

Table Of Contents

Name	Page Number
Pre Mid 1	2-20
Mid1-Mid2	21-35
Mid2-End	36-54

1) Processor:

* Interrupts:

- Hardware signals (INTR signal)
- They call Interrupt Service Request (ISR, subroutine) in the code. (Each interrupt calls a specific ISR)

* IVT (Interrupt Vector Table) ($\times 36$ entries) It stores all the ISR's addresses in order (so they can be addressed) (This IVT is present at lower end of memory) 0x00... to ...;

* Note; Polling

→ Continuously checking status registers to check if it's free. ~~to accept~~

Interrupt driven:

→ Unlike polling it doesn't waste clock cycles.

2) Memory:

- Supplies instructions & data to the processor.

3) I/O Devices:

- Hard drives etc

Note: BIOS is present in our memory which has the bootloader (microcode), so the system can be boot strapped.
ROM or Flash RAM

Note: Boot Sector has the loader which loads the microloader in BIOS loads the microloader in (sector 0) (level 0) It resides at the top of hard disk. (On Boot, the boot loader is loaded into memory first by BIOS, then the boot loader loads rest of kernel)

Note: Read real mode vs protected mode of x86 processors.

Which ISR to call?
** Note: There is only one INTR pin (hardware). So we need to know which ISR needs to be executed, when we get an interrupt request on that pin (i.e., which device generated the interrupt).

For this we give all our devices control to the address bus of the processor, so when a device sets the INTR, it puts its address on the address bus, so we can figure out which ISR to use.

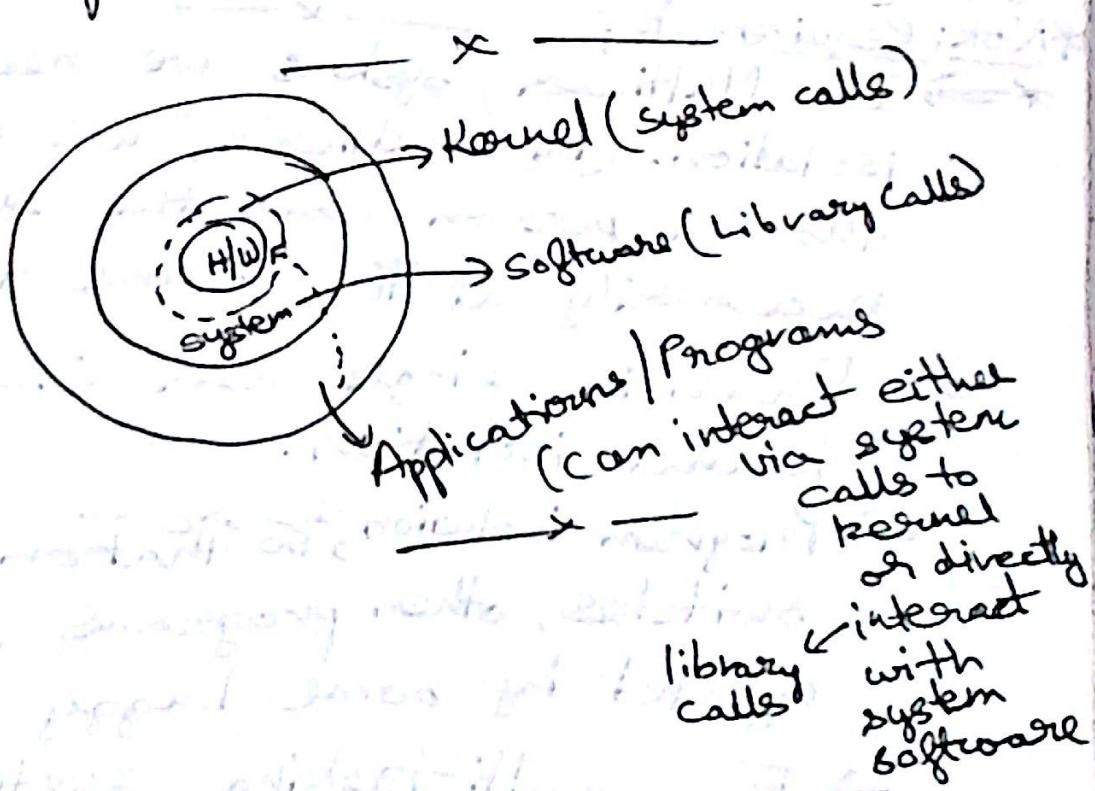
There exists a master ISR, which will tell us which ISR to use.
If multiple devices interrupt, then we use a bus arbitration protocol to decide which device to process first.

Note: OS = Kernel, Programs talk to the kernel only via system calls.
(Kernel has control of hardware) → (BASH) → (Shell does this for us)

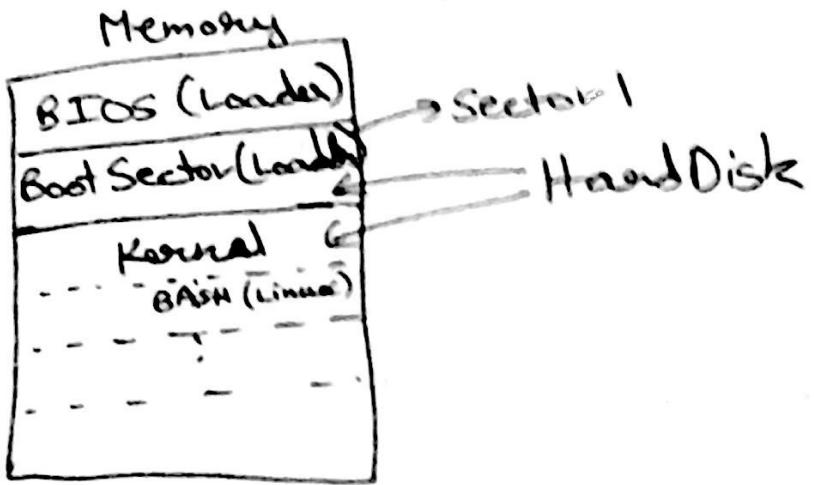
- Multitasking OS; (Mechanism to do this is a context switch)
 (Policy is how we assign priorities)
- When a program requires data ~~read~~ ^{out} from hard disk & other needs the processor, then they can both be run simultaneously & when the data read is done, an interrupt is generated & the processor attends the first program again.

(If we only had polling, then the processor would necessarily be busy, so it couldn't run the 2nd program on the processor).

Note;



Bootstrapping:



(MS DOS)

- Note: Micro kernel → ~~bootstrap~~ shell is within kernel
Monolithic kernel → Shell separate from kernel (Linux)

* (OS Requirements)

* Note: Requirements;

⇒ In Multiuser systems we need user isolation. (By multiuser, we mean more than 1 person uses the system, not necessarily at the same time).

⇒ Even in single user O.S, we need kernel isolation.

