

## \* What is an OS

Shrinithi  
Narayanan

- does not have a "main()" and doesn't really exit
- not like a normal program, more like a superative library
- called by programs and hardware
- bootstrap sets OS up, then used/called however
- acts as an intermediary system between app program + hardware  
e.g. who has permissions to access certain files
- provides an environment in which to run programs, provides service

HW ↔ OS ↔ Apps ↔ User

## \* Goals of an OS

- easy, consistent to use computer (better than barebones)
- efficient use of resources; e.g. abstraction of space and time ("simultaneous" processes)

## \* What things from vendor

- compiler? window system? web browser? X NONE (by our definition)
- INSTEAD: Everything that was in ~~hardware~~ hardware "kernel" mode ("core" (like UNIX "root" or Windows "admin"))
- other essentials? init processes?

## \* OS does 2 things

1. Create abstractions
  - disk → files, directories
  - CPUs → processes, context-switching
  - memory → address spaces, virtual memory
2. Manages resources
  - hardware and software
  - increase hardware utilization =  $\frac{\text{useful time}}{\text{total time}}$

## \* How does the OS do this

fool the user, fool the hardware. (hardware doesn't need to be aware of OS abstractions)

## \* History of OS

- Early history = no OS, just toggled-in instructions, time slots 1940's - 1950's  
punch-cards! inefficient
- Early batch systems - non interactive, patching many programs together 1950's - 1960's
  - a simple Resident Monitor — run program together
  - a job control language

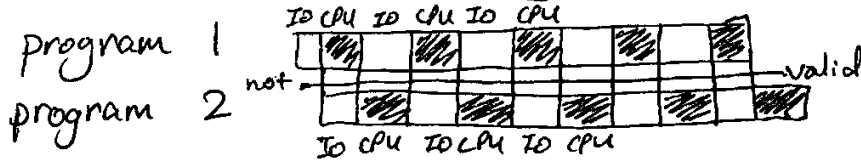
Multiprogramming - multiple "processes" time into memory  
memory split among the many programs (for each process)  
context switch only when needed (no need for extra overhead)

→

1960's - present

# efficient context switching usage of CPU

Shrini Narayan



Multix!  
Unix!

timesharing - interactive multiprogramming  
- context switching fast enough to simulate "timesharing", more frequent (100x frequenter) Batch: turn-around time  $\rightarrow$  time to completion

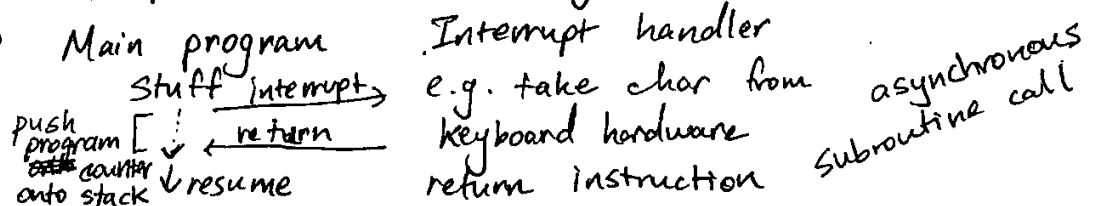
Hardware Support for OS Interactive: response time  $\rightarrow$  time to first response

problem: IO is slow

busy waiting: ~~while~~ while (no char yet?) "necessary" but unproductive  
how to know when to switch back here?

new hardware feature: interrupts

definition: an interrupt is an external signal (outside CPU) to the CPU (is async) Main program Interrupt handler

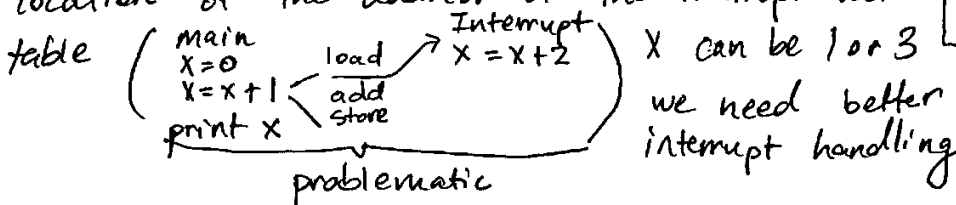


the Interrupt vector table

in memory, an array of pointers to functions/procedures  
on interrupt, external hardware gives interrupt type location? (hardware needs to know where)  $\rightarrow$  new special register! <sup>TM</sup> =

definition: the vector base register holds the location of the address of the interrupt vector table

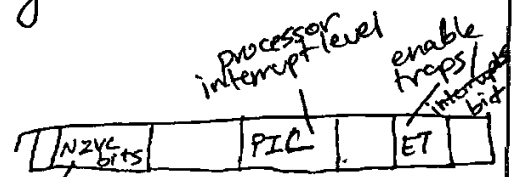
vector base register  
vector here means the ptrs to the other functions



hardware synchronization for interrupts

two mechanisms: disable + enable interrupts

example: SPARC processor storage register (it disables all interrupts)



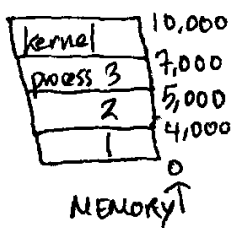
when "disable" pending interrupts are remembered: but only one of each type of interrupt

interrupt properties: new interrupt allowed only if it's a greater priority interrupt than current level priorities for interrupts

## Protection

Shrinithi  
Narayanan  
3

must be built into hardware (software too slow, insecure)



- Two new special registers
  - BASE: First legal address (process 3: 5,000)
  - LIMIT: number of bytes allowed (process 3: 2,000)
- Address  $X$  is legal iff  $0 \leq X \leq \text{LIMIT}$ 
  - IF legal, hardware then accesses  $X + \text{BASE}$
  - Actual BASE is hidden from process (for security)

## Dual-Mode Execution

Example: SPAR, ~~PSE~~ PSR



only allow changes when in kernel mode

"~~privileged~~ instructions"

"supervisor" bit (kernel mode on/off)

- e.g. changing BASE or LIMIT, changing vector base (interrupts), changing PIL, S, ET, perform IO, halt instruction

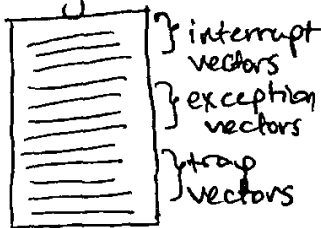
Definition: ~~An~~ An interrupt is an async signal from external source

Definition: ~~An~~ An exception is an error in current instruction execution

Definition: A trap is a special instruction used by user mode to request service from the kernel

all three used operations are similar:

usage of interrupt table for all 3



- however, only interrupts can be disabled (not ~~traps~~ traps)
- all 3 goes into kernel mode
- all push current PSR and PC onto stack (for returning)
  - processor status register
  - program counter
- \* these are the only 3 ways into kernel mode \*

- these 3 act as the only, controlled ways into kernel mode
- the clock interrupt: generates periodic interrupts so that kernel mode would have some consistent control (eg per 10 ms)
  - allows OS to regain control when needed eg for context-switching
  - can be used to count time, down time-outs, CPU usage

## Synchronization/Race Issues

Def: process is a program in execution; has own address space, but can be "multi-threaded"

- threads (also called "lightweight processes") can share same address space

race condition: different "processes" sharing data;  
outcome dependent on order of events

the outcome/order is not determined by the program

• maybe can solve with disabling/enabling interrupts, but problematic and inefficient (messes up clock, only works w/ single processor)

• high level tool for process synchronization

monitor foo { define shared vars: int x  
entry inc-one()  $x = x + 1$ ;  
entry inc-two()  $x = x + 2$ ;

essentially a

lock mechanism

when entry starts, acquire mutual exclusion over whole monitor.

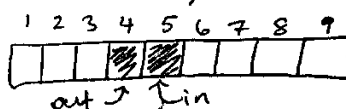
• sync enforced by compiler/runtime  
• only 1 process active at a time in monitor

• can only call entries from outside  
• ~~when~~ upon return, release mutual exclusion (a "lock")

## History of monitors

Three different types, Per Brinch Hansen 1973, CAR Hoare 1974, Mesa language 1974

## The Bounded Buffer Problem



circular buffer w/ in, out ptr  
counter as well

a buffer of N items, a producer, a consumer

producer: while (1) {  
    produce item  
    place item in buffer }

consumer: while (1) {  
    take item from buffer  
    consume item }

need a way to view buffer as full or not w/o it changing

• new variable type: no value type, only used inside monitor (not handlable by monitors)

two operations: X wait - blocks current process until...

