



UNIVERSITY OF
CAMBRIDGE

Cloud Computing

Virtualization

Anil Madhavapeddy

anil@recoil.org

Contents

- Virtualization.
- Layering and virtualization.
- Virtual machine monitor.
- Virtual machine.
- x86 support for virtualization.
- Full and paravirtualization.
- Xen.

- Resources:
 - Book and,
 - VMware White paper: “*Understanding Full Virtualization, Paravirtualization, and Hardware Assisted*” <https://www.vmware.com/techpapers/2007/understanding-full-virtualization-paravirtualizat-1008.html>

Motivation

- Three fundamental abstractions are necessary to describe the operation of a computing systems:
 - (1) interpreters/processors, (2) memory, (3) communications links
- As the scale of a system and the size of its users grows, it becomes very challenging to manage its recourses (see three points above)
- Resource management issues:
 - provision for peak demands → **overprovisioning**
 - **heterogeneity** of hardware and software
 - machine failures
- **Virtualization is a basic enabler of Cloud Computing, it simplifies the management of physical resources for the three abstractions**
- For example, the state of a virtual machine (VM) running under a virtual machine monitor (VMM) can be saved and migrated to another server to balance the load
- For example, virtualization allows users to operate in environments they are familiar with, rather than forcing them to specific ones

Motivation (cont' d)

- ***"Virtualization, in computing, refers to the act of creating a virtual (rather than actual) version of something, including but not limited to a virtual computer hardware platform, operating system (OS), storage device, or computer network resources."*** from Wikipedia
- Virtualization abstracts the underlying resources; simplifies their use; isolates users from one another; and supports replication which increases the elasticity of a system

Motivation (cont' d)

- Cloud resource virtualization is important for:
 - Performance isolation
 - as we can dynamically assign and account for resources across different applications
 - System security:
 - as it allows isolation of services running on the same hardware
 - Performance and reliability:
 - as it allows applications to migrate from one platform to another
 - The development and management of services offered by a provider

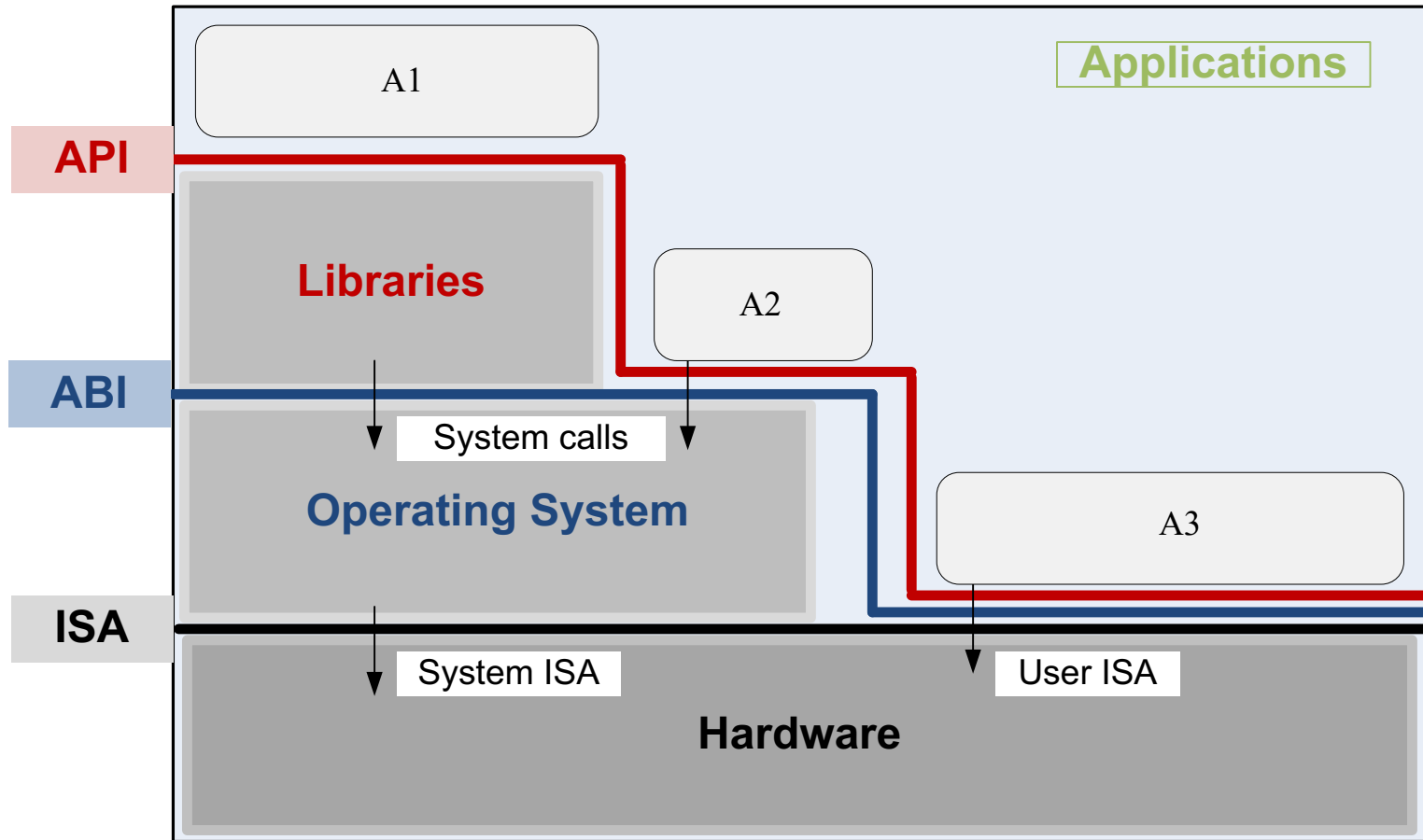
Virtualization

- Virtualization simulates the interface to a physical object by:
 - Multiplexing: creates multiple virtual objects from one instance of a physical object. Many virtual objects to one physical. Example - a processor is multiplexed among a number of processes or threads.
 - Aggregation: creates one virtual object from multiple physical objects. One virtual object to many physical objects. Example - a number of physical disks are aggregated into a RAID disk.
 - Emulation: constructs a virtual object of a certain type from a different type of a physical object. Example - a physical disk emulates a Random Access Memory (RAM).
 - Multiplexing and emulation. Examples - virtual memory with paging multiplexes real memory and disk; a virtual address emulates a real address.

Layering and Virtualization

- Layering – a common approach to manage system complexity:
 - Simplifies the description of the subsystems; each subsystem is abstracted through its *interfaces* with the other subsystems
 - Minimises the interactions among the subsystems of a complex system
 - With layering we are able to design, implement, and modify the individual subsystems independently
- Layering in a computer system:
 - Hardware
 - Software
 - Operating system
 - Libraries
 - Applications

Layering and Interfaces



Application Programming Interface (API), Application Binary Interface (ABI), and Instruction Set Architecture (ISA). An application uses library functions (A1), makes system calls (A2), and executes machine instructions (A3) (from book)

