

OS - layer of software that provides sw access to hardware resources

- abstraction, resource management
- provide VM interface for apps

Program VM - supports execution of single program

System VM - execution & entire OS and apps

VM - gives program illusion it owns machine hardware has features you want

PS VM - each process thinks!

- owns all memory/ CPU time
- owns all devices

- all devices have same interface
- Fault Isolation
- Protection/ Portability

SyVMM - layers of OSes

- OS crash restricted to 1 VM

OS:

- manage sharing of resources, protection, isolation
- provide abstraction, easy to use
- abstraction of physical resources
- common services: networking, storage

4 OS Concepts:

1. Thread

- single unique execution context, fully describes program state
- PC, registers, execution flags, stack

2. Address Space

- programs execute in address space that is distinct from memory space of physical machine

3. Process

- instance of executing program
- consists of address space, ≥ 1 thread

4. Dual Mode Operation/Protection

- only system can access some resources
- controlling translation from program virt addrs to machine physical addrs protects user programs from each other, OS, hw

Thread:

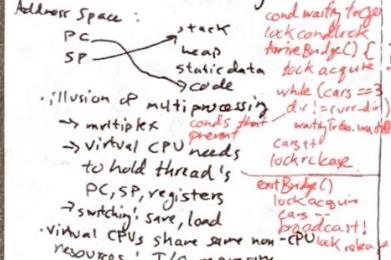
- register hold thread context:

- SP = top of stack

- when resident in processor registers, executing

- PC = addr of executing instr

Address Space:



Process: execution env w restricted rights

- uni memory, file sys context

- protected from each other

- memory protection

Threads: concurrency !!; add space, protection

Dual Mode:

- 1) Kernel (protected)

- before fork point do something, i.e. after fork they

- need state (which mode)

- some ops only allowed in Kernel /etc, open fd after fork

- User → Kernel: set sys mode, save PC, off + put aside user state

- Kernel → User: clear sys mode, restore PC

Kernel → User: clear sys mode, restore PC independently

Base + Bound

- translate address when program loads

- relocate loader

- basically shift program to execute somewhere in memory

Address Translation

Processor → translator → memory

- virtual addr → physical addr

- translate address on fly

- allocate some part of memory to program. Ensure program adds within

size of space allocated

Mode Transfer

Syscall

- ps registers system service
- doesn't have addrs
- marshall syscall ID and args in register

Interrupt

- external asynchronous event triggers context switch
- e.g. Timer/I/O device independent of user ps

Trap, Exception

- internal synchronous event in process triggers context switch

Unprivileged Control Transfer gets target address from interrupt vector

vector: #⇒ handler

Process Control Block (PCB)

- represents managed process
- status
- registers, SP
- PID
- execution time
- memory space, translation table

Safe Kernel Mode Transfers

- controlled transfer into kernel's CS/CSN table
- separate kernel stack
- syscall handler copies user args to kernel space before invoking specific function

Interrupt Control

- interrupt not visible to user process
- no change to process state
- kernel may enable/disable interrupts
- 3 non-markable interrupts
- cannot disable

Safe Interrupts

- handler works regardless of user code state
- atomic transfer of control

Process creates a process

- fork(): - create copy of current process w new PID

- > O (return value)

- main is parent process
- RV = PID of child

- < O

- running in child
- waitpid(): pid, status, options

- error!

- running in dg PS

State of original process duplicated

exec(): change program being run by current ps

Wait(): wait for ps to finish

Signal(): sys call to send notification

to another process

Shell: job control system

User level equivalent of interrupt

= signal

I/O Design Concepts

- Uniformity

- file I/Os, device I/O, interps common thru open, read/write, close

- Open Before Use

- access control

- Byte oriented

- addressable in bytes

- Kernel Buffer Reads

- streaming + block devices look the same

- read blocks process

Kernel buffered writes

- completion of outgoing transfer decoupled from application, allowing it to continue
- explicit close

I/O, storage

High Level I/O streams

Low Level I/O handles FILE DESCRIPTOR# (cont.)

Syscall registers FILE DESCRIPTOR#

File System descriptors

I/O Driver commands/data transfers disks, flash, controllers

File System

- files live in hierarchical namespace of filenames
- File: data, metadata (info about file)
- Directory: contains files, directories

C File API: streams (sequence of bytes)

- FILE: stdIn, stdOut, stdErr
- File Descriptor: OS object representing State of file

- user has "handle" on descriptor
- LOW LEVEL I/O

APIs

Lower Level API

- associated w HW device
- registers/unregisters w kernel
- handler functions for each of file ops

Device Driver

- device specific code in kernel that interacts w device
- standard internal interface

2 pieces:

Top half: accessed by syscalls

- kernel/interface w device driver
- starts I/O to device

Bottom half: run as interrupt

- gets input or transfers next block of output
- may wake sleeping threads

Communication like File I/O

- connected queues over internet
- e.g. client write → buffer → server read

- client write → buffer → server read

Socket: abstraction of network I/O queue

- emulates one side of communication channel

Data transfer like files → fd

Ports:

- 0-1023 system ports
- 65536 privilages

- 1024-49151 registered ports

- 49152-65535 dynamic/private ports

Client: connect() → Server Socket

accept() ← ↓

→ Connection new socket

\* new socket for each unique connection

Server has socket that listens for connection

Upon request:

Listen socket → connection socket

- - - - - - - - - - Fork() - - - - - - - - - -

- close listen socket
- parent socket

- close connection socket

- returns to listening

\* only one thread can access shared vars

context Switch - CPU switch from PS A to PS B

- overhead: min practical switching time

Process States:

- new, ready, running, waiting, timing, admitted

- interrupt exit

- ready → running

- running → term

- I/O event waiting

- I/O event wait

PCBs more queues as they change state

Thread State

- shared by all in process

- content of memory

- I/O state (FD, network conn)

- private to each

- kept in TCB

- CPU registers (PC)

- execution stack

- parameters, temp vars, return PCs

Shared: Heap, global vars, code

Per-Thread: TCB, stack, stack info, saved registers, stack metadata

Dispatcher Loop of OS;

Run Thread()

ChooseNextThread()

SaveStateOfCPU (currTCB)

LoadStateOfCPU (newTCB)

Dispatcher gets control back

- internal events, e.g. blocking on I/O, waiting on signal from other thread, yield()

