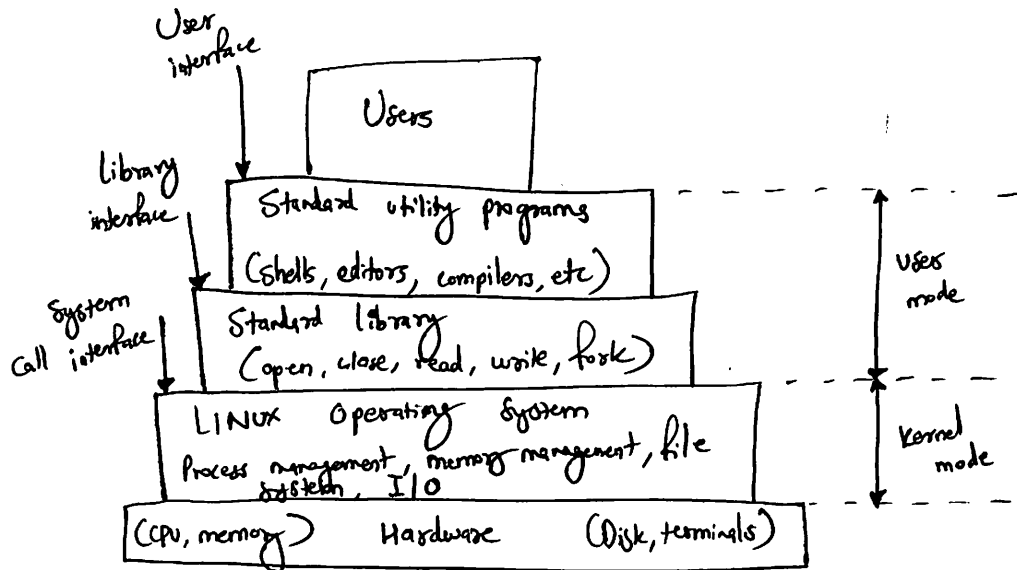Linux — • principle of least surprise (eg.— ls A°  rm A°)

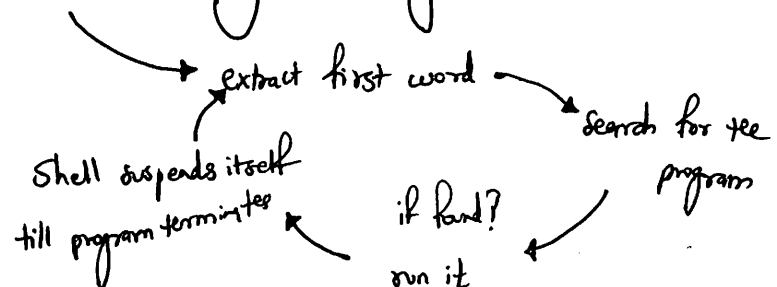• power — basic elements — ∞ variety of combinations

Layers in linux system



Hardware — • control hardware

• system call interface to all programs

How system calls — 
are made()

1. Putting the arguments in registers (or sometimes stack)

2. Issuing trap instructions to switch from user mode to kernel mode

3. To issue trap instr. → library is provided → 1 procedure/produce/system call

Synchronous interrupt

as they are produced by the CPU while executing instructions (done after terminating the execution of the (instruction)

Shell → ordinary user program

→ extract first word → search for the program

Shell suspends itself till program terminates ← if found? ← run it

② Shell — when it starts up — access to ————→ Standard input | File descriptors
                                                              | 0
                              ————— Standard output | 1

                              ————— Standard error | 2

redirection:- Stdin <
             Stdout >

( run process in background — & )

Everything in LINUX is a file, including directory

Linux kernel

————→                    System calls

I/O component          Memory management        Process management
  └•VFS ?               • Virtual Memory          • Signal handling
   • Terminals                    page
   • Socket             • Paging replacement      • Process/thread creation
   • Network              ^                          and termination
   protocols
   • Network            • Page cache              • CPU scheduling
   device
   drivers
   • File systems

   • Generic
     block layer

   • I/O schedulers

   • Block device
     drivers

——→ Interrupts                    ——→ dispatch latency
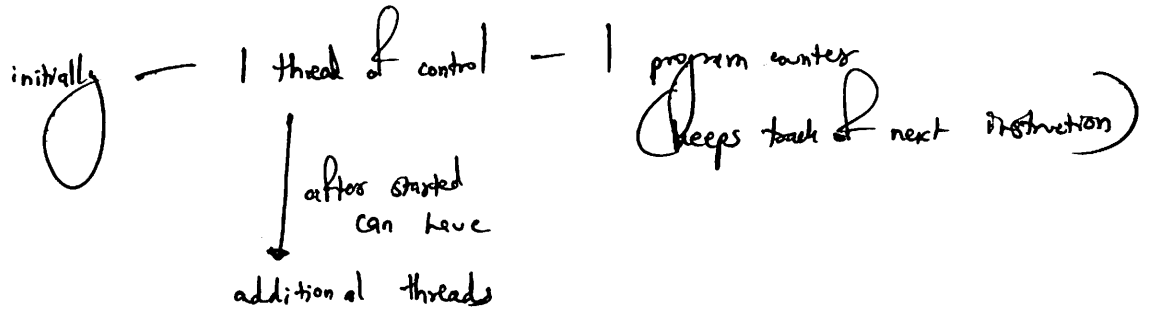
——→ Dispatcher — module that gives control of the CPU to the process selected
                by the short-term scheduler
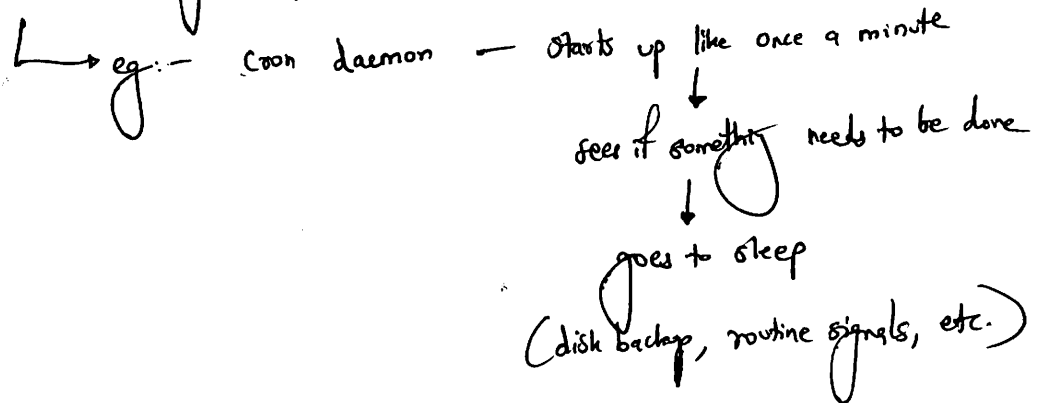          ( It takes care of context switches, switching to user mode )

