

# The Turtles Project: Design and Implementation of Nested Virtualization

Muli Ben-Yehuda<sup>†</sup> Michael D. Day<sup>‡</sup> Zvi Dubitzky<sup>†</sup> Michael Factor<sup>†</sup>  
Nadav Har'El<sup>†</sup> Abel Gordon<sup>†</sup> Anthony Liguori<sup>‡</sup> Orit Wasserman<sup>†</sup>  
Ben-Ami Yassour<sup>†</sup>

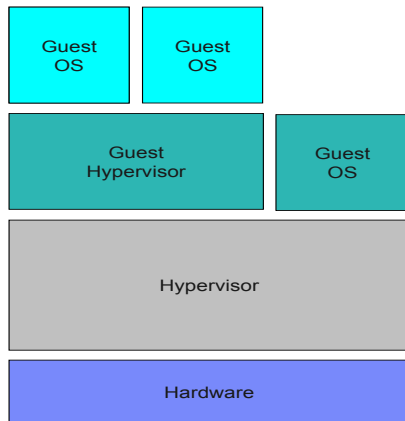
<sup>†</sup>IBM Research – Haifa

<sup>‡</sup>IBM Linux Technology Center



# What is nested x86 virtualization?

- Running multiple **unmodified** hypervisors
- With their associated unmodified VM's
- Simultaneously
- On the x86 architecture
- Which does **not support nesting in hardware**...
- ...but does support a single level of virtualization



# Why?

- Operating systems are already hypervisors (Windows 7 with XP mode, Linux/KVM)
- To be able to run other hypervisors in **clouds**
- Security (e.g., hypervisor-level rootkits)
- Co-design of x86 hardware and system software
- Testing, demonstrating, debugging, live migration of hypervisors



# Why?

- Operating systems are already hypervisors (Windows 7 with XP mode, Linux/KVM)
- To be able to run other hypervisors in **clouds**
- Security (e.g., hypervisor-level rootkits)
- Co-design of x86 hardware and system software
- Testing, demonstrating, debugging, live migration of hypervisors



# Why?

- Operating systems are already hypervisors (Windows 7 with XP mode, Linux/KVM)
- To be able to run other hypervisors in **clouds**
- Security (e.g., hypervisor-level rootkits)
- Co-design of x86 hardware and system software
- Testing, demonstrating, debugging, live migration of hypervisors



# Why?

- Operating systems are already hypervisors (Windows 7 with XP mode, Linux/KVM)
- To be able to run other hypervisors in **clouds**
- Security (e.g., hypervisor-level rootkits)
- Co-design of x86 hardware and system software
- Testing, demonstrating, debugging, live migration of hypervisors



# Why?

- Operating systems are already hypervisors (Windows 7 with XP mode, Linux/KVM)
- To be able to run other hypervisors in **clouds**
- Security (e.g., hypervisor-level rootkits)
- Co-design of x86 hardware and system software
- Testing, demonstrating, debugging, live migration of hypervisors



- First models for nested virtualization [[PopekGoldberg74](#), [BelpaireHsu75](#), [LauerWyeth73](#)]
- First implementation in the IBM z/VM; relies on architectural support for nested virtualization ([sie](#))
- Microkernels meet recursive VMs [[FordHibler96](#)]: assumes we can modify software at all levels
- x86 software based approaches (slow!) [[Berghmans10](#)]
- KVM [[KivityKamay07](#)] with AMD SVM [[RoedelGraf09](#)]
- Early Xen prototype [[He09](#)]
- Blue Pill rootkit hiding from other hypervisors [[Rutkowska06](#)]





## Related work

- First models for nested virtualization [[PopekGoldberg74](#), [BelpaireHsu75](#), [LauerWyeth73](#)]
- First implementation in the IBM z/VM; relies on architectural support for nested virtualization ([sie](#))
- Microkernels meet recursive VMs [[FordHibler96](#)]: assumes we can modify software at all levels
- x86 software based approaches (slow!) [[Berghmans10](#)]
- KVM [[KivityKamay07](#)] with AMD SVM [[RoedelGraf09](#)]
- Early Xen prototype [[He09](#)]
- Blue Pill rootkit hiding from other hypervisors [[Rutkowska06](#)]



## Related work

- First models for nested virtualization [[PopekGoldberg74](#), [BelpaireHsu75](#), [LauerWyeth73](#)]
- First implementation in the IBM z/VM; relies on architectural support for nested virtualization ([sie](#))
- Microkernels meet recursive VMs [[FordHibler96](#)]: assumes we can modify software at all levels
- x86 software based approaches (slow!) [[Berghmans10](#)]
- KVM [[KivityKamay07](#)] with AMD SVM [[RoedelGraf09](#)]
- Early Xen prototype [[He09](#)]
- Blue Pill rootkit hiding from other hypervisors [[Rutkowska06](#)]



## Related work

- First models for nested virtualization [[PopekGoldberg74](#), [BelpaireHsu75](#), [LauerWyeth73](#)]
- First implementation in the IBM z/VM; relies on architectural support for nested virtualization ([sie](#))
- Microkernels meet recursive VMs [[FordHibler96](#)]: assumes we can modify software at all levels
- x86 software based approaches (slow!) [[Berghmans10](#)]
- KVM [[KivityKamay07](#)] with AMD SVM [[RoedelGraf09](#)]
- Early Xen prototype [[He09](#)]
- Blue Pill rootkit hiding from other hypervisors [[Rutkowska06](#)]



## Related work

- First models for nested virtualization [[PopekGoldberg74](#), [BelpaireHsu75](#), [LauerWyeth73](#)]
- First implementation in the IBM z/VM; relies on architectural support for nested virtualization ([sie](#))
- Microkernels meet recursive VMs [[FordHibler96](#)]: assumes we can modify software at all levels
- x86 software based approaches (slow!) [[Berghmans10](#)]
- KVM [[KivityKamay07](#)] with AMD SVM [[RoedelGraf09](#)]
- Early Xen prototype [[He09](#)]
- Blue Pill rootkit hiding from other hypervisors [[Rutkowska06](#)]



## Related work

- First models for nested virtualization [[PopekGoldberg74](#), [BelpaireHsu75](#), [LauerWyeth73](#)]
- First implementation in the IBM z/VM; relies on architectural support for nested virtualization ([sie](#))
- Microkernels meet recursive VMs [[FordHibler96](#)]: assumes we can modify software at all levels
- x86 software based approaches (slow!) [[Berghmans10](#)]
- KVM [[KivityKamay07](#)] with AMD SVM [[RoedelGraf09](#)]
- Early Xen prototype [[He09](#)]
- Blue Pill rootkit hiding from other hypervisors [[Rutkowska06](#)]



- First models for nested virtualization [[PopekGoldberg74](#), [BelpaireHsu75](#), [LauerWyeth73](#)]
- First implementation in the IBM z/VM; relies on architectural support for nested virtualization ([sie](#))
- Microkernels meet recursive VMs [[FordHibler96](#)]: assumes we can modify software at all levels
- x86 software based approaches (slow!) [[Berghmans10](#)]
- KVM [[KivityKamay07](#)] with AMD SVM [[RoedelGraf09](#)]
- Early Xen prototype [[He09](#)]
- Blue Pill rootkit hiding from other hypervisors [[Rutkowska06](#)]



