

Virtualization

Dave Eckhardt and
Roger Dannenberg
based on material from:

Mike Kasick
Glenn Willen
Mike Cui

April 10, 2009

Synchronization

- **Memorial service for Timothy Wismer**
 - **Friday, April 17**
 - **16:00-18:00**
 - **Breed Hall (Margaret Morrison 103)**
 - **Sign will say “Private Event”**
 - **Donations to National Arthritis Foundation will be welcome**

Outline

- **Introduction**
- **Virtualization**
- **x86 Virtualization**
- **Paravirtualization**
- **Alternatives for Isolation**
- **Alternatives for “running two OSes on same machine”**
- **Summary**

What is Virtualization?

- **Virtualization:**
 - Process of presenting and partitioning computing resources in a *logical* way rather than partitioning according to *physical* reality
- **Virtual Machine:**
 - An execution environment (logically) identical to a physical machine, with the ability to execute a full operating system
- The *Process* abstraction is related to virtualization: it's at least similar to a physical machine

Process : Kernel :: Kernel : ?

Advantages of the Process Abstraction

- **Each process is a pseudo-machine**
- **Processes have their own registers, address space, file descriptors (sometimes)**
- **Protection from other processes**

Disadvantages of the Process Abstraction

- **Processes share the file system**
 - **Difficult to simultaneously use different versions of:**
 - **Programs, libraries, configurations**
- **Single machine owner:**
 - **root *is* the superuser**
 - **Any process that attains superuser privileges controls all processes**
 - **Other processes aren't so isolated after all**

Disadvantages of the Process Abstraction

- Processes share the same kernel
 - Kernel/OS specific software
 - Kernels are *huge*, lots of possibly buggy code
- Processes have limited degree of protection, even from each other
 - OOM (out of memory) killer (in Linux) frees memory when all else fails

Why Use Virtualization?

- “Process abstraction” at the *kernel* layer
 - Separate file system
 - Different machine owners
- Offers much better protection (in theory)
 - Secure hypervisor, fair scheduler
 - Interdomain DoS? Thrashing?
- Run two operating systems on the same machine!

Why Use Virtualization?

- **Huge impact on enterprise hosting**
 - No longer need to sell whole machines
 - Sell machine **slices**
 - Can put competitors on the same physical hardware
Can separate instance of VM from instance of hardware
- **Live migration of VM from machine to machine**
 - No more maintenance downtime
- **VM replication to provide fault tolerance**
 - “Why bother doing it at the application level?”

Disadvantages of Virtual Machines

- **Attempt to solve what really is an abstraction issue somewhere else**
 - **Monolithic kernels**
 - **Not enough partitioning of global identifiers**
 - **pids, uids, etc**
 - **Applications written without distribution and fault tolerance in mind**
- **Provides some interesting mechanisms, but may not directly solve “the problem”**

Disadvantages of Virtual Machines

- **Feasibility issues**
 - Hardware support? OS support?
 - Admin support?
 - Popularity of virtualization platforms argues these can be handled
- **Performance issues**
 - Is a 10-20% performance hit tolerable?
 - Can your NIC or disk keep up with the load?

Outline

- Introduction
- **Virtualization**
- x86 Virtualization
- Alternatives for Isolation
- Alternatives for “running two OSes on same machine”
- Summary

Full Virtualization

- **IBM CP-40 (1967)**
 - Supported 14 simultaneous S/360 virtual machines
- **Later evolved into CP/CMS and VM/CMS (still in use)**
 - 1,000 mainframe users, each with a private mainframe, running a text-based single-process “OS”
- **Popek & Goldberg: Formal Requirements for Virtualizable Third Generation Architectures (1974)**
 - Defines characteristics of a *Virtual Machine Monitor* (VMM)
 - Describes a set of architecture features sufficient to support virtualization

Virtual Machine Monitor

- **Equivalence:**
 - Provides an environment essentially identical with the original machine
- **Efficiency:**
 - Programs running under a VMM should exhibit only minor decreases in speed
- **Resource Control:**
 - VMM is in complete control of system resources

Process : Kernel :: VM : VMM

Popek & Goldberg Instruction Classification

- ***Sensitive instructions:***
 - Attempt to change configuration of system resources
 - Disable interrupts
 - Change count-down timer value
 - ...
 - Illustrate different behaviors depending on system configuration
- ***Privileged instructions:***
 - Trap if the processor is in user mode
 - Do not trap if in supervisor mode

Popek & Goldberg Theorem

“... a virtual machine monitor may be constructed if the set of sensitive instructions for that computer is a subset of the set of privileged instructions.”

