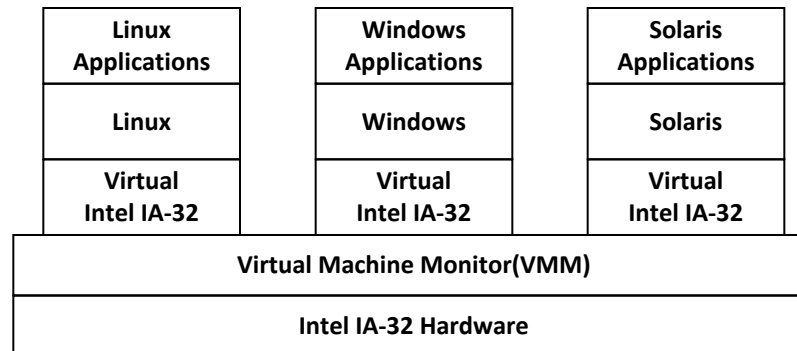


Lecture 12

Virtual Machines

System Virtual Machines

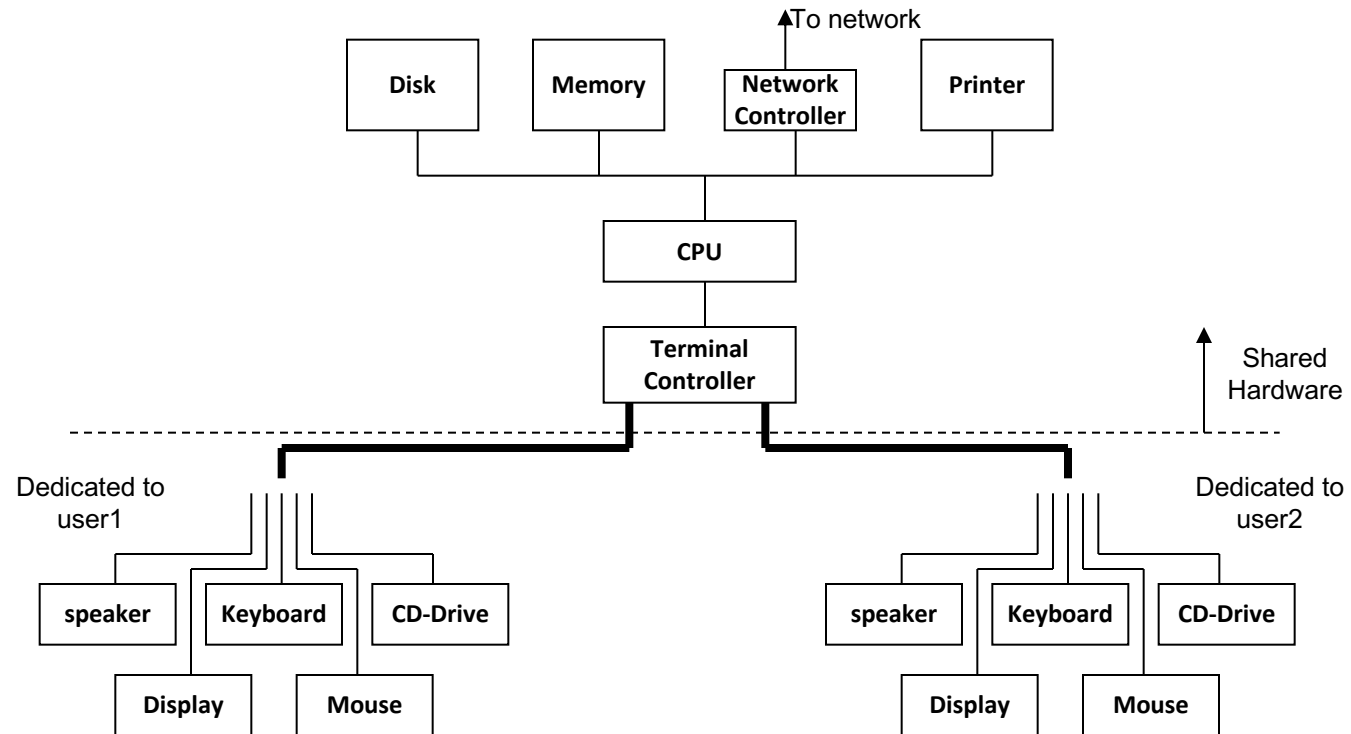
- A **(full) system virtualization** environment is capable of supporting multiple system images simultaneously, each running its own operating system and associated application programs
- Real resources of the *host* platform are shared among the *guest* system with the virtual machine monitor (VMM, a.k.a. **hypervisor**)
- Focus on VMs where ISA of the host and guest are the same



Other types of virtualizations

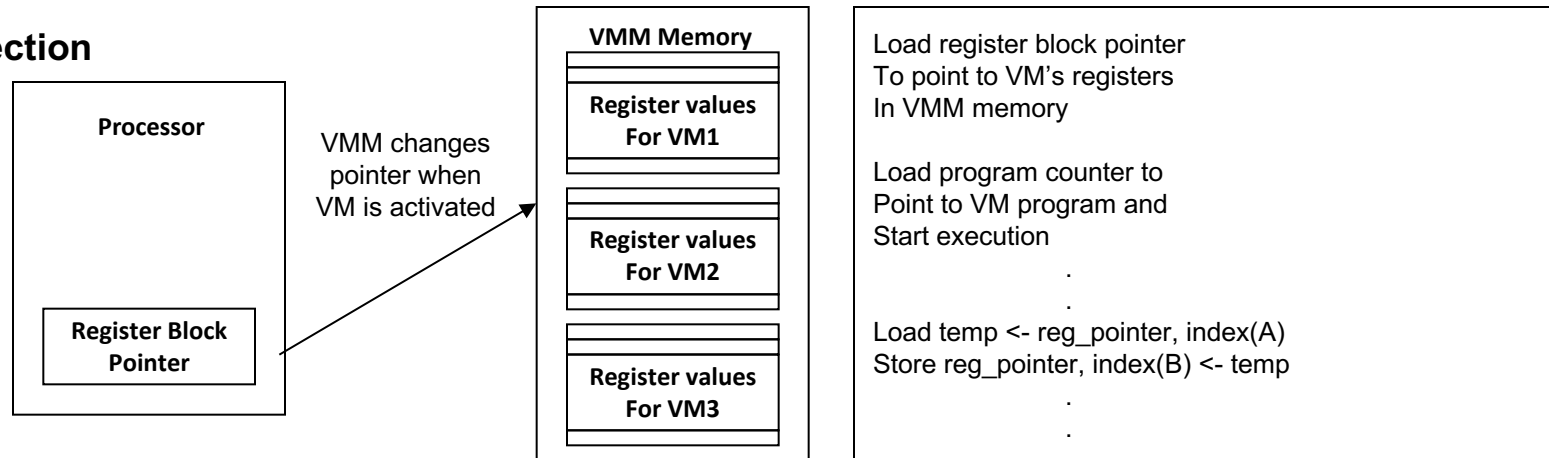
- High level language virtualization
 - For cross-platform compatibility, “write once, run everywhere”
 - Java Virtual Machine
- Cross platform emulation
 - Qemu, Apple Rosetta
- OS-level virtualization
 - Multiple isolated user space instances
 - Docker