# What is the Turtles project?

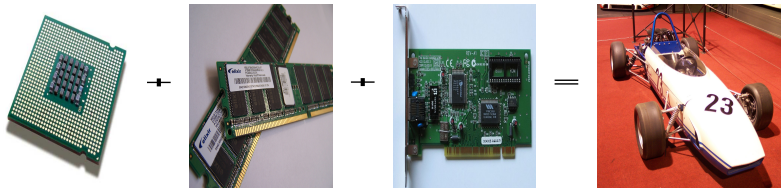


- [Efficient nested virtualization for Intel x86](#) based on KVM
- Multiple guest hypervisors and VMs: VMware, Windows, ...
- Code publicly available



# What is the Turtles project? (cont')

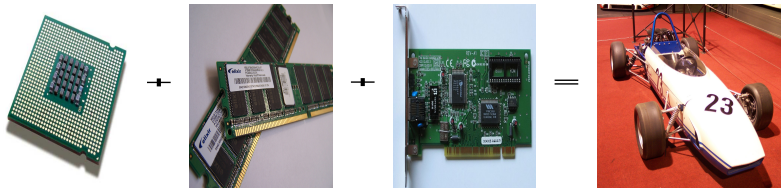
- **Nested VMX virtualization** for nested **CPU** virtualization
- **Multi-dimensional paging** for nested **MMU** virtualization
- **Multi-level device assignment** for nested **I/O** virtualization
- **Micro-optimizations** to make it go **fast** (see paper)





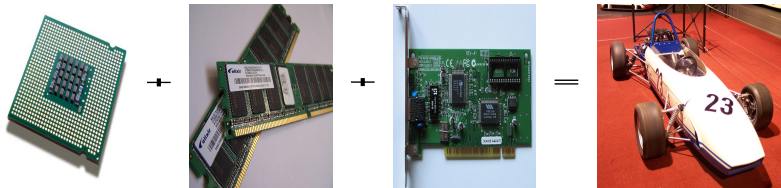
# What is the Turtles project? (cont')

- Nested VMX virtualization for nested CPU virtualization
- Multi-dimensional paging for nested MMU virtualization
- Multi-level device assignment for nested I/O virtualization
- Micro-optimizations to make it go fast (see paper)



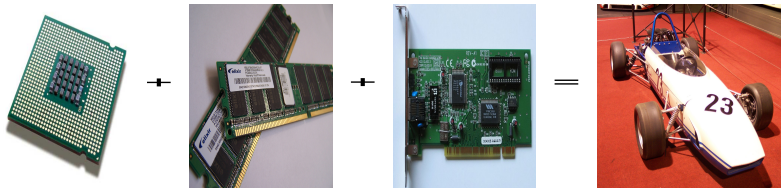
# What is the Turtles project? (cont')

- Nested VMX virtualization for nested CPU virtualization
- Multi-dimensional paging for nested MMU virtualization
- Multi-level device assignment for nested I/O virtualization
- Micro-optimizations to make it go fast (see paper)



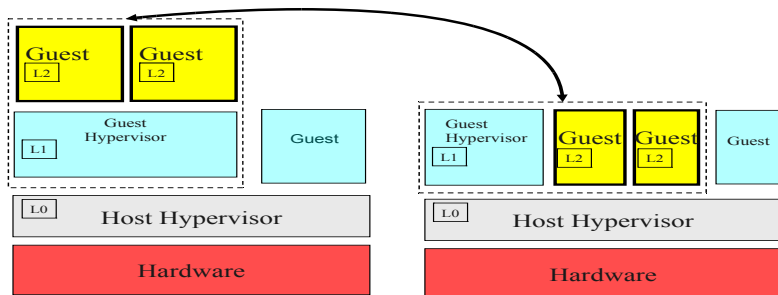
# What is the Turtles project? (cont')

- Nested VMX virtualization for nested CPU virtualization
- Multi-dimensional paging for nested MMU virtualization
- Multi-level device assignment for nested I/O virtualization
- Micro-optimizations to make it go fast (see paper)



# Theory of nested CPU virtualization

- **Single-level** architectural support (x86) vs. **multi-level** architectural support (e.g., z/VM)
- **Single level**  $\Rightarrow$  one hypervisor, many guests
- Turtles approach:  $L_0$  **multiplexes** the hardware between  $L_1$  and  $L_2$ , running both as guests of  $L_0$ —without either being aware of it
- (Scheme generalized for  $n$  levels; Our focus is  $n=2$ )



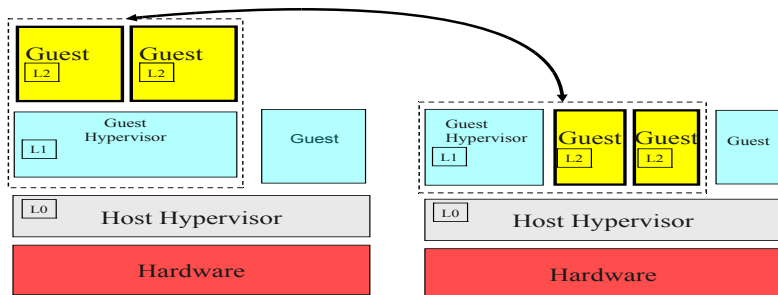
Multiple logical levels

Multiplexed on a single level



# Theory of nested CPU virtualization

- **Single-level** architectural support (x86) vs. **multi-level** architectural support (e.g., z/VM)
- **Single level**  $\Rightarrow$  one hypervisor, many guests
- Turtles approach:  $L_0$  **multiplexes** the hardware between  $L_1$  and  $L_2$ , running both as guests of  $L_0$ —without either being aware of it
- (Scheme generalized for  $n$  levels; Our focus is  $n=2$ )



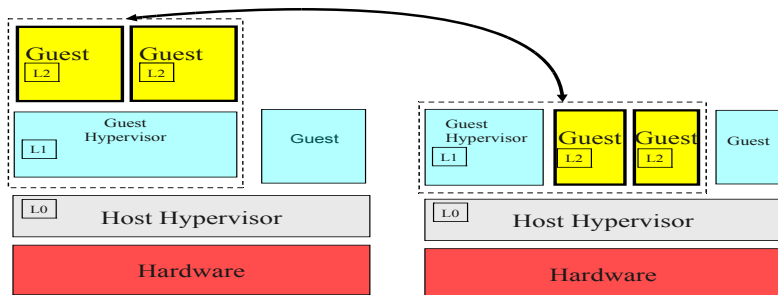
Multiple logical levels

Multiplexed on a single level



# Theory of nested CPU virtualization

- **Single-level** architectural support (x86) vs. **multi-level** architectural support (e.g., z/VM)
- **Single level**  $\Rightarrow$  one hypervisor, many guests
- Turtles approach:  $L_0$  **multiplexes** the hardware between  $L_1$  and  $L_2$ , running both as guests of  $L_0$ —without either being aware of it
- (Scheme generalized for  $n$  levels; Our focus is  $n=2$ )