Context switching mostly depends on cache limits and process/thread

hunger for memory

External Events

- interrupts: stop running code to jump to kernel
- timer

- ensure that dispatcher can regain control

- looks like thread runs synchronously but NOT

- OS keeps track of TCBs in Kernel memory

Thread Fork()

- create new thread and put on ready queue

- 1) sanity check args

- 2) enter kernel mode and check args

- 3) allocate new stack + TCB

- 4) init TCB and put on ready list

Registers:

SP → stack

PC register → Thread Root()

AL, AL → funcPtr, funcArgPtr

Thread Root()

DoStartupHousekeeping()

UserModeSwitch()

Call funcPtr(funcArgPtr)

ThreadFinish()

Multithreaded Processes

PCB → PCB → PCB

↓ TCB → TCB → TCB

TCB → TCB → PCB

①

Process Preemption (REMEMBER) program can exit before context switch

- Switch HIGH
  - CPU state low
  - Memory/I/O HIGH
- Creation HIGH
- Protection
  - CPU ✓
  - Memory/I/O ✓
- HIGH sharing overhead: lock release, sema-up

### Threads

- Switch medium
  - CPU state low
- Creation medium
- Protection
  - CPU ✓
  - Memory/I/O NO
- Sharing low

Kernel threads are expensive  
unlock() → locked = 0

- need to go into kernel mode to schedule

Lighter Weight Option: Userlevel Threads

→ user program provides scheduler, thread phys

→ user threads scheduled non-preemptively  
- cannot interrupt each other

Downside: when one thread blocks on I/O, they all do  
- kernel cannot adjust scheduling

Options: Scheduler activation

→ Have kernel inform user level when thread blocks

Threading Models: (User to Kernel)

- 1-1, many-1, 1-many

Multi-core: run multiple threads each on a different core simultaneously

Multiprocessing: multiple CPUs or cores or hyperthreads

Multiprogramming: multiple jobs or processes

Multithreading: multiple threads per process

Why cooperating threads:

- share resources
- speedup
- modularity

Cooperating threads: shared state between multiple threads

- non deterministic, reproducible

ThreadPool - bounded pool of worker threads

- ensures throughput

\* threads yield overlapped I/O, computation without having to deconstruct code into non-blocking fragments

Atomic Operation - always runs to completion or not at all

- indivisible

- memory references and assignments of words (load, store) are atomic

Synchronization: using atomic operation to ensure cooperation between threads

Mutual Exclusion: ensuring that only one thread does a particular thing at a time

Critical Section: piece of code that only one thread can execute at once

Lock: prevent someone from doing something, e.g. lock before entering critical section and accessing shared data

starvation: thinking other has shared resource when neither does  
solution should protect critical section!  
busy-waiting: consuming CPU while waiting  
- Can avoid context switching by
 

- avoiding internal events
- preventing external events by disabling interrupts

Better implementation of lock:
 

- maintain a lock variable and impose mutual exclusion only during ops on that var.

int value = FREE

Acquire(C) {
 

- enable interrupts
- disable interrupts;

if value = BUSY:
 

- put thread on wait queue
- Go to sleep
- /Enable interrupts?

else:
 

- value = BUSY
- enable interrupts

Release(C) {
 

- enable interrupts
- disable interrupts;

if (anyone on wait queue):
 

- take thread off wait queue,
- place on ready queue;

else:
 

- value = FREE
- enable interrupts;

}

- need to enable interrupts after going to sleep

Alternative: atomic instruction sequence
 

- read/unlock value atomically

Test&Set using value

Lock is free, test&set reads 0 and sets value = 1

Busy: 1, Free: 0

- returns 0 when we complete

If lock is busy, test&set reads 1 and sets value = 1 → while loop continues

↳ BUSY WAITING ⇒ priority inversion

- busy-waiting thread has higher priority than thread holding lock

↳ use while loop for wait

↳ this is to account for cases

in which another thread runs before it and wakes condition

thread 1 was waiting on

Readers/Writers Problem

e.g. database

- readers never modify

- writers read + modify

Want many readers

1 writer

→ give priority to writers

In this example, hold a lock when changing AW, WW, AR, WR

• these vars and checks on them inherently ensure protection

Goal of Scheduling

- minimize response time

- maximize throughput

- fairness

First Come First Serve

↳ thread keeps CPU until it yields

Wait time: amount time job

Completion: amount time job

finishes till

+ simple

- short jobs stuck behind long ones

exec v does not destroy fd, coded happens to fd

Thread Join sema, PC

Finish signals Join

use separate semaphore for each constraint  
Problem: semaphores are dual purpose  
⇒ Locks for mutual exclusion, condition vars for scheduling constraints

Monitor: lock + 20 cond vars for managing concurrent access to shared data

Condition Var: queue of threads waiting for something in a critical section

Key: possible to sleep inside critical section by atomically releasing lock when we sleep  
- cannot do this w sema

Wait(&lock): atomically release lock and sleep  
- reacquire lock later before returning

Signal(&lock): wake up one waiter  
Broadcast: wake up all waiters  
Must hold lock when doing cond var ops!

Construction of lock: --
 

- Separate lock var, use HW mechanisms to protect modification of var

Hedge Monitor: use if

- signaler gives lock, CPU to waiter, waiter runs immediately

- waiter gives up lock, protector back to signaler when it exits critical section or waits

Mesa Monitor: most

- signaler keeps lock to processor

- waiter puts on ready queue

- need to check condition again after wait

↳ use while loop for wait

↳ this is to account for cases

in which another thread runs before it and wakes condition

thread 1 was waiting on

Readers/Writers Problem

e.g. database

- readers never modify

- writers read + modify

Want many readers

1 writer

→ give priority to writers

In this example, hold a lock when changing AW, WW, AR, WR

• these vars and checks on them inherently ensure protection

Goal of Scheduling

- minimize response time

- maximize throughput

- fairness

Round Robin

- each process gets unit of CPU time (10-100 ms); 2 quantum

- after time expires, process preempted

+ put back at q

- no process waits > (n-1) q

+ better for short jobs, fair

- context-switching adds up for long jobs

Today, time slice = 10-100 ms  
context switching = 0.1-1 ms  
~ 1% overhead

Strict Priority Scheduling

- always execute highest priority runnable jobs to completion
- each queue precessed RR

Problems:

- starvation of lower priority tasks
- Deadlock: priority inversion when no priority donation