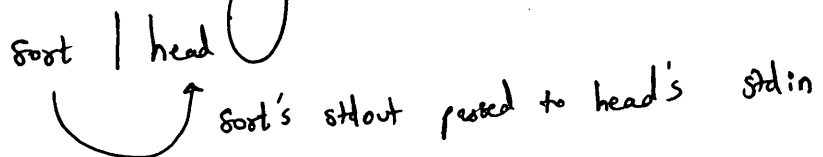Block device drivers — seeks and random accesses are allowed

## Processes

initially — 1 thread of control — 1 program counter

(keeps track of next instruction)

| after started can have
↓
additional threads

daemons — background processes
↳ eg:- cron daemon — starts up like once a minute
↓
sees if something needs to be done
↓
goes to sleep
(disk backup, routine signals, etc.)

fork() — memory images, variables, registers — identical to parent and child
pid different — 0 for child ; non-zero for parent
↳ this is child's PID
(getpid())

Process communication can be done using Pipes ⓘ

sort | head
↳ sort's stdout pasted to head's stdin

Ⅱ Also using signals

process decides what will happen when signal arrives
(terminate, catch, ignore)
↳ if catch → signal handling procedure

④

Process can only send signals to its process groups
↳ Parents, siblings, children

Ctrl^C
↓
SIGINT
(terminate the
application)

---

• SIGALRM
↓
the alarm clock
has gone off

---

Ctrl^D

register a EOF on standard input

---

like waitid()
↓
return
a
siginfo_t

Signals — 
• SIG CHLD - child terminates process
• SIG KILL - send to kill a process
    (cannot be ignored or killed)
• SIG PIPE - process has written to a pipe which has no readers
• SIG SEGV - process has referenced an invalid memory address
• SIG TERM — Used to request a process terminate gracefully
• SIG USR₁, SIG USR₂ — available for application-defined purposes

waitpid() — waitpid (pid, &statloc, opts)
    → it is andpassed (&)
    → Caller blocks or return if no child is already terminated (usually 0)

allows caller to wait for specific child
(if -1 → any old child)

address of variable that will be set to child's exit status
↓
to know the return status

---

execve (name, argv, envp)
    → pointer to the environment array
    ( used to pass stuff like
      • terminal type
      • home directory )

name of the file to be executed

pointer to the argument array

eg:- execve("/usr/bin/ls", ["ls", "-l"], 0x7fffa3746e78) = 0
    ↳ strace of 'ls'

pause ( ) — suspend the process till the next signal arrives

↳ when a program has nothing to do but stay IDLE and wait for a signal to start action

linux kernel — represents tasks — task_struct

task_struct → • scheduling (priority, CPU time, time spent sleeping) — which to run next
• Memory image
• Machine registers
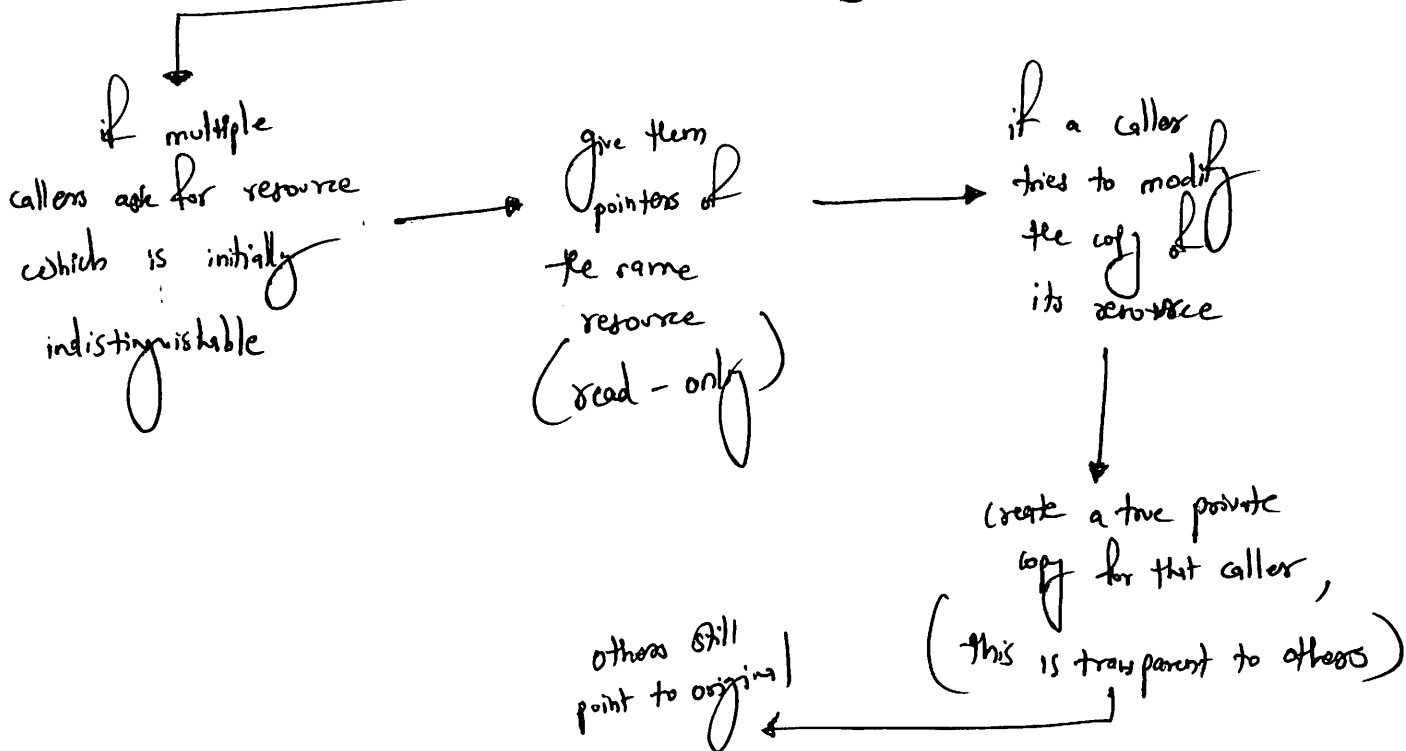ls ←——————— • file descriptor table
• Signals (which are ignored, caught, temporarily blocked)
• system call state
• kernel stack