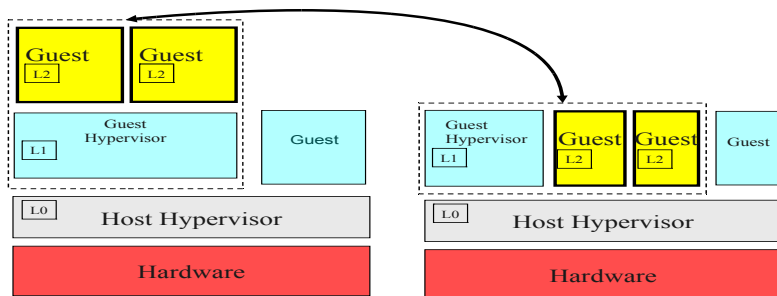
Multiple logical levels

Multiplexed on a single level



# Theory of nested CPU virtualization

- **Single-level** architectural support (x86) vs. **multi-level** architectural support (e.g., z/VM)
- **Single level**  $\Rightarrow$  one hypervisor, many guests
- Turtles approach:  $L_0$  **multiplexes** the hardware between  $L_1$  and  $L_2$ , running both as guests of  $L_0$ —without either being aware of it
- (Scheme generalized for  $n$  levels; Our focus is  $n=2$ )



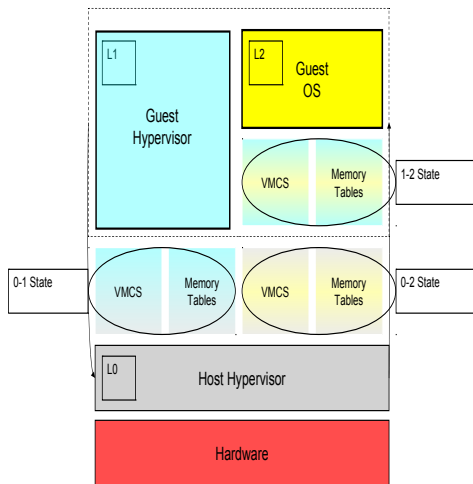
Multiple logical levels

Multiplexed on a single level



# Nested VMX virtualization: flow

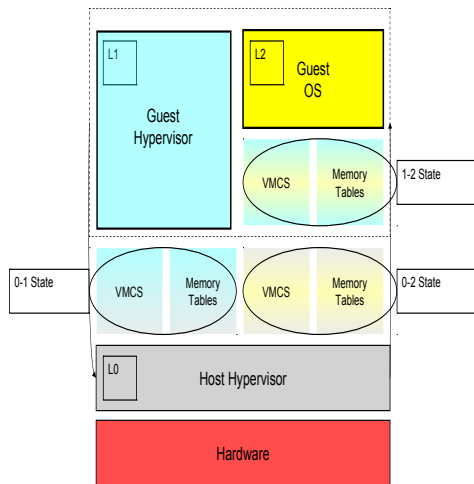
- $L_0$  runs  $L_1$  with  $VMCS_{0 \rightarrow 1}$
- $L_1$  prepares  $VMCS_{1 \rightarrow 2}$  and executes `vmlaunch`
- `vmlaunch` traps to  $L_0$
- $L_0$  merges VMCS's:  $VMCS_{0 \rightarrow 1}$  merged with  $VMCS_{1 \rightarrow 2}$  is  $VMCS_{0 \rightarrow 2}$
- $L_0$  launches  $L_2$
- $L_2$  causes a trap
- $L_0$  handles trap itself or forwards it to  $L_1$
- ...
- eventually,  $L_0$  resumes  $L_2$
- repeat





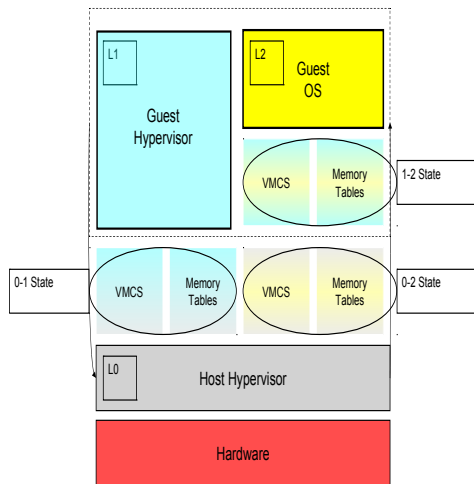
# Nested VMX virtualization: flow

- $L_0$  runs  $L_1$  with  $VMCS_{0 \rightarrow 1}$
- $L_1$  prepares  $VMCS_{1 \rightarrow 2}$  and executes `vmlaunch`
- `vmlaunch` traps to  $L_0$
- $L_0$  merges VMCS's:  $VMCS_{0 \rightarrow 1}$  merged with  $VMCS_{1 \rightarrow 2}$  is  $VMCS_{0 \rightarrow 2}$
- $L_0$  launches  $L_2$
- $L_2$  causes a trap
- $L_0$  handles trap itself or forwards it to  $L_1$
- ...
- eventually,  $L_0$  resumes  $L_2$
- repeat



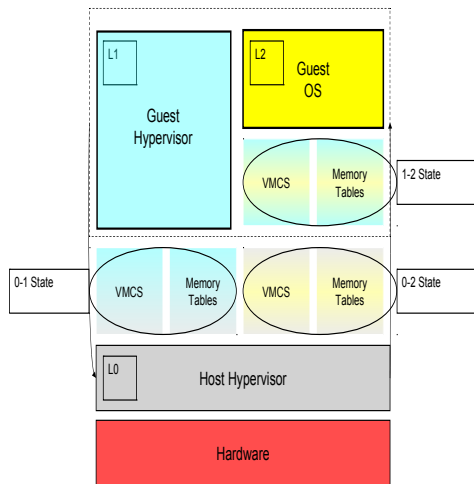
# Nested VMX virtualization: flow

- $L_0$  runs  $L_1$  with  $VMCS_{0 \rightarrow 1}$
- $L_1$  prepares  $VMCS_{1 \rightarrow 2}$  and executes `vmlaunch`
- `vmlaunch` traps to  $L_0$
- $L_0$  merges VMCS's:  $VMCS_{0 \rightarrow 1}$  merged with  $VMCS_{1 \rightarrow 2}$  is  $VMCS_{0 \rightarrow 2}$
- $L_0$  launches  $L_2$
- $L_2$  causes a trap
- $L_0$  handles trap itself or forwards it to  $L_1$
- ...
- eventually,  $L_0$  resumes  $L_2$
- repeat



