

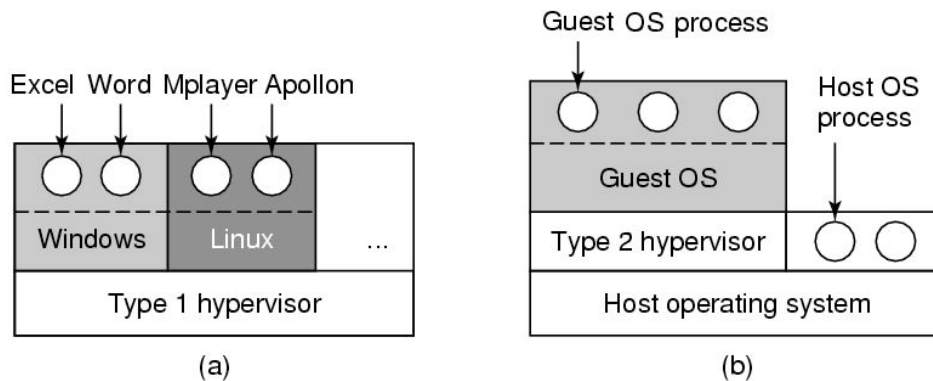
Computing Parable

- The Lion and the Fox

- Courtesy: S. Keshav



Types of Hypervisors



- Type 1: hypervisor runs on “bare metal”
- Type 2: hypervisor runs on a host OS
 - Guest OS runs inside hypervisor
- Both VM types act like real hardware

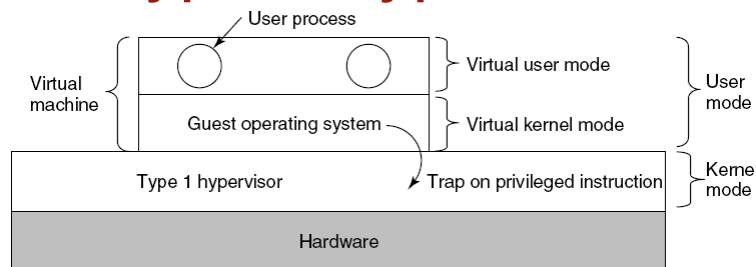


How Virtualization works?

- CPU supports kernel and user mode (ring0, ring3)
 - Set of instructions that can only be executed in kernel mode
 - I/O, change MMU settings etc -- *sensitive instructions*
 - Privileged instructions: cause a trap when executed in kernel mode
- Result: type 1 virtualization feasible if sensitive instruction subset of privileged instructions
- Intel 386: ignores sensitive instructions in user mode
 - Can not support type 1 virtualization
- Recent Intel/AMD CPUs have hardware support
 - Intel VT, AMD SVM
 - Create containers where a VM and guest can run
 - Hypervisor uses hardware bitmap to specify which inst should trap
 - Sensitive inst in guest traps to hypervisor



Type 1 hypervisor



- Unmodified OS is running in user mode (or ring 1)
 - But it thinks it is running in kernel mode (*virtual kernel mode*)
 - privileged instructions trap; sensitive inst-> use VT to trap
 - Hypervisor is the “real kernel”
 - Upon trap, executes privileged operations
 - Or emulates what the hardware would do

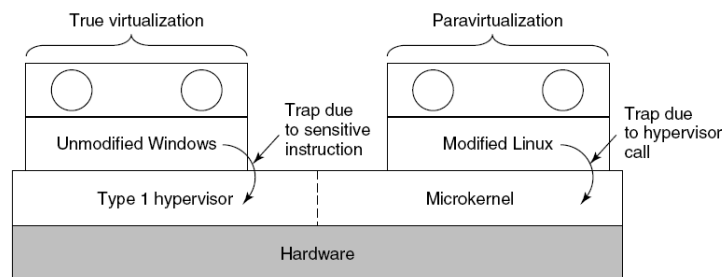


Type 2 Hypervisor

- VMWare example
 - Upon loading program: scans code for basic blocks
 - If sensitive instructions, replace by Vmware procedure
 - Binary translation
 - Cache modified basic block in VMWare cache
 - Execute; load next basic block etc.
- Type 2 hypervisors work without VT support
 - Sensitive instructions replaced by procedures that emulate them.



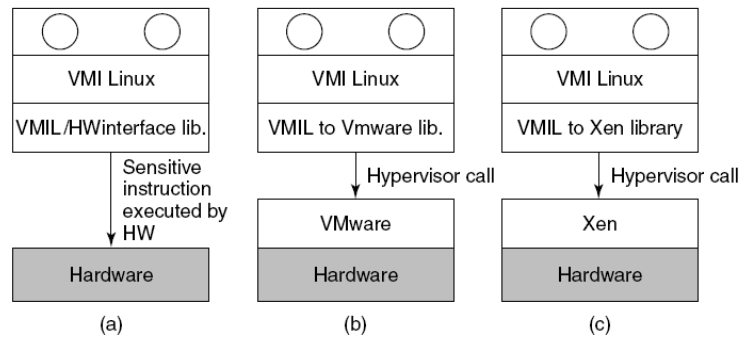
Paravirtualization



- Both type 1 and 2 hypervisors work on unmodified OS
- Paravirtualization: modify OS kernel to replace all sensitive instructions with hypercalls
 - OS behaves like a user program making system calls
 - Hypervisor executes the privileged operation invoked by hypercall.



