

ADVANCED TOPICS: VIRTUAL MACHINES

Shivaram Venkataraman

CS 537, Fall 2024

ADMINISTRIVIA

Project 5 happened? → last week

Project 6 – last project!

- Early deadline this week! → Tomorrow
- Final deadline end of next week → Friday

Shivaram office hours

- TODAY at 1pm!

↳ ~ 2hrs

AGENDA / LEARNING OUTCOMES

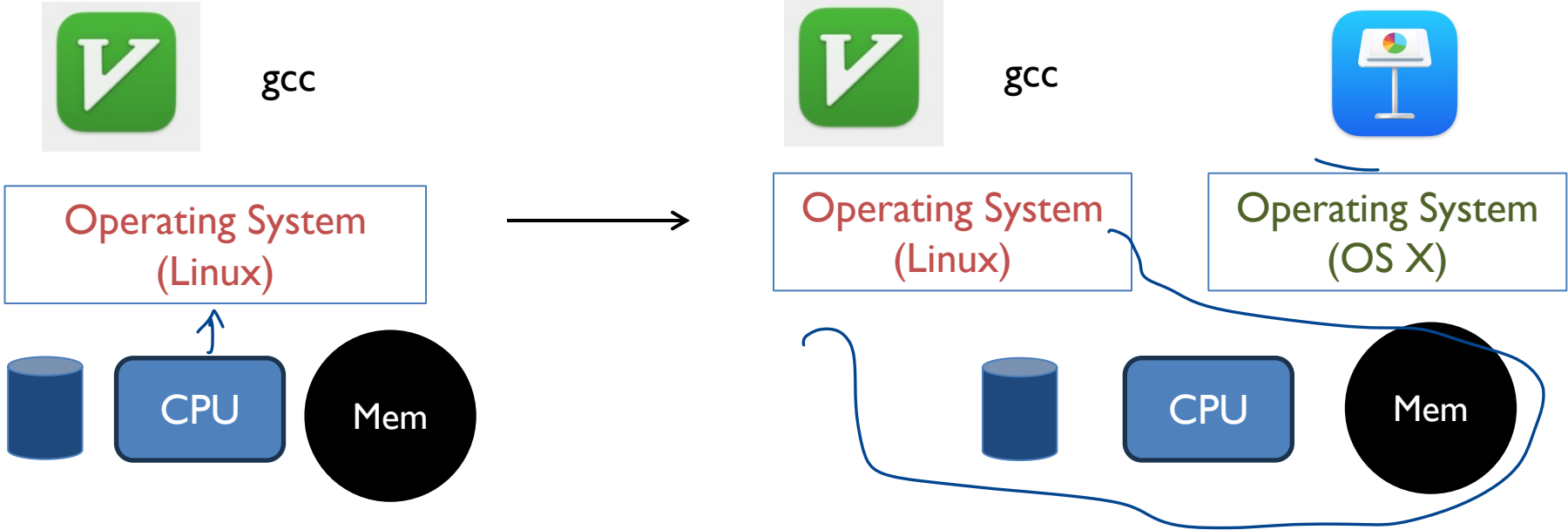
How to virtualize a machine underneath the OS?

PERSISTENCE RECAP

- Managing I/O devices significant part of OS
- Disk Drives, SSDs (pages, blocks)
- File Systems: OS provided API to access disk
- Simple FS: FS layout with superblock, bitmaps, inodes, datablocks
- Fast File System: Key idea – put inode & data close together, namespace locality
- FSCK, Journaling – Handling/Preventing data inconsistencies
- Log Structured File System - Organize data based on writes

} → Project 6

VIRTUAL MACHINES



VIRTUAL MACHINE USE CASES

Share mainframe systems (1970s)

Cloud Computing

- Consolidate multiple tenants running different OS
- Strong Isolation

↳ tenants do not interfere with each

Run applications that only exist for specific OS

Testing, Debugging

↳ want to check if your code works on a diff OS.

