



BITS Pilani

Cloud Computing

CC ZG527

Shwetha Vittal
shwetha.vittal@pilani.bits-pilani.ac.in

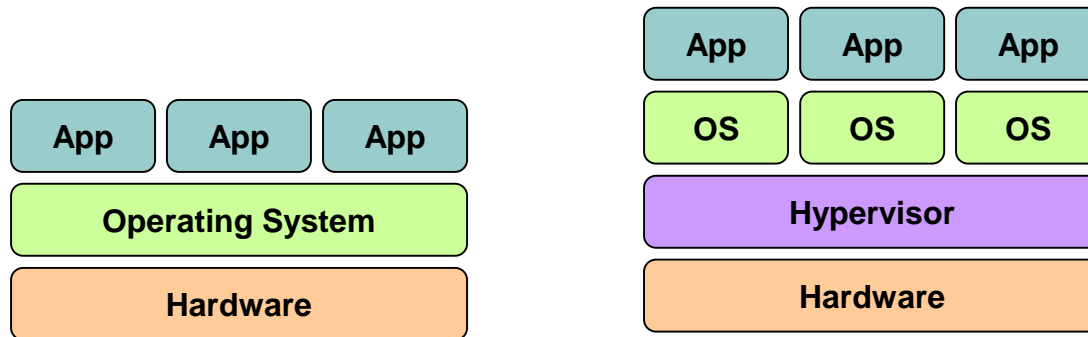
Agenda

Virtualization Techniques and Types

- ❑ Introduction to Virtualization
- ❑ Emergence of Virtualization
- ❑ Types of Virtualization
- ❑ Types of Hypervisors
- ❑ Use & demerits of Virtualization

Technology That Made Cloud Possible

Key Technology is Virtualization



Virtualization gives:

- An enabling technology for datacentre implementation
- An abstract compute, network, and storage service platforms from the underlying physical hardware

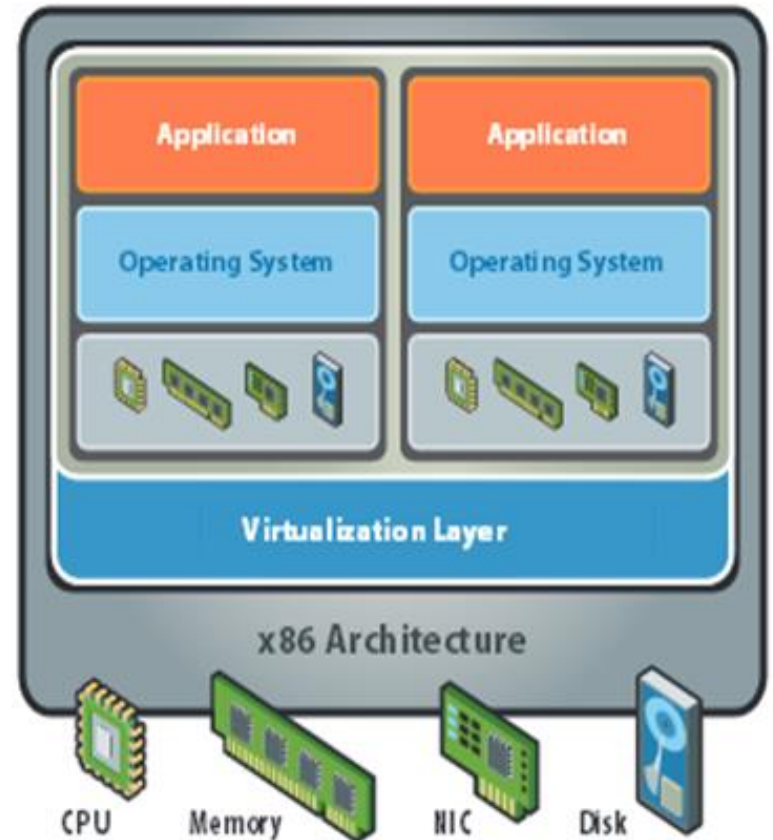
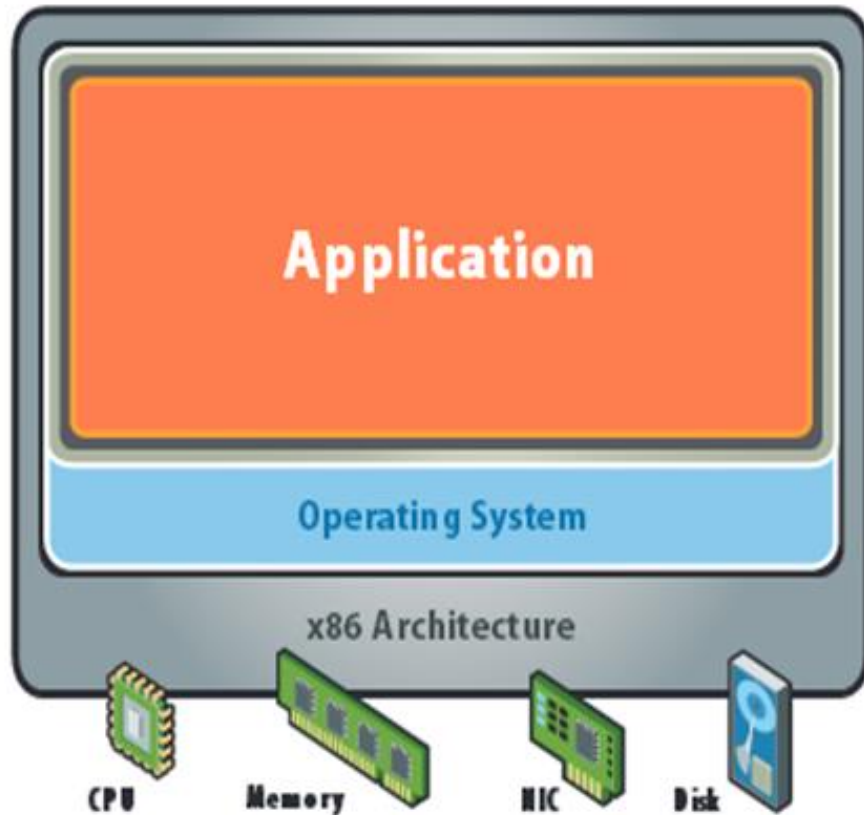
Importance of Virtualization in Cloud Computing

- Cloud can exist without Virtualization, although it will be difficult and inefficient.
- Cloud makes notion of “Pay for what you use”, “infinite availability- use as much you want”.
- These notions are practical only if we have
 - lot of flexibility
 - efficiency in the back-end
- This efficiency is readily available in Virtualized Environments and Machines

Virtualization Re-emergence

- Server Sprawl: Multiple redundant servers present with low utilization
- Goal was to improve server utilization with application isolation
- Server Consolidation
- Platform independence by Virtual Machine Encapsulation
- Improved Scalability and availability with per application isolation
- Faster provisioning for newer applications

What is Virtualization?



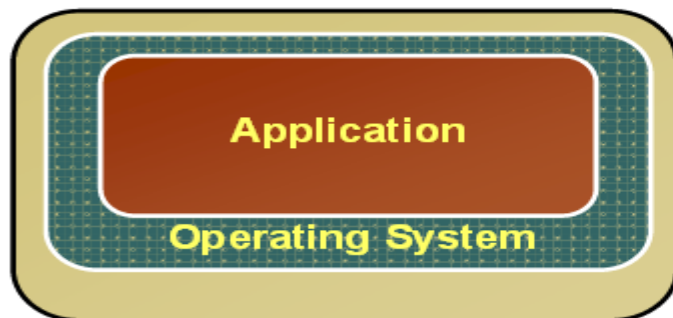
What does Virtualization do?