→ FAIRNESS

- gained by limiting avg response time
- could give each queue a fraction of CPU
- could increase priority of jobs that don't get service

Lottery Scheduling:

- give each job some tickets
- at each time, pick a ticket
- CPU time proportional to # tickets given to each job
- to approximate SRTF
  - short running jobs get more, long ones less
  - behaves gracefully as load changes

To evaluate a scheduling algorithm:

- Deterministic modeling
  - compute performance of alg on predetermined workload
- Queueing Model
  - mathematical approach for handling stochastic workload
- Implementation / Simulation
  - build system to run algorithms against it

Mix of different apps

- Burst time used to figure out whether app is interactive or high throughput
- sleep/lat/short bursts ⇒ interactive
- less bursts ⇒ not interactive

Minor-Best FCFS - Shortest Remaining Job First

SJF! no job w/ least computation

SRJF! preemptive version - if job arrives w/ shorter completion time than running job, preempt
 

- ↳ get short jobs out of system
- may lead to starvation of longer jobs

Make rule: round-robin error task-one target echo recipe task-two task-one dependency echo 2

strace -p pid syscalls

Trace - p pid library calls

- thread creates: creates + immediately starts new thread
- thread, buffer, start, run to stop
- thread join: waits for thread to stop

- execve: current offset + offset
- END: size of file + offset
- dup2: create fd, close old
- dup2: create fd, close old
- signal: signal handler
- bind: socket\_fd, address

- socket: create
- listen: accept
- bind: socket\_fd, address

- send: write
- sendto: write
- sendfrom: write
- sendmsg: write
- connect: connect
- accept: read
- read: read
- readv: read
- readfrom: read
- readmsg: read

- SRTF is not realistic  
 - can predict length of next CPU burst (adaptive)  
 → programs have predictable behavior

### Multilevel Feedback Scheduling

- multiple queues, each w/ own algorithm/CPU time
- jobs go between queues depending on if timeout expires (no bound stays near top)
- serve all from highest priority queue, ↓ real time scheduling

- enforcing predictability of worst case response times
- Hard Real Time (meet all deadlines)
- Soft Real Time (high probability)

Workload: task (deadline, known, computation)

Earliest Deadline First: schedule active tasks w/ closest absolute deadline

response time  $\rightarrow$  utilization  $\rightarrow$  100%.

Deadlock: Mutual Exclusion

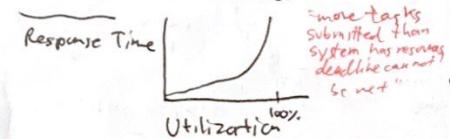
Hold + Wait  
No Preemption

Circular Wait

- detect deadlock, then! (Recovery)
- terminate thread, force resource give up
- preempt resources don't kill thread
- reused actions of deadlocked threads

### Preventing Deadlock:

- infinite resources
- no sharing of resources
- don't allow waiting (retry)
- make all threads request everything they need at beginning
- force all threads to request resources in particular order



Working state of process is defined by data in memory

Memory Multiplexing: protection, controlled overlap TRANSLATION

Locality = locality of instructions, data to memory

Addresses can be bound at compile, link/load execution time

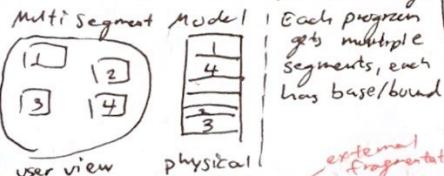
Uniprogramming - no translation/protection - gets entire physical memory

Base/bound: add base to every address during execution (dynamic address translation)

- error if address > bound

- program gets continuous region of memory

Issue: memory gets fragmented  
 no support for sparse address space  
 hard to do interprocess sharing



Segment Map: Seg #  $\Rightarrow$  base, limit, valid  
 - base added to offset  
 as many chunks of physical memory as entries in segment map

e.g. 4 segments

| Seg | Offset | Add base to offset |
|-----|--------|--------------------|
| 1   | 1413   | 0                  |
| 2   | 0      | 0x1000             |
| 3   | 0      | 0x1400             |
| 4   | 0      | 0x1000             |

| Seg ID | Base   | Limit  |
|--------|--------|--------|
| 0      | 0x4000 | 0x0800 |
| 1      | 0x4000 | 0x1400 |
| 2      | 0xF000 | 0x1000 |
| 3      | 0x0000 | 0x3000 |

check offset w/ limit

Virtual Addresses stored in registers

virt addr space has holes - efficient segmentation for sparse addr space

stack/heap grow by addressing outside valid range

need r/w bits in segment table

segment table stored in CPU on context switch

problems: must fit varying size chunks into physical memory

may move up multiple times to fit - limited options for swapping to disk

wasted space (fragmentation)  
 - external: between chunks  
 - internal: within chunks

Solution to Fragmentation:  
Pages - fix sized chunks of physical memory

- pages are small

- use vector of bits to handle allocation: 1  $\Rightarrow$  allocated, 0  $\Rightarrow$  free

One Page Table per process  
 - in physical memory  
 - contains physical page #, permission bits

Virtual Addr [Virt Pg #] offset

2 offset bits = page size  
 e.g. 10 offset bits, page size = 1024

Check Virtual Pg # against page table bounds / permission

Virt Pg #  $\Rightarrow$  (Phys. Pg #, permission)  
 offset - where within page?

Sharing:  
 - page tables for different processes can point to the same physical pg

- if stuff grows in virtual addr space, allocate new pages in physical memory where room

Fix for Sparse Address Space:  
 - two level page table

Virt Addr: [Virt P1 Index | Virt P2 Index | Offset]

PT 1  
 Virt P1  $\rightarrow$  which 2nd level PT  
 Use Virt P2 index in 2nd level PT

Segments = Pages

| Seg # | Pg # | Offset |
|-------|------|--------|
|       |      |        |

Seg  $\Rightarrow$  Pg Table

On context switch, contents of top-level segment register need to be switched

Sharing: Segments point to same Page Table

Page Table Entry

Bits

31-12: Physical Pg #

#11-9: Fine

L: 6-1  $\Rightarrow$  4 MB page

0: Dirty

A: Accessed (recently)

PCD: Page cache disabled

PLW: Page write transparent

U: User accessible

W: writeable

P: present (valid)

Hardware/software translation

HW: fast, inflexible

SW: flexible, every translation must make fault

CPU  $\rightarrow$  TLB  $\rightarrow$  Cache  $\rightarrow$  Memory

Virt  $\rightarrow$  Phys Addr

Phys Addr  $\rightarrow$  Data

Parallellizing TLB, cache access

Virt Addr [VPN | Offset]

Phys Addr [PPN | Tag | Index | Byte]

- check TLB w/ VPN

$\rightarrow$  returns PPN/Tag

- index into cache w/ Index

$\rightarrow$  returns PPN/Tag

- check if return the same

Virtual Address

TLB

Page Table = cached mapping

$\downarrow$  via phys from Virt  $\rightarrow$  phys

Physical Address

$\downarrow$  cache

Physical Memory - address = physical address

What happens when: cannot translate to physical address?