Outward Appearance

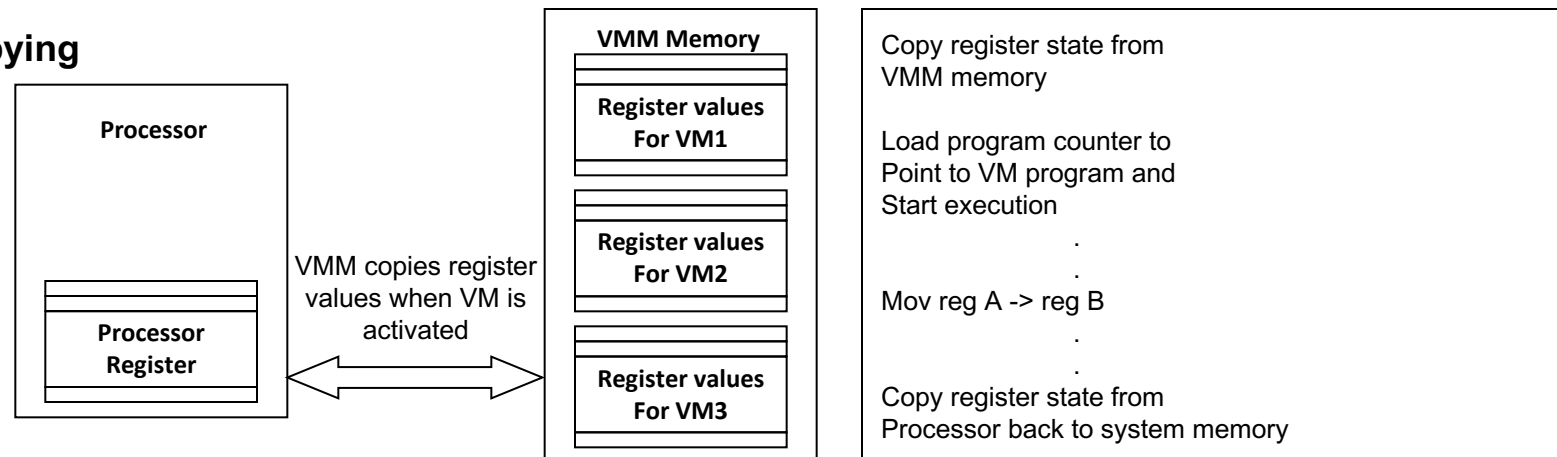


State Management

Indirection

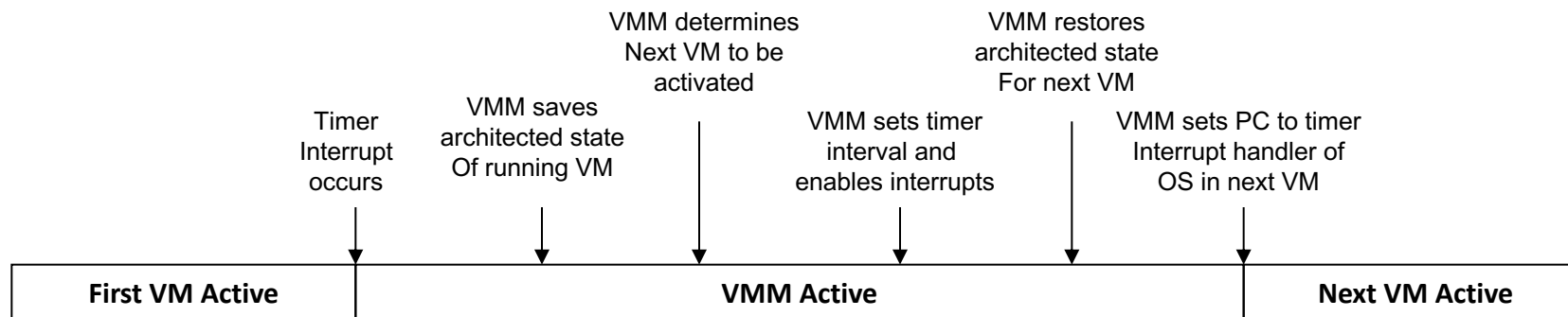


Copying



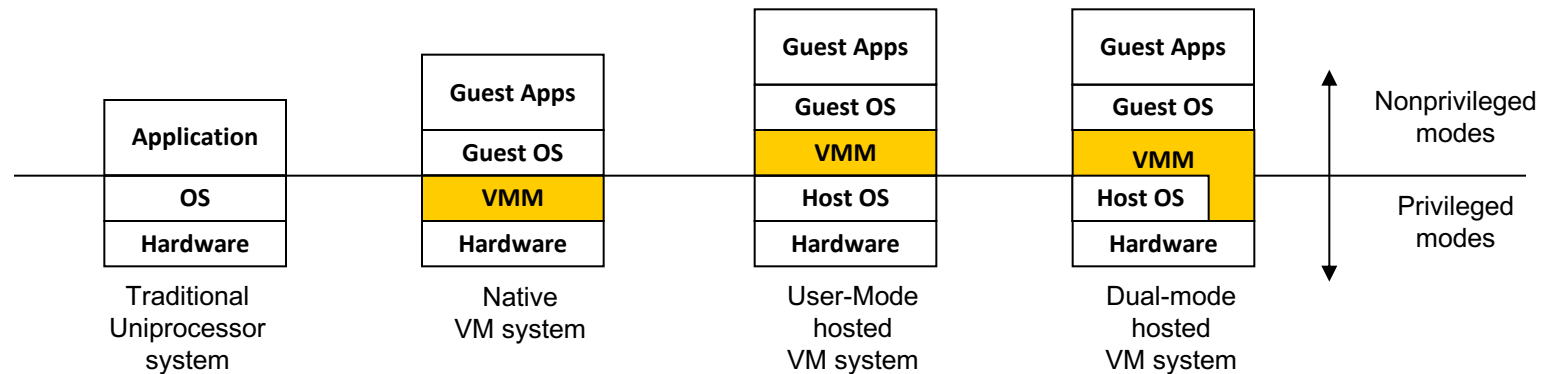
Resource Control

- Adopt time-sharing systems
- The VMM maintain overall control of all the hardware resources
- The interval timer interrupt
 - Guarantee that control is transferred back to VMM and the VMM handles the interrupt itself



Native and Hosted Virtual Machine

- Native VM system
 - A system where a VMM operates in a privilege mode higher than the mode of the guest virtual machines
 - The privilege level of the guest OS is emulated by the VMM
- Hosted VM system
 - The VMM is installed on a host platform that is already running an existing OS
 - The VMM utilizes the functions already available on the host OS to control and manage resources desired by each of the virtual machine



Resource Virtualization - Processor

- The key aspect of virtualizing a processor is the execution of the guest instructions, including both system-level and user-level instruction
- Processor virtualization method
 - Emulation
 - Interpretation, binary translation
 - Emulation is the only processor virtualization mechanism for different ISA between the guest and the host
 - Direct native execution
 - Only if the ISA of the host is identical to the ISA of the guest
 - The goal is to gain the same performance as it runs on the virtual machine

Resource Virtualization - Processor

Conditions for ISA Virtualizability

- We restrict the discussion here to native system VMs
- In a native system VM, the VMM runs in system mode, and all other software runs in user mode
- The VMM keeps track of the intended mode of operation of a guest virtual machine but it will always use user mode in executing the instructions from the guest virtual machine

Resource Virtualization - Processor Conditions for ISA Virtualizability

- A **privileged instruction** is defined as one that traps if the machine is in user mode and does not trap if the machine is in system mode
 - Example: WRMSR, WRGSBASE

WRMSR—Write to Model Specific Register

Opcode	Instruction	Op/En	64-Bit Mode	Compat/Leg Mode	Description
0F 30	WRMSR	Z0	Valid	Valid	Write the value in EDX:EAX to MSR specified by ECX.

Instruction Operand Encoding

Op/En	Operand 1	Operand 2	Operand 3	Operand 4
Z0	NA	NA	NA	NA

Protected Mode Exceptions

#GP(0)

If the current privilege level is not 0.

If the value in ECX specifies a reserved or unimplemented MSR address.

If the value in EDX:EAX sets bits that are reserved in the MSR specified by ECX.

If the source register contains a non-canonical address and ECX specifies one of the following MSRs: IA32_DS_AREA, IA32_FS_BASE, IA32_GS_BASE, IA32_KERNEL_GS_BASE, IA32_LSTAR, IA32_SYSENTER_EIP, IA32_SYSENTER_ESP.

Resource Virtualization - Processor

Conditions for ISA Virtualizability

- To specify instructions that interact with hardware, three categories of special instructions are defined
 - Control-sensitive instruction
 - Attempt to change the configuration of resources in the system
 - For example WRMSR
 - Behavior-sensitive instruction
 - Behavior or results produced depend on the configuration of resource
 - For example Pop Stack into Flags Register (POPF)
 - Innocuous instruction
 - An instruction which is neither control sensitive or behaviour sensitive

Intel x86 POPF instruction

Table 4-21. Effect of POPF/POPFD on the EFLAGS Register

Mode	Operand Size	CPL	IOPL	Flags																		Note
				21	20	19	18	17	16	14	13:12	11	10	9	8	7	6	4	2	0		
				ID	VIP	VIF	AC	VM	RF	NT	IOPL	OF	DF	IF	TF	SF	ZF	AF	PF	CF		
Real-Address Mode (CR0.PE = 0)	16	0	0-3	N	N	N	N	N	0	S	S	S	S	S	S	S	S	S	S	S		
	32	0	0-3	S	N	N	S	N	0	S	S	S	S	S	S	S	S	S	S	S		
Protected, Compatibility, and 64-Bit Modes (CR0.PE = 1 EFLAGS.VM = 0)	16	0	0-3	N	N	N	N	N	0	S	S	S	S	S	S	S	S	S	S	S		
	16	1-3	<CPL	N	N	N	N	N	0	S	N	S	S	N	S	S	S	S	S	S		
	16	1-3	≥CPL	N	N	N	N	N	0	S	N	S	S	S	S	S	S	S	S	S		
	32, 64	0	0-3	S	N	N	S	N	0	S	S	S	S	S	S	S	S	S	S	S		
	32, 64	1-3	<CPL	S	N	N	S	N	0	S	N	S	S	N	S	S	S	S	S	S		
	32, 64	1-3	≥CPL	S	N	N	S	N	0	S	N	S	S	S	S	S	S	S	S	S		
Virtual-8086 (CR0.PE = 1 EFLAGS.VM = 1 CR4.VME = 0)	16	3	0-2	X	X	X	X	X	X	X	X	X	X	X	X	X	X	X	X	X	1	
	16	3	3	N	N	N	N	N	0	S	N	S	S	S	S	S	S	S	S	S		
	32	3	0-2	X	X	X	X	X	X	X	X	X	X	X	X	X	X	X	X	X	1	
	32	3	3	S	N	N	S	N	0	S	N	S	S	S	S	S	S	S	S	S		
VME (CR0.PE = 1 EFLAGS.VM = 1 CR4.VME = 1)	16	3	0-2	N/ X	N/ X	SV/ X	N/ X	N/ X	0/ X	S/ X	N/X	S/ X	S/ X	N/ X	S/ X	S/ X	S/ X	S/ X	S/ X	S/ X	2,3	
	16	3	3	N	N	N	N	N	0	S	N	S	S	S	S	S	S	S	S	S		
	32	3	0-2	X	X	X	X	X	X	X	X	X	X	X	X	X	X	X	X	X	1	
	32	3	3	S	N	N	S	N	0	S	N	S	S	S	S	S	S	S	S	S		

NOTES:

1. #GP fault - no flag update
2. #GP fault with no flag update if VIP=1 in EFLAGS register and IF=1 in FLAGS value on stack
3. #GP fault with no flag update if TF=1 in FLAGS value on stack

Key	
S	Updated from stack
SV	Updated from IF (bit 9) in FLAGS value on stack
N	No change in value
X	No EFLAGS update
0	Value is cleared

Handling Problem Instructions

- The POPF in IA32 instruction is sensitive but not privileged
 - It is a **critical instruction** (sensitive but not privileged)
 - It does not generate a trap in user mode
 - It violate the virtualizability condition of Theorem 1
- An additional set of steps must be taken in order to implement a system virtual machine (with possible loss of some efficiency)
 - It is possible for a VMM intercepts POPF and other critical instructions if all guest software were interpreted instruction by instruction

Processor Virtualization in IA-32

- IA-32 architecture is not efficiently virtualizable
 - There are 17 instructions that are critical (sensitive, but not privileged)
- The VM needs to be hybridized
 - Runs natively, but needs scanning and patching of critical instructions first

IA-32 Critical Instructions

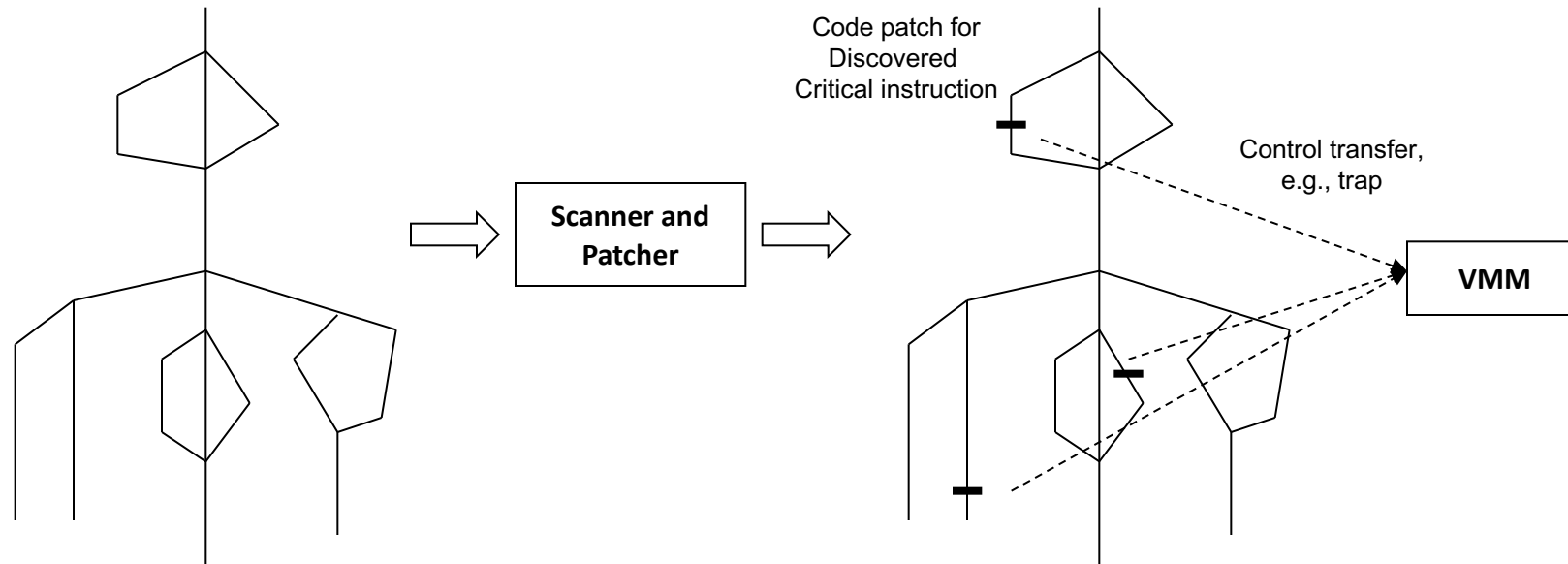
- The 17 critical instructions fall into 2 broad categories:
 - Protection system references – instructions that reference the storage protection system, memory system, or address relocation system
 - E.g., `mov ax, cs` in user mode generates a no-op
 - Sensitive register instructions – instructions that attempt to read or change resource-related registers and memory locations
 - E.g., `popf`, which could potentially change IF flag