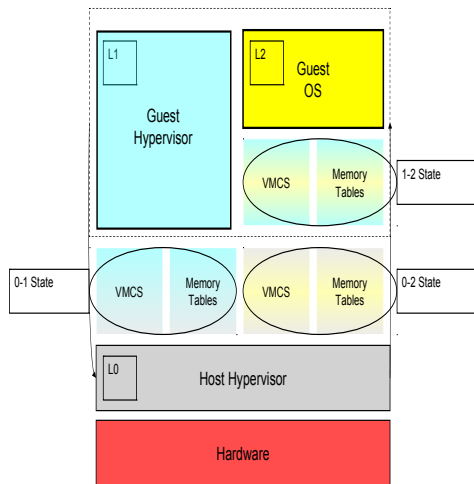
# Nested VMX virtualization: flow

- $L_0$  runs  $L_1$  with  $VMCS_{0 \rightarrow 1}$
- $L_1$  prepares  $VMCS_{1 \rightarrow 2}$  and executes `vmlaunch`
- `vmlaunch` traps to  $L_0$
- $L_0$  merges VMCS's:  $VMCS_{0 \rightarrow 1}$  merged with  $VMCS_{1 \rightarrow 2}$  is  $VMCS_{0 \rightarrow 2}$
- $L_0$  launches  $L_2$
- $L_2$  causes a trap
- $L_0$  handles trap itself or forwards it to  $L_1$
- ...
- eventually,  $L_0$  resumes  $L_2$
- repeat



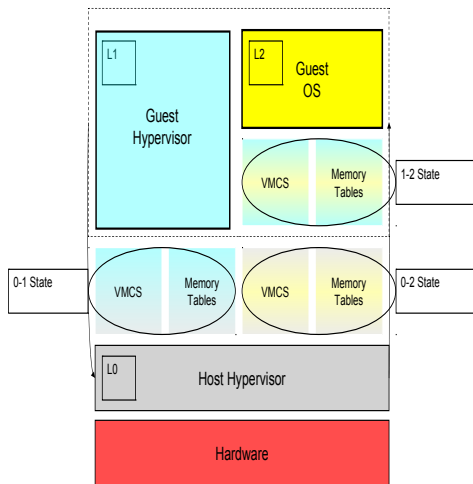
# Nested VMX virtualization: flow

- $L_0$  runs  $L_1$  with  $VMCS_{0 \rightarrow 1}$
- $L_1$  prepares  $VMCS_{1 \rightarrow 2}$  and executes `vmlaunch`
- `vmlaunch` traps to  $L_0$
- $L_0$  merges VMCS's:  $VMCS_{0 \rightarrow 1}$  merged with  $VMCS_{1 \rightarrow 2}$  is  $VMCS_{0 \rightarrow 2}$
- $L_0$  launches  $L_2$
- $L_2$  causes a trap
- $L_0$  handles trap itself or forwards it to  $L_1$
- ...
- eventually,  $L_0$  resumes  $L_2$
- repeat



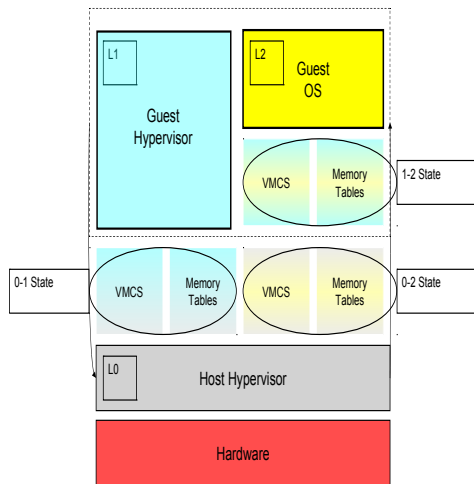
# Nested VMX virtualization: flow

- $L_0$  runs  $L_1$  with  $VMCS_{0 \rightarrow 1}$
- $L_1$  prepares  $VMCS_{1 \rightarrow 2}$  and executes `vmlaunch`
- `vmlaunch` traps to  $L_0$
- $L_0$  merges VMCS's:  $VMCS_{0 \rightarrow 1}$  merged with  $VMCS_{1 \rightarrow 2}$  is  $VMCS_{0 \rightarrow 2}$
- $L_0$  launches  $L_2$
- $L_2$  causes a trap
- $L_0$  handles trap itself or forwards it to  $L_1$
- ...
- eventually,  $L_0$  resumes  $L_2$
- repeat



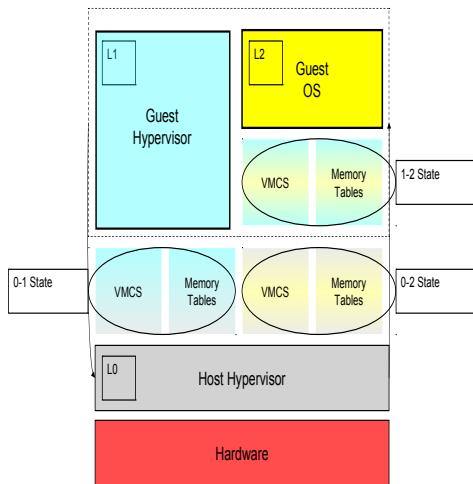
# Nested VMX virtualization: flow

- $L_0$  runs  $L_1$  with  $VMCS_{0 \rightarrow 1}$
- $L_1$  prepares  $VMCS_{1 \rightarrow 2}$  and executes `vmlaunch`
- `vmlaunch` traps to  $L_0$
- $L_0$  merges VMCS's:  $VMCS_{0 \rightarrow 1}$  merged with  $VMCS_{1 \rightarrow 2}$  is  $VMCS_{0 \rightarrow 2}$
- $L_0$  launches  $L_2$
- $L_2$  causes a trap
- $L_0$  handles trap itself or forwards it to  $L_1$
- ...
- eventually,  $L_0$  resumes  $L_2$
- repeat



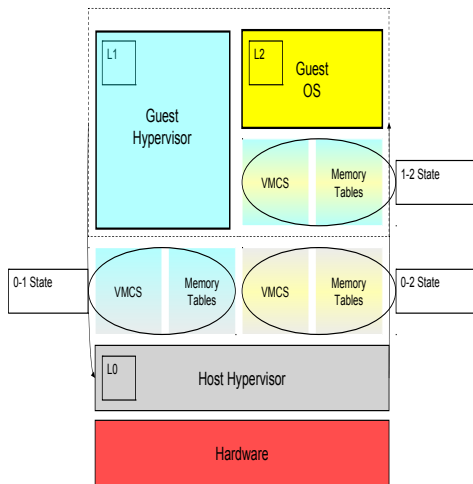
# Nested VMX virtualization: flow

- $L_0$  runs  $L_1$  with  $VMCS_{0 \rightarrow 1}$
- $L_1$  prepares  $VMCS_{1 \rightarrow 2}$  and executes `vmlaunch`
- `vmlaunch` traps to  $L_0$
- $L_0$  merges VMCS's:  $VMCS_{0 \rightarrow 1}$  merged with  $VMCS_{1 \rightarrow 2}$  is  $VMCS_{0 \rightarrow 2}$
- $L_0$  launches  $L_2$
- $L_2$  causes a trap
- $L_0$  handles trap itself or forwards it to  $L_1$
- ...
- eventually,  $L_0$  resumes  $L_2$
- repeat



