

L1: Introduction to Operating System.

→ OS is a program which manages computer system resources

What's an OS?

- middleware between user programs & hardware.
- manages hardware:

CPU, main memory,

IO devices (disk, NIC, mouse etc)

System
Hardware

- CPU
- RAM
- ROM
- various registers
- external devices via drivers.

* When we run a program:-

1. Compiler translates high level program into executable.
about .obj or .exe
2. the .exe contains instructions & data of the program.
good! All addresses!
3. CPU hardware is instructed via ISA.
4. CPU has registers:-
PC
instruction Code
Memory address

* When we run a program:-

- 1) read instruction at PC
 - 2) load data
 - 3) execute instruction
 - 4) store results.
- most recent data & instructions are cached at CPU, for faster results.

→ So what OS do?

- 1) manages program memory

- loads program executable from disk to memory - (code, data) ram!

- 2) manages CPU

- initializes PC & other registers.

- 3) OS manages external devices

- Rd/wr from files from disk.

I) OS manages the CPU:-

- * Operating System provides process abstraction.
- * OS creates & manages processes.
- * each process has the illusion of having the CPU to itself.
 OS virtualizes CPU.
- * Timeshares CPU among processes.
- * Enables co-ordination b/w processes.

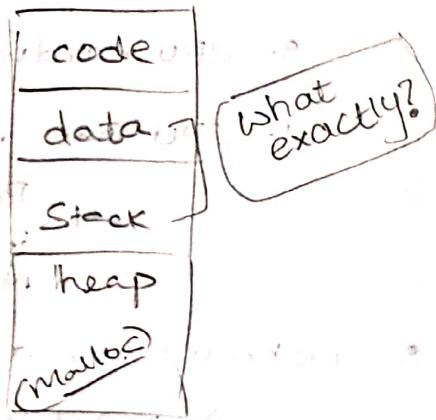
II) OS manages the memory:-

- * OS manages the memory of the processes:

code, data, Stack, heap

- * each process thinks it has separate space, numbers code and data starting from 0.

Virtual addressing.



- * OS abstracts the actual placement in memory. Translates from virtual addresses to actual addresses.

III) OS manages the devices:-

- * OS has code to manage disk, NIC, other devices — device drivers.

Eg: persistent data organized in disk.

- * OS evolved from running a single program, to multiple processes simultaneously.

b-21 The process abstraction

- OS provides a process abstraction.
- OS has a CPU scheduler that picks up one from many active processes to execute on a CPU:
 - Policy: which process to run
 - mechanism: how to "context switch" between processes.

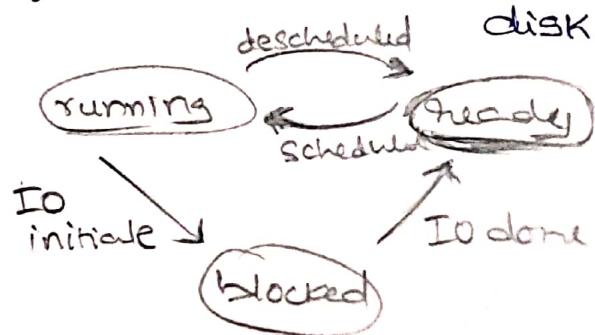
* A process constitutes:-

- a unique process identifier PID.
- memory image
 - code & data (static)
 - stack & heap (dynamic)
- CPU context (registers)
 - PC register
 - current operands
 - stack pointer.
- File descriptors
 - pointers to open files & devices.]

STDIN
STDOUT
STDERR

* State of a process:-

- Running
- Ready (Waiting to be scheduled)
- Blocked : suspended, not ready to run.
 - waiting for some event; like issued a read from disk.
- New : being created.
- Dead : terminated



* OS data structures. Maintains a data structure (list etc) of all active processes.

- Each process's info in a PCB (process control block).

- PID

- Process State

- Pointers related to other processes (parent process)

- CPU context of the process

[PC, Stack pointer, current operand]

- pointers to memory locations

(like the absolute position

- pointer to open files.

position of memory

cool. file B into B exec - image?)

* This datastructure won't contain the memory image of process **dummy!** that's in the memory itself.

This contains position of that.

L21 - xv6 introduction :- & x86 background.

O.S.

{
ISA family.

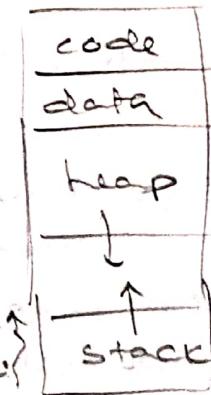
→ xv6 is a O.S. used for teaching

- has 2 versions ; for x86 hardware and one for RISC-V hardware
- we'll learn x86 version. (lets see x86 basics).

→ memory image of process:

* consists of compiled code:-

- compiled code
- Global / static variables
(memory for these is allocated at compile time)
- Heap (grows on demand)
- Stack (local vars. temp. storage during fun' calls etc.) grows 'up' towards lower addresses
- Others like shared libraries



• x86 registers & example instructions:

1) general purpose eax, ebx, ecx ..

2) eip

3) esp, ebp

 |
 stack
 ptr.
 |
 base
 ptr.

4) cr3 → metadata; pointer to page table.

5) Segment registers cs, ds, es, ..

mov %eax, %ebx

mov (%eax), %ebx

push %eax onto stack.

pop %eax
? send stack.top() into eax.] both of these
modify esp.

jmp %eax

* levels of privilege. (0 to 3)

 |
 OS
 code
 |
 Eg:
 setting CR3
 reg,
 io devices
 |
 User
 code
 |
 load value to eax.

* user should request
OS services (syscall)
for high privilege instruc^u

→ Function calls and stack

* What happens:-

1. push arguments to stack
2. "call" fn (this pushes current eip onto stack & jumps)
3. Allocate local vars & complete function.
4. "ret" (this pops the return address & jumps back)

* Register values get clobbered na!

OF CPU

1) caller saved registers:-

- saved on stack by caller before invoking the fn.
- callee code can freely change them.
- caller restores these registers after return.

2) callee saved:-

- caller expects these registers to have same value before & after function invocation.
- saved by callee function & restored as the function ends:
automatically done by C-compiler.

L22:- processes in xVG

PCB → process control block

- * the list of all PCBs is critical kernel datastructure & maintained in kernel memory.

struct proc in xVG

task_struct in Linux

- * in xVG; process states are

UNUSED

EMBRYO

SLEEPING ✓ (blocked)

RUNNABLE ✓

RUNNING ✓

ZOMBIE ✓ after killing a process

Formation - Formation

PCB has:-

- size of memory for proc
- pageable pointer.
- kernel stack pointer
- state of process
- PIB
- list of open files.
- process name (for debugging)
- curr dir.

1) Kernel Stack:- * for syscalls from process

- * when a function is called in process; user stack stores stuff.

- * but when process calls syscalls to run kernel code;

CPU context stored on kernel stack (security).

- this separate area for each process on kernel stack not accessible by users.
- the link to this KStack is through Struct proc of process.

2) List of open files:-

- * Array of pointers. When process opens a new file; a new element is created; & its index is passed as FD. (file descriptor).

- * First 3 elements of list are open by default:

stdin, stdout, stderr.

3) Page table:-

- * every entry in memory image has an address.
 - virtual address starting from 0.
 - actual physical address will be different.
- * Page table maintains mapping from Virtual address to physical address.
nice!
(more on this later)

→ Process table in xv6:-

```
struct {  
    struct Spinlock lock;  
    struct proc proc[NPROC];  
} ptable;
```

- * fixed size array of all procs.
(practically; might have dynamic size).
- * CPU Scheduler loops through ptable & sets a runnable process to running.

→ Process state transition:-

- i) A process that needs to sleep will set its state SLEEPING.
& invoke scheduler.
- ii) A process which has ran for its fair share will set its state to RUNNABLE & invoke scheduler.

L3:- PROCESS API

- * OS API is provided by a set of "system calls".
 - syscall is a function call into OS code which runs at higher privilege levels.
 - access to hardware, is allowed only at higher privilege levels.

* POSIX API:

in order to maintain code portability over various OS; all OSs have to be POSIX compliant.

* programming lang. libraries hide the details of these POSIX calls

Eg:- printf calls the write systemcall

→ process related System calls (in UNIX):-

- fork() creates a new child process
 - All processes are created by forking their parent except init process.
init is ancestor of all processes.
- exec() replaces the complete memory image of process
- exit() terminates a process.
- wait() causes a parent to block, until child terminates.

→ fork:-

- * A new process is created by making a copy of parent's memory image. This new process is added to ptable & set to runnable.
- * the return value from fork() = pid of child ; in parent process = 0 ; in child process.

Eg:-

```
int pid = fork();
if (pid==0){
    print "child";
} else {
    print "parent";
}
```

child prints this. our parent prints this.

→ wait for children to die...

- * process termination scenarios:-
 - 1) exit() syscall. (called automatically, when end of main)
 - 2) OS terminates misbehaving process.
- Terminated processes are state - ZOMBIE. Waiting to be reaped by their parent process.
 - * when parent calls wait(); its zombie are cleaned.
 - * if a parent terminates before child; the init adopts the child.

→ exec():

- After forking, parent and child are running the same code. meh!
(parent) Not too useful.
- * A process runs exec to load another executable into its memory.

→ How does a shell work?

- init process created just after boot up.
- init → forks shell/bash process
- & then this shell, after taking user command, forks a child
it uses exec() to run command executables.
the commands like wc
ls
are all executables.
- * if we wanna redirect output to a txt file;
 - spawn a child
 - close Std-OUT & open the file.
 - call exec().
this'll just update mem image.
won't change open files.
nice! stored in struct procna!

L23: System calls for xv6:-

Boopathi
L23

→ what happens on a system call?

- * System calls are available to user programs, defined in `userlib.h`.

header "user.h":

- * System call implementation, invokes a special ISA instruction

called "trap" instruction, called "int" in x86 ISA.

(in file usys.S)

- This 'int' instruction causes a jump to Kernel code that handles the system call.

System call number is stored in `eax` register.

1) fork() call:-

- new process, set to runnable, returns PID to parent, 0 to child.

achieved by
Setting
`$eax` to 0.
in child.

2) exec():-

- copy new executable into memory
- new Stack, heap.
- Switch process page table to use new memory img.

3) Exit():-

- i) exiting process cleans up state. (e.g. close the files)

* ii) Pass abandoned children to init.

orphans) → the non-zombie ones...?

iii) mark itself as zombie & invoke sched().

xv6 code in
Slides.
L23.

4) wait():

- search for dead (zombie) child in ptable & clean up.
- if no zombie, wait for one to die.
- if no children exist, return -1.