80s/90s

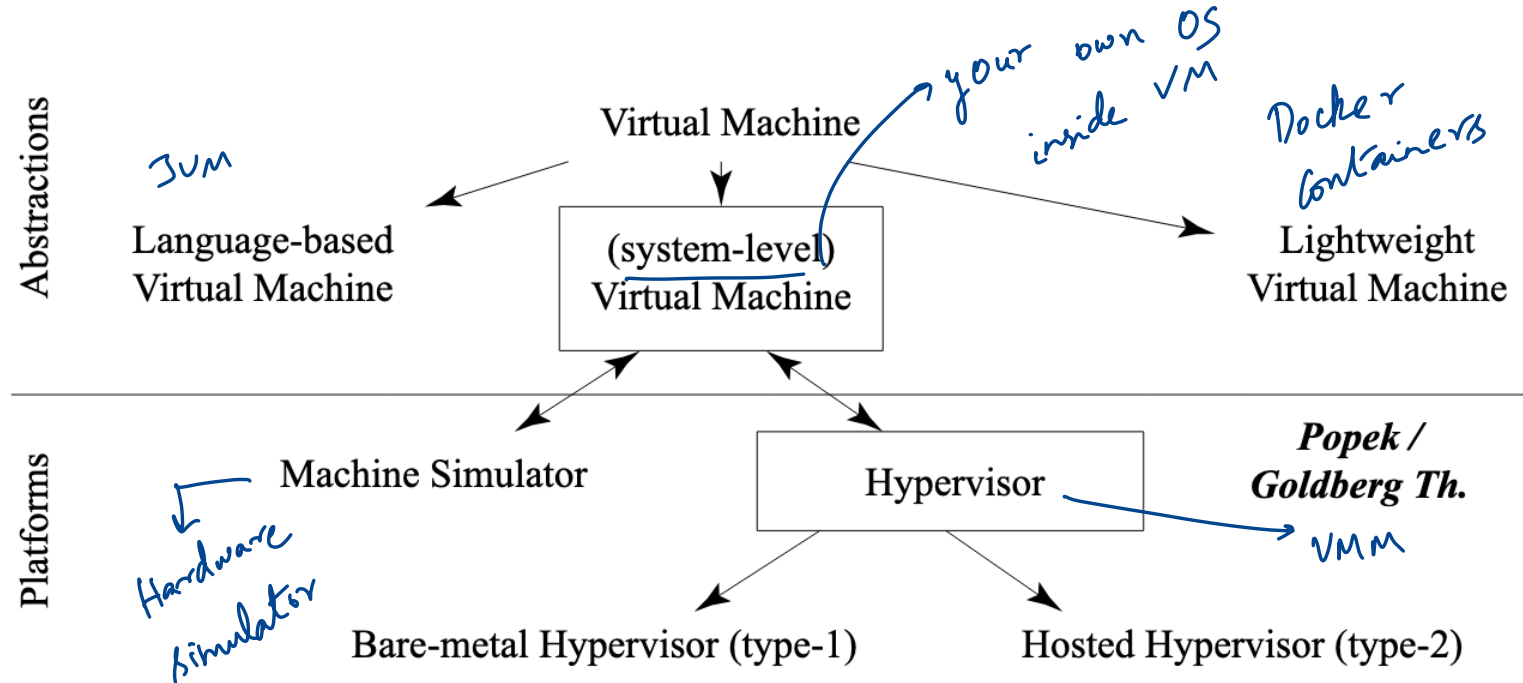
personal computing

tenants
↓
computer



DEFINITIONS

A virtual machine is a **complete compute environment** with its own isolated processing capabilities, memory, and communication channels.