⇒ Program isolation, so that on context switches, other programs don't get affected by some buggy program.

⇒ From multi-tasking systems, we try move to time sharing ~~systems~~ systems -

⇒ File system is required.

⇒ Device abstraction (Makes it similar to files)

* Time Sharing System;

- Similar to a multitasking system but the context switch is done based on an interrupt timer (So all processes are given equal chance to execute)

* Process Abstraction; (Virtualization)

- Done mainly via virtual memory abstraction.
- Each process also thinks the entire processor is for itself because of CPU multiplexing (Time Sharing System)

Note: Each process has a page table (per process page table), which maps virtual memory to actual memory map. This page table is protected, so only the kernel can change it.

* Inter Process Communication

- A process can tell the O.S that it wants to be mapped to some physical memory to which other program is mapped to, but the program cannot specify the location directly in the ~~actual~~ ^{physical} memory to map to.

* Note; Each program has a ~~copy~~
page table entry
~~page~~ in the virtual memory
that corresponds to the kernel
in the physical memory.
(Since every program needs
access to the kernel)

* OS Requirements:

- 1) Multi user, Multi process, reactive
- 2) Process isolation
- 3) Kernel isolation
- 4) File system
- 5) Device abstraction
- 6) Efficient & fair ~~resource~~ allocation
- 7) Inter-process communication.

Note: Useful Commands

- 1) htop {Memory usage}
- 2) sar -B f → time {To see page faults}
- 3) sszie --format=sysV a.out

~~(Base-Bound Approach)~~ (Does not answer how all processes reside in memory)

* Memory Management, (Without Virtual Memory) (Old Systems) ~~so we go for~~

- ⇒ Every core has a base & offset register, when a context switch occurs & a process is loaded, the two registers are set ~~and~~ from the ^{the proc structure} of the process.
- ⇒ The hardware checks for every data instruction if we are within range else it rises an exception.
- ⇒ Only the kernel has access to base & offset registers.

~~Privilege levels~~ ~~in a memory structure~~ ~~and system~~

- 1) User Mode (Programs run in this mode (General))
- 2) Privileged mode → Can manipulate hardware registers etc
(Context switches are done in this mode)

* Note; The proc structure remains in the kernel (It has ^{the} number of ~~processes~~, that a processor can run at max) (Array).

(So for every process it can set values in the proc structure, like the "nice" value)

Priority values for Process

Note: The proc structure also has two variables base & offset assigned from hardware.

Note: (Memory management) *

When we swap out a few process to create a new one, then it would be more efficient to swap out just a part of the process to be more efficient.

e.g. A = 100 MB & B = 100 MB

We need to add C = 150 MB
then it would be efficient to swap A & half of B out.

(on context switches) *

* Note: Schedule function is when kernel modifies base & offset registers of each process which its bringing into memory.

- When we schedule, we also store the base & offset registers so that when we get back, we can acquire the base & offset registers again.

* Program Relocation Problem; (~~Static~~)

1) Relocation Table; (static)

- After bringing a process back after a context switch if we place it in a different location, then we need to modify base &

offset registers, and the data in proc structure.

(This is solved if we use PC relative addressing mode, but we need a solution with absolute addressing).

→ Solved with Relocation Table;

- Relocation Table stores a patching offset and when the program comes back in, we use that to maintain our registers. (For all instructions we have to patch
Software patches all instructions → offset)

2) Segment Addressing Mode; (^(dynamic) Segmentation registers)

- ~~All~~ Data, Stack, Heap and Instructions segment have ~~base~~ a base & offset, so when we move ~~out~~ out & back in, then only corresponding segment registers are modified appropriately.

^{Hardware} (Here instead of patching all instruction while executing, it just adds the value of Data Segment | whichever segment needed from hardware registers & executes instruction).

(Due to all these problems, and inefficiency we move over to Virtual memory & paging).

* Virtual Memory;

- Each ~~process~~ ^{process}, assumes it has the entire memory. ^(address space) or more for itself as ~~a~~ virtual memory. This memory is divided as pages.
- Physical memory is also divided into pages.
- Process shares address space with kernel.
- ~~Process~~ ^{Kernel} isolation must happen with virtual memory.

* Process Memory Map; (Segments)

- 1) Code
 - 2) Heap { Dynamic variables }
 - 3) Stack → Base address is randomized
 - 4) Data { Global Variables }
- Static
(Does not change)
(Base address are static)

Note;

* With this virtual memory, we won't have problems like

- a) Over-allocate
- b) Under-allocate

c) Granularity of process (Swapping out by in at context switch)

of memory for processes.

- * d) External fragmentation

*Locate Free Memory; (In Actual Memory)

- 1) First-Fit } Just traverse all your list of free memory locations.
- 2) Worst-Fit } Heap to just check & allocate maximum free space location.
- 3) Best-Fit } Nearest but greater than required.

*External Fragmentation; (This problem occurs when we don't use virtual memory)

- When small chunks of free memory is present here & there, we do Compaction so that we move the free memory to a certain location.

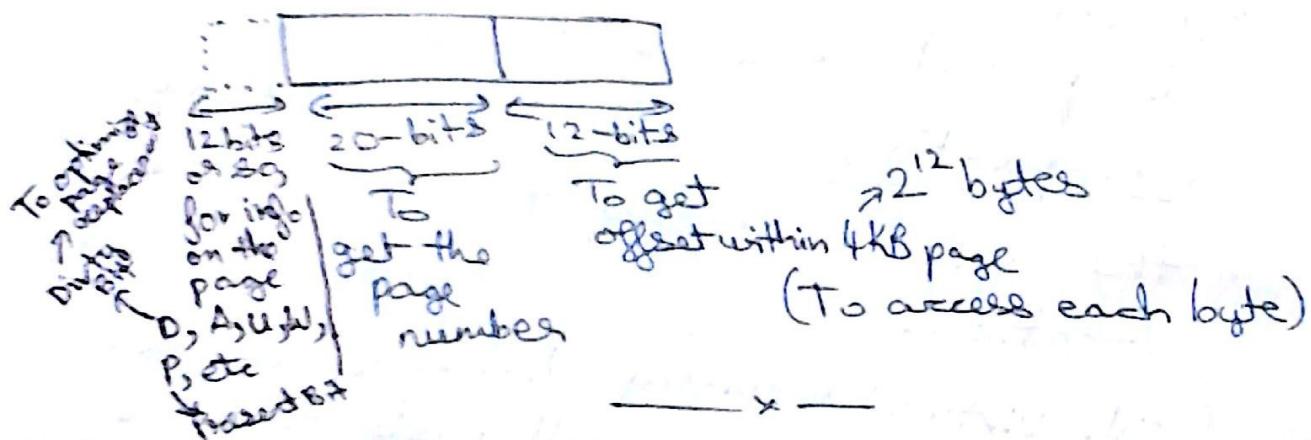
(With virtual memory however, we can use all those small parts of free memory without the need of compaction) (Due to paging)

(These are all the issues we face, so we move on to virtual memory)

* Paging; (Virtual Memory)

- Logical memory is divided into pages (in order of KB).
 - Physical memory is divided into frames \rightarrow of KB.
- Outside processor, memory is always (1MB) physical address based (so we can access memory properly), but within processor, they stay as virtual & physical address.
- ∴ \Rightarrow Page-Table Base register holds physical address.