When fork() — child has to copy memory of parent

BUT — copying is expensive, so — Copy on Write

if multiple callers ask for resource which is initially indistinguishable → give them pointers of the same resource (read - only) → if a caller tries to modify the copy of its resource ↓ create a true private copy for that caller, (this is transparent to others)

others still point to original ←

⑥ clone() → Historically (before linux), file descriptors (open files), signal handlers, alarms, etc. were per process, not per thread

with clone, it can be either per process or per thread

$$pid = clone(function, stack\_ptr, sharing\_flags, args)$$

Linux now scheduling — threads. not processes

real time threads
{
0 — highest priority
⋮
99 — lowest priority
}

non real time threads
{
100
⋮
139
}

$$Time = no. \ of \ clock \ ticks$$
(in linux)      (older — clock ran 1000 Hz — each 1 ms)

nice — each thread — default - 0
range — -20 to +19
high priority

Static priority of each thread

Use system call : nice(value) to change it
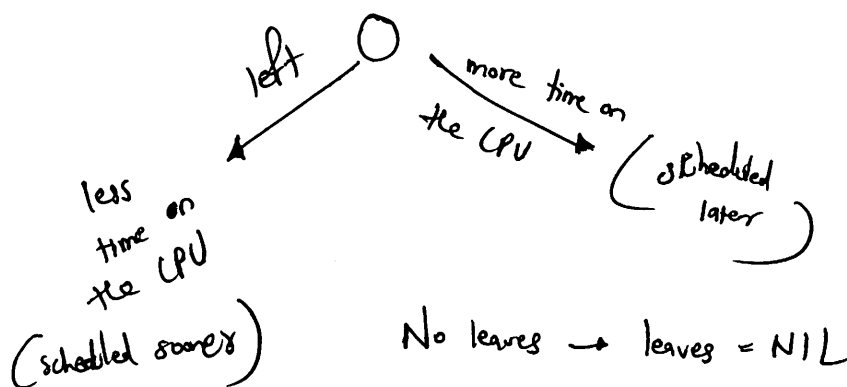
How are PIDs allocated

↳ hash table entry for PIDs — if collision — there is chaining

counter from X to y — increment each time; ~~sometimes~~ if pid is assigned, it is skipped

Completely fair scheduler

red black trees

(tasks are ordered ⟶ amount of time they spent running on CPU)
"V runtime" ⟶ internally called

Internal node ⟶ task



left / more time on the CPU (scheduled later)

less time on the CPU (scheduled sooner)

No leaves ⟶ leaves = NIL
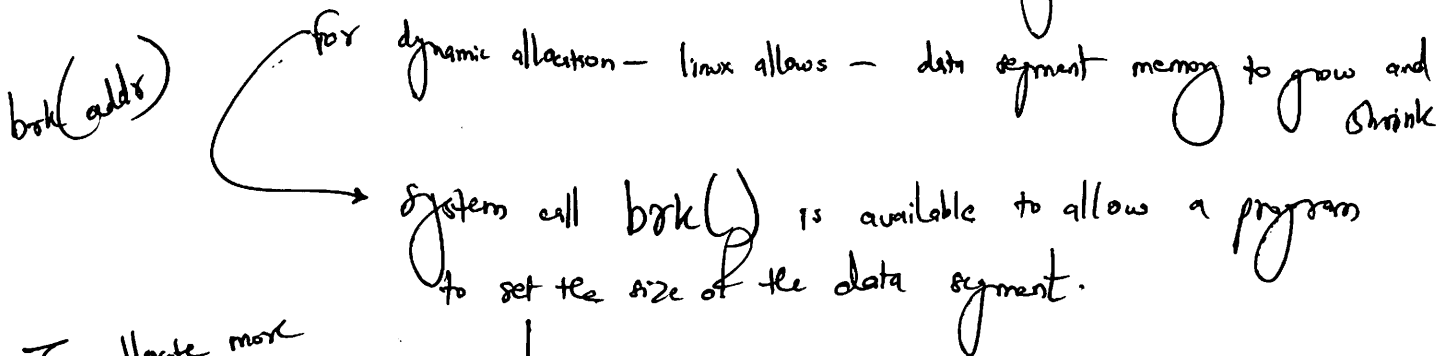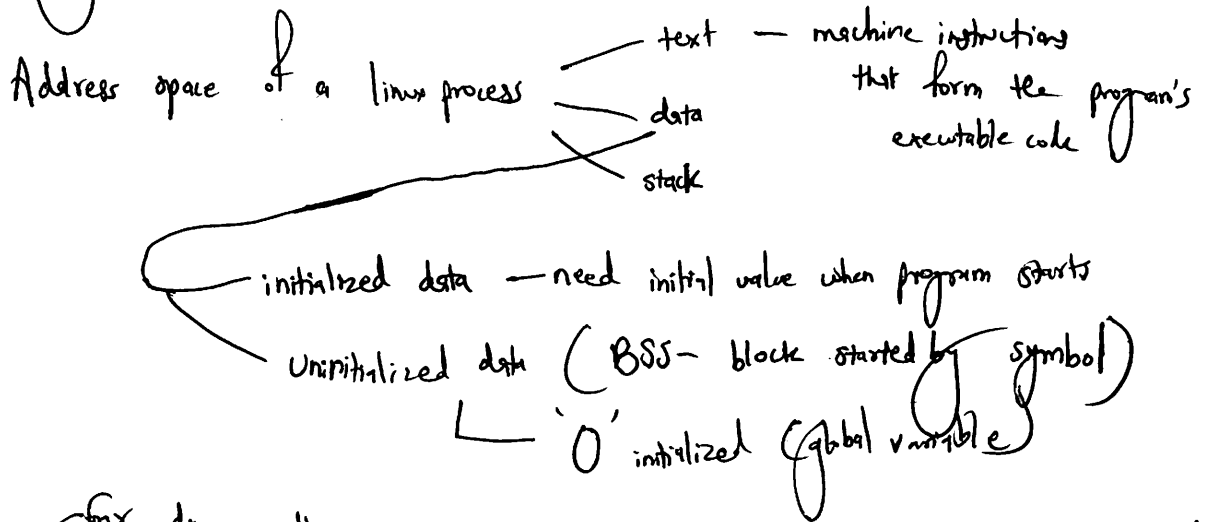
Algorithm ⟶
- Schedule the leftmost node in the tree $[O(1)]$
- Increment that vruntime
- Compare it with the leftmost node in the tree
- If less, continue to run
- If more, go to the appropriate place in tree

$$[O(n \log n)]$$

⑧

1. data structure name

→ **runqueue** —— contains the runnable tasks (ready queue)

→ **waitqueue** —— ~~associated~~ associated with events that tasks may weight on

→ **spinlock** → for preventing concurrent ~~modification~~ modification
^ to datastructures like waitqueues

# Memory Management — Linux

Address space of a linux process
- text — machine instructions that form the program's executable code
- data
- stack

initialized data — need initial value when program starts

Uninitialized data (BSS — block started by symbol)
└─ `0` initialized (global variable)

brk(addr) ⟍ for dynamic allocation — linux allows — data segment memory to grow and shrink

→ System call brk() is available to allow a program to set the size of the data segment.