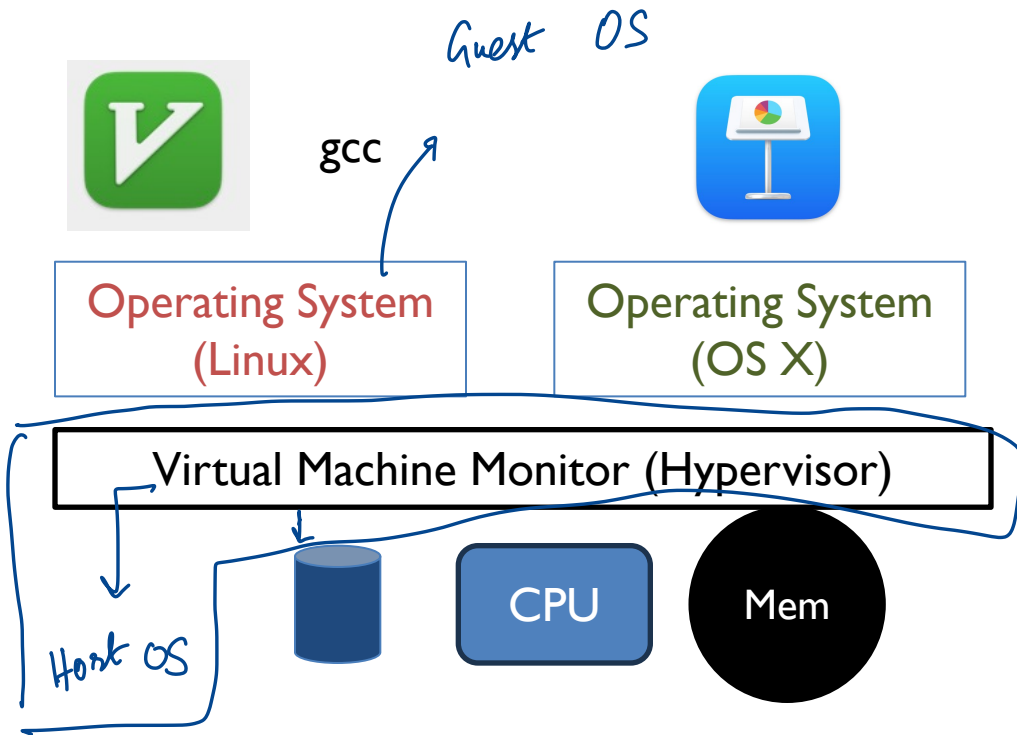
VIRTUAL MACHINE MONITORS

Bare-metal Hypervisor (type-1)
direct control of all resources

xen

Hosted Hypervisor (type-2)
operates as part of or on top of an existing host OS

KVM → part of linux



GOALS

- Equivalence – The exposed resource is equivalent with the underlying computer.
- Safety – Isolation requires that the virtual machines are isolated from each other as well as from the hypervisor.
Diff guest OS
- Performance – The virtual system must show at worst a minor decrease in speed.

CAN WE VIRTUALIZE? (POPEK GOLDBERG 1974)

The processor's system state, called the processor status word (PSW) consists of the tuple $(\underline{M}, \underline{B}, \underline{L}, \underline{PC})$:

- the execution level $M = \{s, u\}$ (superuser or user mode)
- the segment register (B, L) ; (Segmented Memory Model) and
- the current program counter (PC) , a virtual address

kernel or user mode

*B → base
L → limit*

A virtual machine monitor may be constructed if the set of sensitive instructions for a computer is a subset of the set of privileged instructions.

$\{control-sensitive\} \cup \{behavior-sensitive\} \subseteq \{privileged\}$.

*↓
instructions which
can change
state*

*↓
can be
modified by
state*

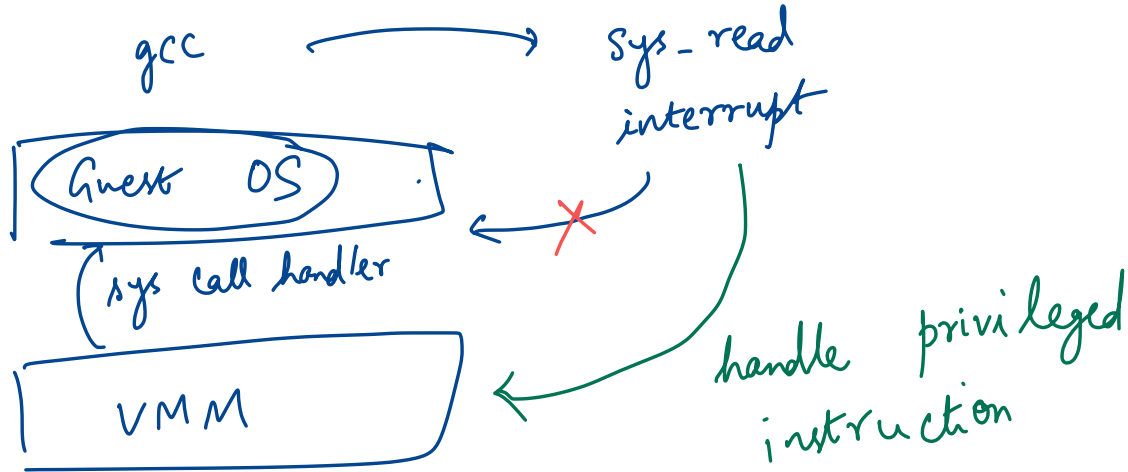
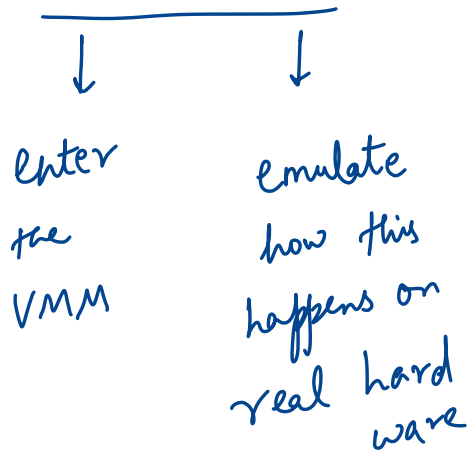
*↳ can only be
run in
kernel mode*