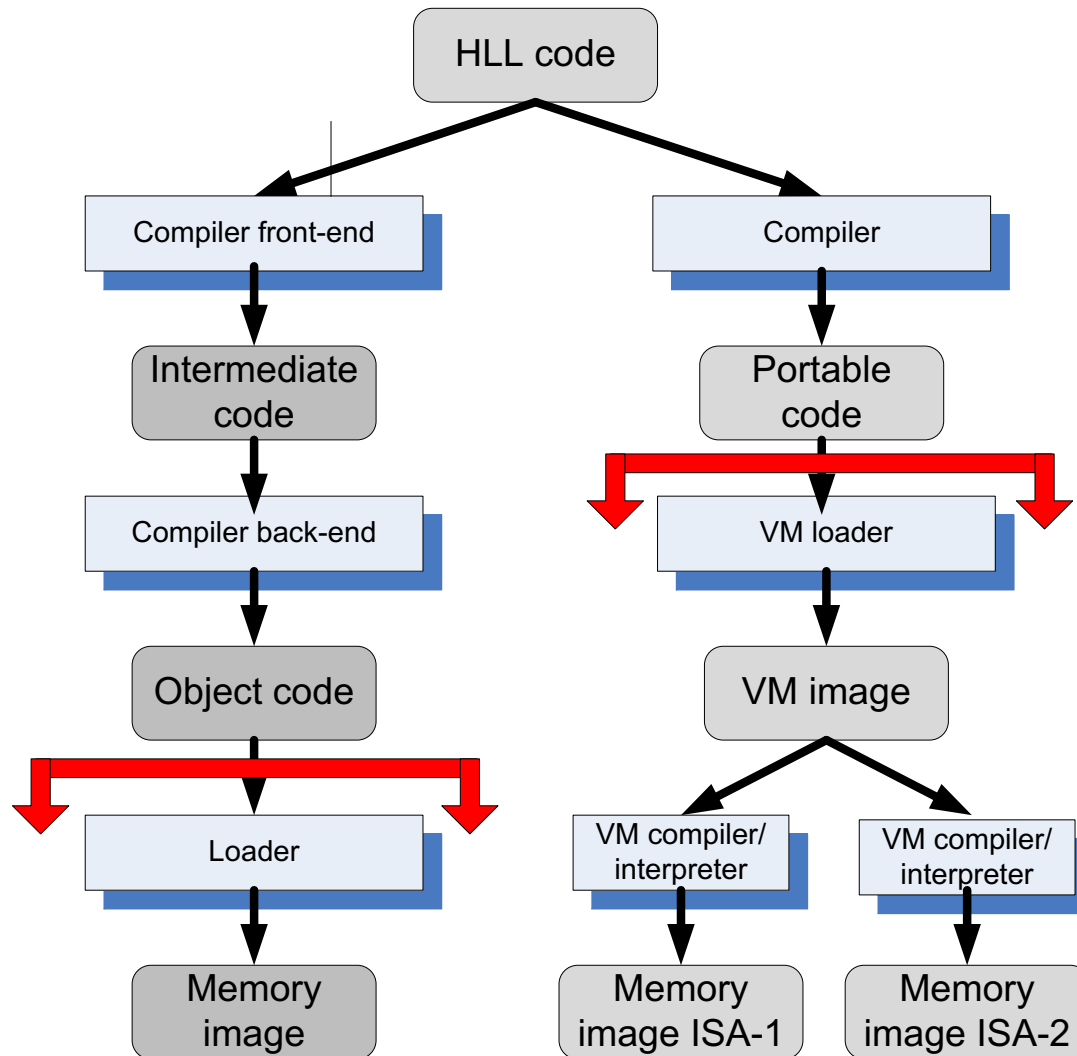
Interfaces

- **Instruction Set Architecture (ISA)** – at the boundary between hardware and software.
- **Application Binary Interface (ABI)** – allows the ensemble consisting of the application and the library modules to access the hardware; the ABI does not include *privileged* system instructions, instead it invokes system calls.
- **Application Program Interface (API)** - defines the set of instructions the hardware was designed to execute and gives the application access to the ISA; it includes high-level language (HLL) library calls which often invoke system calls

Code portability

- Binaries created by a compiler for a specific ISA and a specific operating systems are not portable
- It is possible, though, to compile a HLL program for a virtual machine (VM) environment where portable code is produced and distributed and then converted by binary translators to the ISA of the host system
- A **dynamic binary translation** converts blocks of guest instructions from the portable code to the host instruction and leads to a significant performance improvement, as such blocks are cached and reused

HLL Language Translations



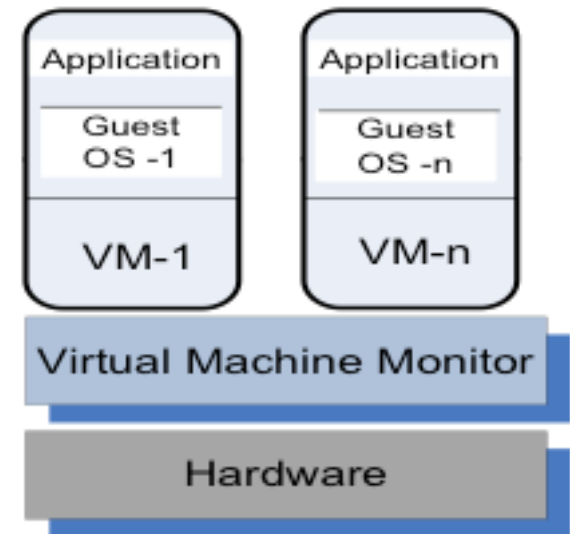
History of Virtualization

(from "Modern Operating Systems" 4th Edition, p474 by Tanenbaum and Bos)

- **1960's, IBM: CP/CMS** control program: a virtual machine operating system for the IBM System/360 Model 67
- **2000, IBM: z-series** with 64-bit virtual address spaces and backward compatible with the System/360
- **1974: Popok and Golberg** from UCLA published "*Formal Requirements for Virtualizable Third Generation Architectures*" where they listed the conditions a computer architecture should satisfy to support virtualization efficiently. The popular x86 architecture that originated in the 1970s did not support these requirements for decades.
- **1990's, Stanford researchers, VMware:** Researchers developed a new hypervisor and founded VMware, the biggest virtualization company of today's. First virtualization solution was in 1999 for x86.
- Today many virtualization solutions: Xen from Cambridge, KVM, Hyper-V, ...
- IBM was the first to produce and sell virtualization for the mainframe. But, VMware popularised virtualization for the masses.

Virtual Machine Monitor (VMM / Hypervisor)

- A **virtual machine monitor (VMM/hypervisor)** partitions the resources of computer system into one or more **virtual machines (VMs)**. Allows several operating systems to run concurrently on a single hardware platform
- A VM is an execution environment that runs an OS
- VM – an isolated environment that appears to be a whole computer, but actually only has access to a portion of the computer resources
- A VMM allows:
 - Multiple services to share the same platform
 - Live migration - the movement of a server from one platform to another
 - System modification while maintaining backward compatibility with the original system
 - Enforces isolation among the systems, thus security
- A **guest operating system** is an OS that runs in a VM under the control of the VMM.

