



*Simple OS Definition:

→ A special piece of software that

- abstracts and
 - ↳ to simplify what the hardware looks like
- ~~enforces~~ arbitrates ^{or the hardware}
 - ↳ to manage, oversee, control

the use of a computer system.

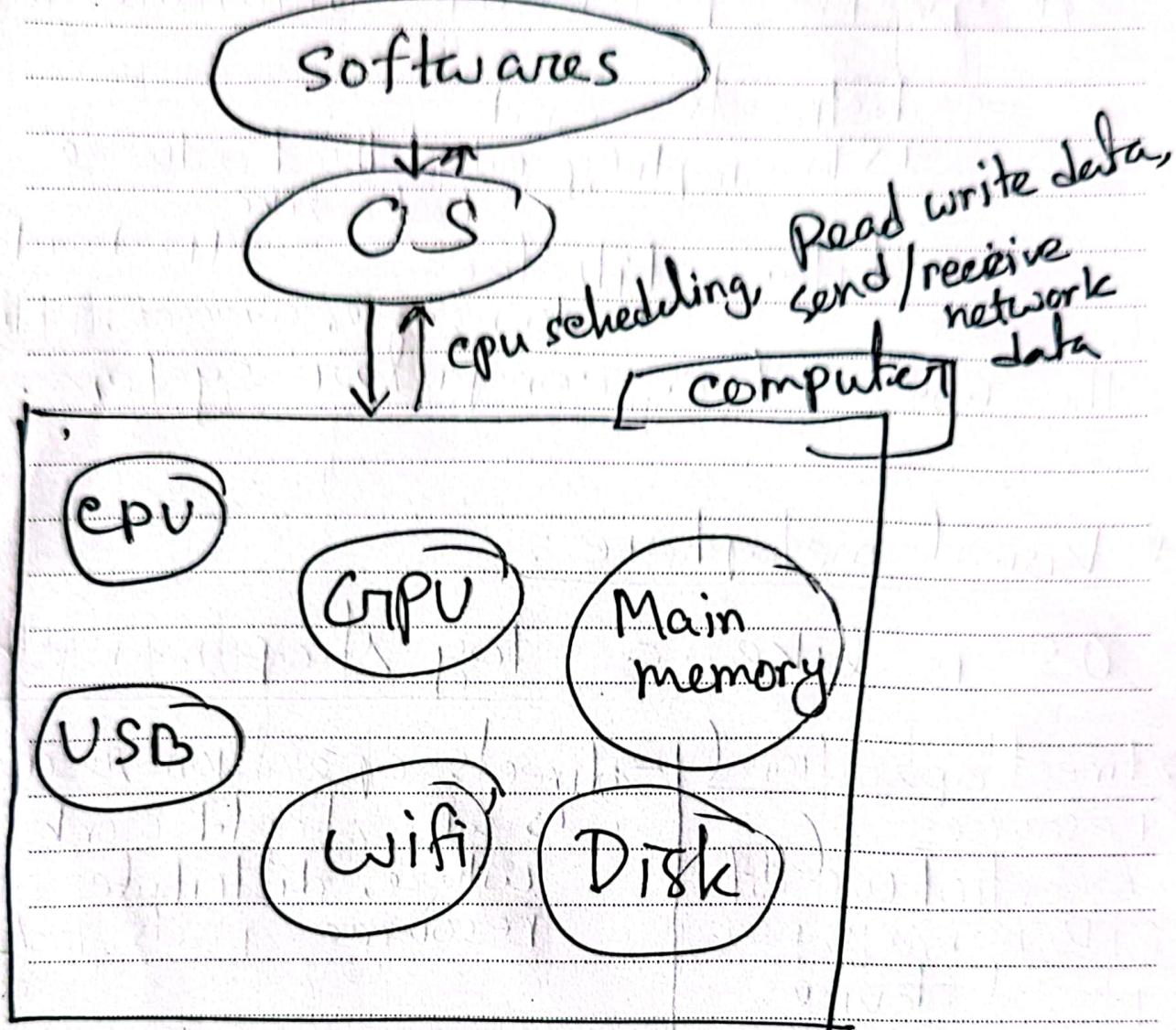
*Visual metaphor:

OS is like a toy shop manager

- | | |
|-------------------------------------------------------------------------------------------------|----------------------------------------------------------------------------------------------------|
| → <u>direct operational resources</u>
↳ Control use of CPU, memory, peripheral devices | → <u>directs operation resources</u>
↳ who should work where, distribute resource, parts, tools |
| → <u>Enforce working policies</u>
↳ fair resource access, limits to resource usages | → fairness, safety, clean |
| → <u>Mitigate difficulties of complex tasks</u>
↳ simplifies operation, optimize performance | |
| → <u>abstract hardware details.</u> | |



* What is OS?



→ hide software hardware complexity

→ resource management

→ Provide isolation & protection

(e.g. one software can't access another's memory)





So, An OS is a layer of system software that :

- ① directly has privileged access to the underlying hardware.
- ② hides hardware complexity
- ③ manages hardware on behalf of one or more applications according to some predefined policies
- ④ In addition, it ensures that applications are isolated and protected from one another.





① Abstractions (OS executes)

eg: Process, Thread, file, socket,
memory page

② Mechanisms: (run on the CPU)

eg: create, schedule, open, write,
allocate

③ Policies (memory manager underlying hardware)

eg: least-recently used, earliest
deadline first (EDF)

④ Memory management example:

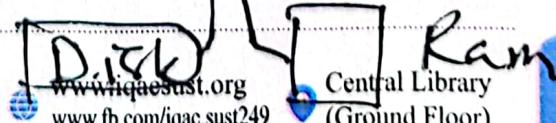
Abstraction: memory page

Mechanisms: allocate, map to a process

Policies: LRU - cache

Process

page





OS Design Principles

Date.....

① Separation of mechanism & policy

→ implement flexible mechanisms to support wide range of policy

like. LRU, LFU, Random

② Optimize for common case

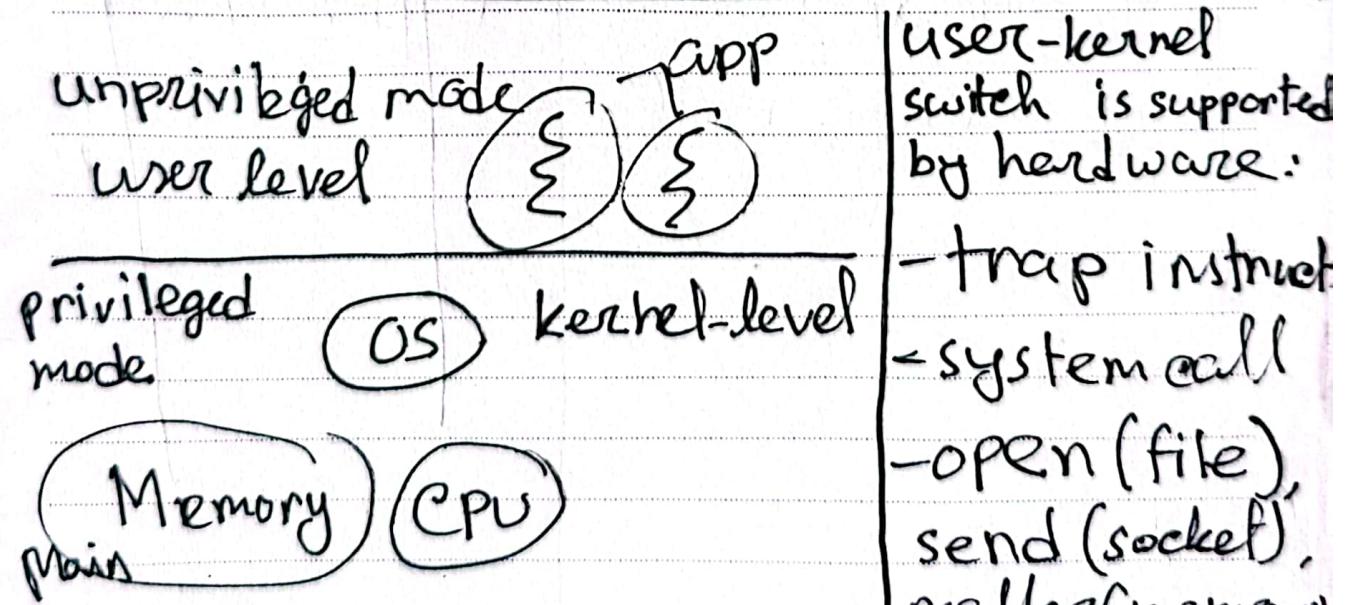
→ Where the OS will be used

→ What will the user want

to execute on that machine

→ what are the workload requi.