Working set: subset of address space program uses to run

Model of Locality: Zipf

- long tail  
 $\rightarrow$  substantial value from tiny cache  
 $\rightarrow$  substantial misses from very large cache

Demand Paging

- program spends 90% of time in 10% of memory  
 - waste if put all of code in memory

solution: use memory as cache for disk

Transparent Level of Indirection!

- flexible placement of user data

- 1 page block size  
 - fully associative - expensive conflict miss

- find page: TSB, page table traversal  
 - on miss, go to disk to grab page  
 - write back

Handling Page Fault: No physical page for 1) reference missing page = invalid 2) temp (not in page table)

3) OS grabs page from disk 4) bring page into phys memory 5) reset page table 6) restart instruction

PTE helps w/ demand paging  
 valid: page in memory, PTE points at physical page

invalid: page not in memory, use info in PTE to find on disk

Page Fault Handling:  
 - replace old page  
 - if dirty write contents to disk  
 - change PTE and cached TLB to invalid

"invalid" means not in memory

Effective Access Time:

- Hit Rate  $\cdot$  Hit Time + Miss Rate  $\cdot$  Miss Time  
 $\rightarrow$  Hit Time + Miss Rate  $\cdot$  Miss Penalty

Replacement Policies for Demand Paging:

- FIFO: minimum (replace page that will be used for longest)  
 - random  
 - LRU

LRU performs badly in cyclic references

Belady's anomaly: adding memory doesn't always help fault rate  
 E.g. FIFO: ABC DABEAB CDE

Implementing LRU  
 - Ideally? timestamp each reference  
 order by timestamp  
 - expensive

• Clock Algorithm  
 - set/use bit on reference  
 - advance clock hand

Valid  
 - advance clock hand  
 - clear and advance  
 - or replace

Valid

Page Fault:  
 - advance clock hand  
 - clear and advance  
 - or replace

Valid

$N^{\text{th}}$  clock  
 - give page  $N$  chances  
 - counter per page: # sweeps  
 - 1: clear use, counter  
 - 0: increment counter; if  $\geq N$ , replace  
 For dirty pages: give extra chance before replacing

- Clean: use  $N=1$   
 - Dirty:  $N=2$  (WR on  $N=1$ )  
 Second Chance List, FIFO LRU  
 - memory in 2: Active List (RW), SC List (invalid)  
 - access pages in active list at full speed  
 - on page fault, access SC list  
 Alternative to SC list: free list  
 - filled in background by clock algo  
 - dirty pages carry back to disk when enter list  
 - faster for page fault

Page Frame Allocation  
 Global: process selects replacement from set of all frames  
 Local: selects from only its frames

Types of Allocation:  
 Equal: every process gets same amount  
 Proportional: allocate according to size of PS  
 Priority: proportional according to priority - take lower priority process' frames on page fault

PS not enough pages  $\Rightarrow$  high page fault rate  
 - low CPU utilization  
 - OS spends lots of time swapping to disk

Thrashing: PS busy swapping pages in + out  
 Core Mapping: physical page  $\rightarrow$  virtual page  
 - when you push physical page out of page table to disk, must invalidate all PTEs

## I/O

interact w devices using I/O controllers  
 r/w to I/O registers like they're memory

Byte or Block data transfer

sequential vs random access

Transfer Mechanism: I/O and DMA

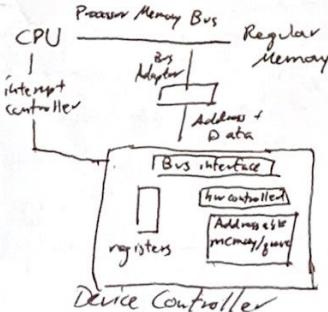
Uniform interface despite different classes  
 - block devices: disk drives, DVD-RW  
 - character devices (byte): keyboards, mouse  
 - network devices: socket interface

## Timing:

- blocking interface (wait)  
 - sleep process till data ready  
 - non-blocking (don't wait)  
 - immediately return from r/w request w # bytes successfully transferred  
 - asynchronous (tell me later)  
 - get pointer to user's buffer  
 - take data and notify when done

CPU interacts w controller

- controller has registers that can be read or written  
 - may have memory for request queues



Processor Accesses Registers:

I/O instr: in/out instr  
 Memory Mapped I/O: word/byte  
 - device registers/memory in physical/addrs space

Transferring Data To/From Controller:

- Programmed I/O:  
 - use instruction + easy, simple hardware  
 - consumes processor cycles proportional to data size

- Direct Memory Access  
 - give controller access to memory bus  
 - transfer directly

$\hookrightarrow$  device driver  $\rightarrow$  disk controller  
 - initiate DMA transfer  
 - DMA controller transfers byto buffer

- DMA interrupt CPU to signal completion  
 - interrupt (high overhead)  
 - Polling: OS checks specific device status + low overhead  
 - wastes cycles

Device Driver in Kernel  
 - device-specific code that interacts w device hardware

Top half: Kernel's interface to device

Bottom:  
 - does actual transfer

$\hookrightarrow$  wake sleeping threads on completion

Latency: time to perform an operation

Bandwidth/Throughput: rate at which ops are performed (ops/s)

Overhead: time to initiate an operation

Latency( $n$ ) = Overhead +  $\frac{n}{\text{Bandwidth}}$

Bus speed, device transfer bandwidth determine peak bandwidth for I/O

Storage Devices

Magnetic Disk:  
 - transfer track (cyl) = sector

R/W:

seek time: position arm over track  
 rotational latency: wait for sector to rotate under r/w head

transfer time: transfer block after

Disk Latency = seek time + controller time +  $\frac{1}{\text{bus speed}}$

request  $\rightarrow$  SW  $\rightarrow$  HW  $\rightarrow$  media  $\frac{1}{\text{bus speed}}$

drive controller time  $\frac{1}{\text{bus speed}}$

Key to using disk effectively is minimizing seek and rotation delays

3 error correcting codes, bad sector remapping  
 - rarely corrupted  $\rightarrow$  large capacity, cheap access

$\rightarrow$  slow for random  $\rightarrow$  better for sequential

Flash Memory

$\rightarrow$  capacity at intermediate os +

$\rightarrow$  block level random access (both)

$\rightarrow$  good for reads, slower for writes

$\rightarrow$  erases sequentially in large blocks

$\rightarrow$  wear patterns issue  
 - no worthy parts (no seek/retention delay)

read is fast, write slower

bc: can only write empty pages in a block

$\rightarrow$  may need to erase a block ...