(different rules for how to handle X signal, usually followed) X signal - wakes up one process if there is one waiting

• Hoare monitor's P is blocked, Q does X signal (signal and wait)

Q is blocked, transfer automatically to P

Q continues to run later after P leaves monitor or P waits on some other condition

(signal and continue)

## Mesa monitors

Q continues to run; P runs later after Q leaves or wait on condition var

## Brinch-Hansen monitors

• two ways to leave: return, or signal and return

(signal and return)



(s) to Bounded Buffer Problem (using condition vars under Shrivathi  
for bounded { struct item buf[N]; Hoare monitor) Narayanan

assuming int in = 0, out = 0, count = 0; condition empty, full; 5

problems can arise  
with multiple  
producers/consumers  
(who gets woken up?)

entry add-item(struct item \*data) {  
while (count == N) {full.wait;}  
buf[in] = \*data; in = (in+1) % N;  
count++; empty.signal;}  
entry remove-item(struct item \*data) {  
while (count == 0) {empty.wait;}  
\*data = buf[out]; out = (out+1) % N;  
count--; full.signal;}  
}

Starvation: infinitely unable to get resources while others keep succeeding

## The Readers-and-Writers Problem

one shared object, any # of 2 kinds of processes: readers and writers  
any # of readers active at once; if reader active, no writer allowed and vice-versa

### Nesa Monitor Solution

Monitor RW { int num-reader = 0, num-writer = 0; condition reader, writer;  
entry begin-read() {  
while (num-writer != 0) reader.wait;  
num-reader++;  
} can add reader.signal  
entry end-read() {  
num-reader--;  
if (num-reader == 0) writer.signal;}  
entry begin-write() {  
while (num-writer != 0 or num-reader != 0) writer.wait;  
num-writer++;  
entry end-write() {  
num-writer--;  
reader.signal; writer.signal;}  
}

### Hoare Monitor Solution # introduce accounting of bypassing to prevent starvation

Monitor RW { int num-readers, num-writers  
condition reader, writer  
int waiting-readers, waiting-writers  
int bypassed = 0  
entry begin-read() {  
if (num-writers > 0 || bypassed >= bypass-limit)  
waiting-readers++;  
reader.wait;  
waiting-readers--;  
num-readers++;  
if (waiting-writers > 0) bypassed++;  
if (waiting-readers > 0 and bypassed < bypass-limit) reader.signal;  
entry end-read() {  
num-readers--;  
if (num-readers == 0 and waiting-writers > 0)  
writers.signal;  
entry begin-write() {  
if (num-readers > 0, num-writers > 0)  
waiting-writers++;  
writer.wait;  
waiting-writers--;  
bypassed = 0;  
num-writers++;  
entry end-write() {  
num-writers--;  
if (waiting-readers > 0)  
reader.signal;  
else if (waiting-writers > 0)  
writer.signal;  
}

# Dining Philosopher's Problem

each philosopher can only use a fork to the left or right  
 can only pick up 1 fork at a time; must hold both forks to eat  
 must be possible for 2 to eat at once

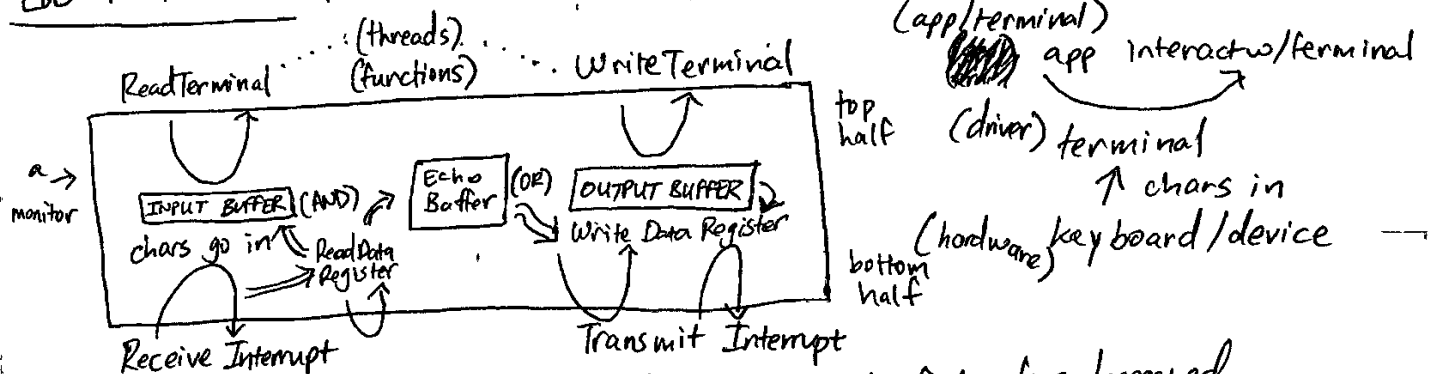
Dijkstra 1965

## Hoare Monitor

monitor dining { int state[5] = "thinking" THINKING, EATING, HUNGRY }  
 condition philcond[5]  
 entry pickup-forks(int i) {  
     state[i] = HUNGRY  
     test-forks(i)  
     if (state[i] != EATING)  
         philcond[i].wait  
 }  
 entry put-down-forks(int i) {  
     state[i] = THINKING  
     test-forks((i+1)%5)  
     test-forks((i+4)%5)  
 }  
 test-forks(i) {  
     if (state[(i+4)%5] != EATING and  
         state[(i+1)%5] != EATING and  
         state[i] == HUNGRY) {  
         state[i] = EATING  
         philcond[i].signal  
     }  
 }

Shrinini  
 Narayan

## Lab 1 Review Device driver for a set of terminals written as Mesa monitor



Flag to keep track of whether a ReadData or WriteData has happened  
 keep cycle of transmit interrupt

## Input

Keyboard: ReceiveInterrupt

• "\r" or "\n"  
 (^M or ^J)

• "\b" or "\177"  
 backspace delete

must preserve "scripts" (e.g. WriteT("hello") and WriteT("world") => "worldhello")

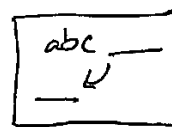
block until everything written or everything read  
 external functions for monitors (not "static")

input buffer  
 "\n"

non-empty last line  
 remove last char  
 else, do nothing

echo buffer  
 "\r\n"

non-empty last line  
 "\b\b"  
 else, optional 17



terminal

abc -> abc -> abc -> abc  
 "bell"

"helloworld" or

"worldhello"

extern foo(arg) {  
 Declare-Monitor-Entry-procedure()  
 condition X;  
 X = cond create();  
 condSignal(X);  
 }  
 typescript entry foo(arg) {  
 condition X;  
 X.signal;  
 }

... - new variable type, like int, but after initializing, Shrinishi Narayanan 7

2 operations for semaphores

$P(S) \Rightarrow \begin{cases} \text{while } (S \leq 0) \\ \text{do nothing;} \\ S = S - 1 \end{cases}$  atomic except "down" (allow context switch when waiting)

$V(S) \Rightarrow S = S + 1 \leftarrow$  atomically "up"

Careful use required (can deadlock)

Example use: Bounded Buffer Problem

Semaphore  $\text{mut-ex} = 1$ ,  $\text{empty} = N$ ,  $\text{full} = 0$

producer: while (1) {  
    produce item  
    P(empty)  
    P(mut-ex)  
    put item in buffer  
    V(mut-ex)  
    V(full)  
}

consumer: while (1) {

    P(full)  
    P(mut-ex)  
    take item from buffer  
    V(mut-ex)  
    V(empty)  
}

so, also need mutex to keep "empty" and "full" synced

need to for limiting both the producer and consumer

mutex is "binary" 0 or 1  
empty, full is "general" 0 to N  
defined by usage

Readers and Writers Solution

int readcount = 0;  
semaphore mutex = 1, writing = 1

Still has same problem of starving writers; add another bypass or semaphore

Dining Philosophers Solution

Semaphore  $\text{fork}[5] = 1 \leftarrow \text{table} = 4$

philosopher (i) {  
    while (1) {

        think...

        P(fork[i])

        P(fork[(i+4)%5])

        eat...