* Note; We can use a 32-bit identification for accessing pages



* Address Translation;

- The offset bits remain the same.
- The other bits are fed into the page table, and at that location we get corresponding bits to the page number of physical address.

(So we can use this page number & offset, to access the byte in physical memory).

* Note; Page table is stored in our physical memory.

→ So in the proc structure,
 (CR3 register) contains the base physical address of the page table.

Always has correct base physical address of P.T

As we context switch a process, then the CR3 register is overwritten with the base physical address of the page table of this process (new).

* Note; Each process has its own page table & its base physical address of P.T is stored in the proc structure.

Stored in kernel space of the physical memory (RAM)
(When more kernel space is needed, other memories can be marked as kernel space)

* * Note; Page Fault; (Page number in page table entry)
• When the present bit is 0 or the page number is invalid in the page table for some virtual address, then a page fault occurs & the kernel allocates a physical page for the virtual page.

(It can be invalid even if it got swapped into the swap space).

* (If this happens, the process is put in blocked state & the kernel fetches the info back (then running state))

* Swap Space; (In hard disk)

- When we require more pages than the physical memory, or the physical memory is full, we swap out some pages into the swap space of our hard disk.

————— x ———

: Note; In x86 processors,

Physical Address Extension (PAE) is used to access more than 4GB gram (2^{32} bits exceeded)

To address data.

————— x ———

* Note; Every context switch, we will have to erase the cache. Since the new process doesn't use the cached data of the previous process

Note;

————— x ———

* Process Sharing; Virtual page addresses of different processes can point to the same physical page address if they need to share memory.

————— x ———

* Internal Fragmentation; (A small downside of virtual memory)
(At max 1 page is wasted per process)

- When the fragmentation occurs not with pages but within pages (ie. of 4KB page only a part of it is used)

————— x ———

Note; a.out files etc are ELF type files.

Note; bss section contains uninitialized variables.
* data section contains initialized variables.

Note; a.out files become large if there are lots of initialized variables/memory.
(size --format=SYN a.out to check memory map).

* Note; Proc filesystem is not actually present on the hard disk, but its generated as a request is made to it on the fly (no file made). Therefore we can't modify these "magic" files.

~~Note; Physical Address Extension (For larger memory access (RAM))~~

Note; When page fault occurs, the same instruction is run again and the executing instruction is aborted.

* Working Set of Process;

- The working set of a process is the minimum amount of memory needed for the process to run. (without swaps to disk).

* Thrashing;

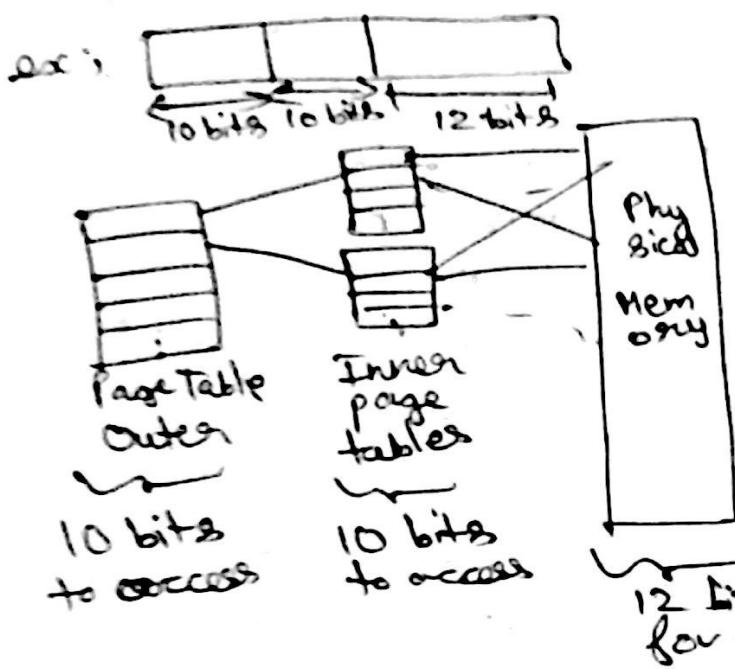
- When $\sum_{i=1}^n W_i \geq \text{Physical Memory}$ Working sets

- Then many swaps to disk happen so the system is slowed down a lot, this is known as thrashing.

Note: Swap space is optimised to access by the OS. (So pages can be accessed quickly & efficiently). The file system structure of swap space is different

* Two level page table; (Disadvantage is 3 memory access)

- Page the page table itself,



(Outer page table
= page directory)

(Inner page table
= page table)

- This helps us save space since we need not store some inner level page tables if they are not needed. (holes are present, ie. they don't use those addresses), so we don't need to store page tables at those addresses)
- We can also swap inner page tables into swap space if needed.

* Note: Inverted Page Tables (For M.d.l)

* Note: In Intel processors, Segmented paging is done (virtual address has segment & offset)
→ We get a translated 32-bit address from segment translation and then get physical address from page table.
⇒ The paging is optional (since there is no paging on boot time).

* Translation Lookaside Buffer; (TLB) (access in 1 clock cycle) (on chip structure)

- It's a small buffer which stores some page table entries, for fast access.
- If TLB has mapping, we use it else we use the page table and we also store that in the TLB.

* Note;

TLB Miss: When a mapping of virtual → physical address is not found.

Note: TLB's are associative memory.

* Note:
At a context switch, the TLB is flushed since there is only 1 TLB for all processes, and the mapping is different from one process's virtual addresses to another process's virtual address.

Page Replacement Policies

⇒ On a page replacement to swap space, the kernel stores where in the swap space the page is present in.

1) First in - First Out ;

⇒ This is bad as from Belady's anomaly, increasing the page table size (number of page frames) increases page faults.

* 2) LRU Replacement,

⇒ Replace page which has not been accessed for the longest time.
(Least Recently used)

⇒ this is an online algorithm

(There is another offline algo (we have entire reference string), so we remove ~~from~~ page which is ~~least~~ ^{farthest} used in the future).



* Note: Swap Daemon | Page Daemon

⇒ Pro actively swaps pages into swap space.

