

## Memory Management Basics



stack, heap exists only at run time of program.

- Desired from a multi-tasking system
- Multiple process in RAM at same time (multi-programming)
- process should not be able to see/make change in each other's code.
- Advanced requirements.
  - process can be put anywhere in RAM (pos<sup>n</sup> should not be fixed)
  - gives flexibility i.e. process can exists anywhere in RAM
  - Process need not be continuous in RAM
  - Parts of process could be moved anywhere in RAM.

### Different Times

- compile time
  - Load time
  - Run time
- } Different action at different time depending on different OS & processor.  
depending on that this leads to diff types of address binding

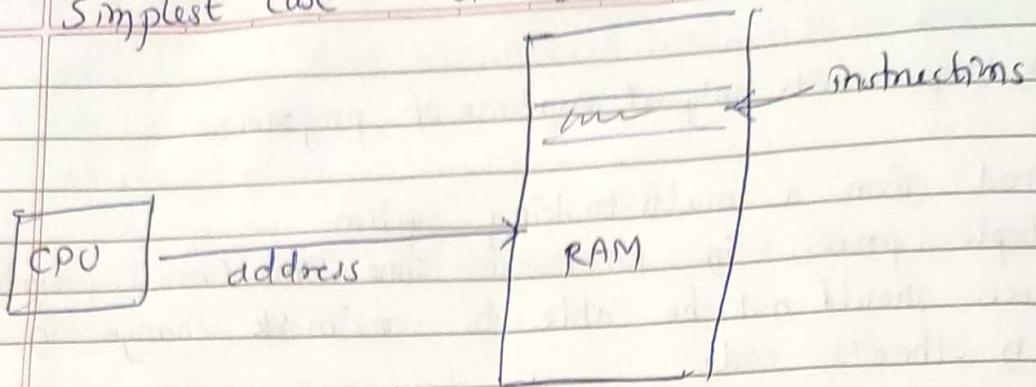
### Diff types of Add binding (determining address)

#### Run time add binding

Address of codes/variables is fixed at the time of executing the code.

It may be decided at load time but it can be change every program is running.

Simplest Case ex: DOS



Example:

address : instruction

1000 : mov \$0x300, r1

1004 : add r1, -3

1008 : jnz 1000  
(1000 add CPU send kernel)  
hor1.

start la 1000 add issued kernel mg tyache contenue fitch  
horas ani te execute horas  
mg 1004 add Bsuad

here \$0x300 address is used while executing instruction

How does this impact compiler or OS?

→ The kernel process runs hot at the same address  
Bsuad je hotat te Ram la directly miltat.

→ So, exact add of globals; add in "jmp" and "call"  
must be part of the machine instructions generated  
by compiler

→ How compiler know addresses at compile time?  
It assumes some fixed addresses for global code.

OS loads the program exactly at the same  
add. specified in the executable file.  
Non-relocatable code

Now program can execute properly.

As here compiler assumes addresses, so they it may will generate overlapping.  
so support only one program (ie, only one program in memory at a time)

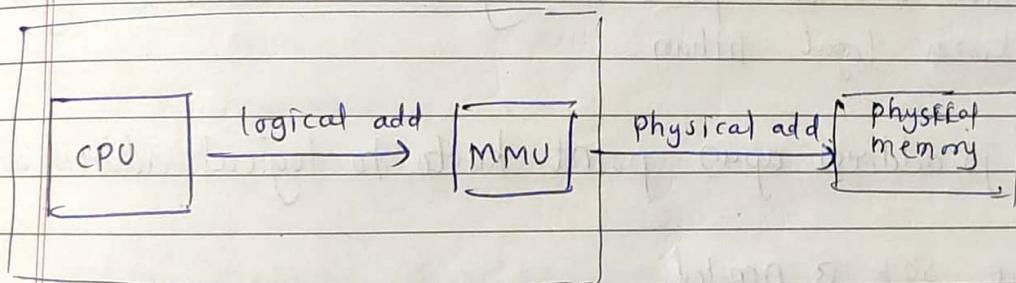
### Base / Relocation + Limit scheme

Log from memory add { limit register ( ~~base register~~)  
relocation register (Base register)

Int value check against limit reg.

If process is in RAM

- location of process → relocation register (Base)
- Total size of process → limit reg.

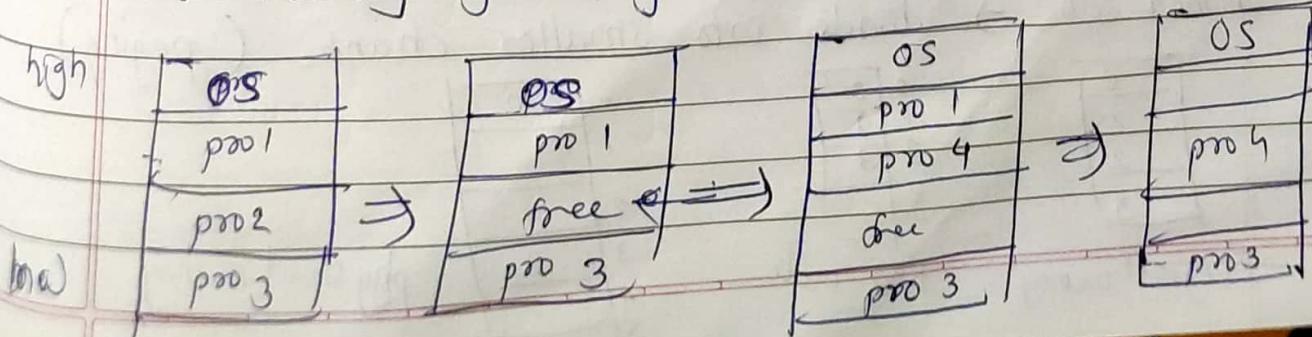


You could access code of your process & can make change in it.

Only one check → limit

Here memory get fragmented.

| Variable partition scheme



Segmentation : Multiple base + limit pairs.

Etb add issue chya reles kathal base & limit reg  
ahe yaeb pn reference geto.

text sathi reg segment, data sathi reg, -  
means text regali pair, data sathi regali, heap, stack

No. of reg increased.  $\rightarrow$  cost  $\uparrow$

Multiple base + limit pairs, with further indirection

seg reg  
used as  
indices  
here.

Base + limit pairs in memory, one reg in CPU to point  
table.

Thod slow hot as pairs memory madhe astat ani tyanna  
access karav lagat tithan.

$\rightarrow$  & ne jo address apna point karto to logical address asto

Suppose 50k is needed

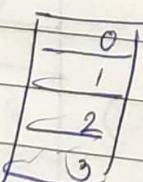
and available 30k, 30k, 40k  $\rightarrow$  they can't allocate  
(needed continuous 50k)

External fragmentation.

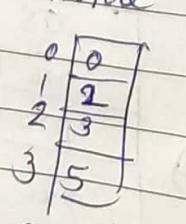
solt

compaction - chunk move karayach

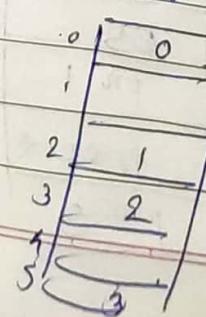
another soln  $\rightarrow$



logical memory

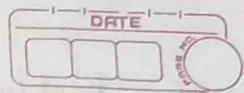


page table



physical memory

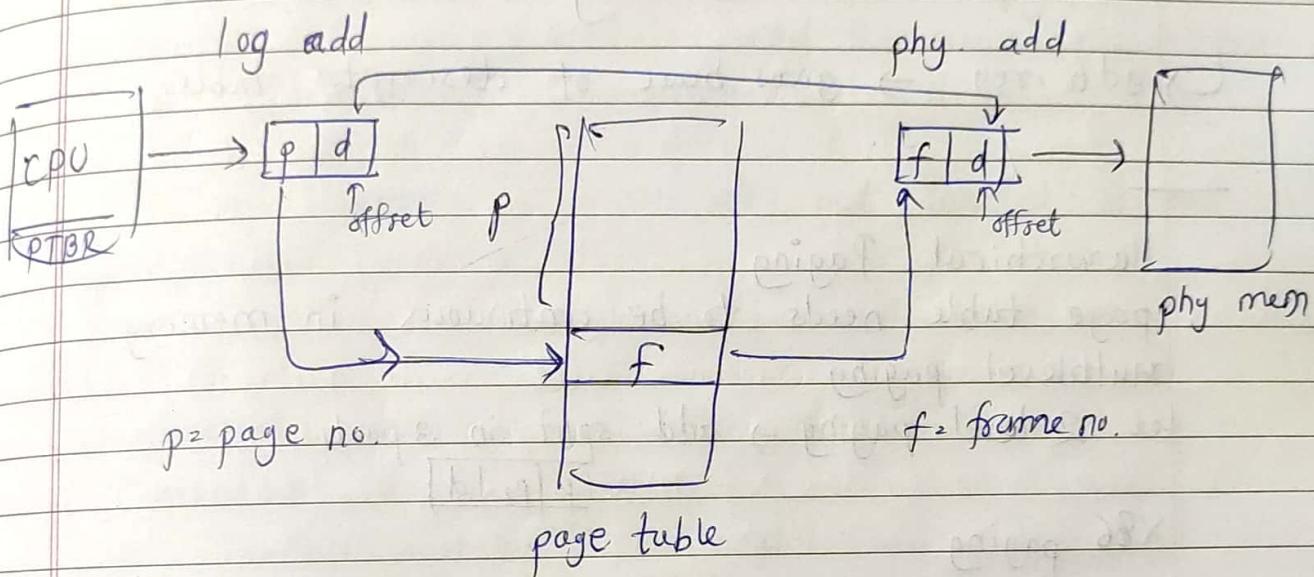




one page is continuous. (further not divided into smaller chunks)

page table located in memory set up by OS

CPU has PTBR (page table base reg) (gives location of page table to CPU)



Process composed of size  
equally divided pages

Actual memory divided into page frames  
size of frame = size of page

SS, ES → they are used as indices in segmentation table.

logical add = seg. reg. + add. compiler provides.