To allocate more memory, program can increase size of brk()

↓

malloc() makes heavy use of it

stack — top of the virtual address space and goes downwards

• Linux has shared text segments for two programs running the same code

• Data and stack segments are never shared except after a fork

Memory mapped file — if a file is in disk

    ↳ we use the file descriptor in the open' sys call to
      open the file and then :—

        read
        write
        seek  —  move the pointer to and fro in the file
        close
      Operations are done

However, you can use mmap(), to map the physical file to the
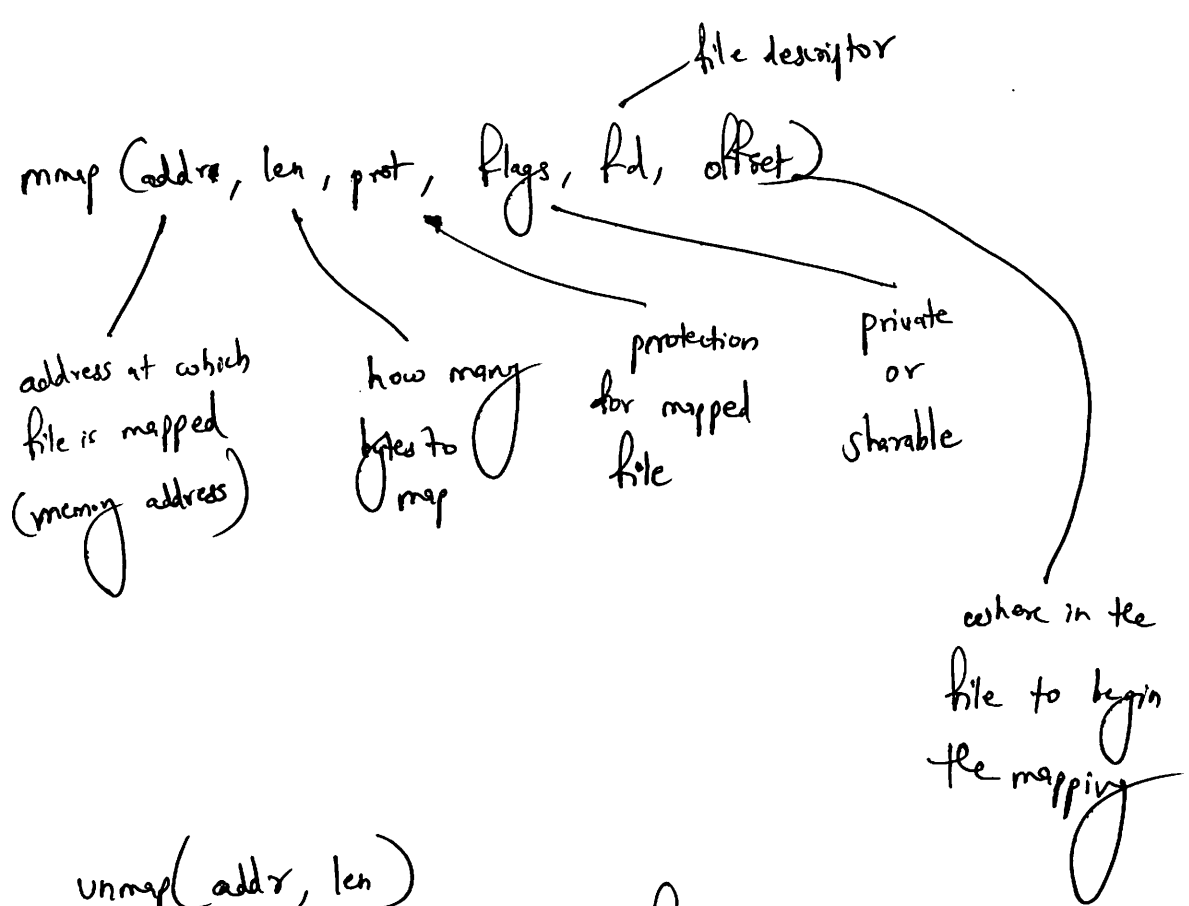address space ———— entire copy is not done initially

   ↳ On Demand copying mechanism using the page table
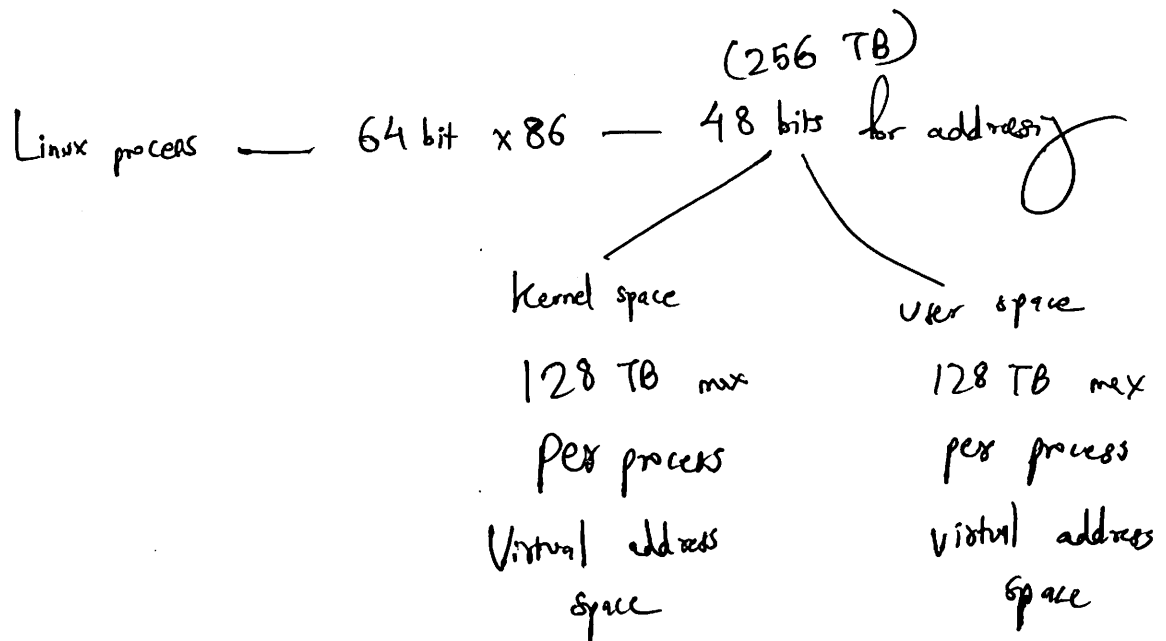     ↓