Resource Virtualization – Processor Handling Problem Instructions

- The VMM scans the guest code stream to discover all critical instructions and replace them with a trap or jump to VMM
- The process above is called *patching*

Resource Virtualization – Processor Handling Problem Instructions

- Patching

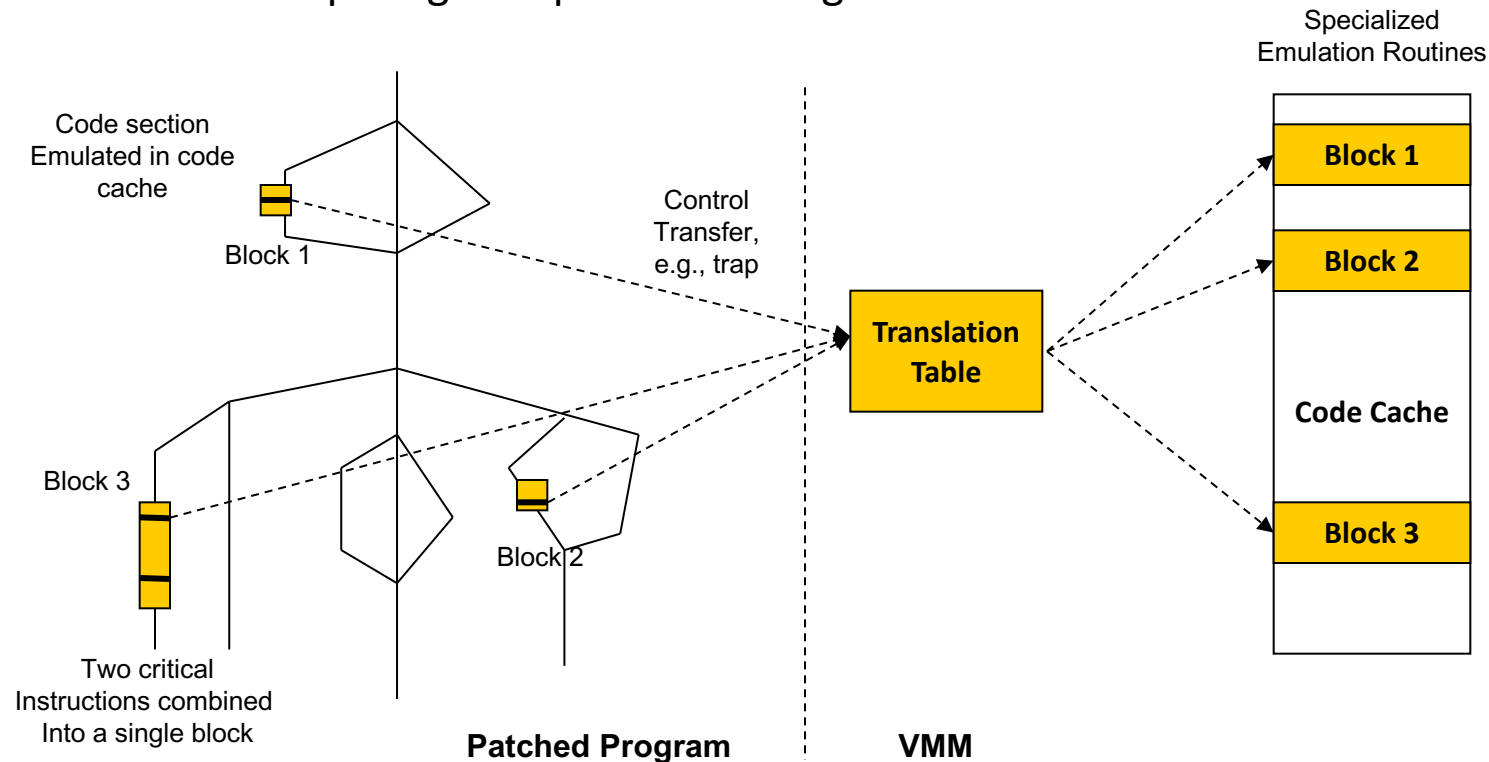


Resource Virtualization – Processor Patching of Critical Instructions

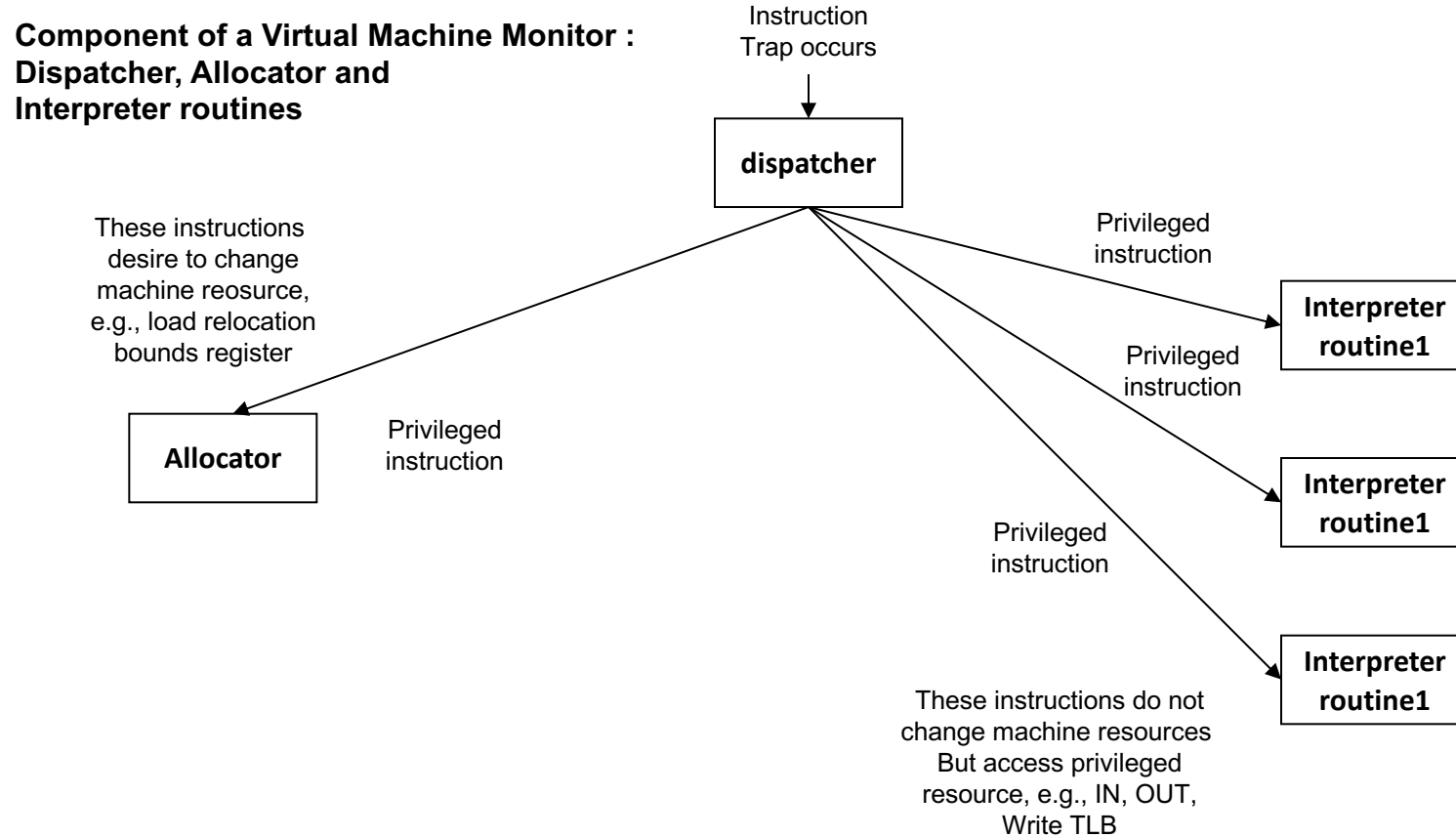
- One way to discover critical instructions
 - The VMM takes control at the head of each guest basic block and scan instructions in sequence until the end of the basic block is reached
 - If a critical instruction is found, it is replaced with a trap to the VMM
 - Another trap back to the VMM is placed at the end of the basic block
- To reduce overhead, the trap at the end of a scanned basic block can be replaced by the original branch or jump instruction

Resource Virtualization – Processor Caching Emulation Code

- Reduce the overhead of VMM interpretation when the frequency of sensitive instructions requiring interpretation is high



Conditions for ISA Virtualizability



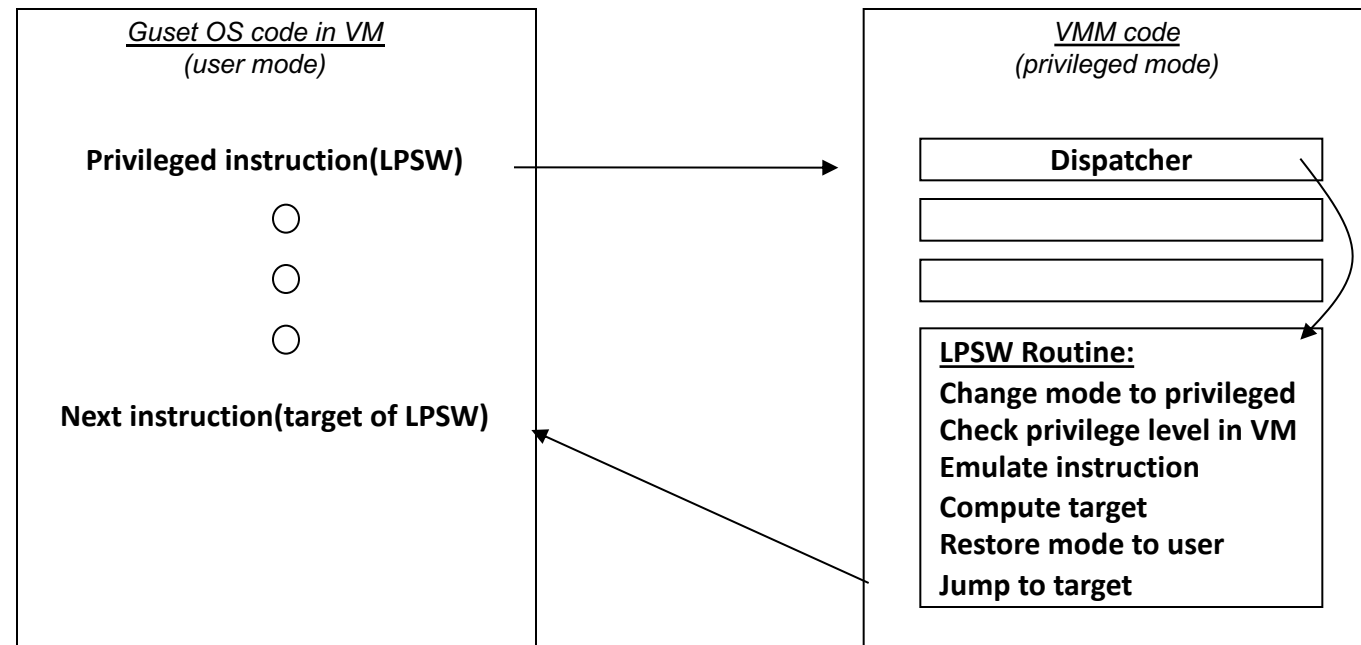
Popek and Goldberg Virtualization Requirements