P

In intel 32 support both segmentation & paging.

Processor gives feature to use only one at time

## xv6 bootloader

Load boot sector in RAM at 0x7c00

The add: documented by xv6 & qemu

.code 16 → ~~(directive)~~ (Directive to assembler)  
(saying code for 16 bit)

descriptor table → array of base, limit

↳ gdtr reg → gives base of descriptor table

## Hierarchical Paging

page table needs to be continuous in memory

Multilevel paging

Let 2 level paging → add split in 3 parts

[ $p_1 | p_2 | d$ ]

## x86 paging

page directory | page table | offset

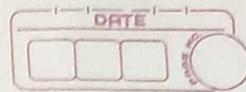
2 level → 4 kB (2<sup>12</sup>)

1 level → 4 MB

page directory, page table → in memory

PS (Page size) 0 = 4 kB, 1 = 4 MB

CR3 (Reg) → gives the base page directory table



## Segmentation + Paging

Kernel role → creating GDT/LDT, making gdtr points to GDT, page table setting / ~~creating~~ <sup>upating</sup> page directory, making CR3 points to base of page directory.

Once it is done, everything will be managed by hardware kernel doesn't need to come into picture again in memory management.

Here all memory violation also get detected in the segmentation + paging unit.

### Segmentation + Paging setup of xv6

Base = 0, Limit = 4 GB

so, logical add = linear add (always)  
(known as identity mapping)

Limit is 4 GB so set to maximum possible (length of RAM)

32-bit processor  $\rightarrow 2^{32} \rightarrow$  ~~4 GB~~ 4 GB  
Bcz did not <sup>setup</sup> any segmentation, didn't divide in any chunk.  
So segmentation is practically off.

Does segmentation setup to do identical mapping.

segmentation table B setup using global descriptor table (GDT) entry.

GDT entry  $\rightarrow$  64 bit

GDT → Global      }    id x86 only GDT  
LDT → Local

gdtdesc

.word (gdtdesc - gdt - 1)  
.long gdt

# 2 byte }  
# 4 byte } total 6 bytes  
→ # add of gdt

(Macro) SEG ASM ( we to do bit-wise operation & create GATE table entry )

SEG\_ASM ( STA X | STA\_R, 0x0, 0xffffffff )  
↑              ↑              ↑  
permitIn        base        limit  
20              = 4FB3

with this type of setup, ( $\text{base} = 0$ ,  $\text{limit} = 4\text{GB}$ )  
every address in kernel will have identical mapping  
(remain as P3)

CRO  $\rightarrow$  registers

last bit PE = 1 they protecting mode enabled

jmp \$(\_SEG\_KCODE << 3), \$start32  
↑  
B) code seg ↑

EIP  
from now code onwards run in 32 bit mode.  
previously running in 16-bit mode.

left shift by 3  $\rightarrow$  Bcz there are 3 bits in every seg selector at the right side & those 3 bits are not part of seg. selector.

These 3 bits  $\rightarrow$  TI (2) Table index  
RPL (10) Request privilege

As go in 32 bit

Here Base will be added in address translation going thru GDT table, But still it will be identical mapping

mov \$start, %esp // start  $\rightarrow$  0x7C00  
esp = 0x7C00

stack goes downwards, so we will use stack from 7C00 to 0.

C code kernel from disk  
read code  $\rightarrow$  put in RAM  $\rightarrow$  do paging setup  $\rightarrow$  run kernel

readreg  $\rightarrow$  read from disk into RAM.

func bootmain  $\rightarrow$  job B to load kernel from disk in memory  
xv6.img  $\rightarrow$  concatenation of bootblock kernel

ELF

unt entry;  $\rightarrow$  1st instruction to be executed in ELF

ph off  $\rightarrow$  program header offset (ptrs add)

ph num  $\rightarrow$  ph number (no. of entries)

we will load kernel at this  
virtual add

vaddr → virtual add → 0x80100000  
paddr → phy → 0x00100000

we can't like this any particular other file  
but kernel can demand, that he wants to occupy particular  
region.

EXTMEM = 0x100000

KERNBASE = 0x80000000 (First kernel virtual add.)

KERNLINK = (KERNBASE + EXTMEM) → add where kernel  
is linked.

elf = (struct elfhdr \*) 0x100000

here we can write this add (points) because we already  
set MMU properly (0 to 4GB)

SECTSIZE = 512 (normally) (Default)

readsect → read single sector at offset into destination

ELF\_MAGIC → 0x464C457FU // \x7FELF is little  
endian

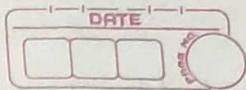
storb → feels memory location with particular value

Eth ph → memsz > ph → filesz

mg storb la call karun o fill karun

entry = (void (\*) (void)) (df → entry);

↓  
type casting particular value in function ptr.



entry.s → code from kernel (upto thr paying 3 off)  
 4KB → 2 level  
 4MB → 1 level

Setup Page Table

only two entries of 512 → both will point to frame 0.  
 $4\text{MB} \times 512 = 2\text{GB}$

→ [dir | offset] →

NPDENTRIES 1024

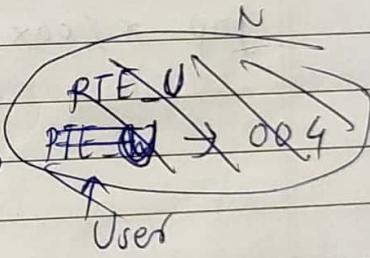
RPT 1024

pde.t - 4 Bytes

PDXSHIFT = 22

Flags set.

PTE\_P → 0001 , PTE\_W → 002 ,  
 ↑ present              ↑ writeable



PTE\_PS → 80 (7th bit=1) page size 4MB

CR4\_PSE 0x 0000 0010 // → enable page size extension  
 which is 4MB

entrypgdir → part of kernel code

kernel code compiled linked by linker  
 so the add all generated in ELF file  
 are going to be all virtual add.

kernel ld → directive to link the kernel

virtual to phy

V2P\_WD(x)  $(x - \text{KERNBASE})$

P2V\_WW(x) c + )

Two flags for enabling paging

CRO\_PG 0x8000 0000 // paging

CRO\_WP 0x0001000 // write protection to enable

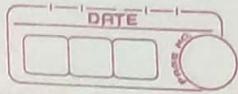
KSTACKSIZE → 4096 // size pre-process kernel stack

.comm stack, KSTACKSIZE // directive to assembler  
allocate data of size KSTACKSIZE

kvmalloc → reinitialize page table (own)  
segment → segment descriptor (own) kern

jmp \*%eax → writing like <sup>wants</sup> bez jmp to particular  
fix phy add.

## Process



Process related data structure in kernel code,  
→ different for diff kernel.

PCB (Process Control Block) (common in all kernel)  
1 PCB for each process

PCB identifier for process in kernel data structure

struct task\_struct → in linux kernel

struct proc → in xv6

Fields also in PCB vary for diff kernel.

Process state → running, sleeping, terminated etc.

List of open files in PCB  
fd[] → array of pbs.

fd (file descriptor) → index in the array of pbs maintained in PCB.

(kernel)  
(in linux) in array 1st 3 are already pointing to files which are opened by kernel.

0 → std::in

1 → std::out

2 → std::error

So 1st file db we get 3

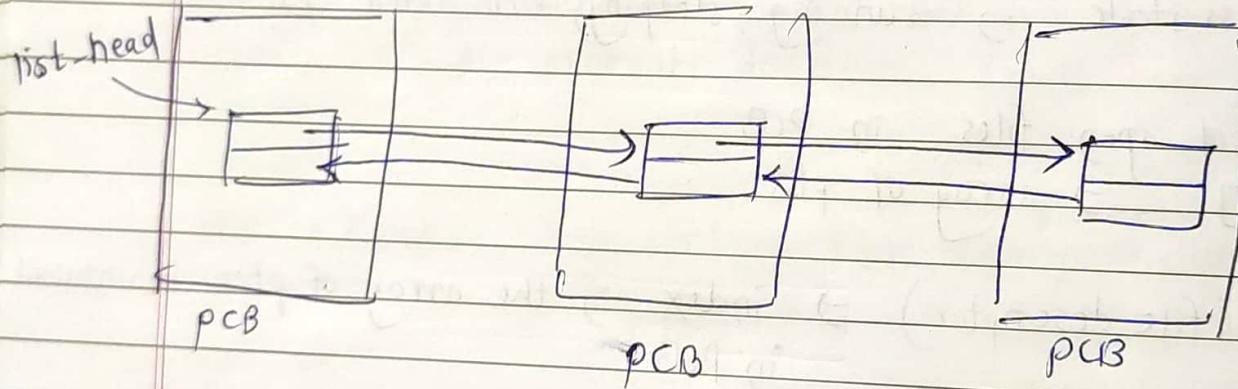
All these data structure in kernel space

all process in xv6 in ptable (global variable)  
maintaining array of `proc[NPROC];`  
`ptable.proc[ ]`

Process Queue / Lsts inside OS

in xv6 only array so scheduler has to do more work  
look for all process and decide which one to schedule

In Linux

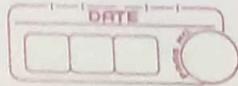


like doubly linked list  
(linking external members)

Process suru astanna  $\setminus \cup$  request ali (like scanf)  
fir mahit nahi ki keyboard  $\setminus \cup$  key kadhhi press  
karnar mhanuy tyala queue madhe takar kahi  
mean nahi karat.