!!

usys.S

```
#define SYSCALL(name) \ * every trap instruction  
    .globl name; \ has  
    name: \ param n in "intr"  
        movl $SYS,%eax; \ param in $eax.  
        int $T_SYSCALL; \ n → unique to a device  
    ret \ T-syscall etc..  
SYSCALL(fork)
```

name variable.
into eax
trap instruction
\$eax → more specific call

SYSCALL(fork)

L4 - Mechanism of process execution:-

Low-level mechanisms (Bye bye 10K feet view)

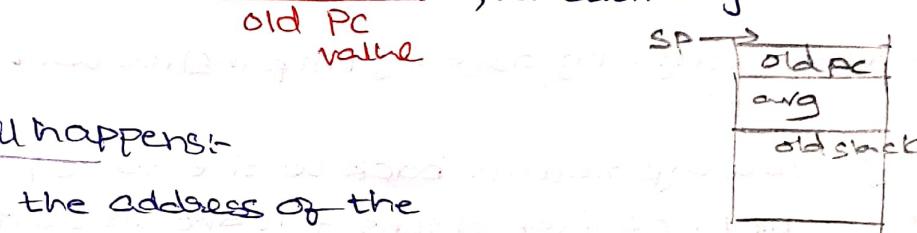
- how OS runs a process
- how OS handles a system call
- how OS context switches from one process to another.

1) OS handling a process:-

sequential instructions via PC have \$SP, \$BP

→ when function call happens:-

- 1) function call means jump instruction.
- 2) A new stack frame pushed onto stack & SP updated
- 3) old value of PC (return value) pushed to stack & PC updated.
- 4) Stack frame contains return value, function arguments etc.



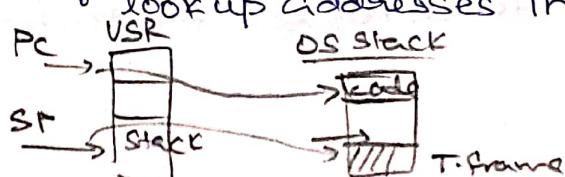
→ when system call happens:-

"we don't know the address of the sys-call instructions!"

- * Kernel doesn't trust user stack - uses a separate Kernel Stack in Kernel mode. (Separate for each process.)
- * Kernel doesn't trust user provided address to jump to
 - Kernel sets up Interrupt Descriptor Table (IDT) at boot time.
 - IDT has addresses of Kernel function instructions to go to, for system calls.
- * When syscall is made, a trap instruction is run. with syscall Baked into Silicon.

Trap instruction:-

- Moves CPU to higher privilege level.
- Switches to Kernel Stack
- Save trapframe (old PC, registers) on Kernel Stack.
- look up addresses in IDT & jump to trap handler function in OS code



→ Trap instruction:-

- trap instruction is executed on hardware in the case of

- System call (by user program)

Program fault (program access illegal memory) **segfault**.

* **Interrupt** (external device needs attention of OS)

or timeshare Eg: network packet arrived on NIC.
interrupt from O.S.

* IDT has many entries:

syscall/interrupt store an ID in CPU register before calling trap;

→ Returning from trap:-

Baked in Silicon

- * when OS is done with syscall/Interrupt; it calls an instruction called return-from-trap.
 - reverse everything done by trap instruction.
- * Must you always return back to the same process?
NO! · OS first checks if it should switch.

why Switch?

- Sometimes OS can't return back to same process it left
 - process has exited or terminated (seg fault)
 - 1. syscall 2. program fault.
 - process made a blocking system call.
- Sometimes we want to timeshare CPU.

In such cases OS performs **CONTEXT SWITCH**.

OS scheduler:-

Policy: which to pick.

Mechanism: how to pick

Policy:-

- Non-preemptive (co-operative) schedulers: switch, Only when process is running
- Preemptive (non cooperative) schedulers: CPU generates periodic timer
- after servicing the interrupt, OS checks if process has ran for too long.

Mechanism:-

Context switch happens; when both processes are in **Kernel mode**.

Now;

- Save context (PC, register, kernel stackpointer) of A on Kernel Stack.
 - Switch SP to Kernel stack of B.] stored in struct proc.
 - Restore context from B's kernel stack.
- [B's format would be similar; which was previously set up, by OS itself or even allocproc()].
- call return-from-trap instruction.

* Trapframe vs switch contexts

when going user code \rightarrow kernel code :-

trap frame stored by trap instruction.

When switching :-

context switching code stores context

then due RFI traps update the additional part because of the previous RFI instruction (which caused RFI segmentation fault).

problem goes with RFI, since marking cache

problem is how to handle multiple pages back to back

pages are merged with memory, and hence it will go to page frame after this new stack layout information will be

Segmentation

multiple memory and application will be mapped to one

memory and hence all memory get mapped to one segment of memory.

multiple memory and application will be mapped to one

L24:- Trap handling in xv6

* Traps:- when OS wants to switch; it traps user process.

- System calls
 - Interrupts | every device has IRQ number
 - ↳ interrupt request
 - program faults.
- * in usys.S, "int" instruction is invoked.] for syscalls.
- for hardware interrupts, device sends signal to CPU & CPU executes "int".
* Trap instruction has a param n) to indicate type of interrupt.
syscall & keyboard interrupt have different value for n.

- * The following happen as part of int(n):-
- eip & esp are pointing to user code & user stack previously
- 1) Fetch nth entry from IDT. (CPU knows location of IDT)
 - 2) Save esp into internal register.
 - 3) Switch esp to kernel stack of current process. (CPU knows this)
 - 4) On kernel stack, save old esp, eip
 - 5) Load new eip from IDT to kernel trap handler.

Also, syscall trap instruction can access syscall IDT but not disk-interrupt IDT. (made sure via CPU privilege levels)

* Trapframe

State pushed onto ^{kernel} stack during trap handling.

- CPU context of where execution stopped is saved. (so as to resume after trap)
- Some extra information needed by trap handler.

- * The intn instruct has only pushed the bottom few entries.
- The trap handler Kernel code will push the remaining.

<lec24> <pg5>

the C structure just indicates the complete structure of trapframe; bottom half of which is built by int instr.
E. top half built by all traps kernel trap handler.

→ Kernel trap handlers: (alltraps) assembly code.

- * IDT entries for all interrupts will set eip to point to the Kernel trap handler "alltraps".
 - * "Alltraps" assembly code will push remaining registers to complete trapframe on kernel stack.
 - "pushal" pushes all general purpose registers.
 - * invoke C trap handling function named "trap".
 - push pointer to trap frame (esp) as argument to trap().
- we have a mix of assembly code & C code.

→ C trap handler function:- trap(struct trapframe* t)

- * C trap handler; written in C; invoked in assembly.
this was passed via assembly code.
- * if `in int(n); n=="T_SYSCALL"` (as in usys.S); indicating this trap is system call.
- * Trap handler invokes common system call function
 - looks at call no. stored in eax & calls the corresponding func.
 - return value of syscall stored in eax.
(fork or exec or....)

hence in fork(); we make eax 0. in child.

`myproc()` returns the current struct proc.

check `myproc() -> if -> eax` & call the corresponding Syscall.

↳ store return value in same eax.

- * Separate code to handle interrupt from devices.
 (each device has different handler for INTN instruction)
 (each device has different interrupt priority level)
- * Timer is special hardware interrupt, & is generated periodically to trap to kernel.

On timer interrupt, a process yields CPU to scheduler.

```

if (myproc() == myproc() -> state == RUNNING &&
    tf -> trapno == T_IRQ0 + IRQ_TIMER)
    yield();
}

void
yield(void)
{
    acquire(&ptable.lock);
    myproc() -> state = RUNNABLE;
    sched();
    release(&ptable.lock);
}

```

→ Return from trap: trapret:

- pop all state from Kernel stack.
- instruction "iret" does the opposite of "int" instruction.
 - changes privilege levels
 - pops values pushed by int.
- execution of user code resumes.

L28: Context Switching in xv6:-

- * every CPU has a scheduler thread (special process, running the scheduler code)
- * After running for some time, a process switches to scheduler when
 - process terminated
 - process wants to sleep (blocking sys call)
 - process **yields** after running for a longtime.
(timer interrupt)
- * context switch happens only in kernel mode.



→ scheduler() and sched() :-

* Scheduler thread shifts to a process via Scheduler().

* user process shifts to scheduler thread via Sched().

Both have `switch(&context*, context*)`...
invoked from exit, sleep, yield.

→ Struct context:-

struct context {

 uint edi;
 uint esi;
 uint ebx;
 uint ebp;
 uint esp;

 set of registers
 to be saved;
 when switching
 processes.

* In both Scheduler() & Sched();
switch() switches between two "contexts".

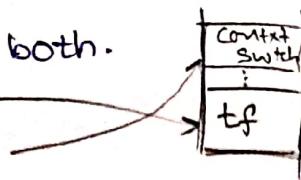
* context is pushed onto kernel stack.
Struct proc maintains a pointer to
the context structure on stack.
 $p \rightarrow \text{context}$.

→ Trapframe vs context:-

- trapframe saved, when CPU switches from user to kernel mode.
eip in trapframe is when syscall was made in usermode.
- context structure is saved, when CPU switches from process to scheduler.
eip in context structure is when switch() is called.
(In Kernel code)

→ Struct proc has pointers to both.

Struct trapframe *tf;
Struct context *context;

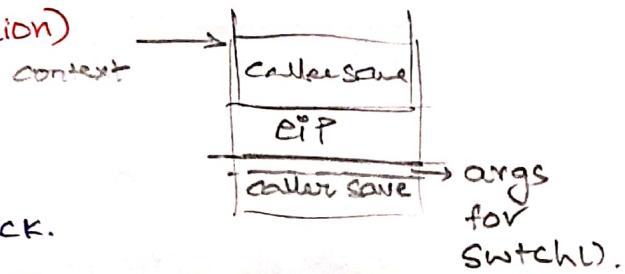


→ switch function:

- * This is the one, which actually creates the context struct.
- * Both CPU struct & procstruct maintain a context struct pointer
(struct context*)
- * Switch takes 2 arguments
 - since we need to update this pointer.
address of old context pointer, to switch from;
 - new context pointer to switch to; → we just need to read this context.
- * When invoked from scheduler();
 Switch(&p->scheduler), p->context);
When invoked from sched();
 Switch(&(p->context), mycpu()->scheduler);

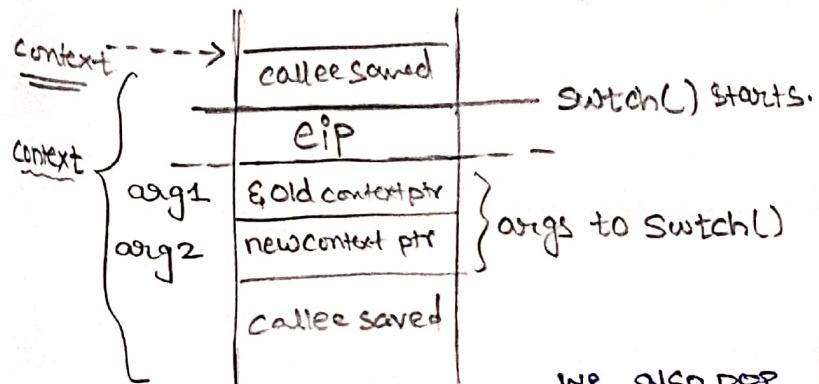
1) What is on Kstack, when a process has just invoked switch();

- caller save registers (C calling convention)
- Return address (eip)



2) What does switch do?

