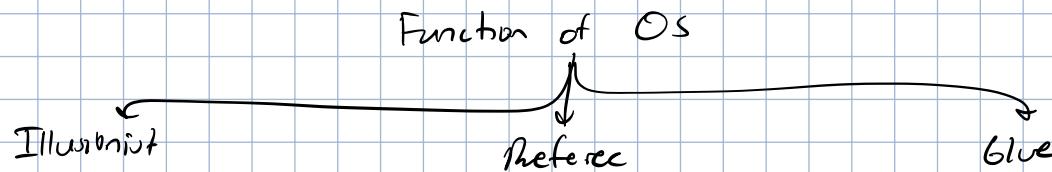


## INTRODUCTION

What is an OS?

- What is shipped w/ device

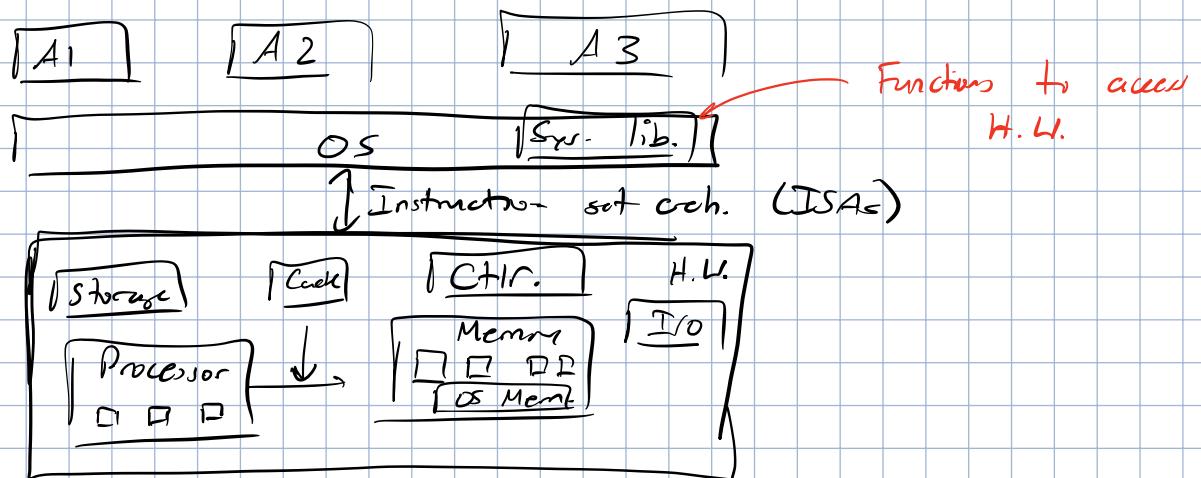
What does OS do?



Illusionist:

Clean abst. of H.W.

↳ No worries about resources, mem.

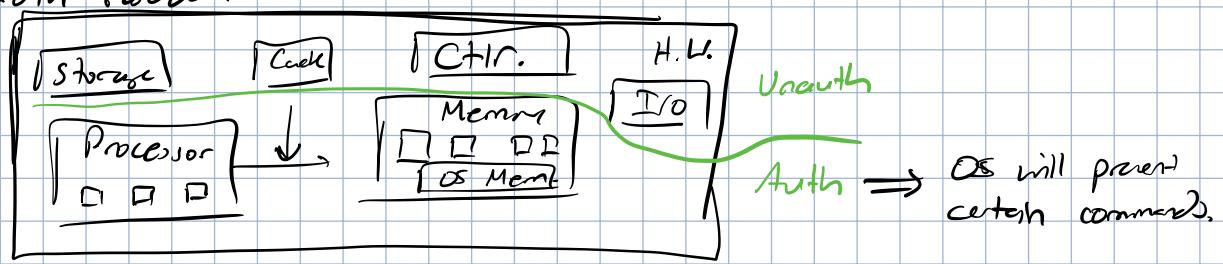


Examples of abstr.:

- Processor → threads
- Memory → address spaces / virtual mem.
- Storage → file system
- Networks → sockets

Reference:

1. Schedule resource allocations:
2. Apps cannot access unauthorized resources.
  - Unauth resources...



- o Exit If accessed unauth. memory  $\Rightarrow$  seg. fault.
- o Resources allocated to 1 process cannot be accessed by another process

Glue:

Common interface for H.W.  $\Rightarrow$  pretty easy to "access"

OS is not more intrusive!

$\hookrightarrow$  How?

1. Pair up L1 H.W. (not necessary, but easier & more lightweight)
2. Algos & D.S.

Kernel: runs all the time on computer

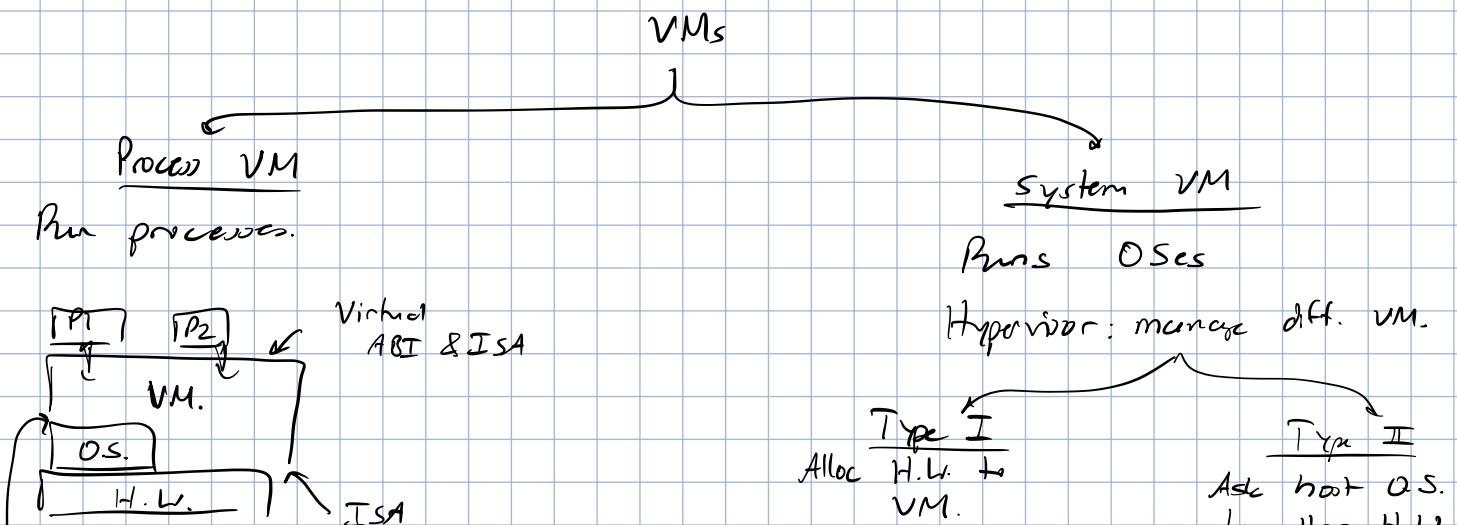
$\hookrightarrow$  Anything else is either an app or system app

## Virtual Machines

Defn: software that emulates H.W.

$\hookrightarrow$  Doesn't depend on H.W.

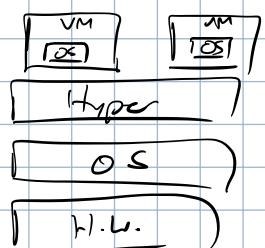
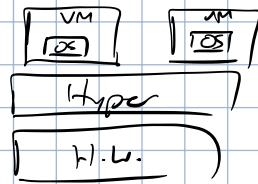
ILLUSION  $\Rightarrow$  isolated H.W.



ABI: app. binary interface.

Ex: Java

$\hookrightarrow$  Run on any computer.

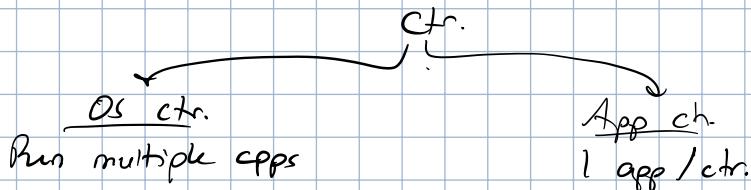
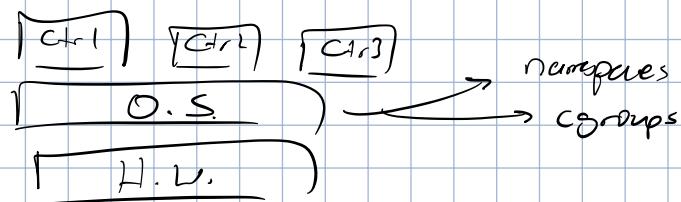


Container:

Defn: emulate O.S.

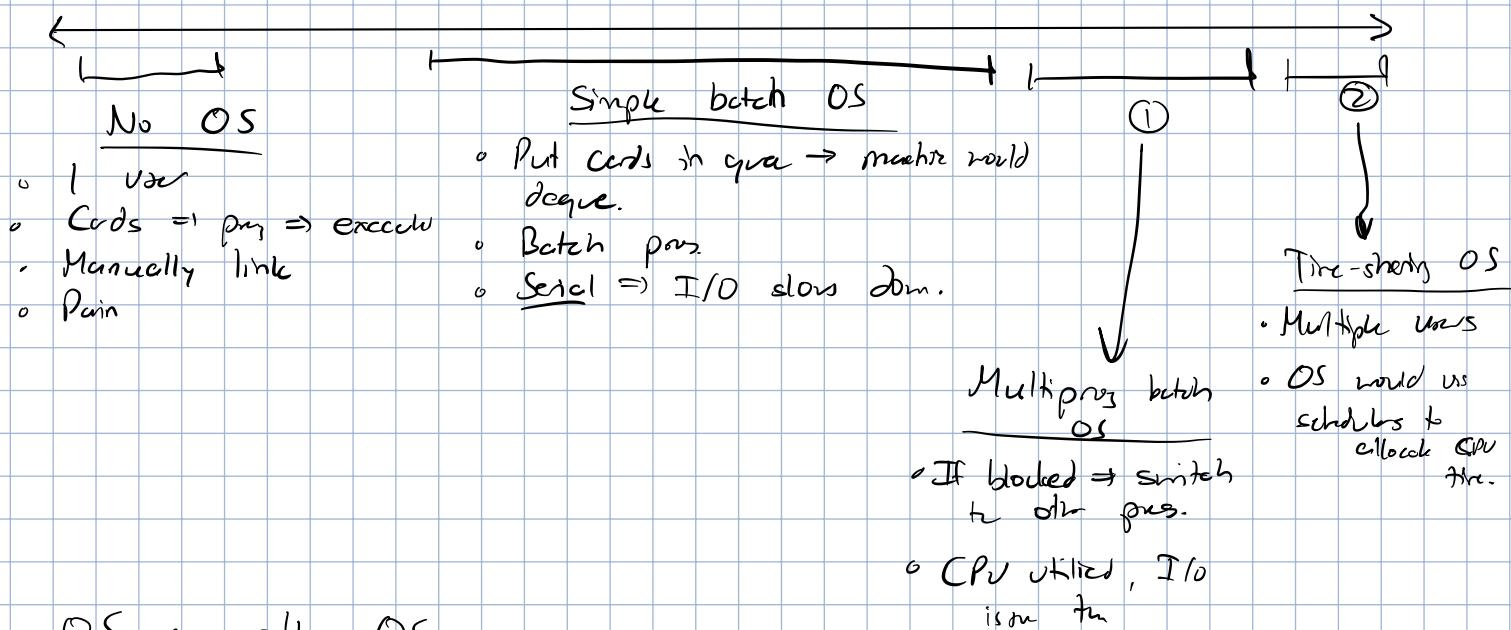
ILLUSION  $\Rightarrow$  isolated O.S.

But only emulating OS  $\Rightarrow$  perf. benefit.



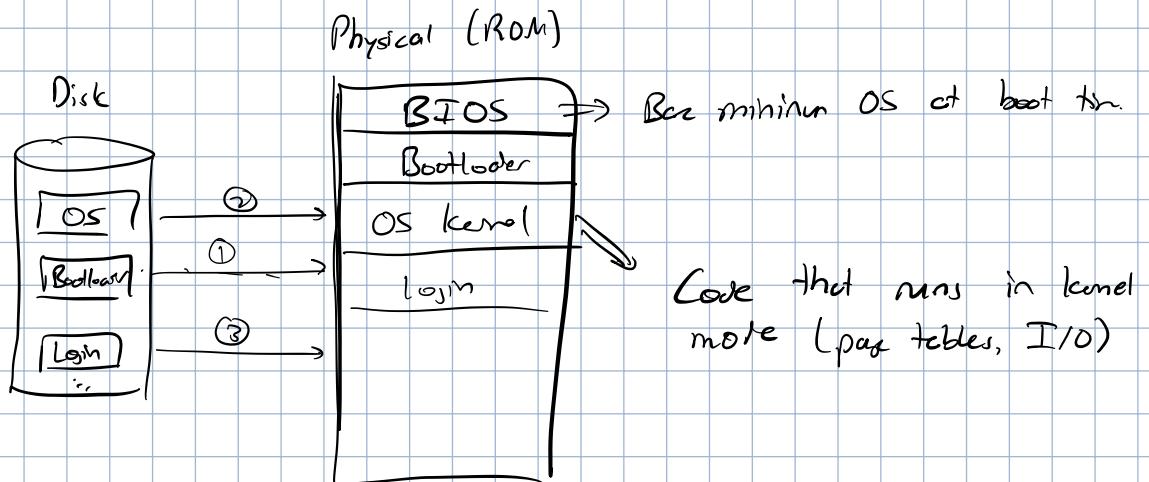
## OS CONCEPTS

### Brief History



OS  $\leftarrow$  other OS

### Bootstrapping



① BIOS loads bootloader

② Bootloader  $\Rightarrow$  local OS

② OS kernel → load / load / start

Q: Why do OS load / load?

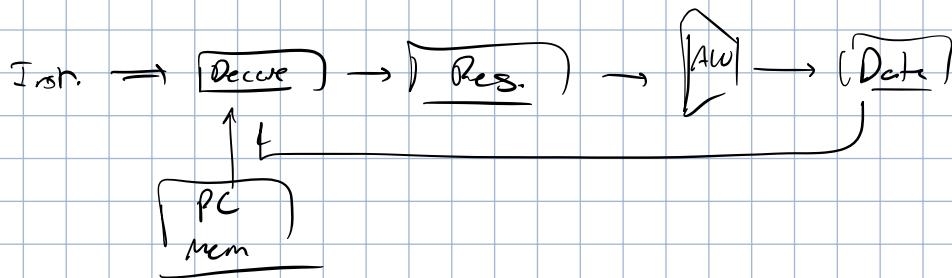
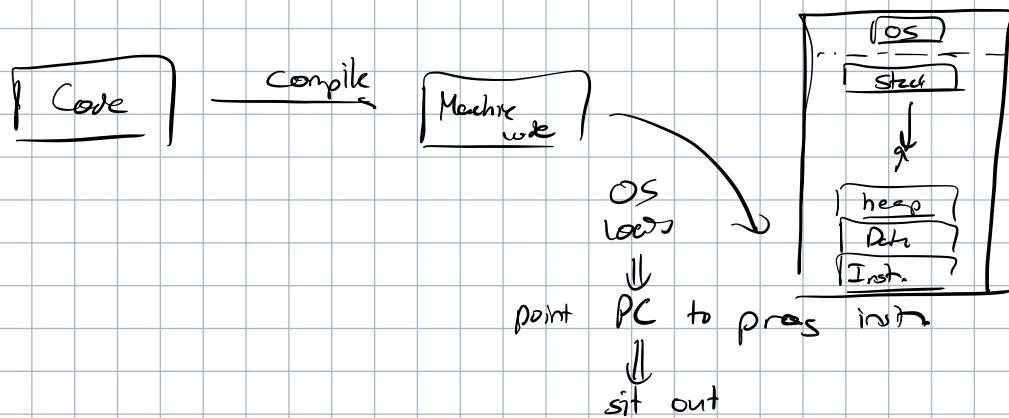
a) ROM is expensive ⇒ kernel is big

b) Kernel is dynamic (update) (had to update if in ROM)

Q: Why boot loader before OS?

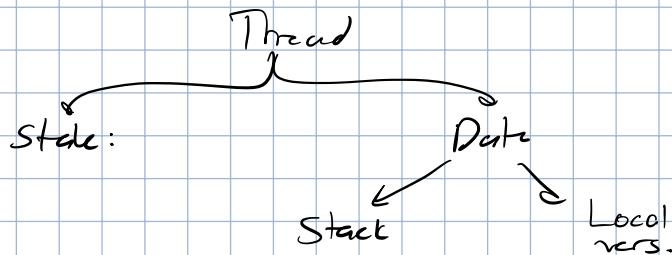
a) Choose on OS

Execution of Programs.



