



COMP9242 Advanced Operating Systems S2/2012 Week 4: **Virtualization**



Australian Government

Department of Broadband, Communications and the Digital Economy





THE UNIVERSITY OF SYDNEY



Queensland Government

NICTA Funding and Supporting Members and Partners



Griffith





QUT





Copyright Notice



These slides are distributed under the Creative Commons Attribution 3.0 License

- You are free:
 - to share—to copy, distribute and transmit the work
 - to remix-to adapt the work
- under the following conditions:
 - Attribution: You must attribute the work (but not in any way that suggests that the author endorses you or your use of the work) as follows:
 - "Courtesy of Gernot Heiser, [Institution]", where [Institution] is one of "UNSW" or "NICTA"

The complete license text can be found at http://creativecommons.org/licenses/by/3.0/legalcode



Virtual Machine (VM)



"A VM is an efficient, isolated duplicate of a real machine"

- → Duplicate: VM should behave identically to the real machine
 - Programs cannot distinguish between execution on real or virtual hardware
 - Except for:
 - Fewer resources available (and potentially different between executions)
 - Some timing differences (when dealing with devices)
- → Isolated: Several VMs execute without interfering with each other
- → Efficient: VM should execute at speed close to that of real hardware
 - Requires that most instruction are executed directly by real hardware

Hypervisor aka virtual-machine monitor: Software implementing the VM





© 2012 Gernot Heiser UNSW/NICTA. Distributed under Creative Commons Attribution License



Why Virtual Machines?



- Historically used for easier sharing of expensive mainframes
 - Run several (even different) OSes on same machine
 - called *guest operating system*
 - Each on a subset of physical resources
 - Can run single-user single-tasked OS in time-sharing mode
 - legacy support
- Gone out of fashion in 80's
 - Time-sharing OSes common-place
 - Hardware too cheap to worry...





Why Virtual Machines?

- Renaissance in recent years for improved isolation
- Server/desktop virtual machines • Improved QoS and Gernot prediction of 2004: 2014 OS textbooks will be - Uniform view identical to 2004 version Complete en except for replication Apps VI2 s/process/VM/g • migration checkpointing Guest Guest OS OS debugging Different concurrent OSes • eg Linux + Windows Virt RAM Virt RAM Total mediation \bigcirc Hyperv sor Would be mostly unnecessary ٠

Mem. region

– ... if OSes were doing their job!



Mem. region

RAM



Why Virtual Machines?



- Embedded systems: integration of heterogenous environments
 - RTOS for critical real-time functionality
 - Standard OS for GUIs, networking etc
- Alternative to physical separation
 - low-overhead communication
 - cost reduction





Hypervisor



- Program that runs on real hardware to implement the virtual machine
- Controls resources
 - Partitions hardware
 - Schedules guests
 - "world switch"
 - Mediates access to shared resources
 - e.g. console
- Implications
 - Hypervisor executes in *privileged* mode
 - Guest software executes in *unprivileged* mode
 - Privileged instructions in guest cause a trap into hypervisor
 - Hypervisor interprets/emulates them
 - Can have extra instructions for hypercalls







Native/Classic/ Bare-metal/Type-I



Hosted/Type-II



- Hosted VMM can run besides native apps
 - Sandbox untrusted apps
 - Convenient for running alternative OS on desktop
- → Less efficient
 - Twice number of mode switches
 - Twice number of context switches
 - Host not optimised for exception forwarding





- → Traditional "*trap and emulate*" approach:
 - guest attempts to access physical resource
 - hardware raises exception (trap), invoking hypervisor's exception handler
 - hypervisor emulates result, based on access to virtual resource
- → Most instructions do not trap
 - makes efficient virtualization possible
 - requires that VM ISA is (almost) same as physical processor ISA









Definitions:

- → **Privileged instruction**: executes in privileged mode, *traps in user mode*
 - Note: trap is required, NO-OP is insufficient!
- → **Privileged state**: determines resource allocation
 - Includes privilege mode, addressing context, exception vectors, ...
- Sensitive instruction: control-sensitive or behaviour-sensitive
 - control sensitive: changes privileged state
 - behaviour sensitive: exposes privileged state
 -includes instructions which are NO-OPs in user but not privileged mode
- → Innocuous instruction: not sensitive

Note:

- Some instructions are inherently sensitive
 - -e.g. TLB load
- Others are sensitive in some context

 e.g. store to page table





Trap-and-emulate virtualizable if all sensitive instructions are privileged

- → Can then achieve accurate, efficient guest execution
 - by simply running guest binary on hypervisor
- → VMM controls resources
- → Virtualized execution is indistinguishable from native, except:
 - Resources more limited (running on smaller machine)
 - Timing is different (if there is an observable time source)
- → Recursively virtualizable machine:
 - VMM can be built without any timing dependence





Impure Virtualization

- \rightarrow Used for two reasons:
 - Architecture not trap-and-emulate virtualizable
 - Reduce virtualization overheads
- → Change the guest OS, replacing sensitive instructions
 - by trapping code (hypercalls)
 - by in-line emulation code
- → Two standard approaches:
 - binary translation: modifies binary
 - para-virtualization: changes ISA









- → Locate sensitive instructions in guest binary and replace on-the-fly by emulation code or hypercall
 - pioneered by VMware
 - can also detect combinations of sensitive instructions and replace by single emulation
 - doesn't require source, uses unmodified native binary
 -in this respect appears like pure virtualization!
 - very tricky to get right (especially on x86!)
 - "heroic effort" [Orran Krieger, then IBM later VMware ;-)]
 - needs to make some assumptions on sane behaviour of guest

Para-Virtualization

- → New name, old technique
 - Mach Unix server [Golub et al, 90], L⁴Linux [Härtig et al, 97], Disco [Bugnion et al, 97]
 - Name coined by Denali [Whitaker et al, 02], popularised by Xen [Barham et al, 03]
- → Idea: manually port the guest OS to modified ISA
 - Augment by explicit hypervisor calls (*hypercalls*)

 Use more high-level API to reduce the number of traps
 Remove un-virtualizable instructions
 - -Remove "messy" ISA features which complicate virtualization
 - · Generally out-performs pure virtualization and binary-rewriting
- → Drawbacks:
 - Significant engineering effort
 - Needs to be repeated for each guest-ISA-hypervisor combination
 - · Para-virtualized guest needs to be kept in sync with native guest
 - Requires source







Virtualization Overheads



- → VMM needs to maintain virtualized privileged machine state
 - processor status
 - addressing context
 - device state...
- → VMM needs to emulate privileged instructions
 - translate between virtual and real privileged state
 - e.g. guest ↔ real page tables
- → Virtualization traps are be expensive on modern hardware
 - can be 100s of cycles (1150 cycles round-trip on latest Intel x86 processors)
- → Some OS operations involve frequent traps
 - STI/CLI for mutual exclusion
 - frequent page table updates during fork()...
 - MIPS KSEG address used for physical addressing in kernel



Virtualization Techniques



- → Impure virtualization methods enable new optimisations
 - due to the ability to control the ISA
- → E.g. maintain some virtual machine state inside VMM:
 - e.g. interrupt-enable bit (in virtual PSR)
 - requires changing guest's idea of where this bit lives
 - hypervisor knows about VMM-local virtual state and can act accordingly -e.g. queue virtual interrupt until guest enables in virtual PSR





Virtualization Techniques



 \rightarrow E.g. lazy update of virtual machine state

- virtual state is kept inside hypervisor
- keep copy of virtual state inside VM
- allow temporary inconsistency between local copy and real VM state
- synchronise state on next forced hypervisor invocation -actual trap
 - -explicit hypercall when physical state must be updated
- Example:guest enables FPU
 - -no need to invoke hypervisor at this point
 - -hypervsior syncs state on virtual kernel exit







Must implement with single MMU translation!



Virtualization Mechanics: Shadow Page Table











Virtualization Mechanics: Real Guest PT







Virtualization Mechanics: Optimised Guest PT



COMP9242 S2/2012 W04 23 © 2012 Gernot Heiser UNSW/NICTA. Distributed under Creative Commons Attribution License







Virtualization Mechanics: Emulated Device











Virtualization Mechanics: Driver OS (Xen Dom0)