- push callee saved regs onto the Kstack.
- save pointer to this context in the struct proc->context.
- switch esp from old kernel stack to new Kernel stack.
- pop callee saved registers from new stack
- return from function call.



movl 4(%esp), %eax;
movl 8(%esp), %edx;

we also pop the callee regs
& return from switch()
in the new process.

Tough to write in C. ☺

L26: process creation in xv6:

* what's this?

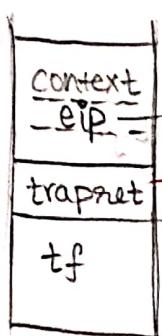
- * init process: first process created by xv6 after boot up.
 - This init process forks shell process; which in turn forks all other user processes.
 - init is ancestor of all other processes, in unix-like systems.
- * allocproc() is called during both init process creation & in fork system
 - Allocates new struct proc, PID etc.
 - [IMP] sets up kernel stack of this new process so that it is ready to be context-switched in by scheduler. COOL.

→ Allocproc:-

- 1) Find unused entry in ptable(). mark it as embryo.
Status = UNUSED. (mark as runnable, after creation completed)
- 2) New PID allocated
- 3) New memory for kstack allocated.
Kalloc() will return first byte address. I think fork() fills in.
- 4) Go to bottom of stack. Leave space for trapframe (move_on_this later)
- 5) push return address of "trapret".
no segment table changes for now
- 6) push context structure, with eip pointing to "forkret"

- When this new process is scheduled; it begins execution at forkret (in kernel code), then returns to trapret (in kernel code); then returns from trap frame to user code.

"we created a hand-crafted kstack, to look like the process was trapped & context-switched in post."



this is a piece of code
in alttraps.S

which takes us from
kernel to user code
after popping some &
calling iret instruction!

L26

Page 3

Awesome

Allocproc code.

intertwine assembly and c.

→ Init process creation:-

→ Allocproc() has created a new process.

Trapframe of process set to make process return to first instruction
of initcode.S in userspace.
But isn't initcode.S kernel?

* the initcode.S simply performs "exec" to run the init program.

* init program opens STDIN,STDOUT,STDERR files.

- inherited by all subsequent processes; as child inherits parent's files.

- forks a child, execs a shell.

- reaps dead children (its own or other orphan descendants).

<Read L26, Pg 435 code>

→ Forking a new process:- (more technical view)

• fork() allocates new process via allocproc() - parent number same.

2) parent memory & files copied:

3) Trap frame of child, copied from parent.

- Hence, child returns to the exact same line of code as parent.

- different physical mem. but same virtual memory address.

- only the value in eax (syscall return values are stored in eax)
is set to 0 in child.

4) set the child to runnable.

5) parent returns from fork() syscall & runs normally.

<L26, Pg 6> fork() code beautiful.

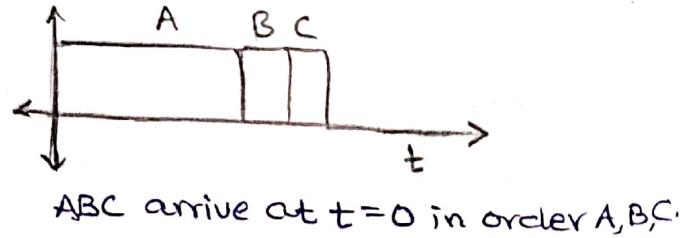
L5- Scheduling Policies:-

preemptive: willing to stop in middle
& large process.

- * Scheduler has 2 parts — policy now!
↳ mechanism
- which process to run next, from a set of processes.
- * OS scheduler schedules the CPU requests (bursts) of processes
 - CPU burst = CPU time used on a single stretch, by a process.
- * What are we trying to optimize?
 - Maximize Utilization (= fraction of time CPU is used)
easy. Always.
 - Minimize average turnaround time
(= time from process arrival to completion)
 - Minimize average response time
(= time from arrival to first execution)
 - Fairness: all processes must be treated equally
 - Minimize Overheads: run processes long enough to amortize cost of context switch ($\approx 1 \text{ ms}$)

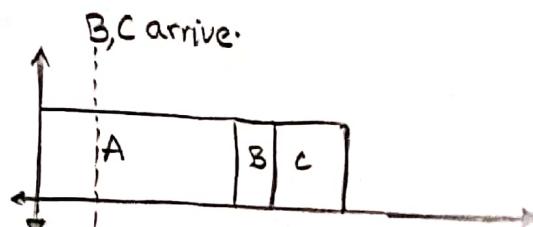
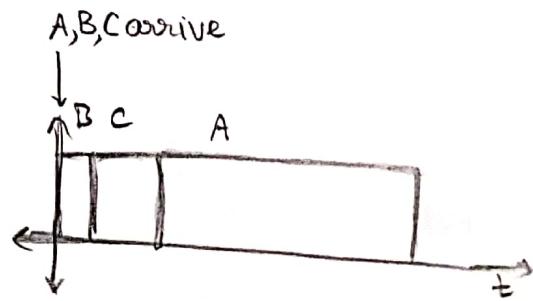
a) FIFO scheduling:-

- * Schedule in the order arrived.
- * Problem: convoy effect
 - A is too big, B,C must wait
 - High turnaround times.
Also response times!



b) Shortest Job first (SJF)

- * provable optimal avg. turnaround time, when all processes arrive together.
- * SJF is non-preemptive. So short jobs stuck behind longer ones; if they arrive after longer ones.

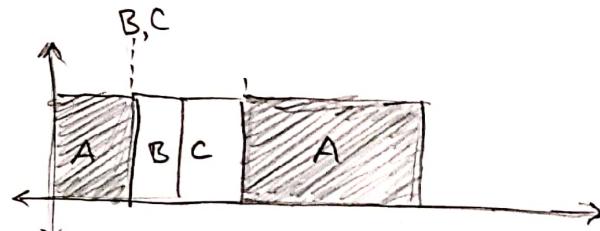


c) Shortest time to completion first :-

(STCF)

* ALSO, shortest remaining time first (SRTF)

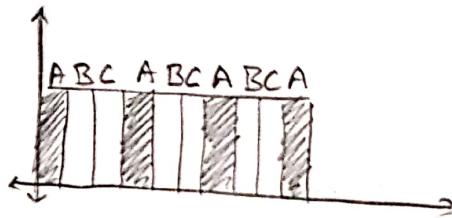
- preemptive scheduling.
- preempts running task, if time left is more than new arrivals.



Checking shortest when new procs arrive.

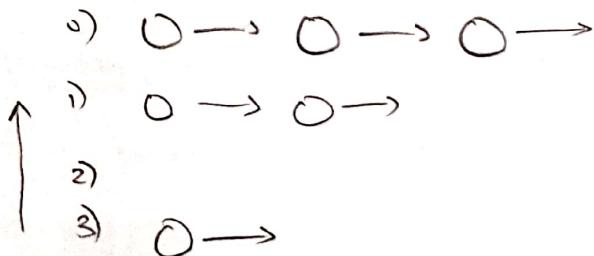
d) Round Robin:-

- * every process executed for a fixed quantum slice.
(keep slice big enough to amortize context switch cost)
- * preemptive.
- * Good for response time & fairness
- * Big blow on turnaround time.



e) Unix:-

- * Schedulers in real system are more complex.
- * Unix uses a multilevel feedback queue (MLFQ)
 - many queues, in priority order.
 - processes from highest queue scheduled first
 - within same priority, any algo like RR
 - priority of process decays with age.



L6: Inter-process Communication (IPC):-

- each process has its own unique memory image.
- If two processes want to work together, they need to use IPC mechanisms.

1. Shared memory:-

Shared memory get.

- * `shmget()` System call.

`int shmget (key_t Key, int size, int shmfld);`

- same key means, both processes have same segment of memory.

- Need to take care one process is not overwriting another.

2. Signals:-

- * Either the OS or a process, can send signal to another process.

`Ctrl+C` → SIGINT.

* Some signals can't be overridden. Eg: SIGKILL. Some can be.

* Signal handler: every process has default code for signal handling.

3. Sockets:-

- * can be used by processes on same machine or different machine to communicate.

- TCP/UDP sockets across PCs
- Unix sockets in local machine

* OS transfers data across socket buffers.

4. Pipes:-

- * Pipe system call returns 2 file descriptors.

- read handle & write handle

- A pipe is a half-duplex communication.



- * Regular pipes: both fds are in same process.

- parent & child fd after fork.

- * Named pipes: two endpoints can be in different processes.

- * OS buffers pipe data b/w read & write.

S message queues:

- * mailbox abstraction
- * process can open a mailbox & send/receive messages from it.
- * OS buffers the intermediate mails.
- * Each system call read/write have blocking & non-blocking versions.
 - read from empty queues
 - write to full queues.

read/write return value depends on whether the queue is full or empty.
returns with a value or an error code.

posting entries in the queue and removing them from the queue.

Termination -> 24.12.93

and now the mailboxes are implemented as shared memory areas in

processes using shared memory and memory protection facilities.

and now the shared memory areas are implemented as shared memory areas in processes using shared memory and memory protection facilities.

2012 version shared memory -

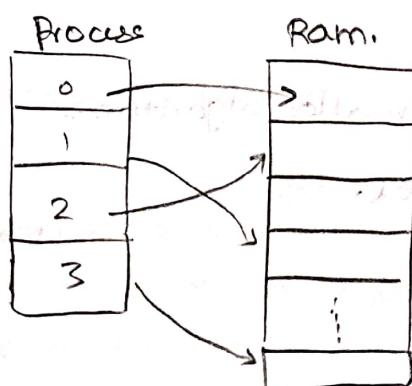
Byte-oriented shared memory -

L7: Intro to virtual memory:-

- * Why virtualize memory?
 - Bcoz real view of physical memory is messy!
 - Multiple processes are stored, all jumbled up.
 - Need to hide this complexity from user.
- * Every process thinks as if it has access to a large space of addresses from 0 to MAX.
- * CPU issues load store instructions with Virtual addresses.
- * MMU memory management unit translates virtual addresses to present on CPU itself!
 - OS makes the pagetable available to MMU.

→ The concept of paging:-

- * The virtual address space is divided into fixed size segments called pages. Also the physical addresses are divided into page frames.
- * To allocate memory to a process, pages are mapped to frames.
- * The pagetable stores mapping from virtual page number to physical frame number.



first 12 bits or so...
of address.

* Goals of memory virtualization:-

- Transparency: user not aware of messy details
- Efficiency : minimize overheads in terms of memory & access time.
- Isolation & protection : user should not access anything outside his address space.