    To write back changes

      ↳ It will be upto the kernel to decide when to
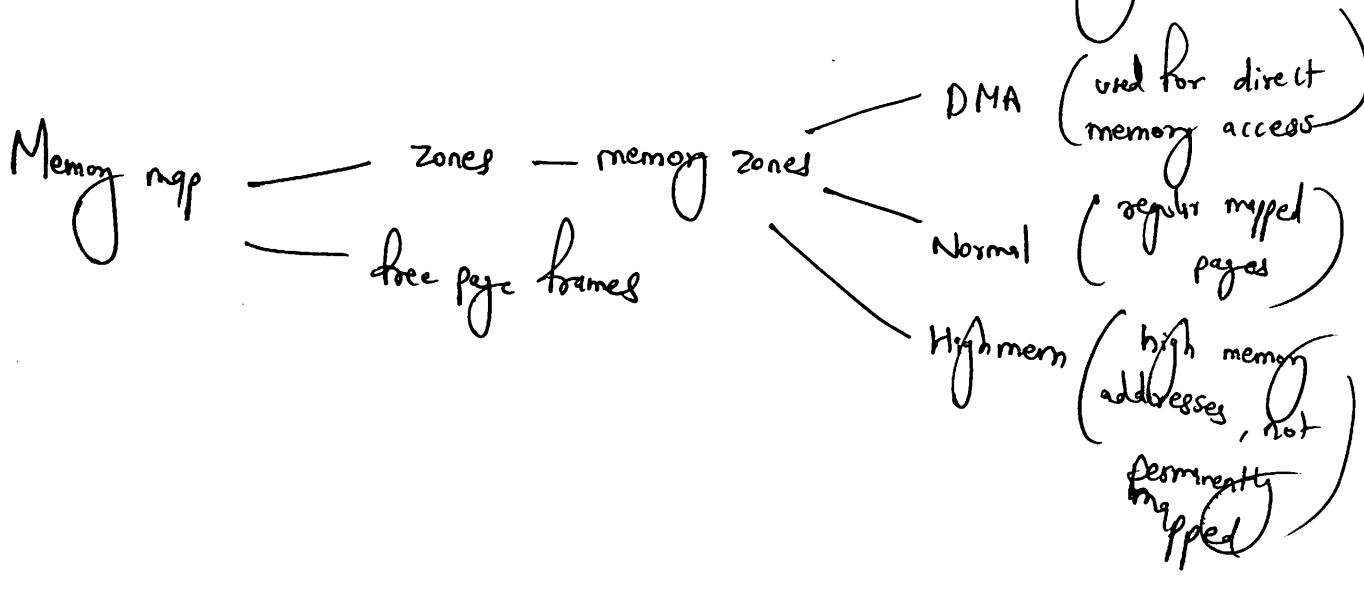        schedule a write back to disk

         ↳ the page table entry has a modified
          bit which tells if an entry is dirty or
        not

file descriptor

mmap (addrs, len, prot, flags, fd, offset)

address at which
file is mapped
(memory address)

how many
bytes to
map

protection
for mapped
file

private
or
sharable

where in the
file to begin
the mapping

unmap(addr, len)
└──→ removes a mapped file

(256 TB)

Linux process ── 64 bit x86 ── 48 bits for addressing

Kernel space
128 TB max
per process
Virtual address
space

User space
128 TB max
per process
virtual address
space

Main memory on linux — 3 parts ———— kernel

memory map

} pinned in memory (never paged out)

page frames

text, data, or stack page

page table page

free list

(saves space)
Linux supports kernel modules ——— pieces of the kernel that reside on disk till needed

loaded from disk and integrated to kernel ——— as kernel needs functionality

live in

"/lib/ modules"

⌐ not used

unloaded from memory ("autocleaning")

Memory map ——— zones — memory zones

free page frames

DMA (used for direct memory access)

Normal (regular mapped pages)

Highmem (high memory addresses, not permanently mapped)

page global directory   page upper directory   page middle directory   Page   offset