So, tyala tithan cut karun wait queue la fakkat  
ans complete ghal ki ethay cut karun punha  
queue madhe

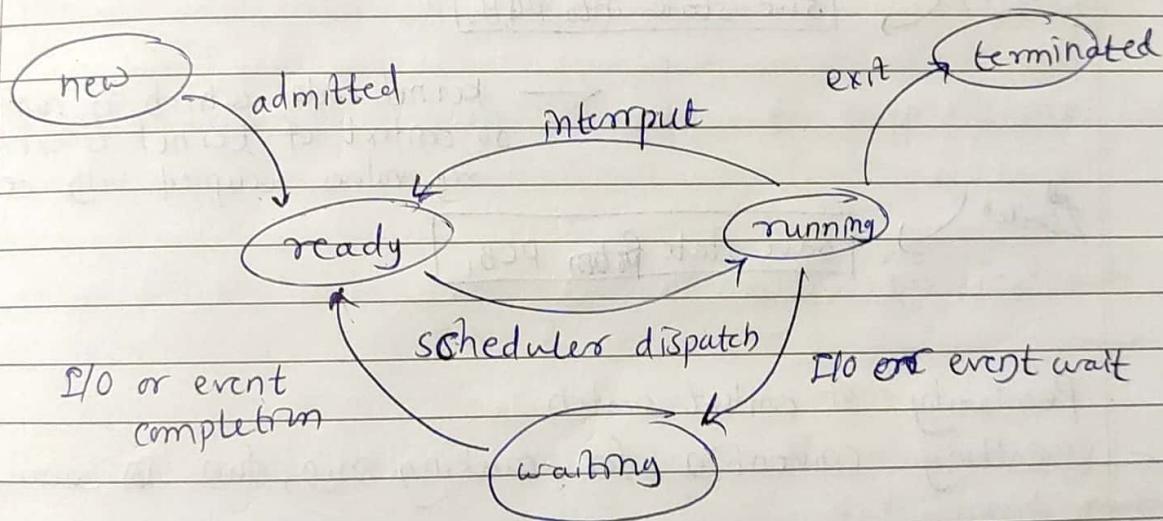
see diagram of process state



only PCB B completely constructed by the kernel  
PCB B moved to ready queue & state also marked as ready (means can be picked up scheduler for scheduling)

If now scheduler selected, change state to running  
then pass off control to process  
It has to be execute on particular processor

For terminating it must call exit  
~~after it will scheduler and will new process~~



ctrl + C → will generate signal

when signal occurs kernel will notify process & process will its own signal handler and if not kernel will run its own default sig handler.

Here kernel code will run, looks for current process & terminate it.

(ETH process exit ne terminate nahi hota  
kernel code ne hote)

Ata illegal memory access ka seg fault yerun termi  
hot tevha, hardware interrupt hot mg te run hodon  
check kare ki kathali process ne illegal memory  
access keli ani tyala terminate kare.

"Giving up" CPU by a process or blocking

sys\_reld

file f = current  $\rightarrow$  fdarray [fa];  
    ↑  
    PCB

Context switch

1<sup>st</sup> switch

[Save state into PCB]

← kernel code which is running  
so context of kernel is executing  
reg values occupied with kernel

2<sup>nd</sup> switch

[Save state from PCB]

Peculiarity of context switch

calling current for making my, fun in same process

Context switch one process take place of another

→ Eth caller sangal mala he he reg use karayachet  
remaining callee use kernel

po context switch madhe tas nast purn use hotat  
by another process



So if code of context switch made the library  
for the calling convention need to,  
then the assembly machine code also

Process in xv6

Max possible virtual add  $\rightarrow$  KERNBASE (0x8000000)  
2 GB

2 GB + 1 MB  $\rightarrow$  onwards  $\beta$  kernel code/text, data

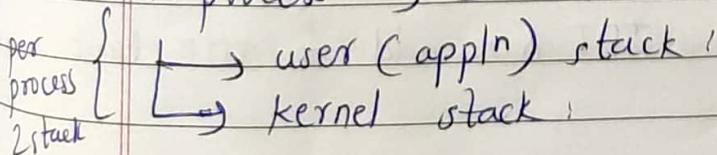
PHYSADDR  $\rightarrow$  20 to 24 MB

Kernel needs to create mapping to map kernel code,  
data, process data

Kernel loaded at 0x100000 (1 MB) physical add.  
0 to 0x100000  $\rightarrow$  BIOS & devices

Kernel  $\beta$  loaded at lower region of phy memory  
as some systems may not have that much memory

A process  $\rightarrow$  2 stacks



3rd stack

used by kernel kernel stack used by scheduler  
Th3  
3 not for  
per process  
(each process)

Generic Stack

struct proc

umt sz; // highest vir. mem add available of page  
code

shouldn't, sz < 2GB

guard page → 4K, stack → 4K

The whole setup has to be done as a part of fork(),  
pgdir → pointing to page dir

## Handling Traps

Kernel needs to get all related to event → that information has to be made available to kernel by hardware.  
Some kind of data passing needs to be done b/w kernel & processor.

Current privilege is → 1. cs (in 2 bits) in the field (1)

If your in protected mode

cs will not be multiplied by 16 (left shift by 4)

IDTR & IDT

Entries in IDT called Gate (Gate is code seg rip pair)  
in memory

IDTR register gives IDT base add & IDT limit

umt var1, 2;

var2 defined of 2 bytes

value for CS      offset  
 /      \  
 SETGATE (gate, Btrap, sel, off, d)

SETGATE (idt[], 0, SEG\_KCODE << 3, vectors[i], 0);  
 (particular entry we trying to set)  
 add of idt  
 Btrap

SEG\_KCODE = 1 << 3

actual add of fun  
to be executed

vectors[] → in vectors.S with 256 entries pointers.

DPL\_USER → 3 (~~4~~)

in xv6 same for all

global vector 64

vector 64

pushl \$0                // push 0

pushl \$64                // push no.

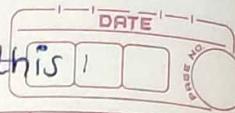
jmp alltraps                // jump to all traps

CPL m % cs B <= DPL

all levels set in descriptor  
all levels set to 0, except sys call level 3.

checking that you should go always go on low level  
not on higher.

On int instruction / interrupt CPU does this



Save %esp & %ss only if PL (privilege level) < CPL

processor gets kernel stack using task segment descriptor. (set of extra reg stores particular values, esp for kernel)

IF → interrupt enable flag

set %cr & %rip values to the descriptor  
changes to refer table (set to value from vector table)

sys call has to return error no to the calling process

struct trapframe → struct of all reg pushed in stack in reverse order.

so we can access content of stack by pointing to trapframe tf. and we access all data by structure trapframe tf → pointing to relevant location in kernel stack

SEG\_KDATA = 2

Trapframe is a part of kernel stack

It can access return value, trap no, old value of reg, ss, esp, everything

int trapframe in struct trapframe {

uint esp;

ushort ss;

ushort padding;

// 2 bytes

// 2 bytes

These 2 bytes

are additionally added

to insure that the content of struct actually matches the size of trapframe built with above code



actual size is 2 bytes

2 bytes was pushed as 4 bytes, so access it as 4 bytes  
so padding.

Boot loader call main

main → call tinit → set up all the gates in idt array  
also set particular value of syscall (says it should  
jump from level 3 to 0)

→ goto vectors.s

Do similar job

push 0

push current value of interrupt

jmp to all traps

In all traps push more things (reg)

→ call trap

Now we have trapframe already on stack.

In trap fun

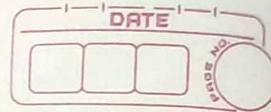
If sys call then

long proc() → tf = tf; // ~~particular~~ currently  
running process

calls syscall(); //

Who sets trap no → It is on trapframe. It was push  
on trapframe by vectors.s (as each level pushing  
its own number)

struct trapframe{} in c file was definition of struct,  
while in assembly code we created a instance of  
that struct.



INT make processor look up in IDT, jump to location given by IDT.

pushing ss, esp (in initialize phase of trapframe)  $\rightarrow$  B esp  
(optional register)

If switch from 3 to 0 they push, then only first will pop.