Untested reads, erasure 10x writes

+ low latency, high throughput

+ lightweight, low power, silent

+ read & memory speeds

- small storage, expensive

- asymmetric block write performance

Effective BW per Op:

$\frac{S}{B} + \frac{n}{B}$ ,  $n = \# \text{ ops}$ ,  
 S = fixed overhead  
 B = time per op  
 transfer size  
 response time

Model burstiness of arrival  
 using exponential distribution  
 $f(x) = \lambda e^{-\lambda x}$



Likelihood of an event occurring is independent of how long we've been waiting  
 Poisson - completely random

Queuing Theory:

Little's Law:  
 in any steady system,  
 average arrival rate = average departure rate  
 $\# \text{ jobs} = \lambda (\text{jobs/s}) \times L(\text{s})$   
 during the  $\rightarrow$  response time

$\lambda$ : mean number of arrivals/s

Tser: mean time to serve

C: standard coefficient of variance

$\mu$ : service rate =  $1/T_{\text{ser}}$

$\zeta$ : server utilization =  $\lambda/\mu = \lambda \cdot T_{\text{ser}}$

$T_g$ : Time spent in queue

$L_g$ : length of queue =  $\lambda \cdot T_g$

Memory latency time distn ( $C=1$ )

$T_g = T_{\text{ser}} \cdot \alpha/(1-\alpha)$

General Service Time Distr

$T_g = T_{\text{ser}} \cdot \frac{1}{2}(1+C) \cdot \alpha/(1-\alpha)$

lock outside while loops that checks condition

- MUTEX: protect access to shared variables

make scheduler fair!

downgrade when thread uses up quota, ~~dequeue~~

upgrade when it yields or gets blocked

Page Table

# entries determined by size of virtual address + size of page

each entry consists of PPN and valid bits

e.g. virt mem size doubles

# entries doubles

possible - exchanging data in memory for data in disk

Banker's J

- look at remaining resources & check against what each thread needs

phys mem larger than aggregate working set sizes  $\Rightarrow$  thrashing

memory isn't actually allocated until you begin using it

Memory Protection - use base/bound for each process

In Pintos, all processes share a single interrupt

Banker's: unsafe state doesn't necessarily imply deadlock  
 - evaluate each request  
 - grant if some ordering of grants makes it still deadlock-free afterwards

Remember! killing something will leave system in inconsistent state, e.g. won't use "semaphore". Given 64 bit addresses space + 12 bit offset, we have  $2^{64}/2^{12} = 2^{52}$  pages.

Multi-level page table will NOT always take less storage space than single level  
 e.g. all pages in virtual memory allocated

In  $\# n$ -way set associativity, num index bits decreases by  $\log_2 n$ . compared to direct mapped.

1 KB Cache w 32 B blocks DM: tag index byte select

$2^10 \times 2^5 \times 2^3$

2-Way tag index byte select

23 4 5

Full tag byte select

27 5

- copy or write in stale pages as readonly  
 - process try to write  $\rightarrow$  fault  $\rightarrow$  OS creates writable copy in new physical frame  $\rightarrow$  new TLB

# VPNs = # entries

possible for FCFS to have lower turnaround time than RR, e.g. jobs arrive shortest to longest

size of PTE determined by how many physical pages we have

0x0442F - 18 bit addr

first 4 bits are 0001 so we look at index 1 in 0x3000 PTE which is 0x4203. Last 2 bits are 11 so it is valid and a directory.

Top 7 bits are taken C bc # bits of PPN in PTE and we don't offset to go to address 0x4200.

Next 5 bits are 00010 so we look at index 2, which is 0x5601. Last bit is valid, so we take top 7 bits of 0x5601 and add that to offset 0x02F which gets us 0x562F.

- virtual address tracked into register

(4)