Thread

Defn: Seq. of commands on CPU

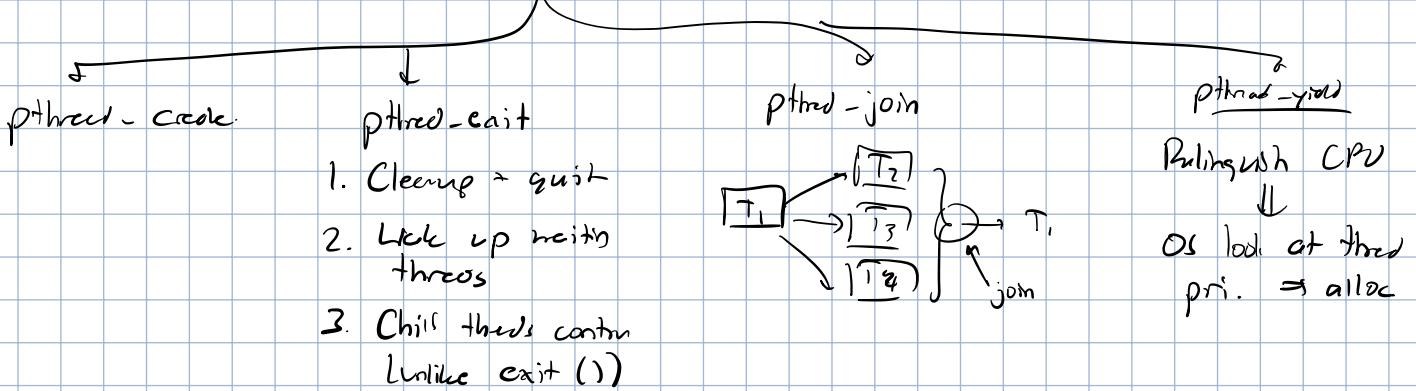


Multithreaded: 1 process using multiple threads?

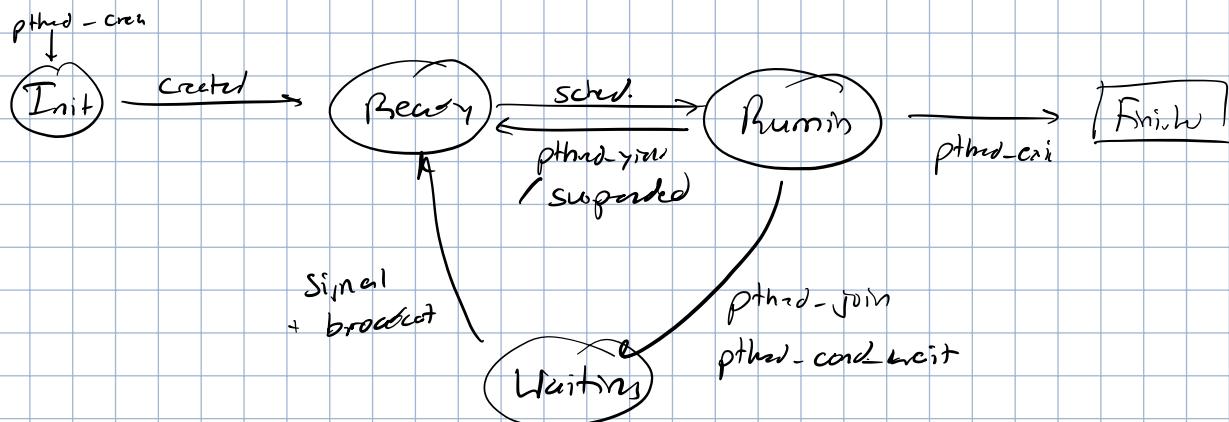
↳ CPU can only 1 thread / time. ⇒ can do max w/ max threads. } thread creates max threads  
↳ can use concurrency to illusion of parallelism } parent + child

pthread: linux lib.

Connex



Thread state:



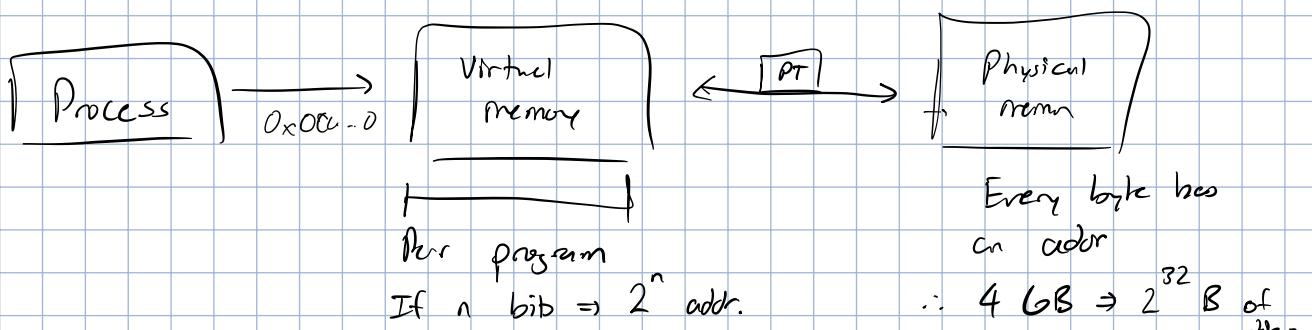
Thread control block (TCB):

Keeps metadata

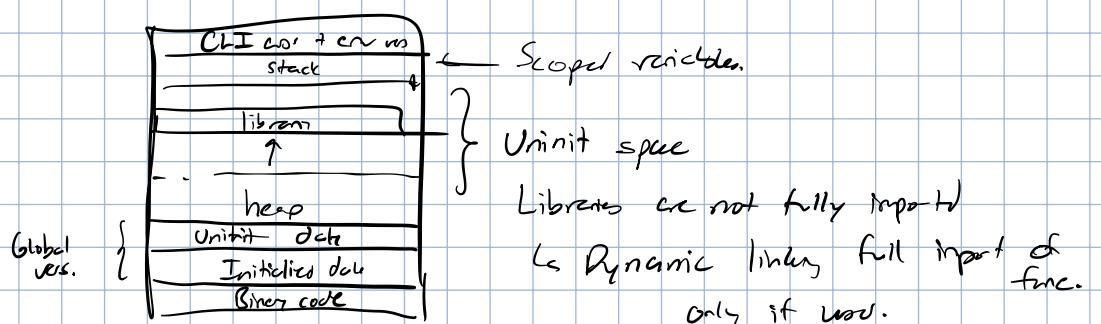
- o TID
- o PC / memory ph + stack
- o Reg values.

Kernel b/c  
concurrency  $\Rightarrow$  switch b/w thds, how do  
+ context switch we restart?

Address Space



Memory organization in virt. mem:



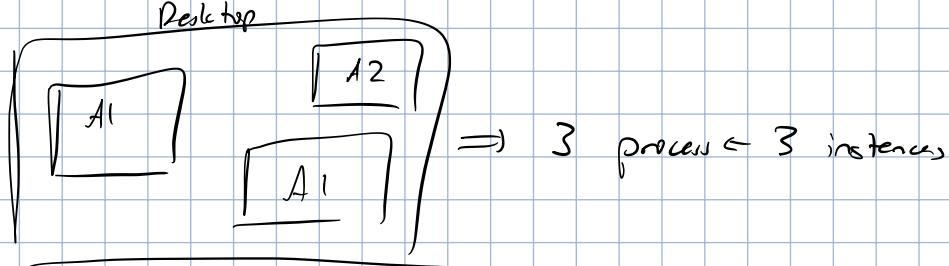
## Stack:

↳ Var values (args are pushed b/f function jump)

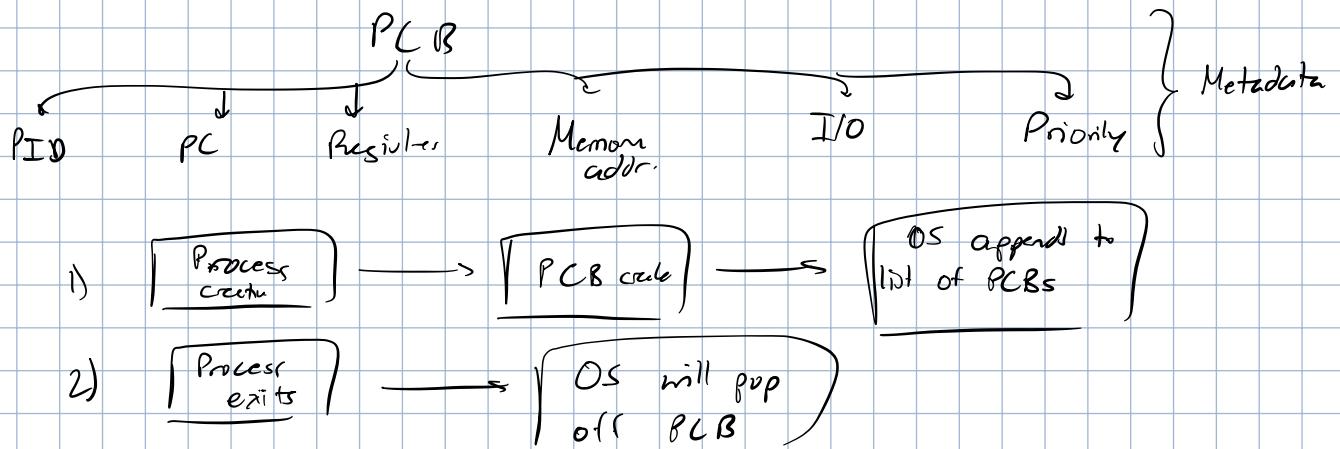
↳ Return address

## Process

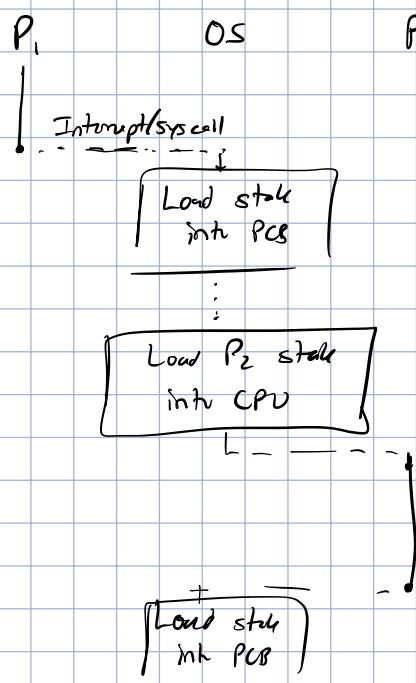
Defn: progs in execution



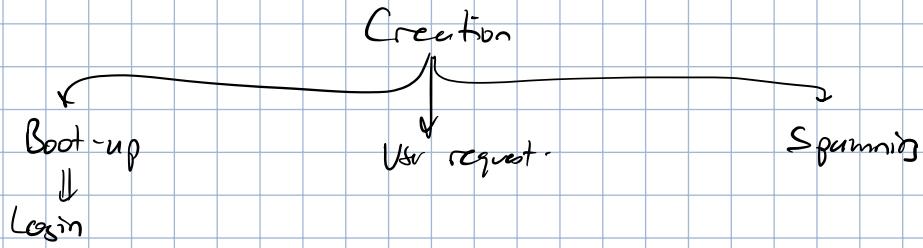
Every process has a process control block (PCB)



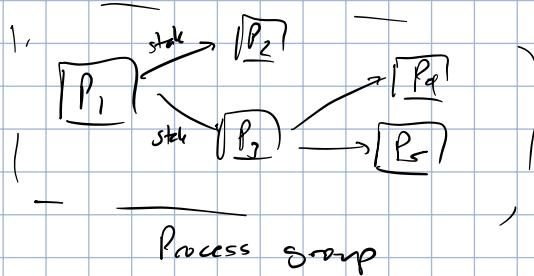
Process context switch: run multiple process on 1 CPU



Process creation:



## Process spawning



Why?

- ① Open up diff program to do smth.
  - ② Organization

How?

Program 1: `int pid = fork();`

if  $p_1 = 0 \Rightarrow$  point process

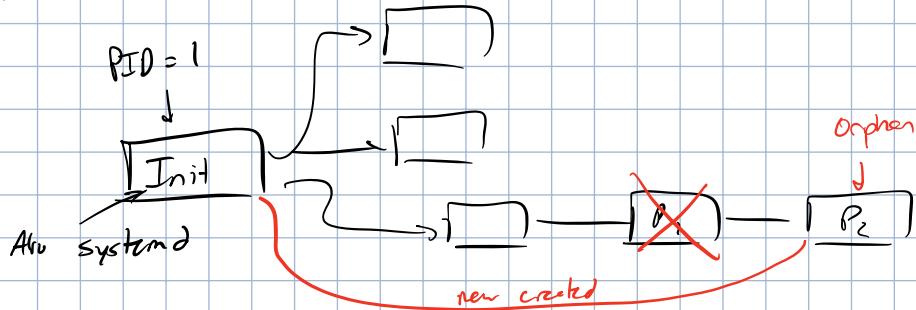
else if  $pid > 0 \Rightarrow$  PID of child proc.

To wait for process to finish: `wait (pid)`

Note: fork bomb

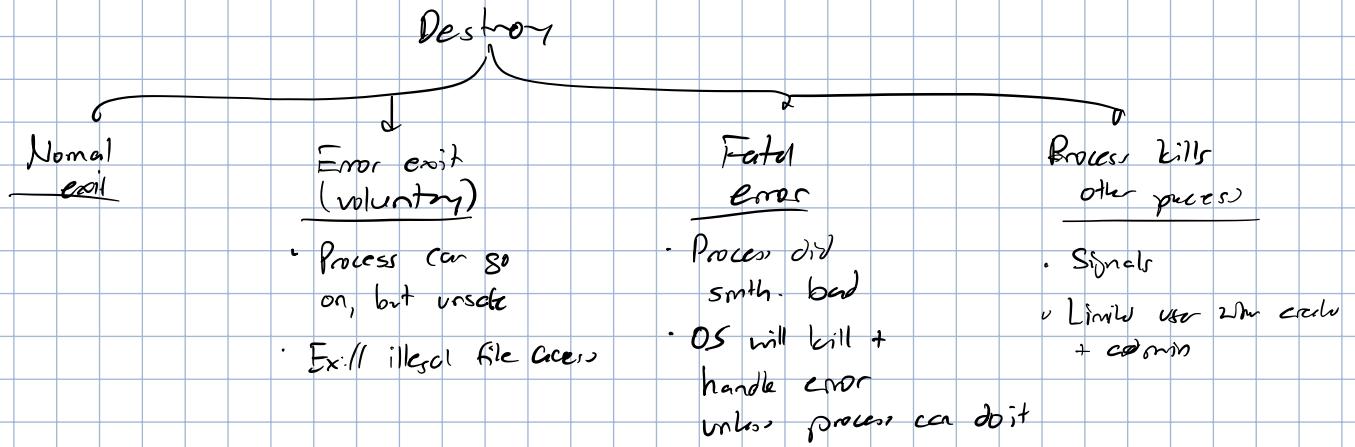
↳ Combat: limit # of process / rate of creation

## Family Tree:

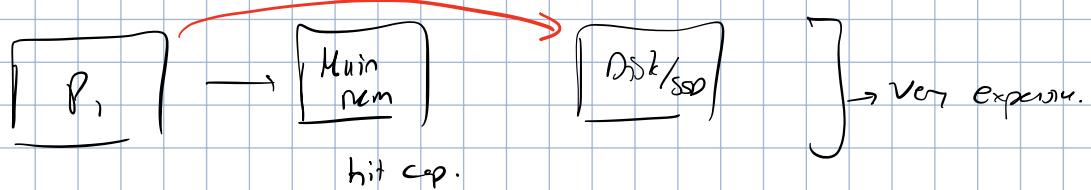


- Func. of init: kernel, pcopy

## Termination:

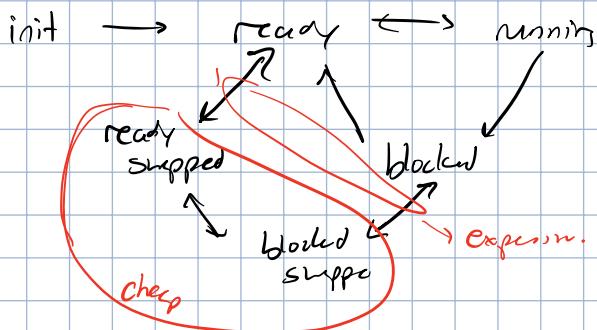


Hau own virt. memory  $\Rightarrow$  request too much mem.  $\Rightarrow$  swapping.