(kswapd man page in linux)

— x —

* Copy-On-Write (COW):

- When a process is forked, then the virtual memory map for child is a replica of parent's page table's ~~physical entries~~. So that we don't need to allocate physical addresses for the child. But they are all marked with COW bits.
- So when either parent or child try to write to that entry, a new location is allocated with the value copied to this new physical page (The old entry is removed in this ^{process} page table)

— x —

* mmap (system call) (Faster than fcreat)

- Given a file descriptor & file size from (fstat) mmap can copy the data onto some memory (addresses) in the kernel & we can access that memory to print to the screen etc.

— x —

* User process data structures

- 1) u-area (only when process is running, this is accessible)
- Has file descriptor tables (Files process is using)
 - Stores process control block (hardware context)
 - when process is not running.
 - Kernel Stack (Per kernel stack) → when kernel executes system calls.
- 2) proc structure

→ Process id

→ Location of u-area.

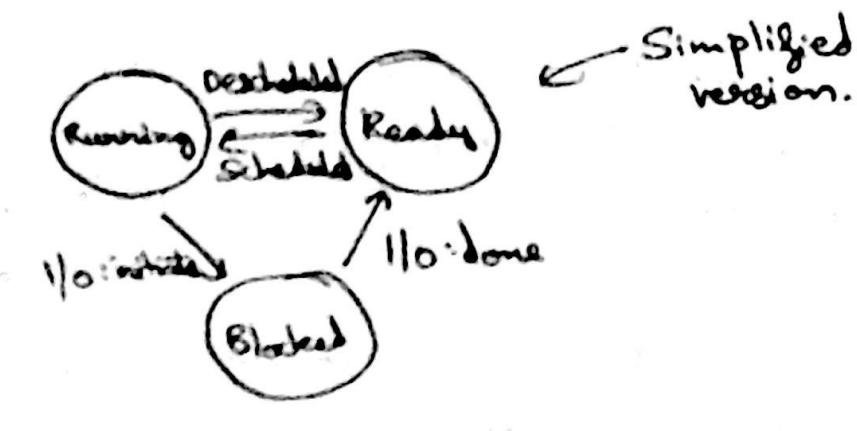
→ Current process state (READY, STOPPED, BLOCKED)

* Mode | Space | Context

Process	System calls, exceptions	Kernel (access system space)	mode
application (user) code			
user mode (access process space)			
not allowed	interrupts, system tasks	System Context	

Q) Why can't kernel also use user stack when executing system calls/etc in a process?

Process State Transitions:



* Note: When a sleep() is present in child after fork(), generally parent is executed first, but it is possible due to high load, ~~and~~ these processes are given the CPU only after α -seconds so now randomly the child can come first as well.

Scheduling; (Context Switches)

i) Preemptive Kernels; (Schedulers)

- When the time slice allocated for kernel is over, we context switch out.
(This has problems)

* (Synchronization problem: Suppose we are removing a node from double linked list, if we get swapped out in between, we might not have completely removed the node yet.)

Locking techniques are used with multicomputer systems.

2) Non-preemptive Kernel, (Scheduler)

- Another process cannot preempt the kernel even if its time slice expires.
- This solves problems seen before (signaling problem).

~~but~~

- But there are better ways to do this, which we will see later.

~~So if interrupt occurs, then the costed context switch is still making, it's interrupt (it's like preemptive).~~

Note: Scheduler is just a program, so who loads the scheduler? (This is the recursive problem (scheduling problem) we encountered before).

- we want to call the scheduler when,
 - a) Time slice expires
 - b) I/O request
 - c) Sleep

Note: In a context switch the scheduler stores process state of p1 & then loads process state of p2. This is the cost of context switch (we want to reduce this time).

* Scheduling Metrics:

1) Interactive Applications, (Shells, Editors etc.)

⇒ Response time ⇒ (first response time - arrival time)

2) Batch Applications (Compilers etc)

⇒ Turnaround time ⇒ (completion time - arrival time)

3) Real-time applications (Video players etc)

⇒ Strict deadlines missile tracking / reaction time

Note: In batch applications, if scheduling is done based on first come first serve, then if a long process comes first, it blocks small process that come later. So turn-around time is not minimized.
(This is the conway effect)

* Shortest Job First Scheduling; (Batch Applications) (Optimize turnaround time)

- We assume we know time each job takes, so we will always schedule the shortest job first.
- In this we will get the minimum turnaround time, so conway effect is not present.
- If a new process comes in which takes less time than the running process remaining time, then we execute the new process. (shortest job first).

* Fairness; (In Scheduling strategies)

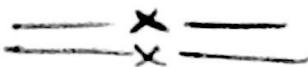
- Shortest job first scheduling is not fair always since a long job might end up starving if small process keep coming in.

* Round Robin Scheduling; (optimizing response time)

- We use time slices to allocate each process some CPU time.
- We access all jobs as soon as they come in so our response time is small, but

the turnaround time for this approach is not good.

- Increasing time slice length, turnaround time decreases but response time rises
- This uses a lot of context switches which wastes time.
- Starving does not occur here generally so it's fair most of the times.



* Multi Level Feedback Queues; (Scheduled in ~~in~~)

- With 2 jobs A & B,
(We don't know if it's batch or interactive)
 - a) If Priority(A) > Priority(B), then the process A is scheduled
 - b) If Priority(A) = Priority(B), then the processes A & B run in round robin scheduling.

* Priority Assignment,

- When job enters, it's given highest priority
 - If job uses entire timeslice, then priority is ~~dropped~~ ~~lost~~ → Batch Application.
 - If job gives up CPU before timeslice is over, it stays in the same priority level
- ⇒ This priority assignment ensures that interactive processes are getting round robin and batch processes are getting shortest job first scheduling (similar).

⇒ We also modify (decrease priority level) if ~~process~~ has ~~a~~ "CPU usage" which is very high, so that starving processes are addressed.

* Priority_{process}^{value} = Base priority + nice value + (CPU Usage)

(lower priority value → more priority)
 $CPU_i = CPU_{i-1}/k$

* Priority_i = $P(i) = \alpha (CPU_i + (1-\alpha)P_{i-1})$

Convex combination of CPU_i & P_{i-1}

Priority value when the process comes in to execute the i^{th} time. → (Lower $P(i)$ → more priority)

Current CPU usage time (Increases as time passes) ↑
 History (Decay Function)
 when process is using CPU History decays.

Note: Credit based scheduling

Schedule function; (of Kernel)

- When a running process goes to sleep or ends then, the processor goes from user mode to kernel mode.
- Then the schedule() function is called from the scheduler (The one with highest priority is scheduled next).

Note: On a context switch all user details of a certain process are stored in one memory & we just store an entry in our page table which we can access later when the process is brought back.

* Note: We can use tools such as 'perf' to see details about page faults of your program etc - (perf stat -e ~~cs~~ migrations file)

IV

* Note: If we don't use CPU consume times with decays, then a batch process can disguise itself as an interactive process by relinquishing the CPU just before its time slice ends so it's treated as interactive & has high priority. (But with accounting CPU consume time, the priority is dropped as more CPU is used (time duration) within timeslice)

1)

2)

3)

4)

5)

6)

Note: An application can either have a single process or a collection of processes.

