

# ID1206 Operating Systems

## 3 tasks of an operating system:

**Abstraction:** Create a layer of abstraction between the hardware and the application.

**Virtualization:** create image for each process, so that it thinks it has the whole CPU & memory while in reality it shares those with other processes.

**Resource management:** How do we share limited resources in a fair way.

**Kernel:** The program that interacts with the hardware.

## POSIX

Portable OS Interface.

Includes:

API: fork, exec, wait...

Process communication: pipes..

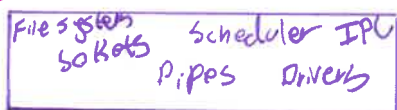
Threads: pthread\_create...

File system.

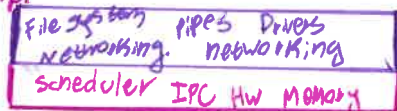
Networking: sockets.

## Monolithic VS Microkernel

A monolithic kernel is a very large system of all parts of the OS with an equal hierarchy of all processes.



A micro kernel is the minimal tools needed and everything else is built atop.



Microkernels are more robust and stable since errors can be handled by the kernel. But they are far slower since everything needs to go through the kernel.

## Standard C (ISO C18)

Memory allocation: malloc, free...

Signal handling: signal, raise, kill...

file operations: fopen, fclose, fread...

## Commandline Interpreter

Shell: text based

Scripting languages:..

## C process

a C process has the following:

- a program
- a stack
- a heap
- program counter
- open file descrip

## The Stack

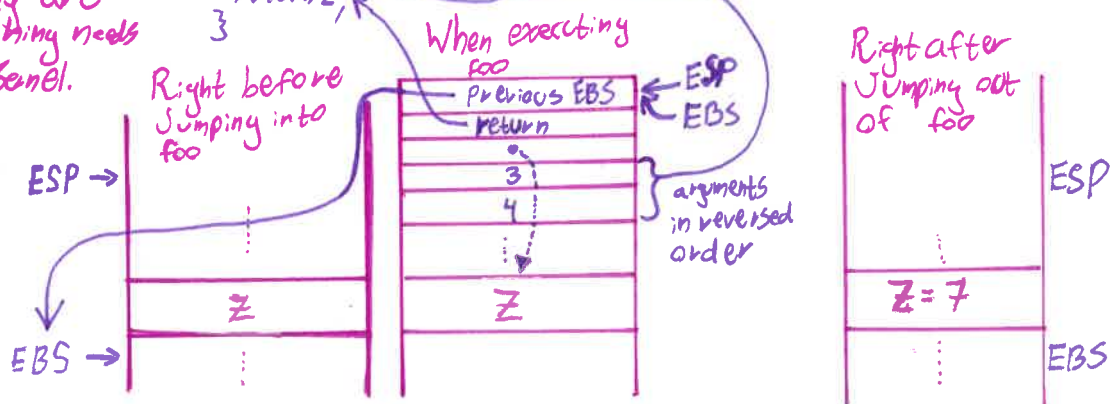
A stack frame contains arguments, local variables, info to be able to exit the frame

Pointers in CPU:

- instruction pointer (EIP) next instruction to execute
- Stack pointer (ESP) top of the stack
- base pointer (EBS) current stack frame.

```
int bar() {
    int z;
    z = foo(3,4);
    return z;
}
```

```
int foo(x,y) {
    return x+y;
}
```



## Returning Datastructures

```
int* foo(int x) {
    int a[5] = {1, 2, 3, 4, 5};
    return a;
}
```

Calling this function will return a segmentation fault. This is because `a` is allocated within `foo`'s stack frame and thus its location is cleared after `foo` is exited. Instead the heap should be used.

## Program to Process

- 1) Find program in persistent storage.
- 2) Allocate to hold the code, statics, heap & stack.
- 3) Process context to hold necessary values.



Context

## Interrupt descriptor table.

Used by kernel to execute kernel level instruction invoked by user level processes.



Protected area of the memory.