- Virtualization allows multiple operating system instances to run concurrently on a single computer
- Each “guest” OS is managed by a Virtual Machine Monitor /Manager (VMM), also known as a hypervisor.
- Because the virtualization layer sits between the guest and the hardware,
 - It can control the guests’ use of CPU, memory, and storage
 - Allows a guest OS to migrate from one machine to another

Changes after Virtualization

Before Virtualization

- Single OS image per machine
- Software and hardware tightly coupled
- Running multiple applications on same machine often creates conflict
- Underutilized resources
- Inflexible and costly infrastructure



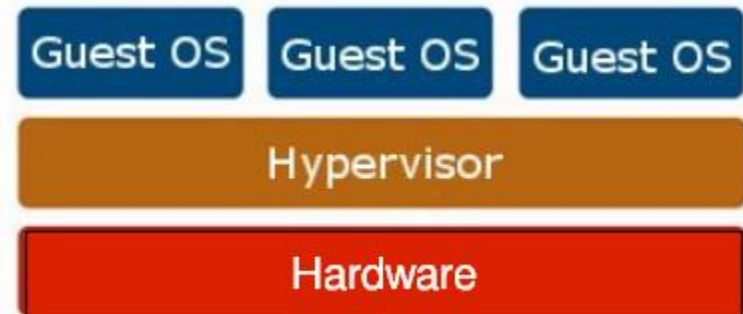
After Virtualization

- Hardware-independence of operating system and applications
- Virtual machines can be provisioned to any system
- Can manage OS and application as a single unit by encapsulating them into virtual machines



Virtualization Architecture

- OS assumes complete control of the underlying hardware.
- Virtualization architecture provides this illusion through a **Hypervisor/VMM**.
- Hypervisor/VMM is a software layer which:
 - Allows multiple Guest OS (Virtual Machines) to run simultaneously on a single physical host
 - Provides hardware abstraction
 - Multiplexes underlying hardware resources



Hypervisor

A layer of software that generally provides virtual partitioning capabilities which runs directly on hardware.

Sometimes referred to as a “bare metal” approach.



Principles of Virtualization

- **Equivalence**
 - Guest OS should run (close to) unmodified
- **Safety**
 - Should (absolutely) not be able to escape its isolated environment
- **Performance**
 - With (close to) native performance

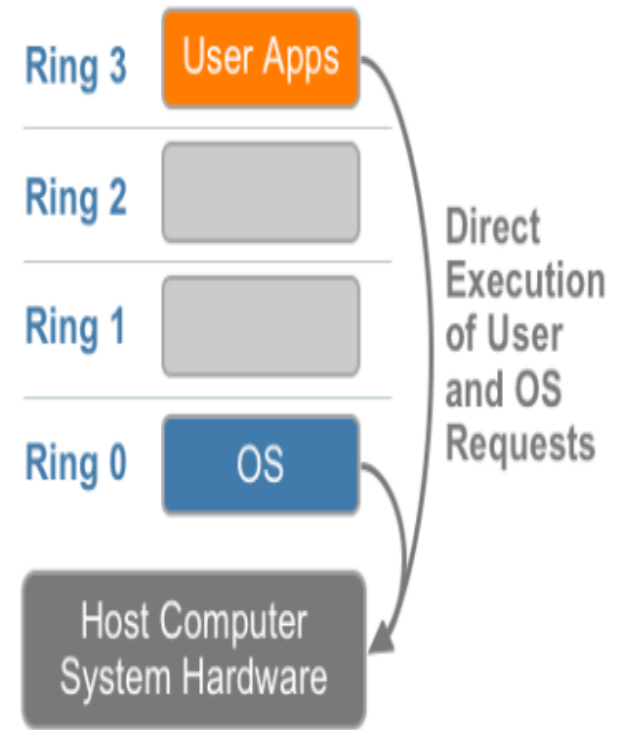
Hypervisor Design Goals

- **Isolation**
 - Security isolation
 - Fault isolation
 - Resource isolation
- **Reliability**
 - Minimal code base
 - Strictly layered design
- **Scalability**
 - Scale to large number of cores
 - Large memory systems

CPU Virtualization - 1

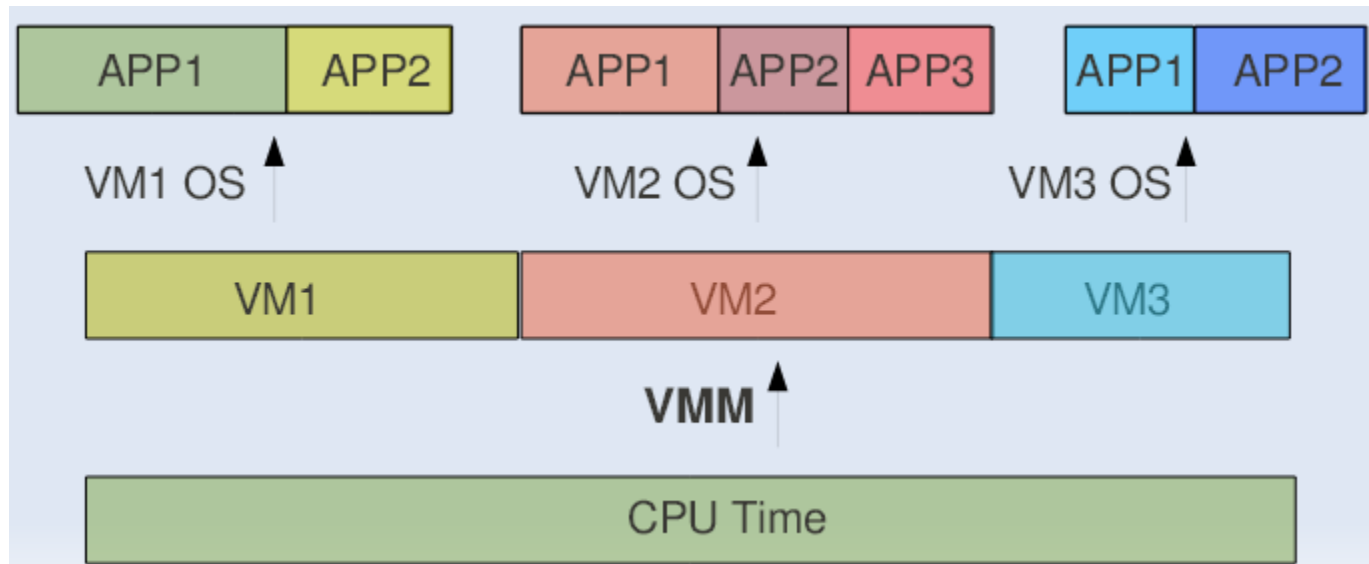
CPUs have multiple privilege levels

- Ring 0,1,2,3 in x86 CPUs
- Virtualizing the x86 architecture requires placing a virtualization layer under the operating system (which expects to be in the most privileged Ring 0)
- Normally,
 - User process in ring 3,
 - Host OS in ring 0
 - Privileged instructions only run in ring 0
 - VMM in ring 0
- Guest OS must be protected from guest apps
 - But not fully privileged like host OS/VMM
 - Run it in ring 1?



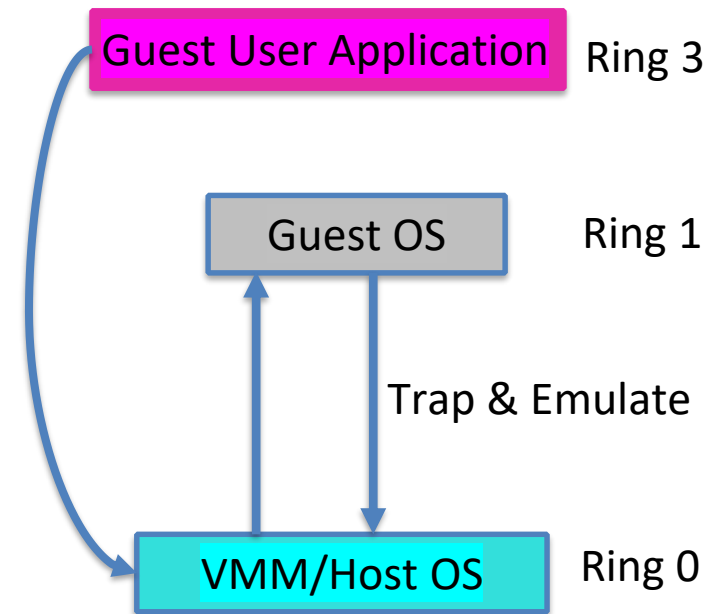
CPU Virtualization - 2

- VMM or Hypervisor provides a virtual view of CPU to VMs.
- In multi processing, CPU is allotted to the different processes in form of time slices by the OS.
- Similarly VMM or Hypervisor allots CPU to different VMs.