        V(fork[(i+4)%5])

        V(fork[i])

        V(table)  
    }

if everyone picks up fork at same time; adding a "table"

semaphore avoids deadlock: a FIFO semaphore eliminates starvation

assumption of fairness in "table" is key

adapt idea for our (Hoare) monitor, again assume fairness

bool avail[5] = 1; condition fork[5]

entry pickup-forks(int i) {  
    if (table == 4) seat.wait  
    table++  
    if (avail[i] == 0) fork[i].wait  
    avail[i] = 0

if (avail[(i+4)%5] == 0) fork[(i+4)%5].wait  
    avail[(i+4)%5] = 0;

putdown-forks  $\Rightarrow$

reader: while (1) {

    P(mutex)

    readcount++

    if (readcount == 1)

        P(writing)  $\leftarrow$

    V(mutex) First or

    read last reader

    P(mutex) "locks/frees"

    readcount-- writers

    if (readcount == 0)

        V(writing)  $\leftarrow$

    V(mutex)

```

entry putdown - forks (int i) {
    avail[i] = 1
    fork[i].signal
    avail[(i+4)%5] = signal = 1
    fork[(i+4)%5].signal
    table --
    seat.signal
}

```

The "General Critical Section Problem" Shrinani

Assumptions: all processes execute critical section at nonzero speeds; no assertions on relative speed of processes

Requirements: mutual exclusion: only one process in critical section at a time

progress: if some process wants to enter

critical section and no process in section, then only processes in enter or exit section influence choice of next process, and choice cannot be indefinitely postponed

bounded waiting: if process wants to enter critical section and has begun entry section, you can't be skipped indefinitely

Implementing Synchronization (for process 0 and 1 for now)

• Shared variables required

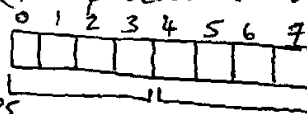
int x automatic assignment to 4 bytes

int \*y pointer to 4 byte space

char z[4] doesn't know if z

y = (int\*)z is multiple of 4

\*y = 5 ← store may thus be problematic



memory spaces

• using 1 shared var for synchronization

int flag = 0

enter while (flag == 0) no mutual exclusion

flag = 1

exit flag = 0

order of these 2 ops

int flag = 1

enter while (flag != 0)

exit flag = 0

deadlock

int turn = 0

entry while (turn != i)

lack of progress

if process doesn't take turn

exit turn = 1 - i

• using 2 shared vars for synchronization

int flag[2] = 0, 0

enter while (flag[i-1])

no mutual exclusion

flag[i] = 1

while (flag[i-1])

deadlock

enter: flag[i] = 1

while (flag[i-1])

flag[i] = 0

wait a little

flag[i] = 1

exit: flag[i] = 0

violates bounded wait, perfect sync, low priority

exit flag[i] = 1

flag[i] = 0

• using 3 shared vars for synchronization

Dekker's algorithm (1965)

int flag[2] = 0, turn = 0

(combo of turn and flags)

See book for full algo

Peterson's Algorithm (1981)

enter: int flag[i] = 1, turn = 1 - i

while (flag[i-1] and turn == 1 - i)

exit: flag[i] = 0

critical point

solution to critical section problem for 2 processes

but not very efficient busy waiting + only 2 processes



Art's Bakery Algorithm (1974) for multiple  $N$  processes Shrinithi Narayanan

9

if choosing  $[N] = 0$ , number  $[N] = 0$

```

enter { choosing[i] = 1
      /* number[i] = max(number[...]) + 1 */
      int x, temp = 0
      for (x = 0; x < N; x++) {
          if (number[x] > temp) temp = number[x]
      }
      number[i] = temp + 1
      choosing[i] = 0
      for (x = 0; x < N; x++) {
          while (choosing[x]) { ; }
          while (number[x] != 0 && (number[x] < number[i] ||
              (number[x] == number[i] && x < i))) { ; }
      }
  }
  exit { number[i] = 0 }
  
```

choose a number larger than the max (this is the ticket number)

no process can have advantage over some process more than once

busy waiting efficiency

## Hardware Support

test and set instructions  
 test\_and\_set x  
 temp = x  
 x = 1  
 return temp

one var solution w/ hardware support, but may still have unbounded waiting acts as a "lock"

Swap a, b

temp = a  
 a = b  
 b = temp

lock = 0  
 local key = 1  
 while (key)  
 swap(lock, key)  
 lock = 0

hardware supported

```

exit { local int x
      x = (i + 1) % N
      while (x != i and waiting[x] != false)
          x = (x + 1) % N
      if (x == i) lock = 0
      else waiting[x] = false
  }
  
```

## Using hardware to reduce starvation

int waiting[N] = false; int lock = 0

entry { local int key  
 enter waiting[i] = true  
 key = 1

while (waiting[i] and key)  
 key = test\_and\_set(lock)  
 waiting[i] = false

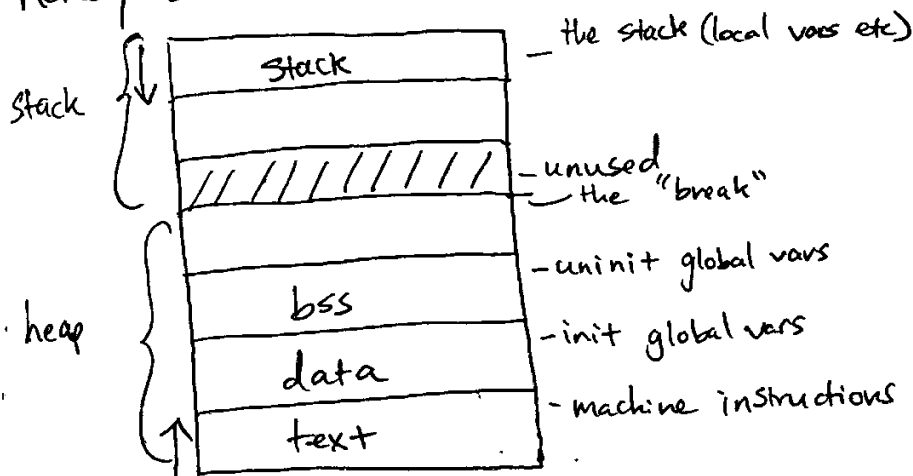
below busy loop, to avoid we if (uniprocessor || lock holder known to hold a "long time") ask OS for context switch else, for (i = 0; i < large # and still can't actually lock; i++) if (still can't acquire lock) ask OS for context switch

sequential loop thru buffer prevents starvation

# Processes, Memory Management

Shrini  
Narayan

## Memory structure



```
int x=1
int y[1000]
main(...) {
    int a=2;
    int b;
    a=b;
    proc();
}
```

register values: must be stored somewhere even when process not used  
 per process OS data structure: process control block (PCB) (track when not being used). PCB holds saved reg. values, pid, user id, uid, group id (gid), accounting info, lots of pointers!

have ready and active queues of PCBs (one for blocked?) <sup>circular linked list</sup>

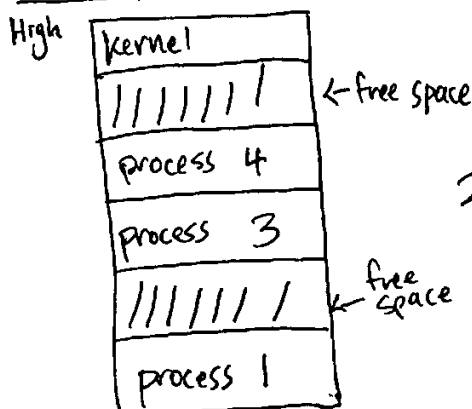
ready  $\square \rightarrow \square \rightarrow \square \dots$  pop and add for context switching

active  $\square \rightarrow \square \rightarrow \square \dots$  each square = PCB

blocked  $\square \rightarrow \square \rightarrow \square \dots$  ← maybe separate queue per reason for being blocked

the "idle" process, X HALT  
 (when ready queue empty) JMP X ← saves power before blocking (interrupt wakes up)  
 active PCB ≠ head of ready queue