These calls are expensive since everything needs to be saved temporarily.

## Fork

Is used to clone a process with a deep copy of the memory (they do not share variables)

```
int pid = fork();
```

It returns 0 once for the child and a new process id for the parent process.

## Heap API

`void *malloc (size_t size):` Allocate size bytes of data on the heap, returns a pointer to the structure.

`free (void *pointer):` Deallocates the structure at ptr

`void *calloc (size_t, nmem, size_t size)` Same as `malloc` but fills the area with `nmem`.

`void *realloc (void *ptr, size_t size)` Changes size of existing `malloc`.

## Indirect Execution

The OS decodes the instructions and runs them like a virtual machine. Java, Python, Erlang.

## Direct Execution

The code is directly executed on the CPU. OS gives up control. Dangerous!!!

## Limited Direct Execution

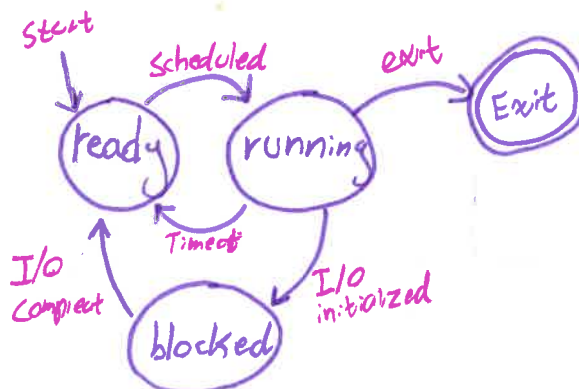
With the use of hardware support code is directly executed with some limitations.

- Will not be able to execute all types of instructions.
- Can only access its own memory
- has limited CPU time.

## Interrupts

These tell the kernel to do something based on a trigger, eg divide by 0, keystrokes, exceptions. They can also be asynchronous, raised by hardware (mouse) or asynchronous, raised by CPU.

## Process State



Many processes can be blocked and ready at the same time but running is limited by the CPU.

## Fork Commands

`exit();` terminates process  
`wait();` wait for other process to terminate.  
`exec();` load different program and start exec.  
`getpid();` gets pid.

`kill();` send signal to process  
`raise();` raise exception  
`sigaction();` sets signal handler.

## Performance Metrics

$$T_{\text{turn around}} = T_{\text{Completion}} - T_{\text{arrival}}$$

$$T_{\text{response}} = T_{\text{First Scheduled}} - T_{\text{arrival}}$$

### MLFQ

round robin with priority  
Priority changes dynamically

Rules:

- 1) if priority(A) > priority(B)  
Pick A
- 2) if priority(A) = priority(B)  
round robin
- 3) New Jobs have highest priority.
- 4) if Job does NOT use an I/O operation  
priority decreases by 1
- 5) After some time of low priority the process is set to the highest priority again.

### Stride Scheduling

Jobs are given a value  
inverse to their priority.  
Every time the job  
is executed this value  
(called stride) is added  
to a value associated  
with the process.  
The process with  
the lowest value  
goes first.

### Real Life

uses stride scheduler  
and keep jobs on the  
same core to  
prevent cache misses.

## Scheduling Strategies

Shortest Job First (SJF)

Shortest time to completion (STCF) \* optimal turnaround

Round Robin (RR)

Multi-level Feedback Queue (MLFQ)

Lottery

Stand in line

Stride Scheduling used in linux

## Realtime Scheduling

Adds a new requirement that jobs should  
be completed before a deadline.

Hard: All deadlines MUST be met, missing a  
deadline is a failure.

Soft: Deadlines can be missed but the application  
must be notified

Tasks are described by a triple (e, d, p)

e: worst case exec time

d: deadline of task

p: how often task should be scheduled.

$d < p$ : Constrained

$d = p$ : default

$d > p$ : Several outstanding

## Rate Monotonic Scheduling (RMS)

assumes default ( $d = p$ )

Rule: Schedule task with shortest period.