Trap & Emulate based Virtualization

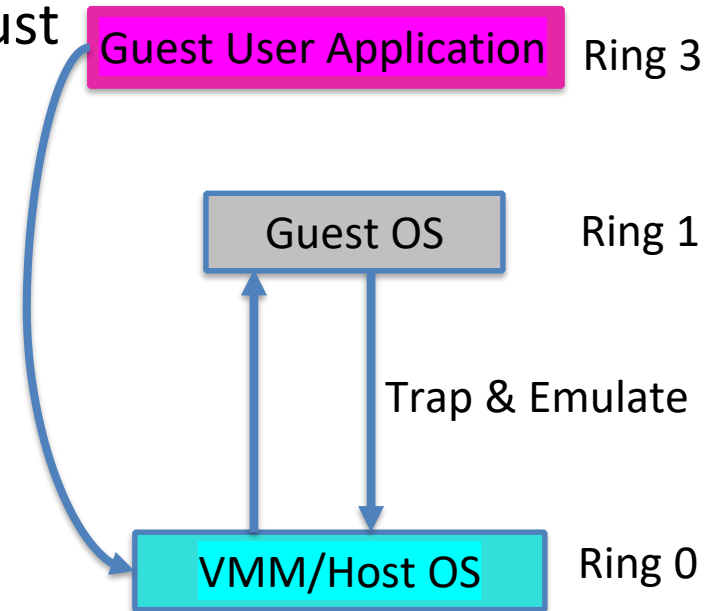
- **Scope of Guest OS**
 - Runs at lower privilege level than VMM
 - Traps to VMM for privileged operation
- **Guest app handling syscall/interrupt**
 - Special trap instr (int n), traps to VMM
 - VMM doesn't know how to handle trap
 - VMM jumps to guest OS trap handler
 - Trap handled by guest OS normally
- **Returning from Trap**
 - Guest OS performs return from trap
 - VMM jumps to corresponding user process



Trap & Emulate based Virtualization

- **Handling privileged instruction**

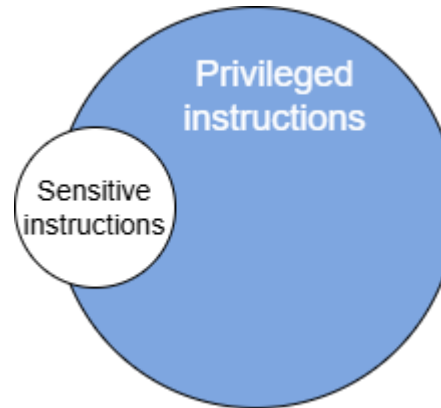
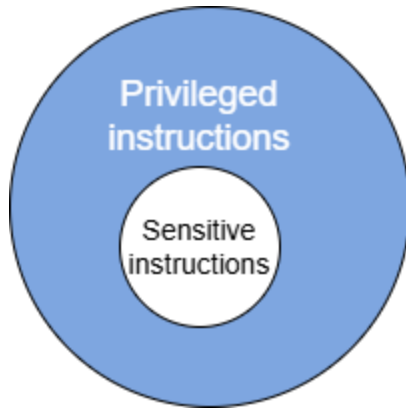
- Sensitive data structures like IDT must be managed by VMM, not guest OS
- Any privileged action by guest OS traps to VMM, emulated by VMM
- Guest OS traps to VMM
- VMM jumps to corresponding user process
 - E.g: Set IDT, set CR3, access hardware



Problem with Trap & Emulate

- OSs not ready /unaware of Virtualization
 - To run at a lower privilege level
 - Machines not designed for virtualization
- Guest OS may realize it is running at lower privilege level
- Some registers in x86 reflect higher privilege level
- **Sensitive instructions:**
 - Can change hardware state, while running in both privileged and unprivileged modes
 - Will behave differently when guest OS is in ring 0 vs in less privileged ring 1
 - OS behaves incorrectly in ring1, will not trap to VMM

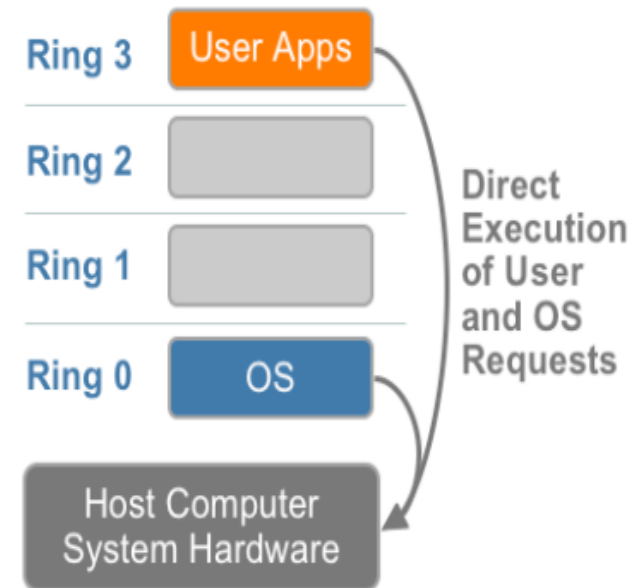
Popek Goldberg theorem



- **Sensitive instruction** may change hardware state
- **Privileged instruction** runs only in privileged mode
 - Traps to ring 0 if executed from unprivileged rings
- **In order to build a VMM efficiently via trap-and-emulate method, sensitive instructions should be a subset of privileged instructions**
- **x86 does not satisfy this criteria, so trap and emulate VMM is not possible**

CPU Virtualization

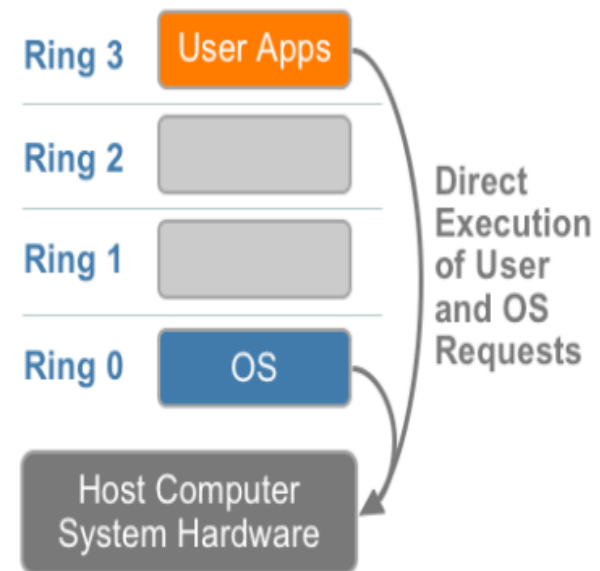
- x86 architecture offers four levels of privilege known as Ring 0, 1, 2 and 3 to operating systems and applications to manage access to the computer hardware
- Virtualizing the x86 architecture requires placing a virtualization layer under the operating system (which expects to be in the most privileged Ring 0)
- Some sensitive instructions can't effectively be virtualized as they have different semantics when they are not executed in Ring 0



Approaches to CPU Virtualization

Three techniques now exist for handling sensitive and privileged instructions to virtualize the CPU on the x86 architecture

1. **Full virtualization** using binary translation
2. **Para-virtualization** or OS assisted virtualization
3. **Hardware assisted virtualization**