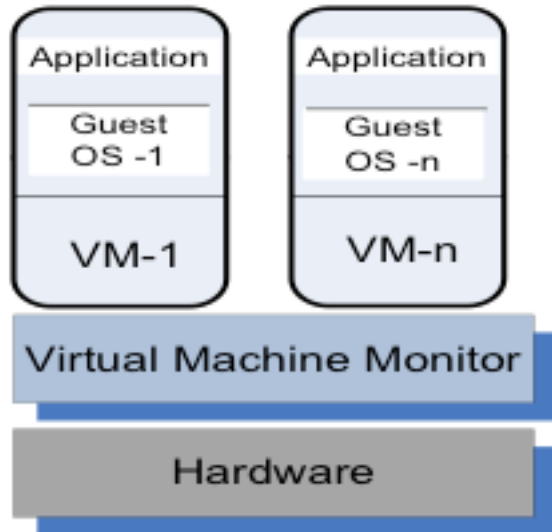


VMM Virtualizes the CPU and the Memory

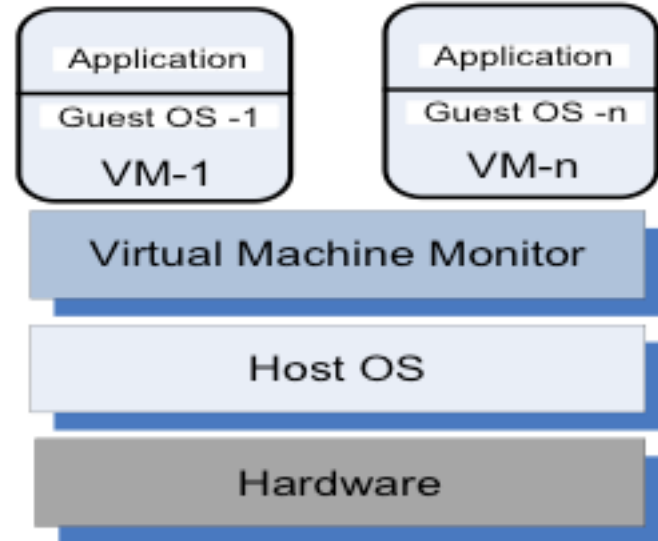
- A VMM (also hypervisor) (howto):
 - Traps the privileged instructions executed by a guest OS and enforces the correctness and safety of the operation
 - Traps interrupts and dispatches them to the individual guest operating systems
 - Controls the virtual memory management
 - Maintains a shadow page table for each guest OS and replicates any modification made by the guest OS in its own shadow page table. This shadow page table points to the actual page frame and it is used by the Memory Management Unit (MMU) for dynamic address translation.
 - Monitors the system performance and takes corrective actions to avoid performance degradation. For example, the VMM may swap out a VM to avoid thrashing.

Type 1 and 2 Hypervisors

Type 1 Hypervisor



Type 2 Hypervisor



■ Taxonomy of VMMs:

1. Type 1 Hypervisor (bare metal, native): supports multiple virtual machines and runs directly on the hardware (e.g., VMware ESX , Xen, Denali)
2. Type 2 Hypervisor (hosted) VM - runs under a host operating system (e.g., user-mode Linux)

Examples of Hypervisors

Name	Host ISA	Guest ISA	Host OS	guest OS	Company
Integrity VM	<i>x86-64</i>	<i>x86-64</i>	HP-Unix	Linux, Windows HP Unix	HP
Power VM	Power	Power	No host OS	Linux, AIX	IBM
z/VM	z-ISA	z-ISA	No host OS	Linux on z-ISA	IBM
Lynx Secure	<i>x86</i>	<i>x86</i>	No host OS	Linux, Windows	LinuxWorks
Hyper-V Server	<i>x86-64</i>	<i>x86-64</i>	Windows	Windows	Microsoft
Oracle VM	<i>x86, x86-64</i>	<i>x86, x86-64</i>	No host OS	Linux, Windows	Oracle
RTS Hypervisor	<i>x86</i>	<i>x86</i>	No host OS	Linux, Windows	Real Time Systems
SUN xVM	<i>x86, SPARC</i>	same as host	No host OS	Linux, Windows	SUN
VMware EX Server	<i>x86, x86-64</i>	<i>x86, x86-64</i>	No host OS	Linux, Windows Solaris, FreeBSD	VMware
VMware Fusion	<i>x86, x86-64</i>	<i>x86, x86-64</i>	MAC OS <i>x86</i>	Linux, Windows Solaris, FreeBSD	VMware
VMware Server	<i>x86, x86-64</i>	<i>x86, x86-64</i>	Linux, Windows	Linux, Windows Solaris, FreeBSD	VMware
VMware Workstation	<i>x86, x86-64</i>	<i>x86, x86-64</i>	Linux, Windows	Linux, Windows Solaris, FreeBSD	VMware
VMware Player	<i>x86, x86-64</i>	<i>x86, x86-64</i>	Linux Windows	Linux, Windows Solaris, FreeBSD	VMware
Denali	<i>x86</i>	<i>x86</i>	Denali	ILVACO, NetBSD	University of Washington
Xen	<i>x86, x86-64</i>	<i>x86, x86-64</i>	Linux Solaris	Linux, Solaris NetBSD	University of Cambridge

Performance and Security Isolation

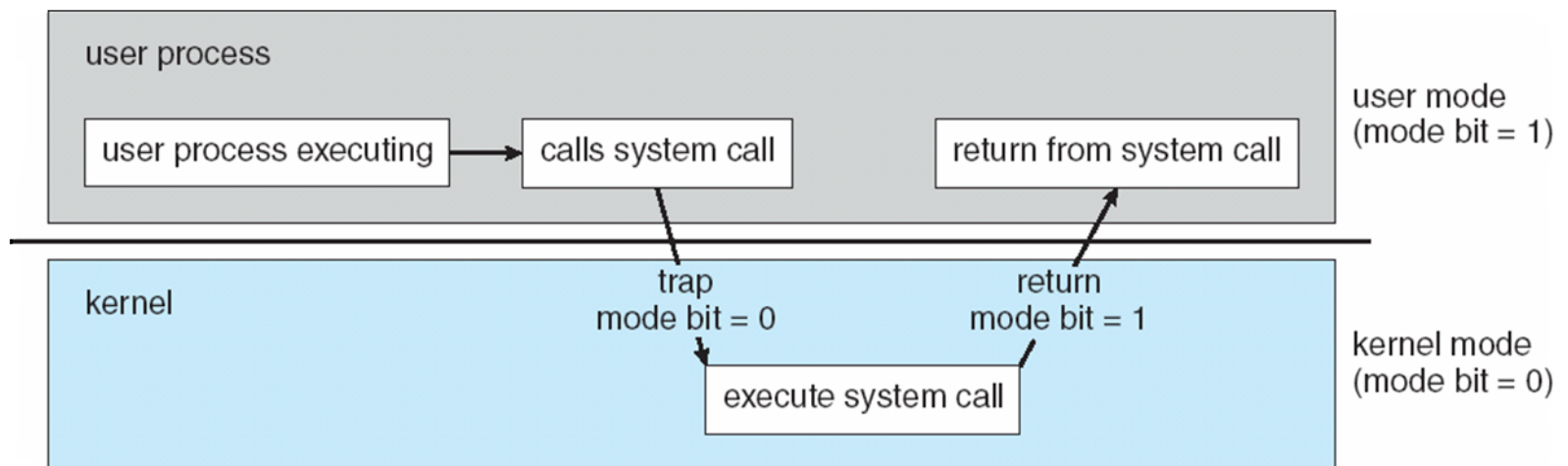
- The run-time behavior of an application is affected by other applications running concurrently on the same platform and competing for CPU cycles, cache, main memory, disk and network access. Thus, it is difficult to predict the completion time!
- Performance isolation - a critical condition for QoS guarantees in shared computing environments
- A VMM is a much simpler and better specified system than a traditional operating system. Example - Xen has approximately 60,000 lines of code; Denali has only about half: 30,000
- The security vulnerability of VMMs is considerably reduced as the systems expose a much smaller number of privileged functions. For example, Xen VMM has 28 hypercalls while Linux has 100s of system calls

Conditions for Efficient Virtualization (from Popek and Goldberg)

- Conditions for efficient virtualization (from Popek and Goldberg):
 1. A program running under the VMM should exhibit a behavior essentially identical to that demonstrated when running on an equivalent machine directly.
 2. The VMM should be in complete control of the virtualized resources.
 3. A statistically significant fraction of machine instructions must be executed without the intervention of the VMM. (Why?)