- A potential virtual machine must satisfy three properties
 - Efficiency
 - Resource control
 - Equivalence
- **Theorem 1** regarding (efficient) VMM construction
 - A virtual machine monitor may be constructed if the set of sensitive instruction is a subset of the set of privileged instructions
 - It means that an efficient virtual machine implementation can be constructed if instructions that could interfere with the functioning of the VMM always trap in the user mode

Resource Virtualization - Processor

Conditions for ISA Virtualizability

- The VMM interprets a sensitive instruction according to the prevailing status of the virtual system resources and the state of the virtual machine



A particularly nasty issue: keeping time

- How to create the illusion of virtual continuous time?
 - Solution: to adjust for the difference
- Example: the **pvclock** protocol of Linux
 - A simple per-CPU structure that is shared between the host and the guest
 - To get the current timestamp counter (TSC) reading, guests must do:

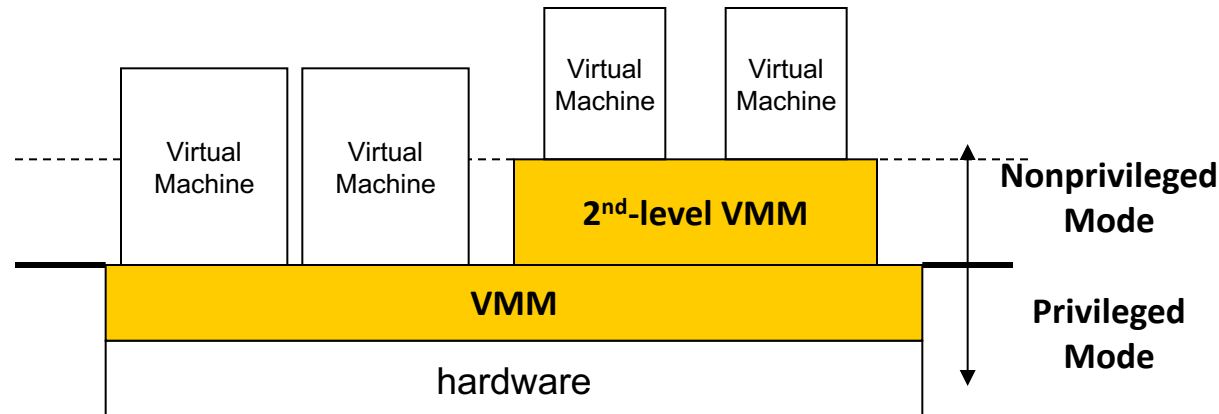
$$\text{PerCPUTime} = ((\text{RDTSC}() - \text{tsc_timestamp}) \gg \text{tsc_shift}) * \text{tsc_to_system_mul} + \text{system_time}$$

```
/ arch / x86 / include / asm / pvclock-abi.h

1  /* SPDX-License-Identifier: GPL-2.0 */
2  #ifndef _ASM_X86_PVCLOCK_ABI_H
3  #define _ASM_X86_PVCLOCK_ABI_H
4  #ifndef __ASSEMBLY__
5
6  /*
7   * These structs MUST NOT be changed.
8   * They are the ABI between hypervisor and guest OS.
9   * Both Xen and KVM are using this.
10  */
11  * pvclock_vcpu_time_info holds the system time and the tsc timestamp
12  * of the last update. So the guest can use the tsc delta to get a
13  * more precise system time. There is one per virtual cpu.
14  *
15  * pvclock_wall_clock references the point in time when the system
16  * time was zero (usually boot time), thus the guest calculates the
17  * current wall clock by adding the system time.
18  *
19  * Protocol for the "version" fields is: hypervisor raises it (making
20  * it uneven) before it starts updating the fields and raises it again
21  * (making it even) when it is done. Thus the guest can make sure the
22  * time values it got are consistent by checking the version before
23  * and after reading them.
24  */
25
26  struct pvclock_vcpu_time_info {
27      u32 version;
28      u32 pad0;
29      u64 tsc_timestamp;
30      u64 system_time;
31      u32 tsc_to_system_mul;
32      s8 tsc_shift;
33      u8 flags;
34      u8 pad[2];
35  } __attribute__((__packed__)); /* 32 bytes */
36
37  struct pvclock_wall_clock {
38      u32 version;
39      u32 sec;
40      u32 nsec;
41  } __attribute__((__packed__));
42
43  #define PVCLOCK_TSC_STABLE_BIT (1 << 0)
44  #define PVCLOCK_GUEST_STOPPED (1 << 1)
45  /* PVCLOCK_COUNTS_FROM_ZERO broke ABI and can't be used anymore. */
46  #define PVCLOCK_COUNTS_FROM_ZERO (1 << 2)
47  #endif /* __ASSEMBLY__ */
48  #endif /* _ASM_X86_PVCLOCK_ABI_H */
```

Resource Virtualization – Processor Recursive Virtualization

- The concept of running the virtual machine system on a copy itself



- Two effects that usually restrict the ability to create an efficient recursively virtualizable system
 - Timing dependancies
 - Memory allocation

Processor Recursive Virtualization

Theorem 2

- A conventional (third-generation) computer is recursively virtualizable if
 - it is virtualizable and
 - a VMM without any **timing dependences** can be constructed for it
- “Timing dependences” = none of the instructions’ execution depends on timing constraints, i.e., no real-time-ness

Resource Virtualization – Input/Output

Virtualizing Device

- Dedicated Devices
 - Some I/O device is dedicated to a particular guest VM or at least are switched from one guest to another on a very long time scale
 - The device itself does not necessarily have to be virtualized
 - Requests to and from the device could theoretically bypass the VMM and go directly to the guest operating system
- Partitioned Device
 - A very large disk, for example, can be partitioned into several smaller virtual disk that are then made available to the virtual machine as dedicated devices

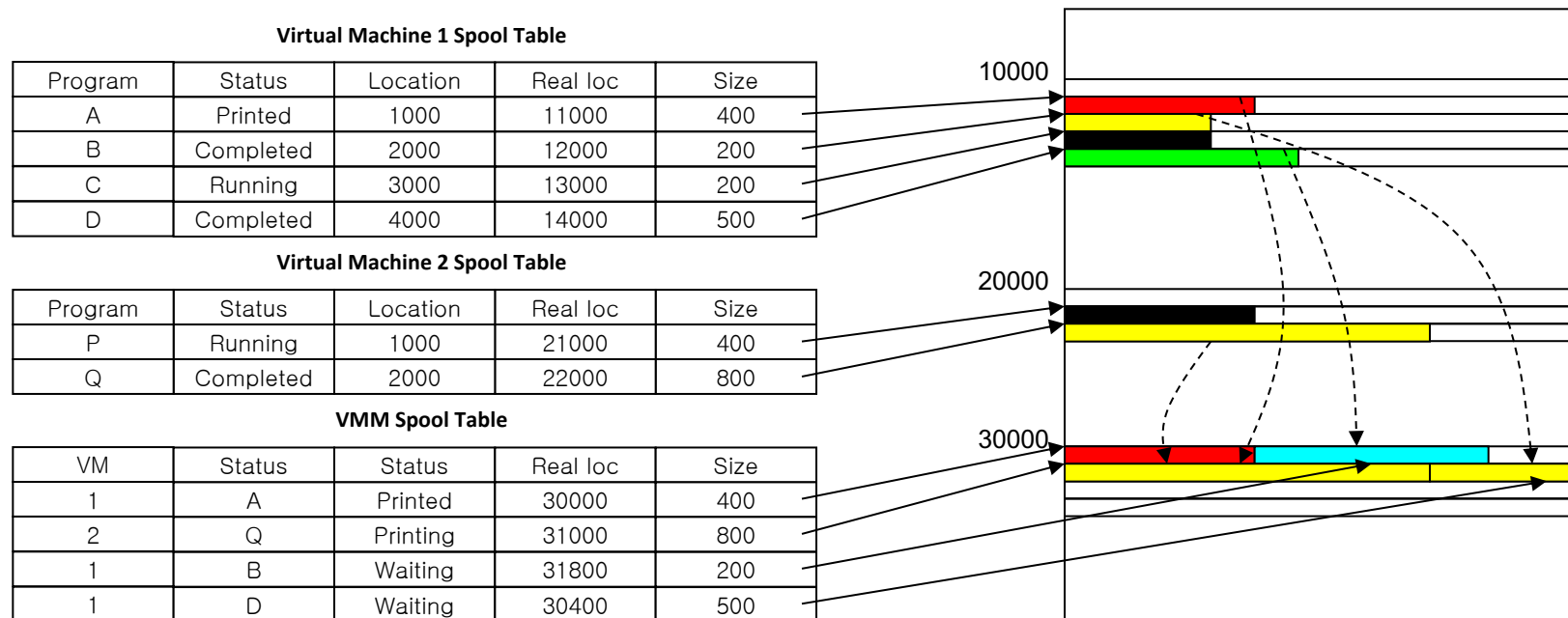
Resource Virtualization – Input/Output

Virtualizing Device

- Shared Devices
 - Some device, such as a network adapter, can be shared among a number of guest VMs at a fine time granularity
 - Each guest may have its own virtual state related to usage of the device, e.g., a virtual network address.
 - This state information is maintained by the VMM for each guest VM
- Nonexistent Physical Device
 - Virtual devices “attached” to a virtual machine for which there is no corresponding physical device
 - For example, a network adapter that is used for communicating with other virtual machines on the same platform

Resource Virtualization – Input/Output Virtualizing Device

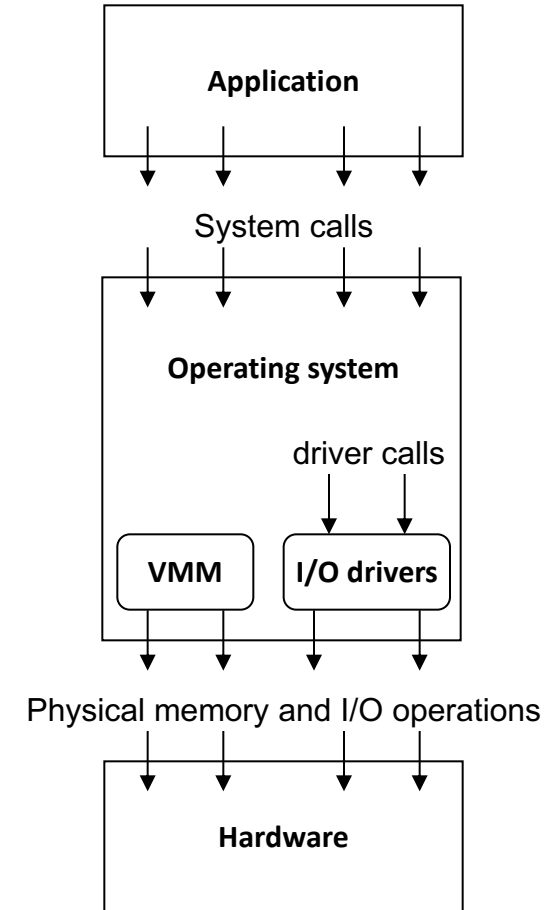
- Spooled Device
 - Virtualization of spooled device can be performed by using a two-level spool table approach



Resource Virtualization – Input/Output

Virtualizing I/O Activity

- An application program makes device-independent I/O request
- The Operating system converts the device-independent request into calls to device driver routines
- A device driver takes care of device-specific aspects of performing an I/O transaction
- The VMM can intercept a guest's I/O action and convert it from a virtual device action to a real device action at any of the three interface
 - The system call interface
 - The device driver interface
 - The operational-level interface



Resource Virtualization – Input/Output

Virtualizing I/O Activity

- Virtualizing at the I/O operation Level
 - The privileged nature of the I/O operations make them easy for the VMM to intercept because they trap in user mode
- Virtualizing at the Device Driver Level
 - If the VMM can intercept the call to the virtual device driver, it can convert the virtual device information to the corresponding physical device and redirect the call to a driver program for the physical device
 - It requires that the VMM developer have some knowledge of the guest operating system and its internal device driver interfaces
- Virtualizing at the System call Level
 - The virtualization process could be made more efficient by intercepting the initial I/O request at the OS interface, the ABI. Then the entire I/O action could be done by the VMM

Input / Output Virtualization and Hosted Virtual Machines

- An I/O request from a guest virtual machine is converted by the native-mode portion of the VMM into a user application request made to the host
- An advantage of a hosted virtual machine
 - It is not necessary to provide device drivers in the VMM
 - the actual device drivers do not have to be incorporated as part of the VMM

Input / Output Virtualization and Hosted Virtual Machines

- A component that form a dual mode hosted virtual machine system
 - VMM-n(native)
 - Intercepts traps due to privileged instructions or patched critical instructions encountered in a virtual machine
 - VMM-u(user)
 - Makes resource requests to the host OS
 - VMM-d(driver)
 - Provide a means for communication between the other two components

Memory Virtualization

Virtual Memory for System VMs

“Any **programming** problem can be solved
by **adding** a level of indirection.”