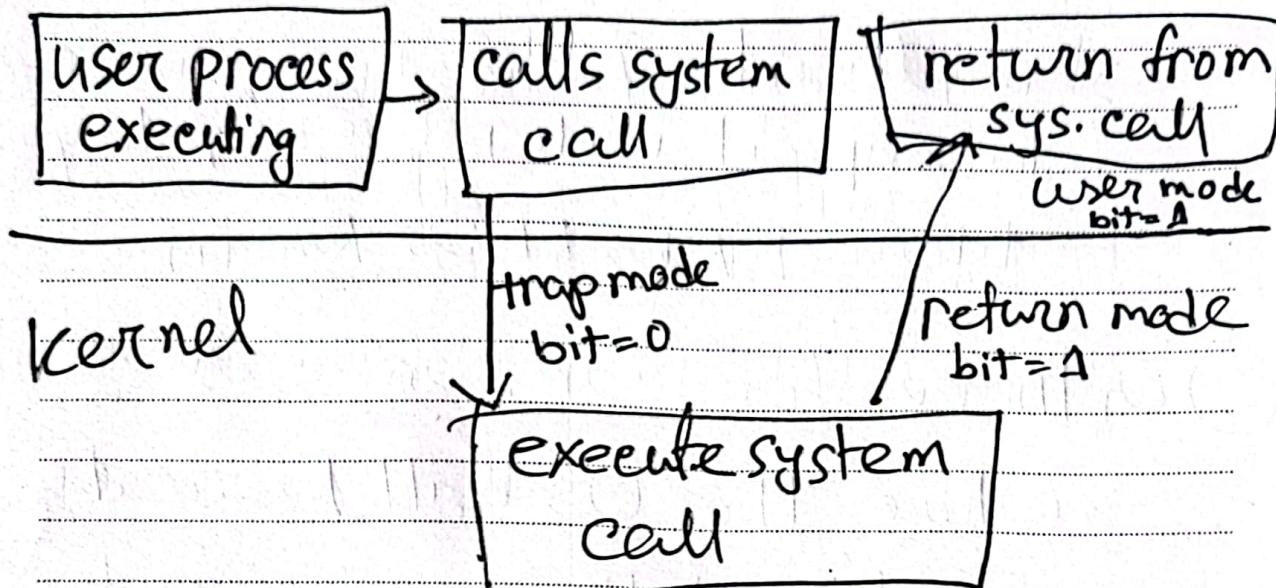
User-kernel protection Boundary





System Call flowchart

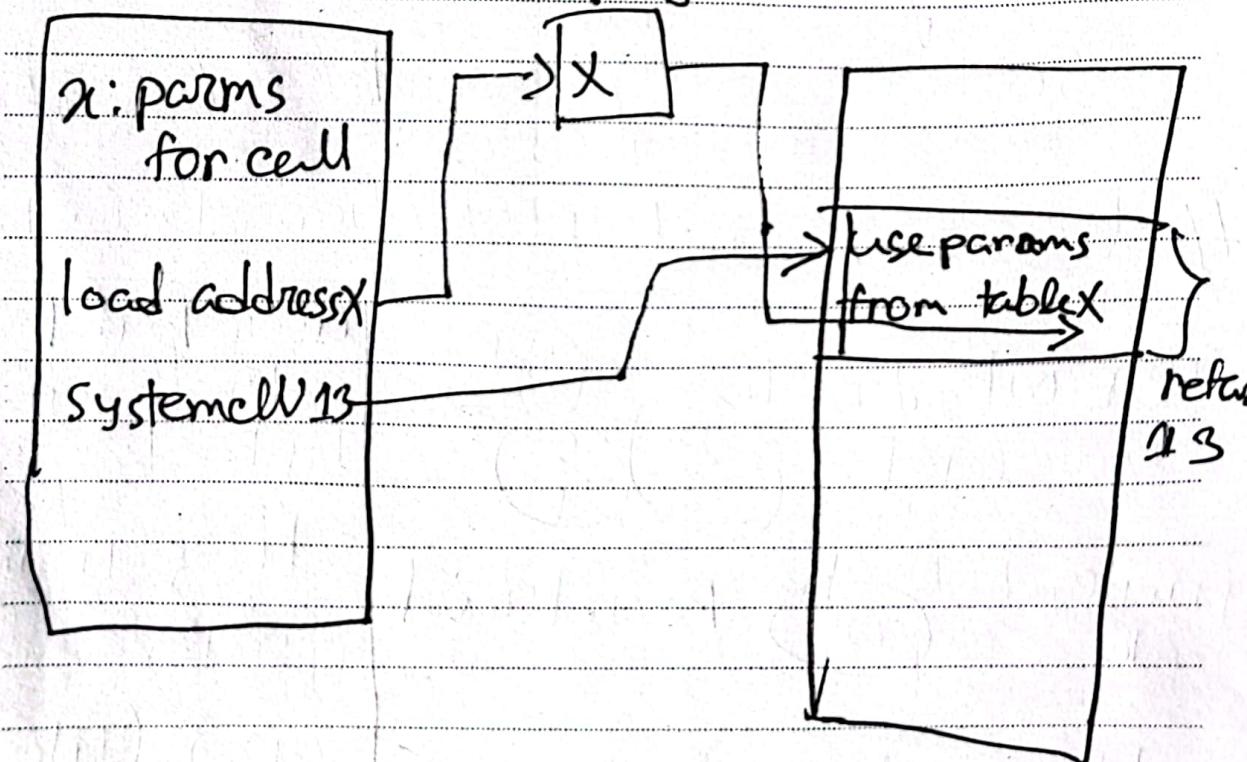
user process



kernel

user program

register





Basic OS Services

- process management
- file management
- device management
- memory management
- storage management
- security . . .

for 64 bit linux:

- send signal to a process → kill
- set group identity of a process → SETGID
- mount a filesystem → Mount
- read/write params → Sysctl





OS organization

* Monolithic OS

→ every possible service that any one of the application can require or any type of hardware will demand is already part of the OS.

Pros: everything included

inlining, compiletime optimization

Cons: customization, portability, manageability, memory footprint, performance

* Modular OS (associated with Linux)

→ has number of basic services and apis already part of it but everything can be added as modular has a module interface that any module must implement in order



pros: maintainability
smaller footprint

less resource needed

cons: indirection can affect
performance

maintenance can still be
an issue.

Microkernel: require the most basic

primitives at sys level such
as: address space, thread.

everything else will run at user
level like: file system, drivers etc.

pros: size, verifiability.

cons: cost of user/kernel crossing
complexity of software development
portability.



* A process is like an order of toys

↳ program counter,
stack

→ State of execution

↳ completed toys,
waiting to be built

→ Parts & temp. holding area

↳ plastic pieces,
containers

↳ data register
state occupies
space in memory

→ may require special hardware

↳ glue, swing machine

↳ I/O devices

⇒ Basically, process represents
the execution state of an active
application.

⇒ What does a Process look like?





Vm

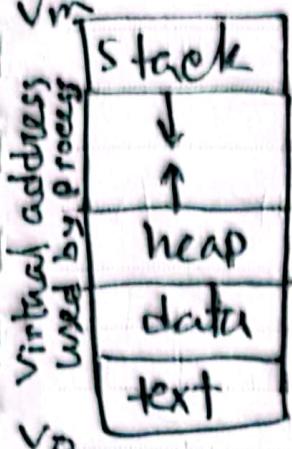


fig : 1

types of state

→ text & data

= static state when process first load

→ heap: dynamically created during execution

→ stack: grows & shrinks

in LIFO Order

eg: call-stack in JavaScript

address space == "in memory" representation of a process

page table == mapping of virtual to physical address

* swapping some virtual address to disk when that's not needed.

example!

if two processes, P1 & P2 are running at the same time, what are the virtual address space ranges they will have

→ both will have access to all the

fig 1





* Process Execution State .

⇒ How does a OS know what a process is doing ?

→ Program Counter (PC): At any given point of time , the CPU needs to know where in the instruction sequence of the process is. (assembly code)

→ CPU registers: PC is stored in maintained CPU while the process is executing in a register. It is necessary to as they may store info like address for data

→ Stack pointer: the last state of the process. Need to know the top of stack .



* Process Control Block (PCB) :

A data structure that the OS maintains for every one of the processes that is manage.

process state
process number
pc
Registers
memory limits
list of open files
priority
signal mask
CPU scheduling info

- PCB is created when process is created
- Certain fields are updated when process state change like mem. increase
- Other fields change too frequently.

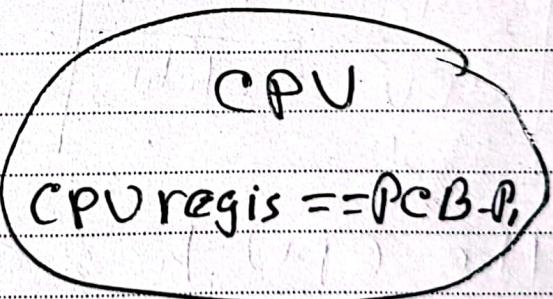
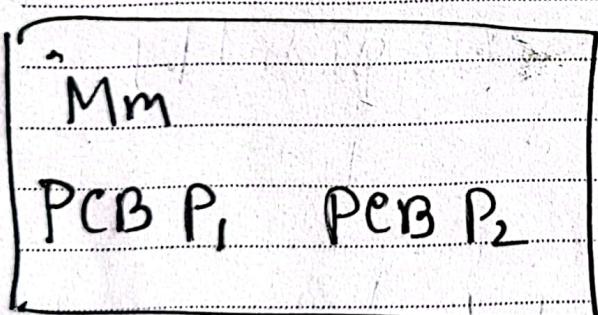
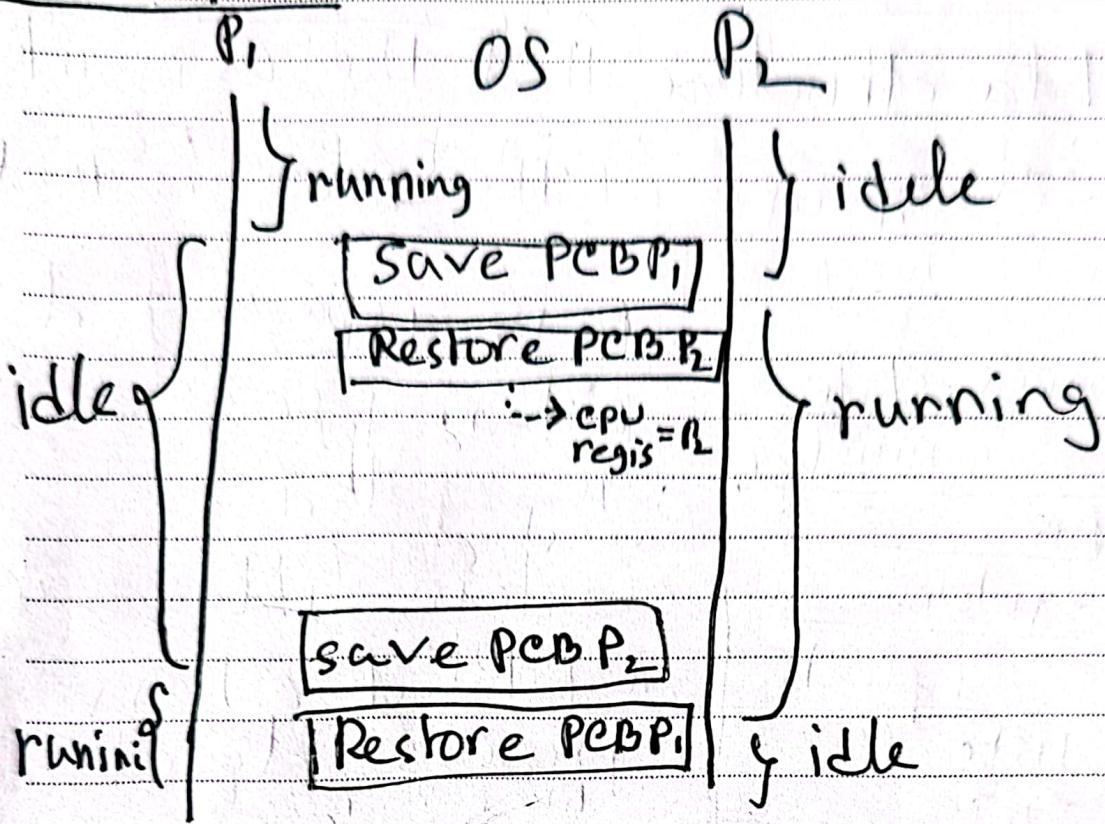
The CPU has a dedicated register which it usages to track the current PC for the currently executing process.

The PC register will auto. upgrade by CPU on every new instruction.





example





Context Switch

→ A mechanism for switching the CPU from the context of one process to context of another.