4-level page table Linux

Virtual address

New page frames & physical memory — page allocator — Buy buddy algorithm
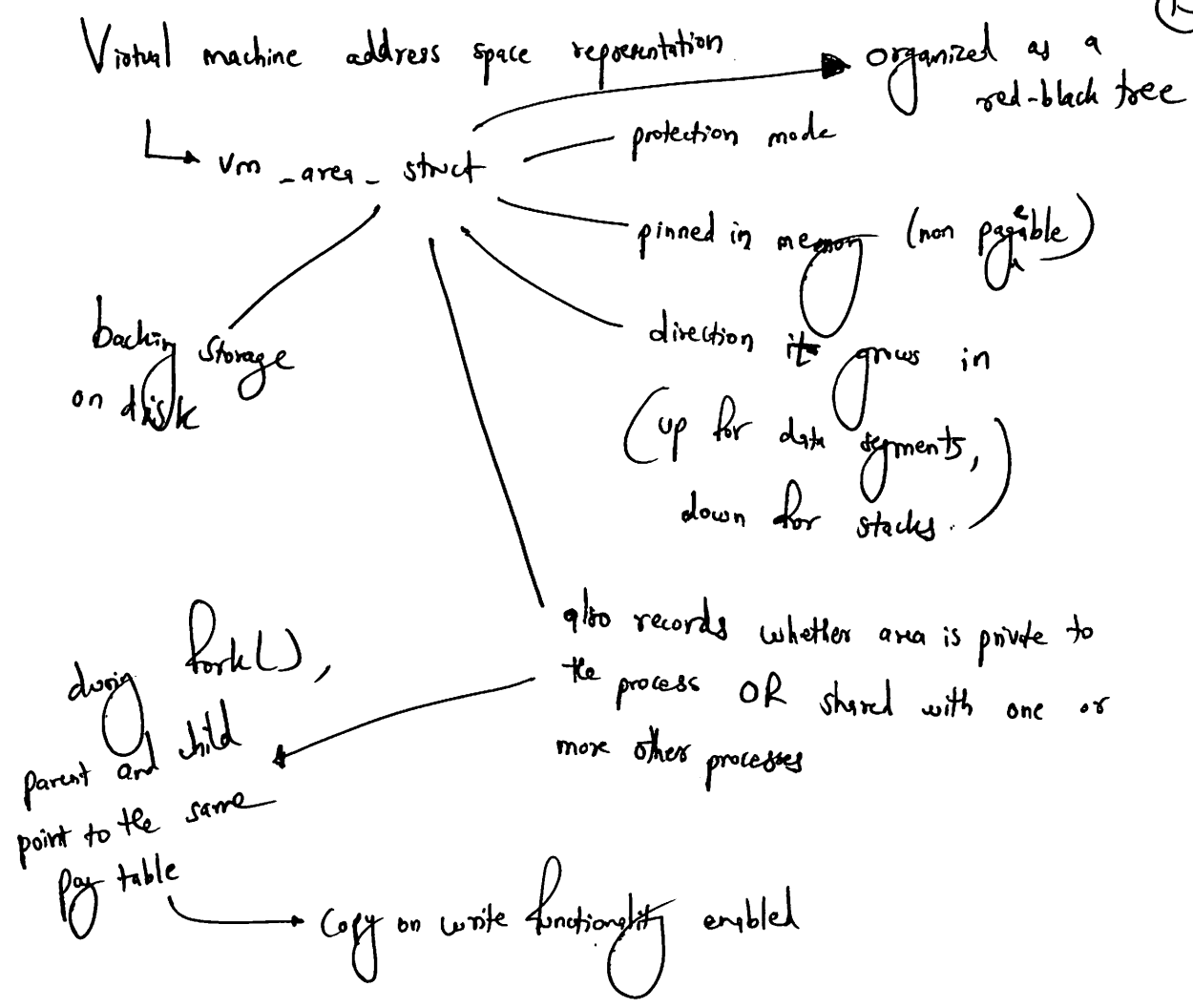
→ divide in power of 2 till you found right size, when freed, combine the buddy of same power ($2^3 + 2^3 \longrightarrow 2^4$)

→ this can cause internal fragmentation

to alleviate this → Slab allocator (carves slabs from these chunks and manages them separately)

Virtual machine address space representation ———→ organized as a red-black tree

└─ vm _area_ struct
    — protection mode
    — pinned in memory (non pagible)
    — direction it grows in

backing storage
on disk

(up for data segments,
down for stacks.)

also records whether area is private to the process OR shared with one or more other processes

during fork(),
parent and child
point to the same
pg table

Copy on write functionality enabled

mm _struct ——story info on——→ • all **VM** areas belonging to the address space

• different segments (text, data, stack)

In linux:-

Main Memory Management unit = Page
(granular)

←WHY?

process need not be entirely in memory in order to run

What is required for a process to run? (Memory Wise)

user structure                    page table

↓ if not there, Swapping needs to be done.

Paging ——————→ implemented ——partly—— kernel

——partly—— page daemon (process 2)

↓

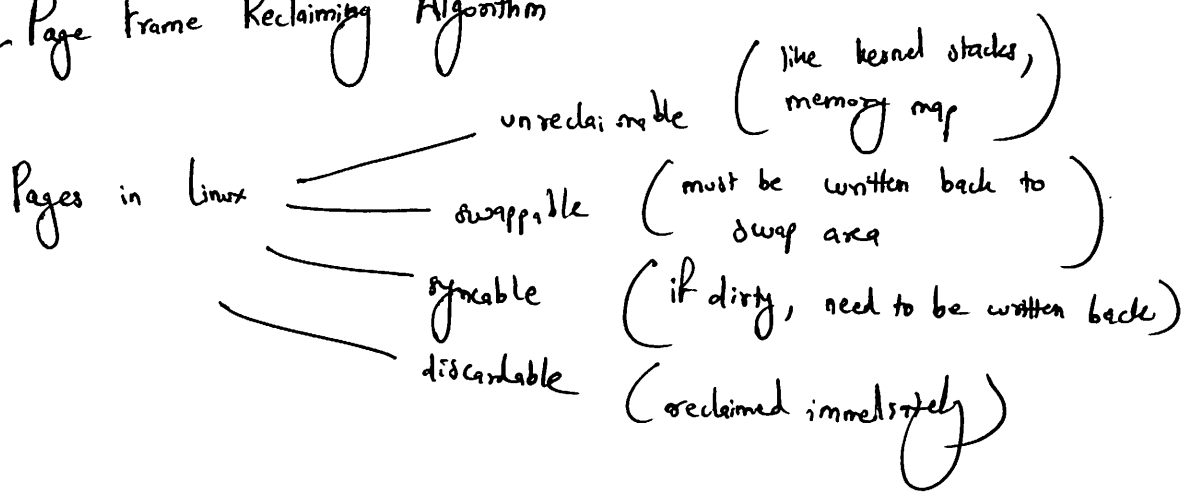if it sees no. of pages on free memory pages is low

↓

free up more pages

Swap area — fixed length paging files

↳ each has a priority

When a page has to be thrown out

↳ the highest priority paging partition or file (still having memory)

is selected

# Page Frame Reclaiming Algorithm

Pages in Linux

- unreclaimable ( like kernel stacks, memory map )
- swappable ( must be written back to swap area )
- syncable ( if dirty, need to be written back )
- discardable ( reclaimed immediately )

The algorithm tries to reclaim easy pages, then proceeds with harder ones.