VIRTUALIZING THE CPU

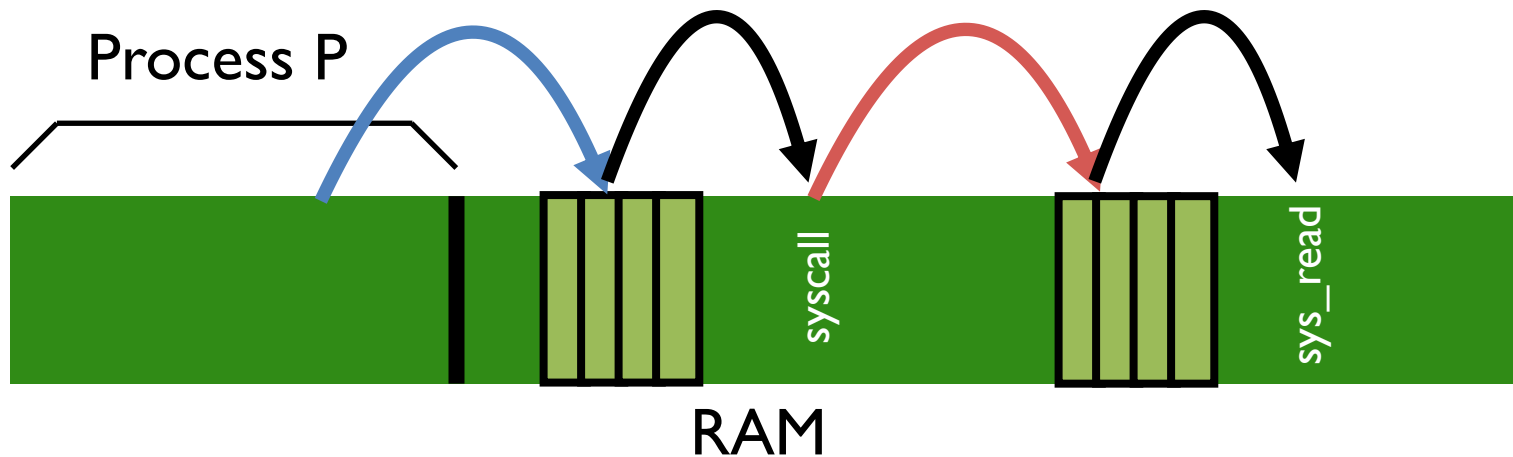
Limited Direct Execution

How to handle privileged instructions (e.g., traps for system calls) ?

Trap and Emulate!