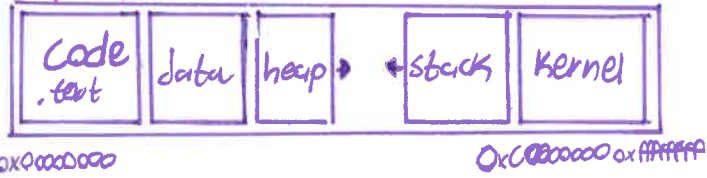
works always when utilization ( $\sum \frac{d}{p}$ ) is less than 69%  
but in general when utilization is less than  $n(2^n - 1)$  n = number of jobs

## Earliest Deadline First (EDF)

Based dynamically on deadline

works always if utilization is less than 100%

# The Process



32 bit Linux

However the process does not actually operate on address  $0x00000000-0xFFFFFFFF$  in memory. instead a so called memory management Unit (MMU) translates between physical and virtual memory addresses.

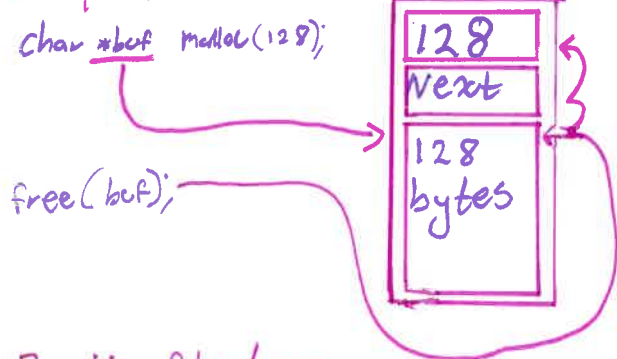
## Brk() & Sbrk()

Change the end of a process heap segment. `brk()` sets the value to the arg and `sbrk()` increases current size by the argument. `sbrk` returns the new address.

these are syscalls while `malloc` and `free` are library routines.

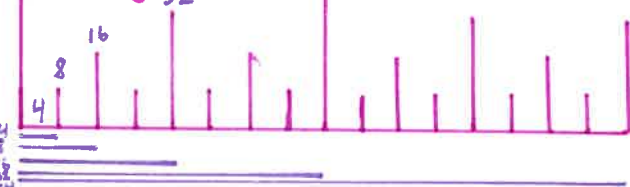
## Malloc & Free

In order to know how big an area is this is stored two addresses above the pointer that `malloc` returns.



## Buddy Strategy

Used in linux for allocation of physical memory.

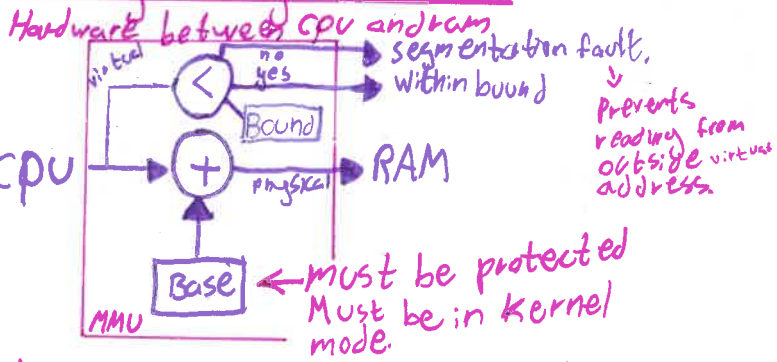


Each memory segment can recursively be split into buddies, each of  $\frac{1}{2}$  the size. they are split until a region fits a node as closely as possible

very efficient  
coalescing efficient ( $\log(n)$ )  
handles fragmentation very well.

internal fragmentation. if 65 bytes is needed 128 is given.