Evolution of CPU Virtualization

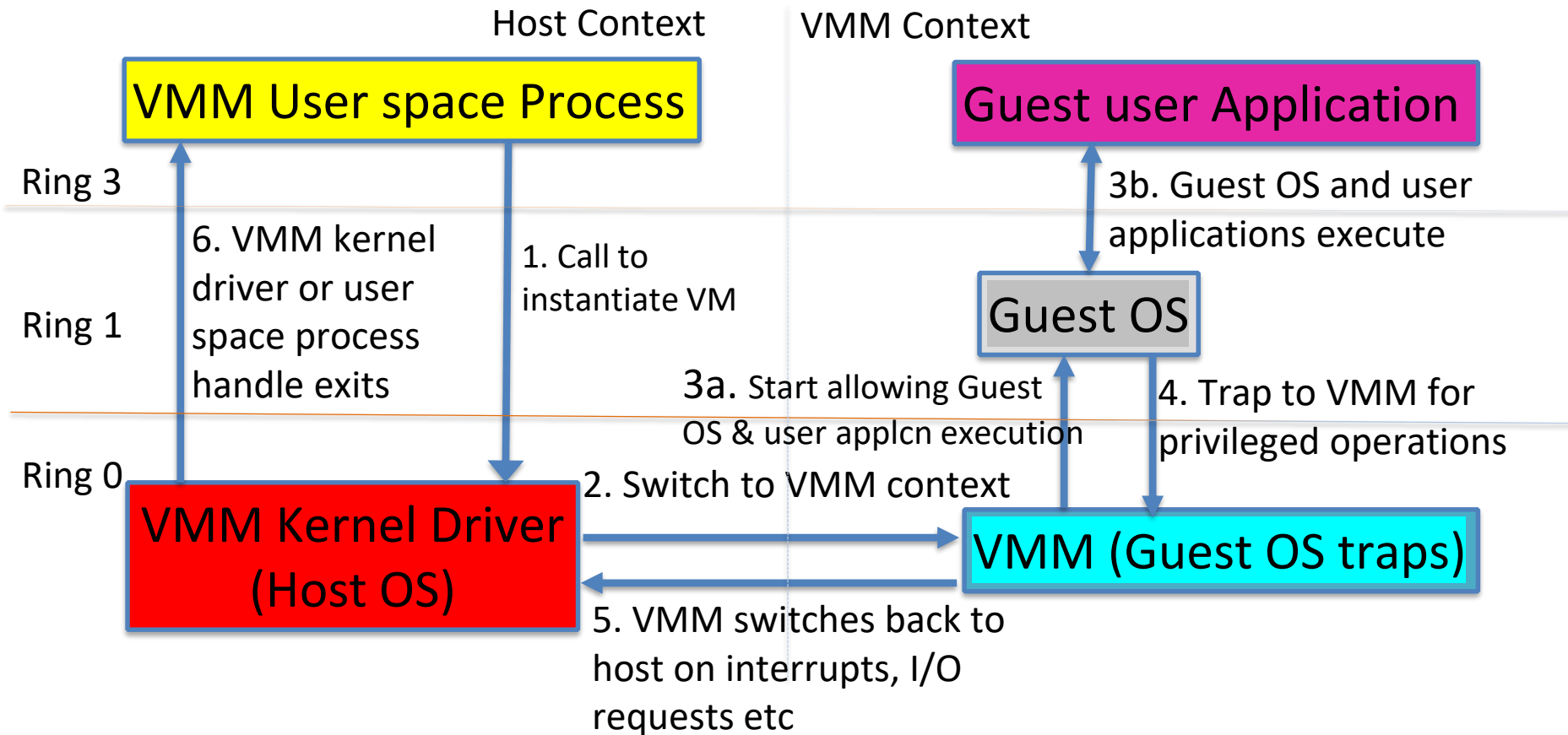
Evolution of Software Solutions

- 1st Generation: Full virtualization (Binary rewriting)
 - Software Based
 - VMware and Microsoft
- 2nd Generation: Para-virtualization
 - Cooperative virtualization
 - Modified guest
 - VMware, Xen
- 3rd Generation: Silicon-based (Hardware-assisted) virtualization
 - Unmodified guest
 - VMware and Xen on virtualization-aware hardware platforms



Virtualization Logic

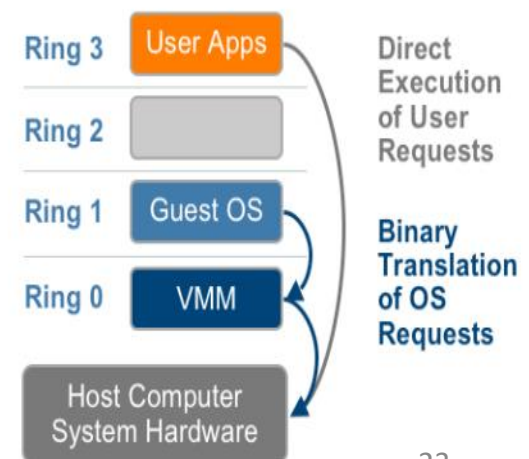
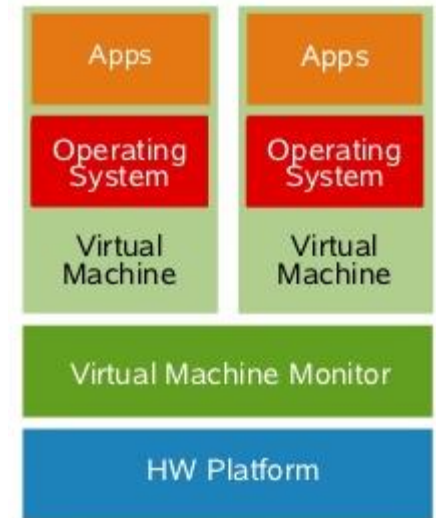
Full Virtualization



“Bringing Virtualization to the x86 Architecture with the Original VMware Workstation”, Edouard Bugnion, Scott Devine, Mendel Rosenblum, Jeremy Sugerman, Edward Y. Wang.

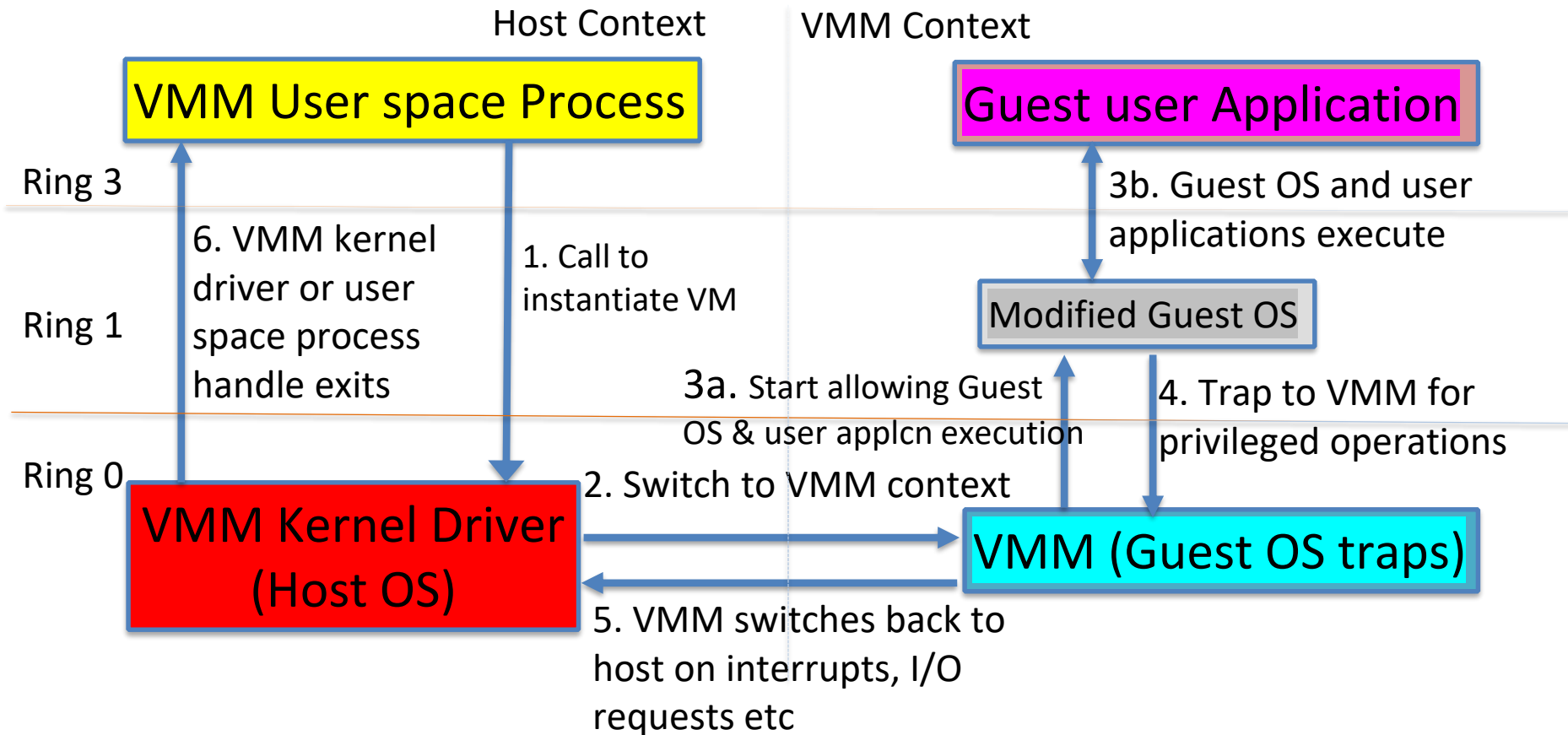
Full Virtualization using Binary Translation