Anonymous

“Any **programming** problem can be solved
by **adding** a level of indirection.”

Anonymous

“Any **performance** problem can be solved
by **removing** a level of indirection.”

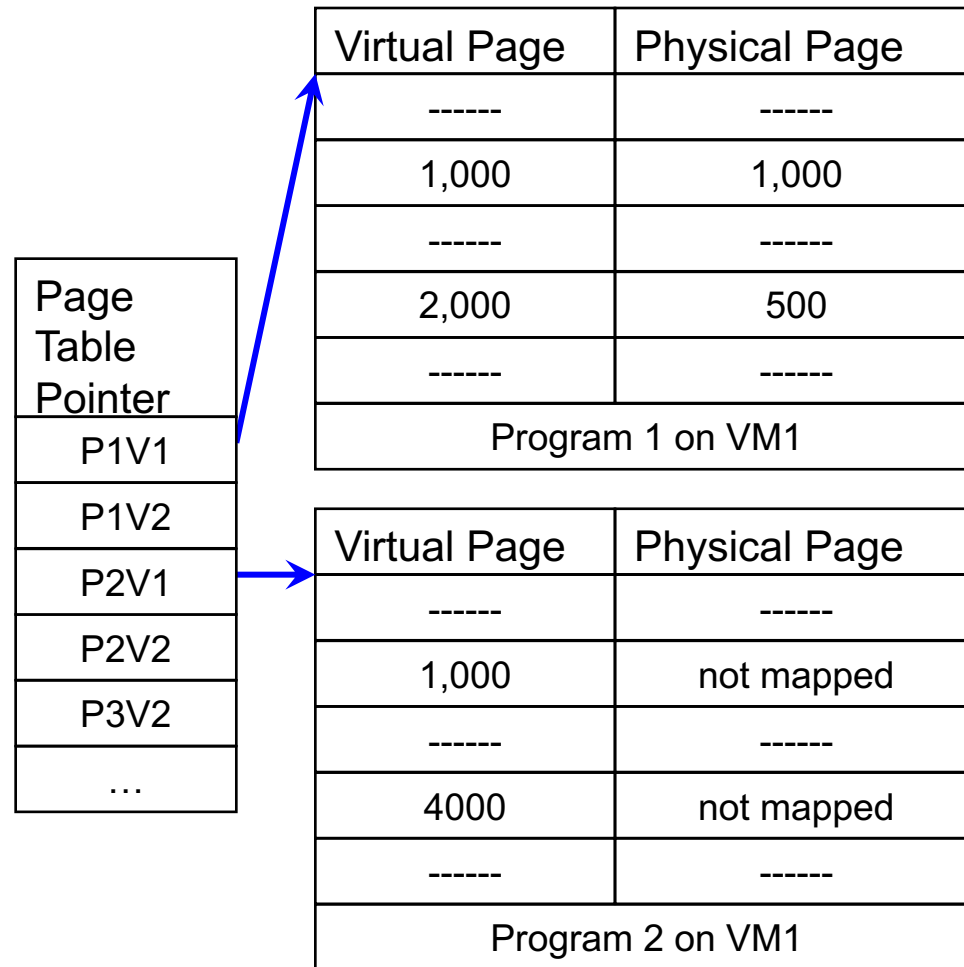
M. Haertel

Implementation

- Translation Lookaside Buffer is used to cache locations of frequently-accessed pages
 - “Architected Page Table” → Hardware maintains page table & instructions regarding it are part of the ISA. The TLB is invisible to ISA and OS.
 - “Architected TLB” → Instructions to manipulate TLB are part of the ISA. OS implements the page table (hardware is unaware) More recent.

Implementation

- Architected Page Tables
 - Guest OS manages virtual-to-real mappings for its processes.
 - VMM manages virtual-to-physical mappings in a *Shadow Page Table*, and virtualizes the page table pointer. Instructions to modify this trap & the VMM updates guest's page table pointer.



Hardware support for virtualization

Intel Virtualization Solution

Source: Y. Chung, National Tsinghua University, Taiwan

What is needed

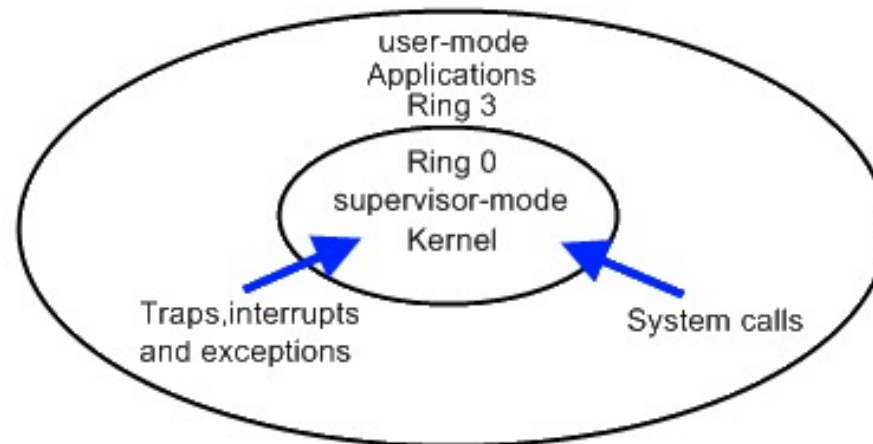
- CPU Virtualization: Intel VT-x
- Memory Virtualization: Extended Page Tables (EPT)
- IO Virtualization: Intel VT-d

Trap and Emulate Model

- We want CPU virtualization to be efficient
 - We should make guest binaries run on CPU as fast as possible.
 - Theoretically speaking, if we can run all guest binaries natively, there will NO overhead at all.
 - But we cannot let guest OS handle everything, VMM should be able to control all hardware resources.
- Solution :
 - Ring Compression
 - Shift traditional OS from kernel mode(Ring 0) to user mode(Ring 1), and run VMM in kernel mode.
 - Then VMM will be able to intercept all trapping event.

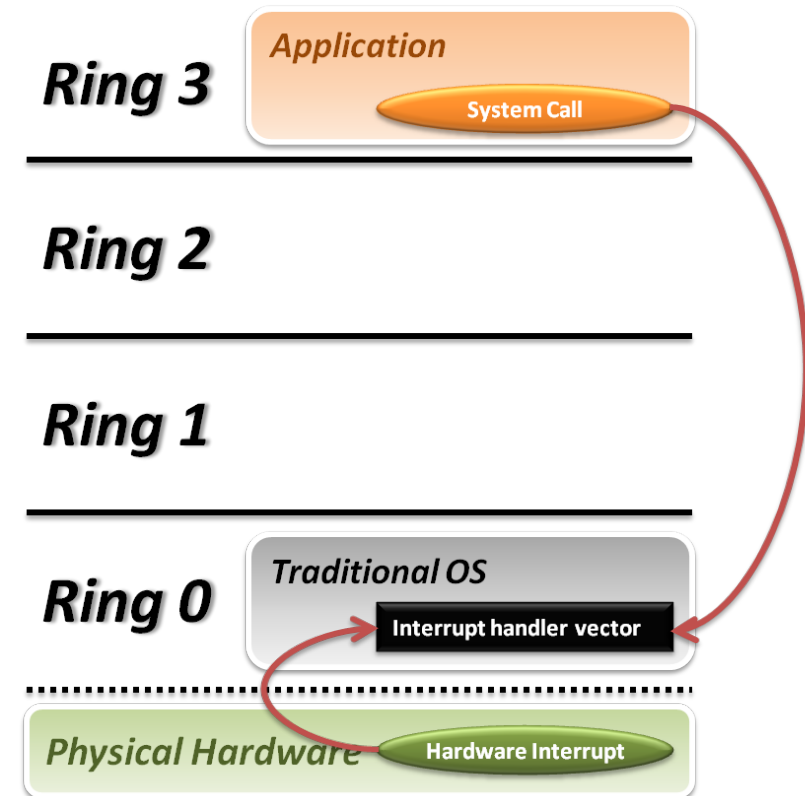
Trap and Emulate Model

- VMM virtualization paradigm (*trap and emulate*) :
 1. Let normal instructions of guest OS run directly on processor in user mode.
 2. When executing privileged instructions, hardware will make processor trap into the VMM.
 3. The VMM emulates the effect of the privileged instructions for the guest OS and return to guest.



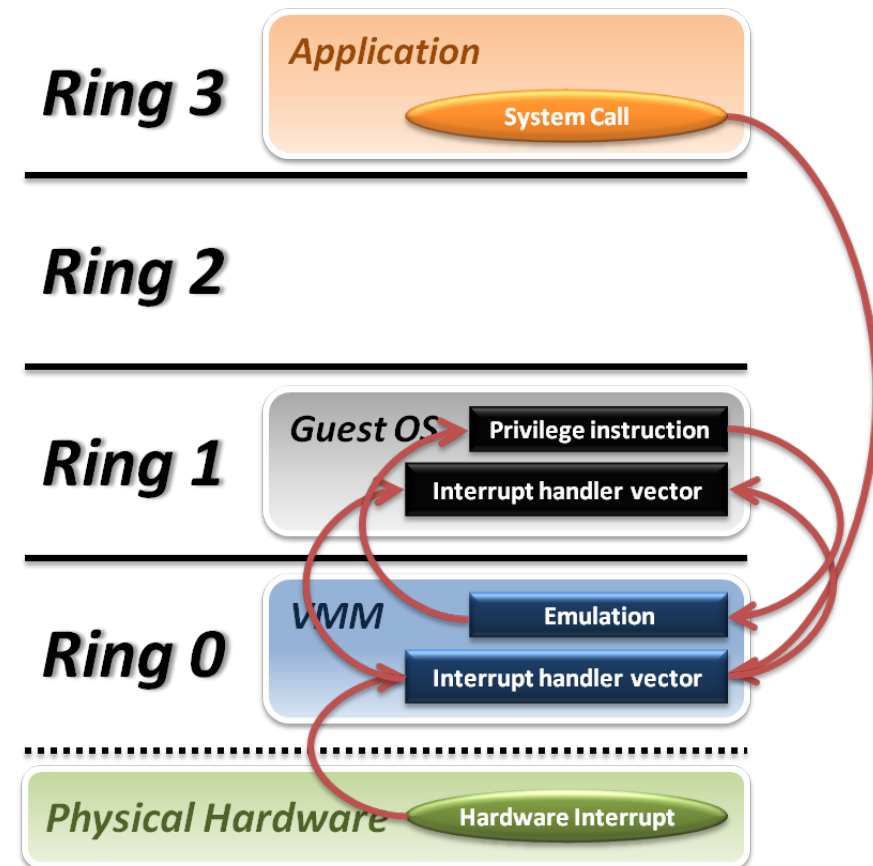
Trap and Emulate Model

- Traditional OS :
 - When application invoke a system call :
 - CPU will trap to interrupt handler vector in OS.
 - CPU will switch to kernel mode (Ring 0) and execute OS instructions.
 - When hardware event :
 - Hardware will interrupt CPU execution, and jump to interrupt handler in OS.