*

* Note: When we execute a set of processes ~~instead of a large single~~ ~~parallel (in multicore)~~, then if one process of the processes is doing a file I/O, the other can continue, ~~executing~~, with the next process and so on. (In one busy process the CPU stays idle on page faults, etc) (So we can overlap I/O with computation).

No

* Note

Inter process communication (IPC)

1) Pipe; when all processes are on the same system.

2) Socket: Processes talk via the ethernet.

↳ Processes can be on different machine

3) Remote Procedure Calls (RPC)

4) Files (Bad performance)

5) Signals

6) Shared Memory

→ Effective alternative of processes (generally).

* Threads: → Output is non-deterministic
(No order threads thread)

- Each thread has its own stack but the same heap, initialized using uninitialized data segment etc.
(Address space is same).
- So they can independently execute the same process.
- Threads are implemented via the clone system call.

Note: Each thread runs on a single core of a processor. Generally ~~threads~~ → threads use core with least load, so so they keep changing cores as needed, but by specifying a certain argument, we can force the threads to run on the same core.

Note: We use threads since they are not as expensive as processes created by `fork()`.

Note: The alone system call is not very expensive, so `pthread-create()` etc use this system call.

- * Note: User Threads $\rightarrow [N:N]$ } The user sees it as N threads but the kernel sees as ^{threads} $[1:1]$
- * Kernel Threads $\rightarrow [1:1]$ We don't even get the advantage of I/O, computation etc.
- Since the kernel views all user threads as one, there are no actual context switches. (No performance benefits since kernel views it as only one thread).
- So nowadays, Linux only uses kernel threads since they get larger time slices compared to user threads.
(And we also have performance)

- Hybrid threads $\rightarrow [N:M]$ } Some O.S's use this.
User threads are managed by user-level library (Not in core).
Note: User threads are faster than kernel threads since they are essentially just a function call.

* Advantages/Disadvantages of Threads:

Processes

- 1) IPC is hard here
- 2) Each process has its own address space, so large address space.
- 3) Offers complete isolation since one other process cannot affect this process.

Threads

- 1) IPC is easy here.
- 2) Since the same address space is used, the stack is shared among threads, so smaller address space.
- 3) If the stack of one of the threads (malicious) can end up affecting other threads.

Since any halting system call from a single thread blocks all other threads since the kernel only sees one thread. (User threads model)

* (So user threads model is useless).

OS might also end up poorly scheduling these threads since it doesn't know about them.

* Note; To the kernel, threads are essentially just processes. (But they were created easily).

• Threads were together forming a process.

* Locking in threads; (Mutex)

- pthread_mutex_lock(&lock) } with a critical section
- pthread_mutex_unlock(&lock) } mutex

(*)

* Note; As we increase the granularity of the locks, (smaller critical sections), then the parallelism increases.

* Note; If a thread puts a lock & then on a context switch, what would the other thread do if there is a lock? (We will see later that there are 2 types of locks) → It just goes to sleep.

* Note: Process context switch is expensive since it changes address spaces of executing process whereas a thread context switch only the registers etc are changed, so its fast.

- * Note critical sections should be as small as possible (least number of them). (We should not have I/O, etc functions in this critical section).

————— x —————
↳ Since they can end up holding lock for too long.

* Coarse Grain vs Fine Grain Locks;

- Coarse Grain → ~~Less~~ parallelism
- Fine Grain → Too many calls to lock & unlock function calls.
(Almost 5x time).
in some cases.

* Variable Privatization (To reduce lock calls).

- * When our global variables used in our threads are reducible, then just use a local variable in each thread & then at last update the global variable.
(This might not always work).
- Compilers might also do this automatically.

* * * (In parallel program) ~~Critical Section~~ Cache Problem).

Note: Write- Invalidate-Store [Ping-pong] problem

- * Suppose we split / parallelize (try this!) a function which sums up integers $\sum_{i=1}^n j$.
- We tend to keep a global array, $lsum[4]$, and split our task to 4 threads, each updates $lsum$, and once all threads are complete

we just print the sum.

→ The problem here is that the entire lsum[4] is stored in a single cache block, so each thread invalidates the cache of other cores on write to its lsum[threadno]. So this wastes a lot of time (Parallel ~~will~~ ^{ending being} worse than serial).

- However if we use separate variables (most likely this problem will not occur. (So parallel ~~will~~ ^{will} be better or equal to serial) (time wise))

→ Or we can use local variables & add to global variable Variable ← at last with privatization. small lock.

=====

* Lock Implementation:

This can work in single core

- 1) If we are using a single core processor, disabling & enabling interrupts might work but it's dangerous since a crash would ruin the system. (This is bad)
⇒ Disabling timer interrupt ^{in program} can stop context switches

• Note: We cannot context switch out of an ISR until it's complete.

* Deadlock; Beware

- If our code locks a variable that an ISR needs to use, then the system goes into a deadlock. (ISR won't finish).
(This generally can happen when programming with kernel code).

*2) In lock call, make it wait until lock is available and then consume the lock.

→ This can be a problem if there are 2 threads waiting for a lock on different cores, so once its available, both the threads can end up entering the critical section.

- So this is a problem,

(Since they are on different cores, they run parallel)

*3) Peterson's Algorithm: (2 threads who want to access lock)

- Each thread tries to be generous.

(Initialize turn=0)

```
void lock() {
```

```
    flag[self]=1;
```

```
    turn=1-self;
```

```
    while((flag[1-self]==1)&&turn==1-self)
        wait here;
```

```
}
```

```
void unlock() {
```

```
    flag[self]=0
```

```
}
```

- Self = 0 for 1st process and self = 1 for 2nd process.

- This solves the problem since if flag[other thread] == 0, then the one thread which requires the lock

can use it as soon as its free.

- If both threads want the lock, then the thread whichever turns it is waits and acquires the lock as soon as its free.

*~~Fairness~~ Fairness; (Or Unfair)
⇒ It's fair since if two threads are trying to acquire a lock, as soon as 2nd process (1) ~~turn~~ acquires it & uses it, then flag(~~turn~~) becomes 0, so the other thread can enter the critical section. Even

if for some reason the 2nd process comes again & makes flag(~~turn~~)=1, it would make turn = 1 - 1 = 0, so the 1st process (0) would use the lock.

*4) ~~(Hardware feature)~~ Test and Set; (Not always fair)

- This works will any number of threads.
- Since the instruction 'Test and Set' is an atomic instruction, only one thread will get the lock at a certain time.

⇒ Test and Set ($\&\text{lock} \rightarrow \text{flag}, 1$) would do; old = $\&(\&\text{lock} \rightarrow \text{flag})$
 $\&(\&\text{lock} \rightarrow \text{flag}) = 1$
return old.