Dual-Mode Operation (recap)

- Dual-mode operation allows OS to protect itself and other system components
 - **User mode** and **kernel mode**
 - Mode bit provided by hardware
 - Ability to distinguish when system is running user or kernel code
 - Some instructions are **privileged**, only executable in kernel mode
 - System call changes mode to kernel, return resets it to user



User-mode vs Kernel-mode (recap)

- Kernel-code (in particular, interrupt handlers) runs in kernel mode
 - the hardware allows all machine instructions to be executed and allows unrestricted access to memory and I/O ports
- Everything else runs in user mode
- The OS relies very heavily on this hardware-enforced protection mechanism

Challenges of x86 CPU Virtualization

- Four layers of privilege execution → rings
 - User applications run in ring 3
 - OS runs in ring 0
- In which ring should the VMM run?
 - In ring 0, then, same privileges as an OS → wrong
 - In rings 1,2,3, then OS has higher privileges → wrong
 - Move the OS to ring 1 and the VMM in ring 0 → OK

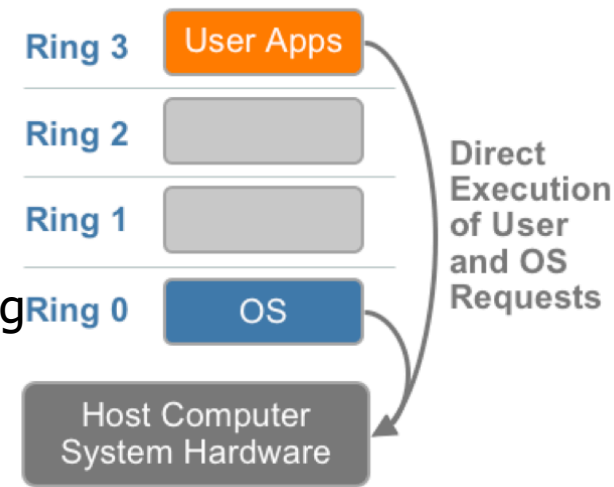


Figure 4 – x86 privilege level architecture without virtualization

- Three classes of machine instructions:
 1. **privileged instructions** can be executed in kernel mode. When attempted to be executed in user mode, they cause a *trap* and so executed in kernel mode.
 2. **nonprivileged instructions** the ones that can be executed in user mode
 3. **sensitive instructions** can be executed in either kernel or user but they behave differently. Sensitive instructions require special precautions at execution time.
 4. sensitive and nonprivileged instructions are hard to virtualize

Techniques for Virtualizing CPU on x86

- 1. Full virtualization with binary translation**
- 2. OS-assisted Virtualization or Paravirtualization**
- 3. Hardware assisted virtualization**

Techniques for Virtualizing CPU on x86

Full virtualization – a guest OS can run unchanged under the VMM as if it was running directly on the hardware platform. Each VM runs an exact copy of the actual hardware.

- **Binary translation** rewrites parts of the code on the fly to replace sensitive but not privileged instructions with safe code to emulate the original instruction
- *“The hypervisor translates all operating system instructions on the fly and caches the results for future use, while user level instructions run unmodified at native speed.”* (from VMware paper)
- Examples: VMware, Microsoft Virtual Server
- Advantages:
 - No hardware assistance,
 - No modifications of the guest OS
 - Isolation, Security
- Disadvantages:
 - Speed of execution

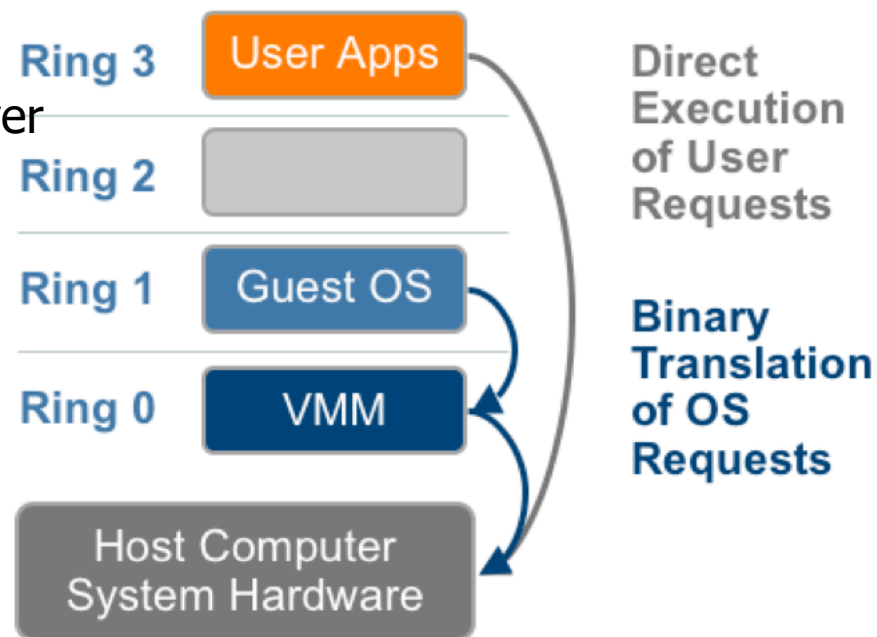


Figure 5 – The binary translation approach to x86 virtualization

Techniques for Virtualizing CPU on x86

Paravirtualization – “involves modifying the OS kernel to replace non-virtualizable instructions with hypercalls that communicate directly with the virtualization layer hypervisor. The hypervisor also provides hypercall interfaces for other critical kernel operations such as memory management, interrupt handling and time keeping.” (from VMware paper)

- Advantage: faster execution, lower virtualization overhead
- Disadvantage: poor portability
- Examples: Xen, Denali

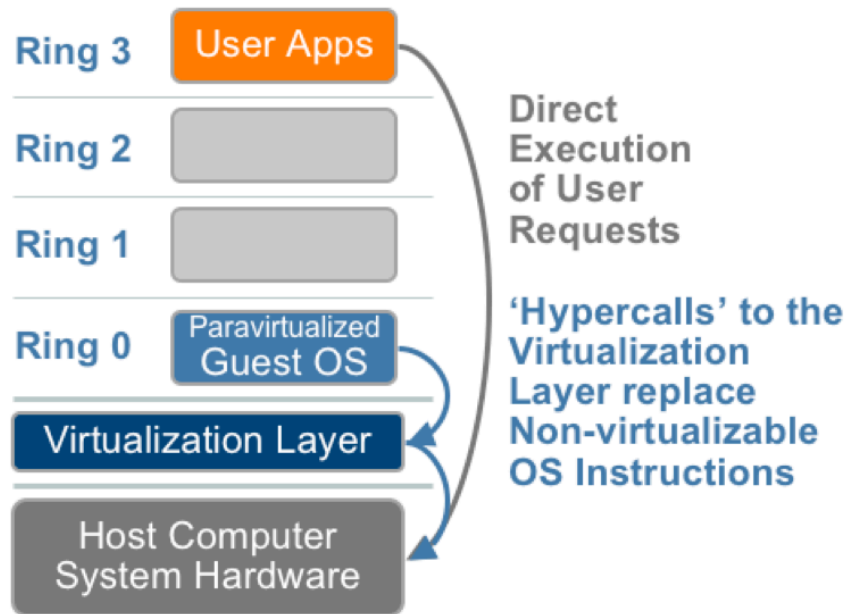
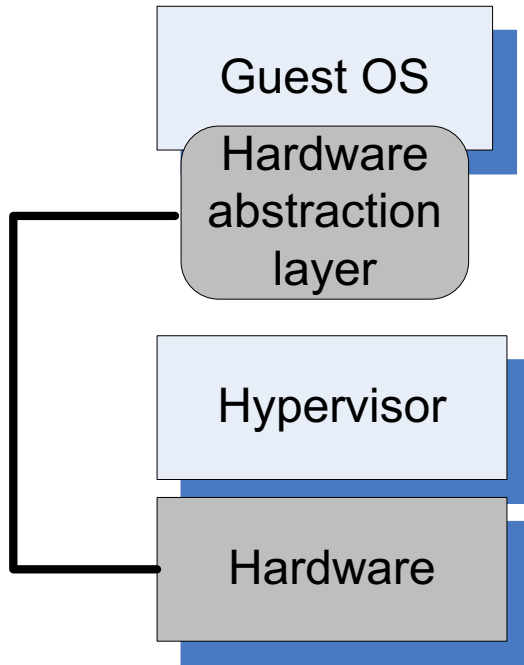
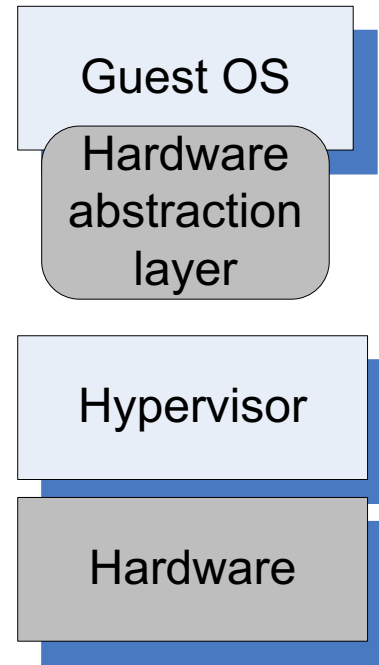


