

① Process abstraction

④ System Resources (CPU, mem, disk, net)

⊕ access only after letting ⇒ but process directly executed.

⊕ Get hardware support

privilege levels

traps

② Virtualization

- time sharing
- resource sharing

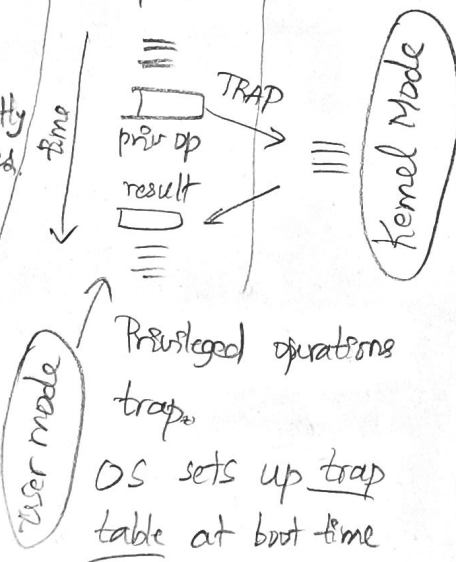
limited direct

↓ execution ↑
Isolation Performance

③ Isolation why?

- process can be buggy
- crash
- infinite loop
- process can be malicious
- take over other processes / users

⑤ trap process



⑥ trap table - hw data structure

trap#	address of fun/handler
0	→
1	→
:	

→ examples Boot time setup

- ① access the storage
- ② execute any privileged operation, in general

⑦ system calls

→ library to do req assembly

⑧ Basic limited direct exec.

(OS (kernel mode) | process (user mode))

Boot time

set up all trap handlers & trap table

Run time

initialize text, data, stack w/ args
save regs
return from trap

(nothing. All this even before process starts)

⑨ sharing between processes

processes

① cooperative approach:

- wait for sys call
- yield()

② Non cooperative approach

Timer Interrupt

hardware support once again.

⑩ Timer Interrupt

boot time kernel
→ start timer
interrupt & set up handler

(whenever choose to move to diff process)

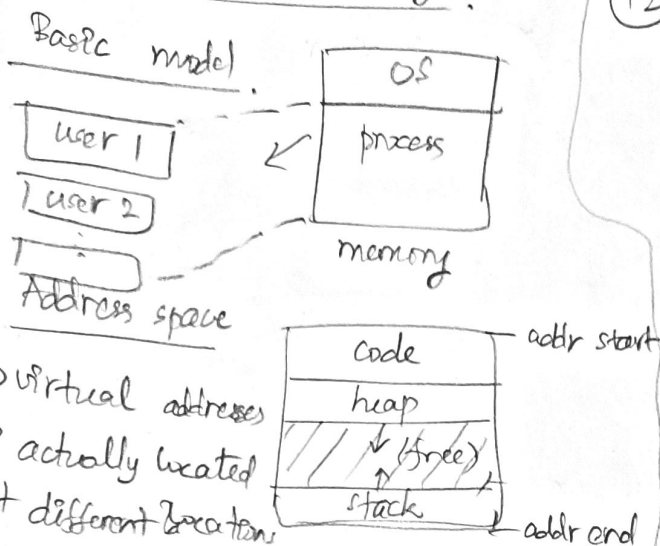
hardware



timer interrupt

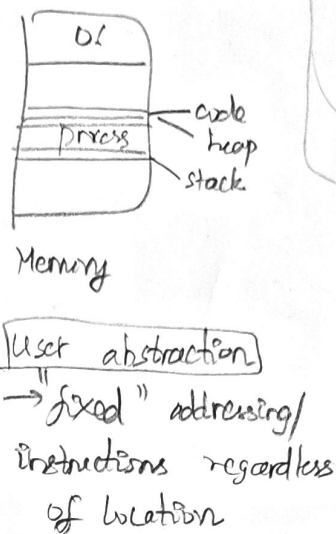
⊕ save regs to kernel stack
move to kernel mode
Jump to handler

11 Virtualizing memory.



- * virtual addresses
- * actually located at different location

12 Dynamic relocation



User abstraction
→ "fixed" addressing/
instructions regardless
of location

13 Hardware translation

* Base & bounds
load register, addr
virtual to phy addr
$$\{ \text{if } \text{base} + \text{vaddr} < \text{bound} \}$$