# Memory Management Unit



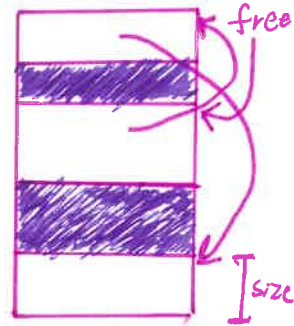
the con is that some memory like libraries must be copied to each process and shared memory is not usable. Problems with fragmentation. Processes are always assigned 4GB even if only a few bytes are needed.

Transparent: Processes are unaware of virtual addresses.  
Protection: Processes cannot interfere with each other.  
Efficiency: Should be as close to real execution as possible.

## Free list

used to find free space in memory

```
typedef struct __node_t {
    int size;
    struct __node_t *next
}
```



When blocks are freed they are added to the start of the list.

When consecutive blocks are freed they must merge (coalesce)

## Free list strategies

If i need to malloc 40 bytes how do i know how big or which block i malloc?

Best fit: smallest possible block

Worst fit: biggest block

First fit: first possible block

## Segregated lists

have separate list of blocks of predefined sizes. (8, 16, 32, 64...) very fast at finding blocks of approximate size.

## mmap

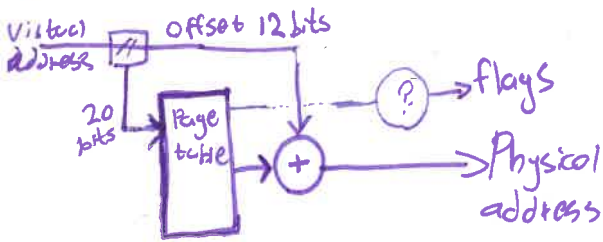
Default use in OSX instead of `sbrk` and `brk`

```
void *mmap(
    void *addr,
    size_t len,
    int prot,
    int flags,
    int fd,
    off_t offset
);
```

Creates virtual mapping  
`addr=NULL` => dont care  
`prot` = protection (allows shared memory)  
Flags, `fd`, offset allows mapping to a file in memory.

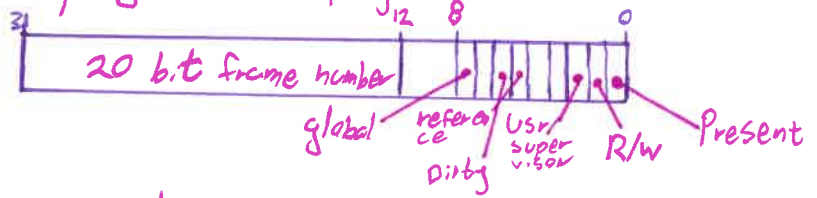
# Paging MMU

Memory is split into pages of equal size. Usually pages are 12 bits large of address space.



# Page Table

Each entry is a map between a virtual and physical address. Or rather a page in the physical address table.



# Physical address extension

In 1995 the x86 extended the frame numbers to 24 bits. Thus 64 GB of physical address could be addressed. Virtual processes still don't know about the extension and still think it's limited to 4 GB. x86-64 has 48 bit virtual address 0xf...fff...

overflow can never happen since a page is exactly  $2^{12}$  bits large. Segmentation and paging is often combined.

## Lookup Procedure:

```
Movl 0x11111222, %eax
      virtual address
```

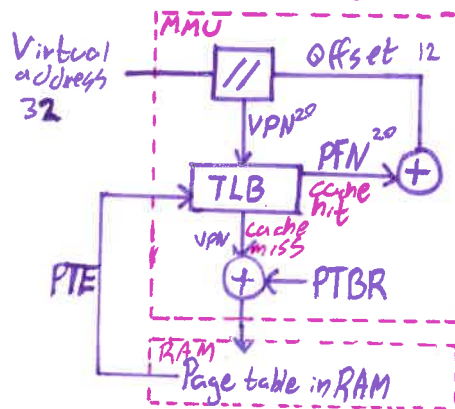
- 1) Page table base register PTBR
- 2) virtual page number VPN 0x1111 (20 bits)
- 3) read from page table  $PTBR + (VPN \ll 3)$
- 4) extract Page frame number from table
- 5) add offset to page frame number.
- 6) read memory at this location  $(PFN \ll 12) + 0x222$

This is too slow since each lookup requires two memory accesses.