Virtual machine Interface



- Standardize the VM interface so kernel can run on bare hardware or any hypervisor



Memory virtualization

- OS manages page tables
 - Create new pagetable is sensitive -> traps to hypervisor
- hypervisor manages multiple OS
 - Need a second shadow page table
 - OS: VM virtual pages to VM's physical pages
 - Hypervisor maps to actual page in shadow page table
 - Two level mapping
 - Need to catch changes to page table (not privileged)
 - Change PT to read-only - page fault
 - Paravirtualized - use hypercalls to inform

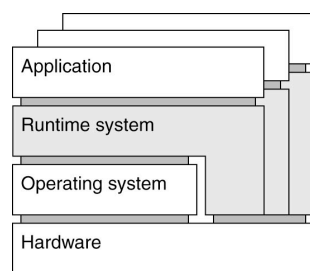


I/O Virtualization

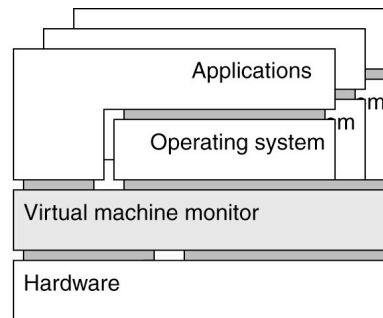
- Each guest OS thinks it “owns” the disk
- Hypervisor creates “virtual disks”
 - Large empty files on the physical disk that appear as “disks” to the guest OS
 - Hypervisor converts block # to file offset for I/O
 - DMA need physical addresses
 - Hypervisor needs to translate



Examples



(a)



(b)

- Application-level virtualization: “process virtual machine”
- VMM /hypervisor



Virtual Appliances & Multi-Core

- Virtual appliance: pre-configured VM with OS/ apps pre-installed
 - Just download and run (no need to install/comfigure)
 - Software distribution using appliances
- Multi-core CPUs
 - Run multiple VMs on multi-core systems
 - Each VM assigned one or more vCPU
 - Mapping from vCPUs to physical CPUs

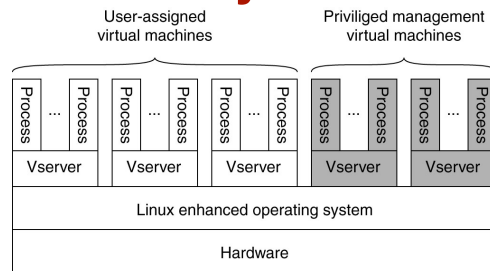


Use of Virtualization Today

- Data centers:
 - server consolidation: pack multiple virtual servers onto a smaller number of physical server
 - saves hardware costs, power and cooling costs
- Cloud computing: rent virtual servers
 - cloud provider controls physical machines and mapping of virtual servers to physical hosts
 - User gets root access on virtual server
- Desktop computing:
 - Multi-platform software development
 - Testing machines
 - Run apps from another platform



Case Study: PlanetLab



- Distributed cluster across universities
 - Used for experimental research by students and faculty in networking and distributed systems
- Uses a virtualized architecture
 - Linux Vservers
 - Node manager per machine
 - Obtain a “slice” for an experiment: slice creation service



Code and Process Migration

- Motivation
- How does migration occur?
- Resource migration
- Agent-based system
- Details of process migration



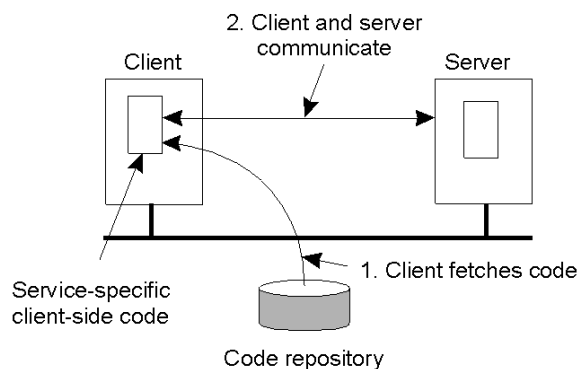
Motivation

- Key reasons: performance and flexibility
- Process migration (aka *strong mobility*)
 - Improved system-wide performance – better utilization of system-wide resources
 - Examples: Condor, DQS
- Code migration (aka *weak mobility*)
 - Shipment of server code to client – filling forms (reduce communication, no need to pre-link stubs with client)
 - Ship parts of client application to server instead of data from server to client (e.g., databases)
 - Improve parallelism – agent-based web searches