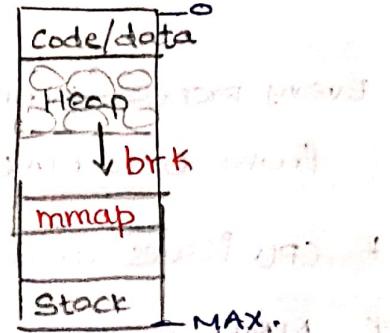
cool!

All 3 are handled.

whole 0 to MAX isn't allotted at process creation.
Waste then!

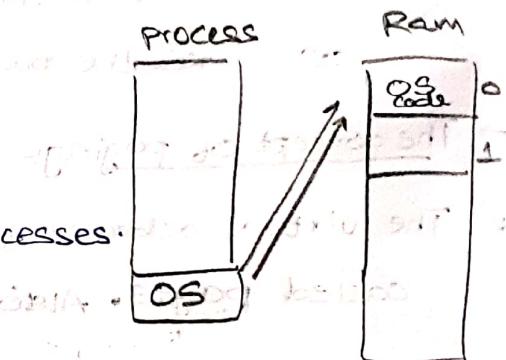
* Memory Allocation:

- * user C-code can allocate more heap memory using `malloc()`.
- * `malloc` implemented in C library-
 - here we get algos for efficient space management.
- * To grow heap,
 - `brk / sbrk system call.`
- * The user program can also allocate a page sized memory using the `mmap()` syscall.
 - get "anonymous" page from OS.



→ The address space of OS code in a process:-

- * OS is not a separate process with its own address space.
- * OS code is part of the address space of all processes.
- * Page table maps OS virtual addresses to real OS code.
 - (only 1 real OS code, like, for all processes).
 - like a header file.



* How does OS allot itself memory?

- large allocations; OS takes up a page.
- small allocations; OS uses various memalloc algorithms.

* Can't use `libc & malloc` in kernel. \textcircled{E}

↳ who will OS syscall to again....

So, in xv6 kernel code

(which is written in C)

We never have

`new ---;`

`malloc();`

hmmmm...
WOW....

L8: Mechanism of Address Translation:-

- * In a simple example of one page translation;
 - OS tells the MMU, the base (start address) and the bound (total size of allocation)
 - then the MMU does:
$$PA = VA + \text{base.}$$

if ($VA \geq \text{bound}$) error!
- * Hence, OS just says the base, bound once. It is not involved in every translation.
- * Hardware role:-
 - ISA has privileged instructions to set translation information.
(base, bound etc)
 - MMU uses this information for every! translation.
 - MMU generates fault & traps into OS when access is illegal!
 - ! Hardware interrupt
 - VA out of bounds.
- * Role of OS:- (in memory management)
 - * OS maintains free list of memory. (In a linked list..)
 - * Allocates frames during process creation (e.g. when requested)
 - * maintains page table in PCB of a process.
 - * Set address translation information in hardware.
 - * Handles traps from Hardware.
- Segmentation:- external fragmentation: - during Context switch ✓
 - during Trap? XX.
- * Say we are using generalized (base, bound) rather than fixed size pages.
- * Variable sized allocations leads to external fragmentation:
 - small holes left out in RAM, between segments.
 - no such issue with fixed size segments.

L9: Paging :-

- * lets allocate memory in fixed sized chunks. → makes memory management easy.
 - Avoids external fragmentation.
 - Has internal fragmentation (partially filled pages)

→ Page table:

- * per-process datastructure
- * Array stores mapping from Virtual page number to Physical frame number

process at boot of OS creates page table for each process.
Index is virtual page number!
not the starting address of page!

- * MMU has access to page table.
- * OS updates page table of MMU upon context switch.

→ Page table entry (PTE):

The OS class scene from

- * page table entry is one per virtual page.
- * VPN is the index into Page table, for its PTE.
- * Each page table entry contains
 - PFN (Physical frame number) & few other bits

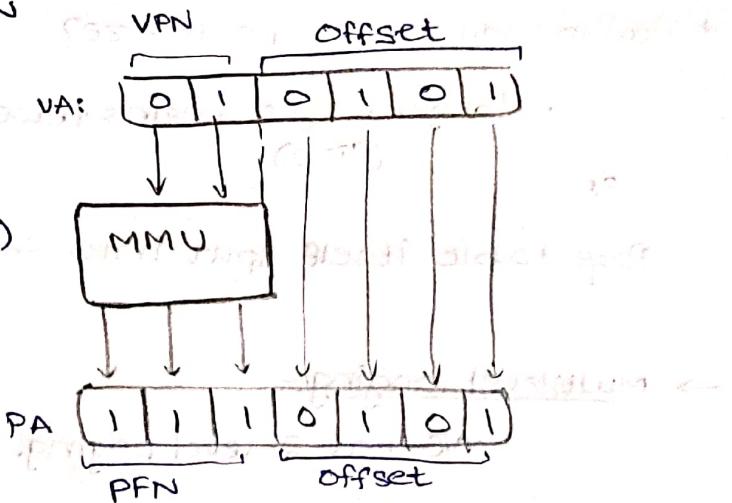
- "The social network".
- (Hacker News)
- valid bit: is the page used by process?
- Protection bits: read/wr. permission.
- Present bit: is this page in memory? Or swap? (even for reading!)
- Dirty bit: has this page been modified
- Accessed bit: has this page been recently accessed?

Status bits

* most significant bits of VA give VPN

* the MMU stores the address of the start of page table.

Need to add (or rather append!) offset.



* MMU needs to translate EVERY address from CPU.

- Paging adds overhead to translation. (could be multi-level page table).

* Hence use a cache for VA-PA translation.

→ Translation Lookaside Buffer TLB:- inside MMU

* A cache of recent VA-PA mappings accessed.

(Since Memory access is slow...)

* MMU first looks inside TLB.

- if TLB hit, PA directly accessed.

- if TLB miss, then MMU walks the page table.

* TLB misses are expensive (multiple memory accesses)

* Locality of reference has a high hit rate.

code usually

asks lws in similar locations.

- for loops
etc.

* TLB entries become invalid on context switch & change of pagetables.

* Page tables typical size:

4KB → first 12 bits of VA is our VPN.

12 bits
Offset

* say each PTE is 4 bytes; then total page table is

$$4 \times 2^{20} B = 4 MB$$

(too large!)
for single process.

* How to reduce page table size?

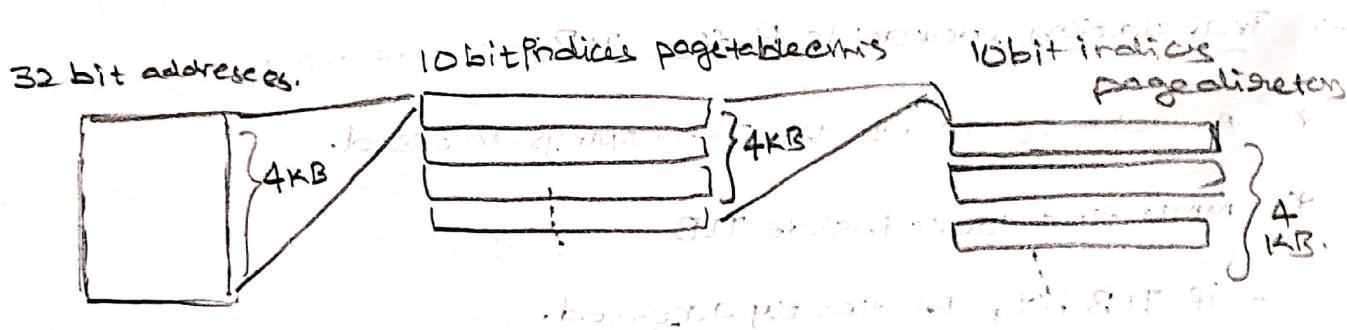
- larger pages means fewer PTE.
(Size)

or
Page table itself split into smaller chunks
(Pages) table.

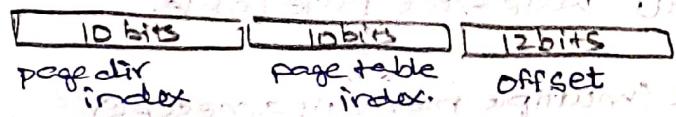
→ Multilevel paging:

X86 has 2-level paging.

- * A page table itself spread over many pages. (All need not be allocated)
→ multilevel PTE based on address (at process creation)
- * An "outer" pagetable or page directory tracks PFNs of pagetable pages.



VA:



$$\therefore 2^{32} \text{ values } \left\{ \begin{array}{l} 4 \times 2^{20} \text{ pagetable bytes} \\ \hline 4M\text{B} \end{array} \right\} \left\{ \begin{array}{l} 4 \times 2^{10} \text{ Pagedirectory bytes} \\ \hline 4K\text{B} \end{array} \right\}$$

All 2^{10} pagetable pages
need not be allocated.
This saves space
on pagetable

* we could need even higher level page tables.

- 64bit arch. - 7 levels.

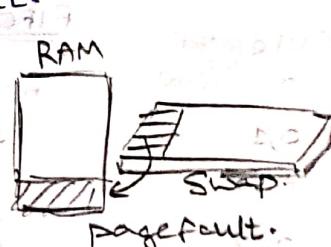
* Like this; in case of TLB miss, we'd need heavy cost.

L10- Demand Paging:-

- * Not all pages of all active processes are present in RAM/memory.
- OS uses a part of disk (swap area) to store pages not in use.

→ Page fault:-

- * present bit indicates if the page is in memory or not:
 - if page not present, MMU raises a trap to OS - page fault.
 - in Kernel mode (after trapping):
 - OS issues read to disk to bring back the page.
 - OS switches to another process. (Disk has a small CPU which works independent of our computer CPU).
 - After disk read completes, disk raises an interrupt & OS updates the page table of our old process & marks it ready.
 - When old process is scheduled again, OS restarts the instruction that caused page fault.



not the next instruction!

Summary:- What happens on memory access by MMU:-

1) CPU issues load to VA =

- checks CPU cache first Ooo....
- goes to main memory in case of cache miss.

2) MMU looks up TLB for VA =

- TLB hit : obtains PA, access memory.
- TLB miss: walks page table & obtains PTE.

When moving a page to swap;
clear TLB entry !!.

- If present bit set, access memory
- If not present, but valid, raise pagefault. Os handles page fault; sets the present bit; restarts instrutn.
- If invalid page, trap to OS for error.

* more complications in page fault:-

- what if OS finds no space to swap in the faulting page.
- OS will readily swap out pages to keep a list of free pages.

→ Page replacement policy:-

1. Optimal: replace page not needed for the longest time in future.
 (not practical...
 we don't know what pages we'll need in future).

2. FIFO: replace page brought in the earliest!
 (but could be a popular page)

3. LRU: replace the page that was least recently (or frequently) used.
LFU: replace the page that was least frequently used.

Eg: 3 frames & 4 pages...
 first few access are cold miss.

Optimal: compulsory → usually worse than optimal.