- VMware(Full Virtualization) can virtualize any x86 operating system using a combination of
 - binary translation
 - direct execution techniques
- User level code is directly executed on the processor for high performance virtualization
- Kernel code is translated to replace non-virtualizable instructions with new sequences of instructions
- This combination of binary translation and direct execution provides Full Virtualization
- Guest OS is fully abstracted (completely decoupled) from the underlying hardware by the virtualization layer
- VMWare Workstation^[1] is an example for full virtualization



“Bringing Virtualization to the x86 Architecture with the Original VMware Workstation”, Edouard Bugnion, Scott Devine, Mendel Rosenblum, Jeremy Sugerman, Edward Y. Wang – 2012 ACM Transactions.

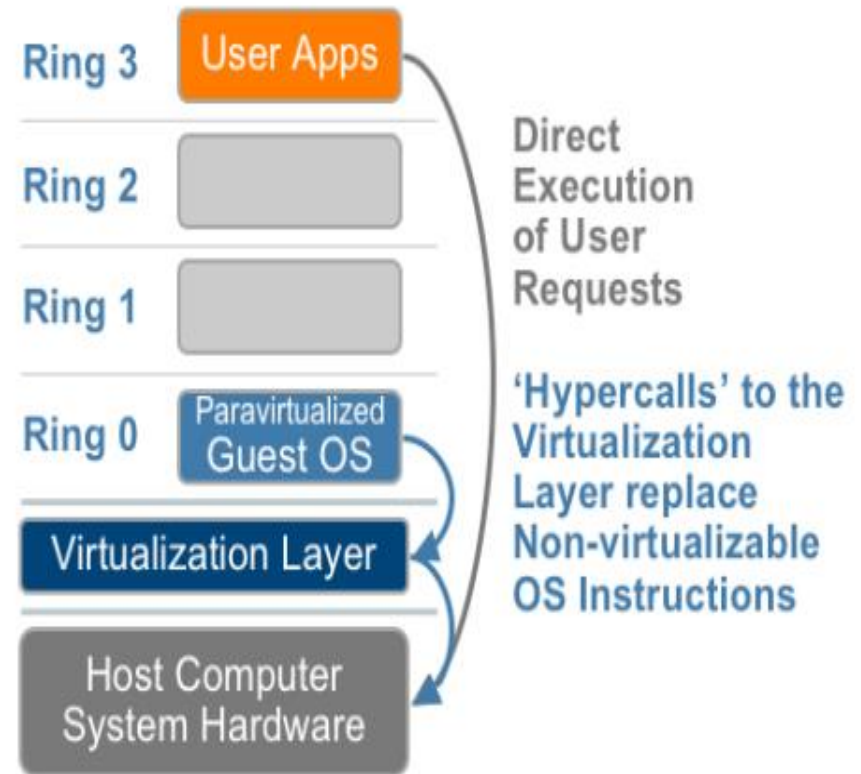
Para Virtualization



“Xen and the Art of Virtualization”, Paul Barham, Boris Dragovic, Keir Fraser, Steven Hand, Tim Harris, Alex Ho, Rolf Neugebauer, Ian Pratt, Andrew Warfield

Para Virtualization

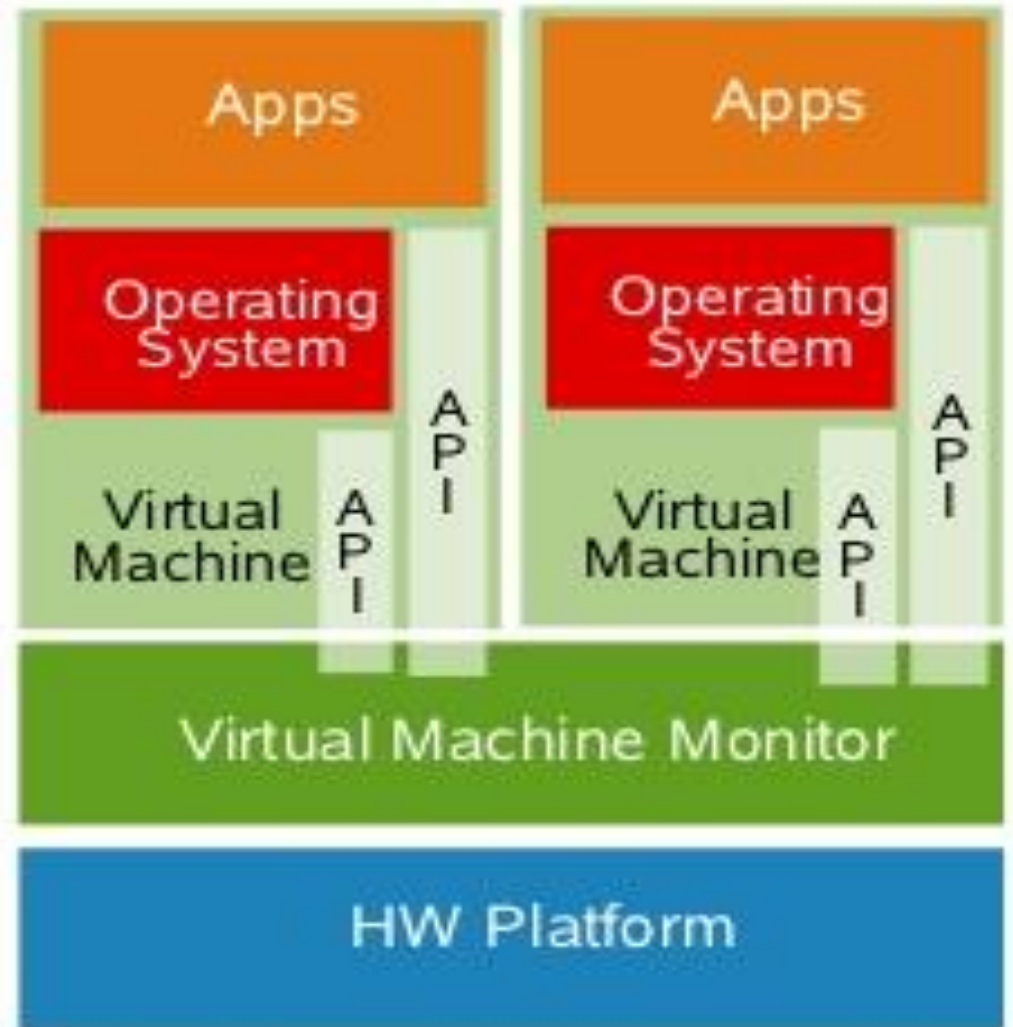
- Para virtualization involves modifying the OS kernel to replace non-virtualizable instructions with **hypercalls**.
- **Hypercalls** communicate directly with the virtualization layer
- The hypervisor also provides hypercall interfaces for other critical kernel operations such as:
 - Memory management
 - Interrupt handling
- Ex: open-source Xen – Type1 Hypervisor