## Speedup using hardware

A copy of each lookup is stored in the CPU cache so sequential lookups are giga fast.

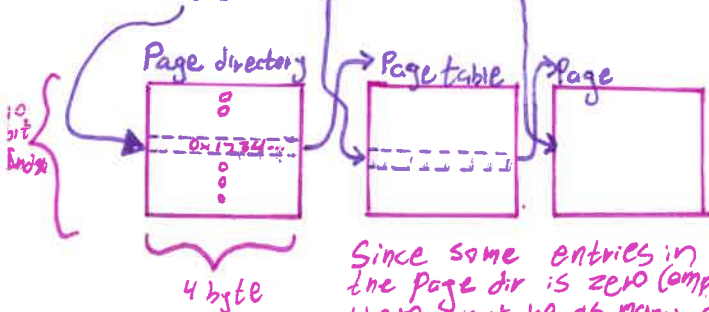
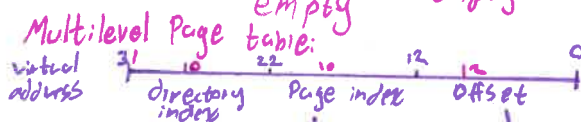
This is called TLB (translation lookaside buffer) Non-global entries need to be flushed between users when switching active processes.



The drawback is that each entry is 4 bytes

Each page table is 4 Mibyte. Each process has its own page table. Thus for 100 processes 400 Mibytes of tables is needed (too much).

Solution: Most pages are actually empty.



Since some entries in the page dir is zero (empty) there won't be as many page tables needed assuming there are gaps.

## Inverted Page lookup:

If we have 8 Gbyte RAM then it'll fit 2M frames of 4Kibyte then a table with 2M entries can map a process to its frames

index	page num	Process id
0		1
2		3
...	...	...
8		4
...	...	...
...	...	...

When you find your entry the index of the table is the actual page number in physical mem.

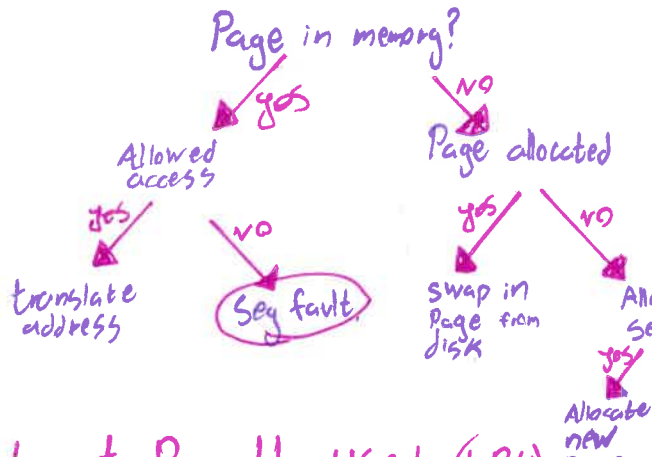
# Swapping

Benefits:

- 1) Give illusion of private memory
- 2) Give illusion of much larger address space.

When the main memory is filled up the unused processes have their pages moved out to external memory i.e. harddrive. later when Paged the OS will fetch the Page from the harddrive.

## Page lookup Process



Least Recently Used (LRU)  
used to decide what page to swap when RAM is full. The idea is that a page that was not used recently is less likely to be used again. It is expensive to keep track of which pages have been used when there are 1000+ pages.

## FIFO

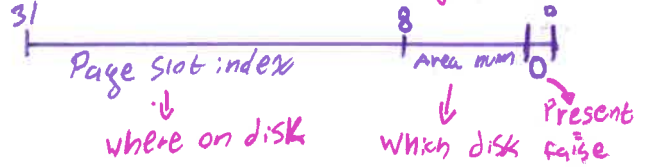
First in first out, approximation of LRU but the que is only updated on misses. much cheaper algorithm since the que is only updated when the processor is already taking time to find page in persistent memory.