Interrupt only can occur at the end of instruction not in betw

for all hardware inturr -  $istrap = 0$   
soft  $istrap = 1$

on 0 clears FLIF (Interrupt flag)

kernel we only general purpose registers

code making system call like read, open, write  
code of read, — in sys.S

make system call number, copied by eax

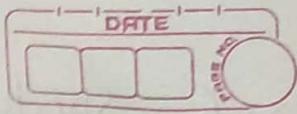
they call int

jump on vectors

push 0

push 69

jump all traps



in all traps  
pushing reg  
pushal (pushall)

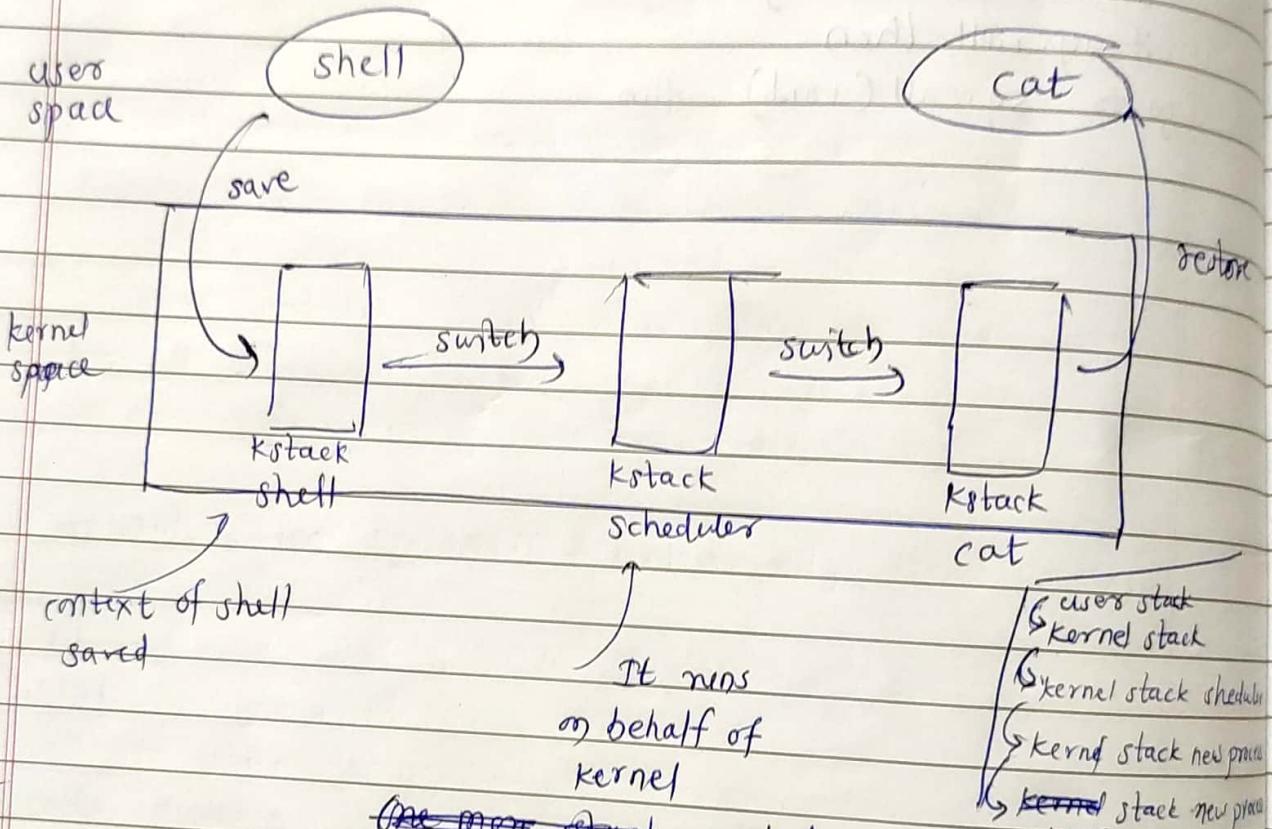
they call trap

they check for sys.call  
if sys.call then  
go to syscall (word) fun

## xv6 Scheduler

Transition from process P1 executing on CPU to process P2 executing on CPU, in between kernel code i.e., scheduler will execute.

4 stacks need to change /

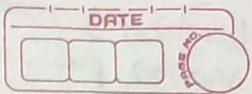


kernel stack of process, it changes to stack of scheduler  
they only scheduler can run.

Again stack needs to change to kernel stack of cat (new)

`mpinit()` → detect other processes  
`mpmain()` → start execution

`mpmain` calls `scheduler()`



mpmain will run on all processor  $\rightarrow$  with scheduler  
will run on all processor

kernel code runs on every CPU independently and  
scheduler process independently

only data structure which represents collection of  
process  $\rightarrow$  struct {

// array of struct & proc  
{ ptable;

scheduler()

sti()  $\rightarrow$  set Interrupt flag (Enable interrupt)

asm  $\rightarrow$  directive to write assembly code in C

Iterate thru proc array

check if p->state  $\neq$  RUNNABLE (ready) ~~if not then~~

c->proc = p; // c pointer to cpu

// remembering process is related to cpu

// per-CPU state

struct context {

edi;

esi;

ebx;

ebp;

esp;

struct cpu {

struct segdesc gdt[NSEGS];

another / ~~one~~ gdt table per-CPU

p->state = RUNNING; // only change state

load  $\rightarrow$  reading from stack put value in proper reg.  
save  $\rightarrow$  push values to stack  
 $\downarrow$   
th B will likely to be read

switch (& (c  $\rightarrow$  scheduler), p  $\rightarrow$  context);

as passing add may it will change later

In switch

switch

p B process to be sched.  
c  $\rightarrow$  struct cpu

~~executed, updated~~

~~pop & put value in esp & returns~~

p  $\rightarrow$  context  $\rightarrow$  struct context  
c  $\rightarrow$  scheduler  $\rightarrow$  also struct context

switch  $\rightarrow$  assembly code in switch.s

// void switch (struct context \*\*old, struct context \*new)

here it will save current registers on stack, in struct cont  
at add in \*old

switch stack to new & pop previous, <sup>saved</sup> reg

eax = old, edx = new

i) switching stack

\*old = updated old stack

esp = new = edx

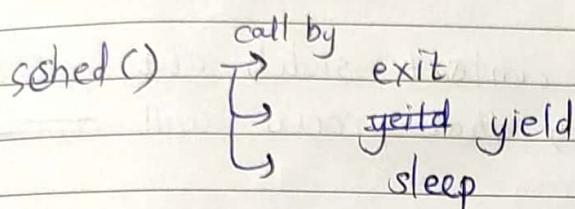
So now stack pointer becomes add given by p->context  
and & (c  $\rightarrow$  scheduler)  $\rightarrow$  the add was save on earlier  
stack.

So now stack started pointing to kernel stack of  
new process



ret → pop from stack one more time (eip will be pop)  
// here will pop from esp → fun where to return  
→ pop & return to that location

& here it will not return to thB (switch to scheduler)



process executing → timer interrupt occurs → came to trap → check condn → call yield.

In yield()

put myproc() → state = RUNNABLE

shed calls switch

switch (&p → context, mycpu() → scheduler);

context of currently running process will be pushed  
on its own <sup>kernel</sup> stack.

so context will be pushed & p → context will  
be made to point it.

mycpu() → scheduler → already pointing kernel's  
it's own stack.

from here you load the context of kernel.

context of kernel → same context that was  
saved when switch was called from the  
scheduler

So calling this, thr will not return here.

It will return to switch in scheduler

Now scheduler fun run in loop again once again call switch

If that process was context switch out as a part of shed (yield), then that process will ~~context~~ return in shed.

So one possibility of returning process B here.

so <sup>one</sup> possibility of start executing ~~star~~ process start from next line of switch in yield

so it will return yield

from yield it will return to trap

so whichever process B schedule will 1st return from kernel & then <sup>start</sup> returning to application code

In sleep

$p \rightarrow state = SLEEPING;$  } put process B sleep  
sched(); } (chang. state)  
call shed

// like blocking

// so loop of scheduler will  
not select it

There B wakeup!

Find out particular process

$p \rightarrow state = RUNNABLE$  } // change state

# IPC (Inter Process Communication)



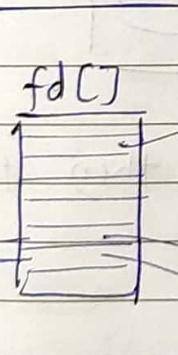
process  
independent → run their own don't need anyone  
cooperating → can affect or get affected by other  
(including data)

Two type → shared memory → common memory  
→ message passing → message pass from  
kernel (kernel has to  
do more work)

In shared memory kernel simply available common memory

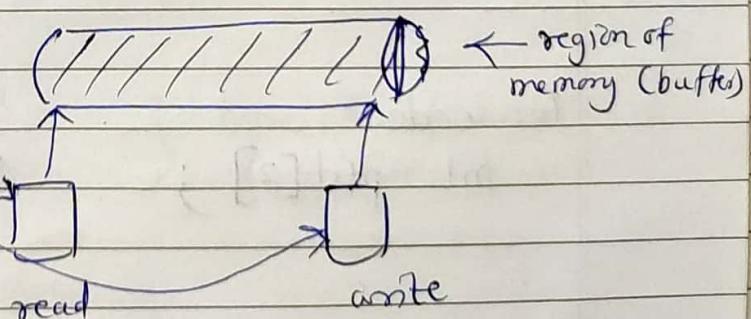
pipes → ordinary / unnamed (also queue FIFO)  
→ named

for process it creates PCB, which has file descriptors:



file (file that are opened)

Now when we make pipe  
sys call



So this  
two indices  
written by  
kernel

also creates  
struct file  
for read & write

& they have some posh.

And kernel make pointer from fd[ ] to pointing  
read write  
So that

In fork PCB gets duplicated.