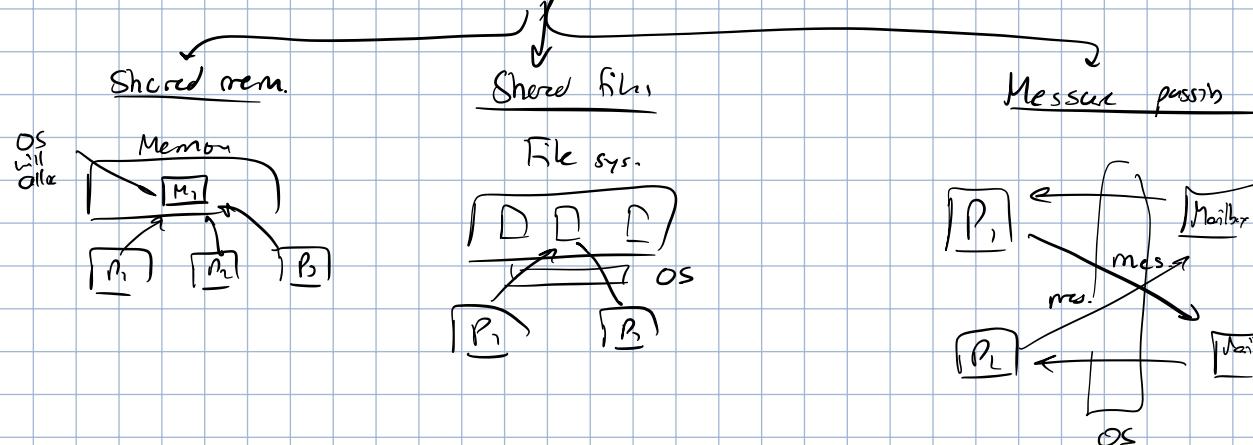
- When: blocked
- If p fails

## State diagram



## Inter-process comms.

### Comms.



## Message passing:

① Direct: use processes to indicate destn. / recip.

send (P<sub>1</sub>, message)      receive (Q, message)

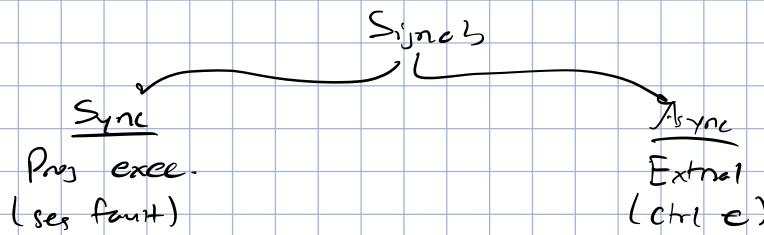
② Indirect: via process mailbox

send (MB, message),      receive (MB<sub>1</sub>, message)

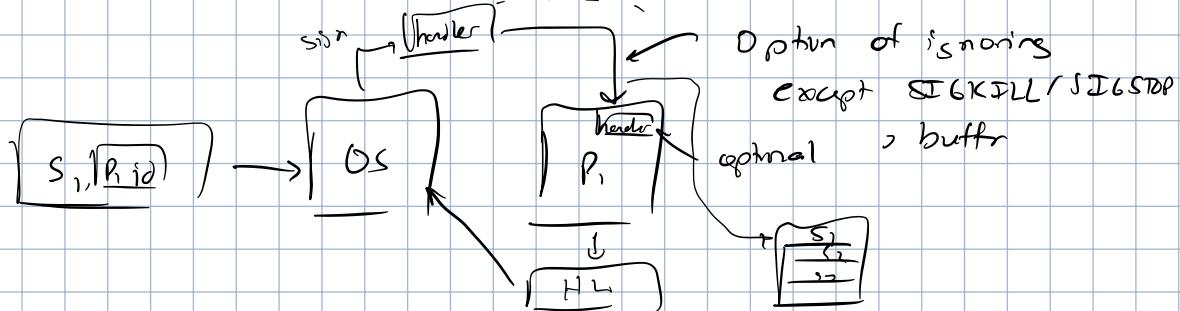
Message can be blocking / non-blocking

Signals: limits direct IPC

↙  
no data      ↗  
needs process  
no.



Workings:



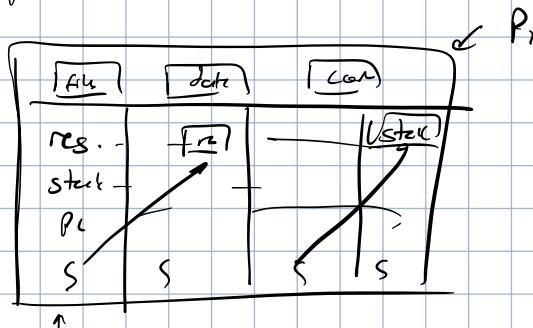
Ex-11 kill -9 PID

↳ SIGKILL signal

Custom headers should be pure  $\Rightarrow$  no side effect (malloc)  
↓  
recurrent functions

## Processes vs. threads

Multithreaded proc:



Thread is  
not isolated w/in proc

↳ Comm. are really easy

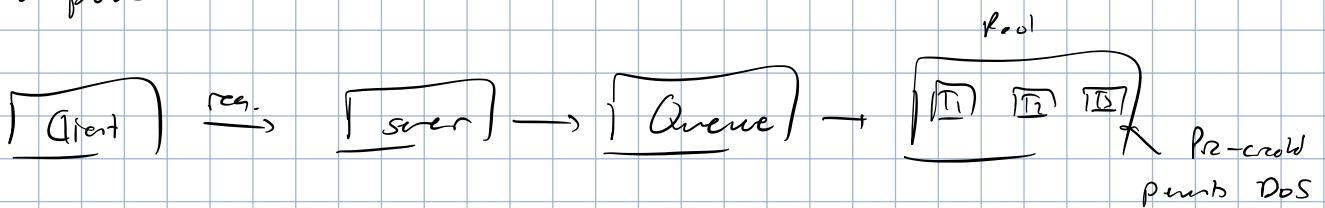
T<sub>i</sub>

CPU can only handle 1 thread / time  $\Rightarrow$  switch b/w thds  $\Rightarrow$  Concurrency

Why multi threads vs. multiprocess.

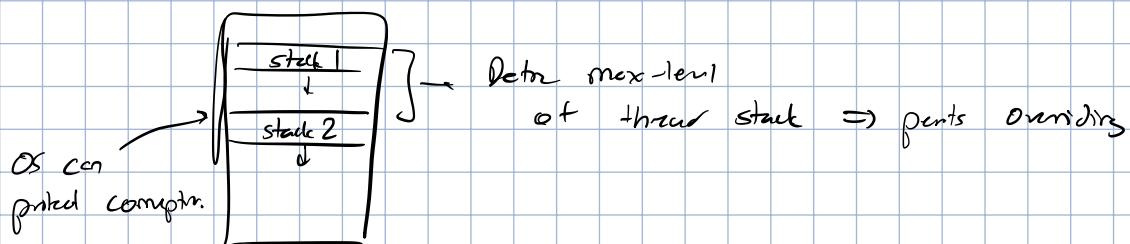
↳ Process are expensive to create, need OS

Thread pools for servers:

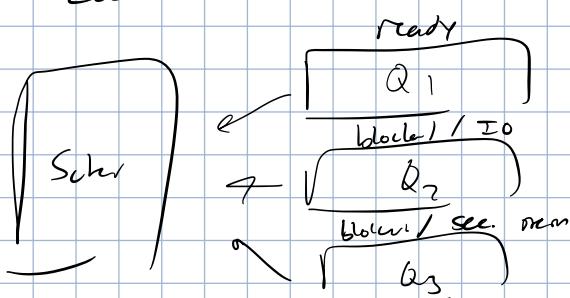
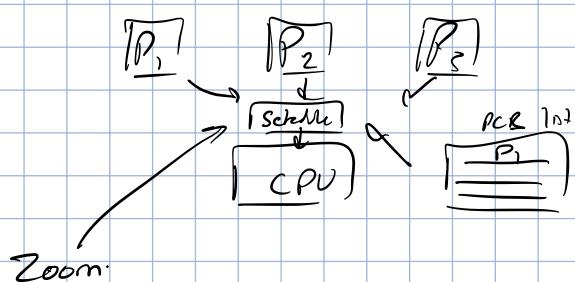
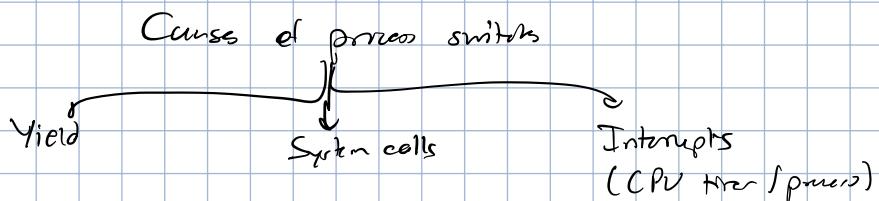


Tradeoff: speed vs. security

1 ton to handle memory of multiple threads in 1 VM:

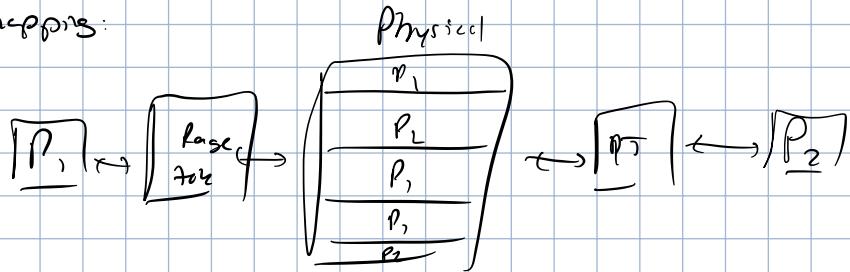


## Multiprogramming

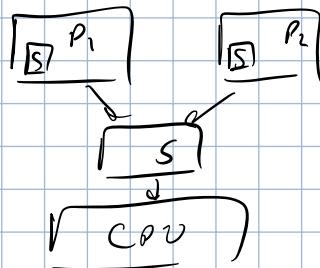


Q: multiple process switches  $\rightarrow$  protect?

Memory mapping:



1<sup>st</sup> Version: basic unithread

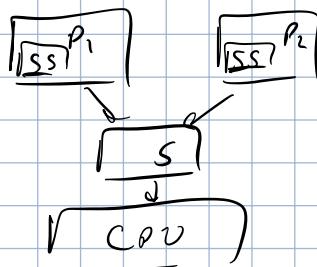


Switch overhead: high

Protection: good

Sharing: not good

2<sup>nd</sup> ver: basic multithread

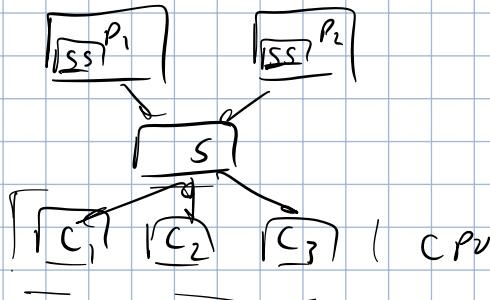


Switch: medium

Protection: medium

Sharing: medium

3<sup>rd</sup> ver: multicore

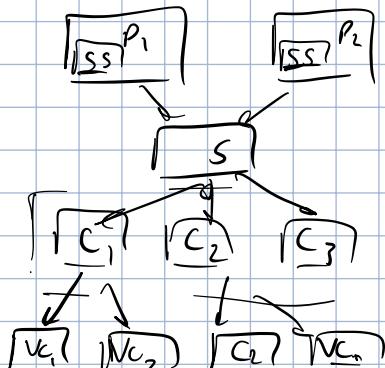


Switch: low

Protection: -

Sharing: low (thread simultaneous)

4<sup>th</sup> ver: hyperthreading:



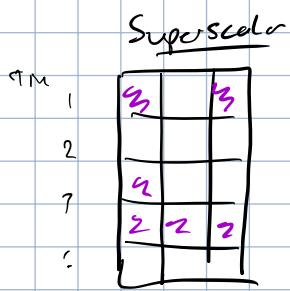
Switch: external low

Protection: -

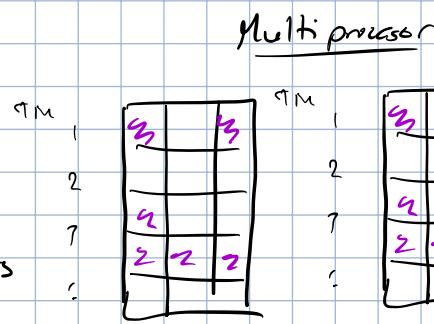
Sharing: -

Performance:  $\uparrow$

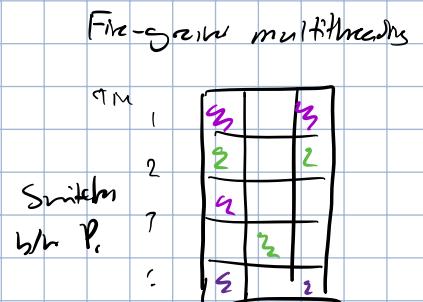
Evolution:



1 process  
Wait b/c blocking



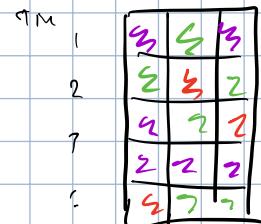
Multi processor



Switch b/w P.

1 processor  
3 ALU

Simultaneous multithreading / hyperthreading



Process executing simu.

High utilization

Hardware is coherent:

OS decides it has virt.  
(more processes to run)

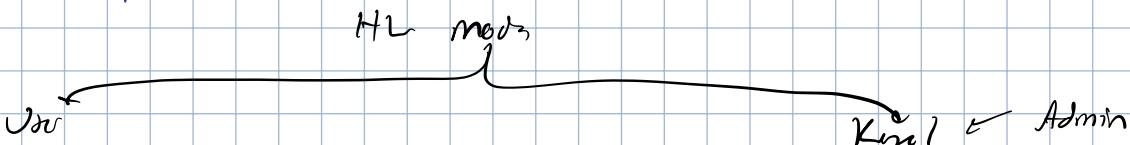
Tradeoff: performance vs. utilization

Only phys. cores  $\Rightarrow$  util  $\downarrow$ , perf (no thread comp.)

if virt. cores  $\Rightarrow$  util  $\uparrow$ , perf  $\downarrow$  (more comp.  $\rightarrow$  delay)

So: 2 VC / physical cores

## Dual-Mode Operation



Why?

◦ User shouldn't have admin privileges.  $\Rightarrow$  uncontrollable

Any code in kernel mode = kernel

Not necessary to run kernel in OS  $\Rightarrow$  OS will need to interpret code

H.W. is needed

↳ Single bit for mode

↳ Transition b/w modes  
↳ Prevent actions

Transition b/w modes

## System call

Requests kernel services

Proc → trap/syscall → kernel

## Processor except.

Internal, sync. event from hardware

↓  
from code

Requirements for safe instr

1. Limited entry into kernel  
↳ APIs few
2. Atomic change in processor state

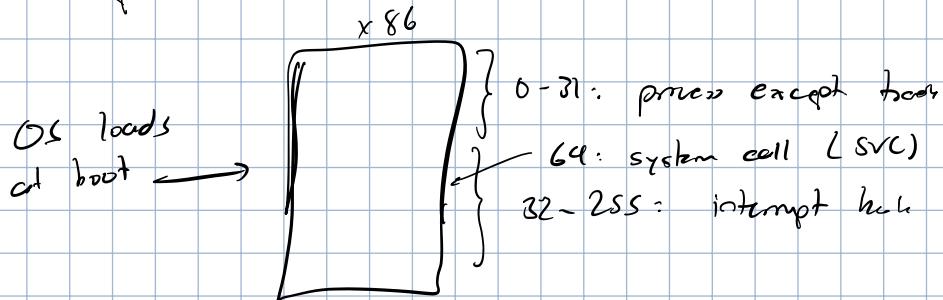
Changes: mode, SP, PC

## 3. Reversible

Easy state restoration after mode trans.

Q: handling?

Interrupt vector table:

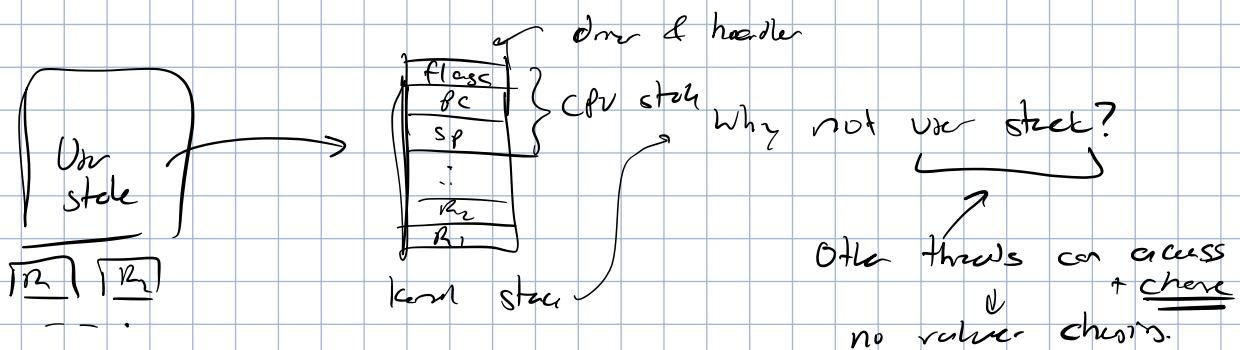


HW event → check IVT → load up handler code → exec.

Q: saving state?

State in user → kernel reversible

We use a kernel stack



1 stack / thread in kernel  $\Rightarrow$  1 k-stack in total

handler () {

  kernel stack push

$\Rightarrow$  check for malicious args

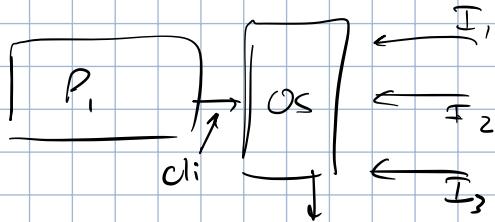
slow!