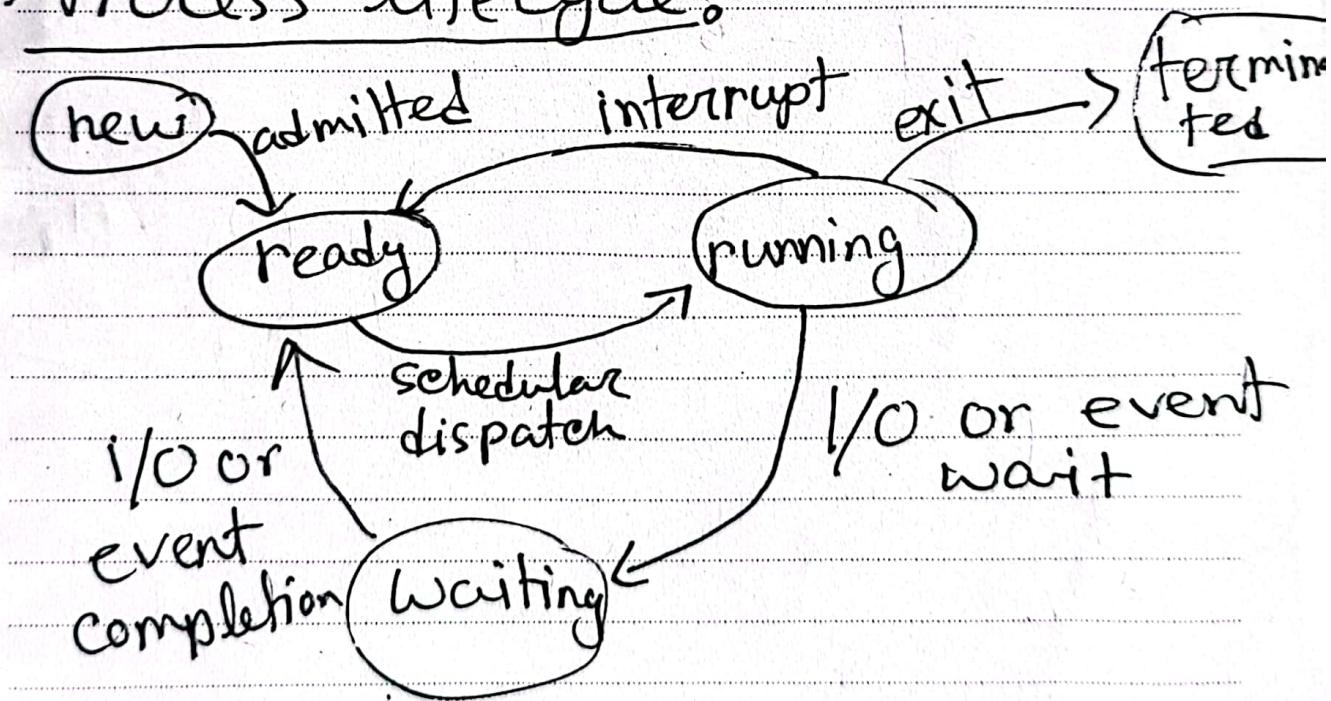
- they are expensive.

- direct costs: number of cycles for load & store instructions
- indirect cost: cold cache! cache miss!

Processor cache: store

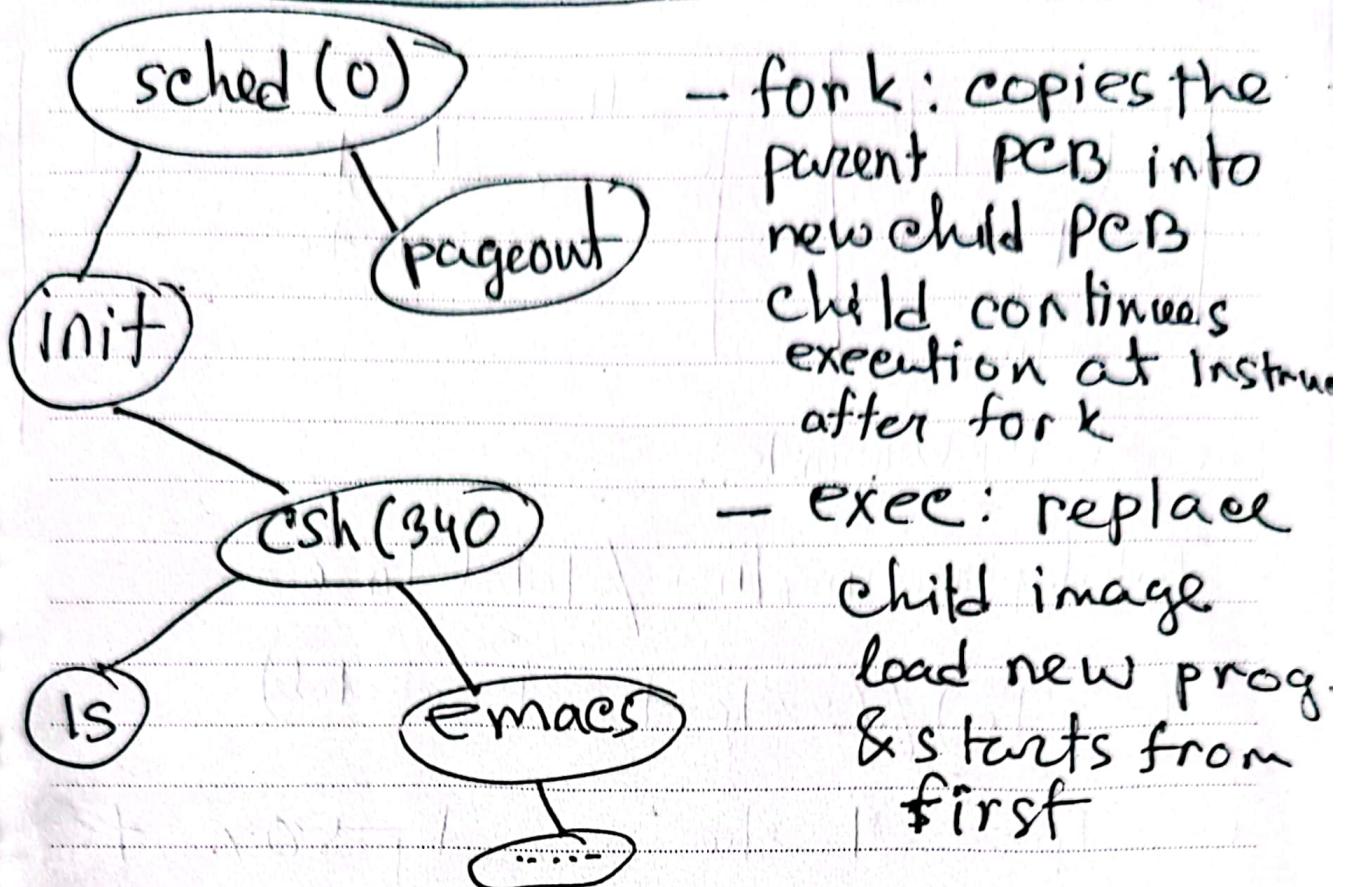
running processes info

Process lifecycle:





Process Creation



Role of CPU Scheduler :

It determines which one of the currently ready processes will be dispatched to the CPU to start running and how long it should run for.

Ready Queue :



Efficient ✓

OS must
 • preempt = interrupt and save current context

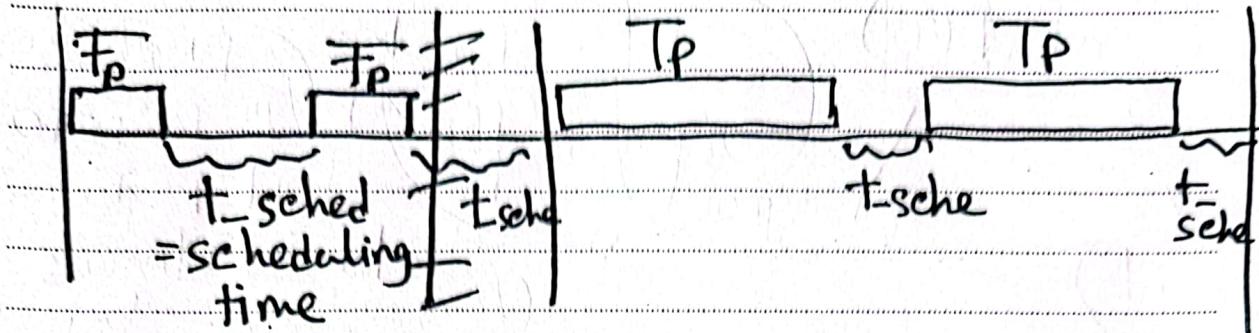
• scheduler : run this to choose next proc.

• Dispatch : dispatch & switch info its context





* Length of a process



useful CPU work:

= Total processing time / total time

$$= (2 \times T_p) / (2 \times T_p + 2 \times t\text{-schedule})$$

if $T_p = t\text{-sche}$ \Rightarrow only 50% of CPU time spent on useful work

if $T_p = 10 \times t\text{-sche}$

\Rightarrow 91% of CPU useful work

Timeslice = time T_p allocated

to a process on CPU

Scheduling Design Decisions:

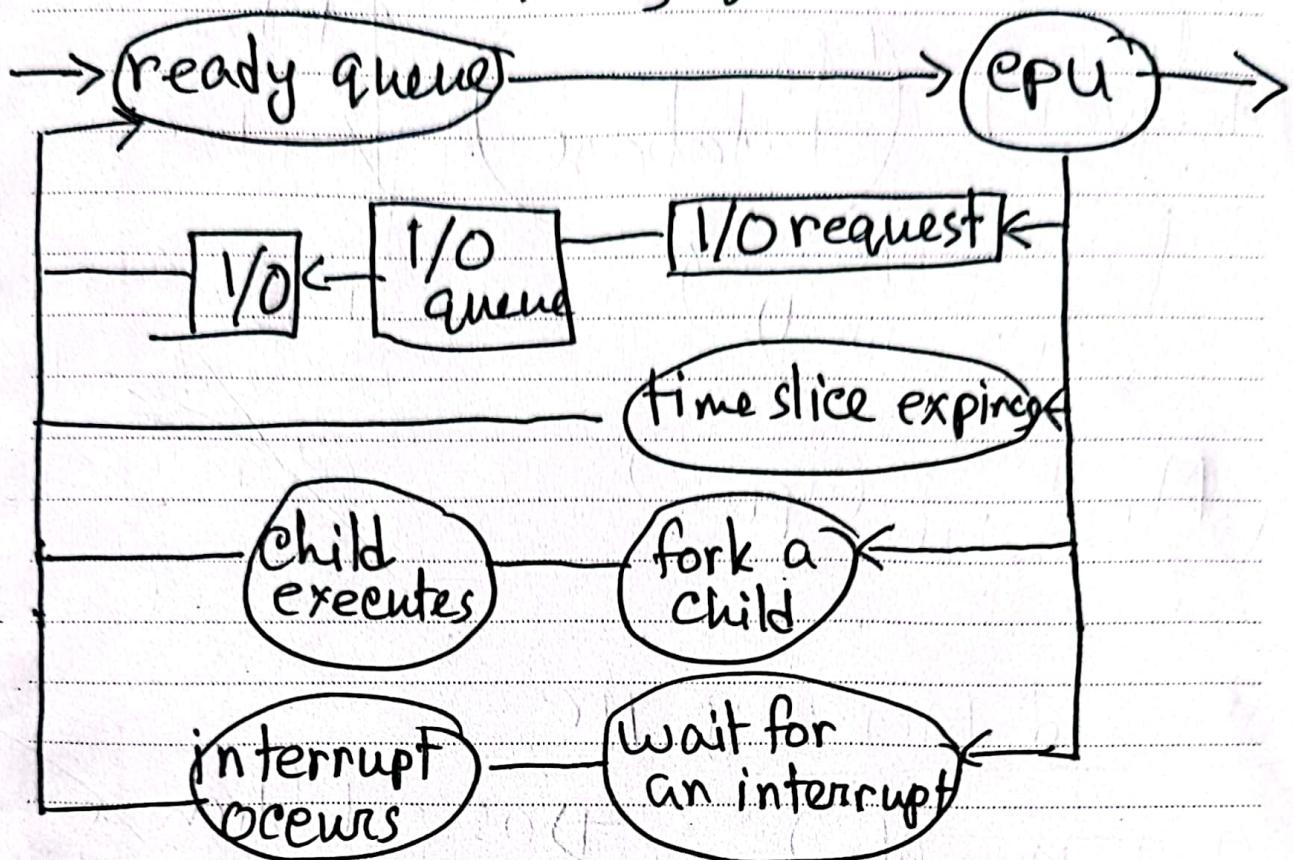
\rightarrow finding out appropriate timeslice

\rightarrow metrics to choose next process to run.