- All this is done atomically.
- ⇒ In the lock function,
while($\text{test\&set}(\&\text{lock} \rightarrow \text{flag}) == 1$)
we wait.

- So now only one of the waiting threads for the lock acquire the lock & set the flag to 1. (Atomically).

⇒ This is not always 'fair' though.

5) FetchAndAdd (This hardware instruction)

is much more expensive to add to hardware

- Another atomic instruction.

⇒ `FetchAndAdd(int *ptr){`
 `int old = *ptr;`
 `*ptr = old + 1;`
 `return old; }`

- On locking, each thread gets its ticket, by , myticket = FetchAndAdd (~~lock~~) (lock)

⇒ And it waits until it gets its turn ~~lock~~



- On unlocking, FetchAndAdd (&lock → turn) happens

⇒ So void lock (...) {
`myturn = FetchAndAdd(&lock->turn);`
`while (lock->turn != myturn)`
`.....(spin)`
`}`

void unlock (...) {

`FetchAndAdd(&lock->turn)`



After Mid2

* Problem with Peterson's Algorithm;

- When instructions are reordered in each of the processor, then the checking of peterson algo can occur before assigning the turn variable (suppose we use 2 variables he-wants & she-wants, then both threads can end up entering critical section).

⇒ There instructions, he-wants = 1

```
while (she-wants && turn == his)
    ||| Get
, she-wants = 1
;
while (he-wants && turn == his)
    ||| Get
```

⇒ They get reordered, so both threads end up entering critical region.

- Individually (within each thread) there are no dependencies b/w the instructions so re-ordering seems reasonable to the processor, but concurrency becomes a problem

* Solution: (Hardware instruction) (Fence)

⇒ We use memory fences, to solve this problem. When we add a memory fence instruction, then this forces the CPU/compila to ensure that no instructions are re-ordered across the fence.

- she-wants = 1;
turn = his;

Fence → sync-synchronize();

while (he-wants && turn == his)

;

- he-wants = 1;

Fence ← turn = hers;

sync-synchronize();

while (she-wants && turn == his)

;

Note: We will also face a problem if the two instructions ($\text{she-wants}=1$, $\text{turn}=\text{his}$) are interchanged as well, so we would think of using another fence, but memory model tells us that writes won't be reordered with other writes. (So we won't get this problem)

=====

* Spin vs Yield;

- If we will get the lock in the near future then yielding is pointless since it wastes context switches. Spinning would be better in this case.

```
1) void lock(){  
    while(i){  
        if(TestAndSet(...)==1) count+=1  
        else return;  
        if(count==1000){  
            count=0; sleep(1);  
        }  
    }  
}
```

- This algorithm is fine but going to sleep no process is sure when it will get the lock - (No fairness)

(Its chance to get the lock might have been lost when it was sleeping)

(Code in slides (Vicking locks using Queues))

*2 To ensure fairness, we add a queue, and, if threads don't get access, we park them in the queue. Then when the lock is released, the unlocking thread unparks the thread at the head of the queue.

⇒ Park() puts the thread to sleep

ready to be awoken when Unpark() is called.

- To maintain consistency in the queue, we use a guard lock when modifying or accessing the queue.

Note: Using a spin lock for the guard is fine since one critical section is very small.

a) Problem with this: If thread is added to queue and guard is set to 0, and then before parking context switch occurs, and another process unparks this and removes it from queue, and the context switching back to the original thread, it parks and the process stays parked forever.

Solution: We add a setpark() before

making m→guard=0 and then park().

(Setpark() tells the kernel that we are

going to park soon so don't context switch until we park()).

Note: These park(), setpark() etc are Solaris based not Linux.

* Condition Variables (Use volatile ints for global) (Only one)

→ We could just spin wait until child finishes execution with a global variable but it wastes time. (Due to spinning).

* 1) We use pthread_cond_signal & pthead_cond_wait since this mechanism is a lot faster. Semantics specify, that we need to use a lock around these functions.

* The pthead_cond_wait releases the lock before going to sleep and then when it is signalled by the child thread, then it will wait until lock is released by the child & then acquires the lock again. ~~(This is done atomically)~~

⇒ So a single function pthead_cond_wait releases lock before sleeping then acquires lock on waking up. (Atomic)

Note: We don't want to call pthead_cond_wait if there is no child to signal, so we use a global variable 'done'. (In code), so that we call the wait only if $\text{done} = 0$.

Note: If we don't use locks with these condition variables, then context switches can mess things up. → See code without locks to understand why.

⇒ Producer Consumer Problem; (1:1)

- When there are 1 consumer & 1 producer, then we can use a single condition variable.
- If item is read, writer is signalled. If writer writes reader is signalled.

* (Linearizability) \rightarrow
 * Linearizability; (Isolability) (Correctness of concurrent programs)

- This is required in concurrent processes.
- The manufacturer specifies memory consistency models, ex: Relaxed, Sequential etc. → These are always satisfied, but we need to add additional synchronization primitives in order for our concurrent programs to work.

~~=) Not reordering all instructions & restricting this reordering is called linearizability.~~

* Note: We say a datastructure etc. is linearizable if at any instant, a snapshot of the system would give us an interleaving of instructions from multiple threads.

We want

valid reordering

(Without unwanted reordering) without reordering which might make our system not work

⇒ We also if any.

For a linearizable, should be able to assume all instructions are atomic in our given ordering, then it's linearizable.

* ⇒ Even if 2 cores run instructions A & B simultaneously we need to guarantee that either A, B happened or B, A.

*Producer - Consumer Problem (N:M) (Code 1)

(Code 1) (No fairness)

(Condition variable like

When we use same conditional Variable

then producers can end up waking other producers or consumers wake up other consumers.

* Note: If there are no waiting processes, signal does nothing.

* * Note: Model checking (Exhaustive checking) is used to verify that our code is correct (Not a mathematical proof) (Exhaustive state space enumeration).

⇒ In code 1, suppose there are 2 consumers & 1 producer, then suppose C1 goes to wait, then P1 signals so C1 goes to ready to run state but waits for lock. Suppose C2 consumes lock after P1 releases it and then reads the data memory it. Then C1 acquires the lock but there is no valid data.
(This is a problem).

The way we fix this is by replacing "if" with "while" in both producer & consumer, so that only one concurrent consumer can consume the data. (So consumer awake consumer at just puts them to sleep again)

But no fairness ←

* Due to this we can end up going into a deadlock state. (when there are multiple producers or consumers).

* Suppose a producer is in wait, and a consumer ^{reads data, then} goes & signals another consumer ^{then goes to sleep}, so that would try to read data, but since there is nothing, it also goes to sleep, so all go to sleep and we get a deadlock.