BEFORE: SYSTEM CALL FLOW

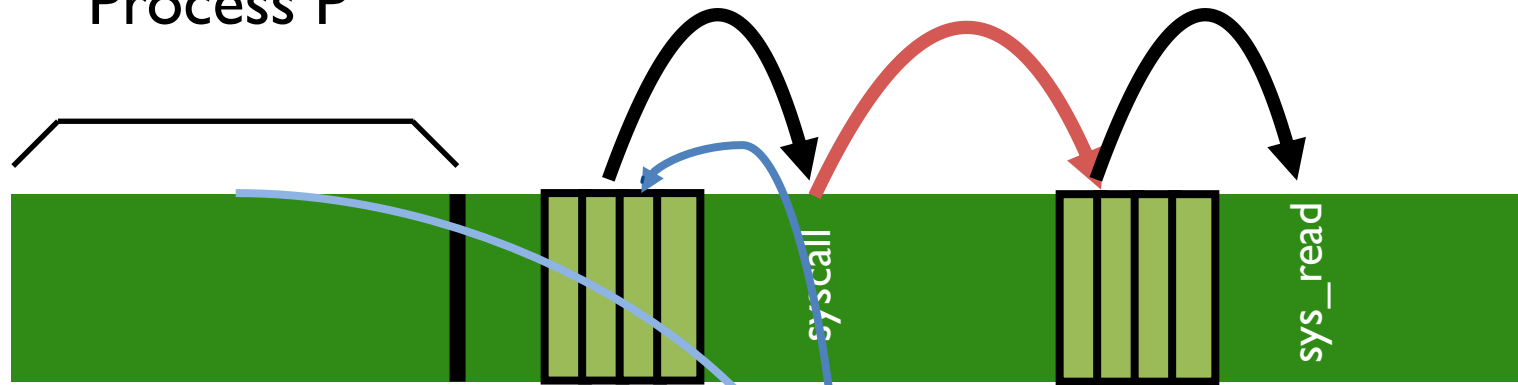


```
movl $6, %eax;   int $64
```

Transfer control to trap handler. Execute appropriate syscall routine

NEW: SYSTEM CALL

Process P



```
movl $6, %eax;
```

```
int $64
```

Guest

OS boot up

↳ privileged inst to install interrupt handler

Virtual Machine Monitor

USER MODE, KERNEL MODE?

MIPS architecture:

- Guest OS runs in “supervisor” mode
- No privileged instructions, some extra memory

→ user -
→ supervisor -
→ kernel -

↳ some extra memory
↳ guest OS data structure

Run Guest OS in user mode

How to protect Guest OS data structures?

↳ from the other processes inside Guest OS

QUIZ 20

Log structured SSD consisting of 3 blocks and 10 pages per block. Each page holds a single character.

The state of each page (i, v, or E), the data stored at each page, and an indicator if a page is currently live (i.e. has a mapping in the FTL).

- read(page#) -- if page is live returns the character at the page, otherwise error
- write(page#,char) -- writes character to logical page #
- erase(page#) -- removes logical page # from the FTL mapping



FTL	0: 15	2: 18	3: 8	4: 4
	14: 16			
Block	0	1	2	
Page	0000000000	1111111111	2222222222	
	0123456789	0123456789	0123456789	
State	vvvvvvvvvv	vvvvvvvvvv	iiiiiiiiiii	
Data	c9XhFAp970	CqFuArsJE		
Live	+ +	++ +		

start

write "t" to
page 4

FTL	0: 15	2: 18	3: 8	4: 19
	14: 16			
Block	0	1	2	
Page	0000000000	1111111111	2222222222	
	0123456789	0123456789	0123456789	
State	vvvvvvvvvv	vvvvvvvvvv	iiiiiiiiiii	
Data	c9XhFAp970	CqFuArsJE	t	
Live	+ +	++ ++		

new

Page 19
has some data
now "t"

If a write(0,'q') is now performed by the OS on the SSD state from the last question, what underlying SSD operations must be performed in order to accomplish this write?

FTL	0: 15 14: 16	2: 18	3: 8	4: 19
Block	0	1	2	
Page	0000000000	1111111111	2222222222	
	0123456789	0123456789	0123456789	
State	v v v v v v v v v v	v v v v v v v v v v	i i i i i i i i i i	
Data	c9XhFAp970	CqFuArsJEt	q	
Live	+	•+ ++	+	

Erase on block 2
Store q in page 20

E E E

FTL	0: 15	2: 18	<u>3: 8</u>	<u>4: 4</u>
	14: 16			
Block	0	1	2	
Page	0000000000	1111111111	2222222222	
	0123456789	0123456789	0123456789	
State	vvvvvvvvvv	vvvvvvvvvE	iiiiiiiiiii	
Data	c9XhFAp970	CqFuArsJE		
Live	+ +	++ +		

initial state

Garbage Collector

FTL	0: 15	2: 18	<u>3: 21</u>	<u>4: 20</u>
	14: 16			
Block	0	1	2	
Page	0000000000	1111111111	2222222222	
	0123456789	0123456789	0123456789	
State	EEEEEEEEEE	vvvvvvvvvE	vvEEEEEEEEEE	
Data	c9XhFAp970	CqFuArsJE	<u>F7</u>	
Live	+ +	++ +	++	

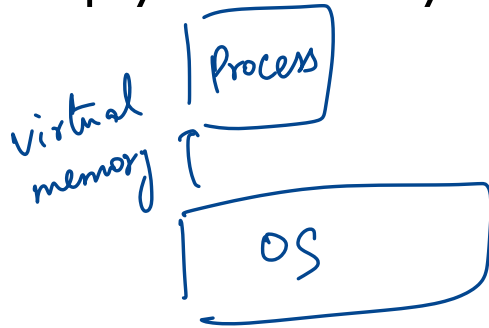
erased

live data

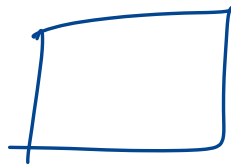
VIRTUALIZING MEMORY

Challenge: Who manages physical memory allocation?