- **All instructions must either:**
 - **Exhibit the same result in user and supervisor modes**
 - **Or, they must trap if executed in user mode**
- **Architectures that meet this requirement:**
 - **IBM S/370, Motorola 68010+, PowerPC, others.**

Outline

- Introduction
- Virtualization
- **x86 Virtualization**
- Paravirtualization
- Alternatives for Isolation
- Alternatives for “running two OSes on same machine”
- Summary

x86 Virtualization

- **x86 ISA does not meet the Popek & Goldberg requirements for virtualization**
- **ISA contains 17+ sensitive, unprivileged instructions:**
 - **SGDT, SIDT, SLDT, SMSW, PUSHF, POPF, LAR, LSL, VERR, VERW, POP, PUSH, CALL, JMP, INT, RET, STR, MOV**
 - **Most simply reveal the processor's CPL**
- **Virtualization is still possible, requires a workaround**

The “POPF Problem”

```
PUSHF                # %EFLAGS onto stack
ANDL $0x003FFDFD, (%ESP) # Clear IF on stack
POPF                # Stack to %EFLAGS
```

- If run in supervisor mode, interrupts are now off
- What “should” happen if this is run in user mode?

The “POPF Problem”

```
PUSHF                # %EFLAGS onto stack
ANDL $0x003FFDFE, (%ESP) # Clear IF on stack
POPF                # Stack to %EFLAGS
```

- If run in supervisor mode, interrupts are now off
- What “should” happen if this is run in user mode?
 - Attempting a privileged operation should trap
 - If it doesn't trap, the VMM can't simulate it
 - Because the VMM won't even know it happened
- What happens on the x86?

The “POPF Problem”

```
PUSHF                # %EFLAGS onto stack
ANDL $0x003FFDFE, (%ESP) # Clear IF on stack
POPF                # Stack to %EFLAGS
```

- If run in supervisor mode, interrupts are now off
- What “should” happen if this is run in user mode?
 - Attempting a privileged operation should trap
 - If it doesn't trap, the VMM can't simulate it
 - Because the VMM won't even know it happened
- What happens on the x86?
 - CPU “helpfully” ignores changes to privileged bits when POPF run in user mode!

VMware (1998)

- Runs guest operating system in ring 3
 - Maintains the illusion of running the guest in ring 0
- Insensitive instruction sequences run by CPU at full speed:
 - `movl 8(%ebp), %ecx`
 - `addl %ecx, %eax`
- Privileged instructions trap to the VMM:
 - `cli`
- VMware performs *binary translation* on guest code to work around sensitive, unprivileged instructions:
 - `popf` ⇒ `int $99`

VMware (1998)

Privileged instructions trap to the VMM:

```
cli
```

actually results in General Protection Fault (IDT entry #13), handled:

```
void gpf_exception(int vm_num, regs_t *regs)
{
    switch (vmm_get_faulting_opcode(regs->eip))
    {
        ...
        case CLI_OP:
            /* VM doesn't want interrupts now */
            vmm_defer_interrupts(vm_num);
            break;
        ...
    }
}
```

VMware (1998)

We wish `popf` trapped, but it doesn't.

Scan “code pages” of executable, translating

```
popf ⇒ int $99
```

which gets handled:

```
void popf_handler(int vm_num, regs_t *regs)
{
    regs->eflags = *(regs->esp);
    regs->esp++;
}
```

Related technologies

Software Fault Isolation (Lucco, UCB, 1993)

VX32 (Ford & Cox, MIT, 2008)

Virtual Memory

- **We've virtualized instruction execution**
 - **How about other resources?**
- **Kernels access physical memory and implements virtual memory.**
 - **How do we virtualize physical memory?**
 - **Use virtual memory (obvious so far, isn't it?)**
 - **If guest kernel runs in virtual memory, how does it provide virtual memory for processes?**
 - **VMM may have to “shadow” page-mapping tables**
 - **Set-CR3 traps, constructs *real* virtual memory**
 - **Writes to page directories and page tables are trapped, mapped to “shadow” tables**

Hardware Assisted Virtualization

- **Recent variants of the x86 ISA meet Popek & Goldberg requirements**
 - Intel VT-x (2005), AMD-V (2006)
- **VT-x introduces two new operating modes:**
 - VMX root operation & VMX non-root operation
 - VMM runs in VMX root, guest OS runs in non-root
 - Both modes support all privilege rings
 - Guest OS runs in (non-root) ring 0, no illusions necessary
- **At least initially, binary translation faster than VT**
 - `int $99` is a “regular” trap, faster than a “special trap”