|  |  |  |
|--|--|--|
| <p>I/O</p> <ul style="list-style-type: none"> <li>- controlled, managed by OS</li> <li>- logic to interact w/ device in I/O controller</li> </ul> <p>Data granularity: block, byte<br/>higher throughput</p> <p>Access Pattern: Sequential vs Random<br/>How to Transfer: Programmed I/O, DMA</p> <p>I/O: Uniform interface for devices</p> <p>Block Devices:</p> <ul style="list-style-type: none"> <li>- open(), read(), write(), seek()</li> </ul> <p>Char:</p> <ul style="list-style-type: none"> <li>- getc(), putc()</li> <li>- read/write to <del>shared</del> file that represents connection</li> </ul> <p>Network:</p> <ul style="list-style-type: none"> <li>- socket interface - socket(), bind(), listen(), select()</li> <li>- select: not block on one socket, wait on all sockets</li> </ul> <p>Dealing w/ Timing</p> <ol style="list-style-type: none"> <li>1) Blocking - wait: puts to sleep till ready</li> <li>2) Non-blocking - return quickly</li> <li>3) Asynchronous - tell me later</li> </ol> <ul style="list-style-type: none"> <li>- kernel fills buffer, notifies user when ready. It's called when op complete</li> </ul> <p>I/O devices connected by controller</p> <ul style="list-style-type: none"> <li>- connected by bus: data moves in/out of PCI bus, connected to CPU by bridge/memory controller</li> <li>- all controllers connected to PCI bus</li> </ul> <p>Processor/Device Interaction</p> <ul style="list-style-type: none"> <li>- inside device controller are hw controllers, memory, bus interface, registers</li> <li>- gets commands, sends info</li> </ul> <p>2 ways to interact + w/ device</p> <ol style="list-style-type: none"> <li>1. I/O instruction       <ul style="list-style-type: none"> <li>- particular register to put instruction for I/O device</li> </ul> </li> <li>2. Memory Mapped I/O       <ul style="list-style-type: none"> <li>- regions of memory reserved for I/O devices</li> </ul> </li> </ol> <p>Transfer To/From Controller (Data)</p> <p>Programmed I/O</p> <ul style="list-style-type: none"> <li>- each byte transferred via processor</li> <li>- Direct Memory Access       <ul style="list-style-type: none"> <li>- give controller access to memory bus</li> <li>- CPU tells I/O controller, DMA how much data, where to store</li> </ul> </li> </ul> <p>Notifying OS</p> <ul style="list-style-type: none"> <li>- interrupts: high overhead, good for unexpected events</li> <li>- polling: OS checks register periodically + low overhead       <ul style="list-style-type: none"> <li>- may waste many cycles on polling if infrequent I/O</li> </ul> </li> </ul> <p>Latency - time to perform operation</p> <p>Throughput/bandwidth - rate at which ops are performed</p> <p>Overhead/Setup - time to initiate an operation</p> <p>Latency: <math>S + \frac{B}{B}</math>, <math>S =</math> startup cost, <math>B =</math> bandwidth</p> <p>Bandwidth: <math>\frac{n}{(C + \frac{n}{B})}</math></p> <p>Peak BW for I/O: (bottleneck)</p> <ul style="list-style-type: none"> <li>- bus speed</li> <li>- device transfer bandwidth</li> </ul> <p>Magnetic Disk:</p> <ul style="list-style-type: none"> <li>- sector/unit of transfer</li> <li>- platter, track, cylinder, head, arm</li> </ul> <p>Seek (position head over proper track)</p> <p>Rotational Transfer</p> | <p>Disk latency:</p> <ul style="list-style-type: none"> <li>- access time + controller time + seek time + rotation time + transfer time</li> </ul> <p>PS: cache mem, fd, filters context</p> <p>Threads: TCB, stack, stackify, cronjobs, stack metadata</p> <p>Shared: heap/global vars, code</p> <p>- minimize seek + rotation delays</p> <p>SSD: no moving parts</p> <p>Queuing &amp; Controller Transfer</p> <p>Segmental, random reads take same time</p> <p>Writing = 10x reads</p> <p>&amp; only write on empty pages</p> <p>Erase = 10x writes</p> <p>Effective BW per Op:</p> <p>transfers/<br/>response time</p> <p>Queuing leads to big increases in latency as utilization increases</p> <p>Service Rate <math>\mu = 1/T_s</math> - the to arrival Rate <math>\lambda = 1/T_a</math></p> <p>Utilization <math>U = \lambda/\mu, \lambda &lt; \mu</math></p> <p>Model Bestness:</p> <ul style="list-style-type: none"> <li>- exponential distribution</li> <li><math>F(x) = 1 - e^{-\lambda x}</math></li> <li><math>\lambda = \text{average arrival rate}</math></li> </ul> <p>internal block arrivals</p> <ul style="list-style-type: none"> <li>- lots of short arrival intervals, few long gaps</li> </ul> <p>Variable Time <math>T</math></p> <p>mean <math>m = \sum p(T) \cdot T</math></p> <p><math>\sigma^2(\text{variance}) = \sum p(T) \cdot (T - m)^2</math></p> <p><math>C = \sigma^2/m^2 = \sum p(T) \cdot T^2 - m^2</math></p> <p><math>C=0 \Rightarrow</math> no variance</p> <p><math>C=1 \Rightarrow</math> memoryless, Poisson completely random</p> <p>Little's Law -</p> <ul style="list-style-type: none"> <li>- in static system:</li> <li>- average arrival rate = avg departure rate</li> </ul> <p><math>N(\text{jobs}) = \lambda \cdot L</math> (queueing time)</p> <p># jobs in system / arrival time</p> <p>throughput</p> <p>arrival rate, how many jobs in sys</p> <p>stable distribution doesn't change over time</p> <p><math>\lambda</math>: arrival rate</p> <p>Tser: mean time to serve a customer</p> <p><math>C = \sigma^2/m^2</math>: standard deviation variance</p> <p><math>\mu = 1/T_{\text{ser}}</math>: service rate</p> <p><math>\omega</math>: server utilization: <math>\lambda/\mu = \lambda \cdot T_{\text{ser}}</math></p> <p>decrease arrival rate <math>L = T_g + T_{\text{ser}}</math></p> <p>Compute: <math>T_{\text{ser}} = \frac{T_{\text{sys}}}{L}</math></p> <p><math>T_g</math>: time spent in queue</p> <p><math>L_g = \lambda \cdot T_g</math> (length of queue by Little's law)</p> <p>M/M/1 queue</p> <p>Poisson arrival, service, losses</p> <p><math>C = 1</math> (memoryless service times)</p> <p><math>T_g = \frac{T_{\text{ser}} \cdot u}{1-u} = \frac{\lambda}{\mu - \lambda}</math></p> <p>General service time:</p> <p><math>1/\mu</math> or <math>1/B</math></p> <p><math>T_g = T_{\text{ser}} \cdot \frac{u}{1-u} (1+C) \cdot \frac{u}{1-u}</math></p> <p>Disk Performance test during sequential reads</p> <ul style="list-style-type: none"> <li>- reorder work during bursts</li> <li>- waste space for speed</li> </ul> <p>Disk Scheduling</p> <ul style="list-style-type: none"> <li>- FIFO - may have long seeks</li> <li>- Shortest Seek Time First - pick request closest on disk</li> <li>- may lead to starvation</li> </ul> <p>SCAN: closest request in direction of travel</p> <ul style="list-style-type: none"> <li>- not fair</li> <li>- requests in the middle get better treatment</li> </ul> <p>C-SCAN - SCAN in only one direction</p> <ul style="list-style-type: none"> <li>- skip requests on "way back"</li> <li>- block: logical transfer unit</li> <li>- sector: physical transfer unit</li> </ul> <p>Fetch All blocks correspondingly to desired bytes i.e. if block contains 1 of desired bytes, fetch File: collection of blocks arranged segmentally in logical space</p> <p>Dir: index mapping names to files</p> <p>Disk: linear array of sectors</p> <ul style="list-style-type: none"> <li>- every sector has integer address</li> <li>- controller translates from address to physical position</li> <li>- need to track free disk blocks</li> </ul> <p>open() has lots of overhead: checking, etc</p> <p>Path</p> <p>inode</p> <p>↳ Directory → file # → block structure offset → index structure</p> <p>File: permanent storage</p> <ul style="list-style-type: none"> <li>- data, metadata</li> </ul> <p>From user space, get file name (file control block) in directory structure in kernel memory</p> <p>jobs corresponding to file structure</p> <p>From SS target paper</p> <p>File control block</p> <p>File Allocation Table</p> <ul style="list-style-type: none"> <li>- which file is allocated</li> <li>- which blocks</li> <li>- each file number is index in FAT</li> <li>- file number is linked list of pointers to nth block</li> <li>- indices of disk blocks</li> </ul> <p>FAT is LL w/ blocks</p> <ul style="list-style-type: none"> <li>- file number is index of root of block list for file</li> <li>- FAT free list - unused blocks</li> <li>- FAT stored in boot code, backup copy on disk</li> <li>- blocks allocated everywhere be new blocks fetched from free list</li> <li>- files usually didn't have blocks in sequential order</li> </ul> <p>Directory: list of entries</p> <p>Access: file header for root</p> <ul style="list-style-type: none"> <li>- read header for root</li> <li>- read blocks for root</li> <li>- file header for "my"</li> <li>- block (first) for "my"</li> <li>- search for blocks</li> <li>- read file header for "book"</li> <li>- etc</li> </ul> <p>FAT:</p> <ul style="list-style-type: none"> <li>- no access rights</li> <li>- no header in file blocks</li> <li>- file metadata in directory</li> </ul> | <p>Unix Filesystem</p> <ul style="list-style-type: none"> <li>- i-node format</li> <li>- file number is index into mode array</li> <li>- metadata w/ file</li> </ul> <p>Inode</p> <p>Inode</p> <p>I data blocks (direct ptr) 48 kB</p> <p>Indirect block 14 MB</p> <p>Double indirect ptr 14 GB</p> <p>Triple indirect ptr 14 TB</p> <p>Worst case: 4 accesses to get to data for random access</p> <p>- Bitmap allocation for free list</p> <p>BSD 4.2: find blocks in contiguous region to expand file</p> <ul style="list-style-type: none"> <li>- limit fragmentation</li> <li>- advantage of bitmap</li> </ul> <p>Skip sector positioning to solve rotational delay</p> <p>Also, prefetching, read-ahead</p> <p>Inode for file stored in same cylinder group as parent directory of file</p> <ul style="list-style-type: none"> <li>- limit seeks</li> <li>- limit fragmentation</li> </ul> <p>How to Expand File:</p> <ul style="list-style-type: none"> <li>- allocate blocks by filling in small holes at start of block group</li> <li>- allows for contiguous region for later - bigger files can use this</li> </ul> <p>Links - DAG</p> <p>link entry in one dir + w/ file in another (file #)</p> <p>Hard - put on disk pointer to where you refer to</p> <p>Maintain ref count for each file</p> <p>Soft Link - contains path + name to file</p> <p>B-Tree - organizing files</p> <ul style="list-style-type: none"> <li>- each entry points to start of range for children OR its children are smaller than that entry</li> <li>- enables quick random access</li> </ul> <p>New Technology File System</p> <ul style="list-style-type: none"> <li>- variable length extents</li> <li>- rather than fixed blocks</li> <li>- everything is sequence of controllable values &gt; pairs</li> </ul> <p>Master File Table</p> <ul style="list-style-type: none"> <li>- database w/ flexible 1KB entries for metadata/data</li> <li>- extend w/ variable size tree</li> <li>- each entry kind of like inode, all in entry itself, can fit</li> </ul> <p>e.g. Record</p> <p>Std. info   File Name   Data   free</p> <p>Tser · u    start length Extent</p> <p>Some amount of pts to extents.</p> <p>- can be indirect, etc</p> |
|--|--|--|