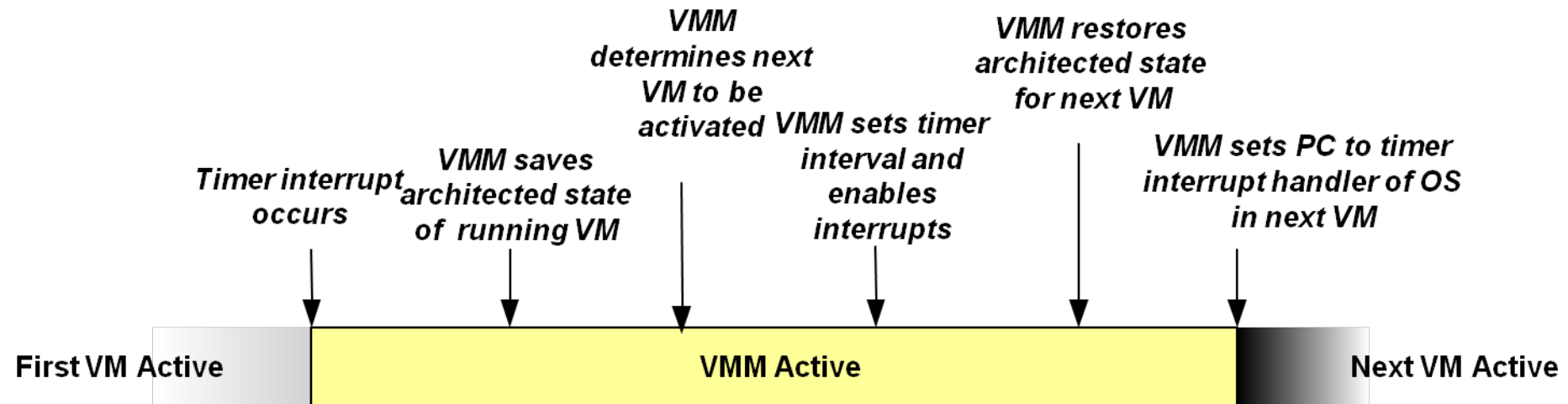


Trap and Emulate Model

- VMM and Guest OS :
 - System Call
 - CPU will trap to interrupt handler vector of VMM.
 - VMM jump back into guest OS.
 - Hardware Interrupt
 - Hardware make CPU trap to interrupt handler of VMM.
 - VMM jump to corresponding interrupt handler of guest OS.
 - Privilege Instruction
 - Running privilege instructions in guest OS will be trapped to VMM for instruction emulation.
 - After emulation, VMM jump back to guest OS.

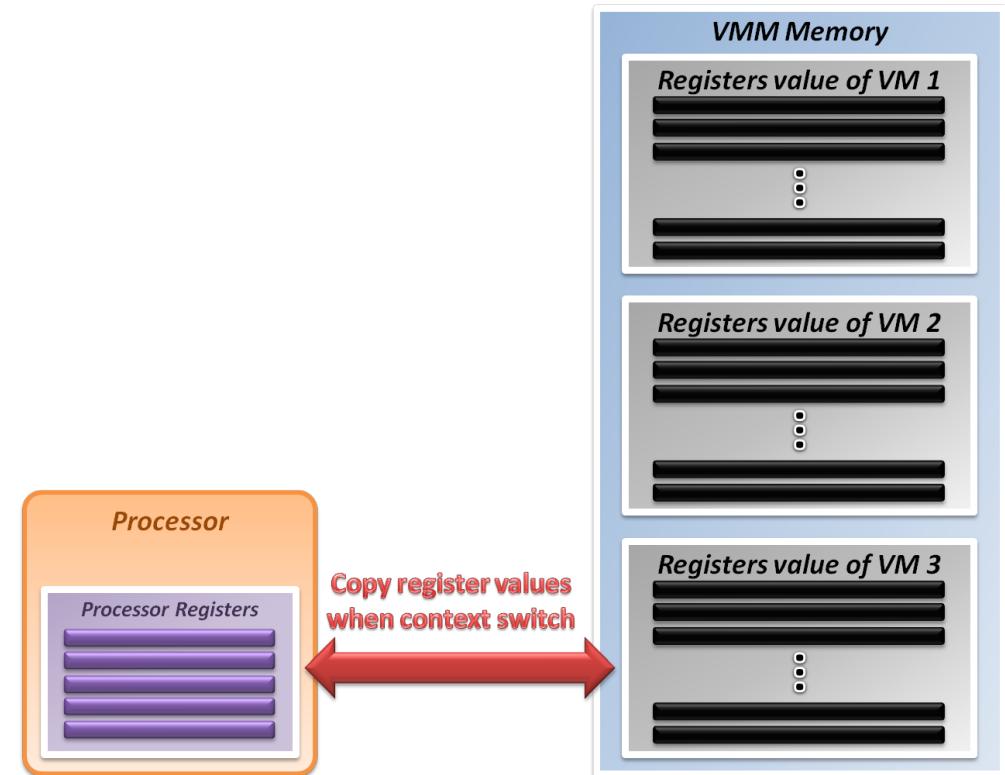


Context Switch



System State Management

- Virtualizing system state :
 - VMM will hold the system states of all virtual machines in memory.
 - When VMM context switch from one virtual machine to another
 - Write the register values back to memory
 - Copy the register values of next guest OS to CPU registers.

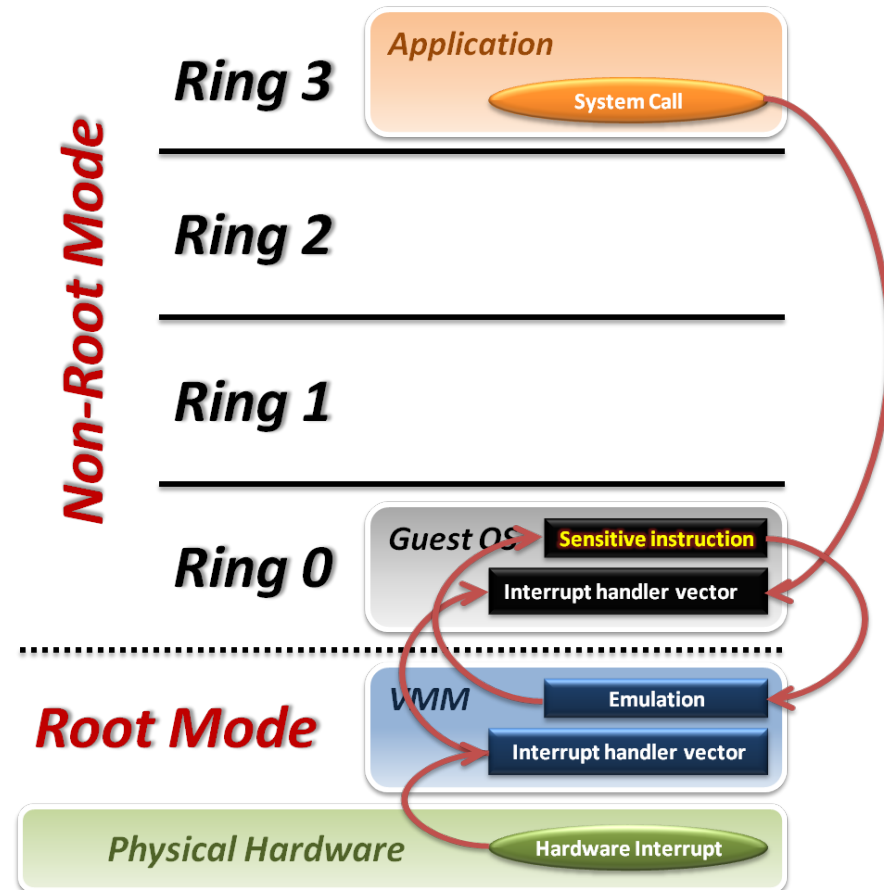


Intel VT-x

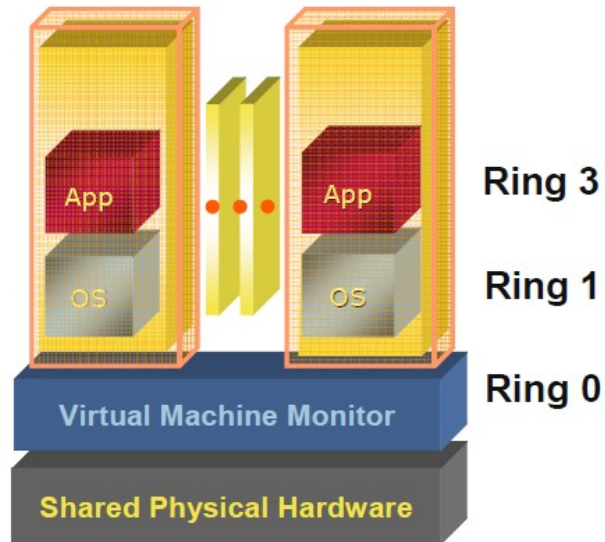
- In order to straighten those problems out, Intel introduces one more operation mode of x86 architecture.
 - **VMX Root Operation** (Root Mode)
 - All instruction behaviors in this mode are no different to traditional ones.
 - All legacy software can run in this mode correctly.
 - VMM should run in this mode and control all system resources.
 - **VMX Non-Root Operation** (Non-Root Mode)
 - All sensitive instruction behaviors in this mode are redefined.
 - The sensitive instructions will trap to Root Mode.
 - Guest OS should run in this mode and be fully virtualized through typical *“trap and emulation model”*.
 - See “Changes to Instruction Behavior in VMX Non-Root Operation” in Chapter 25 of the Intel® 64 and IA-32 Architectures Software Developer’s Manual, Volume 3C

Intel VT-x

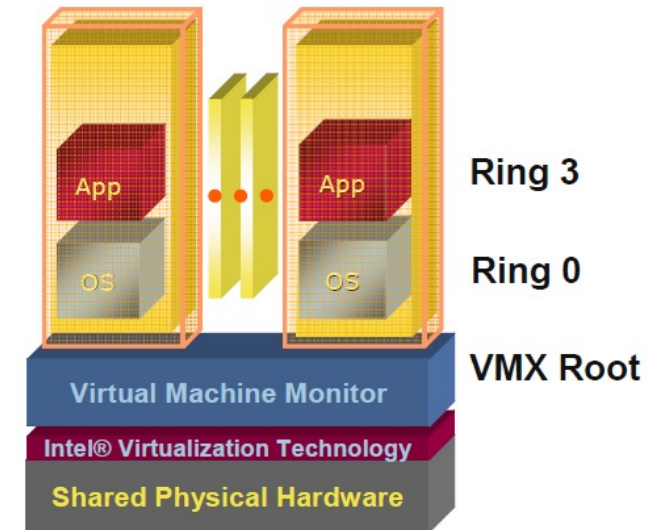
- VMM with VT-x :
 - System Call
 - CPU will directly trap to interrupt handler vector of guest OS.
 - Hardware Interrupt
 - Still, hardware events need to be handled by VMM first.
 - Sensitive Instruction
 - Instead of trap all privilege instructions, running guest OS in Non-root mode will trap sensitive instruction only.