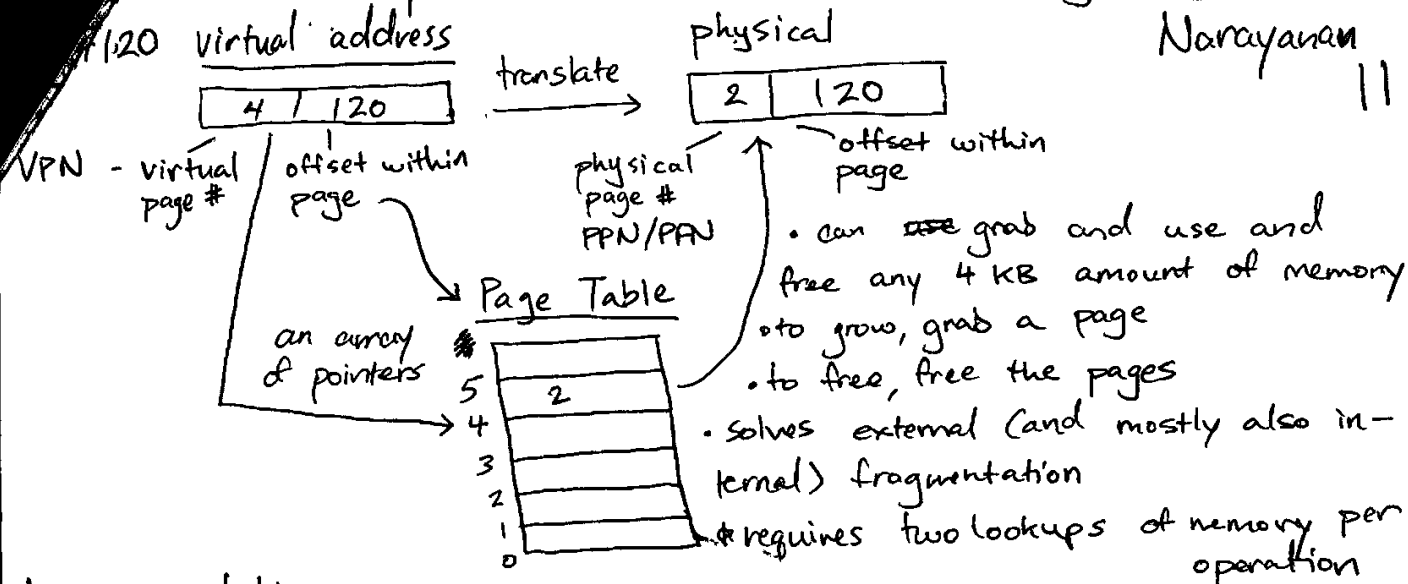
## Memory Management



- "fix" by compaction (copy process memories down to <sup>used free chunks</sup>) ← inefficient
- to allow process growth, must push for space, but over time, may have unused space in each process space

2 problems:  
 external fragmentation: enough free space for new alloc, but ~~not~~ not enough contiguous  
 internal fragmentation: "trapped" memory inside allocated chunks left unused

divide memory into fixed sized chunks, pages (eg 4KB) Shrinithi Narayanan 11



### the page table storage

- special hardware registers (but table can get too large) (PTBR)
- in memory: then, tell hardware location with page table base reg and page table limit register (PTLR). PTBR - physical address of PT array; PTLR - # of entries in PT (physical addresses)
- after one initial physical address read, OS initializes PT, and virtual addresses are used.

### Hardware Cache of Page Table Entries (the TLB)

- high hit rate = less need to fetch something from memory.
- managed almost entirely by hardware (lookaside to TLB before main)

TLB entry: 

vpn	pte
-----	-----

 ; OS flushes TLB on Context switch

valid bit OS can flush single entry of TLB (if present)

### Page Protection PTE

valid bit flag, 0 = not valid  
allowed access

also disallow PTE's out of range of PT range (if  $VPN \geq PTLR$ , disallow)  
protection bits: read/write/exec permissions for kernel + user (valid ≠ 0)

### Kernel: where is it?

> own separate address space? requires context switching to get in/out (expensive!)  
problem: how to access user memory? →

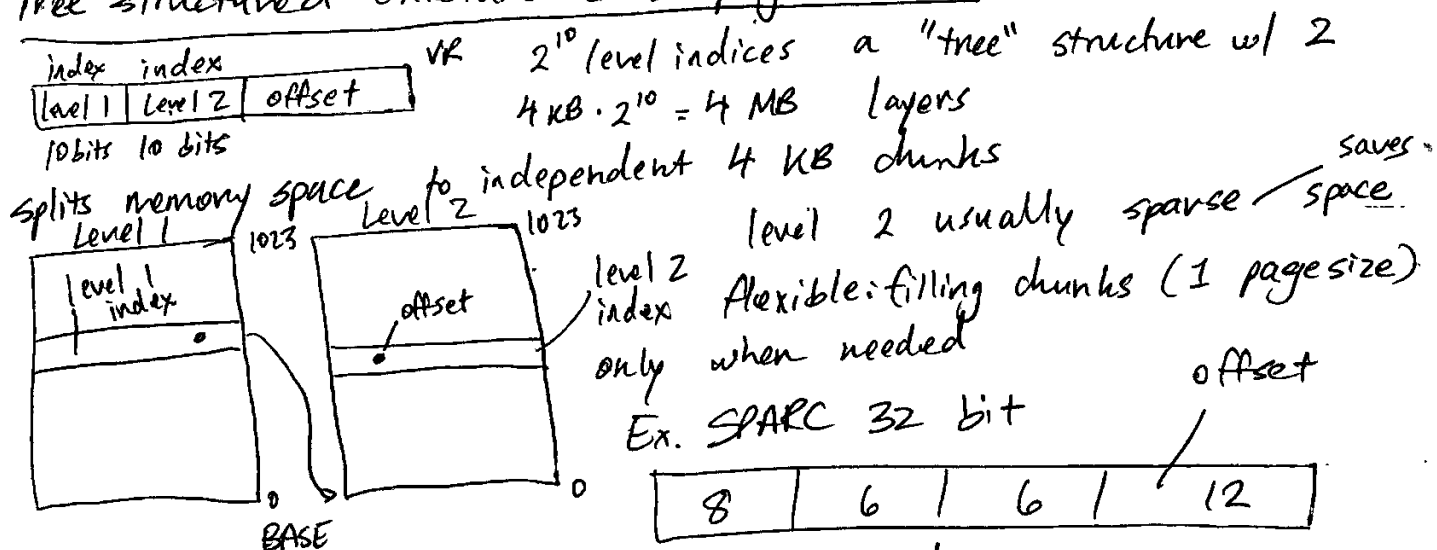
Shrinithi  
Narayanan  
12

> in all process address space  
 problem of differing page numbers per process  
 problems w/ storing data and context switching  
 problem w/ interrupt vectors jumping to wrong spot

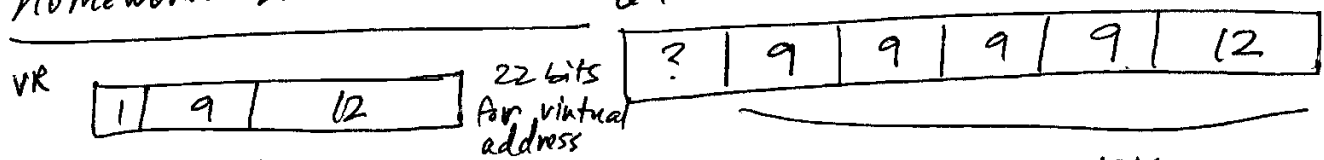
> clever use of virtual address space  
 use different page tables for each region (PTE0 and PTE1)  
 only modify PTE0 during context switch (be careful with protections)  
 kernel stack used (delicately placed in user process region 0)  
 synchronize virtual address space w/ physical address  
 that is,  $vpn \Leftrightarrow pfn$ , before enabling VM

build kernel heap at start (IVT, PTR1, stack pointer), map  
 PTE's correctly, then enable VM

### Tree structured (hierarchical) page table



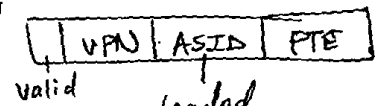
### Homework Lab 2 - RCS 421 64 bit Intel



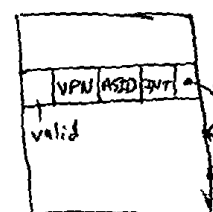
selects 1 of 2 registers, each pointing to a 48 bits level 2 table.  
 level 1 isn't an array in this case, just registers  
 when loading program, one must place PC, SP in kernel stack  
 context switching: may be problematic since kernel stack disappears,  
 instead re-point to shared kernel stack in region of kernel.

have address space ID register (ASID)  
 one for each address space with TLB entries in TLB  
 TLB entry