Memory Mapped Files

- file mapped to file
- write directly to memory
- page fault under hood

Upon mmap()

- Create PTE entries for (file blocks in memory)  
page cache of file  
Read file contents from memory  
mmap() returns addr in memory of file

Availability: probability that system can accept and process requests

Durability: ability to recover data despite faults

Reliability: give me right answers - up + working correctly

- disk blocks contain error correcting codes ( Reed-Solomon )
- use special RAM - NVRAM for dirty blocks in buffer cache
- replicate to ensure that data survives long term

RAID - redundant arrays of inexpensive disks  
- layer of SW on top of HW to manage

RAID 1 - disk mirroring

- duplicate disk
- very expensive
- writes more expensive - duplicated
- reads a little better - can read from either

Recovery: hot spare, idle disk already attached

RAID 5 - High I/O Rate Party

- each disk has blocks, 1 parity block + some data blocks

|     |     |     |     |     |
|-----|-----|-----|-----|-----|
| D0  | O1  | O2  | O3  | P0  |
| D4  | D5  | D6  | P1  | D7  |
| D8  | O9  | P2  | D10 | D11 |
| D12 | P3  | D13 | D14 | D15 |
| P4  | D16 | D17 | D18 | D19 |
| D20 | D21 | D22 | D23 | P5  |

1    2    3    4    5

$P_n = \text{XOR of blocks in row}$

$D2 = D0 \oplus O1 \oplus D3 \oplus P0$

Overhead = parity block

$\frac{1}{n} \text{ disks} - 1 \quad \frac{1}{n}, n = \# \text{ blocks associated w/ 1 parity block}$

Write: ↑

Read: 0 if there 2 overhead fail  
high if reconstruct

Digitalize for read: 80% of time, don't need to reconstruct

= Load Balancing

Replication good for reads

- not as good for write, must wait for slowest one

RPC: client passes to RPC, passes to client stub that marshals args then sends to server  
- shared fate (local)

Write involves updating multiple blocks; mode, indirect block, bitmap  
Guarantee consistency

- threats to reliability:
  - interrupted operation
  - loss of stored data (NVRAM crashes)

Approach 1: careful ordering

- can check if any ops in progress after crash

- all write data block

- write

- allocate inode

- write

- update free bitmap

- update dir w/ file name

- → mode num

- update modify time

Revert/reclaim of crash

Recovery:

- scan inode table

- delete or put unlinked files in lost + found

- compare free block bitmap against inode trees

- scan dirs for missing update blocks

- time

Approach 2: Copy on Write

- create copy, write on that

- Swap new block w/ old

- update in Satch, disk

- writes can occur in parallel

2FS/Open2FS

- Swap out old version

- w/ new versions (ptrs)

in batches

Transactions: atomic updates

- all or none of updates

- from one consistent state to another

BEGIN

UPDATE

- tail or earliest, roll back

COMMIT

Properties:

Atomicty

Consistency - transactions maintain data integrity = constraints

Isolation - execution of one isolated from all others; no consistency issues

Durability: commit = persistent effort beyond crash

Simple Filesystem: disk = 1 big array, files one after another.