## Clock

Uses hardware and is another approximation of LRU. A buffer of all pages is created. When a page is accessed and it is in memory its usage bit is set to 1. If there is a miss a pointer goes through the buffer if the pointer points at a 1 it is set to 0 and pointer is incremented, if it points to a zero, the current page is swapped and pointer is incremented.

# Page Table Entry Pt. 2

As shown on the previous page there is a flag called "present" that indicates if the MMU has an entry for that page. But if it is 0 there will be 31 remaining bits that can be used to find the page on the HDD.



## Cost of page faults

TLB hit address in cache: 1ns  
mem: 10ns

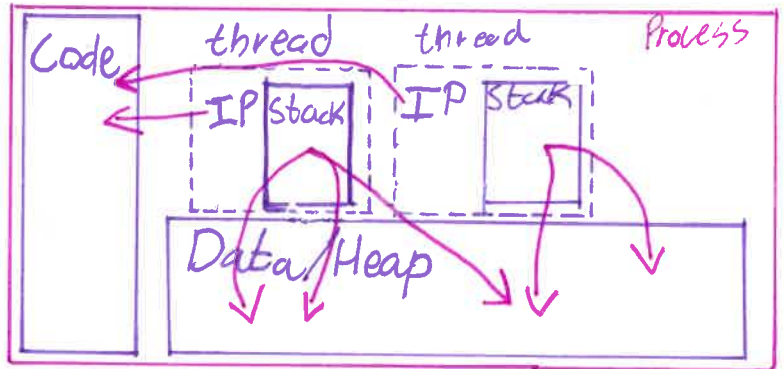
TLB miss page in mem: 100ns  
Disk: 10ms = 10<sup>7</sup>ns

## Concurrency & Parallelism

The illusion of things happening at the same time is concurrency

Parallelism is actually doing things at the same time.

## The Thread



Threads have a shared memory including heap & stack. GCC however may optimize the code, so use the label volatile on shared variables.

Threads can be scheduled by either the process or by the OS

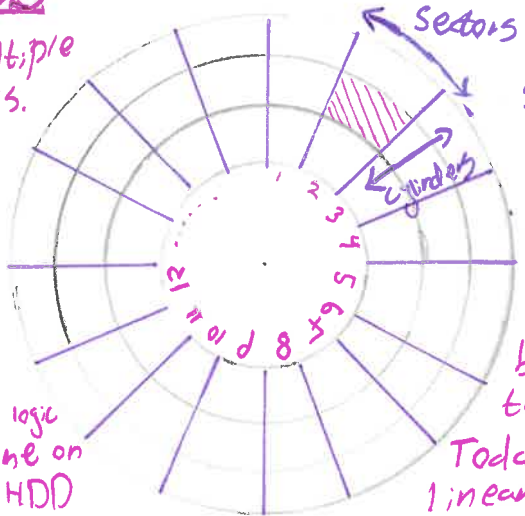
in GNU Linux everything works in kernel space.

```
int pthread_create(...); // create thread
-- thread int local = 4; // create local variable to a thread.
```



# HDD

Multiple disks.  
Each with two sides.



Smallest unit: 4K  
In the past sectors were addressed by disk, cylinders, sector but this limited storage to 500MB.  
Today it is all done with linear addresses.



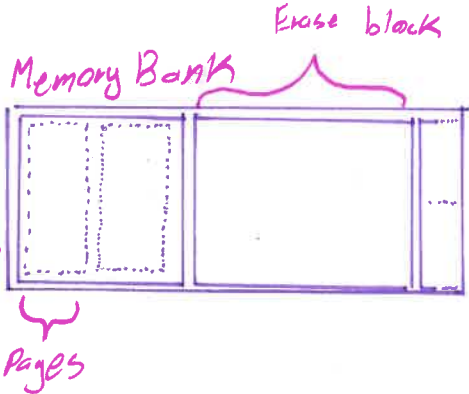
2TB  
3.5"  
7200 RPM  
SATA III  
Cache 64MB  
read 156MB/s  
0.44 K\$/GB

Most logic is done on the HDD

# SSD



500 GiB  
2.5"  
RA 10ms  
read 560MB/s  
SATA III



only a whole erase block can be erased. Thus changes cannot be made to individual pages. Thus to edit the hardware moves a page instead of editing.  
SSD are cool and can be hot on every bus.