Motivation

- Flexibility
 - Dynamic configuration of distributed system
 - Clients don't need preinstalled software – download on demand



Migration models

- Process = code seg + resource seg + execution seg
- Weak versus strong mobility
 - Weak => transferred program starts from initial state
- Sender-initiated versus receiver-initiated
- Sender-initiated
 - migration initiated by machine where code resides
 - Client sending a query to database server
 - Client should be pre-registered
- Receiver-initiated
 - Migration initiated by machine that receives code
 - Java applets
 - Receiver can be anonymous

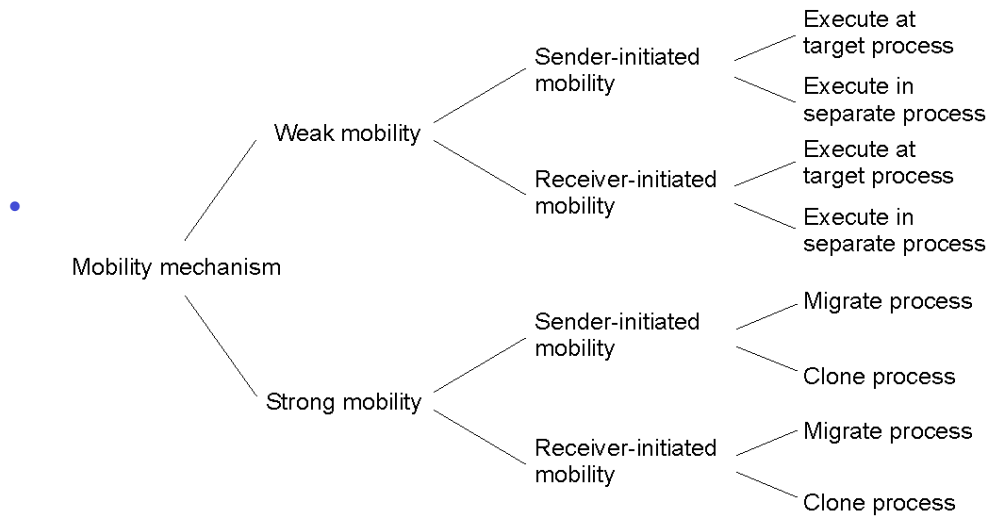


Who executes migrated entity?

- Code migration:
 - Execute in a separate process
 - [Applets] Execute in target process
- Process migration
 - Remote cloning
 - Migrate the process



Models for Code Migration



Do Resources Migrate?

- Depends on resource to process binding
 - By identifier: specific web site, ftp server
 - By value: Java libraries
 - By type: printers, local devices
- Depends on type of “attachments”
 - Unattached to any node: data files
 - Fastened resources (can be moved only at high cost)
 - Database, web sites
 - Fixed resources
 - Local devices, communication end points



Resource Migration Actions

Resource-to machine binding

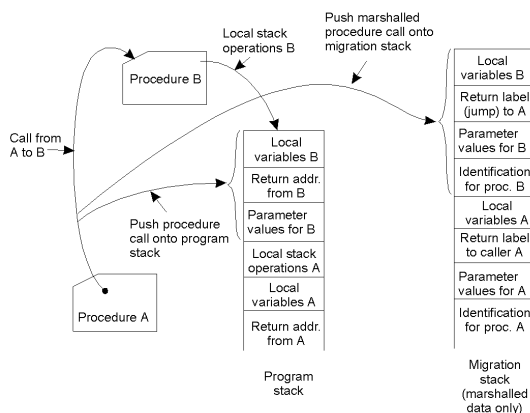
	Unattached	Fastened	Fixed	
Process-to-resource binding	By identifier	MV (or GR)	GR (or MV)	GR
	By value	CP (or MV, GR)	GR (or CP)	GR
	By type	RB (or GR, CP)	RB (or GR, CP)	RB (or GR)

- Actions to be taken with respect to the references to local resources when migrating code to another machine.
- GR: establish global system-wide reference
- MV: move the resources
- CP: copy the resource
- RB: rebind process to locally available resource



Migration in Heterogeneous Systems

- Systems can be heterogeneous (different architecture, OS)
 - Support only weak mobility: recompile code, no run time information
 - Strong mobility: recompile code segment, transfer execution segment [migration stack]
 - Virtual machines - interpret source (scripts) or intermediate code [Java]



Virtual Machine Migration

- VMs can be migrates from one physical machine to another
- Migration can be live - no application downtime
- Iterative copying of memory state



Case Study: Viruses and Malware

- Viruses and malware are examples of mobile code
 - Malicious code spreads from one machine to another
- Sender-initiated:
 - proactive viruses that look for machines to infect
 - Autonomous code
- Receiver-initiated
 - User (receiver) clicks on infected web URL or opens an infected email attachment