  : handler

  : pop off stack

} return (memory, proc, thread)

How to prevent interruption:



$I_1 | I_2 | I_3 \Rightarrow$  Buffer queue

B/c queue size  $\Rightarrow$  can lose interrupts

Kernel has own virtual memory  $\Leftarrow$  stack is in here

↳ Inaccessible from process

Architecture:

Monolithic

Any thing in OS  $\rightarrow$  kernel mode

↳ Pros: failure! It has to restart

modern

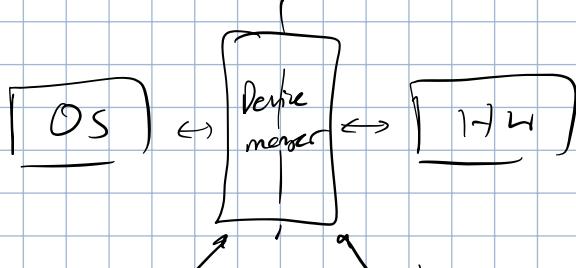
Microkernel

Keep as little as possible  
in kernel mode.

↳ Error handling.

↳ Cons: IPC in OS, not unified.

OS & HW



bottom-halt: interrupt routes, how to spin

↳ Will interrupt when HW is done

Q: how to transfer data

## ① Program I/O

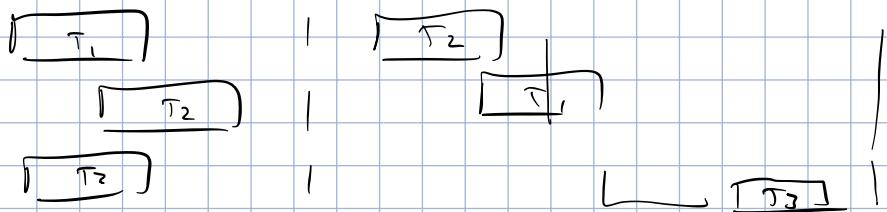
Each byte needs CPU  $\Rightarrow$  CPU time & size of a data.

Slow

② Direct memory access (DMA):

Delegate to controller to transfer data  $\Rightarrow$  free up CPU  
↳ Loss: memory address, size  
Pub it into buffer

# SYNCHRONIZATION & DEADLOCK



Possibility of conflict  $\Rightarrow$  race conditions

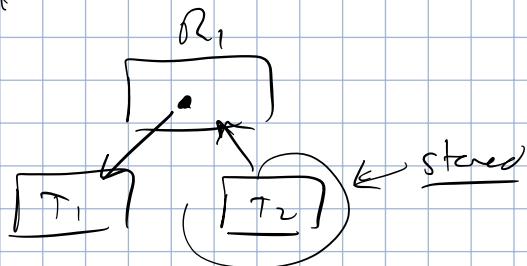
Atomic ops: locs & str of bytes

## Race Conditions

Sharing resources b/w threads  $\rightarrow$  conflicts

↳ Order of three exec matter.

Starvation lockup:



Complication: Compiler reordering instructions

↳ Assure this doesn't happen  $\Rightarrow$  sequential exec.

↑  
No.1 performant

## Mutex

will never  
be reordered

mutex.lock();

critical section  $\Rightarrow$  1 thread runs at a time.

mutex.unlock();

## Tips:

- ① Everything shared in mutex  $\rightarrow$  mutual excl.
- ② Thread that locks also unlocks.

## OS-provided Service

Issue: cannot help w/ scheduling constraints.

## Semaphores

Non-negative int count w/ 2 atomic ops:

P():

1. Wait for semc. true
2. Decrement by 1

} com

} Atomic ops. executed by 1  
threads at a time.

V():

1. Increment semaph by 1
2. Signal any waiting processes in P

} create

No thread will miss a wakeup call

Mutual exclusion: binary semaphore (0 or 1)

sem - p()

crit. sect.

sem - v()

Order matters  $\rightarrow$  Hard to reason

## Monitor

Monitor

Mutex

or condition variables.

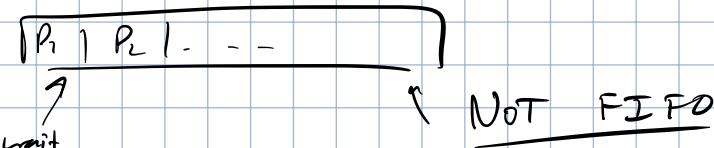
↑  
Mutual excl.

↓  
Scheduling

Condition variable:

- a) CV. wait (& mutex)  $\Rightarrow$    
 { ① Unlock your mutex  
 ② Suspend thread  
 ③ Lock again if woken up.
- b) CV. signal()  $\Rightarrow$  Wakes up anyone in ②

CV:



- c) CV. broadcast()  $\Rightarrow$  wakes up everyone in que.

Memoryless  $\Rightarrow$  no data associated w/ C.V.

Remember: re-enable threads in wait list don't always run immediately.

Design pattern:

method-waits ():

mutex.lock();

while (!testSharedState()) {  
  $\rightarrow$  CV.wait (& mutex);  
 {  
 ;  
 ;  
 ;  
 ;  
 mutex.unlock();

method-signals ():

CV.signal(); }  $\rightarrow$  lock hold

Threads don't run  
immediately  $\Rightarrow$  state could have  
changed before lock acquired.

Implementation

Hoare

- Run as soon as process wakes
- 'if' rather than 'while'

Mesa

- Above, common
- Small bug: process can wake up w/ no signal

Since CV queue is not FIFO, don't depend on FIFO behaviour.

Semaphore useful w/ HW, which cannot access locks (could look itself)

## Communicating Seq. Proc.

Style: no threads share data, use messages b/w threads instead.

↳ Done in channels.

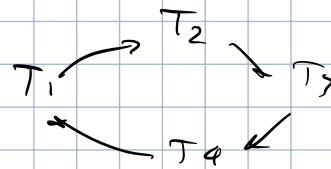
Channels still need sync.

## Deadlocks

Starvation  $\Rightarrow$  thread never gets to access data

Deadlock = starvation for all threads that need shared state

↳ Circular starvation.



Necessary, not sufficient conditions:

- ① Mutual exclusion: resource used by 1 thread at a time.
- ② Hold & wait: thread holds  $R_i$ , but waits on  $R_j$
- ③ No pre-emption for resources: if  $T_i$  holds  $R_i$ ,  $R_i$  can only be released by  $T_i$
- ④ Circular wait

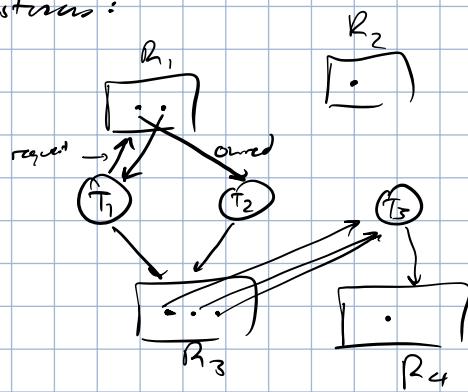
Deadlocks =  $f(\text{timings, threads, resources, order})$

Graphical view: resource allocation graph:

Threads:  $T_1 \dots T_n$

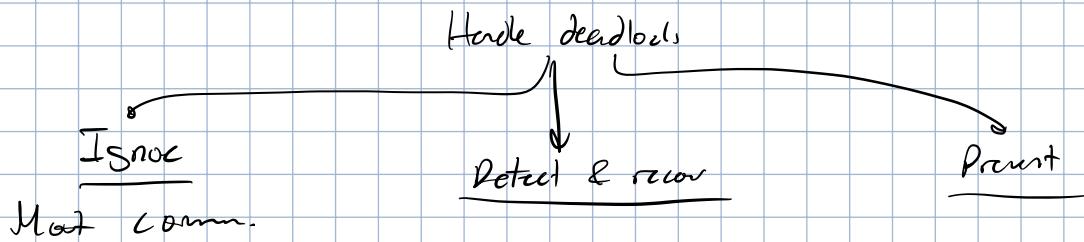
Resource types:  $R_1 \dots R_n$

Each type has instances:



Can help = circular paths  $\Rightarrow$  deadlock.

Deadlocks are dynamics



Deadlock detection:

Assumptions:

$m$  rooms,  $n$  threads

Capacity:  $\{C_1, \dots, C_m\}$

Free rooms:  $\{F_1, \dots, F_m\}$

Reserve allows:

$$\begin{bmatrix} & \xrightarrow{\text{Reserve}} \\ \text{Threads} & \begin{bmatrix} A_{1,1} & \dots & A_{1,m} \\ \vdots & & \vdots \\ A_{n,1} & \dots & A_{n,m} \end{bmatrix} \end{bmatrix}$$

Reserve requests:

$$\begin{bmatrix} R_{1,1} & \dots & R_{1,m} \\ \vdots & & \vdots \\ R_{n,1} & \dots & R_{n,m} \end{bmatrix}$$

Constraint:

$$\forall j, \sum_{i=1}^n A_{i,j} + F_j = C_j$$

Also:

$$\text{unfinished} = [T_1, \dots, T_n]$$

do {

done = true

for each  $i$  in unfinished:

if ( $\{A_{i,1}, \dots, A_{i,m}\} \subseteq \{F_1, \dots, F_m\}$ )

remove  $i$  from unfinished.

$$\{F_1, \dots, F_m\} += \{A_{i,1}, \dots, A_{i,m}\}$$

done = false

} while (done)

If unfinished is non-empty  $\Rightarrow$  deadlock (iters in unfinished strongly)

When to run:

a) Request made for resource.

b) Cannot get resource

c) Periodically

Deadlock recovery:

a) Terminate a thread  $\Rightarrow$  unlock?  $\rightarrow N \rightarrow$  do it cyclic

Issue:

1. Terminate a thread  $\rightarrow$  inconsistent state.

2. Can run into save sit.

b) Proceed until can't go further  $\rightarrow$  exception handling.

c) Rollback to non-deadlock state.

Issue:

1. Recoveres fr checkpoints

2. Same as a)

Deadlock Prevention:

Illums, fix request order, grant resource @ beg. ....

Don't request if deadlock going to occur.

$\hookrightarrow$  Avoid unsafe state: If scheduling order, thread may get blocked if max resources requested.

Ex:// 1 room  $\rightarrow$  10 intra.

3 free visitors:

Thread	Have	Max
A	3	9
B	2	4
C	2	7

Safe / unsafe?

Q: Order that allow each thread to get max resource?

① Allocate free to one w/ least diff. b/w "have" & "max"

Gr 2  $\rightarrow$  B

② Gr back max  $\rightarrow$  free:

F: S room

③ Project: B  $\rightarrow$  C  $\rightarrow$  A

Also: Banker's Also

Same as detector algo except:

$$\text{if } (\sum_{i=1}^n M_{i,1} \dots M_{i,n} - \sum_{i=1}^n A_{i,1} \dots A_{i,n} \leq \sum_{i=1}^n F_{i,1} \dots F_{i,n}) \quad \{ \dots \}$$

Before alloc

Slow & expensive.

## MULTITHREADED KERNELS

① Context switch? User  $\leftrightarrow$  kernel?

When is switch?

- a) System call
  - b) Processor exception
  - c) Interrupts
- } Sync
- } Async

We use handlers:

handler () {

    push regs to kernel stacks

,

    handle

,

    exit:

        pop regs to user stack

    return

}  $\rightarrow$  Switch more efficient.

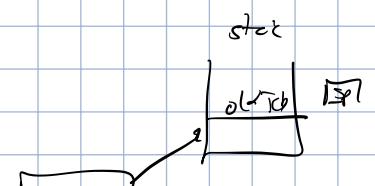
}

Thread switch func:

thread-switch (oldTCB, newTCB):

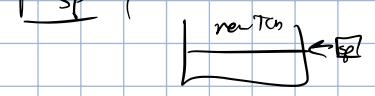
1. push all regs of oldTCB  $\rightarrow$  kernel stack

①



2. oldTCB.sp = sp

②



3. sp = newTCB.sp

③

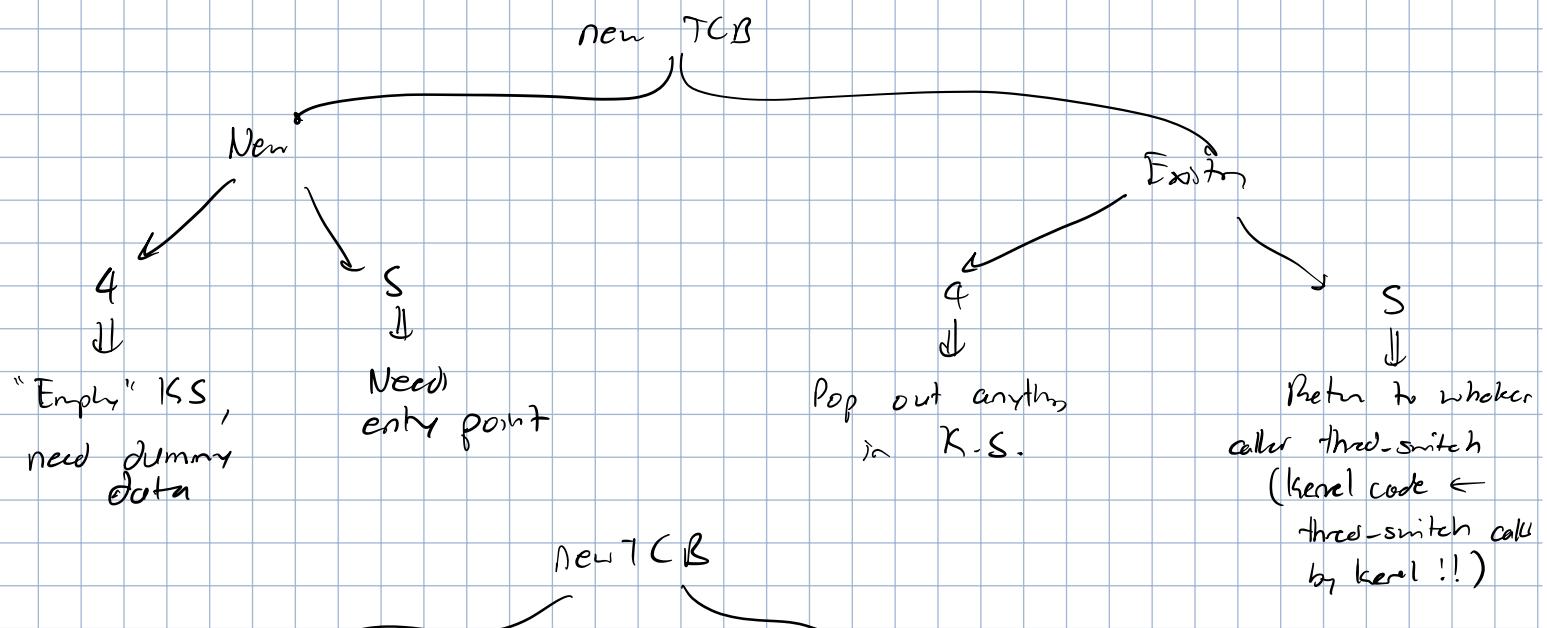


4. pop regs from kernel stack of newTCB

④



5. return

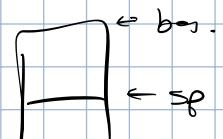


## ② Thread creation

thread-create (func, args):

TCB\* tcb = new TCB();

tcb → sp = new stack + stack-size →



push args

push func.

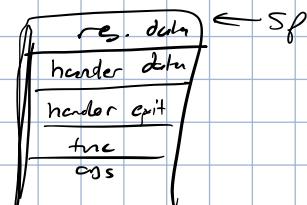
push dummy handle data

push handler exit point

push dummy data for rs.

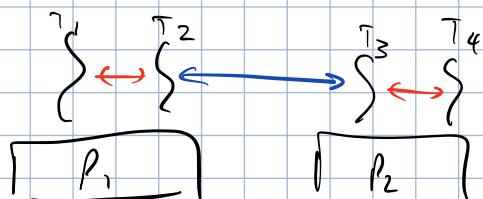
set ston

add to ready que.



Will get out of switch & handle.

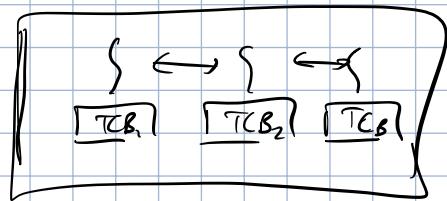
## Kernel vs. User-Managed Threads



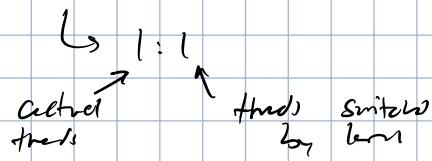
( $\approx$  faster than  $\approx$ ) fast than CPU switch

↳ Reason: caches.