How do we share physical memory across Guest OSes?



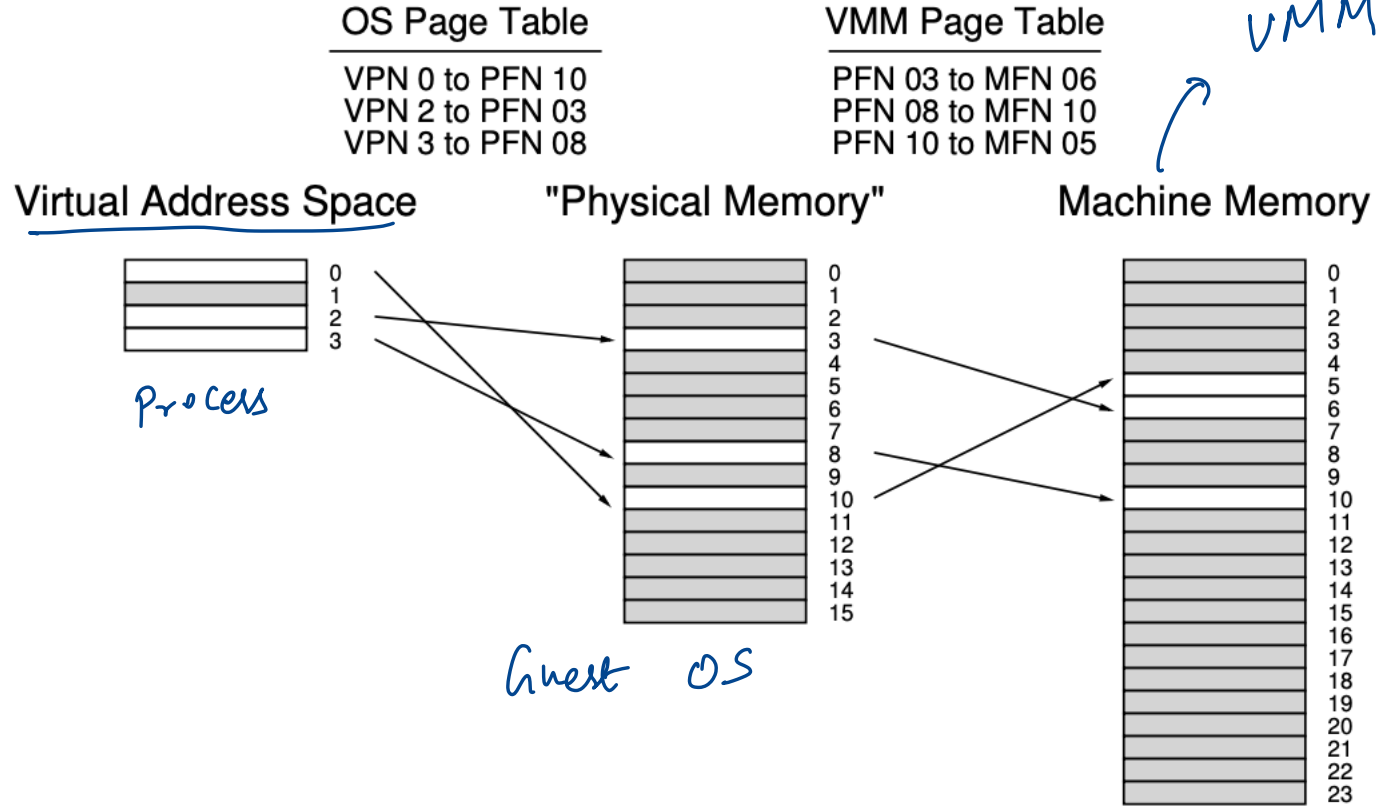
alloc →



Physical memory

Page tables
Virtual → Physical

Extra level of indirection!



