~
lock = false.

while (test & set (& lock));

CS

lock = false,

----->
boolean test & set(
boolean *target)

boolean n = *target
*target = true

return n;

100

FT

LOCK

100

target

F

Among h/w solutions, one of the ways to ensure mutual exclusion is disabling interrupts but this soln is not feasible in a multiprocessor environment as disabling interrupts can be time consuming since the message is passed to all the processes.

- Many modern computer systems therefore provide many special h/w instructions that allow us to either test / modify the content of a word / to swap the contents of 2 words atomically.

(1) Test & Set

- In this solⁿ, we combine line 1 & 2 from Lock Variable solⁿ, so that no preemption can occur & mutual exclusion is ensured. This solⁿ satisfies mutual exclusion & progress but fails on bounded wait. It eliminates formation of race condition.

(2) Compare & Swap

{ lock = 0

while (compare & swap (& lock, 0, 1) != 0) ;

CS ;

lock = 0 ;

{ }) the of last method

- (next method)

int * compare & swap (int * value, int expected, int newValue)

{ int temp = * value ;

 expected

if (temp == 0) { newValue , temp ;

* value = newValue ;

return temp

}

100
lock

100
Value

0

exp

1

newValue

0

temp

This solution satisfies mutual exclusion property as we take a shared variable i.e. initialized to zero. The only process that can enter in the CS is the one that finds lock = 0. All other processes attempting to enter CS go into a busy waiting mode (stuck in loop).

When a process leaves its CS it resets lock to 0 & at this point one of the waiting processes is granted access to CS. This process also satisfies progress but fails on bounded wait.

Atomic function \rightarrow no preemption in b/w

(3) Exchange

do

{

key = true

[while (key == true) {
 swap (&lock, &key); }]

cs

lock = false

swap (boolean * a, boolean * b)

{ temp = * a;

* a = * b;

* b = temp;

}

}

- * Not satisfies bounded wait bcoz no mechanism is there for 1 process to restrict from entering again & again



key



lock

In this solⁿ mutual exclusion is satisfied as a shared variable "lock" is initialized with 0. It excludes all other process from CS by setting lock to 1. Here we have used the swap fn which exchanges the value of key & lock when any process is in CS. Progress is also satisfied but this solⁿ also fails on bounded wait.

Disadv. of above solⁿ(s)

- 1) Bounded wait ↑
- 2) starvation
- 3) Deadlock

Demerits of hardware solutions

1. Busy waiting is employed while a process is waiting for access to a CS, it continues to consume processor time.
2. Starvation is possible

When a process leaves a CS & more than 1 process is waiting the 2nd process could indefinitely be denied access.

3. Deadlock

Suppose P₁ executes a special instruction & enters the CS. P₁ is then interrupted by P₂ which is a higher priority process. If P₂ now attempts to use the CS, it will be denied access bcoz of mutual exclusion.