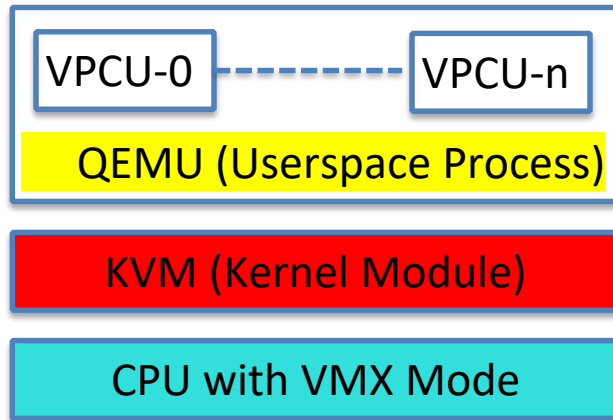
"Xen and the Art of Virtualization", Paul Barham, Boris Dragovic, Keir Fraser, Steven Hand, Tim Harris, Alex Ho, Rolf Neugebauer, Ian Pratt, Andrew Warfield

Para Virtualization

- Para virtualization cannot support unmodified operating systems
- Compatibility and portability is poor
- Ex: open-source Xen



Hardware-assisted Virtualization

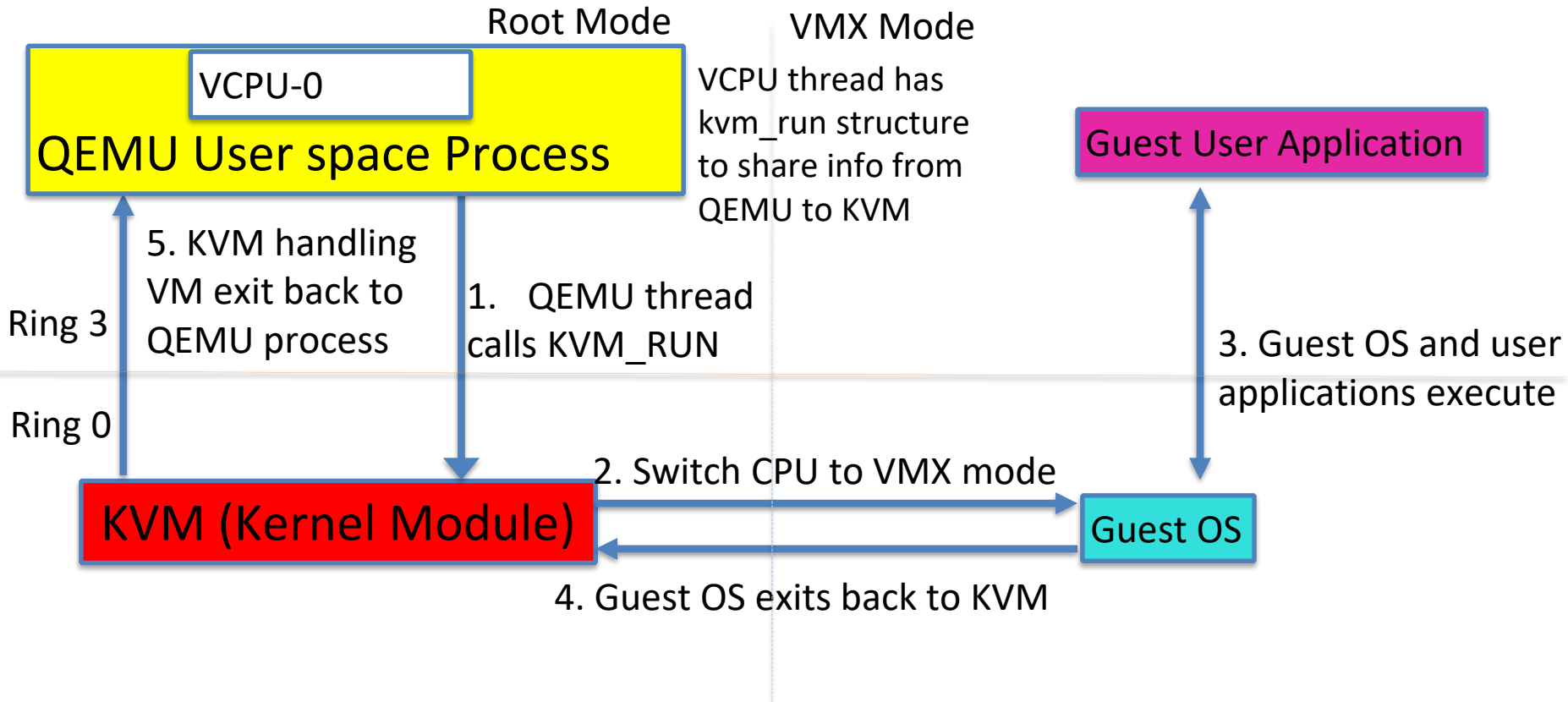


- Intel Virtualization Technology (VT-x) and AMD's AMD-V
- Target privileged instructions with a new CPU execution mode feature that allows the VMM to run in a new root mode below ring 0

Example: QEMU/KVM in Linux

- QEMU is user space process
- QEMU talks to KVM via open/ioctl syscalls
- Host OS sees QEMU as a regular multi-threaded process
- Creates one thread for each virtual CPU (VCPU) in guest
- Multiple file descriptors to /dev/kvm(one for QEMU, one for VM, one for VCPU and so on)
- ioctl on fds to talk to KVM

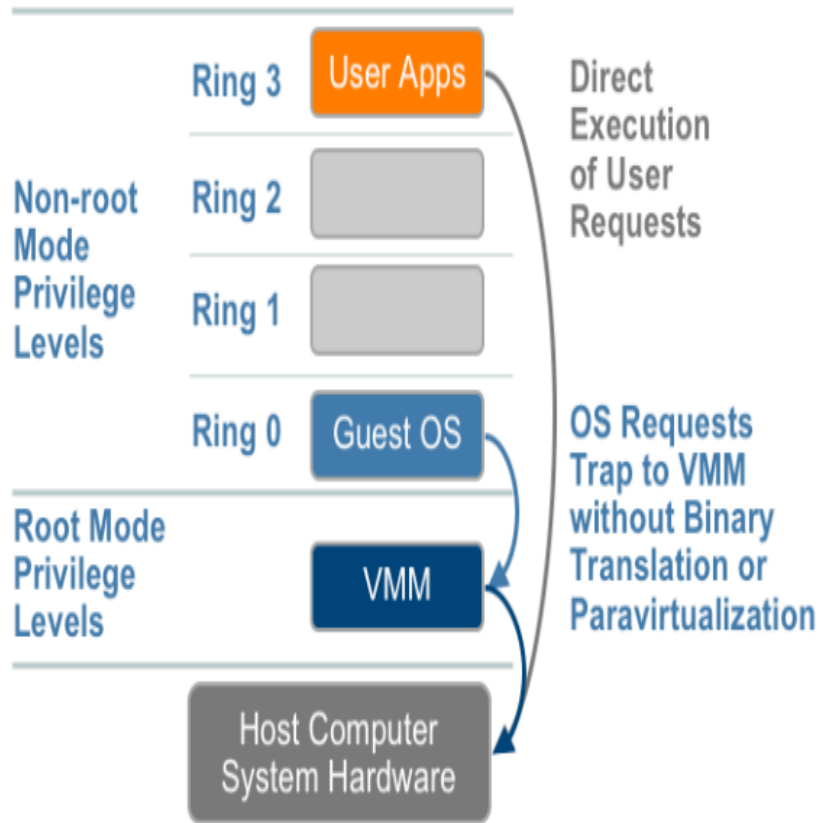
Hardware-assisted Virtualization



VMX Mode Execution

- Special CPU instructions
 - VMLAUNCH, VMRESUME invoked by KVM to enter VMX mode
 - VMEXIT invoked by guest OS to exit VMX mode
- CPU switches context between host OS to guest OS accordingly.
- Hardware manages the switches between these two modes
- VMCS (VM control structure) is for storing CPU context during the switching.
- Page tables (address space), CPU register values etc are switched

Hardware Assisted Virtualization



- Intel Virtualization Technology (VT-x) and AMD's AMD-V which both target privileged instructions with a new CPU execution mode feature that allows the VMM to run in a new root mode below ring 0