Figure 6 – The Paravirtualization approach to x86 Virtualization

Full Virtualization and Paravirtualization



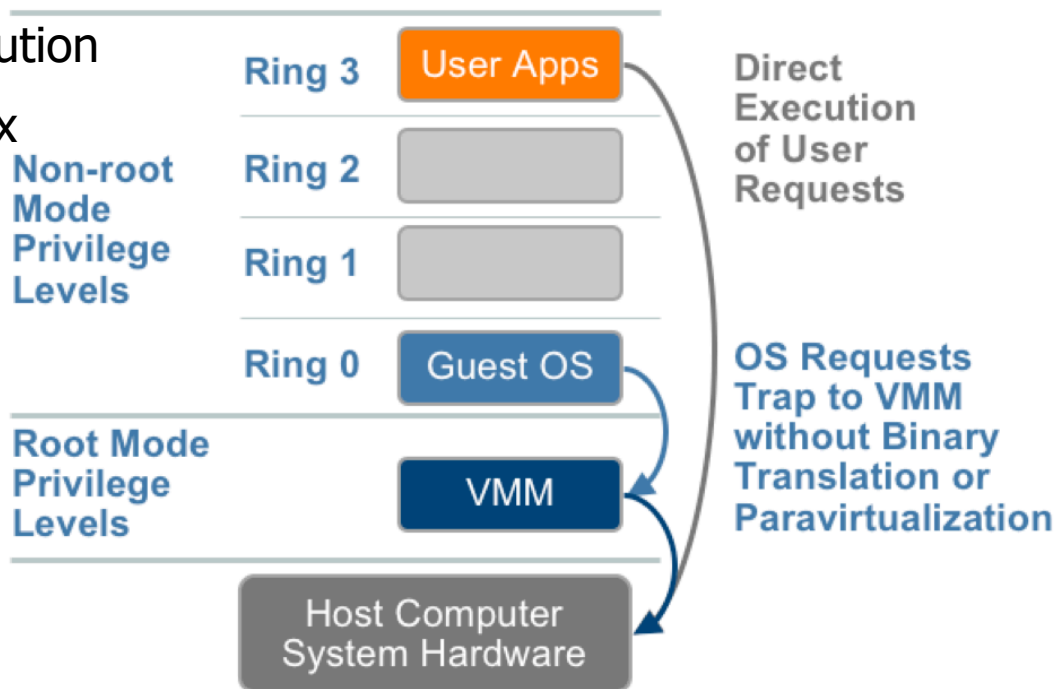
(a) Full virtualization



(b) Paravirtualization

Techniques for Virtualizing CPU on x86

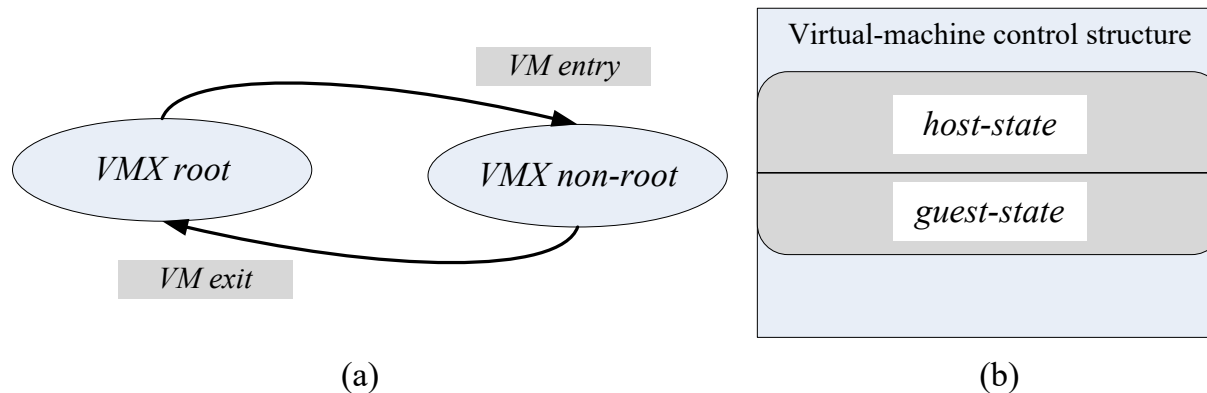
- **Hardware Assisted Virtualization** – “a new CPU execution mode feature that allows the VMM to run in a new root mode below ring 0. As depicted in Figure 7, privileged and sensitive calls are set to automatically trap to the hypervisor, removing the need for either binary translation or paravirtualization” (from VMware paper)
- Advantage: even faster execution
- Examples: Intel VT-x, Xen 3.x



1 Figure 7 – The hardware assist approach to x86 virtualization

VT-x, a Major Architectural Enhancement

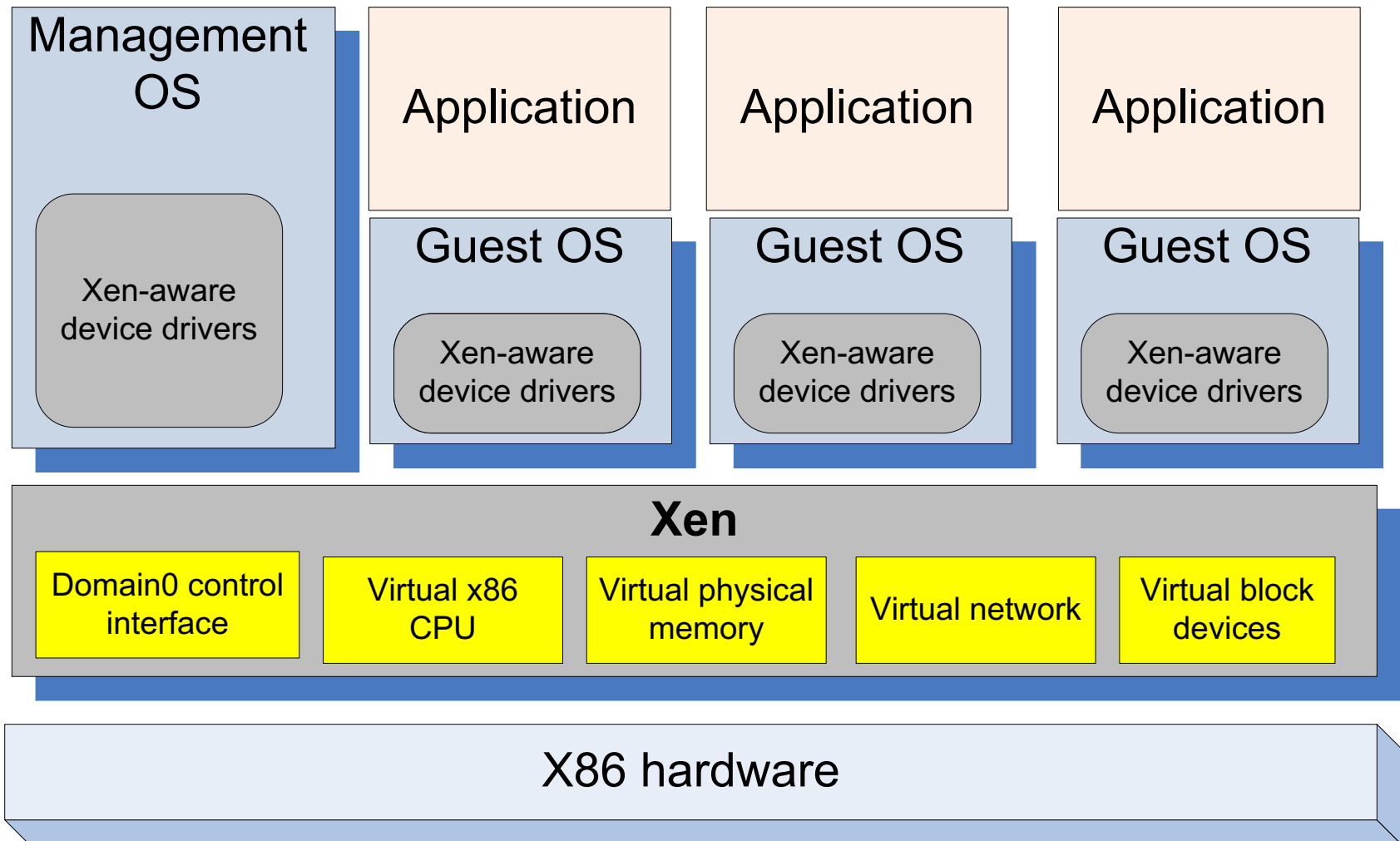
- In 2005 Intel released two Pentium 4 models supporting VT-x.
- VT-x supports two modes of operations (Figure (a)):
 1. VMX root - for VMM operations.
 2. VMX non-root - support a VM.
 - And a new data structure called the **Virtual Machine Control Structure** including *host-state* and *guest-state* areas (Figure (b)).
 - VM entry - the processor state is loaded from the guest-state of the VM scheduled to run; then the control is transferred from VMM to the VM.
 - VM exit - saves the processor state in the guest-state area of the running VM; then it loads the processor state from the host-state area, finally transfers control to the VMM.