Shrinishi  
 Narayanan  
 13  
 for context switch prot



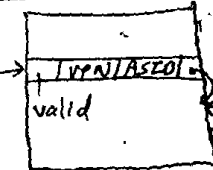
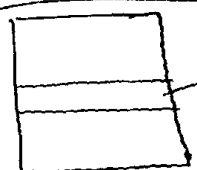
Inverted Page Table - only one per entire system  
 w/ TLB entry, checked when looking up in TLB



check/resolve collisions  
 linked list of image of hash function output

$(ASID, VPN) \rightarrow Hash Func \rightarrow 0 \dots max PPN$

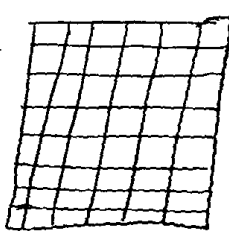
Hash Anchor Table, allows one to go to last seen free space in VPN to reduce collisions



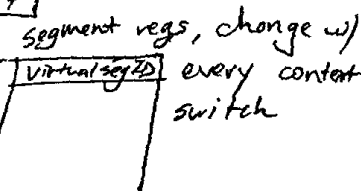
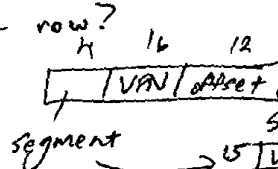
if still collision, go down linked list

Hashed Page Table

VPN, ~~ASID~~ VSID  
 $\hookrightarrow$  hash func  $\rightarrow$  PTE  
 in PTE, VSID + VPN  
 check for actual match



primary  
 alt



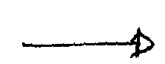
prot, valid, pfn bits as usual

alternate ~~space~~ hash for flowing bits.

storage descriptor #1  
 - size + phys addr of table

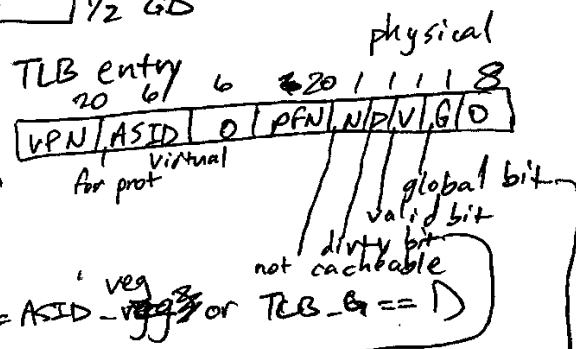
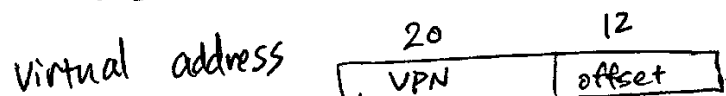
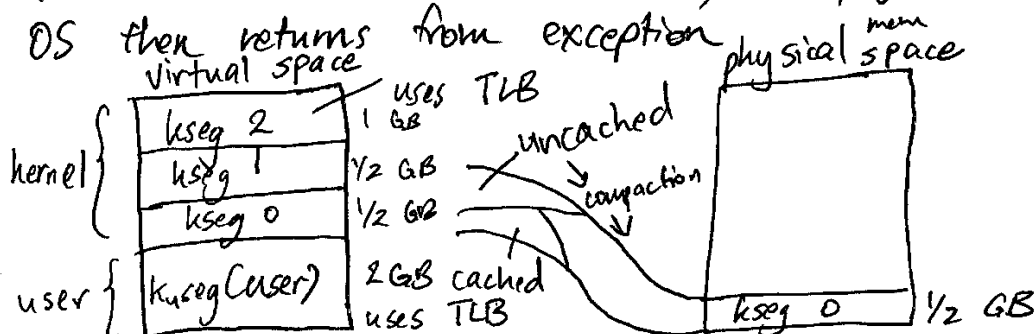
segment regs, change w/ every context switch  
 context switch simple, entries in one compact place

Recap: these are methods where hardware knows page format, TLB, complex structures, can we make the hardware really simple instead?



Shruti

Narayanan



HW searches ~~for~~  $\downarrow$   
 TLB Ar TLB  
 valid == 1  
 and TLB vpn = addr-vpn and (TLB-ASID)

$0 = k_{seg0}$  reading bits  
 $100 = k_{seg0}$  exhibits region  
 $101 = k_{seg1}$  of address  
 $11 = k_{seg2}$

4 CPU registers

TLB entry LO - holds physical stuff for  
TLB entry HI - holds virtual stuff

(writable)  
readability is whether  
the page is in TLB

index - "probe result bit", 6 bits - (TLB subscript)

random - "random number", 6 bits — counts from 8 → 64 counting entries 1-7 safe from over-write

4 instructions

TLBR - read the index<sup>th</sup> entry

TLBP - probe TLB for entry matching TLB entry HI  
(Hardware gets probe result = 1 on success and fills in TLB entry LO)

TLBWI - write into TLB index<sup>th</sup> entry from TLB entry HI and LO

TCBWR - write into random<sup>th</sup> entry (not \*actually "random", but close)

64 bit SPARC

only overwrite  
entries w/ lock bit = 0 enough

"TLB Data-In" register = TLB entry LO and HI  
 a TLB entry has a "lock" bit (writes down the ~~TLB~~ entry in TLB: "locked")  
 removing/adding cached page tables applies to hashed page table  
 hashed page table vs. software TLB

Initial

Shrinithi  
Narayanan

native to paging  
segment: a contiguous chunk of memory w/ hardware enforced 15  
address = (segment #, offset in segment) wraps <sup>beginning + size</sup> around to 0, doesn't  
4 segment registers (since 86 386 FS) affect segment #

CS = code segment register (CPU holds offset)

DS = data segment register (instruction gives only offset)

SS = stack segment register (stack pointer holds offset)

ES = "extra" segment register (override prefix to change default segment register to use format instruction)

### Segment Table

valid	prot	segment size	phys base addr
0			

PBA + offset  $\Rightarrow$  addr

2 tables - local process + global kernel

local table descriptor register (LTDR)  $\rightarrow$  phys

Global table descriptor register (GTDR)  $\rightarrow$  addr #, num entries in table

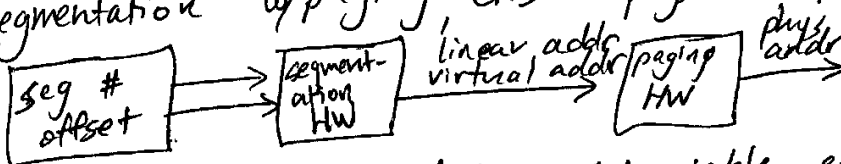
Segment register contains seg # = subscript in table and

table indicator ~~0 = local, 1 = global~~ (0 = global, 1 = local)

segment register has: segment descriptor code regs = copy

of that segment table entry

segmentation w/ paging, CR3 = phys addr page table (in C language) <sup>not really used</sup>



to turn off Segmentation, set table entries to 0, seg registers to 0, limit to MAX

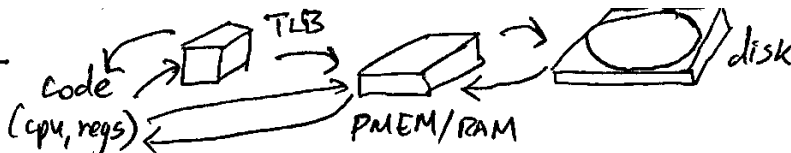
Used sometimes, for ex, for per user shared variables that appear differently for each thread

### MORE ABSTRACTION!

ES register chapter 8: everything in VR space is in PR space (abstract) (physical)

chapter 9: abstraction virtual space VR  $\subseteq$  PR  
may be larger than pmem, physical pages become a cache of virtual pages, other virtual pages are on disk ("backing store")

# Demand Paging



Shrinithi  
Narayana  
16

abstract

virtual space may be larger than pmem

physical pages become a cache of virtual pages, other virtual pages are on disk ("backing store")

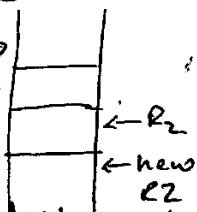
assumes: temporal locality - likely to access same piece of memory again in near future

spatial locality - likely to access same memory nearby in near future (Principle of Locality) victims!