Access	Hit/Miss	Remove	Cache (after page fault)
0	miss		0
1	miss	cold misses.	0,1
2	miss		0,1,2
0	hit		0,1,2
1	hit		0,1,2
3	miss	2	0,1,3
0	hit		0,1,2
3	hit		0,1,2
1	hit		0,1,2
2	miss	3	0,1,2
1	hit		0,1,2

$$\text{net} = \frac{3+2}{2} \text{ misses}$$

cold

Access	Hit/Miss	Remove	Cache
0	miss		
1	miss		
2	miss		0,1,2
0	hit		0,1,2
1	hit		0,1,2
3	miss	0	1,2,3
0	miss	1	2,3,0
3	hit		2,3,0
1	miss	2	3,0,1
2	miss	3	0,1,2
1	hit		0,1,2

$$\text{net} = 3 + 4$$

misses

Buddy's anatomy:-

① Wasted space: performance gets worse, when memory size increases.

② Wasted time: page balancing

LRU policy:-

Access hit/miss evict cache.

0	miss	0
1	miss	0,1
2	miss	0,1,2
0	hit	1,2,0
1	hit	2,0,1
3	miss	2,0,1 0,1,3
0	hit	1,3,0
3	hit	1,0,3
1	hit	0,3,1
2	miss	0 3,1,2
1	hit	3,2,1

* works well due to

locality of reference.

{
- recently used pages
are more popular.

- Hence evict least
recently used
ones.

1. Net = 3 + 2
Same as optimal.

* How is LRU implemented:-

- OS is not involved in each & every memory access. How to implement LRU?
- Hardware helps, some approximation.
- MMU sets a bit "accessed bit" when page is accessed.
- OS periodically looks at this bit; to estimate active & inactive pages.
- To replace; OS tries to find a page that has low access bit
 - may also look for page with dirty bit unset.

* So OS only swaps out
pages of
current process?

classify into
hot/cold
pages

(to avoid swapping out
not to disk).
sure here...
(OS).

* will TLB be cleared, if page swapped out?

L11: memory allocation Algorithms

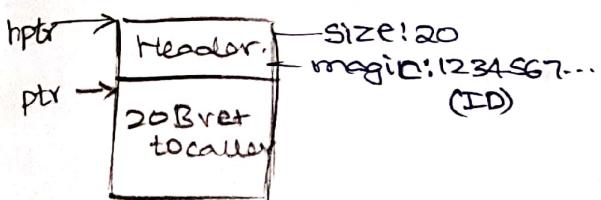
→ Linked list with headers
Data Structure + Algorithm.
↳ first, best, worst.

* Fixed size allocation is straightforward. Let's see for variable sized wast allocations.

* This problem must be solved in C-library in malloc().
↳ also in kernel. Kernel must allocate memory of its internal data structures.

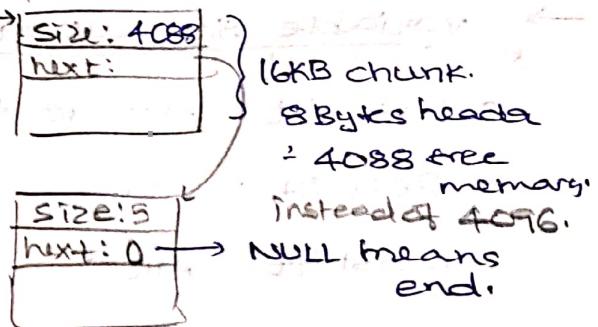
→ Headers:

- consider simple case of malloc()
 - every allocated chunk has a header with info like size of chunk.



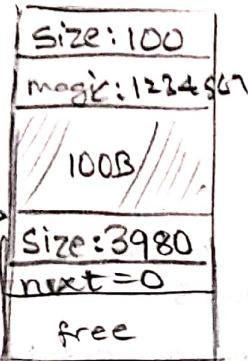
* Size is needed; since when free() called, we'll know how much to free.

- * Free space is managed as a freelist. head



allocating 100B

Initial



head →

final

see how
4088 → 100
occupied + 3980
= 8Bytes

header loss
= 8Bytes

Nice!

* external fragmentation:

- * Say we have 3 100 blocks & 1 3764 block. Now freed middle 100 block. Total we got 3864 B freespace
But fragmented!
- * So, we can't allocate a 3800B request!

Since continuous 3800 B unavailable.

* Splitting & Coalesce:-

* A smart algorithm must coalesce adjacent chunks.

* Must split while allotting requests &

coalesce while freeing requests.

coalescing

also

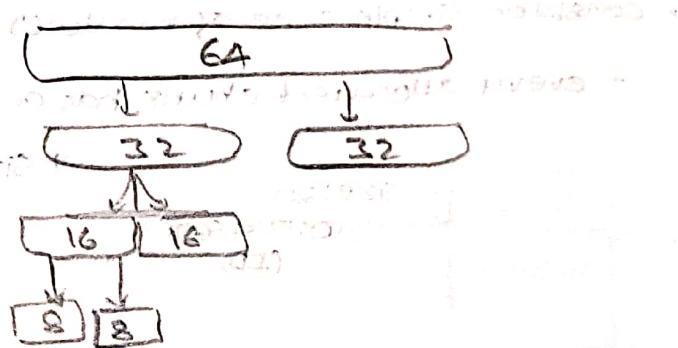
reduces

headers

overhead.

→ Buddy allocation for easy coalescing:-

• Allocate memory in powers of 2.



Eg: for 7KB request,

allocate 8KB

Cause: if this chunk & its buddy

this chunk & its buddy
are free; coalesce
into bigger
chunk.

→ variable size allocation strategies:-

* the policies!

• First fit: Allocate the first chunk that's sufficient.

• Best fit: Allocate the chunk closest in size.

• Worst fit: Allocate the largest chunk in size.

→ fixed size allocations:-

* A bit easier.

- has free list of pages

- pointer to next page stored in this page.

} need a header
na?

* for smaller allocations (like PCBs), Kernel

uses a slab allocator. I don't know if it's needed or not.

- object caches for each type (size) of object.

- within a cache, only fixed size allocations.

- each cache made up of one or more "slabs".

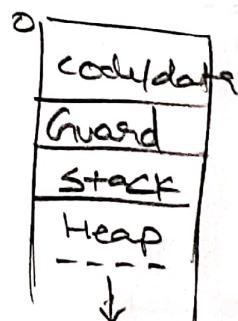
* we could use fixed size allocations in C (instead of malloc...).
But meh!

L27: Paging in xv6

- * 32 bit OS \rightarrow 4GB virtual address space per-process.
 - Page size \rightarrow 4KB
 - No demand paging.
- * Each PTE has :
 - Page table is indexed using 20bit index.
 - * 20 bit PFN (physical frame no.)
 - * Some flags:-
 - PTE-P: present. if not set \rightarrow page fault.
 - PTE-W: writable. if unset \rightarrow only read permitted
 - PTE-U: user : if unset \rightarrow only kernel can access page.
- * 2^{20} PTE can't be stored simultaneously.
 - Two-level page table.
 - 2^{10} inner page-table pages.
 - Outer page directory is 4KB in size.
 - physical address of outer page directory in CPU's cr3 register used by mmu

Process virtual address Space:-

- * From 0:
 - code/data
 - Fixed size Stack (with guardpage)
!xv6 only.
 - expandable heap
- * Kernel code/data:
 - from 2GB onwards
 - KERNBASE.
 - kernel code/data
 - * - free pages maintained by kernel
 - some space reserved for I/O devices.
- * Page table
 - USER PTEs } maps low virtual addresses
 - Kernel PTEs } maps high VAs to real OS code.
Identical in all processes



→ OS page table mappings:

- * maps $\text{VA} \xrightarrow{\text{KERNBASE}} (\text{2GB} \rightarrow \text{2GB+PHYSSTOP})$ to $(0, \text{PHYSSTOP})$
PA in physical memory.
- mapping identical in all processes for kernel code.
- * During trap, we'll use the same pagetable for OS code...
- * $[0, \text{PHYSSTOP})$ has code for
 - Kernel Code/Data
 - I/O devices
 - mostly free pages. In physical Addresses.
! can be mapped to user pages.
- ∴ in kernel VAs;
Phy. frame P has virtual address $P + 2\text{GB}$.
- ∴ same frame has 2 virtual addresses.
- * every RAM byte has 2 bytes in process.
 - ∴ Xv6 process can't use more than 2GB.

* freelist: maintained by OS.

- just a linked list.
 - Kernel maintains a head pointer
- Struct runq
struct runq *next;
};
1st.
- STRUCT {
struct spinlock lock;
int use_lock;
struct runq *freelist;
};
kmem;

→ Alloc & free:

- * who needs a new page: kalloc()
 - returns first free page on freelist
- * who needs to free a page: kfree()
 - Add free page to head of freelist.

Simple!

L-28: memory management of user process:

* New virtual address space for a process during

- init creation

- fork()

- exec()

.....

* existing VA space; modified in Sbrk Systemcall.

→ Building PageTable for a process:-

- Start with one page for directory.

- Allocate inner pages on-demand.

* begins with SetupKvm() (Outer directory allocated)

| Add kernel mappings

mappages() on Kmap[].

after kernel mappings, user page mappings added.

* page table entries added using

mappages(pgd, va, size, pa, int perm)

How many
Bytes.

- for each page,

walks page table; gets pointer to PTE using WalkPgd(...)

& fills it with pa, permissions.

mappages(pgd, va, size, pa, perm)

WalkPgd(pgd, va, alloc)

{

Searches

for

Pgtab in Pgdir

& returns

PTE in Pgtab

{ PDX, PTX are macros.

(page directory index page table index).

→ Fork & copy memory image [implemented in OS, not the user process. Hence, gotta work with PAs.]

- * copyuvm(), called by parent to copy, using with PAs.
 - Create a new Pgtable for child. { setupkvm(): not VAs}
 - Walk through parent VM page by page;
copy to child; add PTE in child.
- * for each page in parent:
 - fetch PTE, get physical address, permissions.
↳ memmove()
 - Kalloc new page for child. needs physical address.
memmove() into new page.
 - Add PTE from Va to Pa in child
not virtual
using mappage().
in both arguments!

* real OS do copy-on-write:

initially child also points to same PA as parent.

duplication occurs if one of them

modifies the page in their code.

→ Considered using EPT as filtering step, makes space efficient.

→ Addressing difficulties of multi-level pgtables

→ Solution: 2-level pgtable

→ 1st level: 1024 entries, 4MB pages

→ 2nd level: 1024 entries, 4KB pages

→ 3rd level: 1024 entries, 4KB pages

→ 4th level: 1024 entries, 4KB pages

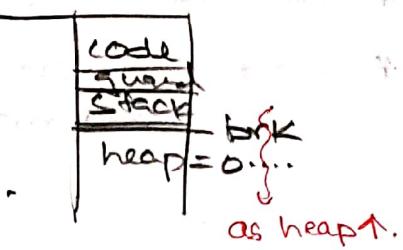
→ 5th level: 1024 entries, 4KB pages

→ Growing memory image: sbrk()

- * initially heap is empty.

program "break" is at end of stack.

- sbrk() is invoked by malloc() internally.