1/0 | How process makes its way to
a ready queue.

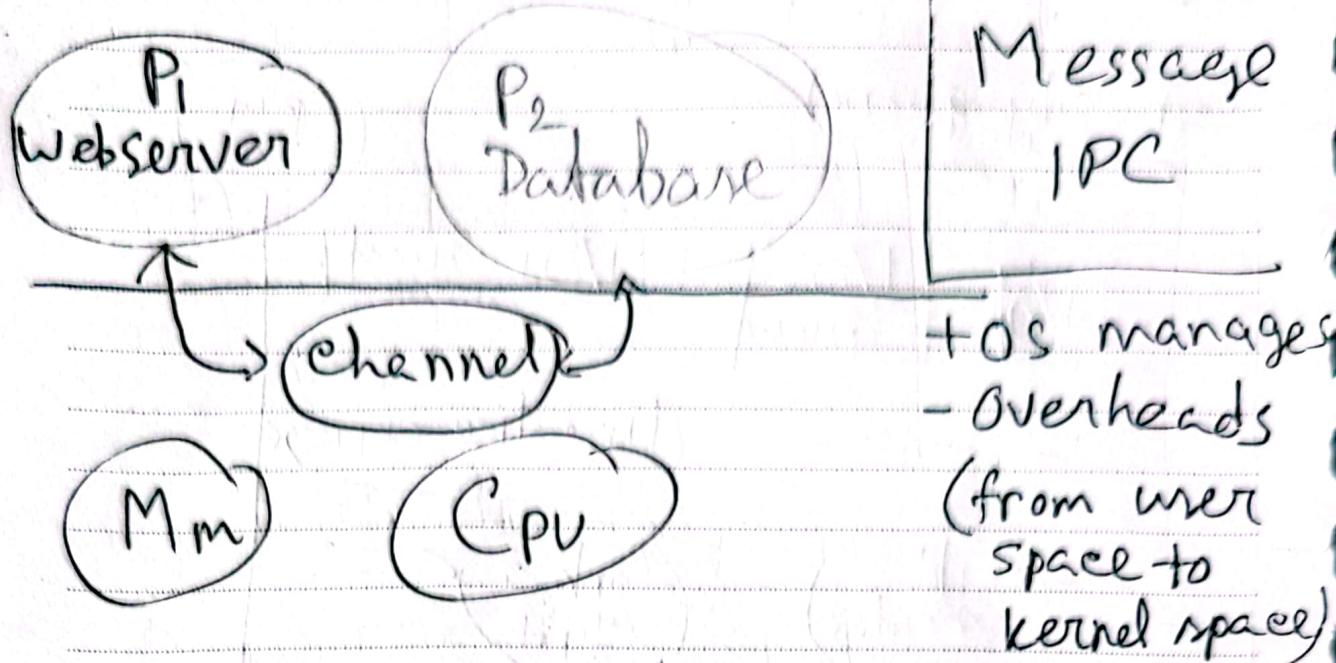


Inter process Communication (IPC)

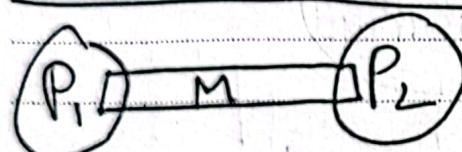
- transfer data/info between address spaces
- maintain protection & isolation
- provide flexibility & performance

Message-Passing IPC

- OS provides communication channel like shared buffer
- Processes read/write to/from



*Shared Memory IPC



→ OS establishes a shared channel and maps it into each process address space

+ Processes directly read/write from memory

+ OS is out of the way (same space)

- (-re)implement code



THREAD



* A thread is like a worker in a toy shop.

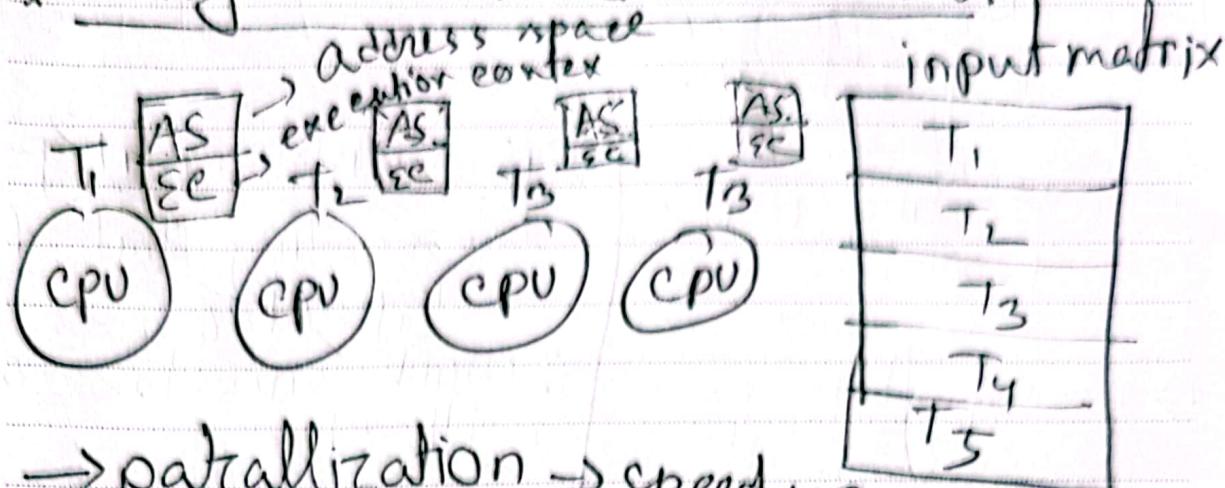
- | | |
|---------------------------------------|-------------------------------------|
| ↳ executing unit of a process | → is active entity |
| ↳ many threads executing | ↳ executing unit of toy order |
| → requires co-ordination | → work simultaneously with others |
| ↳ sharing of I/O device, CPUs, memory | ↳ many workers complete toy order |
| | ↳ sharing tools, parts, workstation |

* difference between thread & Process

- | | |
|---------------------------------------------|-----------------------------------------------------------|
| → Process means a program in execution | → thread means a segment of a process |
| → NOT Lightweight | → Lightweight |
| → takes more time to terminate/create | → less time |
| → takes more time to context switching | → less time |
| → more resource | → less |
| → Different process are treated differently | → All the level peer threads are treated as a single task |



* Why multithread is useful



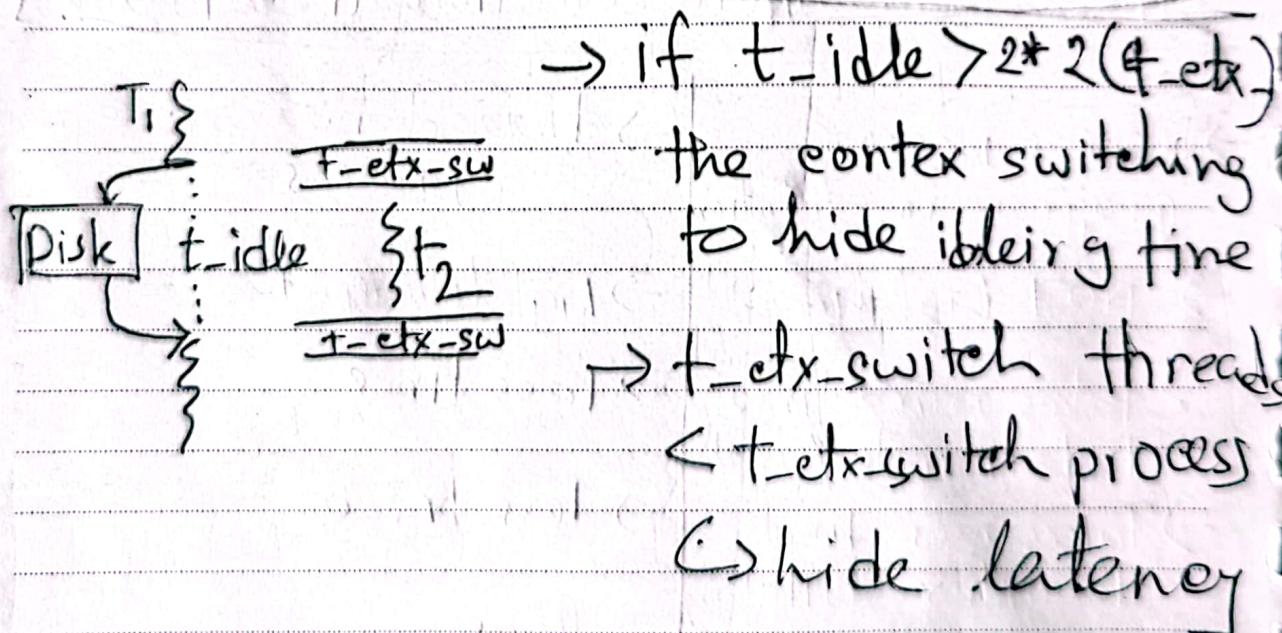
→ parallelization → speed up

→ specialization → hot cache

(as small thread share same memory)

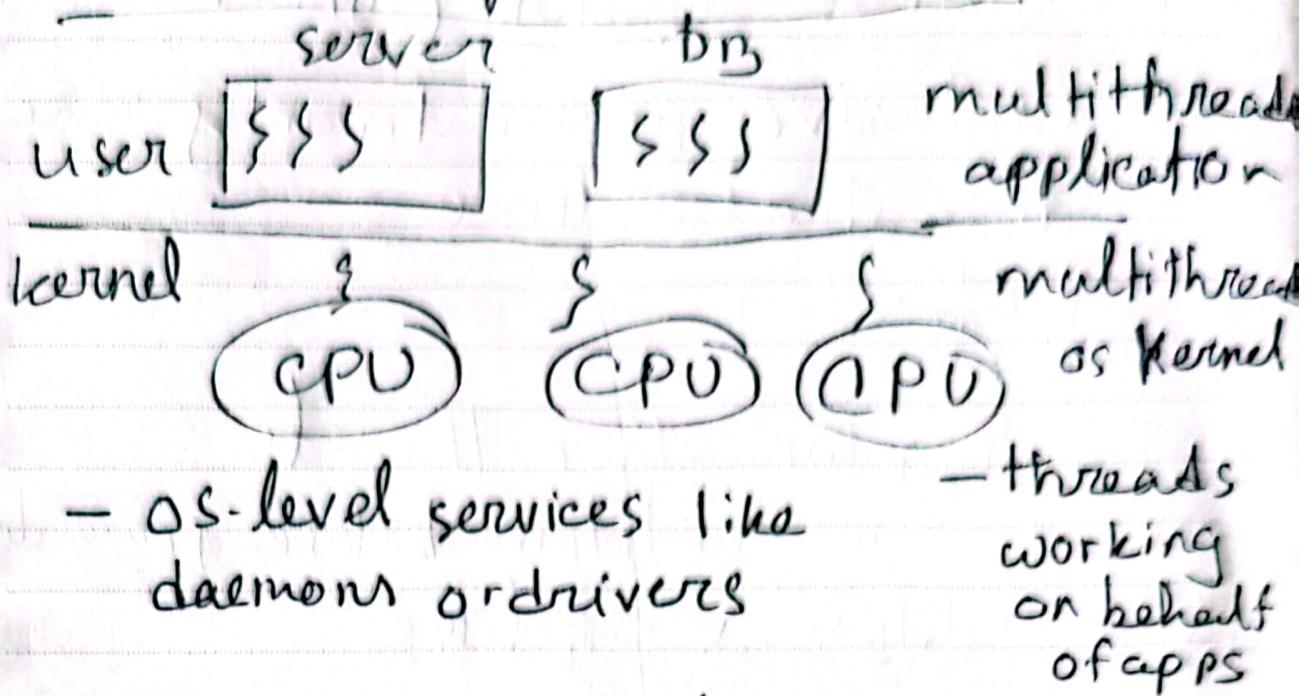
Extension of why multithread is useful

* Are threads useful on a single CPU?