Thus, cache pages when they are accessed

hardware requirement: all CPU instructions must be "restartable" being able to remember and do instructions so to redo instructions, after, eg returning from an exception handler

example: `mov R1, -R2`



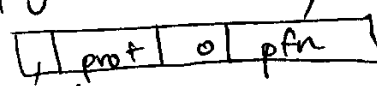
page exception,  
after returning, undo pointer movement

copy buffer to buffer

ex. move 5 src len, src, file, dest len, dest

page table entry: hit or miss

differing buffer sizes = partial copying must be undone

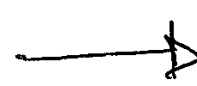


valid hardware generates exception. if because of demand

paging, call a "page fault" - OS finds free physical page on free list, or evict some other page (victim) from phys mem

extra HW support

in PTE, referenced bit: set by HW when HW used its PTE  
dirty (modified) bit: set by HW when HW uses that PTE for a while





# Replacement Algorithms (victim selection)

Shrinithi  
Narayanan

constant (bad idea)  
: random

both ignore principal of locality  
predict future... page default

2: minimum ~~(optional)~~ "page reference string" 2 1 4 2 3 0 4 2 5 4  
victim = page whose next use is farthest into future  $\frac{0}{2}$   $\frac{1}{x_0}$   $\frac{3}{4}$   $\frac{3}{2}$   
depends on size of phys memory and quality of prediction

3: FIFO victim = "oldest" page (time of insertion)

eg 2 1 4 2 3 0 4 2 5 4  
 $\frac{0}{x_0}$   $\frac{1}{x_2}$   $\frac{3}{x_5}$   $\frac{3}{x_4}$

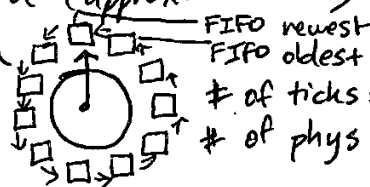
4: LRU victim = last recently-used (time of access) queue overtime

eg 2 1 4 2 3 0 4 2 5 4  
 $\frac{0}{2}$   $\frac{1}{x_0}$   $\frac{2}{4}$   $\frac{3}{x_5}$

2 1 4 2 3 4 2 5 4  
2 1 4 2 0 4 2 5  
2 1 4 3 0 4 2  
1 2 3 0 0

Implementing LRU (approximation)

• Clock algorithm  
gives each page a "second chance" within a cycle, otherwise, victimized  
remove victim from linked list, add new PTE to head of FIFO  
can also try using the dirty bit, but policy is subjective...



while (1) {  
advance hand (clockwise)  
if (PTE[hand].ref == 0) break;  
else clear PTE[hand].ref (set to 0)  
hand points to victim!

enhanced second-chance algorithm  
favor clean over dirty in victim selection

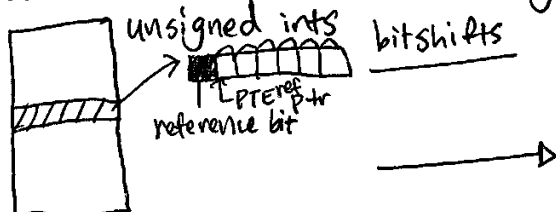
order of	ref	dirty
victim	0	0
pref	0	1
	1	0
	1	1

modified check algo

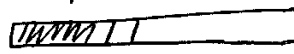
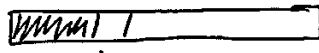
first revolution: only look, don't ~~clear~~ clear ref. bits, but break if (0,0)

second revolution: look, and do clear ref. bits; if (0,0) or if (0,1) break

• additional - reference - bits / aging / reference bit history algorithm



periodically (every few clock ticks), do:  
for each phys page, right shift by 1, entire array  
copy PTE ref to high order bit  
clear ref bit

on page fault, find min value in array; that's the victim Shrinani  
 comparing A and B:  $A < B$  iff A  Narayan  
 Global and Local Replacement B 

global

compete among all frames

Example FIFO: single FIFO list

performs better and more flexibly

(Process-by-process alloc)  
 local

compete among only own frames

Example FIFO: one FIFO list per process  
 more fair?

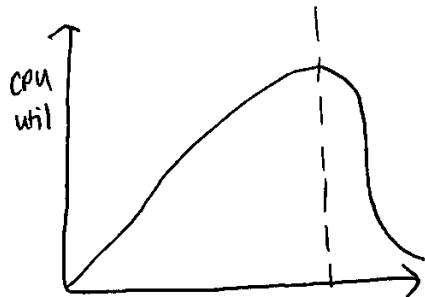
for all  $j > i$ ,  
 $A_j = B_j$  but  $A_i = 0$   $B_i = 1$

The working set: virtual page it is

$WS(T, w) =$  all pages sent over that time

at time  $t$  working at window size

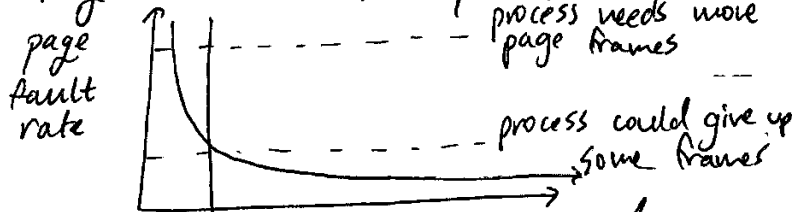
Reference BH History - any phys page  
 while loading  $V$  bits are all 0



degree of multiprogramming

thrashing: spend more time on faults  
 than useful executions  
 which process should be swapped out

Page Fault Frequency Algos



File Systems

• Hard disk

Surface is  
 magnetically coated,

hand reads magnetic signature

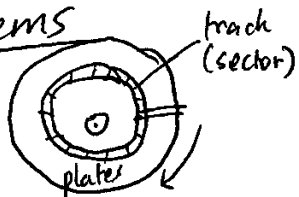
rotational speed not as easily

controlled, read/write head can  
 move in/out, must move seek to place

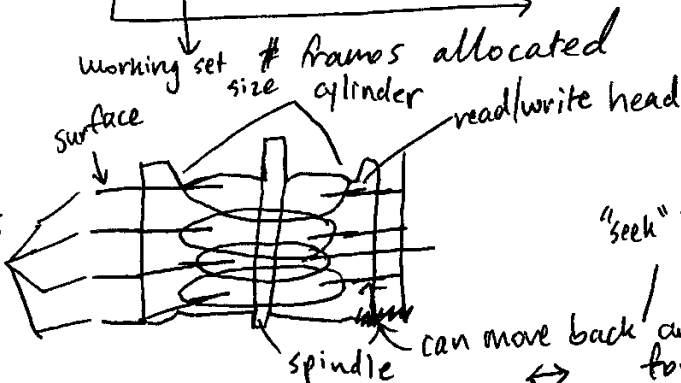
on disk

TODO: seek to correct cylinder + select head (track) +  
 rotational latency + bits of sector

to save time, save time needed for seeking in/out  
 for example



platters



"seek" in/out

	old	common	test
rpm	3600	7200	10k
rev	16 2/3 ms	8 1/3 ms	6 ms
ave. 10% latency	8 1/3 ms	4 1/3 ms	3 ms
seek avg	20 ms	9 ms	4 ms

Geometry (head #) x (cylinder #) x (sector / track #)

Shriniithi  
Narayanan  
19

disk addr: (cylinder #, head #, sector #)

disk request scheduling (no longer just OSS job)

- automatic bad sector forwarding (eg more bits to further cylinder to reduce seeking time)
- variable # sectors/track

- other secret geometries (disk companies competition)

new disk addressing: "Logical Sector Addressing" = sector #

File Systems bits structure (on disk)

put "related" things "nearby" each other

structs are "permanent": bugs cause permanent defect

tradeoff between simplicity and performance

disc sector vs. file system blocks

L512 bytes

L4 KB

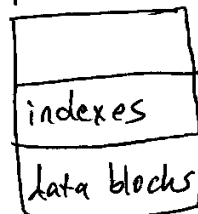
L2<sup>n</sup> bytes ← can lead to internal fragmentation