Outline

- Introduction
- Virtualization
- x86 Virtualization
- **Paravirtualization**
- Alternatives for Isolation
- Alternatives for “running two OSes on same machine”
- Summary

Paravirtualization (Denali 2002, Xen 2003)

- **Motivation**
 - Binary translation and shadow page tables are hard
- **First observation:**
 - *If OS is open, it can be modified at the source level to make virtualization explicit (not transparent), and easier*
- **Paravirtualizing VMMs (hypervisors) virtualize only a subset of the x86 execution environment**
- **Run guest OS in rings 1-3**
 - No illusion about running in a virtual environment
 - Guests may not use sensitive, unprivileged instructions and expect a privileged result
- **Requires source modification only to guest kernels**
 - Abstracting page-tables is a big win
- **No modifications to user-level code and applications**

Paravirtualization (Denali 2002, Xen 2003)

- **Second observation:**
 - **Regular VMMs must emulate hardware for devices**
 - **Disk, Ethernet, etc**
 - **Performance is poor due to constrained device API**
 - **To “send packet”, must emulate many device-register accesses (inb/outb or MMIO, interrupt enable/disable)**
 - **Each step results in a trap**
 - **Already modifying guest kernel, why not provide virtual device drivers?**
 - **Virtual Ethernet could export `send_packet(addr, len)`**
 - **This requires only one trap**
- **“Hypercall” interface:**

syscall : kernel :: hypercall : hypervisor

VMware vs. Paravirtualization

- Kernel's device communication with VMware (emulated):

```
void nic_write_buffer(char *buf, int size)
{
    for (; size > 0; size--) {
        nic_poll_ready();
        outb(NIC_TX_BUF, *buf++);
    }
}
```

- Kernel's device communication with hypervisor (hypercall):

```
void nic_write_buffer(char *buf, int size)
{
    vmm_write(NIC_TX_BUF, buf, size);
}
```

Xen (2003)

- **Popular hypervisor supporting paravirtualization**
 - **Hypervisor runs on hardware**
 - **Runs two kinds of kernels**
 - **Host kernel runs in domain 0 (dom0)**
 - **Required by Xen to boot**
 - **Hypervisor contains no peripheral device drivers**
 - **dom0 needed to communicate with devices**
 - **Supports all peripherals that Linux or NetBSD do!**
 - **Guest kernels run in unprivileged domains (domU's)**

Xen (2003)

- **Provides virtual devices to guest kernels**
 - Virtual block device, virtual ethernet device
 - Devices communicate with hypercalls & ring buffers
 - Can also assign PCI devices to specific domUs
 - Video card
- **Also supports hardware assisted virtualization (HVM)**
 - Allows Xen to run unmodified domU's
 - Useful for bootstrapping
 - Also used for “the OS” that can't be source modified
- **Supports Linux & NetBSD as dom0 kernels**
- **Linux, FreeBSD, NetBSD, and Solaris as domU's**

Outline

- **Introduction**
- **Virtualization**
- **x86 Virtualization**
- **Paravirtualization**
- **Alternatives for Isolation**
- **Alternatives for “running two OSes on same machine”**
- **Summary**

chroot()

- **Venerable Unix system call**
- **Runs a Unix process with a different root directory**
 - Almost like having a separate file system
- **Share the same kernel & non-filesystem “things”**
 - Networking, process control
- **Only a minimal sandbox.**
 - /proc, /sys
 - Resources: I/O bandwidth, cpu time, memory, disk space, ...

User-mode Linux

- **Runs a guest Linux kernel as a user space process under a regular Linux kernel**
- **Requires highly modified Linux kernel**
- **No modification to application code**
- **Used to be popular among hosting providers**
- **More mature than Xen, roughly equivalent, but much slower because Xen is designed to host kernels**

Container-based OS Virtualization

- Allows multiple instances of an OS to run in isolated *containers* under the same kernel
- Assumptions:
 - Want strong separation between “virtual machines”
 - But we can trust the kernel
 - Every “virtual machine” can use the same kernel version
- It follows that:
 - Don’t need to virtualize the kernel
 - Instead, beef up naming and partitioning inside the kernel:
Each *container* can have:
 - User id, pid, tid space
 - Domain name
 - Isolated file system, OS version, libraries, etc.
- Total isolation between containers without virtualization overhead.
- VServer, FreeBSD Jails, OpenVZ, Solaris Containers (aka “Zones”)