Xen - a VMM based on Paravirtualization

- The goal of the Cambridge group - design a VMM capable of scaling to about 100 VMs running standard applications and services without any modifications to the Application Binary Interface (ABI). (2003, Computing Laboratory, Cambridge University)
- Linux, Minix, NetBSD, FreeBSD and others can operate as paravirtualized Xen guest OS running on x86, x86-64, Itanium, and ARM architectures.
- Xen domain - ensemble of address spaces hosting a guest OS and applications running under the guest OS. Runs on a virtual CPU.
 - Dom0 - dedicated to execution of Xen control functions and privileged instructions.
 - DomU - a user domain.
- Applications make system calls using hypercalls processed by Xen; privileged instructions issued by a guest OS are paravirtualized and must be validated by Xen.

Xen



Dom0 Components

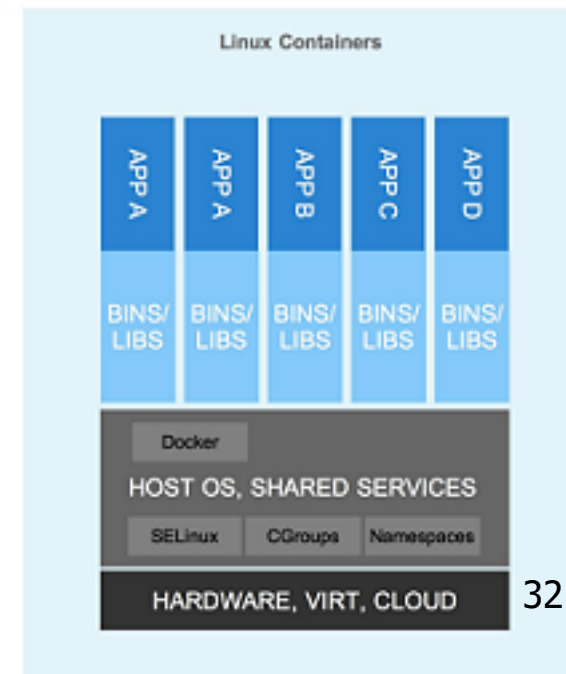
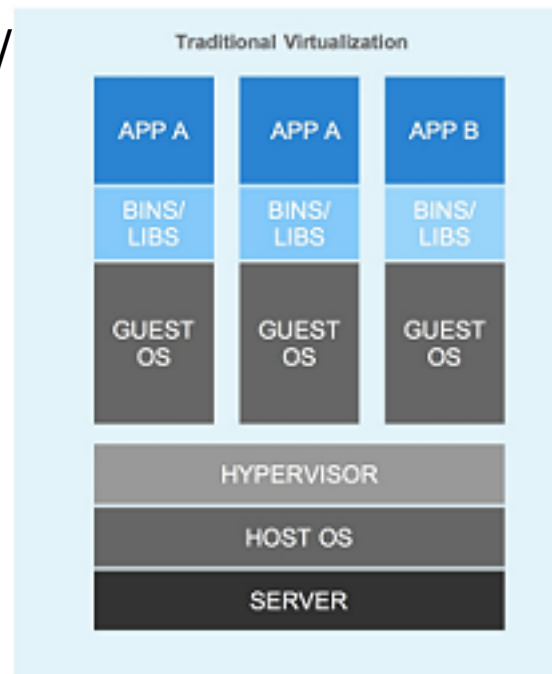
- XenStore – a Dom0 process.
 - Supports a system-wide registry and naming service.
 - Implemented as a hierarchical key-value storage.
 - A watch function informs listeners of changes of the key in storage they have subscribed to.
 - Communicates with guest VMs via shared memory using Dom0 privileges.
- Toolstack - responsible for creating, destroying, and managing the resources and privileges of VMs.
 - To create a new VM, a user provides a configuration file describing memory and CPU allocations and device configurations.
 - Toolstack parses this file and writes this information in XenStore.
 - Takes advantage of Dom0 privileges to map guest memory, to load a kernel and virtual BIOS and to set up initial communication channels with XenStore and with the virtual console when a new VM is created.

Strategies for virtual memory management, CPU multiplexing, and I/O devices

Function	Strategy
Paging	A domain may be allocated discontinuous pages. A guest OS has direct access to page tables and handles pages faults directly for efficiency; page table updates are batched for performance and validated by <i>Xen</i> for safety.
Memory	Memory is statically partitioned between domains to provide strong isolation. <i>XenoLinux</i> implements a <i>balloon driver</i> to adjust domain memory.
Protection	A guest OS runs at a lower priority level, in ring 1, while <i>Xen</i> runs in ring 0.
Exceptions	A guest OS must register with <i>Xen</i> a description table with the addresses of exception handlers previously validated; exception handlers other than the page fault handler are identical with <i>x86</i> native exception handlers.
System calls	To increase efficiency, a guest OS must install a “fast” handler to allow system calls from an application to the guest OS and avoid indirection through <i>Xen</i> .
Interrupts	A lightweight event system replaces hardware interrupts; synchronous system calls from a domain to <i>Xen</i> use <i>hypercalls</i> and notifications are delivered using the asynchronous event system.
Multiplexing	A guest OS may run multiple applications.
Time	Each guest OS has a timer interface and is aware of “real” and “virtual” time.
Network and I/O devices	Data is transferred using asynchronous I/O rings; a ring is a circular queue of descriptors allocated by a domain and accessible within <i>Xen</i> .
Disk access	Only <i>Dom0</i> has direct access to IDE and SCSI disks; all other domains access persistent storage through the Virtual Block Device (VBD) abstraction.