Backing store: — part of hard disk that is used by paging or swapping system to store information not currently in main memory.

$$\text{Swappiness} = \frac{\text{pages with backing store}}{\text{Pages which need to be swapped out selected from PFRA}}$$

loop — scan through — Zone's active and inactive list

reclaim different list with different urgencies

active list
→ used by some process in the system

Inactive list
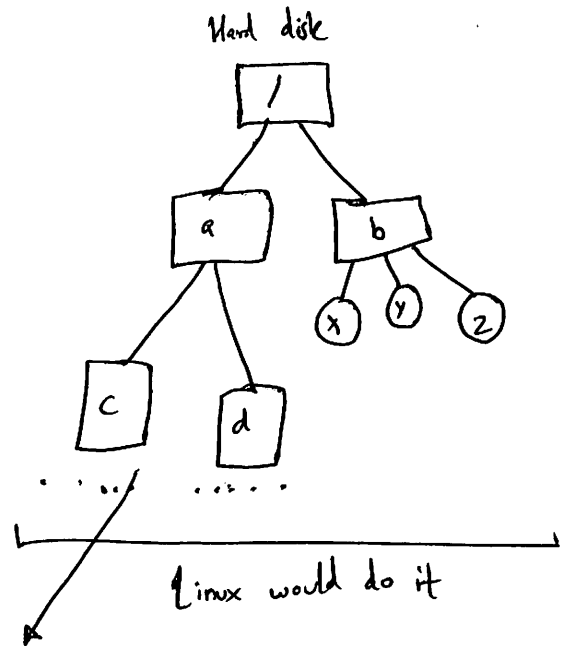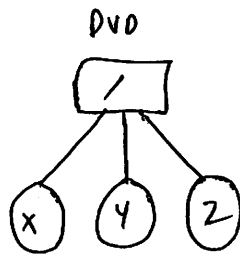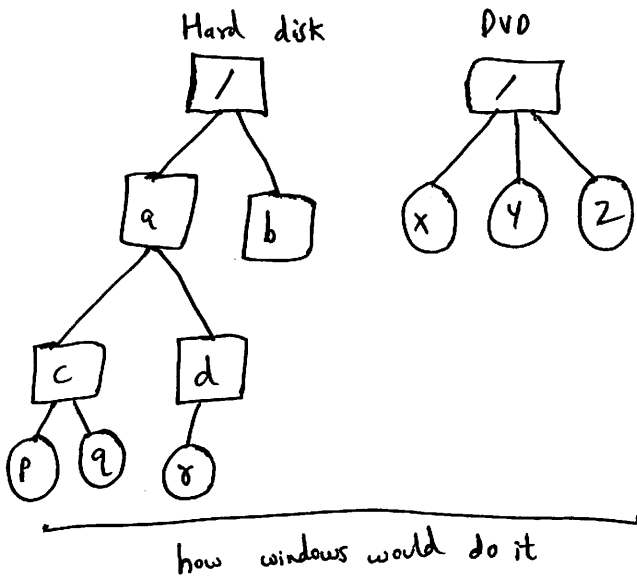→ kernel thinks might not be in use

File — . sequence of zero or more bytes

- File name — 255 characters ≤ ASCII (except NULL)

- Directories are stored as files and treated like files

↳ /bin — binary (executable) programs
Special

/dev — special files for I/O devices

/etc — miscellaneous system files

/lib — libraries

/usr — user directories

/dev/tty — reads from printer

/dev/lp — if write, writes to printer

Hard disk

```
        /
      /   \
     a     b
    / \
   c   d
  /|    \
 p q     r
```

Hard disk

```
        /
      /   \
     a     b
    /  \    / | \
   c    d  x  y  z
```

how windows would do it

Linux would do it

DVD

```
    /
  / | \
 x  y  z
```

no longer has to be aware of
"which file on which device"

Locks — shared — second attempt to place shared lock will work but exclusive lock won't

 — exclusive — all attempts to lock any part of that portion will fail until the lock has been released

creat ( file-name, mode) ——— protection — not only opens but also opens for writing, creates

→ returns fd
(lowest numbered not currently in use)

File System Calls

fd —— 0 (stdin)
   —— 1 (stdout)
   —— 2 (stderr)
} already open ——→ Shell arranges this

read ( fd, buffer, nbytes)
write ( fd, buffer, nbytes) ——————→ no. of bytes to transfer
   ↓                      ↓
file descriptor    where to put the data or
                   get the data from

lseek ( fd, offset, whence) ——— move the file pointer
   ↓                       └——→ relative to what?
returns the            beginning? Current position? or the end of
absolute position after    the file?
the file has been changed

file status information [
   stat ( file-name, &buf)    fstat ( fd, &buf) ——→ creation time
                                                  ——→ file size
            pointer to a structure where the information
            requested is to be put
                                              ——→ group file belongs to
                                    identity of file's owner
device file is on    file mode    no. of links
         I-node number
         (which file on device)
                        time of last access and modification

mkdir (path, mode) — create new directory

rmdir (path) — remove a directory

link (oldpath, newpath) — link to existing file

unlink (path)

chdir (path) — change the working directory  } ls

readdir (dir) — read one directory entry

no lseek inside

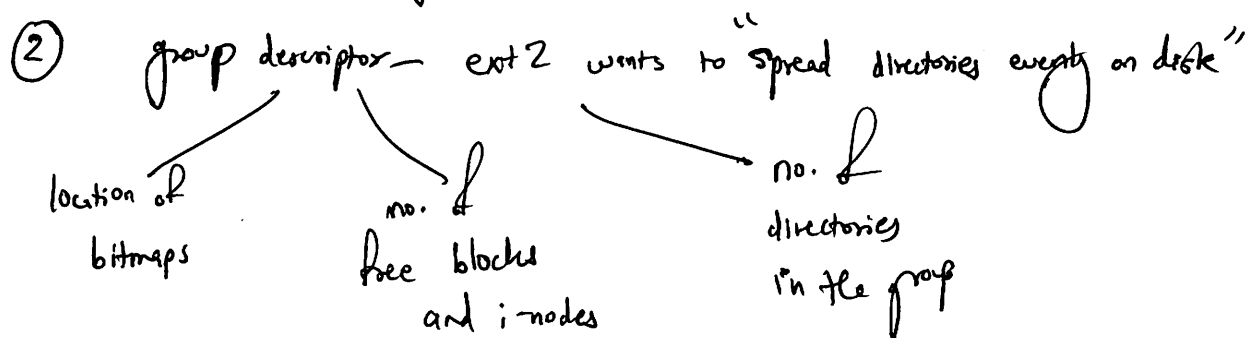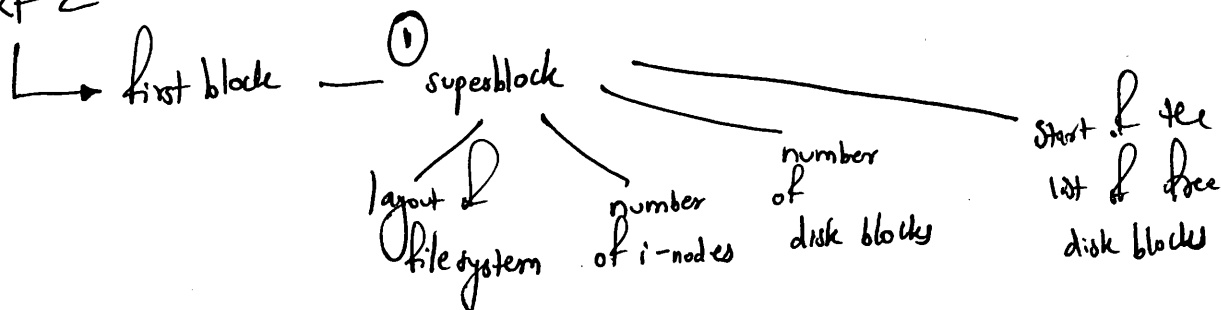## VFS (virtual file systems)

superblock — layout of file system

i-nodes (index nodes) — each describe exactly one file

dentry — directory entry ——→ directory entries are cached in

↓

dentry _ cache

file — normal files to be opened with sys calls like open(), close(), creat()