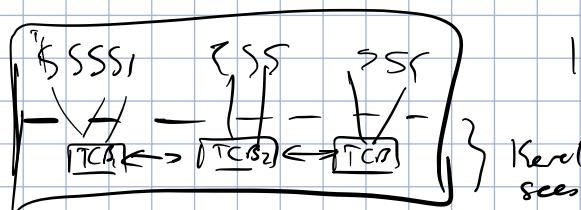
## Kernel-managed threads:



Kernel does all switching:



## User-managed threads:



Kernel does + know # of threads  
 $\therefore N=1$

- Kernel switching has no effect on thread!
- Thread that runs can get blocked if TCB assigned not scheduled by kernel & thread wants I/O
- Cons:
  - Not very parallel in multithread
  - Blocking
    - Scheduler activation: kernel message  $\rightarrow$  user  $\rightarrow$  thread.

## OS Classification:

# of addrs. spaces	1	many
# of threads	1	many
1	early OS	Unix
many	Embedded systems	Modern OS

## Synchronization Implementation

### ① Mutex

Idea #1: disable all interrupts in mutex  $\rightarrow$  prevent context switch

- Bad b/c user can do bad things without interrupt
  - Too much code not under interrupt

Idea #2: restrict no interrupts on mutex:

class Mutex:

private:  
 int value = FREE; // or Busy

Open waitq;  $\Rightarrow$  true to access lock mutex  
public:  
void lock();  
void unlock();

void lock():

disable-interrupts;  
if [value is Busy]:  
    waiting.add (curr TCB)  
    running TCB  $\rightarrow$  state = wait;  
    scheduler TCB = scheduler();  
    thread-switch (running TCB, scheduler TCB);  
    next TCB ready!  
    return  
else:  
    value = Busy  
    enable-interrupts;

void unlock():

enable-interrupts;  
if waitq is not empty:  
    next TCB = waiting.pop();  
    next TCB = READY  
    ready-list.add (next TCB);  
else:  
    value = Free  
    enable-interrupts;

Pros:

Provides section in non-interrupt state.

Cons:

If blocked in lock() / unlock()  $\Rightarrow$  disable interrupt forever

↳ Soln: thread enable interrupt after thread switch

Still relies on disabling interrupt.

Idea #3: HW atomic ops  $\rightarrow$  spinlock  $\rightarrow$  Mutex & scheduler.

test & set (& addr):  
    result = Mem[addr]  
    Mem[addr] = 1  
    return result;

swap (& addr, reg):  
    temp = Mem[addr]  
    Mem[addr] = reg  
    reg = temp

Compare & Swap (&addr, r1, r2):  
    if r1 = mem[addr]?  
        Mem[addr] = r2  
    return  
else:  
    fail

Spinlock:

class Spinlock:

private:

    int value = 0;

public:

    void lock() { while (test & set (value)); }  
    void unlock { value = 0 };

Thus after 1st lock  
continually in while loop  
busy-waits.

Pros: M.E

Cons: waste CPU

## Schedular & mutex:

```
class Mutex:
    private:
        int val = Free;
        Spinlock mutex_spl;
        Queue waiting;
    public:
        void lock();
        void unlock();
```

```
lock():
    mutex_spl.lock();
    if value is Busy:
        waiting.add (curr TCB);
        scheduler->suspend (& mutex_spl);
    else:
        value = Busy
        mutex_spl.unlock();
```

```
suspend (& spinlock);
enable-interrupts;           => prob interrupt harder
scheduler.slock();           from access; spinlock -> infinite loop
spinlock.unlock();
running TCB -> state = WAIT;
switch to next process
scheduler.unlock();
enable-interrupts;
```

Most mutexes free in comp  $\Rightarrow$  can use simple mutex impl. if no conflict, or use Spinlock-based mutexes

## ② Semaphores

```
class Semaphore:
    private:
        int value
        Spinlock semaphore_spl;
        Scheduler scheduler;
        Queue waiting;
    public:
        void P();
        void V();
```

```
void P():
    semaphore_spl.lock();
    if value == 0:
        waiting.add (TCB);
        scheduler->suspend (& semaphore_spl);
    else:
        value --;
        semaphore_spl.unlock();
```

```
void V():
    semaphore_spl.lock();
    if waiting not empty:
        tcb = waiting.pop();
        scheduler->make-ready (tcb);
    else:
        value++;
        semaphore_spl.unlock();
```

Interrupt harder cannot use P() / V()  
blk spinlock

### ③ Condition variables

#### A: Basic implementation.

class CV:

```
Queue waiting;
void wait(mutex*):
    waiting.add(tcb);
    scheduler.suspend(tcb);
    mutex->lock();
    void broadcast():
```

wait(mutex\*):

```
waiting.add(tcb);
scheduler.suspend(tcb);
mutex->lock();
```

signal():

```
if (waiting not empty):
    thread = waiting.pop();
    scheduler.make-ready(thread);
```

broadcast: chan to while

#### B: Using semaphores?

wait(mutex\*):

```
mutex->unlock();
semaphore.P();
```

signal():

semaphore.V();

1 / monitor?

Process will not wait if call wait  $\leftarrow$  state in semaphore!  $\text{Semaphore} = 0$

#### C: Semaphore / waits thread:

wait(mutex\*):

```
ordered
semaphore = Semaphore(1);
→ queue.add(semaphore)
mutex.unlock();
Semaphore.P();
mutex.lock();
```

signal():

```
if queue not empty:
    Semaphore = queue.pop();
    semaphore.V();
```

Guarantees wait!!

## UNIPROCESSOR SCHEDULING

Obj: workload of diff. tasks  $\rightarrow$  optimal sequencing of tasks.

• Minor objectives:

① Minimize overhead

② Work-conserving CPU (don't be idle)

③ Pre-emption

Model:

Task: 

CPU burst	I/O burst	CPU burst	I/O burst
-----------	-----------	-----------	-----------

 ...

↑                   ↑  
don't need  
CPU, longer > CPU burst

Metrics:

① Response time:  $t_{\text{finish}} - t_{\text{ready}}$   
↳ Not  $t_{\text{start}}$

② Throughput:  $\frac{\text{tasks}}{\text{min}}$

↳ Overhead

↳ Efficient resources

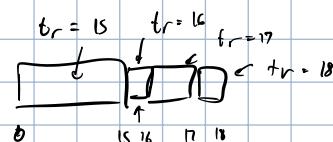
③ Fairness: CPU time shared

First Come First Serve (FCFS) Scheduling:

Pros: simple!

Cons: depends on submit time

↳ Short tasks stuck behind long tasks



But cons: short tasks come first

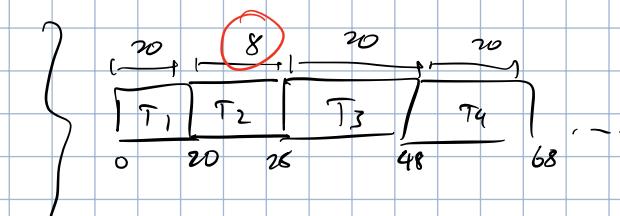
Not: 11: long task 11

Round Robin (RR) Sched:

Time quota  $\rightarrow$  each task gets time quota until expiration.

$\therefore N$  threads,  $q$  is quota: max wait time =  $(N-1)q$  (assuming overhead = 0)

Ex: 11  $T_1 = 53$   $q_1 = 20$   
 $T_2 = 6$   
 $T_3 = 68$   
 $T_4 = 24$



Wait time:  $t_{\text{start}} - t_{\text{process finish}}$

Response time:  $t_{\text{finish}}$

Pros: fair

Cons:

1. Context switching overhead adds up  $\Rightarrow$

$$\left( \frac{1}{q+1} \right) \text{ c.s. time.} < 1 \text{ s.}$$

2. Depends on  $q$ :  $q \uparrow \rightarrow$  response time  $\uparrow$ ,  $q \uparrow \rightarrow$  throughput  $\downarrow$

## Shortest Job First (SJF) Scheduling:

Heuristic: do task w/ least comp. first

Variation: SJTF  
↳ Preemptive!

Pro:

We can prove that SJF minimizes avg. response time.

Cons:

1. Long tasks are starved

2. Context switching

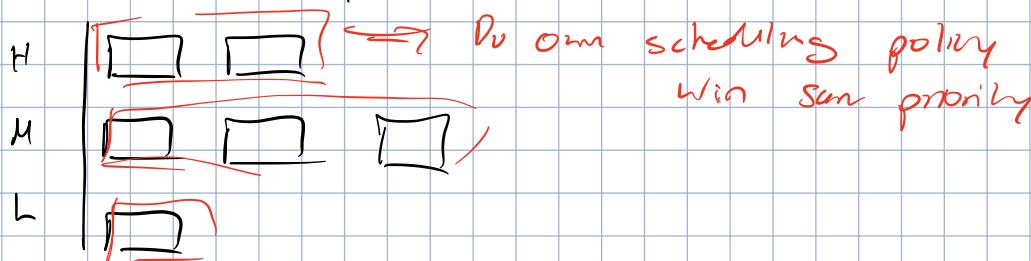
3. Impractical  $\rightarrow$  we don't know length of tasks.

↳ Round robin? Multitasking backoff.

## Strict Priority Scheduling

Prioritize tasks w/ higher priority

↳ Strict: no higher priority waits for lower priority  $\Rightarrow$  pre-emption



Problems:

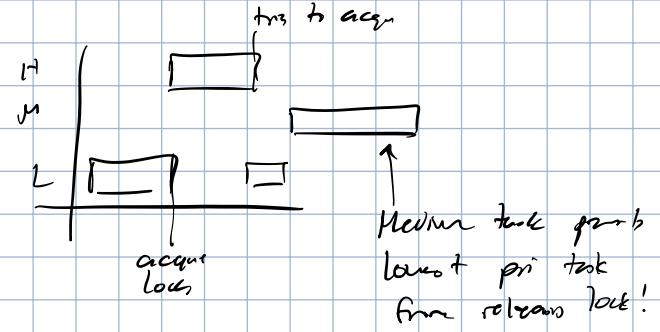
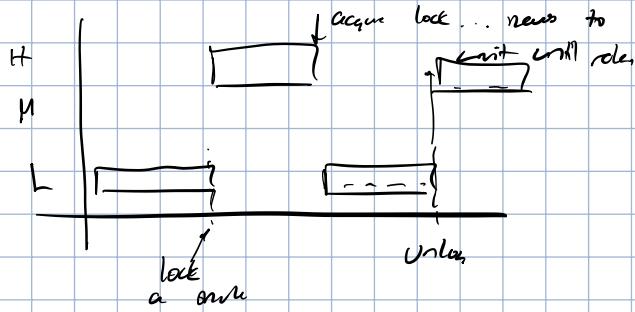
Sol:

1. Starvation of lower prio. task.  $\Rightarrow$

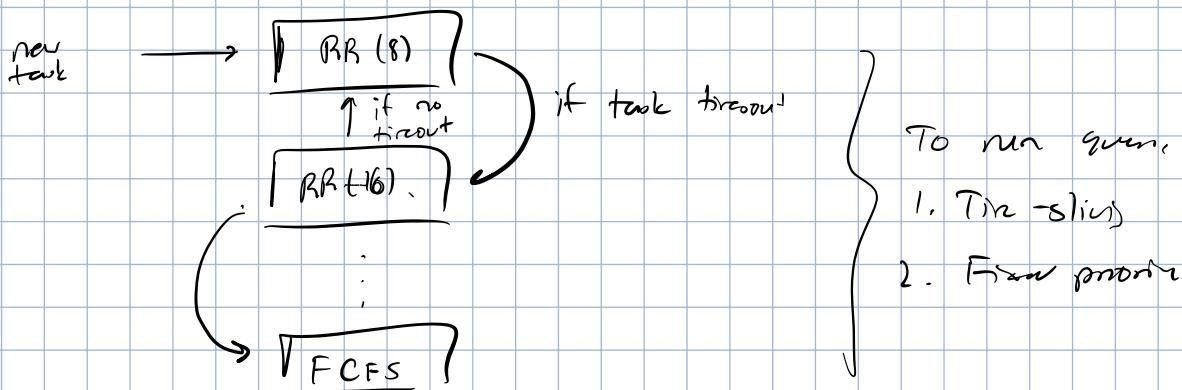
1. Tie slacks b/w priority
  2. Elect long-running task prio
- } will affect response.

2. Priority inversion: low priority delays high priority by holding shared resource.

Ex: //



## Multilevel Feedback Queue



Pros:

1. Fair: short tasks run quickly, long-running tasks will run.

## Lottery Scheduling

Steps:  $\downarrow$  *probation, user*

① Give tasks some lottery tickets  $\Rightarrow$  # of tickets  $\propto$  importance

② At each time slice  $\rightarrow$  run lottery  $\rightarrow$  run task w/ lottery #

Pros:

1. On avg: CPU time for task  $\propto$  property of ticketed tasks
2. Fairness: give every task 1 ticket
3. Works well in cheap OS

Cons:

1. If too many short tasks set # of tickets  $\rightarrow$  RR  $\rightarrow$  response time
2. Randomness

## Strid Scheduling

Obj: proportional time share w/out randomness

Algo:

- ① Stride per thread:

$$\text{Stride: } \frac{W}{N_i} \Rightarrow \begin{array}{l} \text{big #} \\ \# \text{ of latency ticks (weight)} \end{array}$$

② Maintain pass counts for each thread

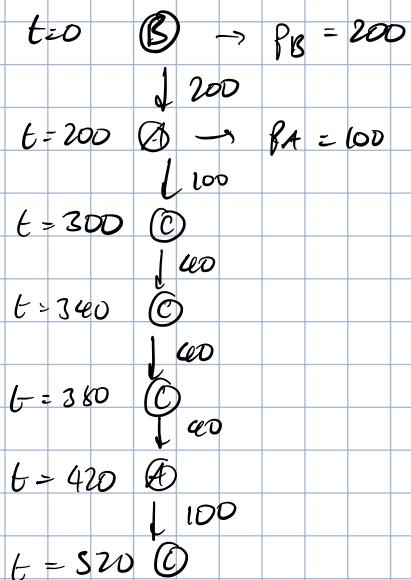
③ On each tick slice:

a: Find thread w/ least pass count (if tie  $\rightarrow$  pick randomly)

b: Increment time w/ start of selected thread

Ex://  $S_A = 100$     $S_B = 200$     $S_C = 40$

$P_A = 0$     $P_B = 0$     $P_C = 0$



Proportional tick share

## Min-Max Fair (MMF) Scheduling

Obj: pick process w/ lowest acc. CPU time

1<sup>st</sup> impl:

- ① Look at specific quota & run  
 ② Switch to thread w/ low. acc. time.

If too many tasks  $\rightarrow$  flk  $\rightarrow$  ~~copy task~~

2<sup>nd</sup> impl: dynamic time quota

- ① Set target latency (interval b/w task runs)

- ② Set time quota:

$$q = \frac{\text{target latency}}{N} \quad \text{N} \leftarrow \# \text{ of thds}$$

If  $N \uparrow \Rightarrow$  too much context switching

- ③ Min. granularity:

pick max ( $q$ , min gran.)

3<sup>rd</sup>: weighted max-min fair schedul.

① Assign  $W_i$  to thread i

② Take quota:

$$q_i = \frac{W_i \times \text{target latency}}{\sum W_i} \Rightarrow \text{Ratio mtn.}$$

③ Keep track of how much the acc.  $\Rightarrow$  switch at end of  $q_i$

↪ Use virtual time:  $\frac{\text{actual time}}{W_i}$

④ Switch according to lowest virtual time.

Why better sched? Dynamic quota  $\rightarrow$  less context switches.

Linux is similar to MMF

## Modelling Scheduling

① Det. modelling: pick workload & alg  $\rightarrow$  compute perf.

Not practical  $\Rightarrow$  can't just "pick" workload IRL

② Queuing modelling: assume events happen in some distn.

Not practical for networked events.

③ Simulation: run system

Bounce of false comparison

## ADDRESS TRANSLATION

### Basics

Two resp:

① Protect memory: share data / excl. data

② Allocate mem.

Units:

$$\text{KiB} \neq \text{KB}, \text{GiB} \neq \text{GB}, \text{MiB} \neq \text{MB}$$

$\frac{1}{2^{10}}$        $\frac{1}{10^3}$

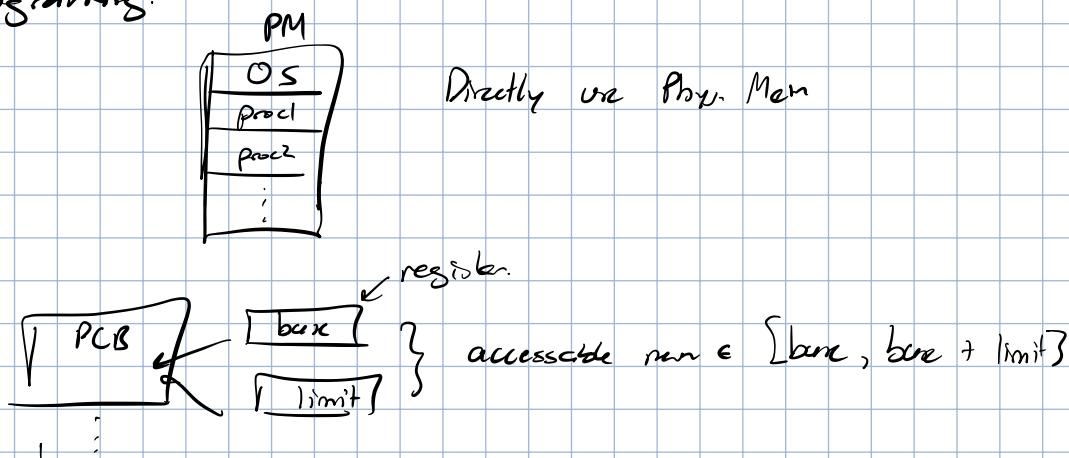
Size of system  $\Rightarrow$  dictate size of addressable contents  
 $n$  bits  $\Rightarrow 2^n$  addr.