* Benefits to application & OS



* Basic Thread Mechanisms

⇒ Threads & Concurrency

- Threads share the same virtual-to-physical address mapping
- it's possible to access same memory space from different threads (not possible for processes)
- violates consistency





* mutual exclusion (synch. mecha.)

→ exclusive access to only one thread at a time

→ mutex

→ waiting for other threads

- specific condition before proceeding

- condition variable

* Thread Creation (accor. Bruell)

Thread type + type

- thread data structure | id, PC, SP, register, stack, attributes

Fork (proc, args)

- create a thread
- not unix thread { To

T_i DS
PC = proc
stack = args

Join (thread)

- terminate
a thread {

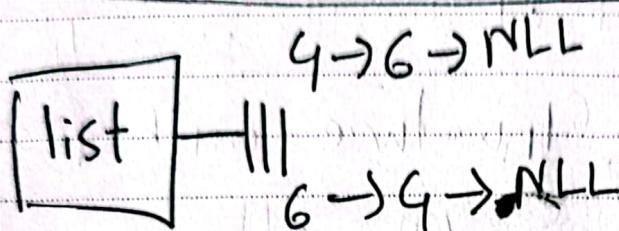
t_i = Fork(proc, args)
↳ T_i

child result { return result
= join(t) end



Example of Thread creation

→ Thread t1;
Shared-list list;
t1 = fork(safe-insert, 4);
safe-insert(6);
join(t1);



{ To

t1 = fork(safe-insert, 4) → t1 {

{ t1, ||| safe-insert(4)

safe-insert(6)

child-result ← join(t1)

* Then discuss mutual Exclusion

Lock(mutex){

critical section

+880-321-717850, 017470
73850 (Ext: 249, 2249)

iqac-off@sust.edu
iqac.sust@gmail.com

Mutex

locked?

owner

blocked -

www.iqac.sust.org
www.fb.com/iqac.sust249

Central Library
(Ground Floor)



Mutex lock with conditional variable

Combine Mutex lock & conditional variable for Read/Write Problem
→ multiple read at the same time.

→ Critical Section Structure with proxy

{ // Enter Critical Section
 | perform critical operation (read/write shared file)
 | // Exit Critical Section

// Enter
Lock(mutex) {
while (! predicate_for_access)
 wait(mutex, cond-var)
update - predicate
} // unlock

// Exit
Lock(mutex) {
update predicate;
signal/broadcast
(cond-var);
} // unlock



* Avoid Common mutex mistakes in mutex

→ keep track of mutex/cond. variables

used with a resource (comment)

→ Check that you are always using lock & unlock

→ use a single mutex to access a single resource.

→ check that you are signaling correct condition

→ check that you are not using signal when broadcast is needed

→ dead lock

→ Spurious wake up → (unnecessary wake up)

// writer

Lock(counter-mutex){

P-C = 0;

Broadcast(read phase);

signal(write phase);

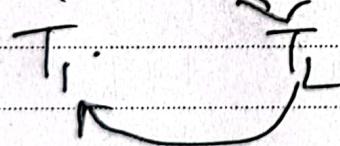
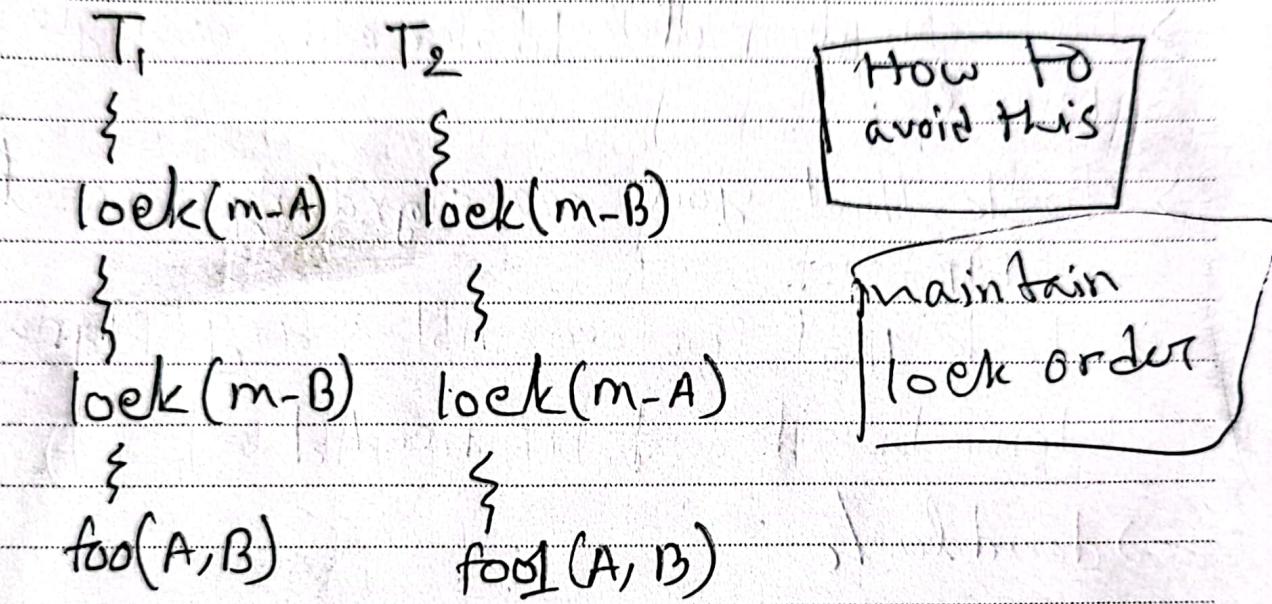
}

// Readers

wait(counter-mutex
read/write phase);



* Deadlock: Two or more competing threads are waiting on each other to complete, but none of them ever do.





* Kernel & User level thread

To execute a user-level thread, first it must be associated with kernel level thread. Then, the OS level scheduler must schedule that kernel level thread onto a CPU.

① One-to-one model: Each user

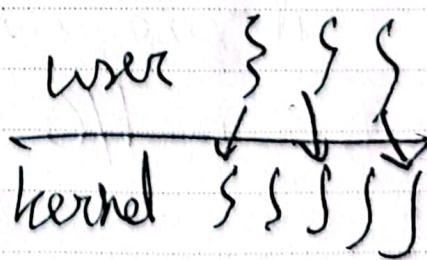
thread is associated with one new/free kernel thread.

⊕ OS understands threads, sync, blocking

⊖ must go to OS for all operations

⊖ OS may have limits on policies
thread

⊖ portability.





⑩ many-to-one model: User level

Threads are mapped onto ~~multiple~~
a single kernel level thread.

In the user level a thread management library that decides which user-level thread will be mapped onto the kernel level thread at any given point of time.

⊕ totally portable, doesn't depend on OS limits.

⊖ OS has no insights into application needs (will see the system as single threaded)

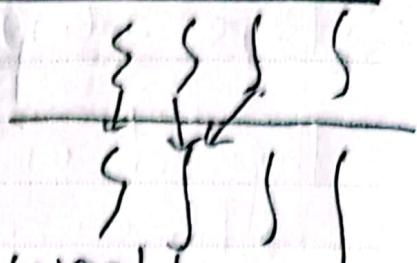
⊖ OS may block entire process if one user thread is blocked on I/O.





(III)

many-to-many:



- + can be best of both world.
- + can have bound or unbounded threads.
- Θ requires coordination between user level and kernel level threads.

→ Process scope: user-level library manages thread within a single process

→ System scope: system wide thread management by OS-level thread manager

*

multithreading pattern

I Boss-workers

II Pipeline

III Layered.

Scenario: for a toy shop.

I accept the order

V paint & decorate

II parse the order

VI assemble the toys

III cut wood parts

VII ship the order





⇒ Boss - Worker Pattern: One boss thread and multiple worker thread.

- boss: assigns work to worker
- worker: performs entire task.

boss: step 1, worker: 2-6 step.

Throughput of the system limited
to by boss thread ⇒ must keep boss
efficient

$$\text{Throughput} = \frac{1}{\text{boss-time-per-order}}$$

boss: - directly signalling specific
worker

⊕ workers don't need to sync.

⊖ boss must keep track what
each worker will be doing

⊖ throughput will go down.

- placing work in producer/consumer
queue.





- ⊕ Boss doesn't need to know the details about workers.
- ⊖ Queue sync. needed.
- ⊕ better throughput
- How many workers?
 - on demand - pool of workers.

Boss + worker variants:

→ all worker created equal

→ worker specialized for certain task

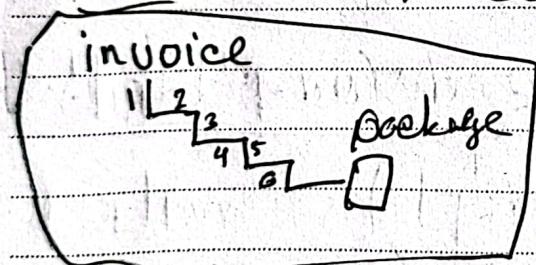
↳ better locality, I/O's managed

↳ load balancing (how many threads should we assign a particular task)



⇒ Pipeline Pattern:

- threads are assigned subtask in the system
- entire tasks == pipeline of threads
- multiple tasks concurrently in the system, in different pipeline stages
- throughput == weakest link
⇒ pipeline stage == thread pool.
- shared buffer based comm.



- sequence of stages
- stage == subtask
- each stage == thread pool
- buffer based comm.

Ⓐ specialization & locality

Ⓑ balancing and sync. overheads.