## ext2

↳ first block — ① superblock

- layout of file system
- number of i-nodes
- number of disk blocks
- start of free list of free disk blocks

② group descriptor — ext2 wants to "spread directories evenly on disk"

- location of bitmaps
- no. of free blocks and i-nodes
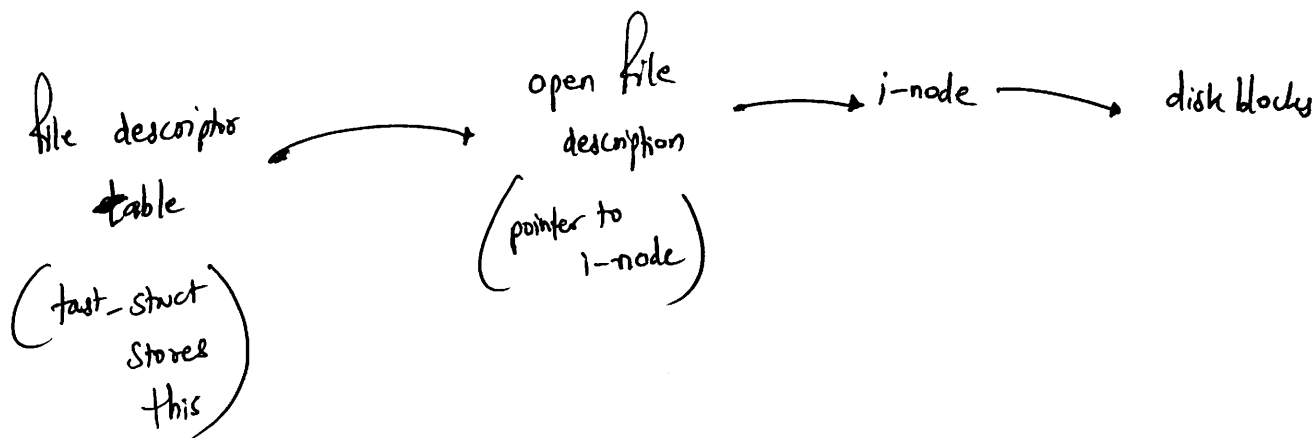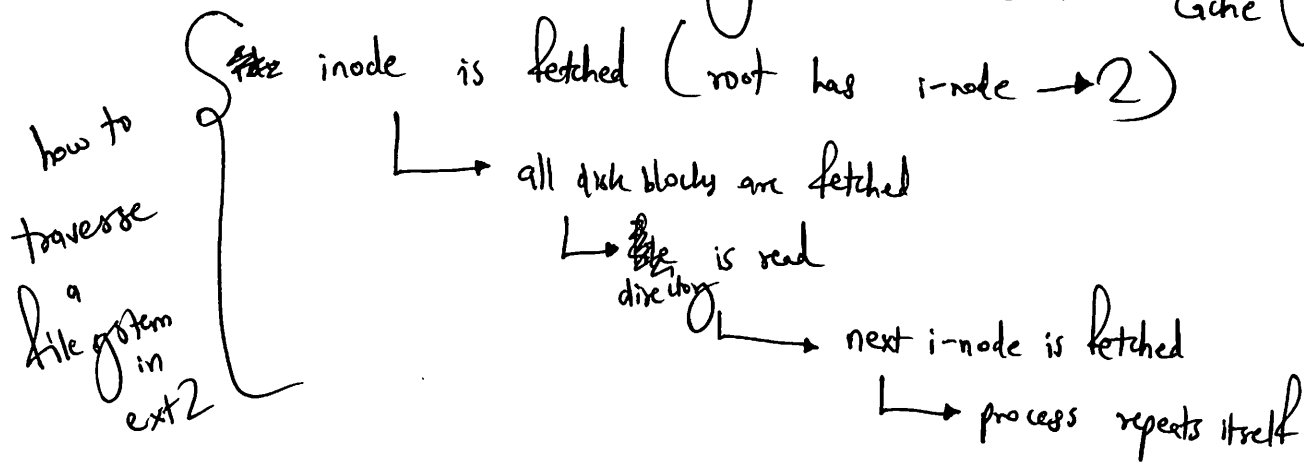- no. of directories in the group

③ i-node — • accounting info
• enough info to locate ALL disk blocks
that hold file's data

④ data blocks — files and directories are stored here

⑤ Bitmaps — used to make fast decisions
↳ where to allocate the data

Directories are searched linearly (also cached) → in dentry cache

how to traverse a file system in ext2 {
the inode is fetched (root has i-node → 2)
↳ all disk blocks are fetched
↳ the is read
directory ↳ next i-node is fetched
↳ process repeats itself

file descriptor table
(task_struct stores this)
→ open file description (pointer to i-node) → i-node → disk blocks

(20) load ⟍— CPU bounded load
        ⟍— out-of memory issues
        ⟍— I/O bound load

Out of memory — System has run out of available RAM
                              ↓
                    Started to go into swap
                              ↓
        ┌───────────── each process gets slowed down as disk gets used
        │
        └──→ TRAP: easy to rule it out as I/O bound load
                        ↓  (as disk is being used as RAM)
                   so ↓
                 Check RAM next

• Don't just look at the free column (see — free + cache)
  WHY? ┌─→ you need file ──────→ Kernel loads it to RAM
       │                                    ↓ after done
                  It tries to      ←─── If enough RAM available,
                  cache as many         kernel leaves it
                  files as it can in
                      RAM
so in some      so ↓ RAM
cases, we might
not need a   ←── this free RAM can get  —so→  free + cache
RAM disk         smaller