# Protection

## ① Uniprogramming

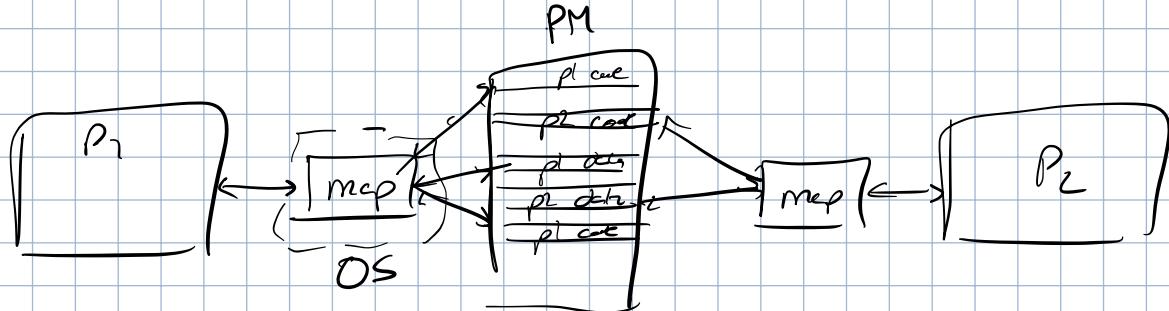
1 prog @ time  $\rightarrow$  entire phys. mem for progs

## ② Multiprogramming:



Issue: no sharing.

## Translation



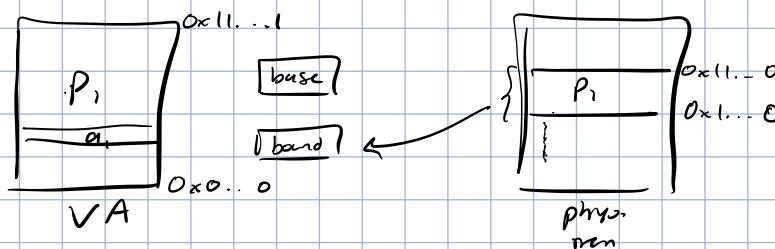
### Pros:

- Illusion  $\rightarrow$  easy dev.
- Protection
- Sharing

### Cons:

- Overhead
- HW support

1<sup>st</sup> impl.: base & bound.



## Algo:

- ① If  $a_i > \text{bound} \Rightarrow$  after add, too big  $\Rightarrow$  raise exc.
- ②  $PA = a_i + \text{base} \Rightarrow$  access

### Pros:

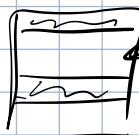
- Protection & isolation.
- Lower overhead

### Cons:

- No sharing.
- Base & bound will shift (happ)  $\rightarrow$  not good.

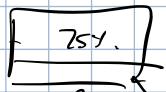
• Frag.

External:



no block  
can be  
allocated

Internal:



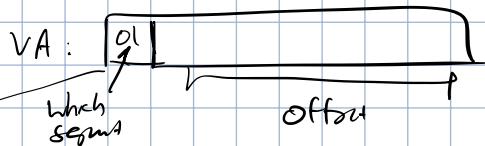
free span  
that  
cannot be  
reused

2<sup>nd</sup> impl: multi-segment addr. trans.

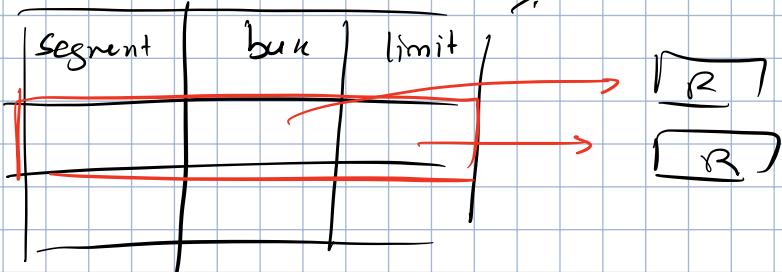
① Give each segment of proc. diff areas  
heap, stack, env vars, ...

② Bound & base for each segment

③ Translation pro:



Multi-segment trn.

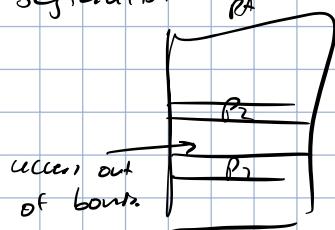


$$PA = \text{base}[\text{segment}] + \text{offset} \quad \text{if } \text{offset} < \text{limit}$$

Pros:

① Reduces frag.  $\leftarrow$  smaller segments

② Segfaults:



③ Easier for growing mem.  $\rightarrow$  new to copy mem.  
 $\hookrightarrow$  will can segfault, but OS handles.

Cons:

① External frags.

② Slapping optim. limits.

Paging

Page: fixed size memory blocks (physical, virtual)

- Somewhat small (1 - 16 KiB)

Translation:

- ① Figure out size of page  $\Rightarrow$  # of bits to address a part of page

Ex: // 1 KiB  $\Rightarrow 2^{10}$  bytes  $\Rightarrow$  10 bits to address byte.

- ② Split virtual address will have page offset & virtual page #



- ③ Table to translate VM  $\rightarrow$  PM page #:

VM P #	PM P #	Permissions	Valid	In memory OS has reference
		(R/L W/R X)	Y N Y	

Per process  $\rightarrow$  virt. mem. changes / proc.

Heuristics:

1. Demand paging: only keep active mem. in table. Inactive  $\rightarrow$  disk, invalid
2. Copy-on-write: on fork  $\rightarrow$  copy PT & set new perm. If write required  $\rightarrow$  copy mem. block
3. Zero-fill-on-demand: set new data alloc to 0

Pros:

- Sharing is easy
- Page faults

Cons:

- Page table is huge: 1 page table  $\Rightarrow$  # of VM pages.

Ex: // 32 bit addr., 4 KiB page.

- ① # of bits to address byte.

$$4 \times 2^{10} = 2^{12} \Rightarrow 12 \text{ bits in addr as offset}$$

- ② # of bits for VM page:

$$32 - 12 = 20 \text{ bits}$$

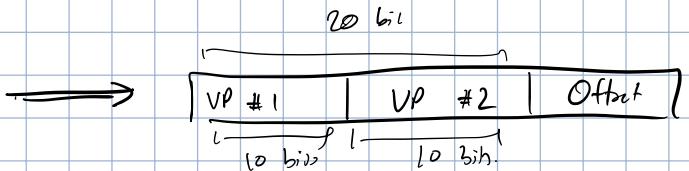
- ③ # of VP:

$2^{20}$  pgs. If each PTE is 4 B

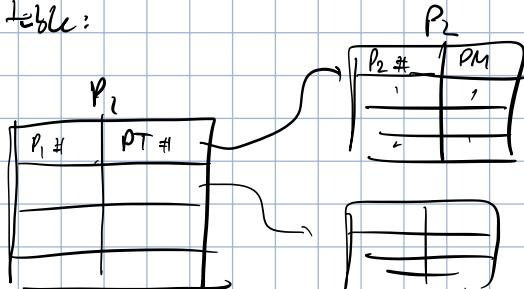
$$\therefore 4B \times 2^{20} = 4 \text{ MiB}$$

## Multilevel Paging

VP #	Offset
------	--------



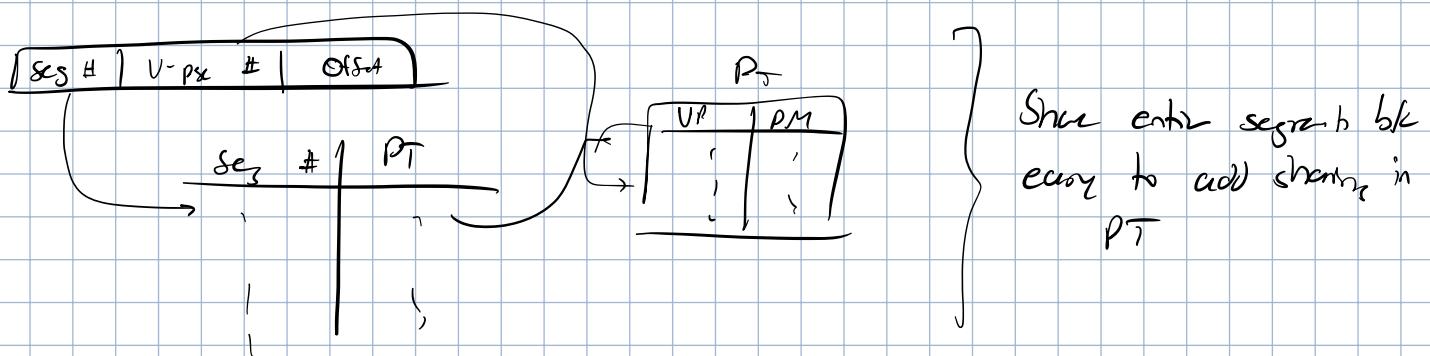
2 level page table:



This saves a lot of space:

If 1<sup>st</sup> level PTE not set  $\rightarrow$  Associate  $P_2$  is not free  $\rightarrow 2^{10}$  pages saved.

Dif. flavor: segment pages



X86:

- 4 levels pages  $\rightarrow 2^4$  entries.
- PTE: 64 bits

Shared libraries

Process should have own global & static variables to be used in shared lib.

- On proc. switch  $\rightarrow$  change value of vars.
- OS creates global-offset table w/ all variables.

No absolute VA in shared library code  $\Rightarrow$  linked dynamically

OS creates procedure-linking table to go to procedures not in shared lib space

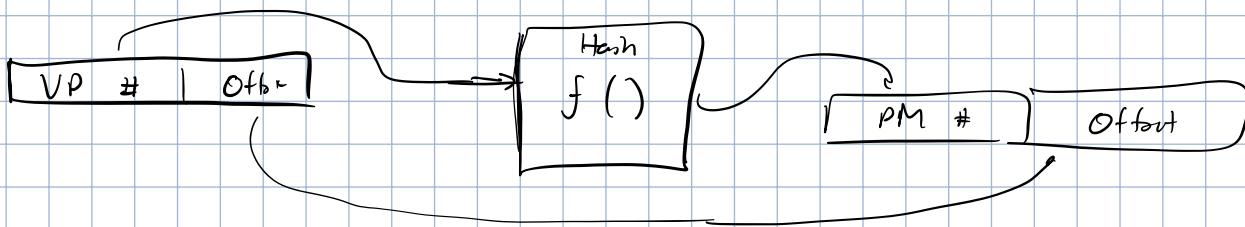
Pros:

- Sharing
- Easy to alloc (page fault)
- Save space.

Cons:

- Overhead
- PT has to be contiguous in PM

## Inverted Page Table



Set size of page table  $\rightarrow$  no need to expand to address entire VM

Con:

- Poor cache locality (spatial)
- Collisions w/ PM

Hardware:

Most translation in TLB  $\rightarrow$  fast (MMU can detect page faults)

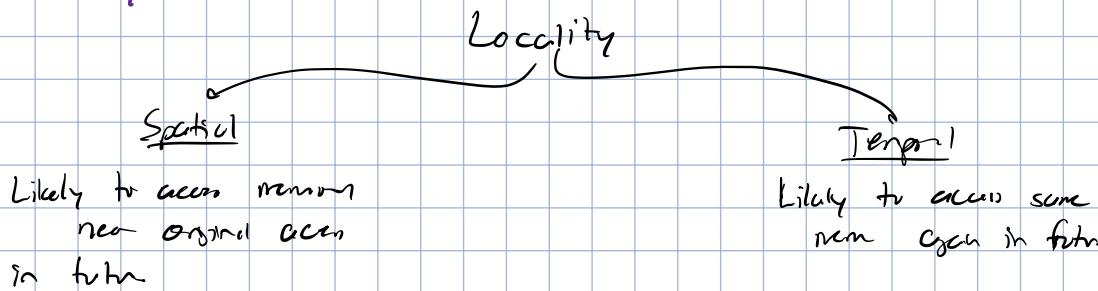
## CACHING

Makes common case fast

2 assumptions:

- ① Fre. case is frequent
- ② Infras. can is not too expensive.

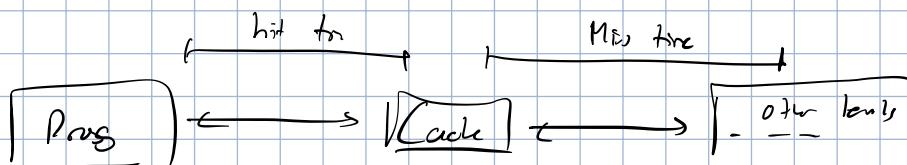
## Locality



In busy sections of memory to cache to prevent page to go to PM

## Terminology:

- Block: consecutive bytes (unit of cache)
- Hit: occurs cache & find item
  - ↳ hit time = time to hit - time to start access
- Miss: occurs cache & fail to find item



o 1st rule:

$$\frac{\# \text{ of hits}}{\# \text{ of hits} + \# \text{ of misses}}$$

o Avg. access time

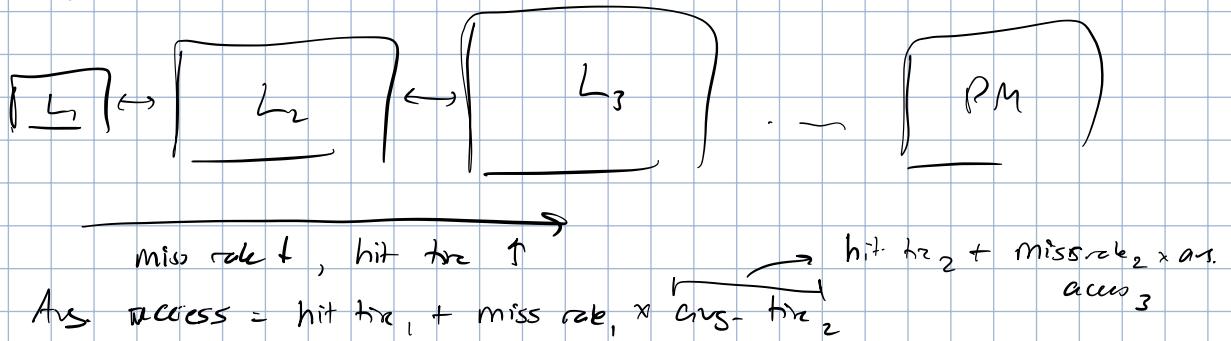
$$= \frac{\text{hit rate} \times \text{hit time} + \text{miss rate} \times (\text{hit time} + \text{miss time})}{\text{hit time} + \text{miss rate} \times \text{miss time}}$$

Tradeoff: low hit time & high miss rate.

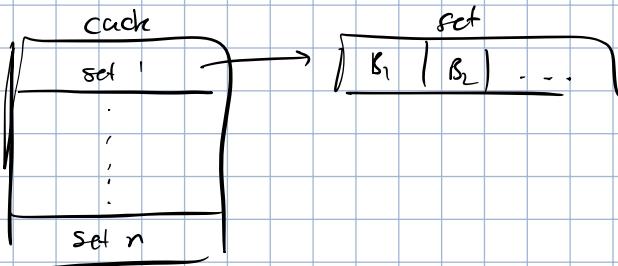
Increase cache: miss rate  $\downarrow$   $\rightarrow$  hit time  $\uparrow$

Decrease cache: hit time  $\downarrow$   $\rightarrow$  miss rate  $\uparrow$

Soln: hierarchy of cache:



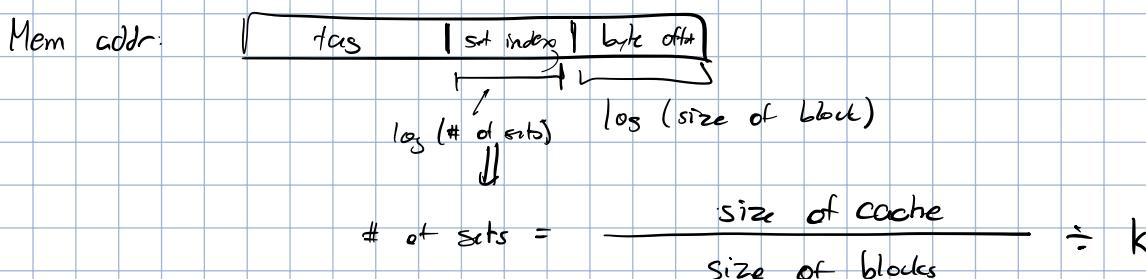
## Cache Mechanics



Associativity: # of blocks / set



To find block in cache:



On access:

- ① Find set
- ② Find block by looking at tags

### ③ Find byte via offset

Ex:// 1 KiB, direct-mapped. 32 B blocks. Mem. addr?