BEFORE: SOFTWARE TLB HANDLER

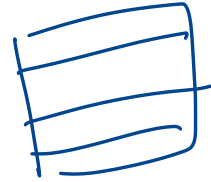
TLB miss in hardware

Trap into OS
OS walks pagetable
Get Virtual → Physical
Update TLB using
privileged instruction

flush the
tlb the / update
 the
 TLB

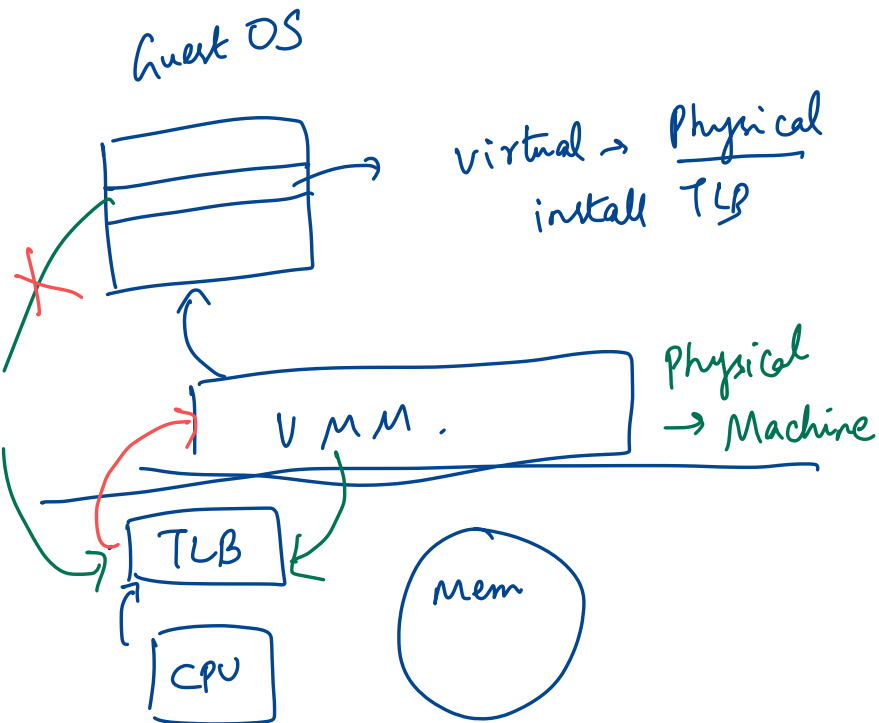
Application

TLB



Virtual
Physical

NEW: SOFTWARE TLB HANDLER



TLB miss

2 traps

Trap into VMM
Call OS Handler

OS walks pagetable
Get **Virtual** → **Physical**
Update TLB using
privileged instruction

Trap handler
Physical → **Machine**
Update TLB

TLB MISS OVERHEADS

Extra trap into VMM for Physical → Machine mapping

Avoid using Software “TLB” in VMM to cache Virtual → Physical

→ Part of Page tables from Guest OS

Hardware managed TLBs

VMM maintains Shadow page table per of Virtual → Machine
Trap when OS tries to update PTE (e.g., lcr3)

→ x86

for every PT in guest OS
create a shadow PT

Guest OS → update PTE
↓ trap into the VMM

SO, CAN WE VIRTUALIZE X86?

expose hardware / processor state

Table 2.2: List of sensitive, unprivileged x86 instructions

Group	Instructions
Access to interrupt flag	pushf, popf, iret
Visibility into segment descriptors	lar, verr, verw, lsl
Segment manipulation instructions	pop <seg>, push <seg>, mov <seg>
Read-only access to privileged state	sgdt, sldt, sidt, smsw
Interrupt and gate instructions	fcall, longjump, retfar, str, int <n>

PARA VIRTUALIZATION, X86 EXTENSIONS

So far: No change to the guest OS. No changes to the hardware.

Downside: Overheads can be quite high?

Para virtualization

Can we make (small?) modifications to the guest OS for efficiency?

Hardware

Instruction set extensions (Intel, AMD)

↳ virtualization friendly x86



XEN

early 2000s

Modify guest OS: simply undefine all of the 17 non-virtualizable instructions!
 Alternate interrupt architecture

Memory Management	
Segmentation	Cannot install fully privileged segment descriptors and cannot overlap with the top end of the linear address space.
Paging	Guest OS has <u>direct read access</u> to hardware page tables, but <u>updates</u> are <u>batched</u> and validated by the hypervisor. A domain may be allocated <u>discontinuous machine (aka host-physical) pages</u> .
CPU	
Protection	Guest OS must run at a lower privilege level than Xen.
Exceptions	Guest OS must register a descriptor table for exception handlers with Xen. Aside from page faults, the handler remains the same.
System calls	<u>Guest OS</u> may install a <u>“fast” handler</u> for system calls, allowing direct calls from an application into its guest OS and avoiding indirection through Xen on every call.