# Nested VMX virtualization: flow

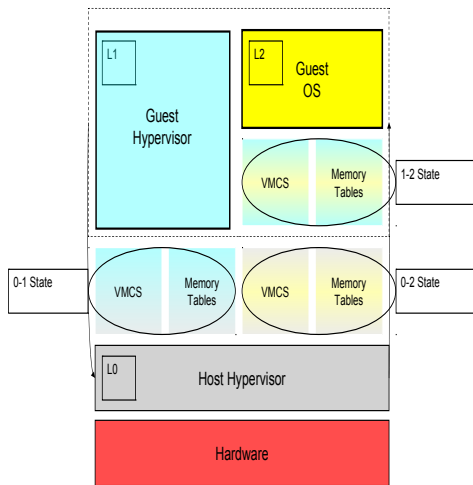
- $L_0$  runs  $L_1$  with  $VMCS_{0 \rightarrow 1}$
- $L_1$  prepares  $VMCS_{1 \rightarrow 2}$  and executes `vmlaunch`
- `vmlaunch` traps to  $L_0$
- $L_0$  merges VMCS's:  $VMCS_{0 \rightarrow 1}$  merged with  $VMCS_{1 \rightarrow 2}$  is  $VMCS_{0 \rightarrow 2}$
- $L_0$  launches  $L_2$
- $L_2$  causes a trap
- $L_0$  handles trap itself or forwards it to  $L_1$
- ...
- eventually,  $L_0$  resumes  $L_2$
- repeat





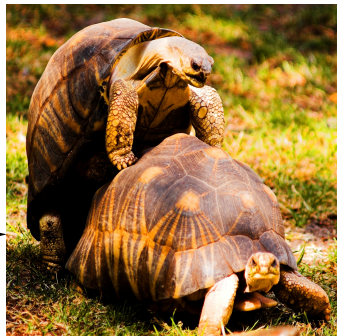
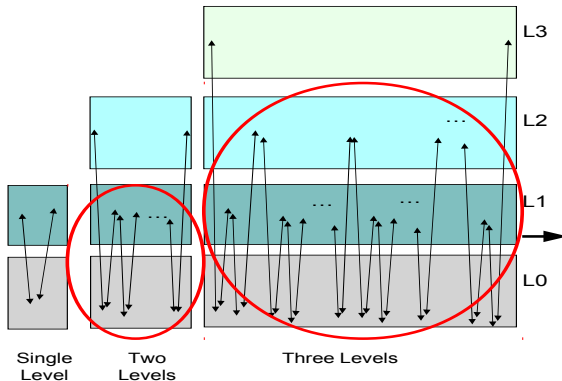
# Nested VMX virtualization: flow

- $L_0$  runs  $L_1$  with  $VMCS_{0 \rightarrow 1}$
- $L_1$  prepares  $VMCS_{1 \rightarrow 2}$  and executes `vmlaunch`
- `vmlaunch` traps to  $L_0$
- $L_0$  merges VMCS's:  $VMCS_{0 \rightarrow 1}$  merged with  $VMCS_{1 \rightarrow 2}$  is  $VMCS_{0 \rightarrow 2}$
- $L_0$  launches  $L_2$
- $L_2$  causes a trap
- $L_0$  handles trap itself or forwards it to  $L_1$
- ...
- eventually,  $L_0$  resumes  $L_2$
- repeat



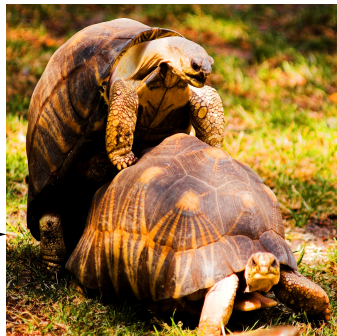
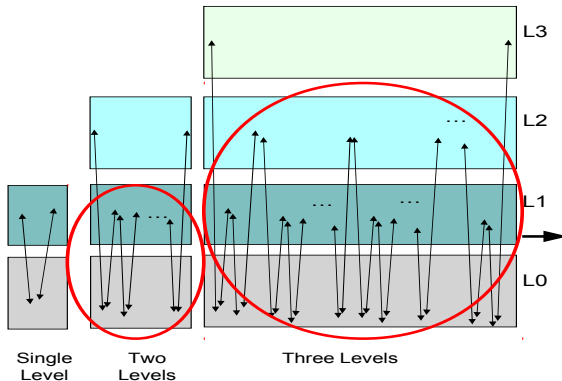
# Exit multiplication makes angry turtle angry

- To handle a single L<sub>2</sub> exit, L<sub>1</sub> does many things: read and write the VMCS, disable interrupts, ...
- Those operations can trap, leading to **exit multiplication**
- **Exit multiplication**: a single L<sub>2</sub> exit can cause 40-50 L<sub>1</sub> exits!
- Optimize: make a **single exit fast** and **reduce frequency of exits**



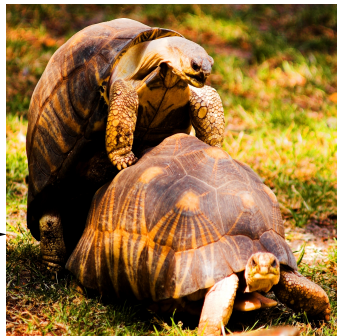
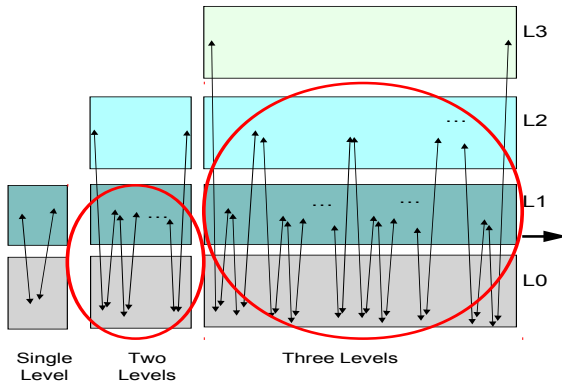
# Exit multiplication makes angry turtle angry

- To handle a single L<sub>2</sub> exit, L<sub>1</sub> does many things: read and write the VMCS, disable interrupts, ...
- Those operations can trap, leading to **exit multiplication**
- **Exit multiplication**: a single L<sub>2</sub> exit can cause 40-50 L<sub>1</sub> exits!
- Optimize: make a **single exit fast** and **reduce frequency of exits**



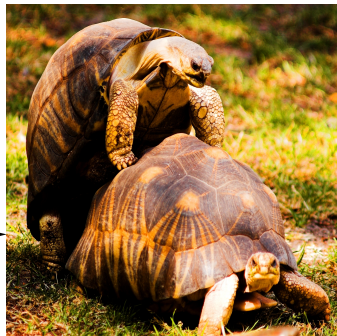
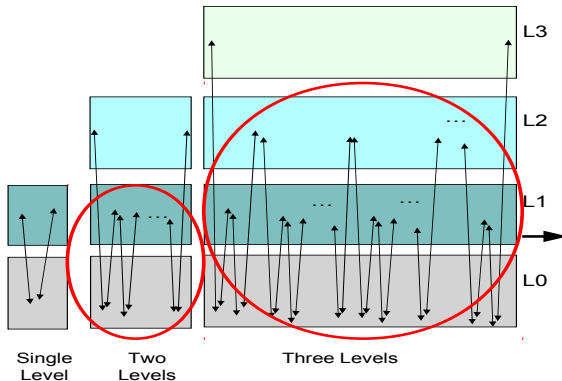
# Exit multiplication makes angry turtle angry