Pre & Post Intel VT-x



- VMM de-privileges the guest OS into Ring 1, and takes up Ring 0
- OS un-aware it is not running in traditional ring 0 privilege
- Requires compute intensive SW translation to mitigate

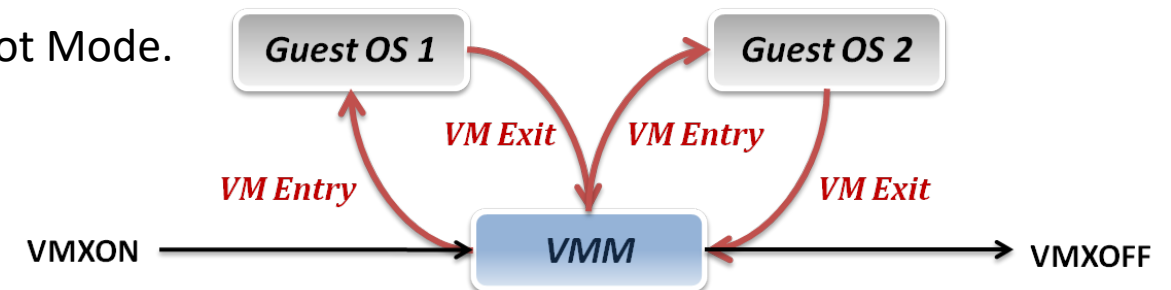


- VMM has its own privileged level where it executes
- No need to de-privilege the guest OS
- OSes run directly on the hardware

CHAPTER 30		
VMX INSTRUCTION REFERENCE		
30.1	OVERVIEW	29-1
30.2	CONVENTIONS	29-2
30.3	VMX INSTRUCTIONS	29-2
	INVEPT—Invalidate Translations Derived from EPT	29-3
	INVVPID—Invalidate Translations Based on VPID	29-6
	VMCALL—Call to VM Monitor	29-9
	VMCLEAR—Clear Virtual-Machine Control Structure	29-11
	VMFUNC—Invoke VM function	29-13
	VMLAUNCH/VMRESUME—Launch/Resume Virtual Machine	29-14
	VMPTD—Load Pointer to Virtual-Machine Control Structure	29-17
	VMPTDST—Store Pointer to Virtual-Machine Control Structure	29-19
	VMREAD—Read Field from Virtual-Machine Control Structure	29-21
	VMRESUME—Resume Virtual Machine	29-23
	VMWRITE—Write Field to Virtual-Machine Control Structure	29-24
	VMXOFF—Leave VMX Operation	29-26
	VMXON—Enter VMX Operation	29-28
30.4	VM INSTRUCTION ERROR NUMBERS	29-31

Context Switch

- VMM switch different virtual machines with Intel VT-x :
 - **VMXON/VMXOFF**
 - These two instructions are used to turn on/off CPU Root Mode.
 - VM Entry
 - This is usually caused by the execution of **VMLAUNCH/VMRESUME** instructions, which will switch CPU mode from Root Mode to Non-Root Mode.
 - VM Exit
 - This may be caused by many reasons, such as hardware interrupts or sensitive instruction executions.
 - Switch CPU mode from Non-Root Mode to Root Mode.

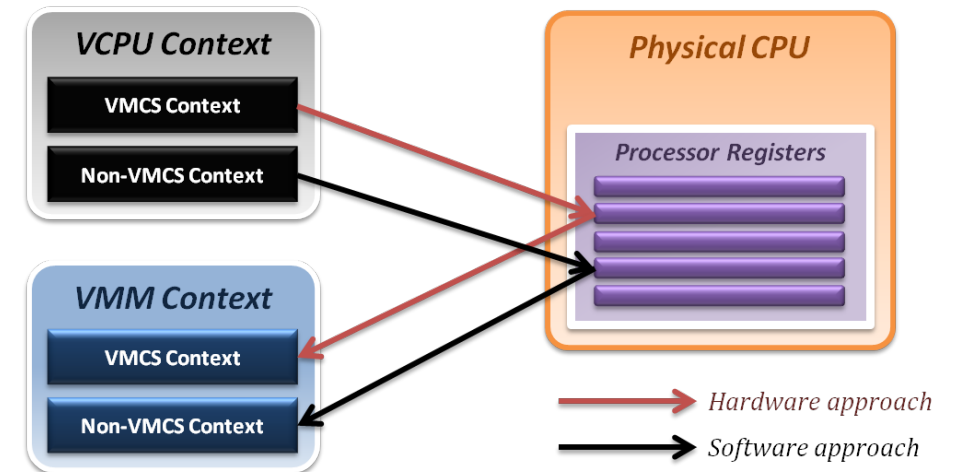


System State Management

- Intel introduces a more efficient hardware approach for register switching, **VMCS** (*Virtual Machine Control Structure*) :
 - State Area
 - Store host OS system state when VM-Entry.
 - Store guest OS system state when VM-Exit.
 - Control Area
 - Control instruction behaviors in Non-Root Mode.
 - Control VM-Entry and VM-Exit process.
 - Exit Information
 - Provide the VM-Exit reason and some hardware information.
- Whenever VM Entry or VM Exit occur, CPU will automatically read or write corresponding information into VMCS.

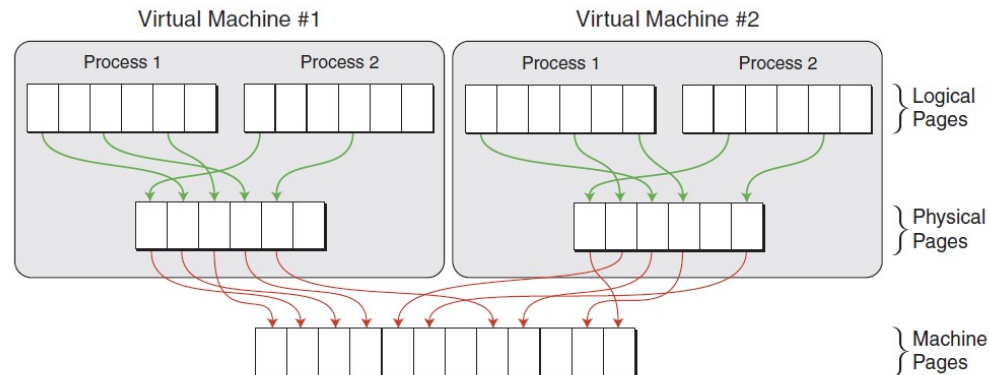
System State Management

- Binding virtual machine to virtual CPU
 - VCPU (Virtual CPU) contains two parts
 - VMCS maintains virtual system states, which is approached by hardware.
 - Non-VMCS maintains other non-essential system information, which is approach by software.
 - VMM needs to handle Non-VMCS part.

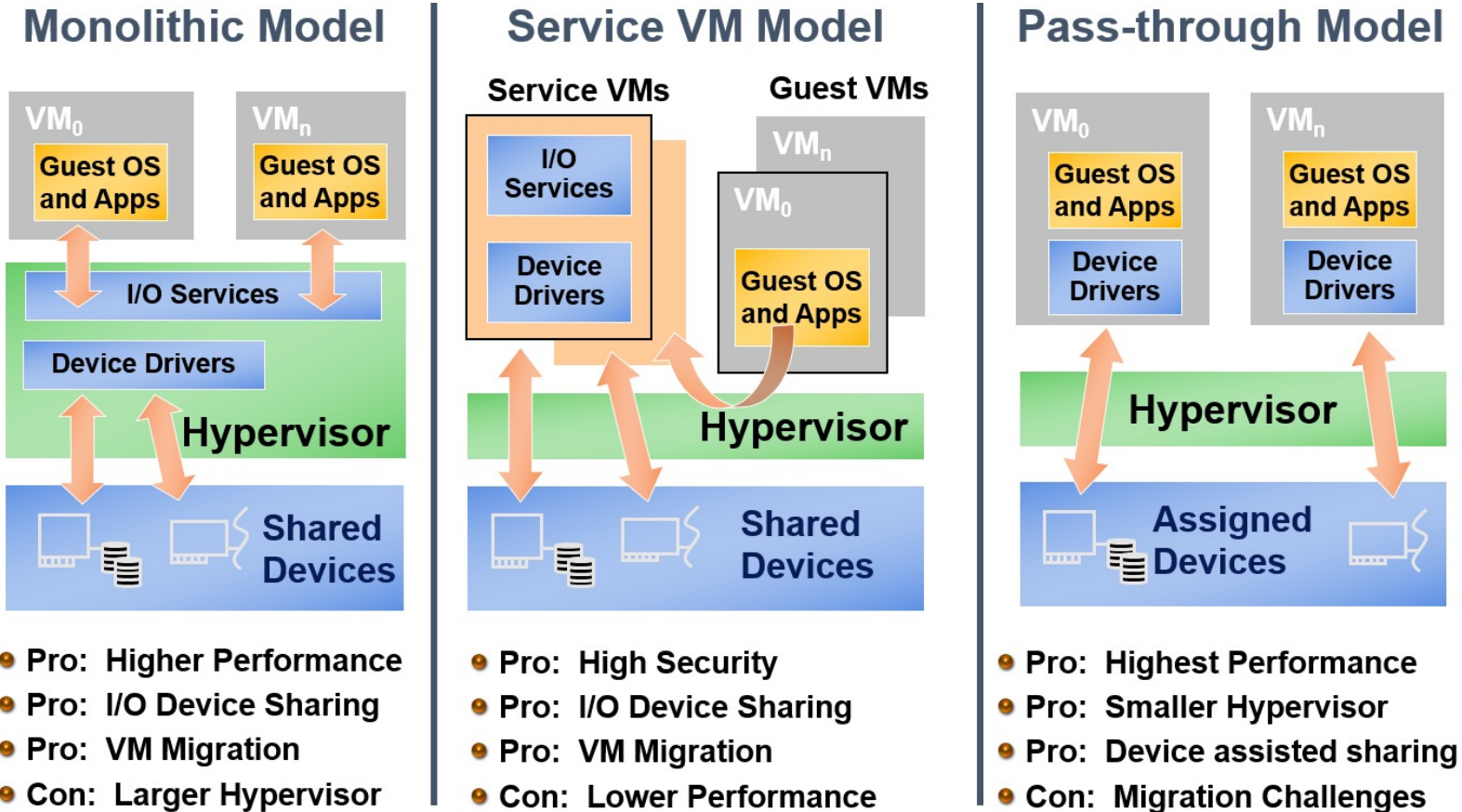


Memory Virtualization: Extended Page Table

- Concept of Extended Page Table (EPT) :
 - Instead of walking along with only one page table hierarchy, EPT technique implement one more page table hierarchy.
 - One page table is maintained by guest OS, which is used to generate guest physical address.
 - The other page table is maintained by VMM, which is used to map guest physical address to host physical address.
 - For each memory access operation, EPT MMU will directly get guest physical address from guest page table, and then get host physical address by the VMM mapping table automatically.