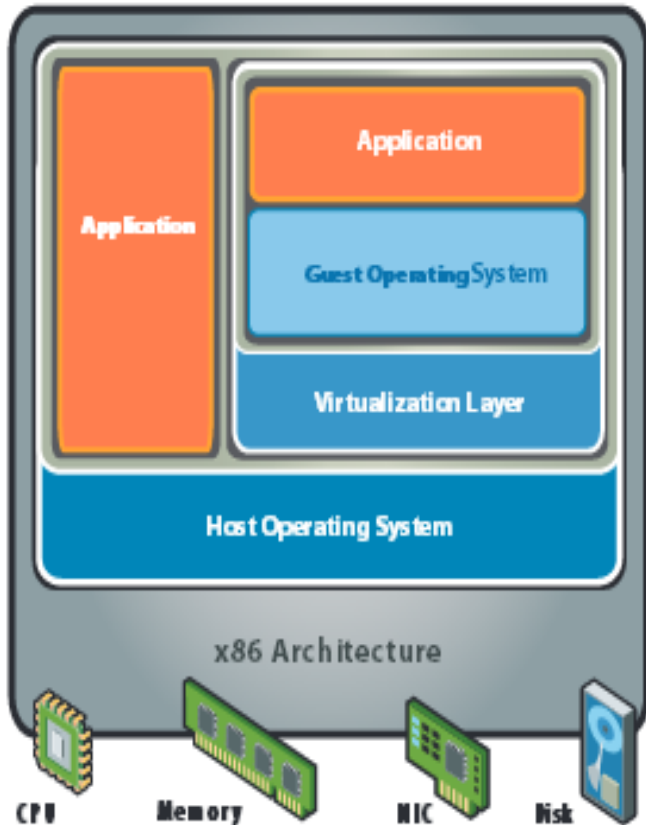
Example: QEMU/KVM in Linux

- Privileged and sensitive calls are set to automatically trap to the hypervisor, removing the need for either binary translation or paravirtualization.
- Underlying hardware provides special CPU instructions to aid virtualization

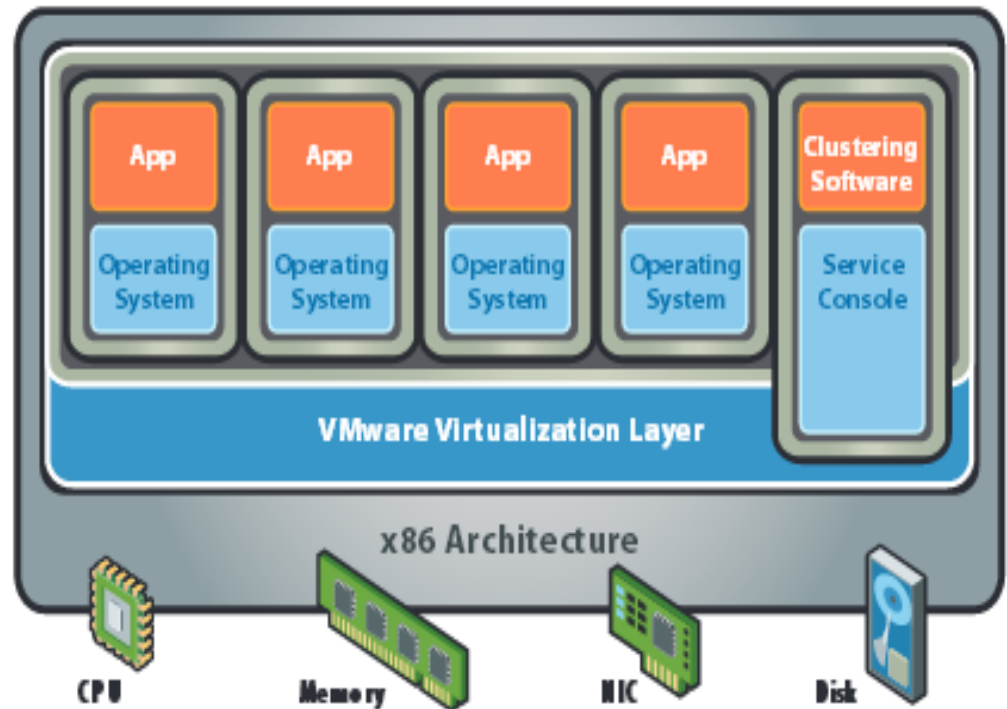
Types of Hypervisors

- For Industry-standard x86 systems, the two approaches typically used with software-based partitioning
 - Bare-metal architectures
 - Hosted
- Bare Metal / Native (Type I)
 - Hypervisor installed on a clean x86-based system.
 - Direct access to the hardware resources, a hypervisor is more efficient than hosted architectures.
 - Enables greater scalability, robustness and performance
 - E.g: VMware ESXi, Citrix XenServer, and Microsoft Hyper-V hypervisor
- Hosted (Type II)
 - A hosted approach provides partitioning services on top of a standard operating system.
 - Supports the broadest range of hardware configurations.
 - E.g: VMware Player or Parallels Desktop