- * to grow memory; allocuvm() allocates new page, adds mappings into Pgdir.

- * whenever page-table updated, must update Cr3 reg & TLB
also done during context switch.

growproc(int n)

→ allocuvm(pgdir, sz, sz+n)

→ deallocuvm if n ≤ 0

→ switchuvm(cuproc); } refreshing

Allocuvm():

- * walk through new VAS to be added

* Alloc new page kalloc()

add to pagetable mappages().

* similarly deallocuvm() to kfree().

Allocuvm()

for (pa < newsz; pa += pg_size)

→ mem = kalloc()

→ mappages(pa, mem).
okay.

Afterwards, the new pages are mapped to the new VAS.

and the old ones are unmapped from the old VAS.

Memory is freed.

allocuvm() and freeuvm() are shared with allocvm() and deallocvm().
they also share the same refresh function.

allocuvm() and deallocuvm()

allocuvm() and deallocuvm()

allocuvm() and deallocuvm()

→ exec system call:

- executable & linkable format.
- * read ELF binary file from disk to memory.
 - * Start with new pgdir.
 - Add mappings to new executable pages & grow virtual addresses.

exec

+ setupKVM

+ readi(ELF)

+ allocUVM

- * After executable is copied (copy code + data hal.)

- Allocate 2 pages - Guard + Stack.

↳ permissions cleared.
Accessing this will trap!

- push exec() arguments onto user stack for main function of new program. (command line args).

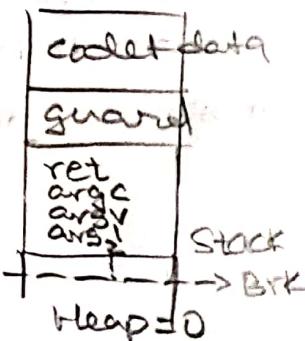
* Stack now has return address, args, argv array & the arguments themselves!!

exec

+ allocUVM (2PAGESIZE)

+ pushing arguments onto userstack.

new memory



- * if no errors so far; switch to new pgdir.

- if any error; don't switch. return exec() with error.

- * Set EIP in trap frame to entry of new program.

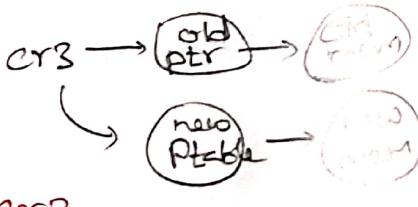
- returning from trap; new code will run.

exec

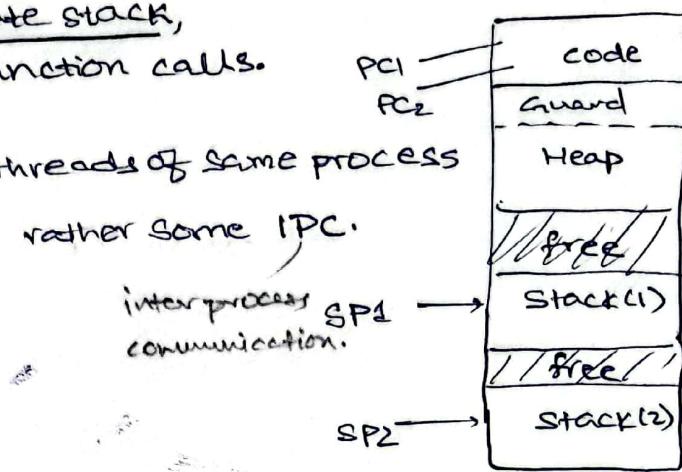
curproc->pgdir = pgdir

curproc->tf->eip = elf.entry

Switch UVM (curproc) } update CR3 & TLB



L12 : Threads and concurrency:

- * Thread is like another copy of a process, that executes independently.
 - Threads share the same address space (code, heap...)
 - Each thread has separate stack, for independent function calls.
- * For communication among threads of same process use Global variables rather some IPC.
 - interprocess SP1 communication.
 - SP2
- * Concurrency vs parallelism:-
 - multiple threads/processes at same time; even on single core; by interleaving their execution.
 - multiple processes/threads in parallel over multiple cores.
- + why threads?
 - single process can effectively use multiple cores parallelly.
 - Even if no parallelism, concurrency of threads ensures one thread runs, when other thread is blocked... imp!
- * Scheduling threads:-
 - The context of a thread (PC, register) is in Thread Control Block. TCB.
 - each PCB has a list of TCBs.
 - The threads those are scheduled independently by kernel are kernel-threads.
 - Linux pthreads.h is a kernel thread.
- * Some libraries provide user threads.
 - user sees many threads
 - library multiplexes large no. of user threads over small no. of K-thread
 - low overhead of switching in user threads.

→ Race Condition:- concurrent execution leads to different results.

- * We really can't predict the execution order of multiple threads.
at assembly level.

say

counter = 0

thread 1:

```
for(i=0; i<1000; i++)  
    counter++;
```

thread 2:

```
for(i=0; i<1000; i++)  
    counter++;
```

By end of both threads, we expect the counter to be 2000,
not anymore...

But we get less, sometimes
counter++; in assembly:-

- 0) mov 0x8abc, \$1
- 1) add \$1, 0x01

- 2) mov \$1, 0x8abc

Thread 1: executes 0,4 & then interrupt!

Thread 2: executes 0,4,8

Thread 1: executes 8.

Now, we finally got (41) instead of (42).

- * Critical Section : portion of code, that can lead to race condition.

- What shall we do? Mutual exclusion! only one thread executes critical section at once.

- What we need? Atomic Instructions & Atomicity of critical section from ISA

Achieve by

using Locks. 😊

LIB: locks:- pthreads library in linux provides such lock.

* We need to update a shared variable...

use a lock: only a thread which holds the lock can change it.

• Goals of building a lock

- Mutual exclusiveness: 1st...

- Fairness

- Low overhead: acquiring, releasing, waiting for a lock; Shouldn't consume many resources.

• locks needed for user programs & kernel programs.
(pthreads)

Implementation of locks needs support from microarchitecture & OS.
(atomic instructions)

→ Disabling interrupts:

* Having such an implementation isn't nice.

really?
maybe...
to send
signals...

good! - Disabling Interrupts is a privileged instruction & cannot give this power to user code.

- In multicore; another thread on another CPU can access critical section.
(unless u disable on all processors--> which is a big waste.)

* can use this in single processor systems inside OS.

→ Lock Implementation:-

* we'll use a flag to track whether lock is taken/available. But simple instructions won't do. Race condition moves to lock acquisition code.

Use atomic Instructions:-

1) test-and-set: update a variable & return old value. All in one instruction.

Ex So

• struct lock_t { int flag=0; }

• void lock (lock_t* lock) {

 while (TestAndSet(&lock->flag, 1) == 1) {

 ;

 } Spin wait. do nothing till condition true.

Spin Lock

Spinning until lock is acquired...

2) compareAndSwap (int* ptr, int expected, int new) {

```
    int actual = *ptr,  
        if (actual == expected)  
            *ptr = new;  
    return actual;  
}
```

Again, can implement a spinlock.

→ Alternative to spinning:- Sleeping mutex.

* A contending thread yields the CPU.

```
while ( testAndSet(&lock->flag, 1) == 1 )  
    yield(); // give up CPU.
```

Nice!

* Most userspace locks → Sleeping mutex kind. Saves lots of resources...

inside OS. → spinlocks only.

Who will OS yield to?

itself no?

Doubt!

* When OS acquires a spinlock!

nice.

It must disable interrupts, when the lock is held. Cuz interrupt handler could also request for the same lock & spin forever.

NO!

maybe scheduler itself needs that lock... again it itself yields

2. Must not perform any blocking operation.

2: recursive loop...

"never go to sleep with a lock". You are O.S!

who'll wake you up?

* When to use locks:

★ 1. A lock should be used before acquiring any shared data structures. "thread-safe data structures".

2. All Shared Kernel DS. should be accessed after locking.

3. Coarse-grained locks vs fine-grained locks:-
one big lock for all shared data

Separate locks...
- allows more parallelism.
- harder to manage.

L14: Condition Variables:- Pthreads provide CV for user prog.

* we done mutual-exclusion. Lets see Waiting-Signalling.

→ Condition Variables:

* A CV is a queue, a thread/process can put itself into; waiting for condition.

- Another thread signals CV to wake up a waiting thread.

Signal → wakeup one thread.

Signal broadcast → wakeup all.

- no flag inside condition variable. just wait(), signal() methods.

- wait() is different from yield() & yield() makes status runnable; where wait() makes status sleeping.

Eg:

Shared value int done = 0

LOCK mutex_t m = init_mutex;

CV: condvar_t c = init_condvar;

```
void thr_exit() {
    lock(&m);
    done = 1;
    Signal(&c);
    unlock(&m);
}
```

```
void* child {
    thr_exit();
    return NULL;
}
```

```
void thr_join() {
    lock(&m);
    while(done == 0)
        wait(&c, &m);
    unlock(&m);
}
```

```
main {
    pthread_t p;
}
```

```
pthread_create(&p, child);
thr_join();
return 0;
}
```

use a 'while' loop; rather than if-condition;
to account for wrong waking of CV.

use lock, while accessing a shared
variable. Otherwise here; race condition
could send parent to permanent sleep.

"Lock must be held while using
wait, signal on conditional variables!"

and so;

The wait(.,.) releases the lock; before
putting thread to sleep. Never skip with lock
Similarly; when thread wakes up; it would
already be holding the lock.

1. CV is a queue

2. Update shared vars inside a
(Read) lock

3. Use lock; before using
wait, on conditional variables."

→ Producer - Consumer Problem:-

* A common pattern in multi-threaded programs:

Setup: One or more producer threads
one or more consumer threads
a fixed buffer for everyone.

* Using 2 CVs:-

```
int count = 0;  
cond_t empty, full;  
mutex_t mutex;  
  
void* producer(void* args){  
    int i;  
    for(i=0; i<loops; i++){  
        lock(&mutex);  
        while(count == MAX)  
            wait(&empty, &mutex);  
        put(i);  
        signal(&full);  
        unlock(&mutex);  
    }  
}
```

no need lock here.
locks already stored in
queue of cv. Those need to
be returned to processes.

```
void* consumer(void){  
    int i;  
    for(i=0; i<loops; i++){  
        lock(&mutex);  
        while(count == 0)  
            wait(&full, &mutex);  
        get(i);  
        signal(&empty);  
        unlock(&mutex);  
    }  
}
```

Again we note:-

- 1) Lock must be held,
during wait, signal on
conditional variables.
- 2) Pass the lock to wait(),
which releases the lock.
- 3) Never sleep with a lock.

L15. Semaphores:-

- * Semaphore is also a synchronized primitive like CVS.
 - Semaphore has an underlying counter.
 - up/post increments the counter.
 - down/wait decrements the counter & blocks the calling thread if the resulting value is negative.

- * Semaphore with value 1, acts like a lock.