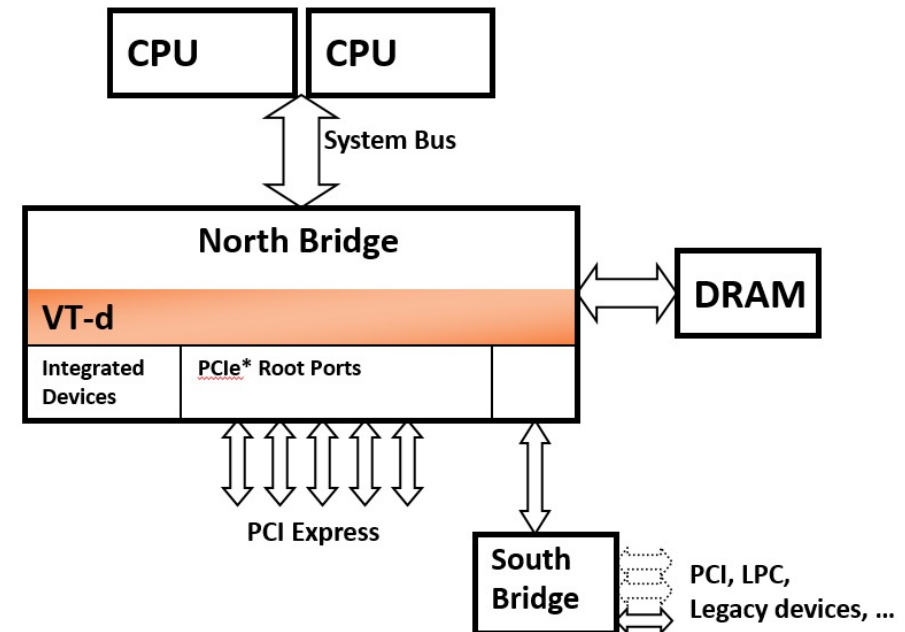
Options For I/O Virtualization



VT-d Goal: Support all Models

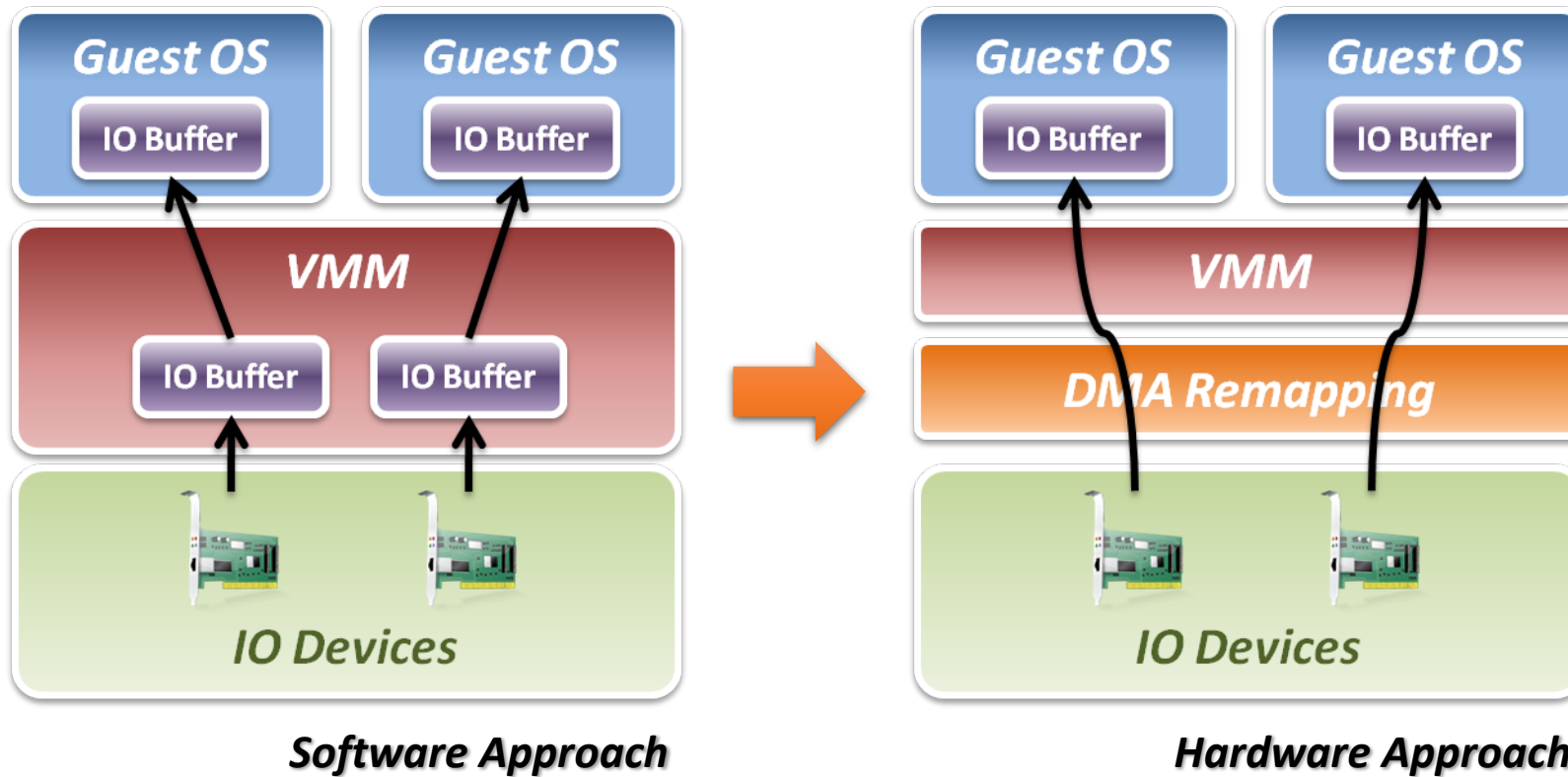
VT-d Overview

- VT-d is platform infrastructure for I/O virtualization
 - Defines architecture for DMA remapping
 - Implemented as part of platform core logic
 - Will be supported broadly in Intel server and client chipsets

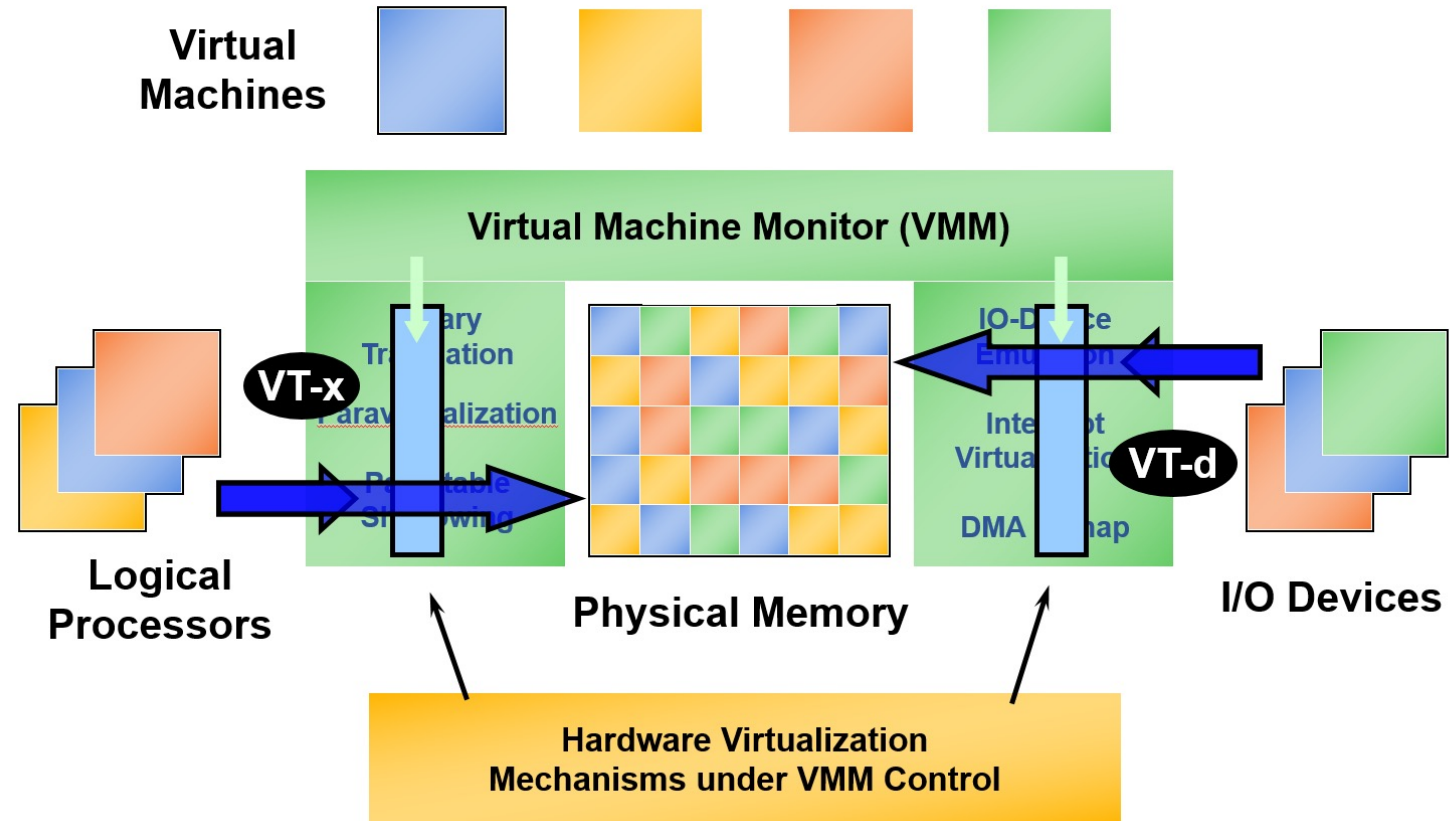


Intel VT-d

- Add DMA remapping hardware component.



VT-x & VT-d Working Together



Summary

- CPU Virtualization
 - Trap and Emulate Model
 - Virtualization technique, VMX Root/Non-Root Operation, VMM and Guest OS, VMCS ... etc.
- Memory Virtualization: Extended Page Tables (EPT)
 - EPT implement one more page table hierarchy
 - MMU virtualize, EPT translation, Memory Operation,... etc.
- IO Virtualization: Intel VT-d
 - Implement DMA remapping in hardware
 - Hardware Page Walk, Translation Caching

Our journey in CS5250

What we have gone through

- We (re)visited the traditional OS concepts such as process management, virtual memory etc.
 - Linked them with actual ISA support
 - Deep dive into an industrial strength OS
 - Studied some of the key data structures and algorithms used

What you have achieved

- The kernel is no longer some mysterious blackbox
 - You actually work some kernel code!
 - Understand the API a bit better
- A better appreciation of how it works and the tradeoffs involved
 - Some issues translate to wider software engineering practices

Looking forward

- No doubt the kernel will change
 - Even as we speak
- But not the underlying principles
 - You will be able to appreciate future changes even better
- Hopefully, the experience will make you a better systems person

END