2.4K\$/GB

# BUS limitation

- SATA III 6 Gb/s
- SAS-3 12 Gb/s
- USB 3.1 10 Gb/s
- PCI Ex. 30x6 128 Gb/s

# RAID

Redundant array of independent disks. Multiple disks seen as one, also has error corrections.

## Raid 0:

stripe file across multiple disks.

## Raid 1:

keep a mirror copy of each file.

## Raid 2-6:

spread a file + parity info across several drives.

# Files

- Persistent, shared, path.

A file is:

- sequence of bytes.
- metadata
  - ↳ size & type etc.
  - ↳ owner & permissions.
  - ↳ author
  - ↳ created, changed
  - ↳ (icons, fonts)

functions:

- create & delete
- finding
- read write
- control authorization.

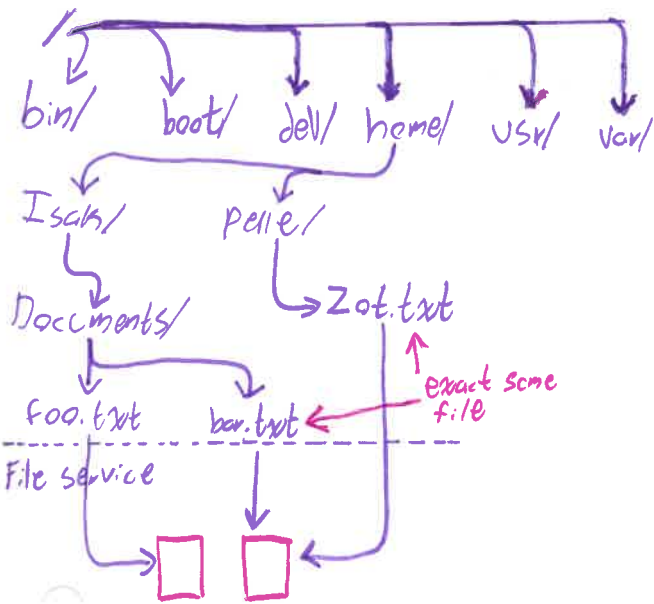
Directory } mapping

File module }  
Access control } file structure  
File operations }

Block operation }  
Device operation. } Drivers



# Tree File Structure



# OS Functions

**creat** - create file

**unlink** - remove link, when last link is gone remove file

**link** - hard link

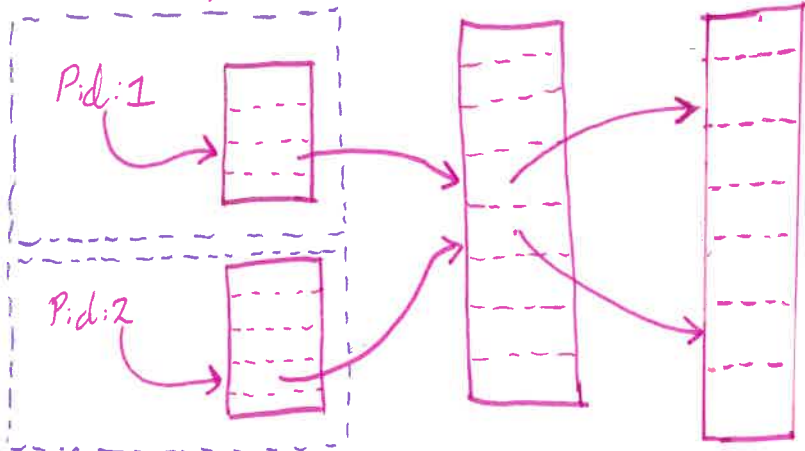
**symlink** - soft link

**stat** - read metadata

The kernel keeps a record of each process accessing a file. By default a process has 3 open files, standard input, standard output, standard error.