Binary Semaphore = mutex.

```
sem_t m;  
sem_init(&m, 1);  
    X = 1  
    for lock.  
wait(&m);  
    //critical section  
post(&m)
```

- * Semaphores for ordering,

```
void* child() {  
    sempost(&s);  
}  
{  
main {  
    thread(&p, child);  
    sem_wait(&s);  
}
```

See how we don't use lock beforehand...
Since if() condition is atomically implemented inside semu.

Here no lock is passed in wait(); unlike CV.wait.

→ Producer Consumer problem again!

- * 2 semaphores for signalling:

- one to track empty slots
- one to track full slots

- * 1 Semaphore as mutex for shared buffer.

```

sem_t empty; sem_init(&empty, 0, MAX);
sem_t full; sem_init(&full, 0, 0);           ↘ some other arg... ← loaded
sem_t mutex; sem_init(&mutex, 0, 1); } a lock...
void* producer() {
    int i;
    for (i=0; i<loops, i++) {
        sem_wait(&empty); ] mind this order!
        sem_wait(&mutex);
        put(i);
        sem_post(&mutex);
        sem_post(&full);
    }
}
void* consumer() {
    int i;
    for (i=0; i<loops, i++) {
        sem_wait(&full); ←
        sem_wait(&mutex);
        get(i);
        sem_post(&mutex);
        sem_post(&empty);
    }
}

```

never sleep with a lock.
"Acquire lock after signalling".

A different conclusion; when CV are used...

Since cv-wait are able to free locks,
whereas sem-wait are not.

no if() - else statement.
here; since inbuilt inside semaphore.

L16:

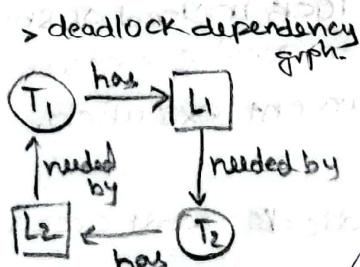
Bugs in concurrent programming:-

Deadlock bugs:-
cannot go further

classic: T₁ has L₁, needs L₂
T₂ has L₂, needs L₁

Condition for deadlock:

See next page.



when our assumption about atomicity of instructions is violated.

when our assumption abt order of exec. of 2 or more threads is wrong.

can't just assume such stuff.

Non deadlock bugs:-
wrong results obtained

Atomicity Bugs:

- fix is: use locks for mutual excl.

Order violation Bugs:

- use CV to impose ordering among various processes.
- Signal in 1st process wait in 2nd process & a bool variable.

→ conditions for deadlock to occur

1. Mutual exclusion: a thread claims exclusive control over a resource.
2. Hold-and-wait: thread holds a resource & is waiting for another.
3. No preemption: Thread cannot be made to give up its resources.
4. Circular wait: There is cycle in resource dependency graph.

* All four above must hold for deadlock.

* preventing circular wait :-

- Acquire locks in a fixed order, everywhere.
- Impose total ordering.

Eg: Acquire the least address lock ...

```
if ( $m_1 < m_2$ ) {  
    lock (& $m_1$ );  
    lock (& $m_2$ );  
} else {  
    lock (& $m_2$ );  
    lock (& $m_1$ );  
}
```

reboot system or

kill deadlocked process

Doubt

will killing work?

lock is still held...

* preventing hold-and-wait:

- Have a big master lock; instead of small ones.
 then acquire smaller ones.
- may reduce concurrent execution.

* OS can schedule so deadlocks won't occur... impractical!

Banker's Algorithm

OS won't know nat.

L29: Locking in xv6 :-

* xv6 has no threads. But supports multiple cores.

So; we need locking in OS code.

* Use spinlocks to protect critical sections.

→ Spinlocks in xv6 :-

* xchg (ELK→locked,1) x86 atomic instruction

Similar to testAndSet.

* must disable interrupts; before spinning for lock.

!

* Interrupts stay disabled; until ALL locks are released.

Maintain a counter for no. of locks held.

mycpu() → ncli;

* pushcli()

```
{  
    disable interrupts  
    if mycpu()→ncli==0;  
    mycpu()→ncli++;  
}
```

* popcli()

```
{  
    decrement ncli;  
    if ncli==0: enable interrupts  
    sti();  
}
```

acquire () {

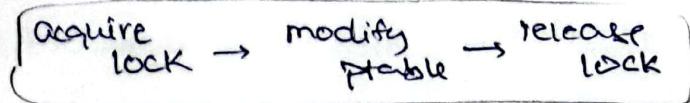
```
    pushcli(); → disable  
    while(xchg(ELK→locked,1)!=0)  
    ;  
}
```

→ ptable.lock:-

* Any access to ptable must be done with ptable.lock held.

But in code, sometimes...

* Normally a process:-



But during context switch :- ptable is changing throughout.

(P₁ holds lock → goes to scheduler → P₂ releases lock.)

* Every function that calls sched(); does so, while holding the lock.

yield()

sleepc()

exit()

yield(), sleepc(), forkret()

For a newly created process.

Every function switch(); returns to, releases the ptable lock immediately.

* Subtle points about Scheduler();

- scheduler runs a loop with lock held.
- Periodically; when a loop through all the processes is over;
- ptable.lock is released & interrupts are enabled.
- what if no process is runnable & interrupts are all disabled....

To avoid this; we enable interrupts periodically.

- * Interrupts during lockheld time are queued up.(not discarded).

```
void
scheduler(void)
{
    for(;;){ // loop
        sti(); // turn on interrupt
        acquire(ptable.lock);
        for(over all processes){
            switch(0,0);
        }
        release(ptable.lock);
    }
}
```

Kernel functions... not userspace...

L-30: Sleep and Wakeup in xv6 :-

How interrupts work... Similar to condition variables
in user codes.

→ Basis:-

- * process P₁ blocks for an event - disk read.

Invokes sleep() function.

- * while P₂ is running; say disk interrupt occurs and wakeup() is run.

How to know which proc? Use channel:

Stored in
Struct Proc 'chan'

- * ALSO; have to use lock for sleeping & wakeup...
 same common value known to both sleep process & wakeup process

- to bypass missed wakeup : location of struct proc
problem. : address of disk read etc...

So: wakeup sent, even before proc. went to sleep...

P₁: lock l; if (!done) sleep(c, l)

P₂: lock l; wakeup(c); unlock l;

release l
unlock l;

① void sleep(void* chan, spinlock *lk){
 p = curproc();
 if (lk != &ptable.lock) {
 acquire(&ptable.lock);
 release(&lk); for wakeup fn!.

P → chan = chan; // stored in struct proc

schel();

P → chan = 0;
if (lk != &ptable.lock) {
 release(&ptable.lock);
 acquire(lk);

} release the lk given; only if its NOT ptble.lock.

(called when lk is ptble. lock)

② void wakeup1(void* chan){
 struct proc* p;
 for (p in proc[]) & p->state = SLEEPING
 P → chan == chan
 P → state == RUNNABLE;

③ if lk is not ptble lock.

void wakeup(void* chan){
 acquire(ptable.lock);
 wakeup1(chan);
 release(ptable.lock);

A wakeup cannot run in below sleep; as sleep holds lk or ptble.lock.

* Examples are pipes:-

Struct Pipe{

```
pipewrite ( pipe* P; char* addr;
            int n);
```

```
piperead ( pipe* P; char* addr;
            int n);
```

```
struct spinlock lock;
char data[PIPESIZE];
uint nread;
uint nwrite;
int readopen;
int writeopen;
```

- one process writes into pipe; another reads from the pipe.
- reader sleeps, if the nread == nwrite; & writer wakes it.
- writer sleeps if nwrite = nread + PIPESIZE; & reader wakes it.
- Channel is addresses of member variables. fine.

* Example is wait and exit! - Here the lock held by both is ptable.lock

- * Parent calls wait when child is running

```
wait() // wait for children
sleep(cuproc, &ptable.lock);
```

In child; exit() has lock & wakes up parent

```
exit() // parent could be sleeping
wakeup1 (cuproc->parent)
minimally directly this called 😊
```

Yay! lets see!

addr. of
chan = struct
ptable of
parent
process.

- * child proc. by itself can't cleanup its complete memory... & its struct

CR3, Kstack
are in use always

proc.

So; wait() by parent needs to do this.

Lecture 17: Communication with IO devices:-

- * IO devices connect to CPU and memory via bus.

Block devices:-

- * Store a set of numbered blocks.
- Eg: Disks

Character devices:

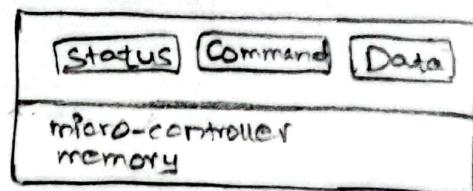
- * produce/consume stream of bytes
- Eg: Keyboard

- * Expose an interface of registers :- status, command, data

- * How OS reads/writes these device regs:-

1) Explicit IO instructions

- In, Out are two such, on x86.
- privileged instructions for OS only.



2) Memory mapped IO:

- Device regs. appear like memory locations.
- Memory hardware routes links these special memory addresses to device regs.

Execution of IO requests:-

```

while (status == BUSY)
;
[write data to data reg
 write command to command reg
 while (status == BUSY)
;
]

```

2) Programmed I/O: explicit copying of data

Direct memory Access (DMA):

- CPU wastes time in copying data.
- Instead, a special hardware DMA engine

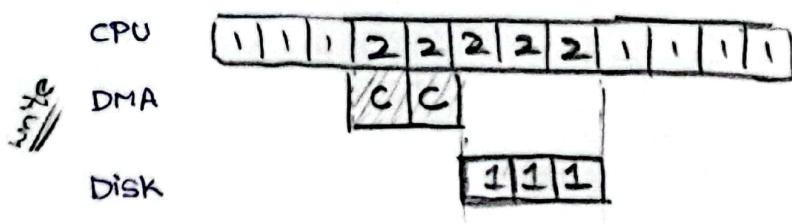
copies from main memory to device and back.
CPU provides memory location to DMA.

read → DMA raises interrupt after copying to mem.

write → disk starts writing after DMA copies to register.

1) Polling status to see if device's ready:-

- wastes CPU cycles
- use interrupts. Send process to sleep. Eg. device at completion raises interrupt.
- IRQ number identifies which interrupt handler function to call to...
- wakes up blocked Process & starts next pending IO request.
- If device is fast, polling better than interrupts as we avoid kernel mode transition overheads. Monitor?

