Distributed and Cloud Computing K. Hwang, G. Fox and J. Dongarra

Chapter 3: Virtual Machines and Virtualization of Clusters and datacenters

Adapted from Kai Hwang University of Southern California March 30, 2012

Copyright © 2012, Elsevier Inc. All rights reserved.

Virtualization for Datacenter Automation to serve millions of clients, simultaneously

- Server Consolidation in Virtualized Datacenter
- Virtual Storage Provisioning and Deprovisioning
- Cloud Operating Systems for Virtual Datacenters
- Trust Management in virtualized Datacenters

Difference between Traditional Computer and Virtual machines



(Courtesy of VMWare, 2008)

Virtual Machine, Guest Operating System, and VMM (Virtual Machine Monitor) :

Virtual Machine

A representation of a real machine using software that provides an operating environment which can run or host a guest operating system.

Guest Operating System

An operating system running in a virtual machine environment that would otherwise run directly on a separate physical system.

The Virtualization layer is the middleware between the underlying hardware and virtual machines represented in the system, also known as *virtual machine monitor* (VMM) or *hypervisor*.

User's view of virtualization



PHYSICAL VIEW



(Courtesy of VMWare, 2008)

Virtualization Ranging from Hardware to Applications in Five Abstraction Levels

Application level		
	JVM / .NET CLR / Panot	
Library (user-level	API) level	
	WINE/ WABI/ LxRun / Visual MainWin / vCUDA	
		\leq
Operating system	level	
	Jail / Virtual Environment / Ensim's VPS / FVM	
Hardware abstract	ion layer (HAL) level	
	VMware / Virtual PC / Denali / Xen / L4 / Plex 86 / User mode Linux / Cooperative Linux	
Instruction set arc	nitecture (ISA) level	
	Bochs / Crusoe / QEMU / BIRD / Dynamo	

Virtualization at ISA (Instruction Set Architecture) level:

Emulating a given ISA by the ISA of the host machine.

- e.g, MIPS binary code can run on an x-86-based host machine with the help of ISA emulation.
 - Typical systems: Bochs, Crusoe, Quemu, BIRD, Dynamo

Advantage:

- It can run a large amount of legacy binary codes written for various processors on any given new hardware host machines
- best application flexibility

Shortcoming & limitation:

- One source instruction may require tens or hundreds of native target instructions to perform its function, which is relatively slow.
- V-ISA requires adding a processor-specific software translation layer in the complier.

Virtualization at Hardware Abstraction level:

Virtualization is performed right on top of the hardware.

- It generates virtual hardware environments for VMs, and manages the underlying hardware through virtualization.
- Typical systems: VMware, Virtual PC, Denali, Xen

Advantage:

- Has higher performance and good application isolation
 Shortcoming & limitation:
- Very expensive to implement (complexity)

Virtualization at Operating System (OS) level:

It is an abstraction layer between traditional OS and user placations.

- This virtualization creates isolated containers on a single physical server and the OS-instance to utilize the hardware and software in datacenters.
- Typical systems: Jail / Virtual Environment / Ensim's VPS / FVM

Advantage:

• Has minimal starup/shutdown cost, low resource requirement, and high scalability; synchronize VM and host state changes.

Shortcoming & limitation:

- All VMs at the operating system level must have the same kind of guest OS
- Poor application flexibility and isolation.

Virtualization at OS Level



Figure 6.3 The virtualization layer is inserted inside an OS to partition the hardware resources for multiple VMs to run their applications in virtual environments

Virtualization for Linux and Windows NT Platforms



By far, most reported OS-level virtualization systems are Linux-based. Virtualization support on the Windows-based platform is still in the research stage. The Linux kernel offers an abstraction layer to allow software processes to work with and operate on resources without knowing the hardware details. New hardware may need a new Linux kernel to support. Therefore, different Linux platforms use patched kernels to provide special support for extended functionality.