- To handle a single L<sub>2</sub> exit, L<sub>1</sub> does many things: read and write the VMCS, disable interrupts, ...
- Those operations can trap, leading to **exit multiplication**
- **Exit multiplication**: a single L<sub>2</sub> exit can cause 40-50 L<sub>1</sub> exits!
- Optimize: make a single exit fast and reduce frequency of exits



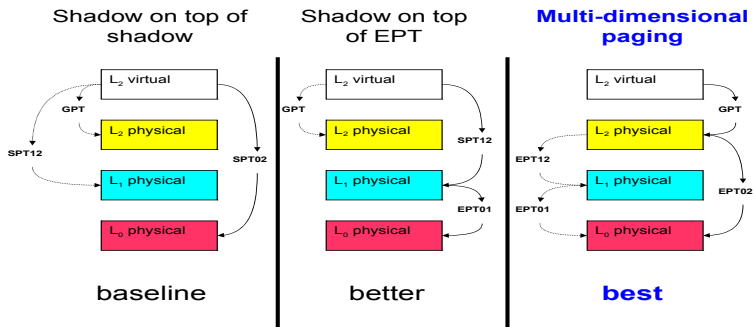
# Exit multiplication makes angry turtle angry

- To handle a single L<sub>2</sub> exit, L<sub>1</sub> does many things: read and write the VMCS, disable interrupts, ...
- Those operations can trap, leading to **exit multiplication**
- **Exit multiplication**: a single L<sub>2</sub> exit can cause 40-50 L<sub>1</sub> exits!
- Optimize: make **a single exit fast** and **reduce frequency of exits**

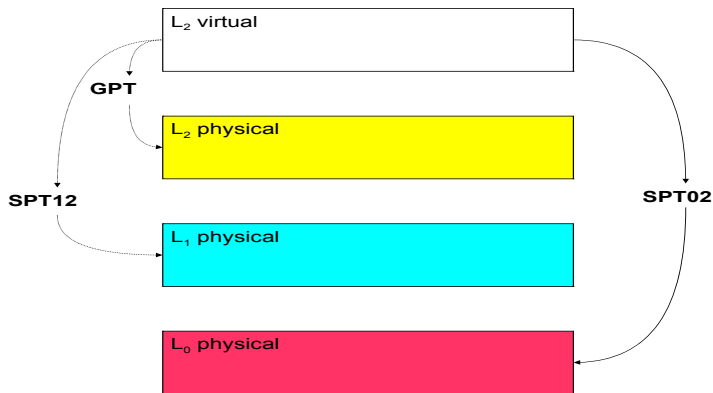


# MMU virtualization via multi-dimensional paging

- **Three logical translations:**  $L_2 \text{ virt} \rightarrow \text{phys}$ ,  $L_2 \rightarrow L_1$ ,  $L_1 \rightarrow L_0$
- Only **two** tables in hardware with EPT:  
virt  $\rightarrow$  phys and guest physical  $\rightarrow$  host physical
- $L_0$  **compresses** three logical translations onto two hardware tables



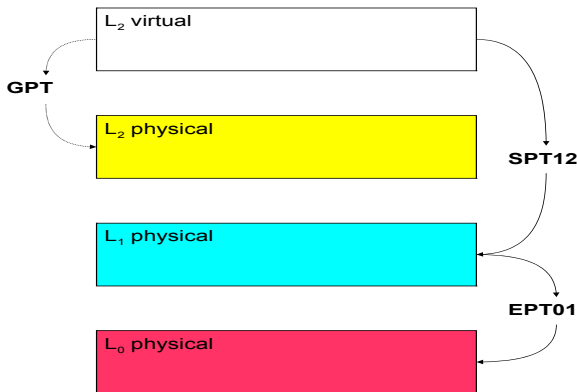
# Baseline: shadow-on-shadow



- Assume no EPT table; all hypervisors use shadow paging
- Useful for old machines and as a baseline
- Maintaining shadow page tables is expensive
- **Compress**: three logical translations  $\Rightarrow$  one table in hardware



# Better: shadow-on-EPT

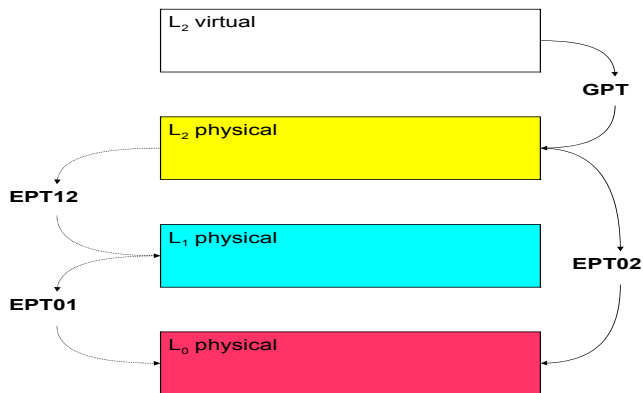


- Instead of one hardware table we have two
- **Compress**: three logical translations  $\Rightarrow$  two in hardware
- Simple approach: L<sub>0</sub> uses EPT, L<sub>1</sub> uses shadow paging for L<sub>2</sub>
- Every L<sub>2</sub> page fault leads to multiple L<sub>1</sub> exits





# Best: multi-dimensional paging

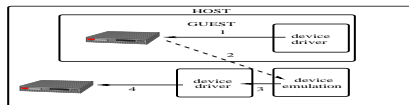


- EPT table rarely changes; guest page table changes a lot
- Again, **compress three** logical translations  $\Rightarrow$  **two** in hardware
- L<sub>0</sub> **emulates** EPT for L<sub>1</sub>
- L<sub>0</sub> uses EPT<sub>0 $\rightarrow$ 1</sub> and EPT<sub>1 $\rightarrow$ 2</sub> to construct EPT<sub>0 $\rightarrow$ 2</sub>
- End result: a lot less exits!