# File Tables

One table per process, copied when process is forked.



**File Descriptor table.**

list of files used by process.

**Open file Table**

keeps track of open(); calls, has a counter for when 2+ processes access a file

**Inode Table**

One Inode entry for each file.

# The Inode

Includes:

- Mode (access rights)
- #links
- User Id (owner)
- group Id
- size
- blocks
- identifiers
- generation.
- access time
- modify time
- change time.

Total 15 total Pointers.

first 13 pointers are normal block pointers

Pointer 14 is a pointer to a block of pointers.

Pointer 15 is a pointer to a block of blocks of pointers.

# The File system

The first few blocks of a HDD is metadata that tell the OS about the HDD.



• Describes file system

• How many Inodes and where

• How many data blocks and where

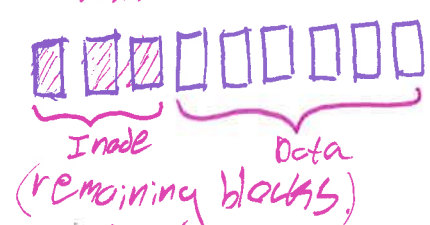
• where bitmaps



• bitmap for which inode is taken



• bitmap for which datablocks are needed.



Directories are Inodes. with a single block.

# File Operators

**Read** - reads bytes into buffer

**Write** - writes bytes from buffer

**Seek** - reposition write head.

# File System errors

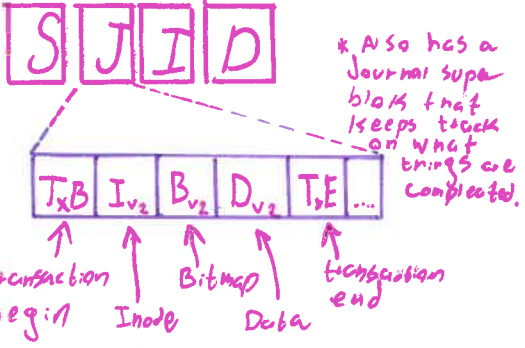
When a file is changed 3 changes must be made.

- update bitmap
- update Inode
- update Blocks

Crashing before doing 3/3 tasks is very bad.

## Journal

Keep track of what we are going to do. Add a new Journal Block.



Either All are complete or none are complete. These can be queued before many are done at once.

## Virtualization

The OS is run with a so called Hypervisor between the OS and the Hardware. It mimics the interrupt table and acts as a man in the middle. This is done in order to emulate different hardware eg. Arm on a x86 machine, or to have different versions on an OS before everything is migrated/updated.

## Approaches:

File system check: flag that tells file system that it did not shut down properly prompting a complete file check.

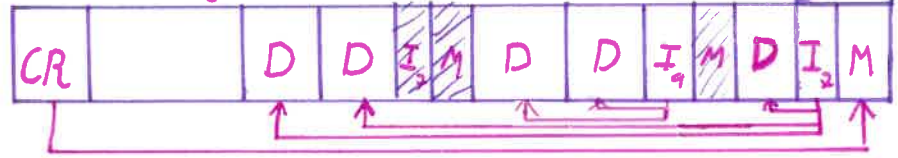
Journal: write down what you want to do before you do it.

log: log all changes.

Copy on write: create a perfect copy and last thing you do is flip a pointer.

## Log Structures

Write sequentially since reads are cached in memory write speed is the most important



Everything is written left to right. Inodes point to data blocks before them. Map point to which Inodes are most up to date. Check region points at the latest Map. This way the system can fail before CR is updated and it will simply be reverted to a previous state. When full the first blocks in the list will be copied to the end, clustering them and freeing consecutive space.

## Hypervisor Guest OS Application