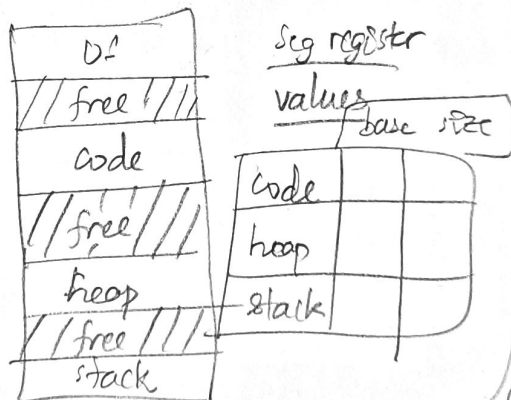
load from base + vaddr
else fault // trap
Memory Management Unit (MMU)

14 primitives for hardware address translation.

- 1 privilege modes
 - 2 base / bound registers
 - 3 check vaddr within bound
 - 4 privileged insn to update base / bounds ⇒ "sensitive"
 - 5 privileged insn to register exception handlers
 - 6 raise exceptions on invalid address
- OS must start managing memory

15 segmentation

empty mostly addr spaces
⇒ generalize to many
base and bounds.



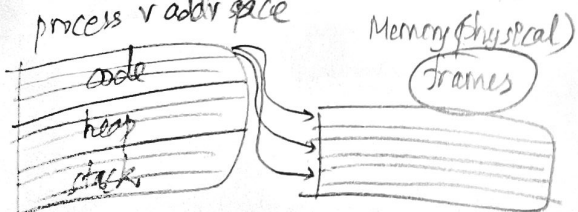
16 hardware for segmentation

- 1 tell hardware which segment?
 - 2 Memory protection bits.
read? write?
execute?
- problems

- 1 context switch
→ save & restore segdesc off- or on-disk
- 2 non-compact memory layouts!

17 Paging

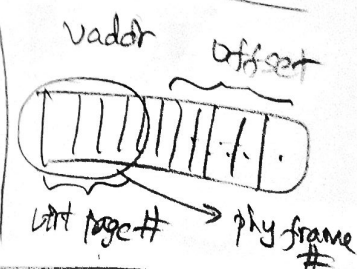
use fixed size chunks. Many of them.



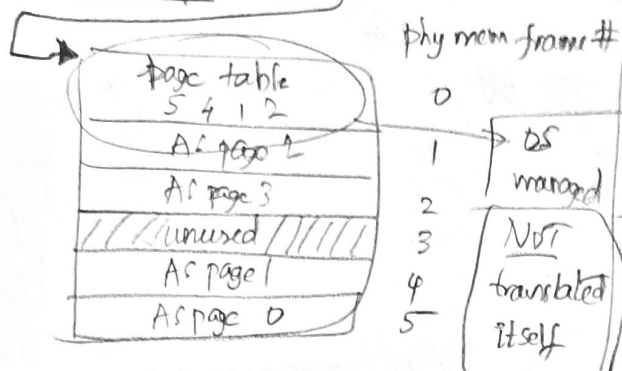
- * can't use dedicated registers any more
⇒ page table in memory.

Per process data structure

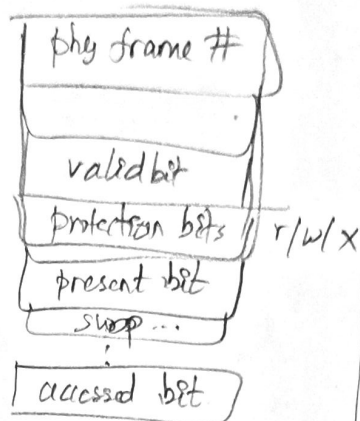
18 Translation



19 example



20 Page table entry



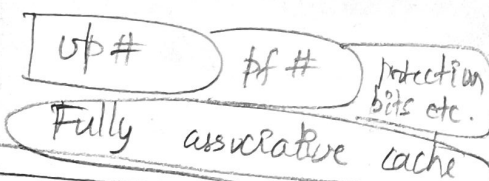
22 translation

- check valid, protection, PTE valid.
- Else trap!
 - Minor fault (retrieve shared)
 - Major fault (retrieve from disk)
 - Invalid fault (vaddr not valid)