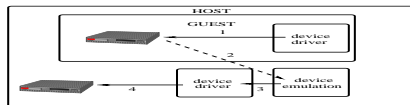
# Introduction to I/O virtualization

- Device emulation [Sugerman01]

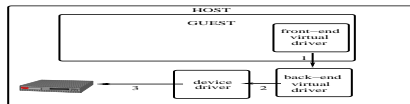


# Introduction to I/O virtualization

- Device emulation [Sugerman01]

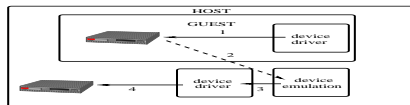


- Para-virtualized drivers [Barham03, Russell08]

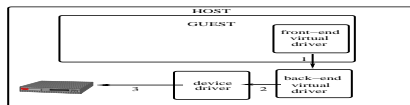


# Introduction to I/O virtualization

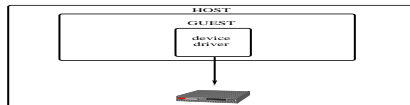
- Device emulation [Sugerman01]



- Para-virtualized drivers [Barham03, Russell08]

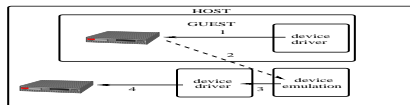


- Direct device assignment [Levasseur04, Yassour08]

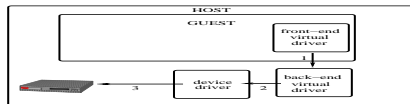


# Introduction to I/O virtualization

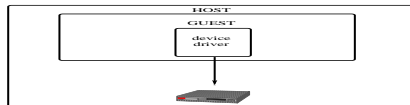
- Device emulation [Sugerman01]



- Para-virtualized drivers [Barham03, Russell08]



- Direct device assignment [Levasseur04, Yassour08]

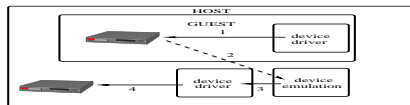


- Direct assignment **best performing option**

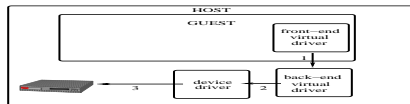


# Introduction to I/O virtualization

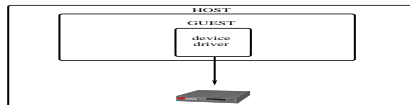
- Device emulation [Sugerman01]



- Para-virtualized drivers [Barham03, Russell08]



- Direct device assignment [Levasseur04, Yassour08]

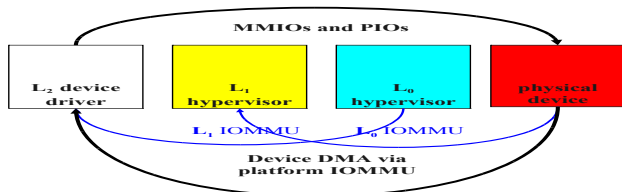


- Direct assignment **best performing option**
- Direct assignment **requires IOMMU** for safe DMA bypass



# Multi-level device assignment

- With nested **3x3** options for I/O virtualization ( $L_2 \Leftrightarrow L_1 \Leftrightarrow L_0$ )
- Multi-level device assignment means giving an  $L_2$  guest direct access to  $L_0$ 's devices, safely bypassing both  $L_0$  and  $L_1$

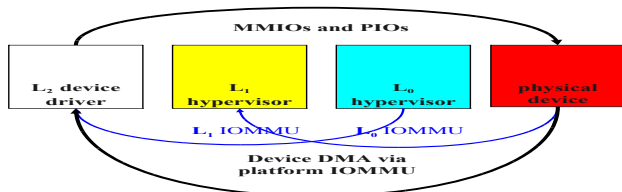


- How?  $L_0$  emulates an IOMMU for  $L_1$  [Amit10]
- $L_0$  compresses multiple IOMMU translations onto the single hardware IOMMU page table
- $L_2$  programs the device directly
- Device DMA's into  $L_2$  memory space directly



# Multi-level device assignment

- With nested **3x3** options for I/O virtualization ( $L_2 \Leftrightarrow L_1 \Leftrightarrow L_0$ )
- **Multi-level device assignment** means giving an  $L_2$  guest direct access to  $L_0$ 's devices, safely **bypassing both  $L_0$  and  $L_1$**



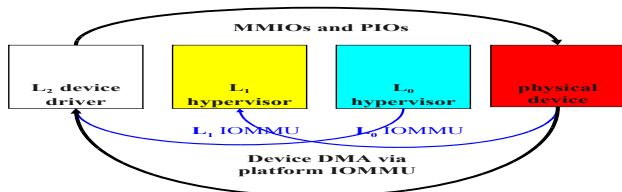
- How?  $L_0$  emulates an IOMMU for  $L_1$  [Amit10]
- $L_0$  **compresses** multiple IOMMU translations onto the single hardware IOMMU page table
- $L_2$  programs the device directly
- Device DMA's into  $L_2$  memory space directly





# Multi-level device assignment

- With nested **3x3** options for I/O virtualization ( $L_2 \Leftrightarrow L_1 \Leftrightarrow L_0$ )
- **Multi-level device assignment** means giving an  $L_2$  guest direct access to  $L_0$ 's devices, safely **bypassing both  $L_0$  and  $L_1$**

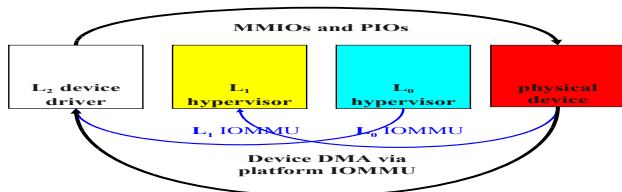


- How?  $L_0$  emulates an IOMMU for  $L_1$  [Amit10]
- $L_0$  **compresses** multiple IOMMU translations onto the single hardware IOMMU page table
- $L_2$  programs the device directly
- Device DMA's into  $L_2$  memory space directly



# Multi-level device assignment

- With nested 3x3 options for I/O virtualization ( $L_2 \Leftrightarrow L_1 \Leftrightarrow L_0$ )
- Multi-level device assignment means giving an  $L_2$  guest direct access to  $L_0$ 's devices, safely bypassing both  $L_0$  and  $L_1$

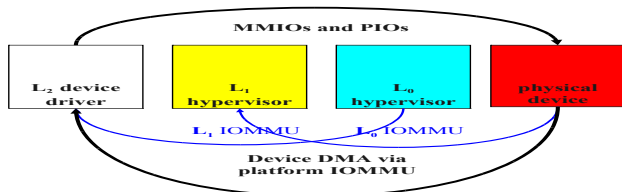


- How?  $L_0$  emulates an IOMMU for  $L_1$  [Amit10]
- $L_0$  compresses multiple IOMMU translations onto the single hardware IOMMU page table
- $L_2$  programs the device directly
- Device DMA's into  $L_2$  memory space directly



# Multi-level device assignment

- With nested **3x3** options for I/O virtualization ( $L_2 \Leftrightarrow L_1 \Leftrightarrow L_0$ )
- **Multi-level device assignment** means giving an  $L_2$  guest direct access to  $L_0$ 's devices, safely **bypassing both  $L_0$  and  $L_1$**

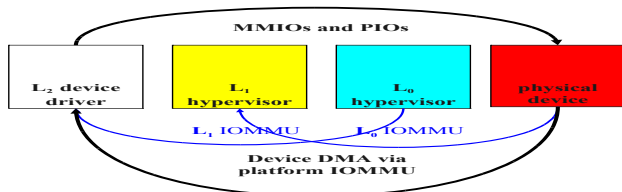


- How?  $L_0$  emulates an IOMMU for  $L_1$  [Amit10]
- $L_0$  **compresses** multiple IOMMU translations onto the single hardware IOMMU page table
- $L_2$  programs the device directly
- Device DMA's into  $L_2$  memory space directly