Linux Containers

- A Linux Container is a Linux process (or processes) that is a virtual environment with its own process network space. (lightweight process virtualization)
- Containers share portions of the host kernel
- Containers use:
 - Namespaces: per-process isolation of OS resources (filesystem, network and user ids)
 - Cgroups: resource management and accounting per process
- Examples for using containers:
 - <https://www.dotcloud.com/>
 - <https://www.heroku.com/>



Xen (old) Implementation on x86 Architecture

- Xen runs at privilege Level 0, the guest OS at Level 1, and applications at Level 3.
- The x86 architecture does not support either the tagging of TLB entries or the software management of the TLB. Thus, address space switching, when the VMM activates a different OS, requires a complete TLB flush; this has a negative impact on the performance.
- Solution - load Xen in a 64 MB segment at the top of each address space and delegate the management of hardware page tables to the guest OS with minimal intervention from Xen. This region is not accessible or re-mappable by the guest OS.
- A guest OS must register with Xen a description table with the addresses of exception handlers for validation.

Xen Abstractions for Networking and I/O

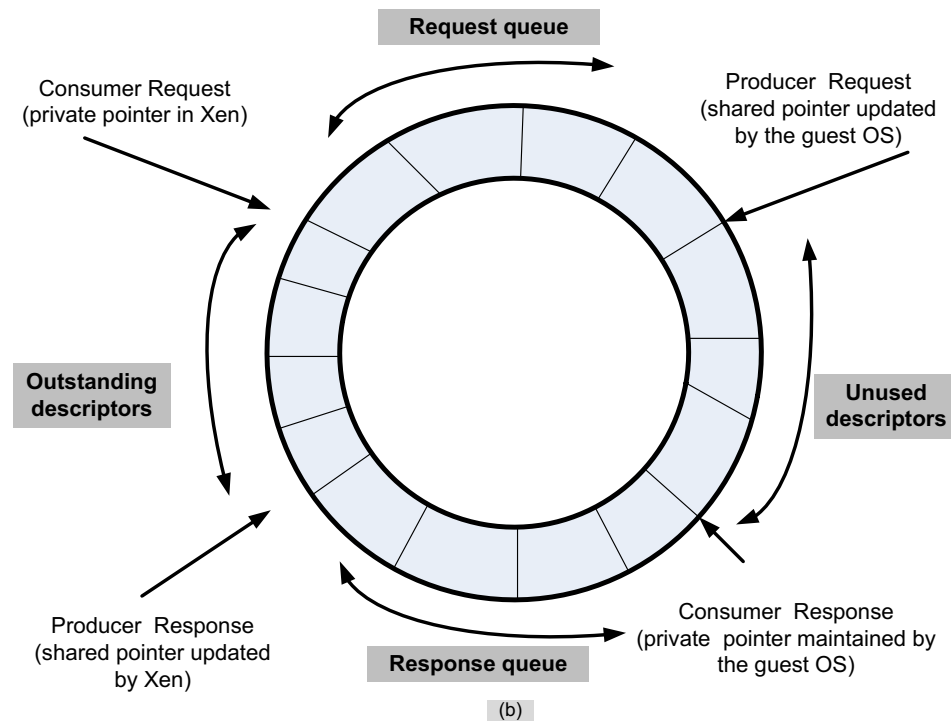
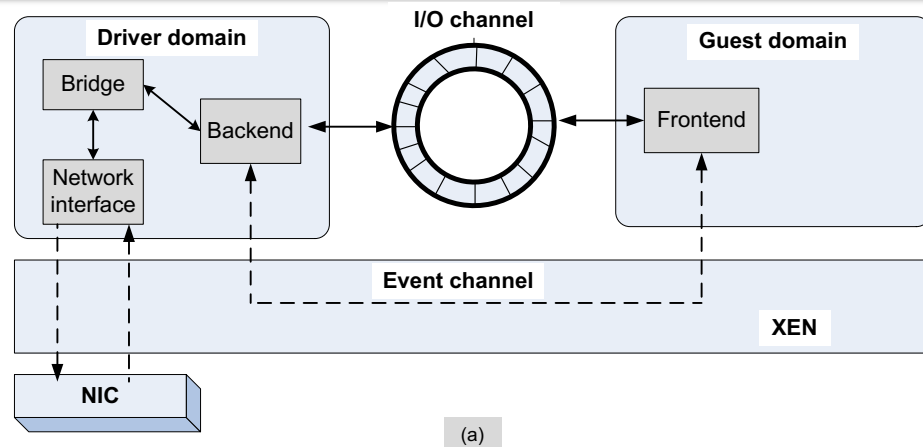
- Each domain has one or more Virtual Network Interfaces (VIFs) which support the functionality of a network interface card. A VIF is attached to a Virtual Firewall-Router (VFR).
- Split drivers have a front-end in the DomU and the back-end in Dom0; the two communicate via a ring in shared memory.
- Ring - a circular queue of descriptors allocated by a domain and accessible within Xen. Descriptors do not contain data, the data buffers are allocated off-band by the guest OS.
- Two rings of buffer descriptors, one for packet sending and one for packet receiving, are supported.
- To transmit a packet:
 - a guest OS enqueues a buffer descriptor to the send ring,
 - then Xen copies the descriptor and checks safety,
 - copies only the packet header, not the payload, and
 - executes the matching rules.

Xen I/O

Xen zero-copy semantics for data transfer using I/O rings.

(a) The communication between a guest domain and the driver domain over an I/O and an event channel; NIC is the Network Interface Controller.

(b) The circular ring of buffers.