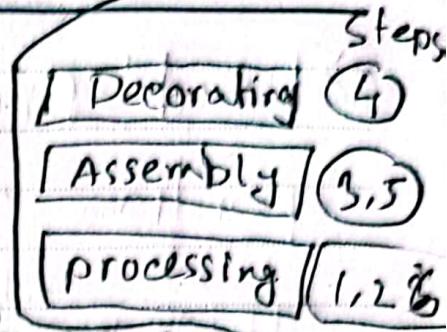




⇒ Layered Pattern:

→ each layer group of related subtasks

→ end-to-end task must pass up & down through all layers.



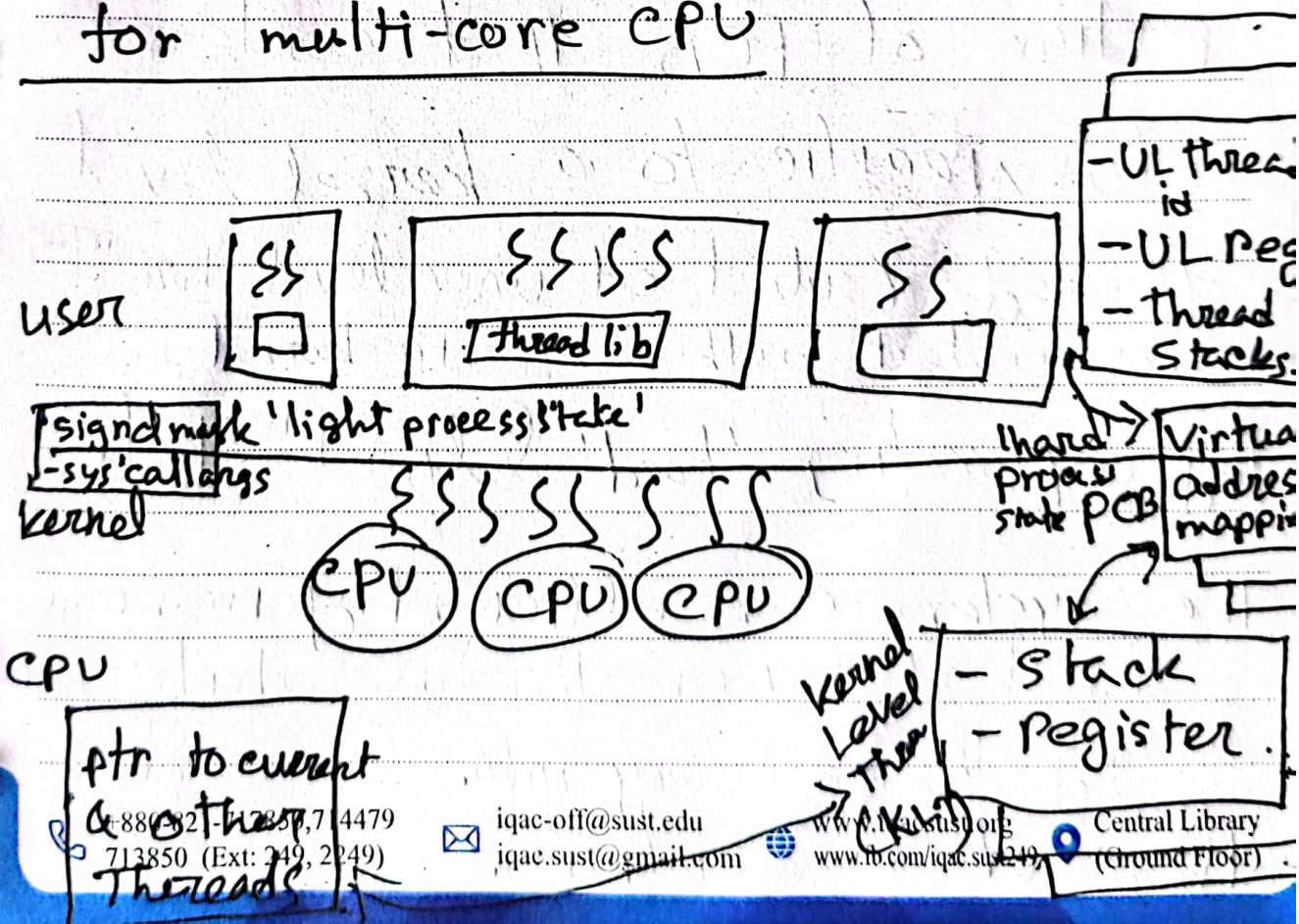
⊕ specialization ⊕ less-fine-grained than pipeline

⊖ not suitable for all applications

⊖ sync.

* Thread Related Data Structure:

for multi-core CPU





→ 'hard' process state: refers to the portion of the PCB that contains information relevant for all user-level threads within a process. This includes virtual address mapping, process id, status etc.

It includes info that needs to be preserved during context-switching.

→ 'light' process state: refers to the portion of PCB that contains info. specific to a kernel level thread. (subset of user-level thread associated with a specific kernel level thread). By separating 'hard' & 'light' state

The system can efficiently manage & switch between threads while preserving the necessary info. for each thread execution.

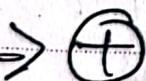


* Rational for multiple Datastructure

⇒ Single PCB: large continuous DS, private for each entity, saved & restored on each context switch, update for any changes.

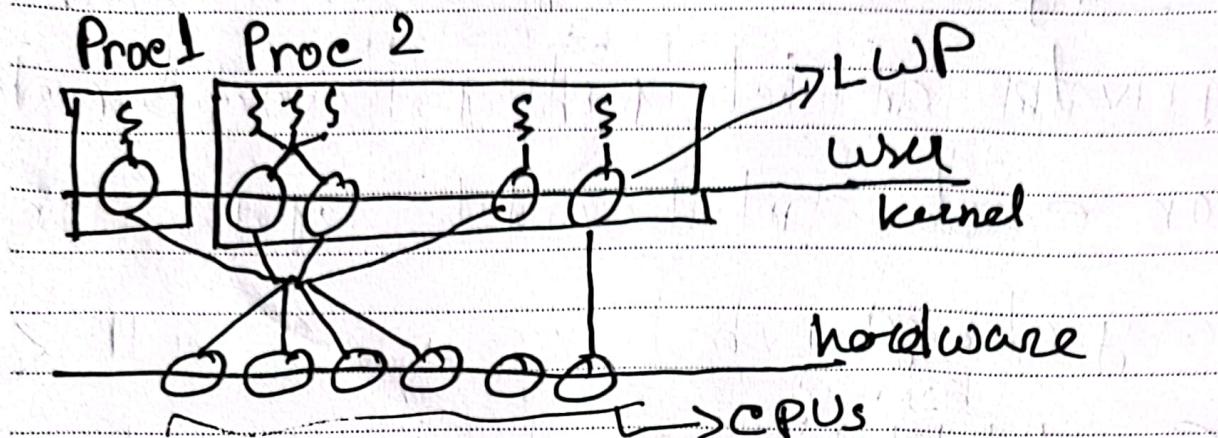
- ⊖ Scalability, ⊖ Overheads, ⊖ performance
- ⊖ flexibility.

⇒ Multiple DataStructure: smaller DS, easier to share, on context switch only saved and restore what needs to change, user-level library need only update portion of the state.





SunOS 5.0 Threading model / Solaris 2.0



→ both many to many & one to one model supported.

→ Each kernel-level Thread that executing a user-level thread has a light-weight process Ds associated with it.

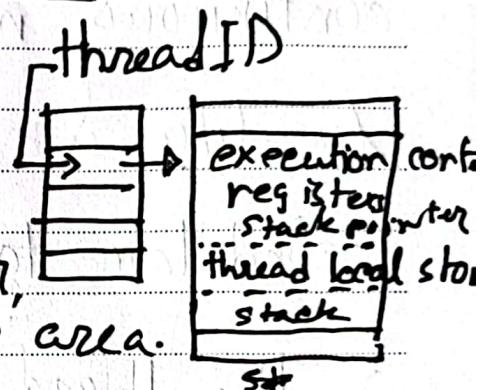
→ from User-level perspective, these LWP represents the virtual CPUs onto which its going to be scheduling the user-level threads.



→ At the kernel level, there will be a kernel-level ~~manager~~ scheduler that will be managing the kLT and scheduling them into physical CPUs.

④ User-level thread DS

→ Contains various fields such as, registers, signal mask, priority, stack pointer, and thread local storage area.



→ When a thread is created Threading lib. returns a thread-ID which is an index in a table of pointers rather than direct pointer to a thread DS.

→ This indexing approach allows for meaningful feedback or error message if there is any issue with thread, as opposed to pointing to a corrupt

~~memory~~



- This also includes The stack, which can be of defined size by the library or the user can provide it.
- This DS is created in a ~~manner~~
~~& layered~~
continuous manner, enabling better locality & making it easier for the scheduler to find next thread
- The threading library does-not control stack growth, the OS is unaware of multiple threads user-level thread, which can lead to overwriting another's DS
- To mitigate this, a 'red zone' is introduced, which refers to a portion of the Virtual add. space that is not allocated. If a thread tries to write



in this zone, the OS will cause a fault, making it easier to identify and debug the problem.

* Kernel-level DS (muliple DS)

Process

- list of kernel-level threads
- virtual address space
- user credentials
- signal handlers

Light-weight Process (LWP)

- user level registers
- system call argument
- resource usages info
- signal mask

Similar to VLT, but visible to kernel, not needed when process not running

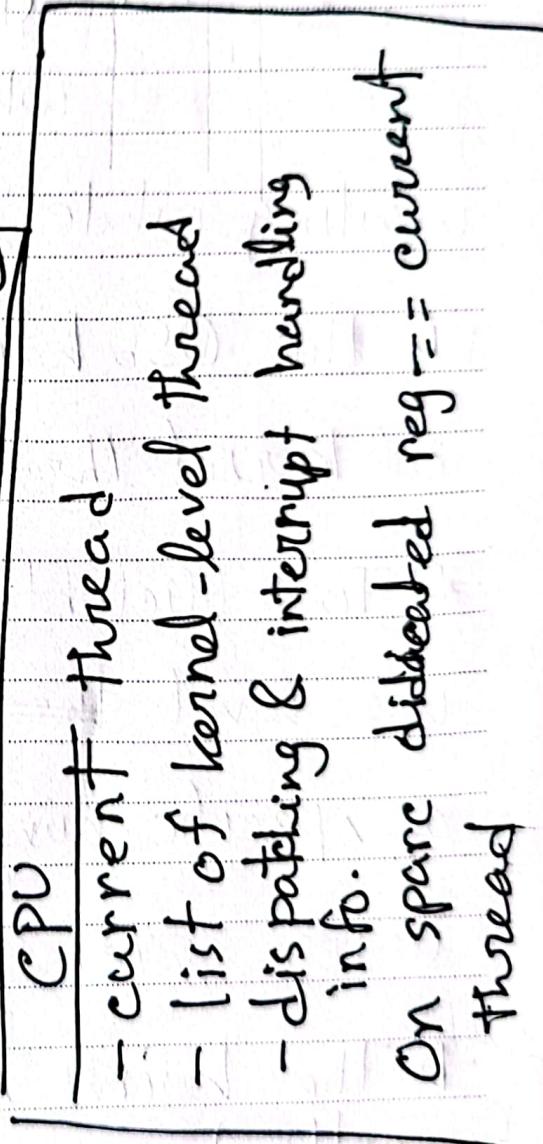
kernel-level threads

- kernel-level registers
- stack pointer
- scheduling info
- pointers to associated LWP, Process, CPU Struc.