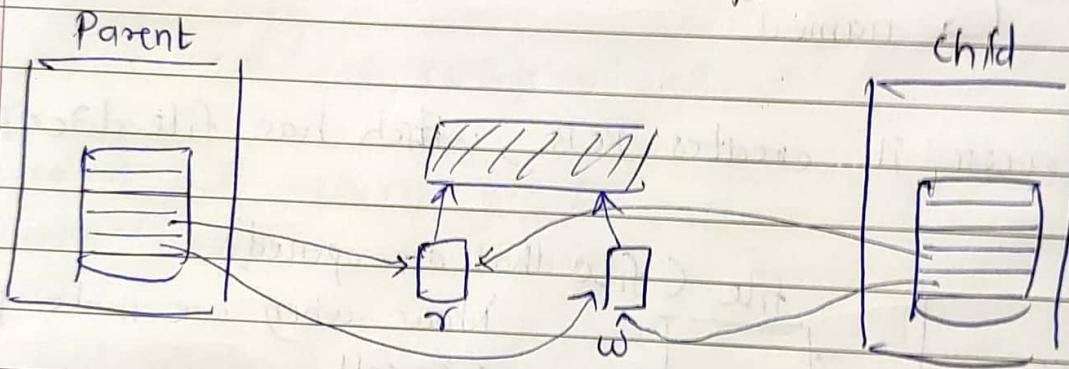
Suppose pipe B called & then fork B called & in fork PCB get duplicated then

Then fd[ ] of child PCB also points to read write from same memory

i.e. buff B shared b/w two (parent & child)

Buff B part of kernel memory not a part of process data. They have only access to read & write thru fd[ ].

OS job to insure read & write happen B FIFO



This happen when we call pipe & then fork

for code

```
int pfd[2];
```

```
Two fd pfd = 0 // for reading  
pfd = 1 // writing
```

```
pipe(pfd); // they will point accordingly to  
           pipe
```

while writing or reading from pipe  
only one end must be open.  
It's unidirectional communication mechanism

So for writing,  
`close(pfd[0]); // require to close unused end`  
`write(pfd[1], string, no.of bytes);`

for reading,

~~read~~ `close(fd[1]);`  
~~read~~ `(pdf[0], buf, no.of bytes to read);`

kernel often does delay write

`close(0);`  
`fd = open(filename, O_RDONLY);`

At a time at least one fd is closed & another file opened  
delay till the file has a file descriptor before.  
As open try to locate 1st available slot in fd.

Now,

`scanf("%s", buf); } here scanf will read`  
`printf("%s\n", buf); } from file.`

~~This is redirection of stdio~~

`ls > /tmp/file` (output of ls stored in /tmp/file)

shell does → creates child but calling fork →  
in child before exec, simply `close(1)` and open  
`/tmp/file` → so output goes written in /tmp/file

<                   ⇒ std i/p redirection  
  >                   ⇒ std o/p redirection

dup → duplicates file descriptor

cat /etc/passwd | grep jay

Here, cat will not write entire output to the pipe & then grep will read

cat B continuously writing to the pipe & the grep B continuously reading from the pipe

Both cat & grep exist at same time

• Named pipes → FIFO

process can create a file cat acts as pipe more powerful than ordinary pipes.

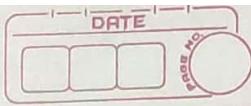
int mkfifo (const char \* pathname, mode\_t mode);

mkfifo /tmp/xyz

ls -l /tmp/xyz

→ prw-rw-r-- 1

at start p which not present for normal file, kernel knows it B special file



## Shared Memory

`shmget (IPC_PRIVATE, size, S_IRUSR | S_IWUSR);`

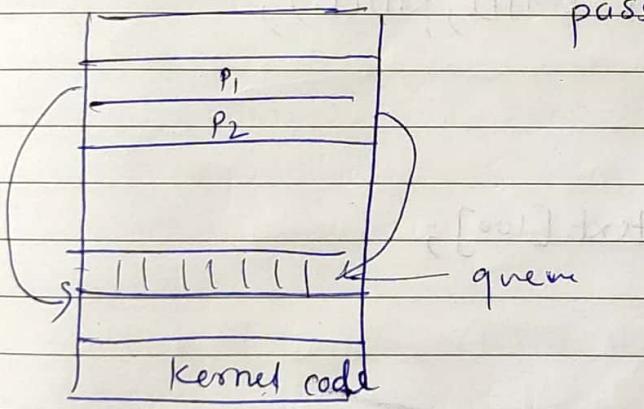
key      size

`ptr = (char*) shmat (id, NULL, 0);`  
 ↑  
 ptr to shared mem  
 Now can read & write from it.  
`sprintf (ptr, "string");`

To detach shared mem from its add space  
`shmdt (ptr);`

## Message Passing

P<sub>1</sub>, P<sub>2</sub> 2 processes and they want to talk using msg passing

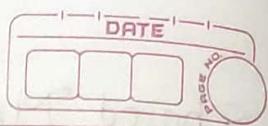


kernel creates / makes available queue of msgs to processes.  
 Now, P<sub>1</sub> & P<sub>2</sub> can access to queue & they can communicate.

Kernel job B to create communicating link bet processes

`key = ftok ("str", no); // creates unique key`

↑  
 To generate same other process should know these two parameters.



msgid = msgget (key, 0666 | IPC\_CREAT);

// msgget creates a msg queue

// with help of id , can access to queue

msgsend // msgsnd to send msg)

msgsnd (msgid, & message, sizeof(message), 0),  
data

In other file reads

msgrecv (msgid, & message, sizeof(message), 1, 0);

saying that want  
msg of type 1.

// To destroy msg queue

msgctl (msgid, IPC\_RMID, NULL);

```
struct msg_buffer {  
    long msg_type;  
    long char msg_text [100];  
} message;
```



syst-pipe)

argptr() // use to fetch arguments passed by user  
in the system call.

argptr(0, (void \*) &fd, 2 \* sizeof(fd[0])) < 0

pipealloc(&rf, &wf) → // modify rf & wf  
// allocate array & make  
// rf & wf point to array

fd = fdalloc(rf) // → will allocate index in file dB  
(table) & make ptr point to rf

pipealloc()

\*fd = filealloc() → gives struct file

kalloc() → give a page frame of 4K

struct file {

enum { FD\_NONE, FD\_PIPE, FD\_INODE } type;

;

struct pipe \* pipe;

struct mode \* ip;

Treating pipe also as  
file and file also as file.

so one of ptr we will  
be use

This for pipe if file is actual pipe  
for file

Type will be decided on variable type

When we read from pipe it will call sys.read()  
sys.read() → fileread()

fileread(f, addr, n)

↑  
fileptr

char \*ptr

no. of bytes to be read

→ fileread(f) {

if ( $f \rightarrow type == FD\_PIPE$ )

return piperead( $f \rightarrow pipe, addr, n$ ); // read from pipe

Similar to next cond) check for FD\_INODE

Then it will read from file actual on disk.

piperead(struct pipe \*p, char \*addr, int n)

// code

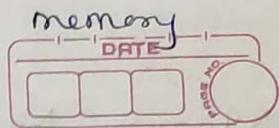
// waiting if the queue is empty

Then in loop reading from pipe & putting one byte  
after one in  $addr[i]$

wakeup(& $p \rightarrow pwrite$ ); // wake up writer  
// bcz there can be writer which was waiting  
bcz queue was fullled.

release(& $p \rightarrow lock$ );

ipcs // shows all existing shared



### pipewrite()

/ Do almost same thing

, wait if pipe was full

// write data into pfp array of pipe

wakeup (if p->nread);

release (if p->lock);

return n;

}

## syswrite → filewrite → <sup>calls</sup> pipewrite

# Memory Management

variable partition scheme

strategies to find free chunk

- Best fit
- Worst fit
- First fit

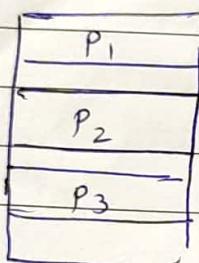
External fragmentation problem

30K, 40K, 20K chunks free

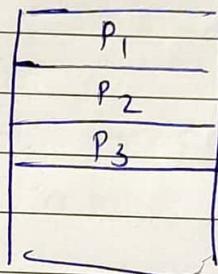
Need 50K

→ Total free = 90K, still can't allocate 50K.

Compaction



compaction



more them  
to be continuous

Fixed sized partitions

- ↳ Lead to internal frag
- ↳ 4K normal size

after During context switch of the process, load the PTBR using PCB. OS

once this done, Then process B ready to run.

All the memory <sup>add</sup> translation for process will be done by hardware. (Process will run as if it is running its own)

OS has to do some global level task.