(Note: We can fix this with 2 condition variables)

* Note: Signal \rightarrow Wait \rightarrow Ready To Run \rightarrow Scheduled \rightarrow Get Lock \downarrow

• The problem occurs when another process gets lock & changes state of shared variable.
Here we solved with while loop.

Come out of wait

* II) Using 2 condition variables & while loops our code works.

- A consumer only signals the producer and vice-versa.

* III) When there are n-locations in a circular queue. (Code 3)

- We use ~~a~~ similar code as II but we just change put(i) & get() functions

* Semaphores (Code in slide) (Reader-Writer problem)

- How are post, wait implemented in semaphores (please read source code).

⇒ We can implement the reader writer problem with semaphores.

* I) (Problem in code 1)

I) The code given must ensure that the put() & get() functions are guarded since they can simultaneously be accessed by multiple producers/consumers.

II) (Problem in code 2)

- If consumer sleeps with guard mutex, the producer cannot get the lock, so we get a deadlock.

* III) (Code 3)

- This code works fine, since we only lock critical section. (Generally it's bad to sleep/wait while holding a lock).

Note: Read about Model Checking.

*Dining Philosophers Problem; (Code in slides)

⇒ The problem is that all philosophers can hold their left fork & wait for right fork, so we get a deadlock.

Solution: Make one philosopher the leader who gets the right fork first whereas all the others get the left fork first.

_____ + _____

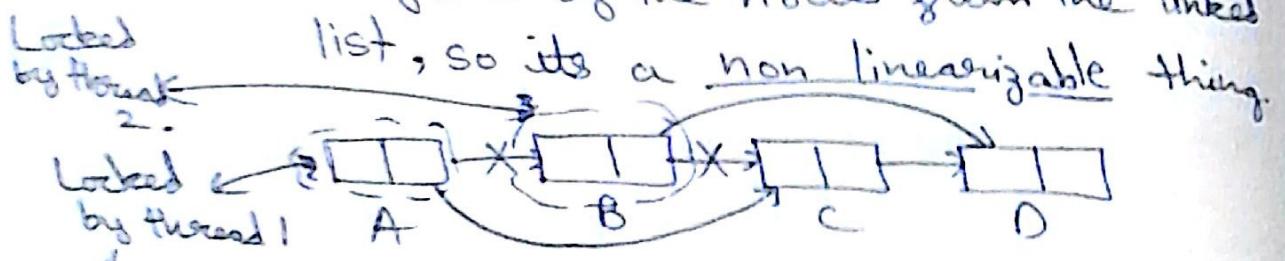
Concurrent Datastructures

* 1) Making Linked Lists concurrent; (with locks)

(Slides have diagrams)

- * 1) Using lock variables in the struct of a node itself. Each thread needs only 1 lock.
- ⇒ Suppose each thread can access functions like add, remove etc on nodes only when they have acquired the locks.

- When 2 threads try to remove adjacent nodes, we might end up removing only one of the nodes from the linked list, so it's a non linearizable thing.



So we can change pointers of A. ⇒ So now our linked list is A, C, D. So C has not been removed.

* 2) Use 2 locks (Hand on hand) (Slides for diagrams)

Same as above, but each thread needs to hold 2 locks to insert or delete a node. (The node & its predecessor)

⇒ So now, one of the thread would have to wait till adjacent deletion/additions have been complete.

(But if they are not adjacent, then both deletions/additions can happen parallelly).

Note: We could just use a global lock so only one thread can modify the list at a time but the level of parallelism would be very low then.

→ In our hand-on-hand method, the threads can act parallelly if they are not modifying adjacent nodes.

=====

* Lock free - Datastructures;

I) Concurrent Stack; (Slides)

- * • We use the CAS (Compare & Swap) operator to implement this.

* 1) Lock-free push:

- First we do $\text{node} \rightarrow \text{next} = \text{head}$;
- We use CAS to compare if $\text{node} \rightarrow \text{next}$ is equal to head, then we make $\text{head} = \text{node}$. Atomic
so we try
use this.

⇒ We perform the above 2 operations in a loop, so if another thread adds a node to stack after we do $\text{node} \rightarrow \text{next} = \text{head}$ then CAS fails, so we go and loop and try to make $\text{node} \rightarrow \text{next} = \text{head}$, CAS succeeds & we make $\text{head} = \text{node}$.

* 2) Lock-free pop:

- First we do $\text{current} = \text{head}$ } temporary pointer.
- Then in a loop we do,
 - ⇒ CAS(head & current) → break on success
 - ⇒ $\text{current} = \text{head}$. On failure.
- So if head & current point to same node, then CAS makes $\text{head} = \text{current}$ next .

⇒ If CAS fails, it loops and tries again.

* Note: CAS is atomic, (execution), but setting up its parameters is not atomic (since we need to put them in registers).

(This is called the ABA problem)

⇒ This can ~~end up~~ mess everything up we used CAS till now.



* Note: In the pop operation, ABA problem will make the operation fail if previous push operations have used recycled memory, so context switches b/w assignment of arguments & execution of CAS can mess up the stack.

(So what's the solution? → ~~Don't know!~~)



We use Loadlinked, Storelinked which would automatically update the registers if we are using recycled memory, so that no problem occurs with the pop operation.

instructions

~~Still this is not very clear.~~

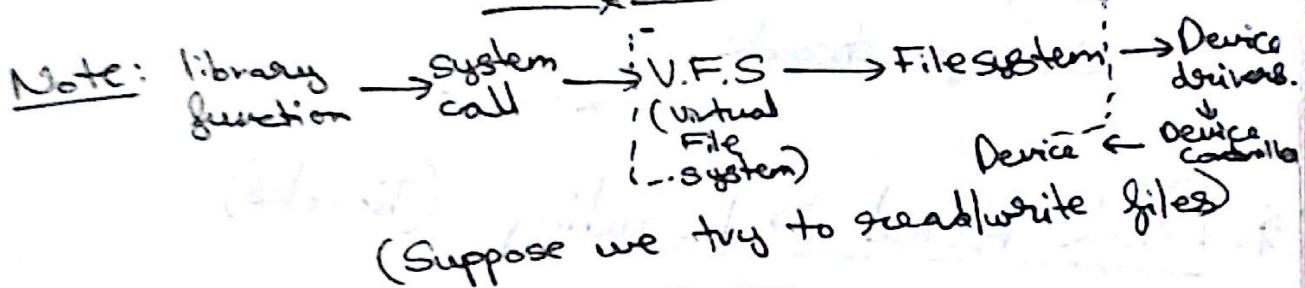


*Recap:

Note: Interrupts & DMA controllers speed up the CPU (ie no wastage of CPU time).

Note: Memory mapped I/O, just assigns parts of memory to I/O device registers, so load/store instructions work by its fast. In just I/O, we use IN, OUT instructions to copy data to/from device registers to our registers, this is also quite fast, but limited in no. of available registers.

- Device drivers must be given along with devices so they can be used.