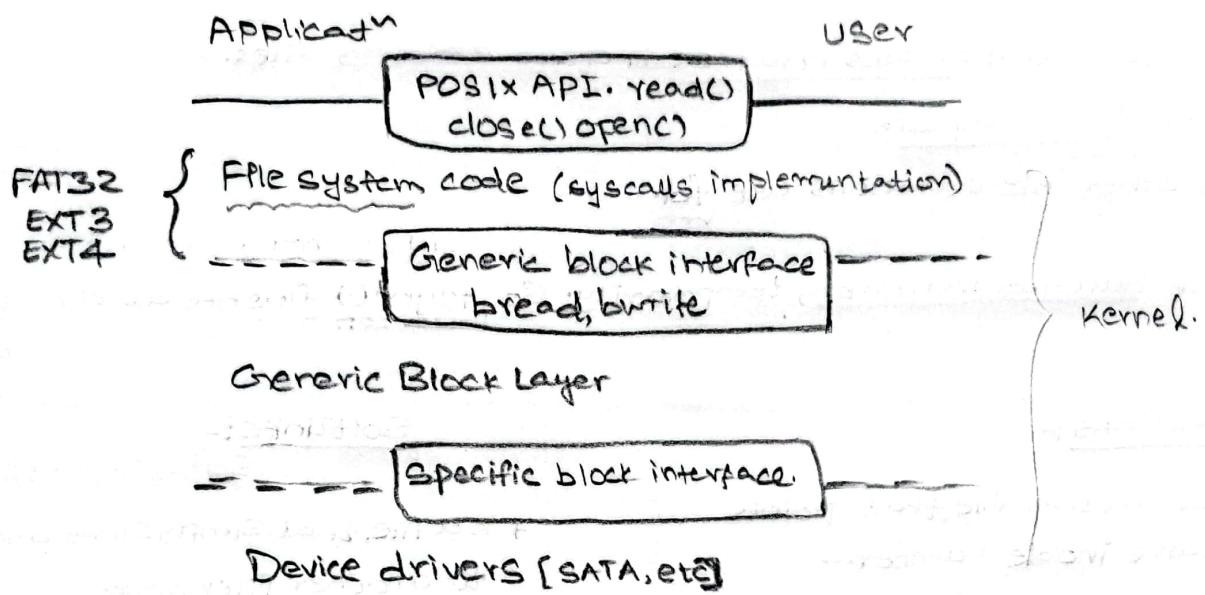


→ Device drivers: Software.

* part of OS code, that talks to device, handles its interrupts etc.

* Most OS code abstracts the device details.

Eg: file system code is written on top of a generic block interface.



L18: Files and Directories:

inode → index node.

- * File: stored persistently.

identified with filename (human readable)

& inode number (OS-level)

for each & every file? Wow!

Directory: A kind of file, whose contents are filename - inode mappings which it contains.

- * open() system call creates new files & opens existing files.

Returns file descriptor.

All other file operations use fd.

- * Files are buffered in memory temporarily. So fsync() flushes all changes to disk.

Hardlinks :-

* creates another file, that points to same inode number...

* when one deleted; we access inode through another.

* inode maintains a link count.

Delete when Links = 0.

Softlinks :-

Symbolic links

* is a file, that simply stores pointer to another filename.

* if main file deleted, then inode gone.

- * mounting a filesystem & devices using mount command in linux.

→ Memory mapping a file:

* Alternative way of file access without fd, read(), write().

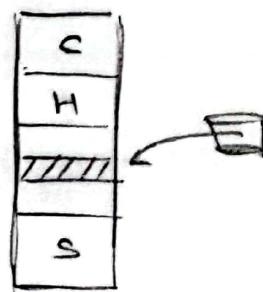
* mmap() allocates a page in virtual address space.

- "Anonymous" page for program data

- file-backed page contains file data.

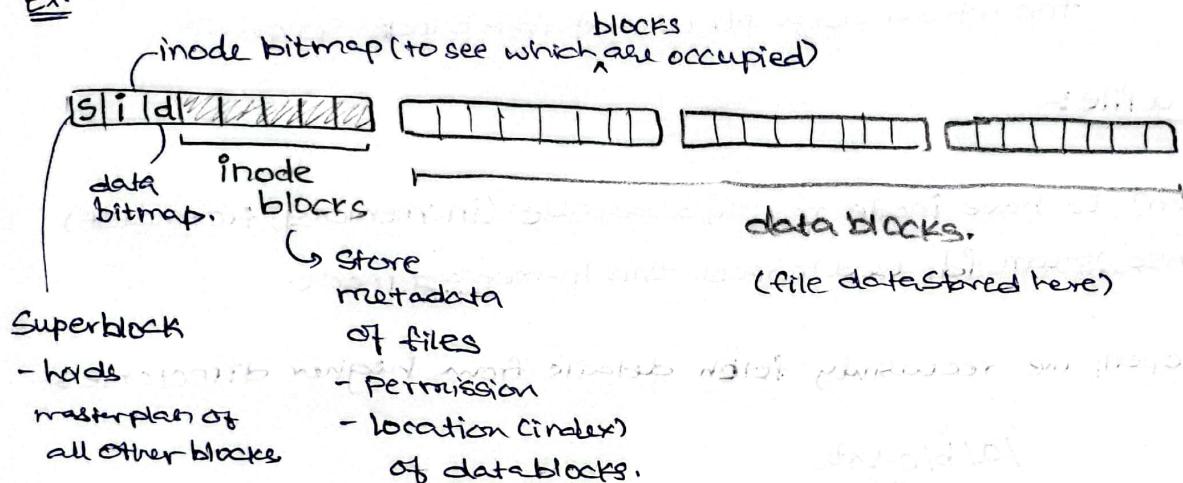
(filename argument to mmap())

* Access file data like any other memory location.



- L19: File System Implementation: File system ≠ device driver.
 I think...
 It's another layer of OS code.
- * Filesystem is an organization of files and directories on disk.
 - * Two main aspects
 - Data structures to organize
 - System call implementations to read, delete etc.
 - * Disks expose a set of blocks. say 512B. Filesystem organize files onto these blocks.

Ex:



→ Inode table:

- * Inode no. is index into this table. each inode block is:
- * Inode stores:
 - 1) file metadata - name, access time.
 - 2) pointers to file data.
(block numbers)

data bit map	0	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16	17	18	19	20	21	22	23	24	25	26	27	28	29	30	31	
	1	0	1	0	1	0	1	0	1	0	1	0	1	0	1	0	1	0	1	0	1	0	1	0	1	0	1	0	1	0	1	0	1

iblock 0 iblock 1

array of inodes.

- * file not stored contiguously on disk. Have to track block numbers.

So, inside inode...

- Direct pointers: no. of first few blocks stored in inode itself.
(enough for small files)
- Indirect pointers:
inode stores no. of block, which in turn has block numbers of file.

* we similarly have double and triple indirect blocks.

multilevel indexing.

- * Or, use File allocation table (FAT); where each block has pointer to next block.
First block add. stored in inode.

file system

a separate
filesystem
na...

* tracking free blocks: (both inode blocks and data blocks)

- Bitmaps

Store one bit per block to indicate free/used.

- Freelist

Superblock stores pointer to 1st free block.

this intern stores ptr to next free block, so on...

→ opening a file:-

- * why open? to have inode readily available (in memory, from disk)
 - Also return fd; used to reach this in-memory inode.

- * during open; we recursively fetch details from higher directories

↓↓↓
/a/b/c.txt
create
or
load inode

* open file table:

Global:-

across all processes.
* one entry for every open file.

(even pipes, sockets)

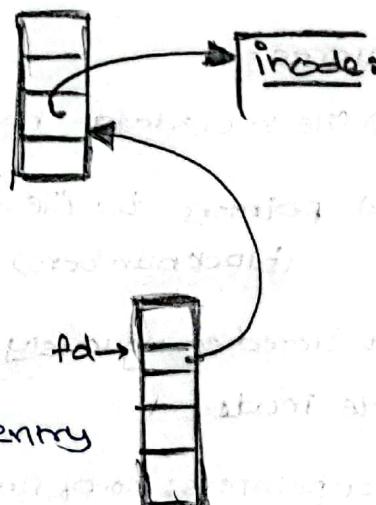
- Entry points to in-memory copy of inode.

Per-process:-

* Array of files opened by process.

- fd is index into this array

- entry points to global-open-file-table entry



- * open() creates entries in both filetables & returns fd.

→ Virtual File System:- (abstraction)

- * File systems differ in implementation of data structures.
- * VFS looks at a file system as objects & operations.
- * Syscall logic is written on VFS.
- * To develop new file system; simply implement functions on VFS objects & provide these to kernel.

* Syscall implementation does not have to change with FS implementation.
Spirit of SW Dev.

→ Disk buffer cache:-

- * recently fetched disk blocks are cached.
 - FS issued read/write are passed onto buffer first.
 - while writing;
 - Synchronous/write-through cache : write to disk imm.
 - Asynchronous/write-back cache : have a dirty bit set & write back when evicted.

* Benefits:

- Improved performance due to no disk I/O
- single copy of block in memory (no inconsistency)
- * Some applications like databases, avoid caching altogether.
to avoid inconsistencies
due to crash.

L31,32: XVG filesystem

Abstractions:

System Call

open()

link()

read()

operation on FS-data Structs

struct dirent{} → disk structure.

struct dirent{} → disk structure.

struct file{} →

struct ftable{} for open files.

struct inode{} → in-memory inode.

struct icache{} →

Block I/O layer

(disk buffer cache)

* buffers disk &

Synchronizes process access
to disk.

* use input to read/write
blocks

device drivers

(communicate with
harddisk)

struct buf{} → disk block.
buffer.

struct bcache is buf[]

bread()

bwrite()

bget()

brelease()

use below, driver

function...

lock! exclusive access.!

idewr() → read/write to disk.

idestart() → many assembly
code.

ideintr() → interrupt handler

→ in, out
instructions

logging in disk: we want atomicity on disk changes. (when crash...)

* logging groups disk changes into transaction.

- later, installs the changes into disk one by one.

- if crash happens after logging; the entries are
replayed on restart of disk.

* link count of inode = no. of directory entries pointing to the inode.

diskinode

* for in-memory structures:-

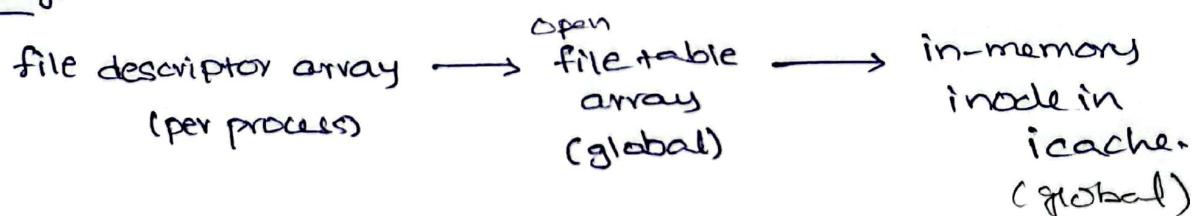
- two processes P,Q open same file, use different file entries in ftable.
 - points to same inode
 - different offsets
- two parent-child process use same ftable entry...
 - Shared offset.

reference no. in file is how many ftable entries point to it.

* on disk:

inodes, datablocks, free bitmap, logs.

in-memory:-



* updates to disk happen via buffercache.

- changes to all blocks in a systemcall are wrapped into log
for atomicity.

we either want
all or none,
in case of crash.

Nothing like,
inode is updated but
data block is trash.