File system format

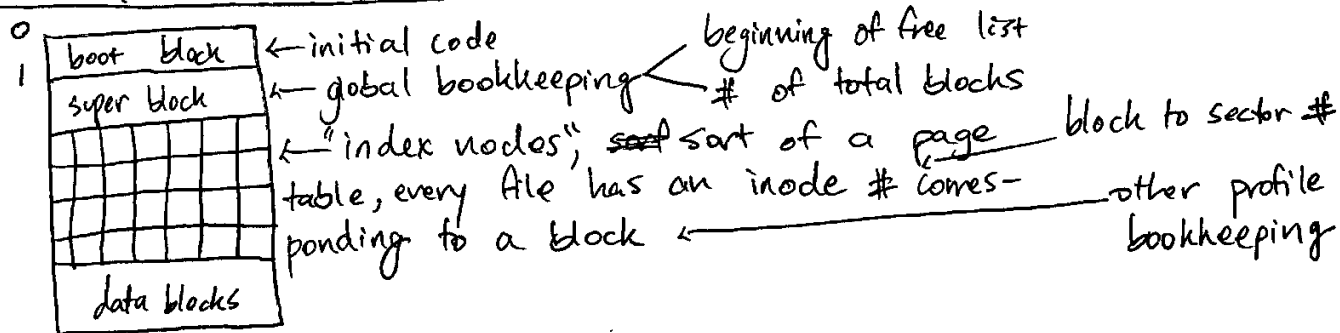
data structs on disk

global: free space list

~~file~~ profile: which blocks in this file



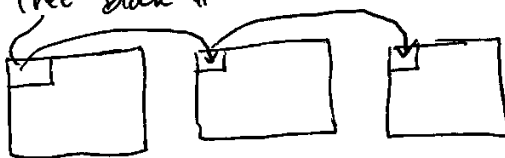
File System Format (classical Unix)



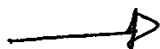
Free space management

1. In superblock: first free block #

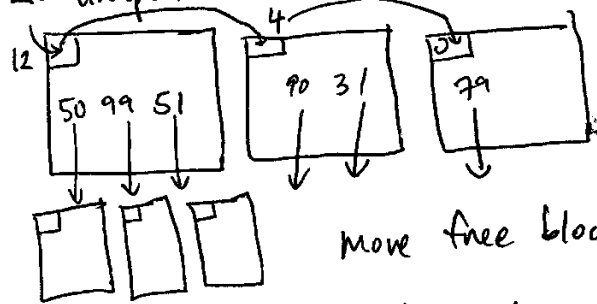
normal linked list



inefficient <sup>seek</sup> lookup to disk  
for next free block  
space taken up



## 2. "Grouped" Linked list



consume same amount of space but fewer disk reads for more free block #'s

Shrinithi Narayana 20

move free blocks once used, value is  $\emptyset$  or 0

3. linked list w/ counts change each entry to (block #, count of contiguous free blocks) but how to handle coalescing?

4. bitmap (1 bit/block)

0 = free, 1 = ~~not~~ not free

4 KB block size = 4096 bytes = 32768 bits  $\cdot$  4096 = (large #)  $\leftarrow$  many blocks represented (1 bit/block)

per file block allocation "mode"

1. contiguous: in inode = (1<sup>st</sup> block #, # blocks, # bytes)

2. linked list: first block # in inode  $\xrightarrow{\text{seek}}$  etc  
inefficient # of seeks to get next block  
ptr takes up space: powers of 2 broken (blocksize = 2<sup>k</sup> - ptr size)

3. file allocation table (FAT): one global table = array of block #'s  
entries correspond to file #, first block # in inode 40, follow table ~~entries~~ entry #'s 40  $\rightarrow$  5  $\rightarrow$  100; alternate file e.g. 3  $\rightarrow$  41  
similar to linked list, but table indexing much faster/efficient.

4. indexed allocation

(a) in inode array of block #'s

(b) in inode: block # of index

(c) in inode: block # of block of block #'s

5. combined scheme: handles files of all sizes w/ not much overhead

Naming: where is the name stored?

in the inode?

symlinks? stored where?

store directory entries same as file data

dir has inode  
dir entries in data blocks

structured (hierarchical) directories  
on hand. for pathnames, suppose lookup for  
a/b/c/d/e/f/g/h

- hard to type, slow lookup for each directory down
  - so, current directory at each process: in PCB: inode # of cur dir
- looking up path name: starts with "/" (absolute pathname) →  
lookup start at constant, root inode, otherwise use cur dir inode  
to allow cd. and cd .., have dir entries for . and ..

• multiple names for same file, references to same file from different locations. keep a reference count and free when ref count is 0. no easy way to check this problem.

no links to directories? link("old", "new") unlink("new")

mounting filesystems: how to name inodes unlink("old")

when one has multiple discs: each disc has its own inode #s, starting at 0 or 1; mount(dev, dir, type, flags, opt)

dir goes to device

symbolic links symlink("old", "new") vs. hard links link("old", "new")

- way to distinguish original name
- way to detect when a cycle would happen when disconnected
- limit # of symlinks when traversing to prevent infinite loops
- when original link is deleted, other symlinks invalid

- no links to discs
- can't tell which is old vs. new link



New Linux Filesystem - divide disk into cylinder groups

Shrinan  
Narayan  
22

- divide disk into cylinder groups

in each: subset of data blocks, subset of inodes,  
list of free blocks, free inodes in cylinder group  
copy of superblock (all read only)

put non dir file inodes in same CG as dir inode

put new dir inode in new CG

put data blocks in same CG as inode (spread lead around)

start new CG after all dir block #'s, and every 1 MB after

> limits # of seeks, backs up superblock

Multiple FS on same disk

how to track all these dynamically allocated inodes  
NTFS (Windows NT Filesystem)

MFT (Master File Table)

First 16 MFT entries

0 MFT itself

3 root dirs

6 bitmap of free space

7 boot block

block #  
within file

virtual cluster #

each MFT entry 1 KB

data block in each file sequence & boot block

of contiguous chunks = "extent", each described by (VCN, LCN, len)

### Lab 3 Yalrix FS

just look at spec, order of operations

important (might overwrite <sup>wrong</sup> file after  
(could have unused space trapped/allocated forever)

block #  
within whole FS

logical  
cluster #

size of  
block

freeing inode w/ links pointing to it  
(could have unused space trapped/allocated forever)

### File Protection

In PCB, uid + gid  
user group

In inode, "mode" protection bits

Ex. file mode 644 or 446

dir mode 755 or 444

Rwx Rwx Rwx

(restricted delete on dir)

1777  
/tmp

Rwx Rwx Rwx  
user group other

"exec" opens dir

ability to create files  
in dir

cannot rename or  
delete files in tmp

sticky bits: ☐ ☐ ☐  
setuid setgid sticky

bits

→

## File Protections

### File

Read read contents, view prot, owner, attrib  
Write change contents, change prot, owner, attrib  
Read/Exec "read" and can exec (pass to kernel call)  
Modify "write" and delete file  
Full control all of above  
list contents X

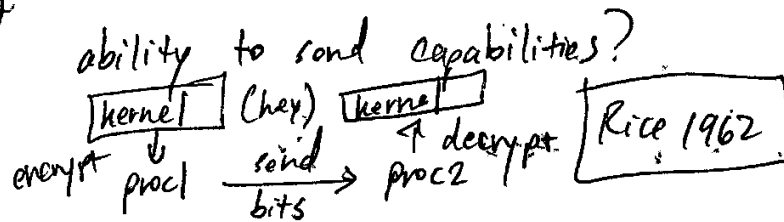
### Dir

"list contents", view prot, owner attrib  
change names in dir, change prot, owner, attr  
"read", list permissions, traverse  
"read+exec"; delete dir  
all above + below  
list contents

Shruthi  
Narayanan  
23

Protection domain = list of all "objects" and, for each, list of "rights" on that object  
could represent as access matrix this matrix would be large + sparse  
column = access control list (ACL) for that "object"  
(what can others do to you) (like window prot., unix a bits)

row = capability list  
(what can you do to other objects)

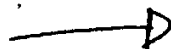


capability: secured name of some object and the rights to do various operations on that object (possession gives you the rights to do things on that object)

## Program Security - Worms and Viruses

worm - independent program, attempts to spread self over network to other computers

virus - code in some program that, when ~~run~~ run, tries to spread itself by embedding in other programs



The Morris Internet Worm Nov 2 1998 from MIT

Three ways of spreading

1. Remote login ~~and~~ - rlogin and rsh try variations on username  
trusted computers, brute force typing passwords

2. Send mail DEBUG command (SMTP)

prompts for password; the email addr = shell command

default = "u" email body = std input

often left unchanged

runs as root  
admin access  
to remote comps

backdoor!