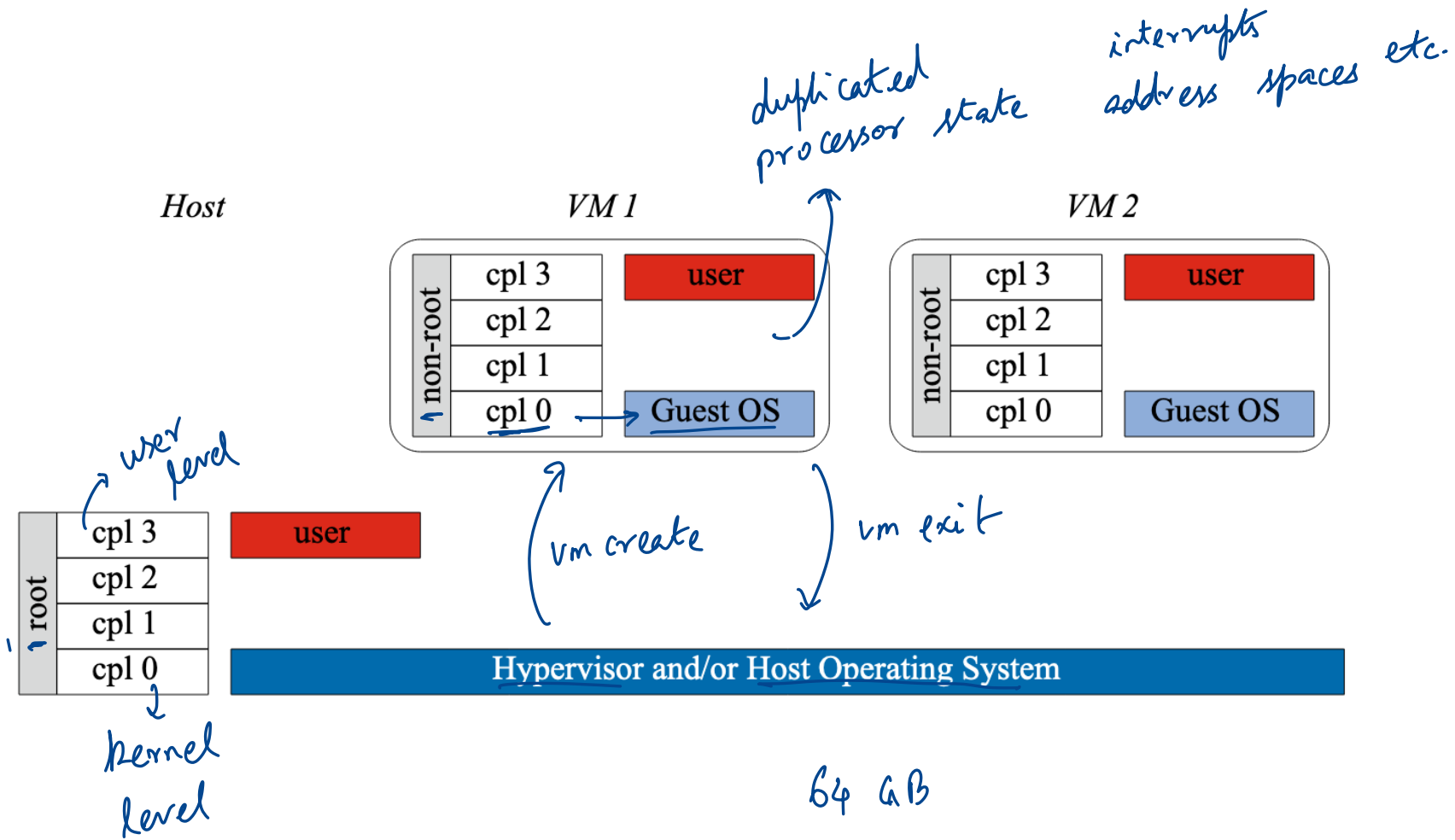
INTEL VT-X EXTENSIONS

True Hardware Support meeting Popek / Goldberg Criteria

Do not change the semantics of individual instructions, instead duplicate the entire visible state and introduce a **new mode of execution**: the root mode.

- Hypervisor is in root mode, Guest OS in non-root mode.
- Special new instructions for detecting mode (only available in root mode, otherwise a trap is caused).
- New mode only used for virtualization
- Each mode has own address space
- Each mode has own interrupt flag

CR3 → each mode has its own PTBR



Next class: Multi-CPU scheduling

Thanksgiving break!