32 B blocks  $\Rightarrow$  5 bits for byte offset

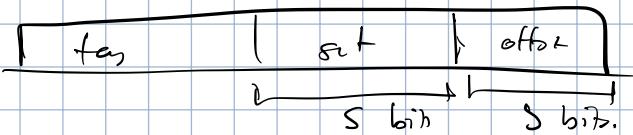
1 KiB cache, direct mapped.

$$\frac{2^{10} B}{\text{size of block}} = \frac{2^{10} B}{2^S B} = 2^S \text{ blocks to str.}$$

Since  $D = 10 \Rightarrow 32$  sets.

$\therefore$  Set index: 5 bits

Rest is tag:



Ex:// Same set., 2-way associat.:

Could store  $2^S$  blocks in tot.

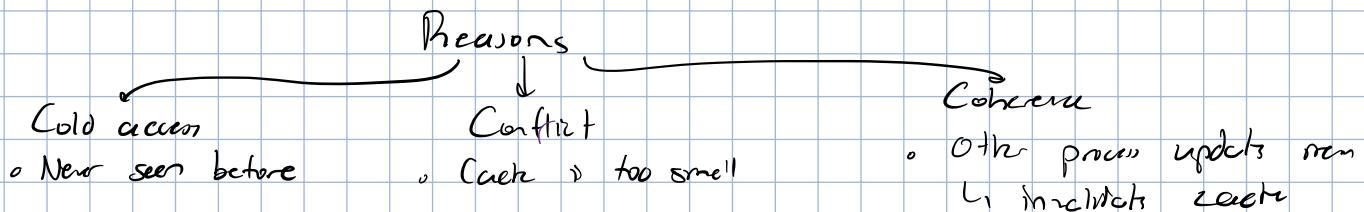
$$2^S \text{ blocks} \times \frac{\text{set}}{2 \text{ blocks}} = 2^4 \text{ sets} \Rightarrow 4 \text{ bits.}$$

Prob: new memory always has same set index  $\rightarrow$  only use 1 set, cache is mostly empty!

L. Soln: coloring  $\rightarrow$  PA is 'colored' so that blocks of same color  $\rightarrow$  same cache set

OS will reduce use seen VA's from diff. colors.

### Cache Miss & Eviction



### Eviction:

- Direct mapped: kick out block
- M-way associat.: LRU, random choose

## Cache Writes

### Write methods

#### Write-through

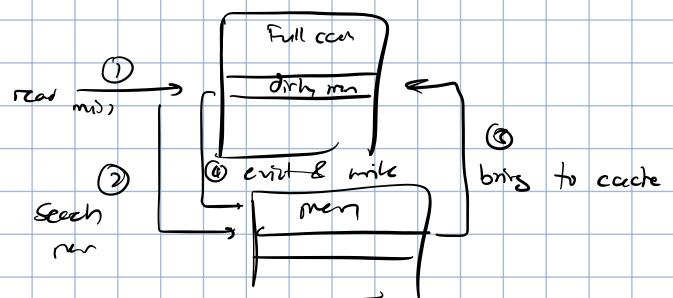


o Pros: read misses will not create writes.

o Cons: processor is slow  $\rightarrow$  write buffer.

#### Write-back

- o On write  $\rightarrow$  modify cache to signal it's dirty
- o On read miss & replace  $\rightarrow$  write to memory



o Pros: no write holdup

o Cons: overhead on read miss

## Address Translation

Cache VA  $\leftrightarrow$  PA translations  $\Rightarrow$  TLB (translation lookaside buffer)

$\hookrightarrow$  If miss, check PTE

Q: What if address from cache

Cache  $\uparrow$   $\Rightarrow$  nothing

Cache  $\downarrow$   $\Rightarrow$  flush out translation.

Q: How are writes handled?

Use a dirty bit in TLB.

If bit == 0  $\Rightarrow$  write through on bit (set dirty)

If bit == 1  $\Rightarrow$  don't have to write through

OS can enable dirty bit

If read  $\rightarrow$  update if necessary

Challenges:

① Homonymy

B/w diff processes: 1 VA  $\rightarrow$  many diff. PA

Soln: add proc. ID on TLB entries to point many PA usage

## ② TLB Shoot down

TLB / processor  $\Rightarrow$  1 processor updates TLB  $\rightarrow$  reflects in all other proc. TLB

Need to send inter-processor signal to remove entry if update.

## ③ TLB miss

Expensive  $\leftarrow$  Page table walks!

Q: What do we do if we access invalid page / page is not in PT?

A: Hardware for PTW

Page interrupt  $\rightarrow$  OS will trap to page fault handler

B: Software for PTW

Page fault raised  $\rightarrow$  kernel will do PTW  $\rightarrow$  page fault handler

Challenge: out of order execution

Ex. 11 1. mul r1, r2, r3

2. bne r1, r4, loop

3 ldr r5, (r6)

Issue: compiler may order 3 before 2  
if 3 faults  $\rightarrow$  should 2  
be done again? Out-of-order  
exception.

Solution: use precise exception

① Keep track of retired instr. (state is in HW)

② Only handle exceptions if instr. that raised is retired  
↳ all prev. instructions should complete.

In architectural  
state (res.,  
mem.)

Speedup of TLB?

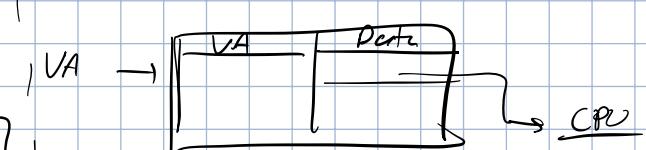
① Multiple levels of TLB



② Virtually-addressed Cache:



VA Cache



Why speedup? No translation?

Challenges:

1. Synonyms: diff VA  $\rightarrow$  same PA

Caché stores multiple KVP w/ same VA  $\Rightarrow$  Aliasing

Soln: tag entry in caché w/ PA & evn some VA

L, How to translate VA  $\rightarrow$  PA via TLB  $\rightarrow$  check if correct data in caché

This problem happens if you depend on VA in any DS.

③ Bring PTE to cache

RAMB Sory to mem.

④ Set associativity:

$$\text{Avg mem access time} = \text{hit time}_{\text{TLB}} + (\text{miss ratio}_{\text{TLB}} \times \text{miss time}_{\text{DRAM}})$$

↓ ↑  
Direct mapped      Fully associative      ↑ ↓

Cost of miss is really high  $\rightarrow$  fully associat.

⑤ Superpages

TLB will now address bigger pages  $\rightarrow$  don't need + have more entries!

L, Cons: more offset bit needed in VA

Need a bit to indicate whether PTE is superpage

# of PTE  $\downarrow$

Linux

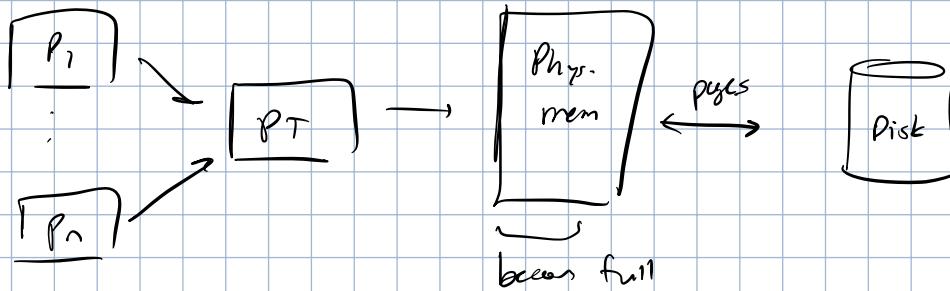
Issue: TLB filler w/ user progr. translators en kernel switch

Soln (superpage): kernel all in same addr, no need to flush TLB

L, Security concerns

## DEMAND PAGING

### Intro



Physical mem. acts like cache!

- Block size: 1 page
- Associativity: fully-associative
- Write-back cache upon eviction

Switching is v. expensive, switching pages is a background

Q: How do we know if page is switch out?

Valid bit in PTE (1 if in PM, 0 if in disk)  $\rightarrow$  cascades to all entries.

This valid bit is updated on eviction

If we try to access  $\rightarrow$  raise page fault  $\rightarrow$  OS will bring back

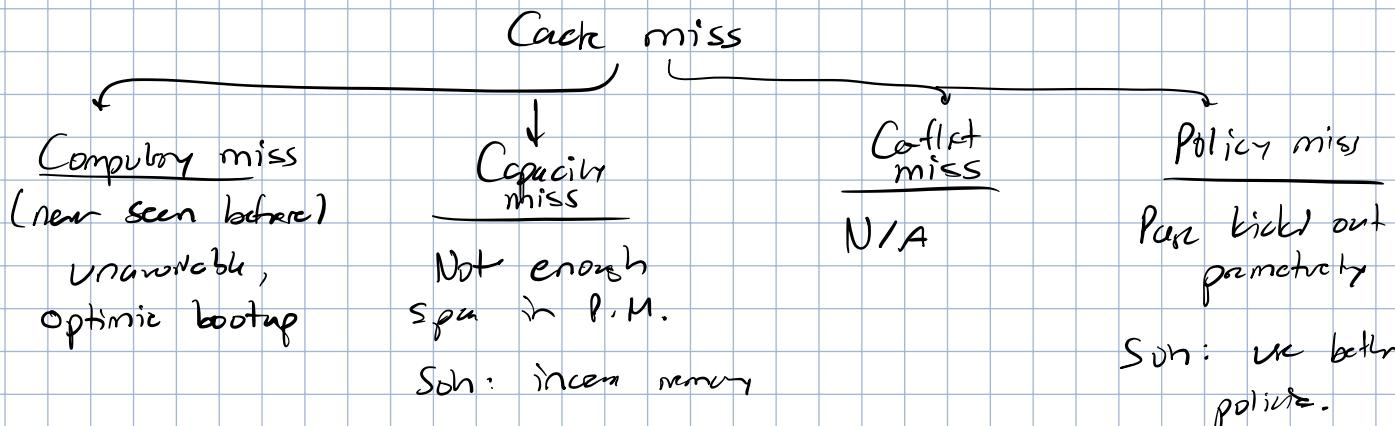
Backing store / swap file:

Disk will mirror all virt. pages

Q: How do we bring back pages into P. M.?

OS will keep track of non-recent V.P.

Q: What causes misses?



## Page Replacement Policies

### ① Random eviction

- Simple ↳ Use in TLR
- Not efficient

### ② FIFO: kick out oldest

- Fair
- Oldest pages could be well used.

### ③ Minimum: kick out pg w/ least in fut.

- Optimal
- No idea about fut.

### ④ LRU: kick out pg w/ least in past

- Good approx. of min., but not necessarily best.
- Overhead

## Belady's Anomaly:

Thought: just incres. amount of memory  $\rightarrow$  drops miss rate?

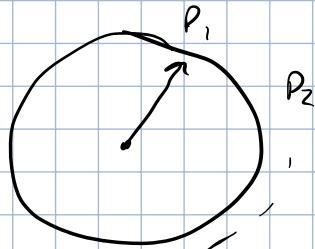
Ans: NO! In certain patterns  $\rightarrow$  incres. size, incres. in perf. fault (FIFO)

## LRU: Clock Algorithm

LRU w/ L.L. has overhead.

Clock algo:

- ① Arran phys. pg in circ.



TLB/PTE  
↓

- ② Whenever we evict, we incant hand & examine "access bit" of pages:

If access = 0 (it is not access)  $\rightarrow$  kick out

If access = 1  $\Rightarrow$  set to 0 and move to next page.

③ If all access bits = 0  $\Rightarrow$  FIFO policy

If hand is moving fast  $\leftarrow$  lot of page faults.

LRU: N<sup>th</sup> Chance Algorithm

- ① Set variable/page to record # of times it could be evicted.
- ② If hand is at pos & access bit is 0  $\rightarrow$  increment variable.
- ③ If variable = threshold  $\rightarrow$  evict
- ④ If access bit is 1, reset to 0

Threshold: # of chances that we will give before eviction

- ↳ Too high: it will take too much time.
- ↳ Too low: not good approach.

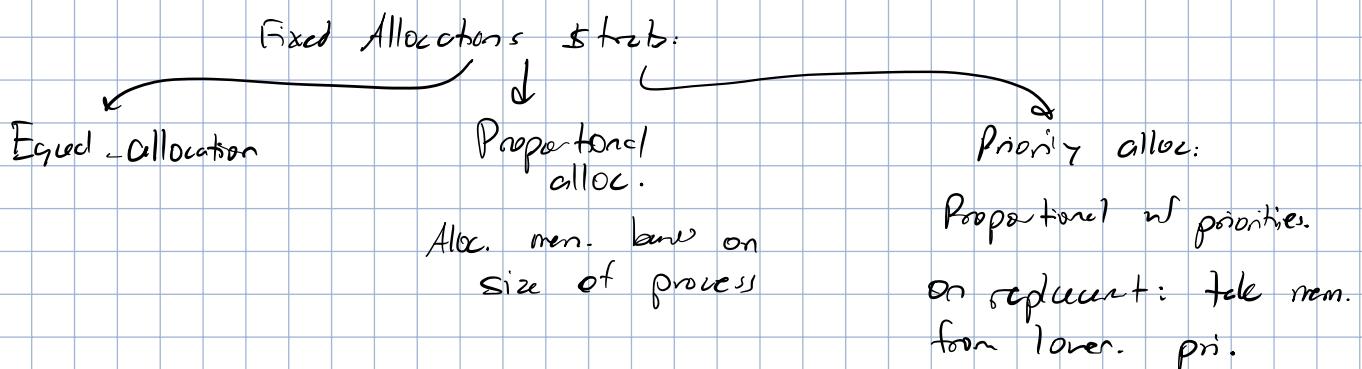
Pros: better approach of which pages are least used.

If page is dirty: write back & give 1 more chance.

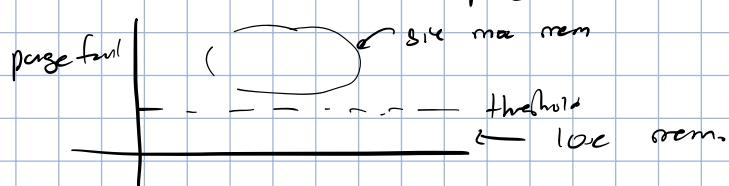
Q: when to run?

- ① On page fault (add. counter on page fault)
- ② from async

Memory Allocation



Adaptive alloc.: give mem. according to page fault



Thrashing: processes busy swapping mem. from each other

Sol:

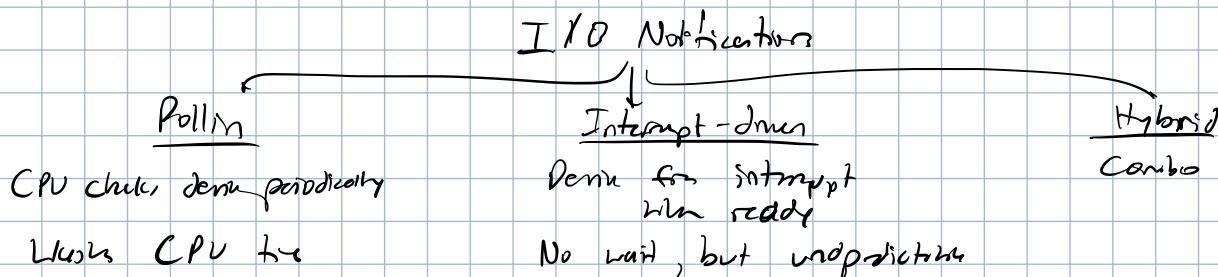
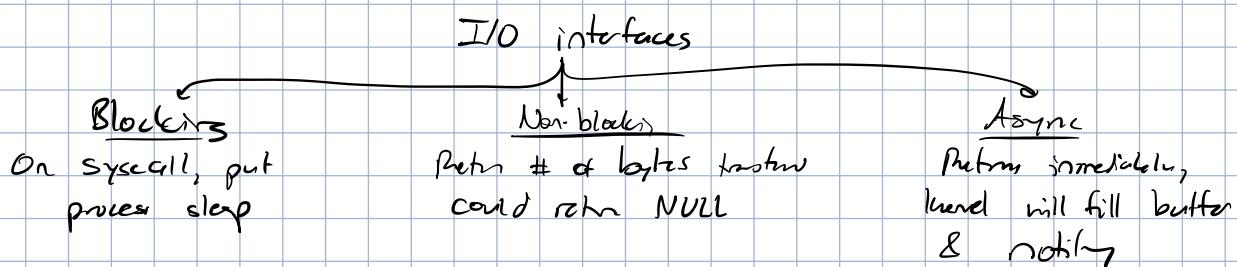
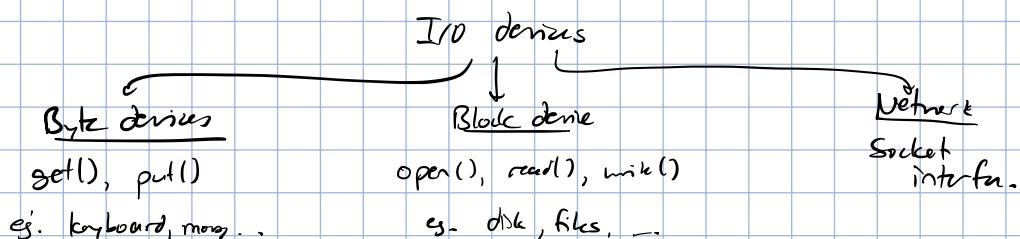
- ① Detect
- ② Swap out entre processes.