3. The "finger" Server

Prevention gets → fgets.

performs a buffer overflow

check for strcpy → strncpy

something C {

char buf[512] buf;

gets(buf);

return;

sprintf → snprintf

Widely used, open ~~source~~

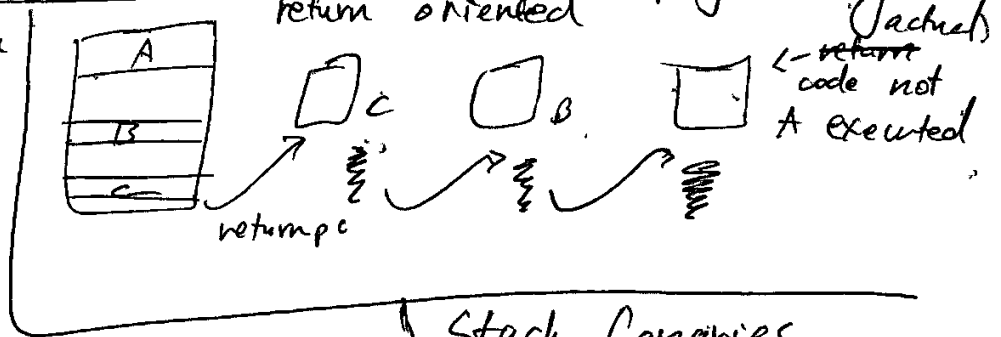
or little snippets

Address Space Layout Randomization

return to lib.c/  
return oriented programming

load program @ random  
address, load libs @  
random address

random space b/w  
stack frame, random  
malloc locations



compile to position independent code (PIC)

Stack Canaries

at program start, generate  
random #, store it "out of  
the way"

inside procedure: push that # onto stack, before return,  
check random # matches (random # relative to segment #, so can't  
know it's there)



# Security Services and Cryptography

Shrinithi  
Narayanan

confidentiality: protection against "eavesdropping"

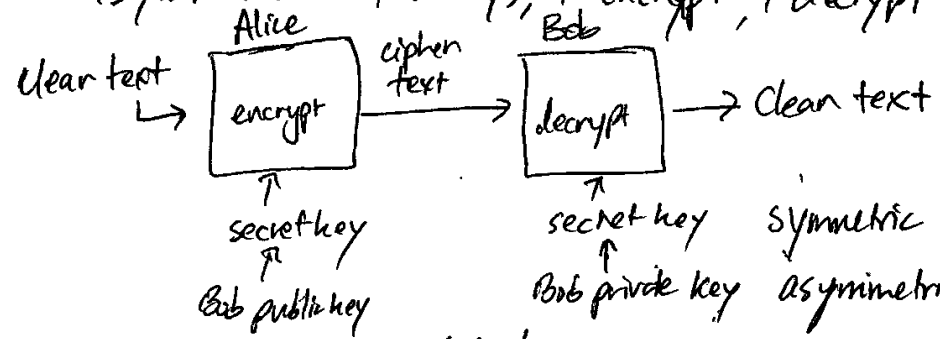
authentication: confirm msg sent by who sender claims to be

integrity: confirm contents not modified

confidentiality - encryption

symmetric - same key used for encrypt and decrypt (secret)

asymmetric - two keys, 1 encrypt, 1 decrypt (public, private)



symmetric "cheap"

asymmetric "very CPU expensive"

authentication

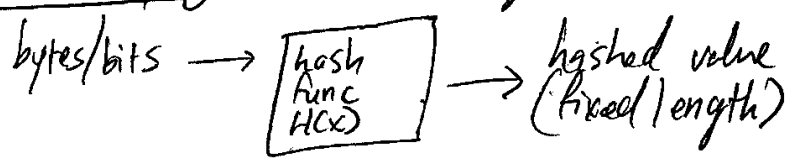
symmetric:  $E_{\text{secret } A+B}(\text{data} \parallel \text{Alice}) \parallel \text{Alice}$   
 ↑ encrypt w/ secret key

check Alice = Alice after decrypt

asymmetric:  $E_{\text{private } A}(\text{data} \parallel \text{Alice}) \parallel \text{Alice}$

decrypt w/ public key  
 (same can be done w/ integrity)

## Secure Cryptographic One-Way Functions

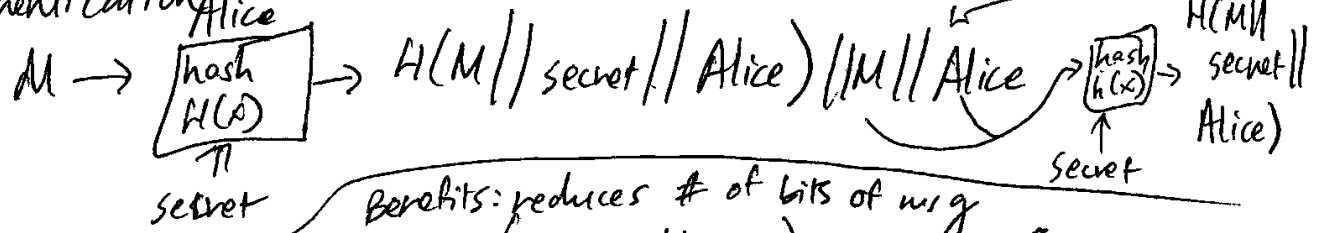


Given  $H(x) = y$ , cannot find  $z \neq x$  such that  $H(z) = y$  as well

Example: MD5(x) = 128 bits

SHA-1(x) = 160 bits

authentication

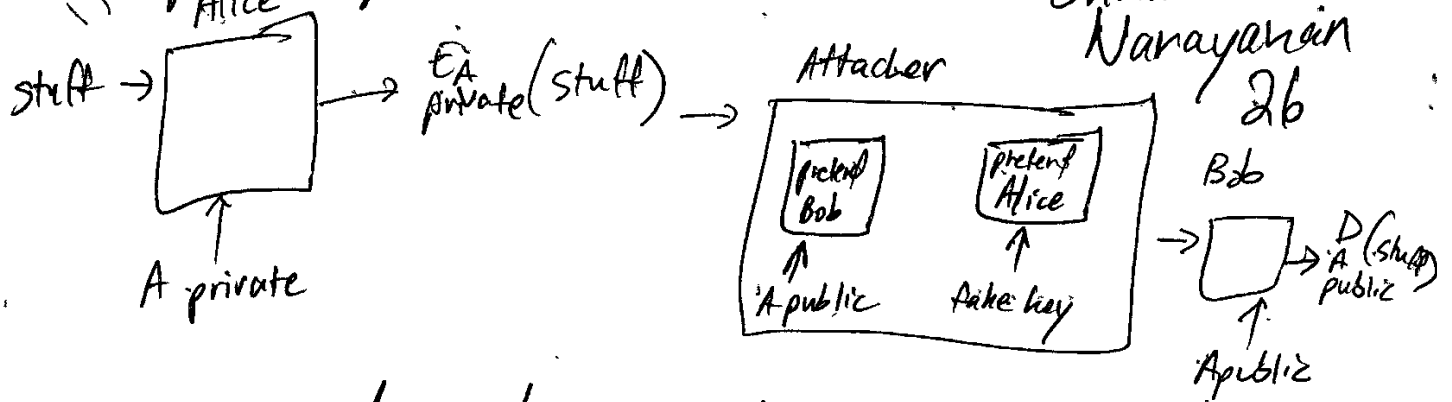


Benefits: reduces # of bits of msg

eg  $E_{\text{private } A's \text{ key}}(H(\text{data} \parallel \text{Alice}) \parallel \text{Alice} \parallel \text{data})$

Shruthi  
Narayanan  
26

## Attaching the Asymmetric Model (MATH)



the public key passed around can be hijacked

## Public Key Certificates

certification authority = comodo, symantec, godaddy

• ~~certificate~~ certificate  $X = (CA, \text{Alice}, \text{Alice public key})$   
 which cert authority

• use :  $E_{CA \text{ private}}(X) \parallel CA$  or  $E_{CA \text{ private}}(H(X) \parallel X \parallel CA)$

Passwords - prevent leaking, avoid expensive crypt/decrypt

hashed passwords

store  $H(\text{password})$

check  $H(x) = H(\text{password})$

add a "salt"

one random salt/user

$H(\text{password} \parallel \text{salt}) \parallel \text{salt}$

But same password means

someone ~~can~~ login

multiple places w/ same pwd