Outline

- **Introduction**
- **Virtualization**
- **x86 Virtualization**
- **Paravirtualization**
- **Alternatives for Isolation**
- **Alternatives for “running two OSes on same machine”**
- **Summary**

Full System Simulation (Simics 1998)

- **Software simulates hardware components that make up a target machine**
- **Interpreter executes each instruction & updates the software representation of the hardware state**
- **Approach is very accurate but very slow**
- **Great for OS development & debugging**
- **Break on triple fault is better than a reset**

System Emulation (Bochs, DOSBox, QEMU)

- **Seeks to emulate just enough of system hardware components to create an accurate “user experience”**
- **Typically CPU & memory subsystems are emulated**
 - **Buses are not**
 - **Devices communicate with CPU & memory directly**
- **Many shortcuts taken to achieve better performance**
 - **Reduces overall system accuracy**
 - **Code designed to run correctly on real hardware executes “pretty well”**
 - **Code not designed to run correctly on real hardware exhibits wildly divergent behavior**
- **E.g. run legacy 680x0 code on PowerPC, run Windows on ??**

System Emulation Techniques

- **Pure interpretation:**
 - Interpret each guest instruction
 - Perform a semantically equivalent operation on host
- **Static translation:**
 - Translate each guest instruction to host once
 - Happens at startup
 - Limited applicability, no self-modifying code

System Emulation Techniques

- **Dynamic translation:**
 - Translate a block of guest instructions to host instructions just prior to execution of that block
 - Cache translated blocks for better performance
- **Dynamic recompilation & adaptive optimization:**
 - Discover what algorithm the guest code implements
 - Substitute with an optimized version on the host
 - Hard

Outline

- **Introduction**
- **Virtualization**
- **x86 Virtualization**
- **Paravirtualization**
- **Alternatives for Isolation**
- **Alternatives for “running two OSes on same machine”**
- **Summary**

Are We Having Fun Yet?

- **Virtualization is great if you need it**
 - If you must have 35 /etc/passwd's, 35 sets of users, 35 Ethernet cards, etc.
 - There are many techniques, which work (are secure and fast enough)
- **Virtualization is overkill if we need only isolation**
 - Remember the Java “virtual machine”??
 - Secure isolation for multiple applications
 - Old approach –Smalltalk (1980)
 - New approach –Google App Engine
- **Open question**
 - How **best** to get isolation, machine independence?

Summary

- **Virtualization is big in enterprise hosting**
- **{Full, hardware assisted, para-}virtualization**
- **Containers: VM-like abstraction with high efficiency**
- **Emulation is a slower alternative, more flexibility**

Further Reading

- Gerald J. Popek and Robert P. Goldberg.
Formal requirements for virtualizable third generation architectures.
Communications of the ACM, 17(7):412-421, July 1974.
- John Scott Robin and Cynthia E. Irvine.
Analysis of the intel pentium's ability to support a secure virtual machine monitor.
In *Proceedings of the 9th USENIX Security Symposium*, Denver, CO, August 2000.
- Gil Neiger, Amy Santoni, Felix Leung, Dion Rodgers, and Rich Uhlig.
Intel Virtualization Technology: Hardware support for efficient processor virtualization.
Intel Technology Journal, 10(3):167-177, August 2006.
- Paul Barham, Boris Dragovic, Keir Fraser, Steven Hand, Tim Harris, Alex Ho, Rolf Neugebauer, Ian Pratt, and Andrew Warfield.
Xen and the art of virtualization.
In *Proceedings of the 19th ACM Symposium on Operating Systems Principles*, pages 164-177, Bolton Landing, NY, October 2003.
- Yaozu Dong, Shaofan Li, Asit Mallick, Jun Nakajima, Kun Tian, Xuefei Xu, Fred Yang, and Wilfred Yu. Extending Xen with Intel Virtualization Technology.
Intel Technology Journal, 10(3):193-203, August 2006.
- Stephen Soltesz, Herbert Potzl, Marc E. Fiuczynski, Andy Bavier, and Larry Peterson.
Container-based operating system virtualization: A scalable, high-performance alternative to hypervisors.
In *Proceedings of the 2007 EuroSys conference*, Lisbon, Portugal, March 2007.
- Fabrice Bellard.
QEMU, a fast and portable dynamic translator.
In *Proceedings of the 2005 USENIX Annual Technical Conference*, Anaheim, CA, April 2005.