- maintain a list of all page frames

TLB entries can be maintained by → least frequent use  
Irrespective of it TLB speeds up. → most frequent use

- DB advantage of page

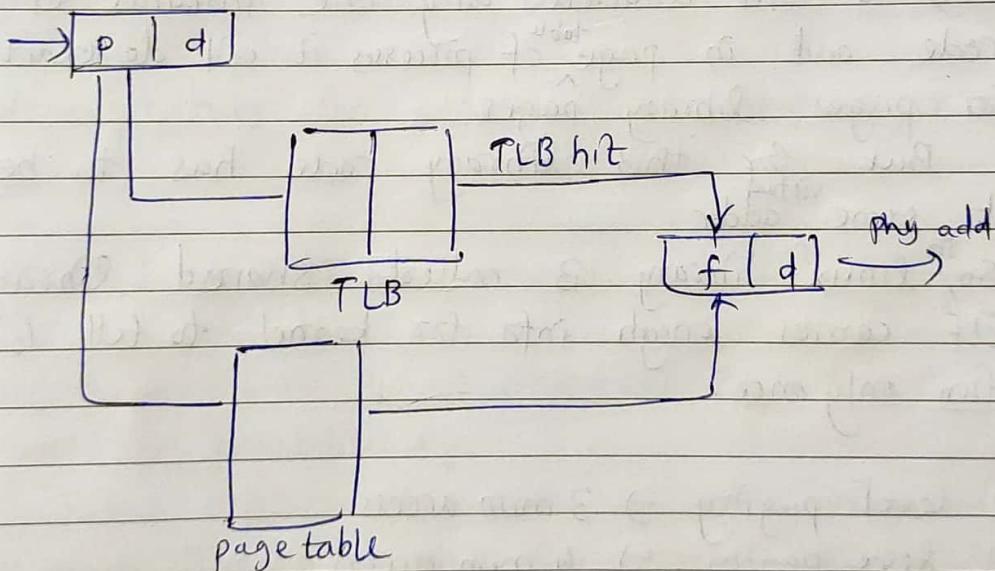
memory access results in 2 memory access

- 1 actual memory access
- 2 page table

Translation Lookaside buffer (TLB) → 2 entries

It's subset of page table entries

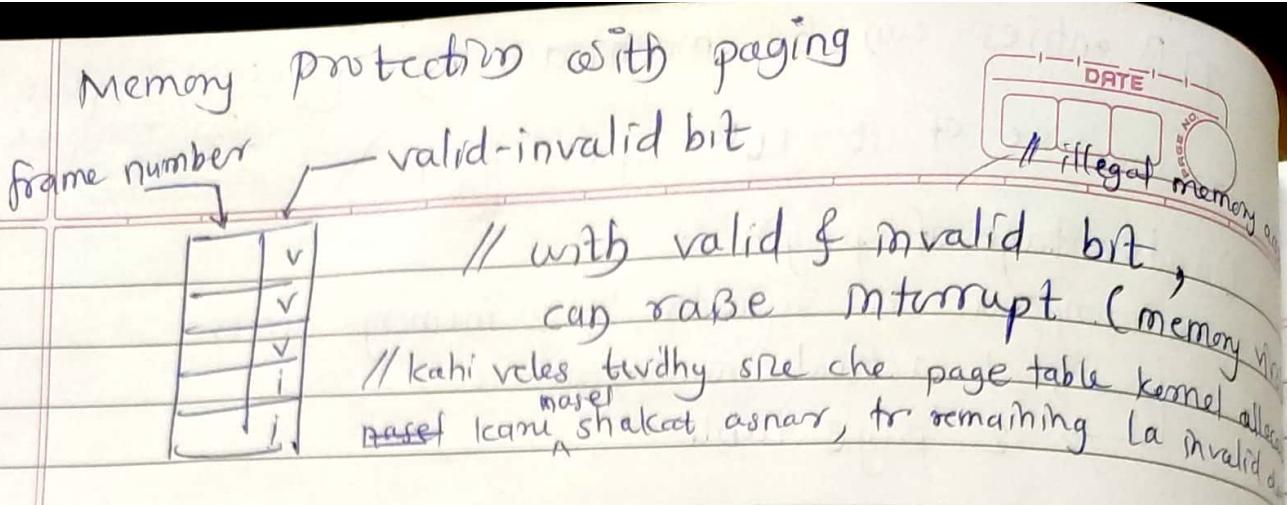
→ doesn't frame no. eq to <sup>number to number</sup> page contains less,  
So it has to maintain both page no. & frame no.



First try to look up in TLB if get it → TLB hit  
& if don't it's TLB MISS, then look it for in  
page table

Now if TLB hit → 1 memory access  
TLB miss → 2 memory access

TLB access & updating is done by hardware  
automatically. There is no kernel involved here.



shared pages (eg library with paging)

suppose  $p_1, p_2, P_3$  uses same library code  
while normally allocating for processes it will allocate  
as separate library for  $p_1$ , for  $p_2, P_3$  (End up allocated  
same library pages multiple times)

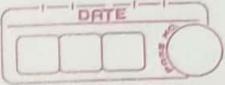
But If kernel can identify which part of ELF file  
is library code and which part of file is not library.  
Then it will allocate only one instance of library  
code and in page <sup>table</sup> of processes it will do exact mapping  
for pages library pages.

But for this library code has to be mapped  
at same virtual addr

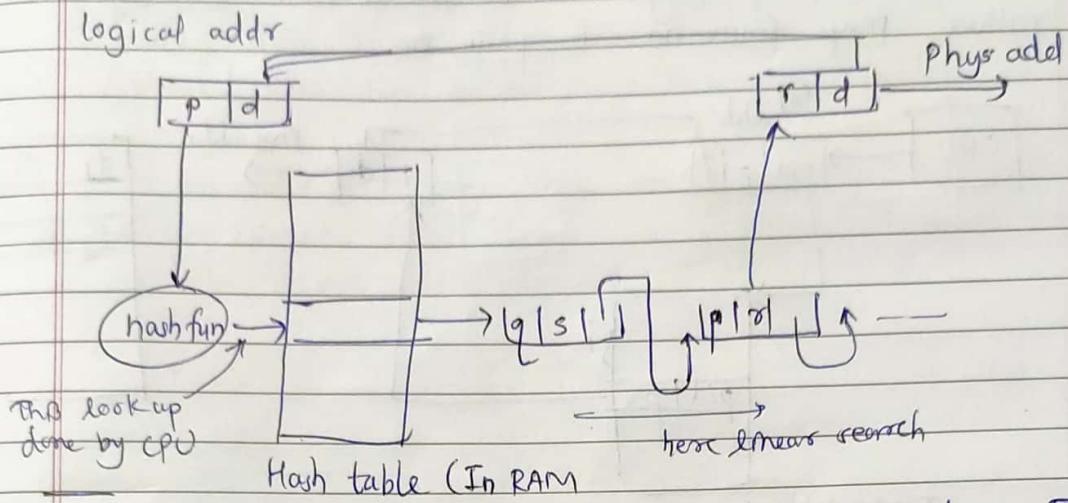
So, In Linux library is called shared library.  
ELF comes enough info for kernel to tell load library  
fun only once.

- 2 level paging  $\rightarrow$  3 mem access
- 3 level paging  $\rightarrow$  4 mem access

Hierarchical Paging



## Hashed page table



$p \rightarrow$  page no. It doesn't use as index. It pass thru hash fun. That will be index in hash fun.  
In hashing for problem of collision  $\rightarrow$  chaining is used.

Here entries for q & p collided. From p, r is used.

## Inverted page Table

we have process  $P_1 \rightarrow$  page table  $P_2$ , pro  $P_2 \rightarrow$  ptable  $P_2$ ,  
pro  $P_3 \rightarrow$  pagetable  $P_3$

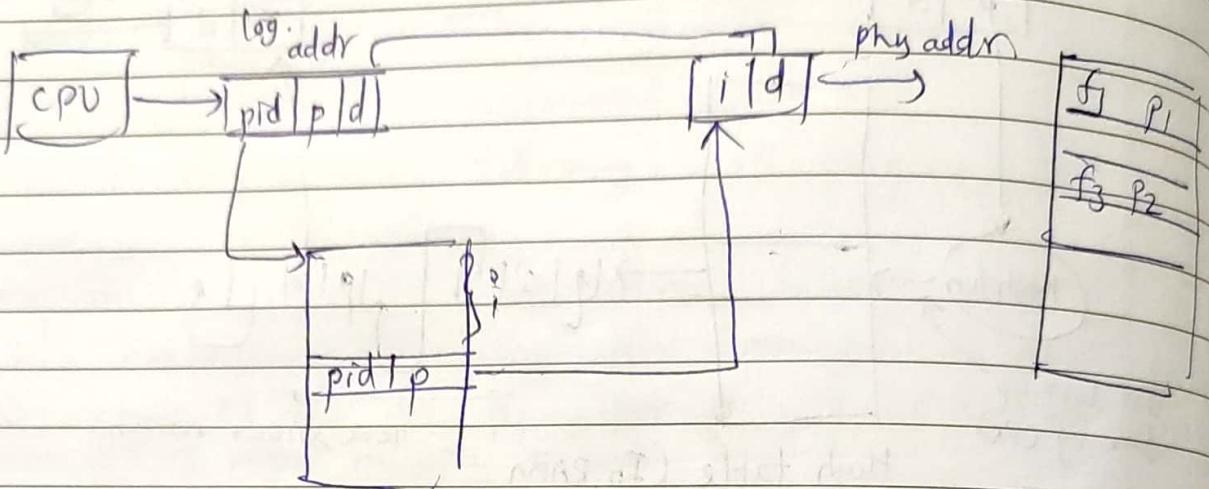
Now all of them think entire ram belongs to them, that's why they have mapping to the same RAM.

So, more no. of pages, mapping all to same RAM.  
All mapping to same page frame in RAM.  
RAM is same.

So think inverted.

No matter how many processes  $\rightarrow$  they still have to be contained within same memory.

so didn't map page no. to frame no. (As page no. for each process)  
 rather map frame no. to page no.