24 Address translation cache (TLB)

- efficient memory, small → cache PTEs
- first check TLB
 - hit → use (protection etc.)
 - miss → walk page table

26 TLB entry



27 TLB + context switches

Per process data structure of pagetables

option 1 TLB flush

option 2 tagged TLB

use address space ID

29 TLBs make paging feasible

physically indexed versus virtually indexed caches

21 translation (hardware)

$$VP\# = (vaddr \& vmask) \gg \text{SHIFT}$$

$$PTE\text{addr} // \text{linear pageable} = \text{Page Table Base Ptr} + (vp\# \times \text{sizeof}(PTE))$$

translate to get phy frame #

offset = vaddr & offset mask

$$\text{Phy addr} = (pf\# \ll \text{SHIFT}) | \text{offset}$$

23 Paging is too slow

- at least two mem access
- large memory just for page tables

25 Handling TLB miss

- hardware managed TLB (Intel x86)
- software handled (use traps) (RISC)

28 smaller page tables?

multilevel page tables

page directory + table

- superpages
- hybrid segment & paging
- swap page tables to disk

30 Principles applied to entire operating systems!

- Xen
- amazon ec2

31 system virtualization

- ① OS assumes it is the most privileged entity.
- ② OS virtualization, why?
 - ⇒ time sharing uses
 - ⇒ unit of software/app deployment — any software & config can be run
- ③ actually properties of both OS & hardware arch matter (together provide isolation)

33 theorem (basic result)

$\text{efficiency virtualized if } S2 \cup S3 \subseteq S1$ } FULL
 → "trap and emulate" virtualization } virtualization
 (for privileged instructions)

- sufficient condition, x86 does not satisfy it (classically) → now Intel VT-x and AMD-V
- unmodified guest OS
- workaround tech

① Dynamic recompilation — replace critical insns at run time.

② {

- cache emulation code
- assist through hardware

② para-virtualization → port before running

③ layering Xen copy on write ← ③ union filesystems

32 virtual machine monitor

(VMM) 1974

→ virtualize all system resources (processor, mem, disk/net/I/O, ...)

3 requirements (popek & goldberg)
of interest

① equivalence

— same as direct execution on underlying phy machine

② Resource Control/safety

→ VMM must have complete control of resources

③ efficiency

→ mostly execute without VMM intervention (not a "requirement" for functionality. But need it in practice)

33 ISA classification

① privileged insns

→ must trap in user mode

② control sensitive

→ can change processor mode or "whether something will trap in 1st place" (eg. change mem bounds)

③ Behaviour sensitive

→ effect of instruction depends on mode or bounds registers

34 Containers

① kernel namespaces

(access isolation) process, network, mem

② Control groups

(resource isolation) process, mem, cpu, blkio

35 Xen

① Cpu → ring 1 of
(else trap ring 0)

② I/O virtualization

Ring buffers

aggr producer-consumer
descriptor-data separation / zero-copy.

③ Memory

types

guest OSes manage page table
writes validated by xen.
avoid effect of refresh on xen

④ Management

← separate policy from mech
dom 0.

synchronous
system calls

easy.
same handlers.

page faults
extended stack frame
(cr3 unrelatable)

hypercalls

33C Behaviour sensitive

effect of execution depends
on values of relocation
bounds R, or mode M

e.g. load directly phys addr;
move data from previous instruction
space.

33D VMM control program

- dispatcher (for traps)
- allocator (what resources to provide)
- interpreter for instructions that trap
(simulate privileged instructions)

⊕ VMM runs as supervisor, all
others run in user mode.

⊕ Hardware assumptions: trap,
relocation, MMU, privilege levels

33A

(E, M, P, R) ← Program counter (vaddr) → Memory-mapped I/O
↓ ↓ ↓
executable mode relocation register (bounds) simple dynamic relocation (base/bounds)

storage (reg+mem)
[memory trap: attempt to OOB access]
Trap: storage left unchanged except

currently active (m, p, r) storage

Privileged: (e, s, p, r) and (e, u, p, r)

↓ ↓
does not trap traps

33B

Control sensitive attempt to change
mem resources (through reloc bounds)
or affects processor mode without
a memory trap (but regular trap OK)
e.g. return to user mode.