Working set: set of pages accessed freq. in past faults

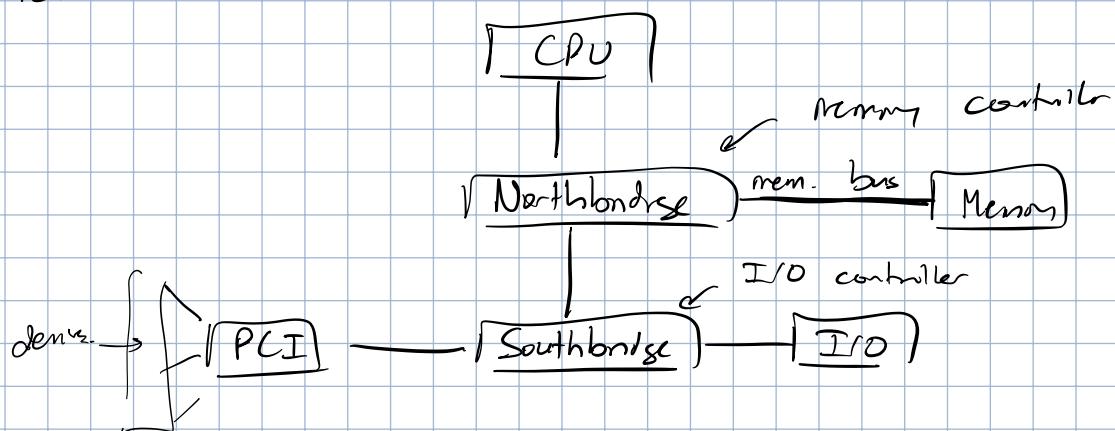
↳ Big entr set into memory  $\rightarrow$  hit rate  $\uparrow$

↳ If too many  $\rightarrow$  thrashing likely so swap out procs.

## I/O



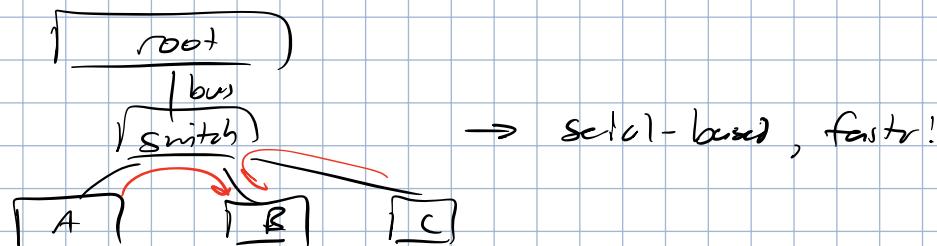
Architecture:



# PCI & Dennis

Originally a parallel bus  $\rightarrow$  limited by speed of slowest dev.

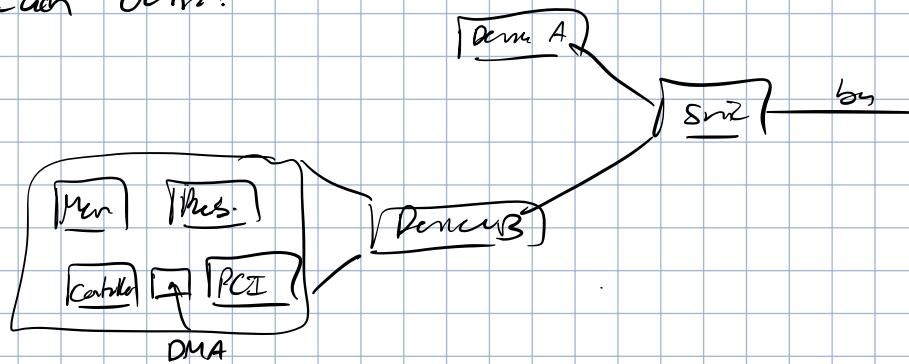
PCIe:



$\rightarrow$  serial-based, faster!

Good backward compatibility

Each dev:



On BIOS, dev will tell OS & auto-reg.

OS  $\leftrightarrow$  I/O acc

Port-mapped

I/O dev has a sep. addr. space.

CPU vms special init.

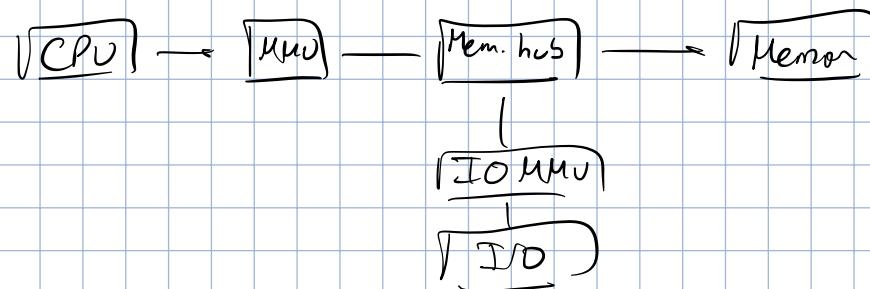
Memory mapped

Port of physical mem  $\rightarrow$  I/O dev

OS can use memory ops.  
needs protection!

We need I/O to interact memory

$\hookrightarrow$  For security  $\rightarrow$  address translation for I/O



2 problems:

① I/O can bypass IOMMU by caching

② JONMU not use for p2p dense coms.

Soln: access control service forces all data packets to go through JONMU

## I/O Performance

Peak bandwidth: max rate of data transfer

↳ depends on HLR

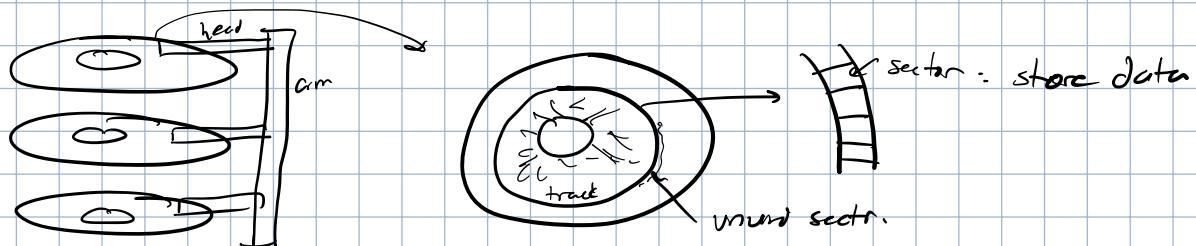
Effective bandwidth: actual rate of data trans.

↳ latency  $\uparrow \rightarrow$  size of data  $\uparrow$

$$\text{Latency } (n) = \text{overhead} + \frac{n}{\text{peak bandwidth}}$$

$$\text{Effective bandwidth } (n) = \frac{n}{\text{Latency } (n)}$$

## Disks



Outer half of disk is used

Process of reading data:

- ① Seek track  $\rightarrow$  move head (seek time)
- ② Rotate disk to locate sector (rotational latency)
- ③ Transfer blocks b/w sector & head (transfer time)

Eg. 11 Seek time is 5ms. Disk rotates at 7200 RPM. Transfer rate of 4 MiB/s w/ sector size of 1 MiB.

a) How long to read sector from random sector?

① Seek time  $\rightarrow$  5ms

② Rotational latency:

$$60000 \frac{\text{ms}}{\text{minute}} \times \frac{\text{minute}}{7200 \text{ rotations}} = 8 \text{ ms/rotation}$$

③ Transfer time:

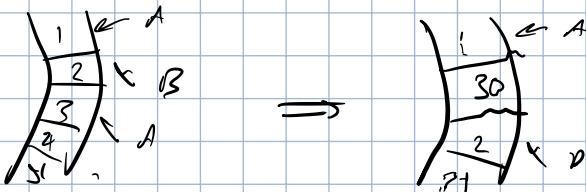
$$2^{10} \text{ B} \times \frac{\text{s}}{2^{22} \text{ B}} = 0.24 \text{ ms/B}$$

$$\text{Add: } S_{\text{ms}} \rightarrow 97_{\text{ms}} + 0,24 = \dots$$

Sequential reads are  $\mathcal{O}(n^2)$  → cut out scale tire & no need to rotate

## Error correction:

- ① Sector spans: send complete data to spec sectors.
  - ② Slip spans: remap sectors  $\rightarrow$  sequential reads
  - ③ Track steering: offset sector # to allow for seq. reads



## Disk Scheduling

- ① FCFS

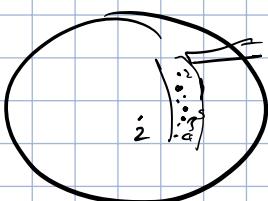
Pros: Simple

Cons: Poor performance  $\rightarrow$  keep switching tracks, create

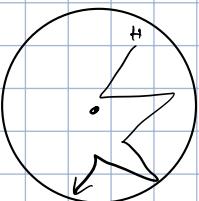
- ② SSTF: pick request w/ shortest seek time (closest to head)

Pros: scale on time

Cons: structure of task, far away from head



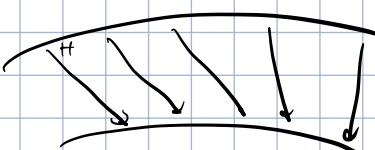
- ③ SCAN: arm moves in 1 dir  $\rightarrow$  spin the disk to pick closest request, move up & down in disk



Pros: low seek time & no start/stop

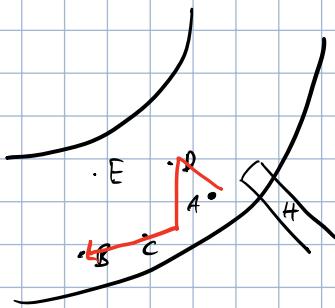
Cons: unfair to inner/outer regns

- ④ CSCAN: arm moves but switches to top



Much faster than SCAN

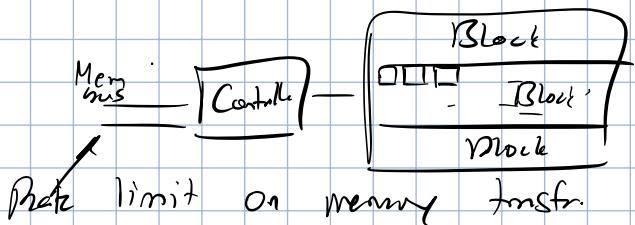
⑤ R-SCAN: take rotational delay



→ Move up & down to read request if < rotational delay to get to next req. in track

## Flash Memory

Better hard drive ← more efficient, no moving parts



Latency overall = queue time + controller + transfer byt.

Con: initing

- ↳ Can only write to empty pages. Clearing empty pages → clear what you
- ↳ Premap memory using flash translator layer (FTL)

Flash has garbage collector → keeps blocks together if stored @ same time.

## FILE SYSTEM

### Intro

Abstraction for persistent, named data

always stored unless deleted

address data by human names

Goals:

- ① Naming
- ② Data-storage management
- ③ Protection
- ④ Reliability/durable

File

Data

Metadata

- Bunch of blocks

Characteristics:

- ① Mostly small files, grows over time
- ② Most space taken by big files

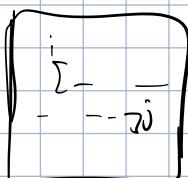
File access

Pattern:

- ① Sequential access: read bytes in order



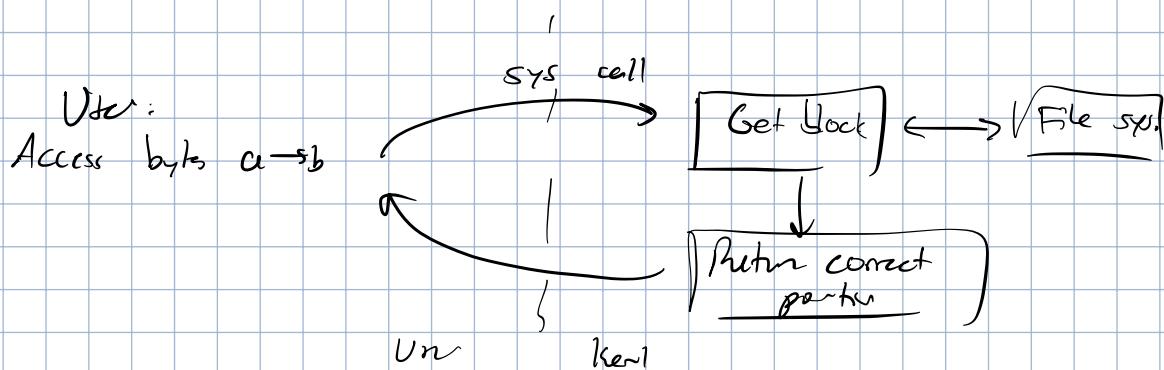
- ② Random access: can read anywhere in file



→ Should be fast

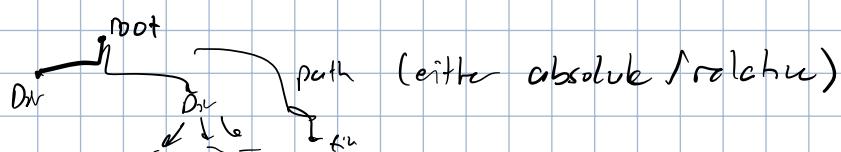
- ③ Content-based access

Eg: grep, filtering → need indexes



Directory

Defn: list of human-readable names  $\leftrightarrow$  files

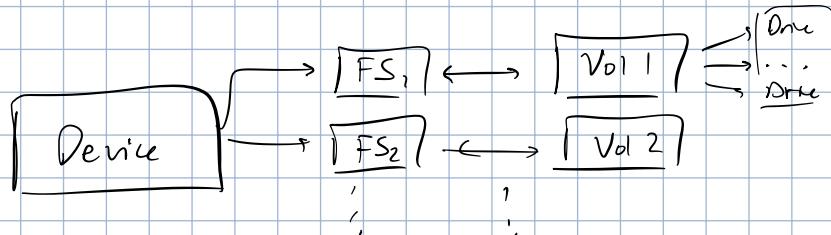


Can create links, which are alt. paths to file. File is not deleted if hard link exists.

1. Alternatively, use soft links which can point to nothing

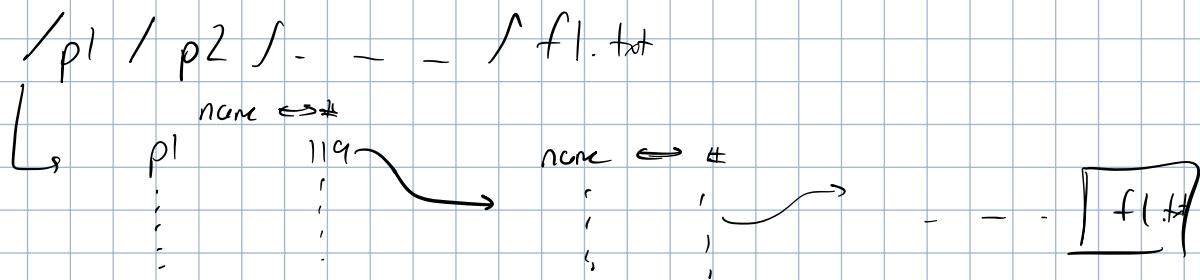
Stored like a file

## Volume



Device volume file systems will handle files by mounting to devia file system

## Search



## Implementation:

### ① Linked list



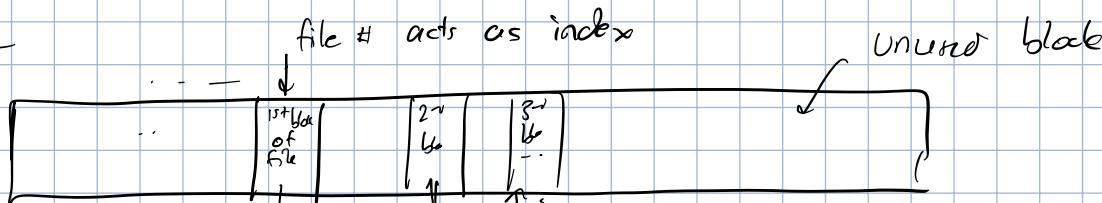
### ② B+ tree

Fast, blk logn search / trav

## Block Access

① Find file # → ② Get block to access (FAT table)

### FAT table



Q: Create a file?

Easy, just link copy block to end block

Q: Where is it stored?

In disk → bring to memory & cache on boot

Q: Format disk?

Just clear FAT table

Pros:

+ Extremely simple

Cons:

- Random access is slow
- Poor spatial locality: blocks can be all over FAT for 1 file
- No hard link support: no metadata stored
- Limited volume & file sizes
- Security issues