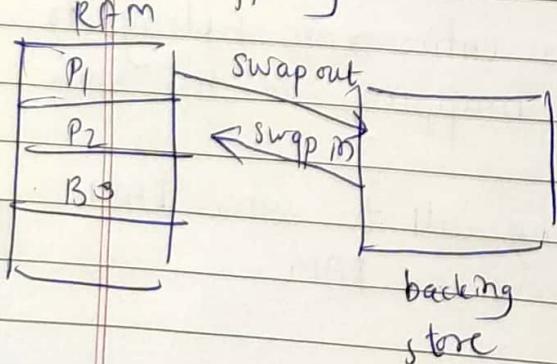
page table (frame table)

Problem B that  $f_1$  can belong to  $P_1$ ,  $f_3$  to  $P_2$   
 so also store pid & p.no.

This is one global table not perprocess table.  
 There are no perprocess table.

pid  $\rightarrow$  not generate by ~~CPU~~ CPU  
 $\rightarrow$  loaded by kernel when process executing

### Swapping



Suppose there is no space for  $P_1$  in RAM, But at a time we can run only one process, so we can swap  $P_1$  & allocate pages for  $P_2$ . So we can allocate pages more than efficiency of RAM.



of process

Entire memory img  $\uparrow$  map into page table (including code, data, stack, heap)

shared also has to be map into page table

As for std c library they mostly don't get removed as every process needs them. Only when there is pressure on mem & you want page out.

→ In day of pte (2<sup>20</sup>)  
→ can do this in (1000 pages \* 2<sup>10</sup>)  
→ enough physical address space

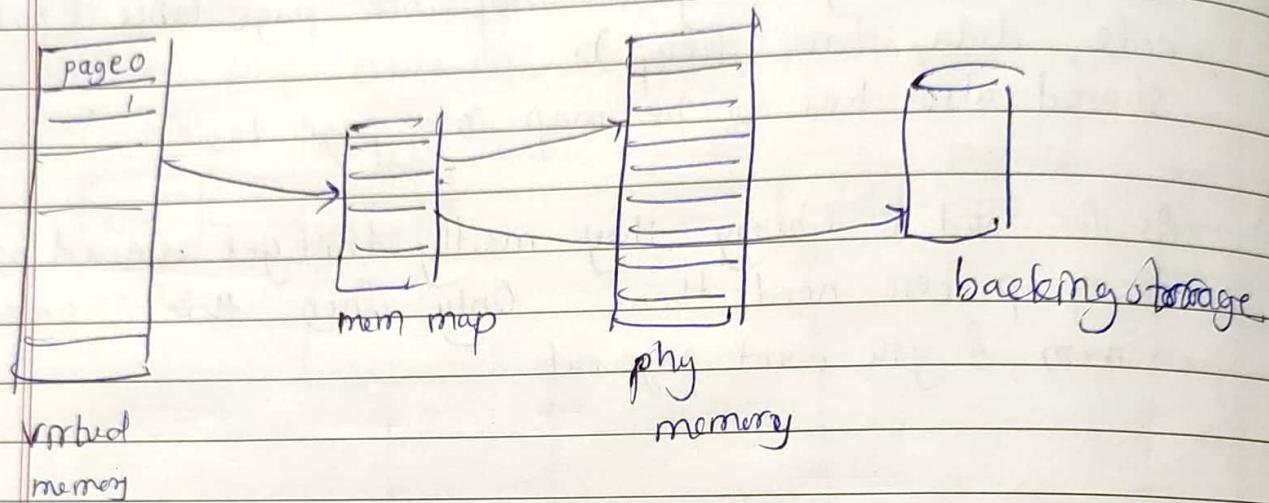
map bld function

→ using memory allocation →  
→ use of desired malloc  
→ and deal with heap  
→ deal with free heap  
(copying of func.)

sdk → sys call → ask kernel to allocate certain region  
virtual mem & return ptr to it.

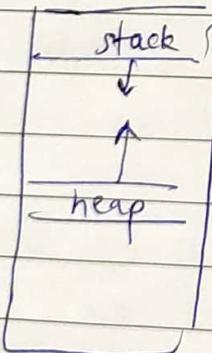


## Virtual Memory



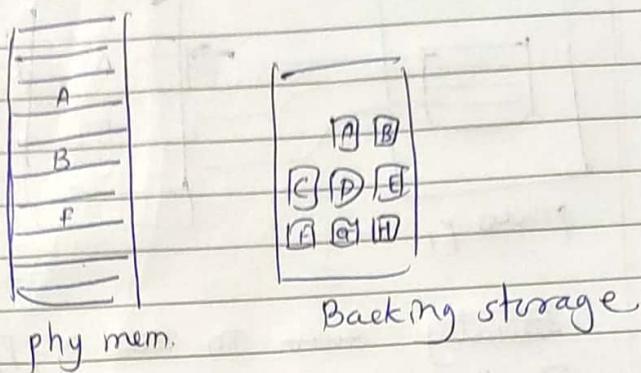
$\text{vir} > \text{phy}$ . So job of OS to creates mem.m (like page table) it will not only map to RAM but some also to backing storage

## Virtual add space



→ Normally yarathi purv space adhish allocate keli asti po in virtual add space stack la 1 page & heap 1 page and they later size can be increased (left for growth)

## Demand Paging



Pages in phy mem also in backing stg as we never know which one go out & which one come back in.  
So in backing stg  $\rightarrow$  all loaded.

valid - invalid bit

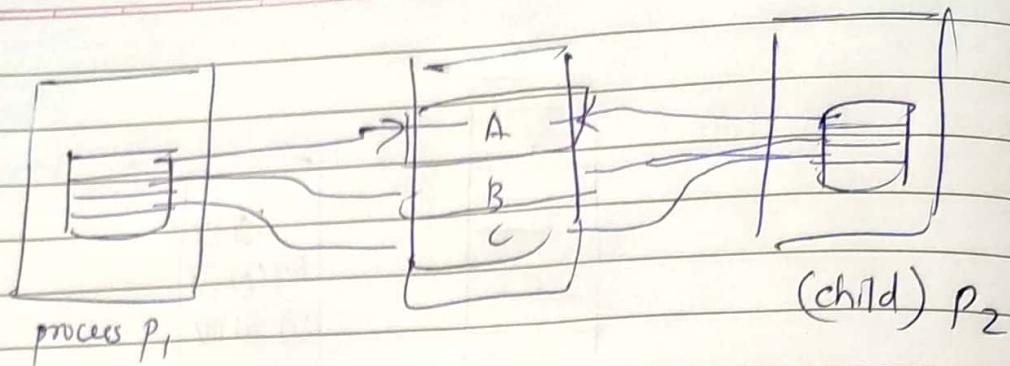
i  $\rightarrow$  not in mem or illegal  $\rightarrow$  raised trap called page fault

movamny 0x100, 0x110, 20  
Yamul overlapping hoil & this will lead to page fault.

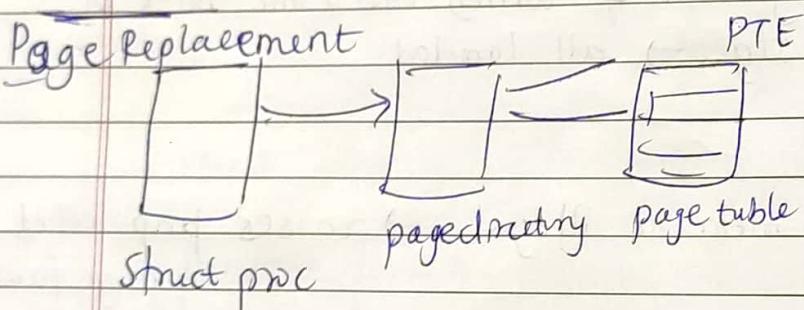
Jr ekhadi pnt partially over asel ani mg page fault  
ala kr te kas revert back honarr.

for demand paging  $\rightarrow$  CPU has to insure with something goes wrong it has to revert back as if nothing happen.

i  $\rightarrow$  while bringing page from backing storage to phy. mem  $\rightarrow$  Invokes read from  $\rightarrow$  takes some milisecond  $\rightarrow$  so Meanwhile scheduler will schedule another process.



Initially save in PCB P<sub>1</sub>, P<sub>2</sub> all the pages get shared.  
In PCB there should one extra bit.  
Indicating there should copy on write.



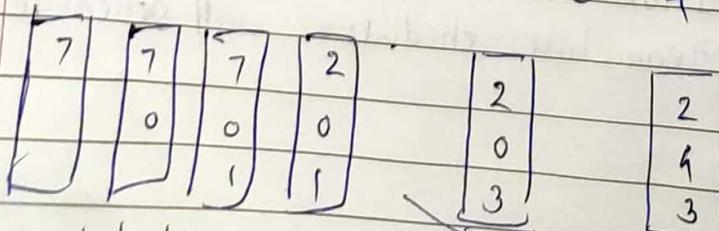
one more bit,  
PTE has ~~dirty~~ dirty bit get set by hardware  
when you modify corresponding page.  
only modify pages anty to disk.