Table 3.3 Virtualization Support for Linux and Windows NT Platforms				
Virtualization Support and Source of Information	Brief Introduction on Functionality and Application Platforms			
Linux vServer for Linux platforms (http://linux- vserver.org/)	Extends Linux kernels to implement a security mechanism to help build VMs by setting resource limits and file attributes and changing the root environment for VM isolation			
OpenVZ for Linux platforms [65]; http://ftp.openvz .org/doc/OpenVZ-Users-Guide.pdf)	Supports virtualization by creating virtual private servers (VPSes); the VPS has its own files, users, process tree, and virtual devices, which can be isolated from other VPSes, and checkpointing and live migration are supported			
FVM (Feather-Weight Virtual Machines) for virtualizing the Windows NT platforms [78])	Uses system call interfaces to create VMs at the NY kernel space; multiple VMs are supported by virtualized namespace and copy-on-write			

Advantages of OS Extension for Virtualization

- 1. VMs at OS level has minimum startup/shutdown costs
- 2. OS-level VM can easily synchronize with its environment

Disadvantage of OS Extension for Virtualization All VMs in the same OS container must have the same or similar guest OS, which restrict application flexibility of different VMs on the same physical machine.

Library Support level:

It creates execution environments for running alien programs on a platform rather than creating VM to run the entire operating system.

- It is done by API call interception and remapping.
- Typical systems: Wine, WAB, LxRun, VisualMainWin

Advantage:

It has very low implementation effort

Shortcoming & limitation:

poor application flexibility and isolation

Virtualization with Middleware/Library Support

Table 3.4 Middleware and Library Support for Virtualization					
Middleware or Runtime Library and References or Web Link	Brief Introduction and Application Platforms				
WABI (http://docs.sun.com/app/docs/doc/802-6306)	Middleware that converts Windows system calls running on x86 PCs to Solaris system calls running on SPARC workstations				
Lxrun (Linux Run) (http://www.ugcs.caltech.edu/ ~steven/lxrun/)	A system call emulator that enables Linux applications written for x86 hosts to run on UNIX systems such as the SCO OpenServer				
WINE (http://www.winehq.org/)	A library support system for virtualizing x86 processors to run Windows applications under Linux, FreeBSD, and Solaris				
Visual MainWin (http://www.mainsoft.com/)	A compiler support system to develop Windows applications using Visual Studio to run on Solaris, Linux, and AIX hosts				
vCUDA (Example 3.2) (IEEE IPDPS 2009 [57])	Virtualization support for using general-purpose GPUs to run data-intensive applications under a special guest OS				

The vCUBE for Virtualization of GPGPU



FIGURE 3.4

Basic concept of the vCUDA architecture.

(Courtesy of Lin Shi, et al. [57])

User-Application level:

It virtualizes an application as a virtual machine.

- This layer sits as an application program on top of an operating system and exports an abstraction of a VM that can run programs written and compiled to a particular abstract machine definition.
- Typical systems: JVM, NET CLI, Panot

Advantage:

• has the best application isolation

Shortcoming & limitation:

• low performance, low application flexibility and high implementation complexity.

Table 3.1 Relative Merits of Virtualization at Various Levels				
Level of Implementation	Higher Performance	Application Flexibility	Implementation Complexity	Application Isolation
ISA	Х	XXXXX	XXX	XXX
Hardware-level virtualization	XXXXX	XXX	XXXXX	XXXX
OS-level virtualization	XXXXX	XX	XXX	XX
Runtime library support	XXX	XX	XX	XX
User application level	XX	XX	XXXXX	XXXXX

More Xs mean higher merit

Hypervisor

A hypervisor is a hardware virtualization technique allowing multiple operating systems, called guests to run on a host machine. This is also called the Virtual Machine Monitor (VMM).

Type 1: bare metal hypervisor

- sits on the bare metal computer hardware like the CPU, memory, etc.
- All guest operating systems are a layer above the hypervisor.
- The original CP/CMS hypervisor developed by IBM was of this kind.