*Disk:

- Platter → Track → Sector {Granularity
⇒ Read / write happens at sector level.

⇒ Parameters affecting read times;

1) Seek time: To bring read / write head onto correct track.

2) Rotational delay: To get to the correct sector (latency)

3) Transfer time.

- ⇒ The device drivers can do all sorts of tricks in assigning sector reads/write so as to minimize average read times.

* Block Size (No. of sectors in a block)

- As we increase block size, fragmentation may occur as blocks may not be full.

* Track Skew

- We keep a little offset of next sectors on a new track ($i+1$) once sectors of track (i) are full since we don't want to pass the sector before our seek head changes track.

* Disk Head Scheduling (Tracks)

- 1) Greedy (Nearest track first) (Unfair)
- 2) First-come - first serve. (Slow)
- 3) SCAN (Go to leftmost & then rightmost & repeat).

Secondary Storage (Non-volatile persistent storage)

- 1) Capacity
- 2) Consistency
- 3) Performance (R/W speeds)
- 4) Price

Working with a single device driver.

Redundant Array of

Inexpensive Disk; (RAID)

I) RAID 0:

- Multiple hard disks
- ⇒ Increased capacity & R/W can happen in parallel.
- (No fault tolerance)

II) RAID 1:

- 2 Disks store copies of the other 2 disks.
- ⇒ Good failure tolerance
- Bad capacity.

III) RAID 34:

- We use a parity disk.



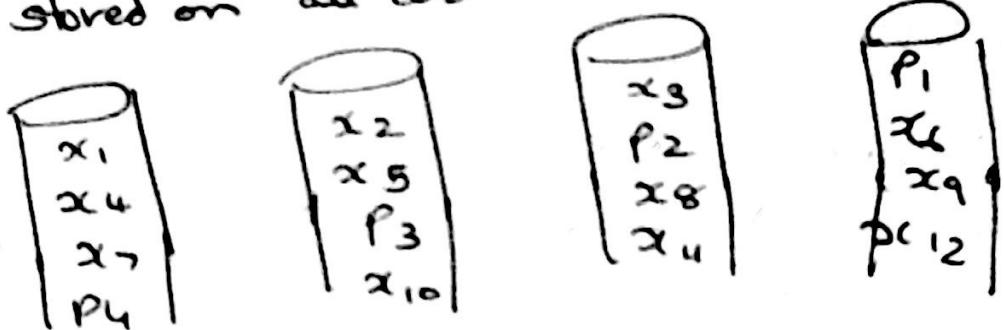
$\Rightarrow x_1 \oplus x_2 \oplus x_3 \oplus P = 0$, if this is false, then this means there is an error.

- Failure tolerance upto one disk

⇒ Bottleneck on parity disk.

III) RAID 5:

- Interleave the parity disk, so its stored on all disks

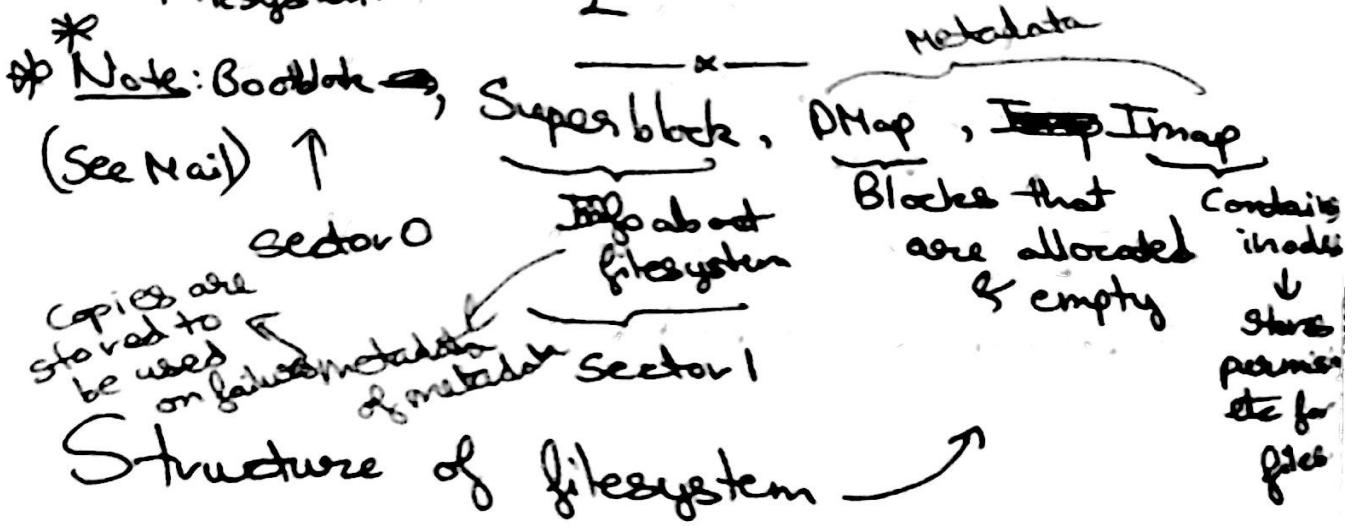
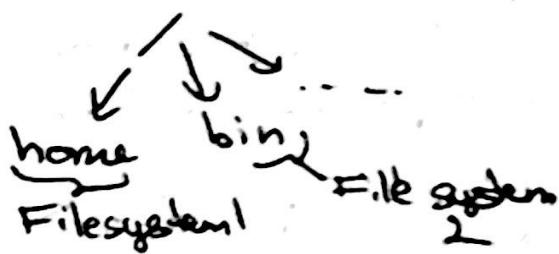


⇒ So the parity bottleneck is gone.

Note: Storage area networks

* Filesystems;

- Note: we can stitch multiple file systems together to make a directory structure.



* The memory allocated to a certain file is stored. These ^{inode index} structures will be small, so we use indirection (or double triple) so that we can store memory allocated to large files.



⇒ Hard links count is also stored in inode.

⇒ Different ~~files~~ ^{users}, can give hard links to a certain inode file, so that it can be shared by multiple users. (Essentially a copy)

* Note: Hard links cannot be made across file-system boundaries. straddle across.

— x —

Note: Symbolic links can be made across file system boundaries. (inode hard-link count does not change).

* To add memory to a file:

- 1) Update DMap & get a free block.
- 2) Update inode index structure by adding this block
- 3) Write data to the block.

* If a crash occurs after step ①, then we will have a block which will be wasted & never get unallocated. (space leak?)

⇒ We want to do all these three instructions atomically.

- fsck is a file system check which happens on boot of the system (generally time consuming).

*Transactions; (As in DBMS)

- Journaled file systems, which keep track of the steps performed.
- When a crash happens, we try to do partial transactions.

(A transaction consists of the 3 instructions we have seen before).

Note: Memory is freed only when all hard links referring to that memory are removed.