### Page replacement algo.

1) FIFO

→ optimal page replacement

7 0 1 2 0 3 0 4 2 3 0 3



next hya entries made 0 snarv  
→ last ahe so tyda  
replace ket

Eth check ket tr  
next chya reference made 7 last ahe.  
so tyda replace ket.

optimal  $\beta$  not possible

Least recently used (LRU)  $\rightarrow$  approximation of optimal.

stack algorithms

for  $n$  frames

$$\text{No. of page fault} \geq \frac{n+1}{n} \text{ frame no. of page fault}$$

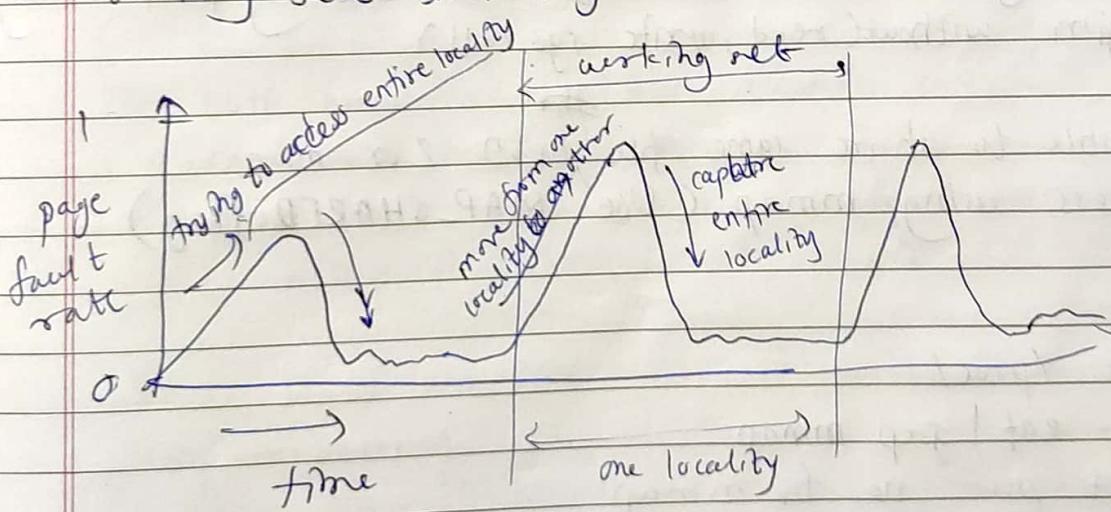
## Thrashing

CPU utilization  $\beta$  low when

$\rightarrow$  less no. of process

$\rightarrow$  more no. of process, but most time is page fault

## working sets and Page fault rates



```
fd = open(argv[1], O_RDWR);
```

```
data = mmap(NULL, 4096, PROT_READ | PROT_WRITE, MAP_SHARED, fd, 0);
```

map 1st 4k bytes into memory region

and returning ptr to region.

can be shared with other process

- When mmap call is did, → mapping done only conceptually
- so basically page table entry created for process
- having storage in the file
- Now get access in memory region

```
munmap(data, 1024); // undo mapping  
close(fd);
```

Here you read & write from file directly from memory region without (read, write sys call).

possible to share same file betw 2 or more process using mmap (Use <sup>#</sup> MAP\_SHARED flag)

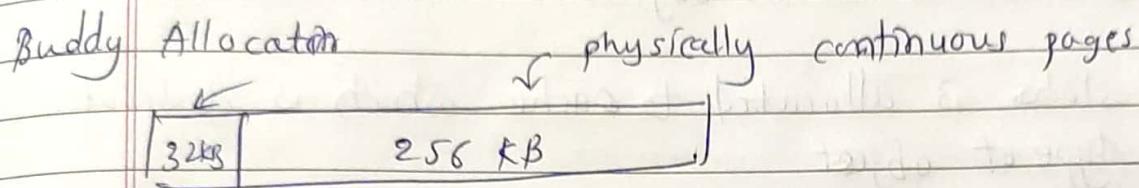
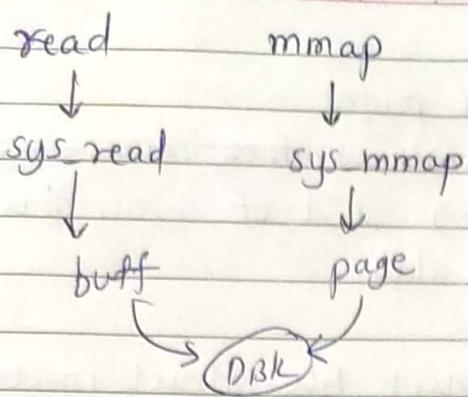
```
cat /proc/
```

```
ps -ef | grep mmap
```

→ get proc no. for mmap

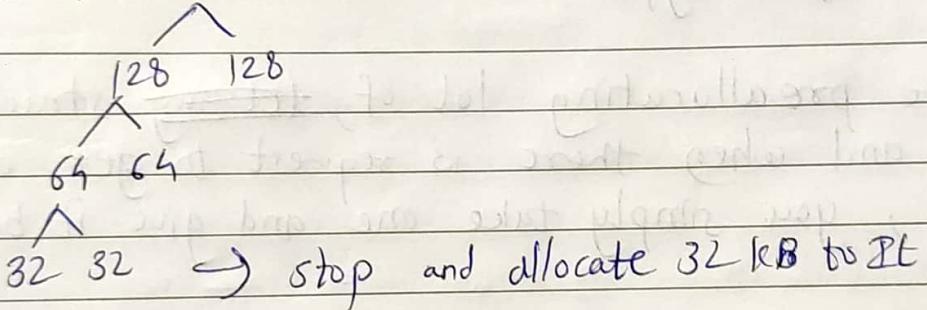
→ cd no.

```
cat maps → memory map of process
```



Suppose need 30 KB

Break 256



7th kuth mem - we ahe kar identify karayach

→ kernel will set set limit

→ suppose here set lower limit to be 1 KB ~~that is~~  
i.e., below it won't get divided.

so,  $256 / 1 = 256$  so it will have 256 bit long bit pattern.

Suppose here 82 KB B at start, so at start all starting 32 bits will be set to 1 from 0.

## Slab allocator

Slabs → set of consecutive pages.

When kernel loaded → combines free frames (~~continuous~~)

→ list of free frames → list of conti free frames

→ combine some frames into slab.

Kernel often needs struct proc, struct file, struct inode, struct buffer  
etc → called as kernel object.

Set of slabs is allocated to cache which is dedicated to one type of object.

Take a slab → divide it into equal size chunks of that object type

Let say, so here preallocating lot of, ~~let say~~ struct proc here and when there is request to get one of them, you simply take one and give it back

Here preallocating → so wastage of memory.

Advantage → caches are already there

when loaded kernel → kernel created slabs

→ from slabs kernel created proc cache, file cache

inode cache, IO cache etc

→ so very fast, very quick

→ No internal fragmentation but wastage if don't use all completely (don't use all cache completely)

## Threading



Thread create →

separate concurrent flow of execution B created

set of inst executed independently.  
process B also thread

concurrency → progress at same time

parallelism → execution at same time

suppose f<sub>1</sub>, f<sub>2</sub> are tr mg f<sub>1</sub> kahil millisee mg f<sub>2</sub> millisee  
sathi as Eth as nahi y ki f<sub>1</sub> purn execute hotay mg  
f<sub>2</sub> hotay.

Eth parallelly do gh sbut <sup>run hotat</sup> ~~start start~~. (multi core)

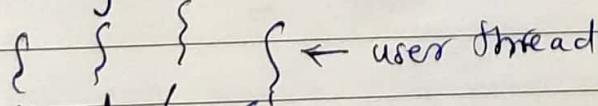
Process creating B typically heavy weight , (amount of time)  
Thread light →

Going do fork, exec → <sup>fork</sup> going to duplicate everything

Part of fun → going to duplicate only some part of it

Pratically it switch to context switch while scheduling  
diff threads.

Many-one model



K ← kernel thread

Trouble with thB B

Suppose ~~thB~~ <sup>thread</sup> user  
thread ~~make~~ <sup>do</sup> scanf,

→ make sys-call → thB  
sys-call and suspend the process

→ go process get blocked → so  
what happened to other threads →  
they also get blocked.

One-one Model  $\rightarrow$  need to many resources

main {  
    pthread create (

}

pthread create ( f - )

clone (f, - ) library

{

f

main

kernel

If don't have clone  $\rightarrow$  then many ~~to~~ one

pthread create ( - f1 - - )

thread {

context c

stack ps

// other things needed

}

}

th = new struct  
setup th  
switch ()

Read  $\rightarrow$  clone

getcontext

set context

, make set context

## Thread pools

→ create threads without scheduling

Scheduler activations for thread.

suppose 3 user threads, 2 kernel threads  
many to many

Application

```
f() {  
    scanf();  
}  
g() {  
    recv();  
}
```

→ after allocating thread to process  
→ when scanf → sys call → thread block  
→ then call sched → goes to scheduler, schedules process

As here we don't know thread get block, so that we can create new thread

On thread B always ON which should not block → scheduler

SIG\_KILL → <sup>rule says</sup> Not possible to handle sigkill

```
int *p = 123;  
main () {  
    Signal ( SIGSEGV, seghandler );  
    *p = 1000; // writing at any random pointer  
    // it will handle in this (segfault)
```

we can modify that fun

but if segfault occurs before th3 func line of code,  
then that segfault will be handled by default signal  
handler.

```

int *p = 1234; // initialize any random loc
void seghandler(int signo) {
    printf("Seg fault occurred \n");
    return;
}
int main() {
    signal(SIGSEGV, seghandler);
    *p = 100; // Try to write at random location.
    return 0;
}

```

Here we are override default handler.

Basic - B

The default action for SIGSEGV is to terminate your program. But here we are override.

So after this line of code, for every segfault our code will execute.

Q- Why signal handler goes to infinite loop?

As we override, so for every instruction that triggers a sigsegv, this (our) handler is called & the instruction is restarted. But here our handler does nothing, it doesn't fix anything, to fix it did nothing to fix what was wrong in the first place with faulting instruction.

When instruction B restarted, it will fault again. That's why it goes in infinite loop.