Type 2: hosted hypervisor

- Run over a host operating system.
- Hypervisor is the second layer over the hardware.
- Guest operating systems run a layer over the hypervisor.
- The OS is usually unaware of the virtualization

Major VMM and Hypervisor Providers

VMM Provider	Host CPU	Guest CPU	Host OS	Guest OS	VM Architecture
VMware Work-station	X86, x86-64	X86, x86-64	Windows, Linux	Windows, Linux, Solaris, FreeBSD, Netware, OS/2, SCO, BeOS, Darwin	Full Virtualization
VMware ESX Server	X86, x86-64	X86, x86-64	No host OS	The same as VMware workstation	Para- Virtualization
XEN	X86, x86-64, IA- 64	X86, x86- 64, IA-64	NetBSD, Linux, Solaris	FreeBSD, NetBSD, Linux, Solaris, windows XP and 2003 Server	Hypervisor
KVM	X86, x86- 64, IA64, S390, PowerPC	X86, x86- 64, IA64, S390, PowerPC	Linux	Linux, Windows, FreeBSD, Solaris	Para- Virtualization

The XEN Architecture (1)



FIGURE 3.5

The Xen architecture's special domain 0 for control and I/O, and several guest domains for user applications.

The XEN Architecture (2)

Xen is an open source hypervisor program developed by Cambridge University. Xen is a microkernel hypervisor, which separates the policy from the mechanism. The Xen hypervisor implements all the mechanisms, leaving the policy to be handled by Domain 0, as shown in Figure 3.5. Xen does not include any device drivers natively [7]. It just provides a mechanism by which a guest OS can have direct access to the physical devices. As a result, the size of the Xen hypervisor is kept rather small. Xen provides a virtual environment located between the hardware and the OS. A number of vendors are in the process of developing commercial Xen hypervisors, among them are Citrix XenServer [62] and Oracle VM [42].

The XEN Architecture (3)

The core components of a Xen system are the hypervisor, kernel, and applications. The organization of the three components is important. Like other virtualization systems, many guest OSes can run on top of the hypervisor. However, not all guest OSes are created equal, and one in particular controls the others. The guest OS, which has control ability, is called Domain 0, and the others are called Domain U. Domain 0 is a privileged guest OS of Xen. It is first loaded when Xen boots without any file system drivers being available. Domain 0 is designed to access hardware directly and manage devices. Therefore, one of the responsibilities of Domain 0 is to allocate and map hardware resources for the guest domains (the Domain U domains).

For example, Xen is based on Linux and its security level is C2. Its management VM is named Domain 0, which has the privilege to manage other VMs implemented on the same host. If Domain 0 is compromised, the hacker can control the entire system. So, in the VM system, security policies are needed to improve the security of Domain 0. Domain 0, behaving as a VMM, allows users to create, copy, save, read, modify, share, migrate, and roll back VMs as easily as manipulating a file, which flexibly provides tremendous benefits for users. Unfortunately, it also brings a series of security problems during the software life cycle and data lifetime.

Full Virtualization vs. Para-Virtualization Full virtualization

- Does not need to modify guest OS, and critical instructions are emulated by software through the use of binary translation.
- VMware Workstation applies full virtualization, which uses binary translation to automatically modify x86 software on-the-fly to replace critical instructions.
- Advantage: no need to modify OS.
- Disadvantage: binary translation slows down the performance.

Para virtualization

- Reduces the overhead, but cost of maintaining a paravirtualized OS is high.
- The improvement depends on the workload.
- Para virtualization must modify guest OS, non-virtualizable instructions are replaced by hypercalls that communicate directly with the hypervisor or VMM.
- *Para virtualization is supported by Xen, Denali and VMware* ESX.

Full Virtualization



Figure 6.9 The concept of full virtualization using a hypervisor or a VMM directly sitting on top of the bare hardware devices. Note that no host OS is used here as in Figure 6.11.

Binary **Translation** of Guest OS Requests using a VMM:



FIGURE 3.6

Indirect execution of complex instructions via binary translation of guest OS requests using the VMM plus direct execution of simple instructions on the same host.

Para- Virtualization with Compiler Support.



The KVM builds offers kernel-based VM on the Linux platform, based on para-virtualization



FIGURE 3.8

The Use of a para-virtualized guest OS assisted by an intelligent compiler to replace nonvirtualizable OS instructions by hypercalls.