# Multi-level device assignment

- With nested **3x3** options for I/O virtualization ( $L_2 \Leftrightarrow L_1 \Leftrightarrow L_0$ )
- **Multi-level device assignment** means giving an  $L_2$  guest direct access to  $L_0$ 's devices, safely **bypassing both  $L_0$  and  $L_1$**



- How?  $L_0$  emulates an IOMMU for  $L_1$  [Amit10]
- $L_0$  **compresses** multiple IOMMU translations onto the single hardware IOMMU page table
- $L_2$  programs the device directly
- Device DMA's into  $L_2$  memory space directly



# Experimental Setup

- Running Linux, [Windows](#), KVM, [VMware](#), SMP, ...
- Macro workloads:
  - kernbench
  - SPECjbb
  - netperf
- Multi-dimensional paging?
- Multi-level device assignment?
- KVM as L<sub>1</sub> vs. VMware as L<sub>1</sub>?
- See paper for full experimental details, more benchmarks and analysis, including worst case synthetic micro-benchmark



# Macro: SPECjbb and kernbench

kernbench				
	Host	Guest	Nested	Nested <sub>DRW</sub>
Run time	324.3	355	406.3	391.5
% overhead vs. host	-	9.5	25.3	20.7
% overhead vs. guest	-	-	14.5	10.3

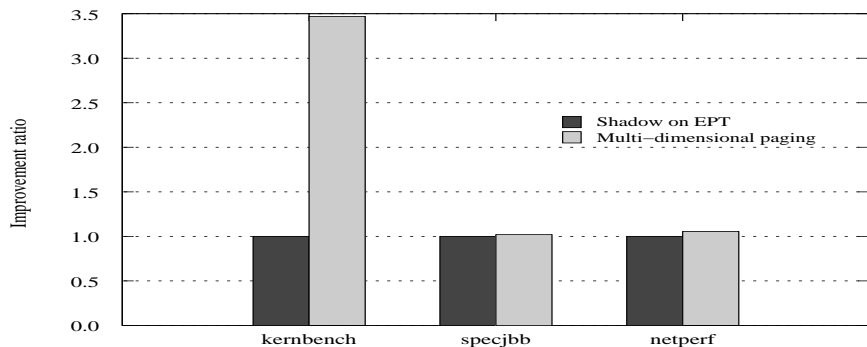
SPECjbb				
	Host	Guest	Nested	Nested <sub>DRW</sub>
Score	90493	83599	77065	78347
% degradation vs. host	-	7.6	14.8	13.4
% degradation vs. guest	-	-	7.8	6.3

Table: kernbench and SPECjbb results

- Exit multiplication effect not as bad as we feared
- Direct `vmread` and `vmwrite` (DRW) give an immediate boost
- Take-away: each level of virtualization adds approximately the same overhead!



# Macro: multi-dimensional paging

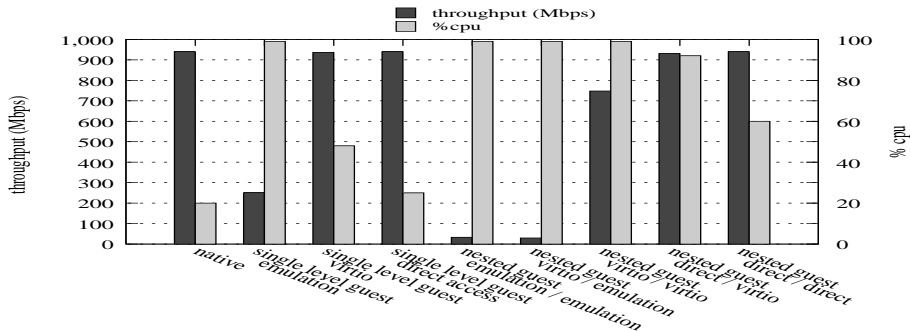


- Impact of multi-dimensional paging depends on rate of page faults
- Shadow-on-EPT: every  $L_2$  page fault causes  $L_1$  multiple exits
- Multi-dimensional paging: only EPT violations cause  $L_1$  exits
- EPT table rarely changes:  $\#(\text{EPT violations}) \ll \#(\text{page faults})$
- Multi-dimensional paging huge win for page-fault intensive

kernbench



# Macro: multi-level device assignment

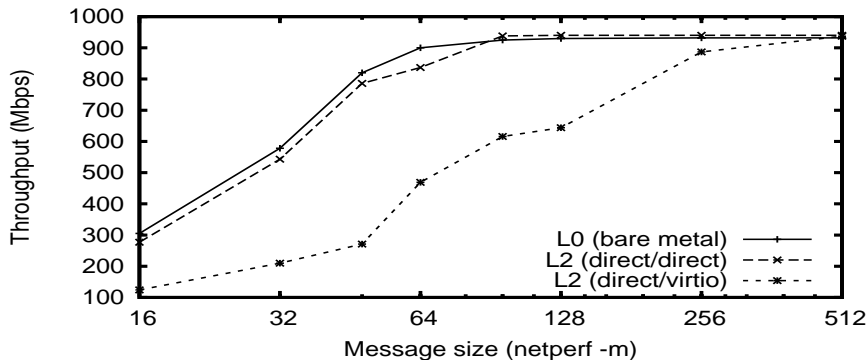


- Benchmark: netperf TCP\_STREAM (transmit)
- Multi-level device assignment best performing option
- But: native at 20%, multi-level device assignment at 60% (x3!)
- Interrupts considered harmful, cause exit multiplication





# Macro: multi-level device assignment (sans interrupts)



- What if we could deliver device interrupts directly to L<sub>2</sub>?
- Only 7% difference between native and nested guest!



# Conclusions

- Efficient nested x86 virtualization is challenging but feasible
- A whole new ballpark opening up many exciting applications—security, cloud, architecture, . . .
- Current overhead of 6-14%
  - Negligible for some workloads, not yet for others
  - Work in progress—expect at most 5% eventually
- Code is available
- It's turtles all the way down



# Conclusions

- Efficient nested x86 virtualization is challenging but feasible
- A whole new ballpark opening up many exciting applications—security, cloud, architecture, . . .
- Current overhead of 6-14%
  - Negligible for some workloads, not yet for others
  - Work in progress—expect at most 5% eventually
- Code is available
- It's turtles all the way down



# Conclusions

- Efficient nested x86 virtualization is challenging but feasible
- A whole new ballpark opening up many exciting applications—security, cloud, architecture, . . .
- Current overhead of 6-14%
  - Negligible for some workloads, not yet for others
  - Work in progress—expect at most 5% eventually
- Code is available
- It's turtles all the way down



# Conclusions

- Efficient nested x86 virtualization is challenging but feasible
- A whole new ballpark opening up many exciting applications—security, cloud, architecture, . . .
- Current overhead of 6-14%
  - Negligible for some workloads, not yet for others
  - Work in progress—expect at most 5% eventually
- Code is available
- It's turtles all the way down



# Conclusions

- Efficient nested x86 virtualization is challenging but feasible
- A whole new ballpark opening up many exciting applications—security, cloud, architecture, . . .
- Current overhead of 6-14%
  - Negligible for some workloads, not yet for others
  - Work in progress—expect at most 5% eventually
- Code is available
- It's turtles all the way down



# Questions?