x86 Hardware Virtualization



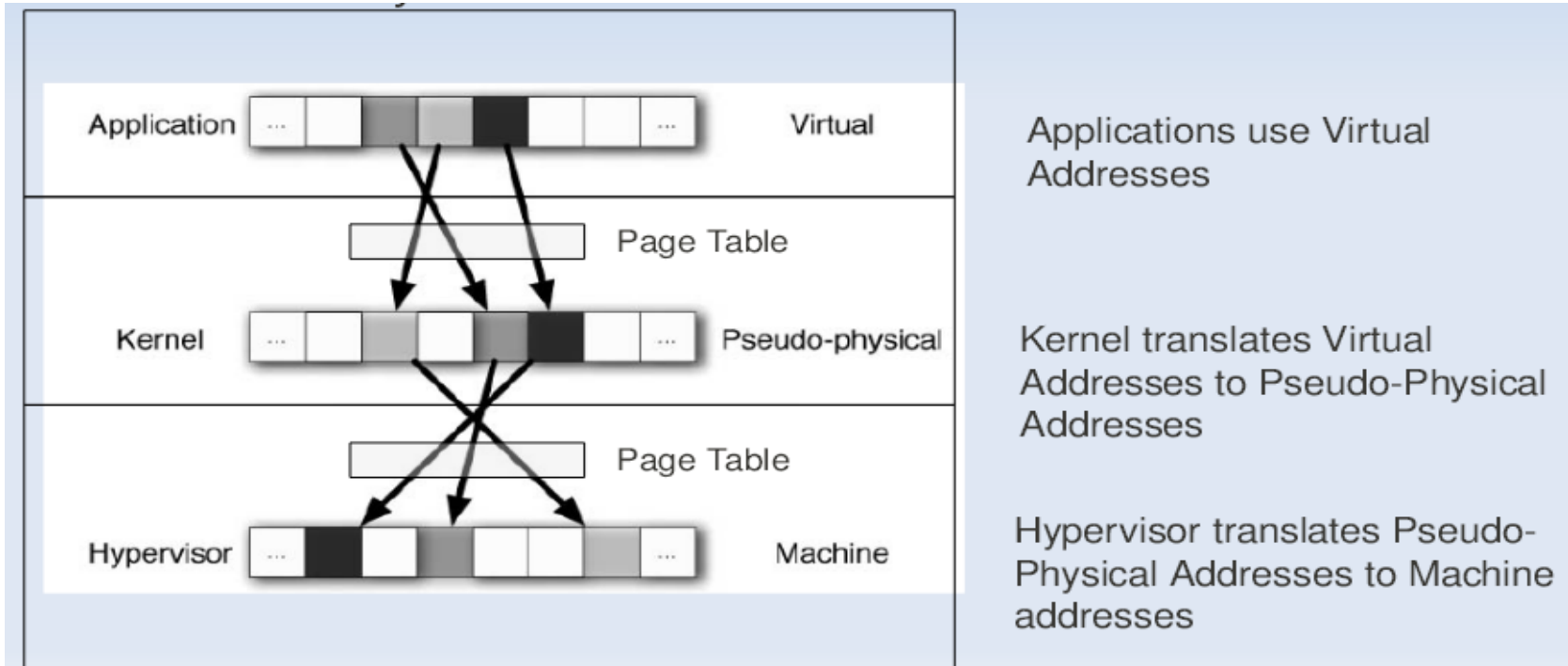
Hosted Architecture



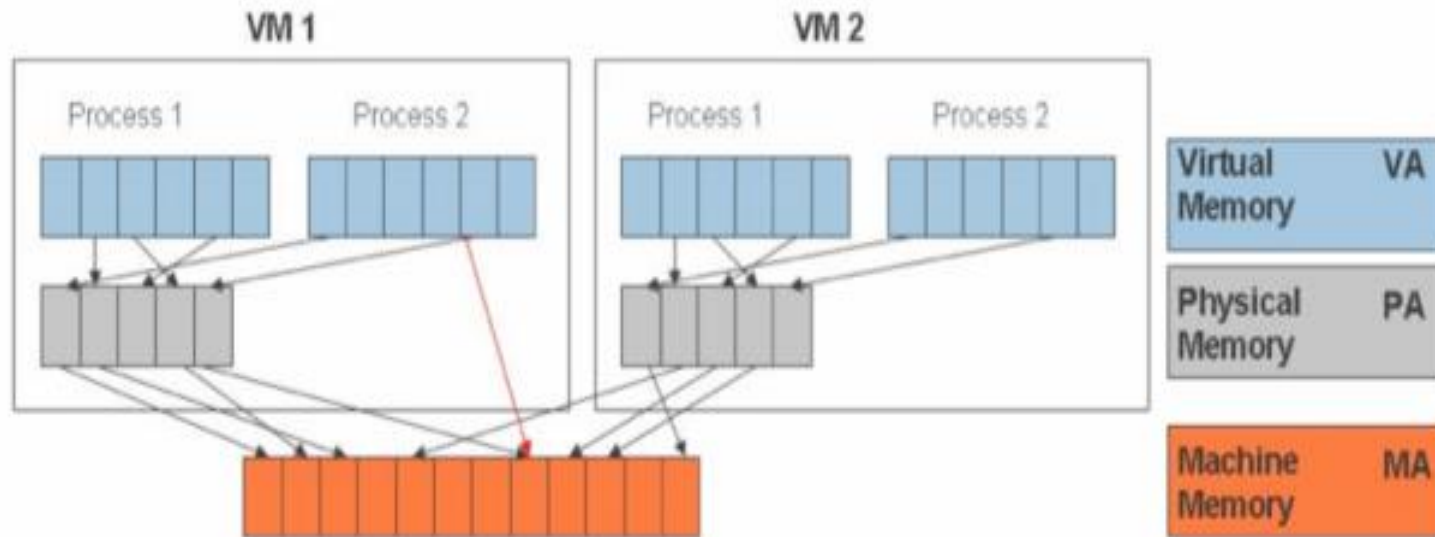
Bare-Metal (Hypervisor) Architecture

Memory Virtualization

- In multiprogramming there is a single level of indirection maintained by Kernel.
- In case of Virtual Machines there is one more level of indirection maintained by VMM



Memory Virtualization



GVA: Guest Virtual Address

GPA: Guest Physical Address

HVA: Host Virtual Address

HPA: Host Physical Address

- Guest page table has GVA->GPA mapping
- Each guest OS thinks it has access to all RAM starting at address 0
- VMM / Host OS has GPA->HPA mapping
- Guest “RAM” pages are distributed across host memory

Memory Virtualization Techniques

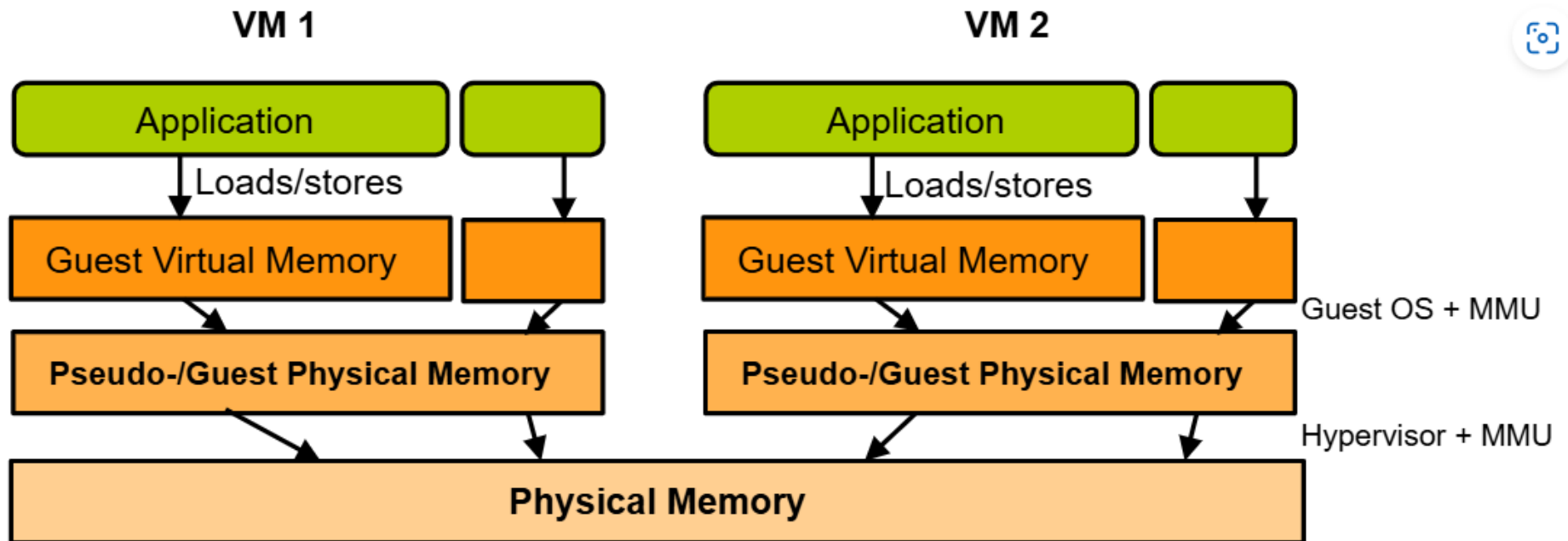
1. Shadow paging

- VMM creates a combined mapping GVA->HPA and Memory Management Unit (MMU) is given a pointer to this page table
- VMM tracks changes to guest page table and updates shadow page table

2. Extended page tables (EPT)

1. MMU hardware is aware of virtualization, takes pointers to two separate page tables.
 2. Address translation walks both page tables
- EPT is more efficient but requires hardware support

Memory Virtualization



Guest virtual memory → (guest) pseudo-physical memory → (host) physical memory

[Virtualization 101 - Introduction to Virtualization \(olivierpierre.github.io\)](https://olivierpierre.github.io)

Summary

- Virtualization is key to exploiting trends
- Virtualization Types
 - CPU Virtualization
 - Problems with standard Trap & Emulate Virtualization
 - Full: Dynamic binary translation
 - Para: Modified guest OS source code
 - Hardware-assisted: CPU has special virtualization mode
 - Memory Virtualization
 - Shadow page tables: Combined GVA->HPA mappings tracked by VMM
 - Extended page tables: Two separate pointers for GVA to GPA and GPA to HPA mappings given to MMU

Further Reading

- J.S. Robin, and C.E. Irvine, “Analysis of the Intel Pentium’s Ability to Support a Secure Virtual Machine Monitor”
- M. Rosenblum, and T. Garfinkel, “Virtual Machine Monitors: Current Technology and Future Trends”
- Bringing Virtualization to the x86 Architecture with the Original VMware Workstation”, Edouard Bugnion, Scott Devine, Mendel Rosenblum, Jeremy Sugerman, Edward Y. Wang.
- S. Devine, E. Bugnion, and M. Rosenblum, “Virtualization system including a virtual machine monitor for a computer with a segmented architecture”, US Patent 6,397,242
- P. Barham, et. al., “Xen and the Art of Virtualization”
- “Xen Developer’s Reference”
- [05_PopekGoldbergTheorem.pdf](#)