VMWare ESX Server for Para-Virtualization



Memory Virtualization Challenges

Address Translation

- Guest OS expects contiguous, zero-based physical memory
- VMM must preserve this illusion

Page-table Shadowing

- VMM intercepts paging operations
- Constructs copy of page tables

Overheads

- VM exits add to execution time
- Shadow page tables consume significant host memory



Current virtual I/O devices



- Guest device driver
- Virtual device
- Virtualization layer
 - emulates the virtual device
 - remaps guest and real I/O addresses
 - multiplexes and drives the physical device
 - I/O features, e.g., COW disks,
- Real device
 - may be different from virtual device

Conclusions on CPU, Memory and I/O Virtualization :

- CPU virtualization demands hardware-assisted traps of sensitive instructions by the VMM
- Memory virtualization demands special hardware support (shadow page tables by VMWare or extended page table by Intel) to help translate virtual address into physical address and machine memory in two stages.
- I/O virtualization is the most difficult one to realize due to the complexity if I/O service routines and the emulation needed between the guest OS and host OS.

Multi-Core Virtualization: VCPU vs. traditional CPU



Figure 3.16 Four VCPUs are exposed to the software, only three cores are actually present. VCPUs V0, V1, and V3 have been transparently migrated, while VCPU V2 has been transparently suspended. (Courtesy of Wells, et al., "Dynamic Heterogeneity and the Need for Multicore Virtualization", *ACM SIGOPS Operating Systems Review,* ACM Press, 2009 [68])

Virtual Cores vs. Physical Processor Cores

Physical cores	Virtual cores
The actual physical cores present in the processor.	There can be more virtual cores visible to a single OS than there are physical cores.
More burden on the software to write applications which can run directly on the cores.	Design of software becomes easier as the hardware assists the software in dynamic resource utilization.
Hardware provides no assistance to the software and is hence simpler.	Hardware provides assistance to the software and is hence more complex.
Poor resource management.	Better resource management.
The lowest level of system software has to be modified.	The lowest level of system software need not be modified.



(Courtesy of Marty and Hill, 2007)

Virtual Clusters in Many Cores Space Sharing of VMs -- Virtual Hierarchy



(Courtesy of Marty and Hill, 2007)

Virtual Cluster Characteristics

- The virtual cluster nodes can be either physical or virtual machines. Multiple VMs running with different OSs can be deployed on the same physical node.
- A VM runs with a guest OS, which is often different from the host OS, that manages the resources in the physical machine, where the VM is implemented.
- The purpose of using VMs is to consolidate multiple functionalities on the same server. This will greatly enhance the server utilization and application flexibility.
- VMs can be colonized (replicated) in multiple servers for the purpose of promoting distributed parallelism, fault tolerance, and disaster recovery.
- The size (number of nodes) of a virtual cluster can grow or shrink dynamically, similarly to the way an overlay network varies in size in a P2P network.
- The failure of any physical nodes may disable some VMs installed on the failing nodes. But the failure of VMs will not pull down the host system.

Virtual Clusters vs. Physical Clusters



FIGURE 3.18

A cloud platform with 4 virtual clusters over 3 physical clusters shaded differently



FIGURE 3.19

The concept of a virtual cluster based on application partitioning.

(Courtesy of Kang, Chen, Tsinghua University 2008)

Live Migration of Virtual Machines



FIGURE 3.20

Live migration process of a VM from one host to another.

(Courtesy of C. Clark, et al. [14])



Effect on data transmission rate of a VM migrated from one failing web server to another.

(Courtesy of C. Clark, et al. [14])

Virtual Cluster Projects

Table 3.5 Experimental Results on Four Research Virtual Clusters				
Project Name	Design Objectives	Reported Results and References		
Cluster-on-Demand at Duke Univ.	Dynamic resource allocation with a virtual cluster management system	Sharing of VMs by multiple virtual clusters using Sun GridEngine [12]		
Cellular Disco at Stanford Univ.	To deploy a virtual cluster on a shared-memory multiprocessor	VMs deployed on multiple processors under a VMM called Cellular Disco [8]		
VIOLIN at Purdue Univ.	Multiple VM clustering to prove the advantage of dynamic adaptation	Reduce execution time of applications running VIOLIN with adaptation [25,55]		
GRAAL Project at INRIA in France	Performance of parallel algorithms in Xen-enabled virtual clusters	75% of max. performance achieved with 30% resource slacks over VM clusters		



(Courtesy of Jeff Chase, et al, HPDC-2003 [12])

Cluster-on-Demand (COD Project) at Duke University

Developed by researchers at Duke University, the COD (*Cluster on Demand*) project is a virtual cluster management system for dynamic allocation of servers from a computing pool to multiple virtual clusters [12]. The idea is illustrated by the prototype implementation of the COD shown in Figure 3.23. The COD

The Duke researchers used the Sun GridEngine scheduler to demonstrate that dynamic virtual clusters are an enabling abstraction for advanced resource management in computing utilities such as grids. The system supports dynamic, policy-based cluster sharing between local users and hosted grid services. Attractive features include resource reservation, adaptive provisioning, scavenging of idle resources, and dynamic instantiation of grid services. The COD servers are backed by a configuration database. This system provides resource policies and template definition in response to user requests.