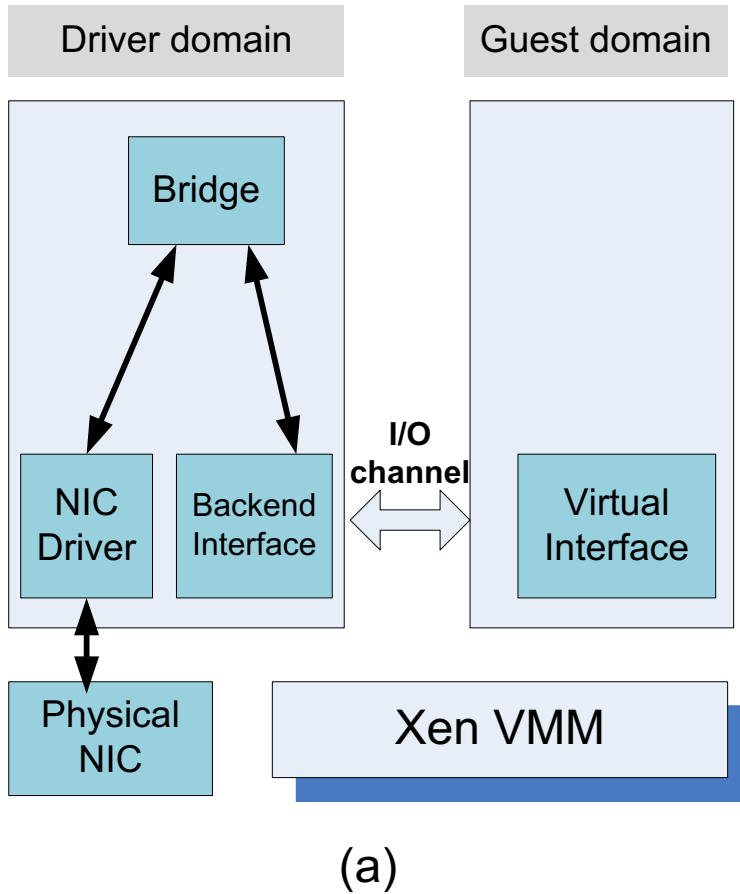


Xen 2.0

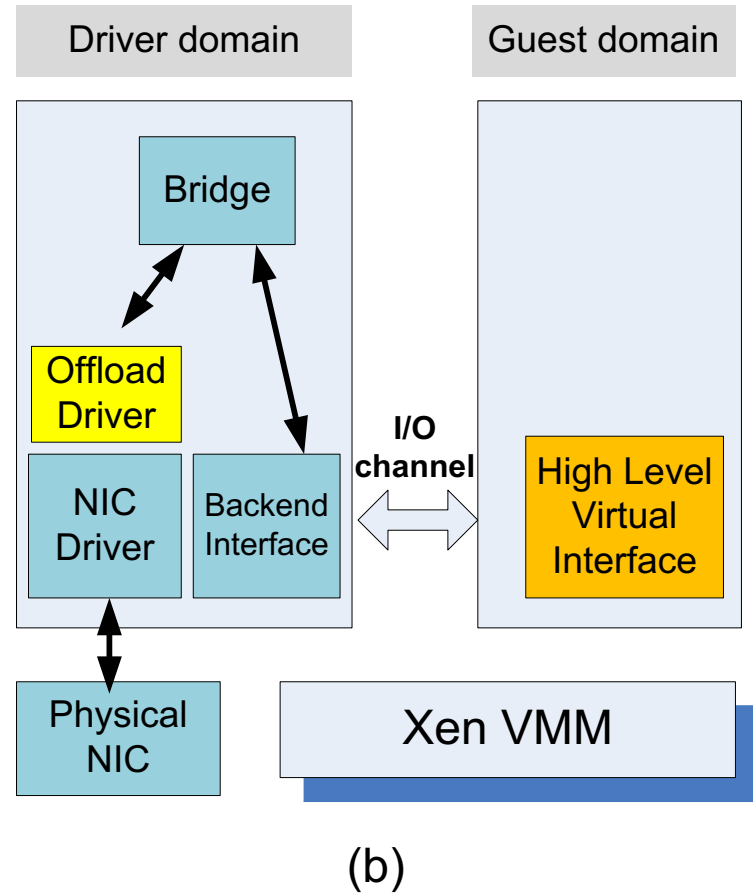
- Optimization of:
 - Virtual interface - takes advantage of the capabilities of some physical NICs, such as checksum offload.
 - I/O channel - rather than copying a data buffer holding a packet, each packet is allocated in a new page and then the physical page containing the packet is re-mapped into the target domain.
 - Virtual memory - takes advantage of the superpage and global page mapping hardware on Pentium and Pentium Pro processors. A superpage entry covers 1,024 pages of physical memory and the address translation mechanism maps a set of contiguous pages to a set of contiguous physical pages. This helps reduce the number of TLB misses.

Xen Network Architecture

The original architecture



The optimised architecture



Performance Measurements

System	Receive data rate (Mbps)	Send data rate (Mbps)
Linux	2 508	3 760
<i>Xen</i> driver	1 728	3 760
<i>Xen</i> guest	820	750
optimized <i>Xen</i> guest	970	3 310

A comparison of send and receive data rates for a native Linux system, the Xen driver domain, an original Xen guest domain, and an optimised Xen guest domain.

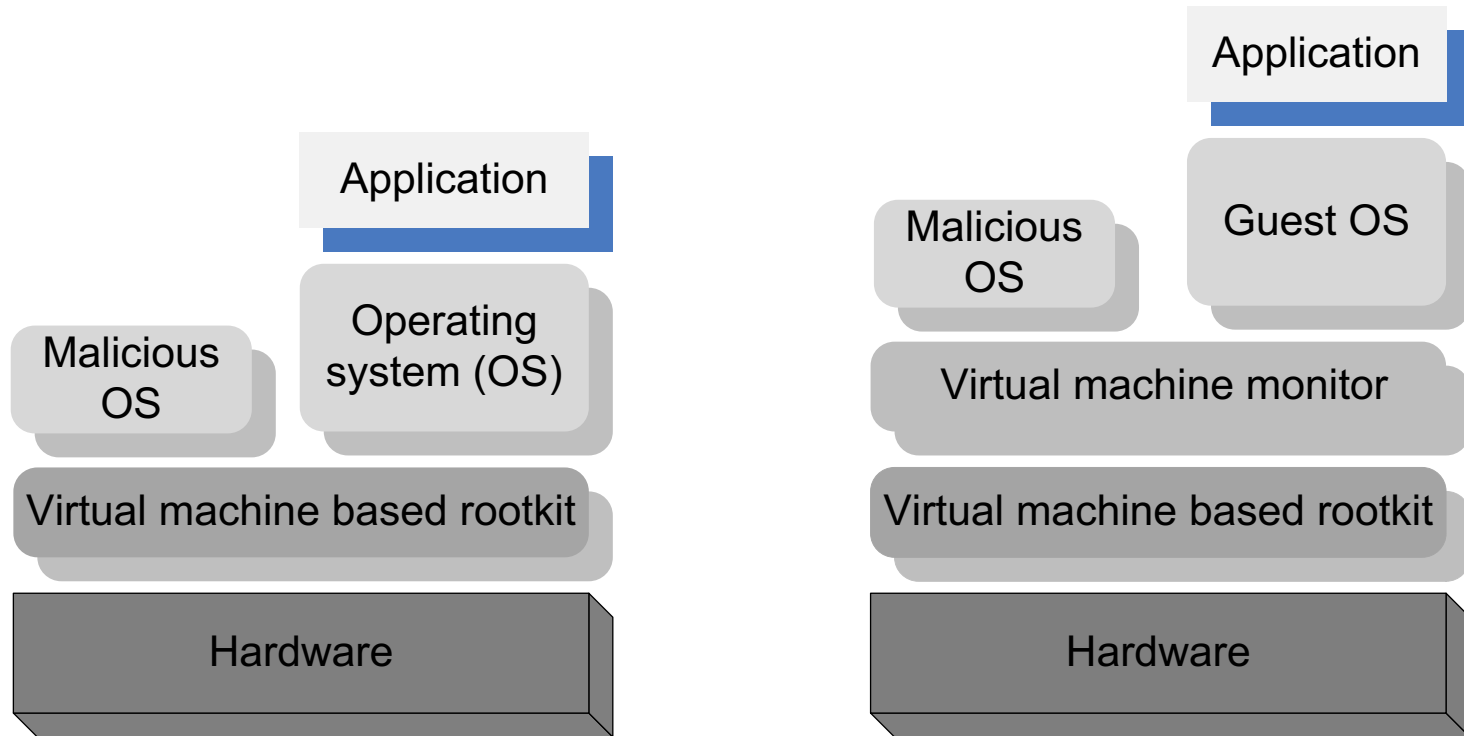
Performance Comparison of Virtual Machines

- Compare the performance of Xen and OpenVZ with, a standard operating system, a plain vanilla Linux.
- The questions examined are:
 - How the performance scales up with the load?
 - What is the impact of a mix of applications?
 - What are the implications of the load assignment on individual servers?
- The main conclusions:
 - The virtualization overhead of Xen is considerably higher than that of OpenVZ and that this is due primarily to L2-cache misses.
 - The performance degradation when the workload increases is also noticeable for Xen.
 - Hosting multiple tiers of the same application on the same server is not an optimal solution.

The Darker Side of Virtualization

- In a layered structure, a defense mechanism at some layer can be disabled by malware running at a layer below it.
- It is feasible to insert a *rogue VMM*, a Virtual-Machine Based Rootkit (VMBR) between the physical hardware and an operating system.
- Rootkit - malware with a privileged access to a system.
- The VMBR can enable a separate malicious OS to run surreptitiously and make this malicious OS invisible to the guest OS and to the application running under it.
- Under the protection of the VMBR, the malicious OS could:
 - observe the data, the events, or the state of the target system.
 - run services, such as spam relays or distributed denial-of-service attacks.
 - interfere with the application.

The Darker Side of Virtualization (con't)



(a)

(b)

The insertion of a Virtual-Machine Based Rootkit (VMBR) as the lowest layer of the software stack running on the physical hardware; (a) below an operating system; (b) below a legitimate virtual machine monitor. The VMBR enables a malicious OS to run surreptitiously and makes it invisible to the genuine or the guest OS and to the application.

Summary

- Virtualization (Chapter 5, Sections 5.1-5.8)
- Layering and virtualization.
- Virtual machine monitor.
- Virtual machine.
- x86 support for virtualization.
- Xen.