- not swapable

1380827478, 1380827479
13850 (Ext: 240, 2240)

always needed even when the process is not running



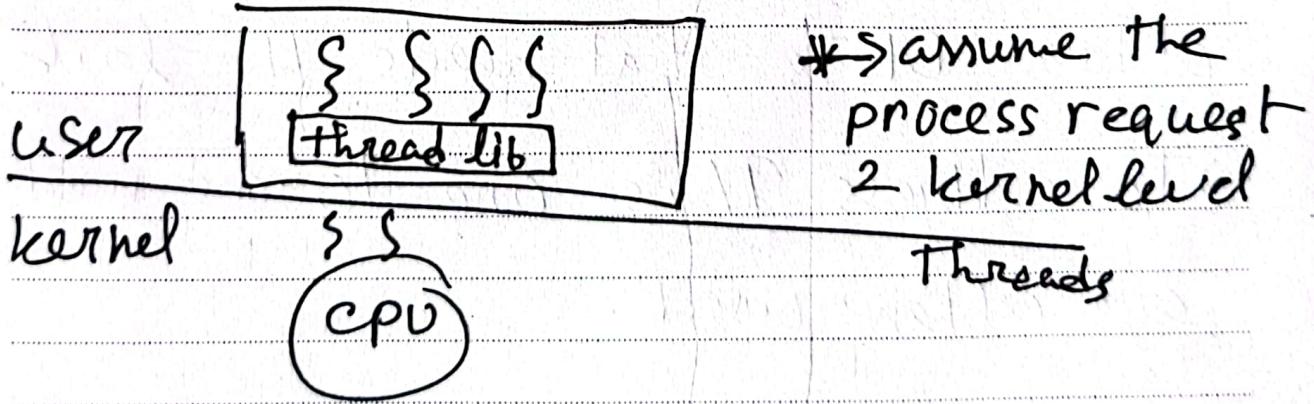


* Basic Thread Management Interaction

- In a multi-threaded process, There can be multiple user-level threads that executes concurrently.
- The actual level of concurrency may be limited by the factors such as I/O Operations, where some thread may be waiting while others are executing
- The OS - has the limit on the number of kernel-thread it can support
- To efficiently manage threads, the user-level ~~process~~ process can request a specific number of kernel-level thread using system called thread-concurrency.
- The kernel responds to this system calls by creating additional threads.

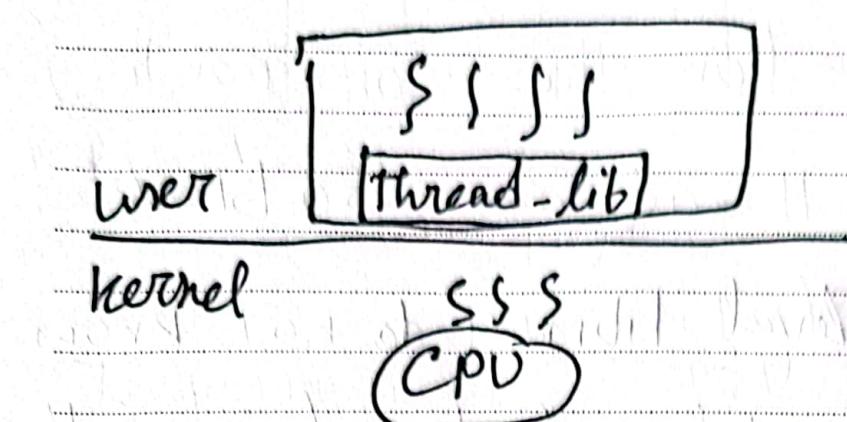


- If the user-level threads are blocked on I/O operation, the corresponding kernel-level threads are also blocked.
- The user-level library does not know the blocking status of kernel-level threads which leads to a situation where the whole process is blocked.
- To address this, the kernel can notify the user-level library before blocking the KL threads, allowing the library to request more KL threads or lightweight processes.





* Thread management Visibility & Design



UL library sees -
- ULTs
- available KLTs

kernel sees
- KLTs.
- CPUs

- KL scheduler

This cause,

* lets assume, if a kernel supporting the execution of a

Critical section preempted by the Kernel

The execution of that critical section

can not continue until the kernel schedules that thread again. This may lead

to delays in the execution of

other ULTs that needs the lock.

→ To address this one-to-one model

are often used.





Process jumps to U Library scheduler

When :

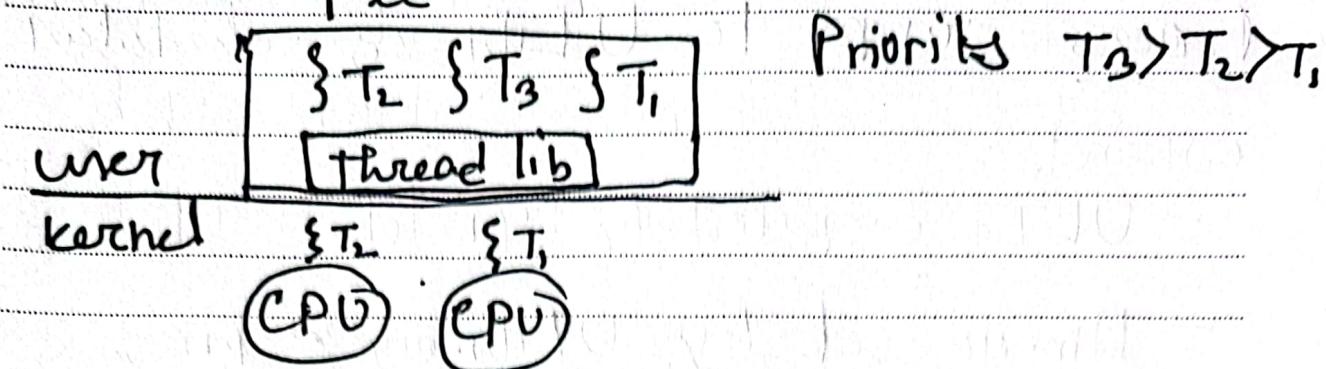
- ULTs explicitly yield
- Timer set by ULibrary expires
- ULTs call library function like lock/unlock
- blocked threads become runnable.

Issues on multiple CPUs

This arise when dealing with multiple CPUs in a system. In a single CPU system user threads run on top of that CPU and any scheduling changes are immediately reflected on that CPU. However, in a multi-CPU system, KLTs support a single process may be running on multiple CPUs concurrently. This introduce complexities in managing user-level threads.



For example



T_2 running on One CPUs and holds a mutex while T_3 waiting on that mutex and are on block. T_1 is running in another CPU.

release the mutex, T_3 become runnable and all 3 threads are now runnable.

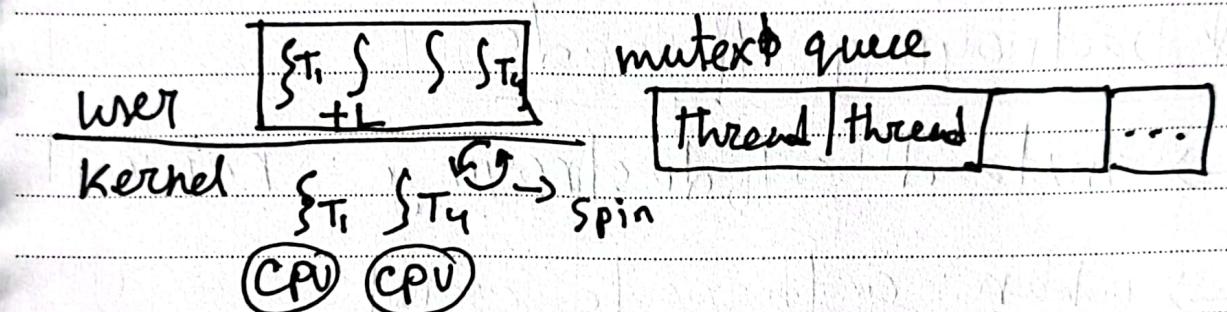
Now, to ensure the priority, T_1 needs to preempted. As T_1 is running on another CPU, we need to notify CPU to update its registers & PC.

This is done by sending a signal or interrupt from one CPU to the other. instructing it to execute necessary library code



* Synchronization-related issue:

In a multi-CPU system, when multiple threads needs to access a shared resource protected by a mutex, it is possible that the owner of the mutex to be running on different CPU.



In such case it may be more efficient for the waiting thread to spin instead of blocking and queuing up on the mutex.

This is effective if CS is short.

→ Adaptive mutex are the type of mutex that dynamically decide whether a thread should spin or block based on the current state of the mutex and



They are designed for multi-CPU system and provide balance between spinning and blocking, depending on the characteristics of the cs. For long critical section apply default blocking behaviour.

* Destroying Threads

→ Instead of destroying - reuse.

→ when a threads exist.

- put on a death row

- periodically destroyed by reaper thread

- otherwise Thread structure / stacks are reused - performance gain!



Interrupts & Signals

Interrupts

Signals

- events are generated → events triggered by externally by components other than CPU (I/O device, timers)
- determined based on physical platform → based on the OS.
- appear asyncro. → appear sync or async.

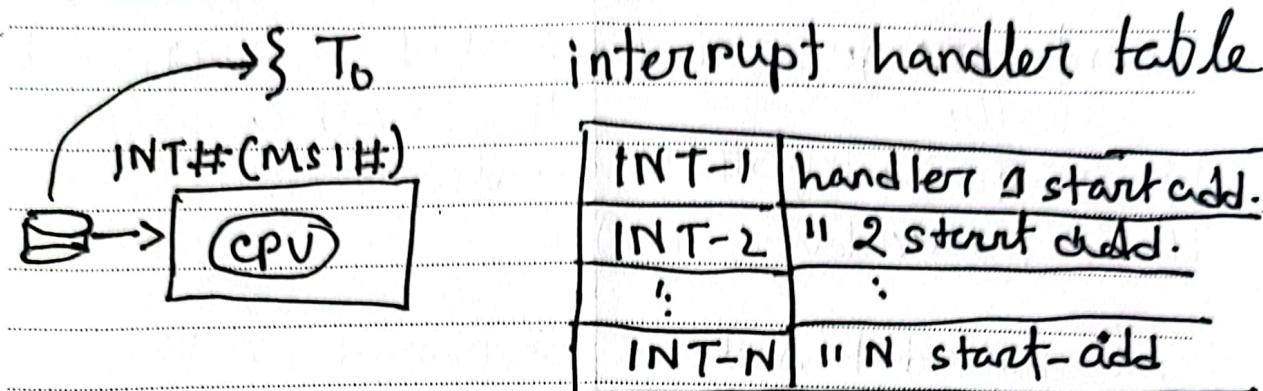
Both have

- unique ID depending on the hardware (intru). or OS.(sig)
- can be masked and disabled / suspended via corresponding mask
 - per-CPU interrupt mask, per-process signal mask
- If enabled, trigger corresponding handlers
 - interrupt handler set for the entire system by OS
 - signal handlers set on per process basis, by process.