FIGURE 3.24

Cluster size variations in COD over eight days at Duke University.

(Courtesy of J. Chase, et al. [12])



FIGURE 3.22

Live migration of VM from the DomO domain to a Xen-enabled target host.

VIOLIN Project at Purdue University

The Purdue VIOLIN Project applies live VM migration to reconfigure a virtual cluster environment. Its purpose is to achieve better resource utilization in executing multiple cluster jobs on multiple cluster

domains. The project leverages the maturity of VM migration and environment adaptation technology. The approach is to enable mutually isolated virtual environments for executing parallel applications on top of a shared physical infrastructure consisting of multiple domains. Figure 3.25 illustrates the idea with five concurrent virtual environments, labeled as VIOLIN 1–5, sharing two physical clusters.

The message being conveyed here is that the virtual environment adaptation can enhance resource utilization significantly at the expense of less than 1 percent of an increase in total execution time. The



FIGURE 3.25

VIOLIN adaptation scenario of five virtual environments sharing two hosted clusters; Note that there are more idle squares (blank nodes) before and after the adaptation.

(Courtesy of P. Ruth, et al. [24,51])

Parallax for VM Storage Management



FIGURE 3.26

Parallax is a set of per-host storage appliances that share access to a common block device and presents virtual disks to client VMs.

(Courtesy of D. Meyer, et al. [43])

Cloud OS for Building Private Clouds

1

Table 3.6 VI Managers and Operating Systems for Virtualizing Data Centers [9]					
Manager/ OS, Platforms, License	Resources Being Virtualized, Web Link	Client API, Language	Hypervisors Used	Public Cloud Interface	Special Features
Nimbus Linux, Apache v2	VM creation, virtual cluster, www .nimbusproject.org/	EC2 WS, WSRF, CLI	Xen, KVM	EC2	Virtual networks
Eucalyptus Linux, BSD	Virtual networking (Example 3.12 and [41]), www .eucalyptus.com/	EC2 WS, CLI	Xen, KVM	EC2	Virtual networks
OpenNebula Linux, Apache v2	Management of VM, host, virtual network, and scheduling tools, www.opennebula.org/	XML-RPC, CLI, Java	Xen, KVM	EC2, Elastic Host	Virtual networks, dynamic provisioning
vSphere 4 Linux, Windows, proprietary	Virtualizing OS for data centers (Example 3.13), www .vmware.com/ products/vsphere/ [66]	CLI, GUI, Portal, WS	VMware ESX, ESXi	VMware vCloud partners	Data protection, vStorage, VMFS, DRM, HA

٦.

Eucalyptus: An Open-Source OS for Setting Up and Managing Private Clouds

Eucalyptus is an open source software system (Figure 3.27) intended mainly for supporting Infrastructure as a Service (IaaS) clouds. The system primarily supports virtual networking and the management of VMs; virtual storage is not supported. Its purpose is to build private clouds that can interact with end users through Ethernet or the Internet. The system also supports interaction with other private clouds or public clouds over the Internet. The system is short on security and other desired features for general-purpose grid or cloud applications.

- Instance Manager controls the execution, inspection, and terminating of VM instances on the host where it runs.
- Group Manager gathers information about and schedules VM execution on specific instance managers, as well as manages virtual instance network.
- Cloud Manager is the entry-point into the cloud for users and administrators. It queries node managers for information about resources, makes scheduling decisions, and implements them by making requests to group managers.



FIGURE 3.27

Eucalyptus for building private clouds by establishing virtual networks over the VMs linking through Ethernet and the Internet.

(Courtesy of D. Nurmi, et al. [45])



Trusted Zones for VM Insulation