External Fragmentation - free space is unusable bc each block is so small

Internal - space wasted in allocated memory blocks bc of restriction on min size of allocated blocks

FAT Dir Only: Num - Addr - Index of 1st block - size

ext2 - top half; kernel starts ZIO ops

free driver bottom! serious interrupts produced by

FAT32 has 4 GiB file size limit

deportant - repeated without effect after first iteration

ZPC: V slave replicas and ACTION from master

Master [Write-request ACTION] → slave

Slave [Write-about/commit] → master

Master [Global-commit/about] → slave

Slave [ACK] → Master

I slave about, master sends a global-ABORT

• CLR, GLOBAL-COMMIT

• response from slave → request has been recognized

• committed

• ACK for global commit ⇒ action has been executed/committed

now fault tolerant

• log of slave contains proposal transaction and COMMIT/ABORT

LogByByteWriter - last byte put in send buffer

Log Byte Sent - last byte sent by sender to receiver

Log Byte Acked - Log byte acknowledged by receiver

Log Byte Read - last byte read by receiver in sequence

Log Byte Expected - last byte expected by receiver

Advertised Window + Max Rec Buffer = (Last Byte Read - Log Byte Read)

Send Window + Advertised Window = (Last Byte Sent - Last Byte Acked)

WriteWindow + Max Send Buffer = (Last Byte Written - Last Byte Acked)

ACK sent and before reading

Recursive Acq - faster but less scalable

R/W ≥ N+1 : Quorum Consensus

Prepare - timeout: global abort

Commit - timeout: resend msg

2PC-only AC

Log Structure - data stored in log form

Journaling - log used for recovery

need transactions for writes

- apply to disk, more tail

i Commit is there

- needs everything

else

discard log

Centralized System

- major func on one computer

Distributed - separate computers working together

Why:

- more compute power, storage

- cheaper/easier to build lots of simple computers

- easier to expand

- user has control over some components

- easier collaboration

Promise:

- higher availability

- better durability: store data in multiple locations

- more security/easier to make secure

Reality: 3 fails

Goal: traits/priority

- mostly complexity behind

simple interface

MAC - 48 bit assigned by card vendor

IP - 32 bit assigned when corp joins new

Endpoint identified by port

Functionalities:

- 1) delivery

- 2) reliability: tolerate pkt loss

- 3) flow control: avoid overflowing rx buffer

- 4) congestion control: avoid overflowing buffer & router along path

Layers provide set of abstractions for nw functionality + technologies

Properties of Layers:

Service - what a layer does

Service Interface - how to access Service

Protocol: how peers communicate to implement service

1-Physical

- more info between 2 sys connected by phys link

- specifying how to send/receive

- coding scheme used to represent bits, voltage/level

2-Data Link

- enable end hosts to exchange frames (atomic msgs) on same link

- send frames to other hosts

CDN - group of hosts on same link

- switch forwards frames to intended recipients

Protocols:

- partition channel

- token ring

- random access (listen then speak)

Switches

3-Network Layer

- deliver packets to specific nw address

- construct forwarding table

WAN - nw that covers broad area, connects multiple link layer nws

- connected by routers

IP - Best effort pkt delivery

4-Transport Layer

- end-to-end communication between processes

- reliability, time, flow/congestion control

- de-multiplexing of communication between host

VDP: multiplexing/demux of processes

TCP: connection set up/tear down, pkt corruption, retransmission, control

5-Application

- many services provided to end user

Lower 3 everywhere, top 2 only at hosts

Hourglass: all nw can interoperate

applications that run IP can use any

simultaneous innovations above/below IP

E2E: where to implement functionality

Conservative:

- default values can be completely implemented

- default unless relieve burden from hosts

Moderate:

- only if performance enhancement

- AND no additional burden on apps that don't need it

server: socket, bind, listen, accept

client: socket, connect

named pipe - unidirectional; UNIX socket - bidirectional

computers can share IP using network address translation

selective ACK (TCP) - adds ability to acknowledge individual pcks beyond current cum ACK

advertised window in TCP for new data

Log Flows

- writer groups writes before flushing to disk, improves performance

- read - must read through log to ensure we have most up-to-date version but can just cache this

buffer cache holds filesystem data in memory

Simple Filesystem: disk = 1 big array, files one after another.

- external fragmentation

External Fragmentation - free space is unusable bc each block is so small

Internal - space wasted in allocated memory blocks bc of restriction on min size of allocated blocks

FAT Dir Only: Num - Addr - Index of 1st block - size

ext2 - top half; kernel starts ZIO ops

free driver bottom! serious interrupts produced by

FAT32 has 4 GiB file size limit

deportant - repeated without effect after first iteration

ZPC: V slave replicas and ACTION from master

Master [Write-request ACTION] → slave

Slave [Write-about/commit] → master

Master [Global-commit/about] → slave

Slave [ACK] → Master

I slave about, master sends a global-ABORT

- CLR, GLOBAL-COMMIT

- response from slave → request has been recognized

- committed

- ACK for global commit ⇒ action has been executed/committed

now fault tolerant

- log of slave contains proposal transaction and COMMIT/ABORT

LogByByteWriter - last byte put in send buffer

Log Byte Sent - last byte sent by sender to receiver

Log Byte Acked - Log byte acknowledged by receiver

Log Byte Read - last byte read by receiver in sequence

Log Byte Expected - last byte expected by receiver

Advertised Window + Max Rec Buffer = (Last Byte Read - Log Byte Read)

Send Window + Advertised Window = (Last Byte Sent - Last Byte Acked)

WriteWindow + Max Send Buffer = (Last Byte Written - Last Byte Acked)

ACK sent and before reading

Recursive Acq - faster but less scalable

R/W ≥ N+1 : Quorum Consensus

Prepare - timeout: global abort

Commit - timeout: resend msg

2PC-only AC