Virtualization Mechanics: Pass-Through Driver







Non-Virtualizable Architectures



- → x86: lots of non-virtualizable features
 - e.g. sensitive PUSH of PSW is not privileged
 - segment and interrupt descriptor tables in virtual memory
 - segment description expose privileged level
- → Itanium: mostly virtualizable, but
 - interrupt vector table in virtual memory
 - THASH instruction exposes hardware page tables address
- → MIPS: mostly virtualizable, but
 - kernel registers k0, k1 (needed to save/restore state) user-accessible
 - performance issue with virtualizing KSEG addresses
- → ARM: mostly virtualizable, but
 - some instructions undefined in user mode (banked registers, CPSR)
 - PC is a GPR, exception return in MOVS to PC, doesn't trap
- → Most others have problems too
- Modern trend are virtualization extensions to ISA
 - x86, Itanium since ~2006 (VT-x, VT-i)
- → Case study: ARM
 - announced '10, samples '11, products '12



ARM Virtualization Extensions (1)

Hyp mode

"Non-Secure" world	"Secure" world			
User mode				
Kernel modes	User mode			
Hyp mode	Kernel modes			
Monitor mode				

- New privilege level
 - Strictly higher than kernel
 - Virtualizes or traps all sensitive instructions
 - Only available in ARM TrustZone "non-secure" mode
- Note: different from x86
 - VT-x "root" mode is orthogonal to x86 protection rings



NICTA

ARM Virtualization Extensions (2)



Configurable Traps





ARM Virtualization Extensions (3)





ARM Virtualization Extensions (3)







ARM Virtualization Extensions (4)

NICTA

2-stage translation





ARM Virtualization Extensions (4)



- On page fault walk twice

2-stage translation cost





ARM Virtualization Extensions (5)

Virtual Interrupts



- ARM has 2-part IRQ controller
 - Global "distributor"

•

- Per-CPU "interface"
- New H/W "virt. CPU interface"
 - Mapped to guest
 - Used by HV to forward IRQ
 - Used by guest to acknowledge
- Halves hypervisor entries for interrupt virtualization



NICTA



Hypervisor	ISA	Туре	Kernel	User
OKL4	ARMv7	para-virtualization	9.8 kLOC	0
Prototype	ARMv7	pure virtualization	6 kLOC	0
Nova	x86	pure virtualization	9 kLOC	27 kLOC

- Size (& complexity) reduced about 40% wrt to para-virtualization
- Much smaller than x86 pure-virtualization hypervisor
 - Mostly due to greatly reduced need for instruction emulation



Overheads (Estimated)



	Pure virtualization		Para-virtualiz.
Operation	Instruct	Cycles (est)	Cycles (approx)
Guest system call	0	0	300
Hypervisor entry + exit	120	650	150
IRQ entry + exit	270	900	300-400?
Page fault	356	1500	700
Device emul.	249	1040	N/A
Device emul. (accel.)	176	740	N/A
World switch	2824	7555	200

- No overhead on regular (virtual) syscall unlike para-virtualization
- Invoking hypervisor 500–1200 cycles (0.6–1.5 μs) more than para
- World switch in ~10 μ s compared to 0.25 μ s for para
- ⇒ Trade-offs differ



Hypervisors vs Microkernels



- Both contain all code executing at highest privilege level
 - Although hypervisor may contain user-mode code as well
 - privileged part usually called "hypervisor"
 - user-mode part often called "VMM"
- Both need to abstract hardware resources
 - Hypervisor: abstraction closely models hardware

Difference to traditional terminology!

- Microkernel: abstraction designed to support wide range of systems
- What must be abstracted?
 - Memory
 - CPU
 - I/O
 - Communication







Closer Look at I/O and Communication



- Communication is critical for I/O
 - Microkernel IPC is highly optimised
 - Hypervisor inter-VM communication is frequently a bottleneck



Hypervisors vs Microkernels: Summary



- Fundamentally, both provide similar abstractions
- Optimised for different use cases
 - Hypervisor designed for virtual machines
 - API is hardware-like to ease guest ports
 - Microkernel designed for multi-server systems
 - seems to provide more OS-like abstractions



Hypervisors vs Microkernels: Drawbacks



- Communication is Achilles heel
 - more important than expected
 - critical for I/O
 - plenty improvement attempts in Xen
- Most hypervisors have big TCBs
 - infeasible to achieve high assurance of security/safety
 - in contrast, microkernel implementations can be proved correct

Microkernels:

- Not ideal for virtualization
 - API not very effective
 - L4 virtualization performance close to hypervisor
 - effort much higher
 - Virtualization needed for legacy
- L4 model uses kernelscheduled threads for more than exploiting parallelism
 - Kernel imposes policy
 - Alternatives exist, eg. K42 uses scheduler activations

More on this later!