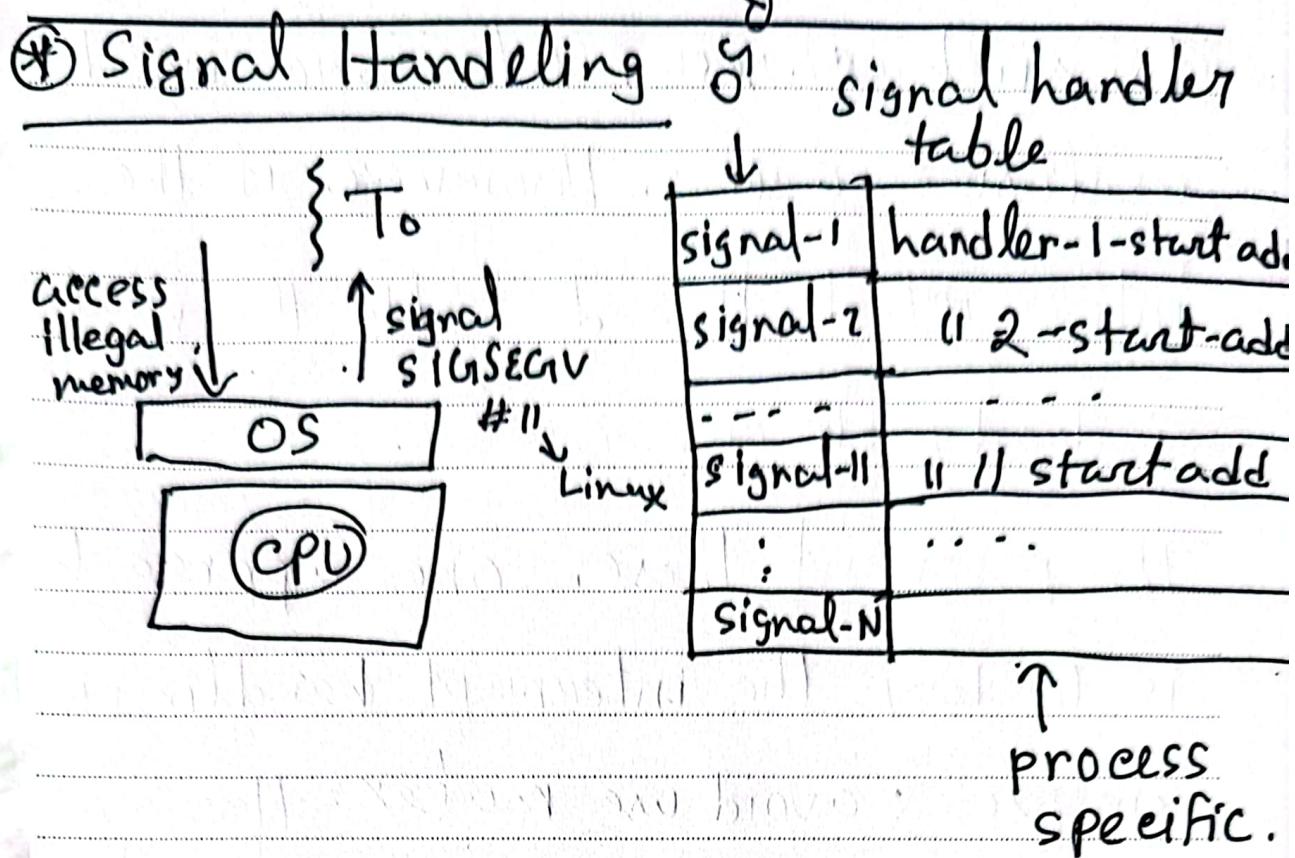


④ Interrupt handling



When a device, such as disks, wants to send a notification to the CPU, it sends an interrupt that connects the device to the CPU. In modern system this can be done using a message-signal-interpreter carried on the same interconnect.

This interrupt stops the execution of the thread that was running on the CPU. If the interrupt is enabled, a table is referenced based on the



⑤ Why disable Interrupts or signals

→ When the interrupt or signal occurs, it is executed in the context of thread that was interrupted, & this can lead to a issue.

→ One example is that, when the interrupt handling code needs to access shared state that other threads in the system may also be accessing.



In such cases mutex are used for exclusive access. However, if the interrupted thread holds the mutex, a deadlock can occur.

To prevent these, one approach is to keep the interrupt handling code simple & avoid use mutex. However,

these can impose limitation on hand.

Here masks can be used. These mask allow dynamic enabling or disabling of interrupts or signals during specific sections of code.

→ Interrupt masks are per CPU

If masks disables interrupt \Rightarrow hardware interrupted routing mechanism will not deliver interrupt to CPU.



\Rightarrow signal marks are per execution context (ULT on top of KLT)

gf marks disables signal \Rightarrow kernel sees the mark and won't interrupt corresponding thread.

④ Interrupts on Multicore System

- Interrupts can be directed to any CPU that has them enabled.
- may set interrupt on just a single core \rightarrow avoids overheads & perturbations related to interrupt handling from any other CPU/cores





④ Types of Signals

① One-shot signals: These are standard signals. They are designed to be handled at least once when multiple instances of the same signal occurs. The handling routine for one shot must be re-enabled every time this is invoked.

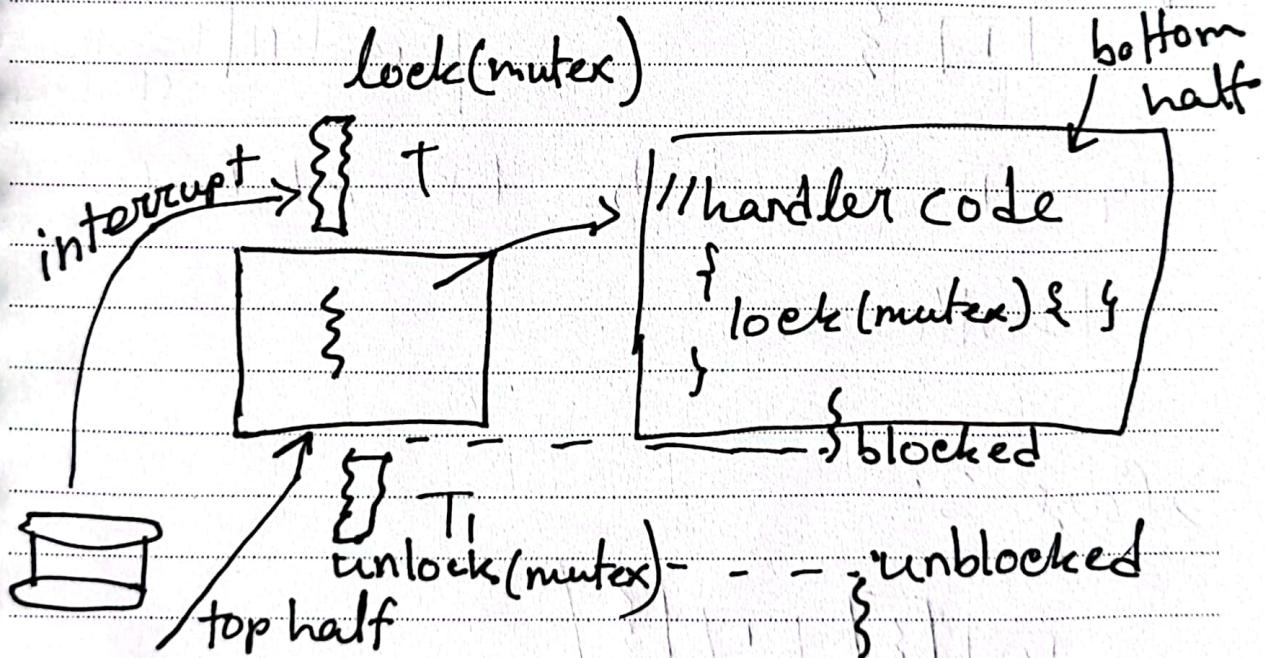
② Real-time signals: Every signal will be handled.
 n signals $\rightarrow n$ handler funcs.

* Top half & bottom half: refers to a technique used in OS to handle interrupts efficiently & avoid potential issues like 'deadlock'.



The interrupt handling process is divided into two parts: The top half & bottom half.

The top half of the interrupt handling code is immediately executed & the bottom half is a separate thread or task that is responsible for performing complex & time consuming operation.

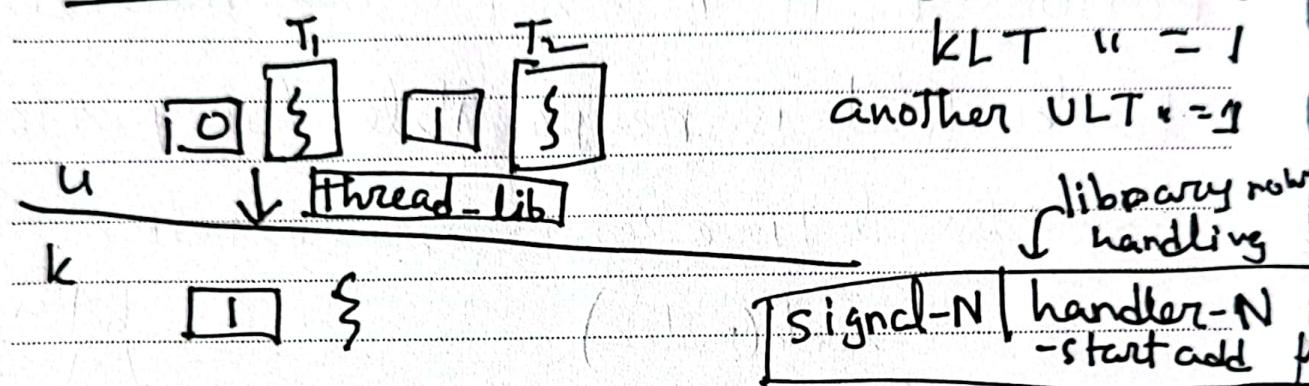




* Threads & signal Handling

\Rightarrow Case 1 : (both mask is one)
 $ULT\ mask = 1$
 $KLT\ \text{in} = 1$
so will interrupt
The process & complete the interrupt.

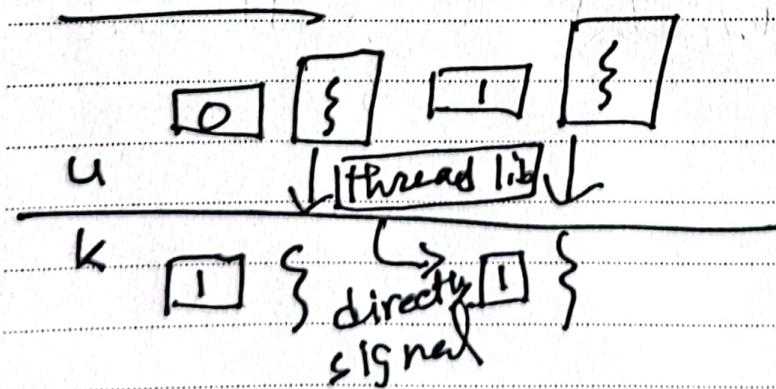
\Rightarrow Case 2 :



stops T_1 & run T_2 on KLT to

handle signal.

\Rightarrow Case 3 :





Case 4 :

