

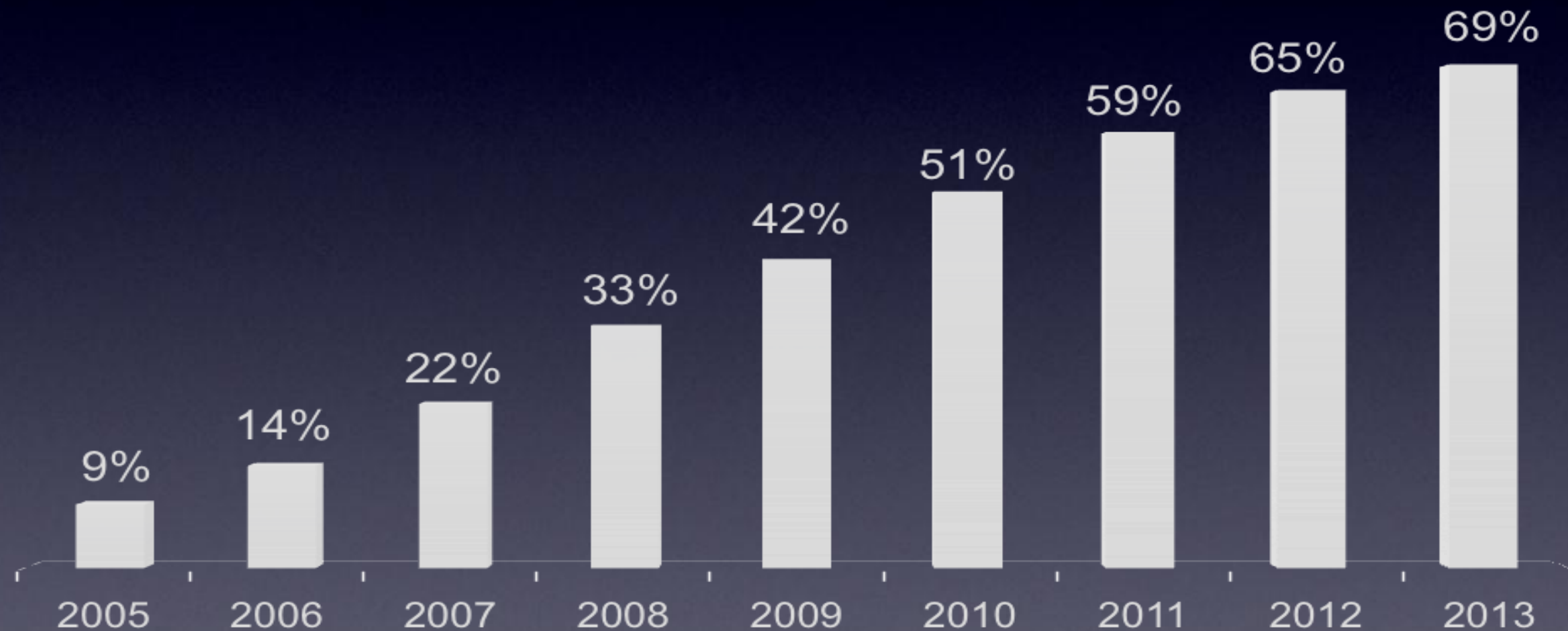
Software Techniques for Avoiding Hardware Virtualization Exits

Ole Agesen Jim Mattson Radu Rugina Jeffrey Sheldon

VMware

Server Virtualization

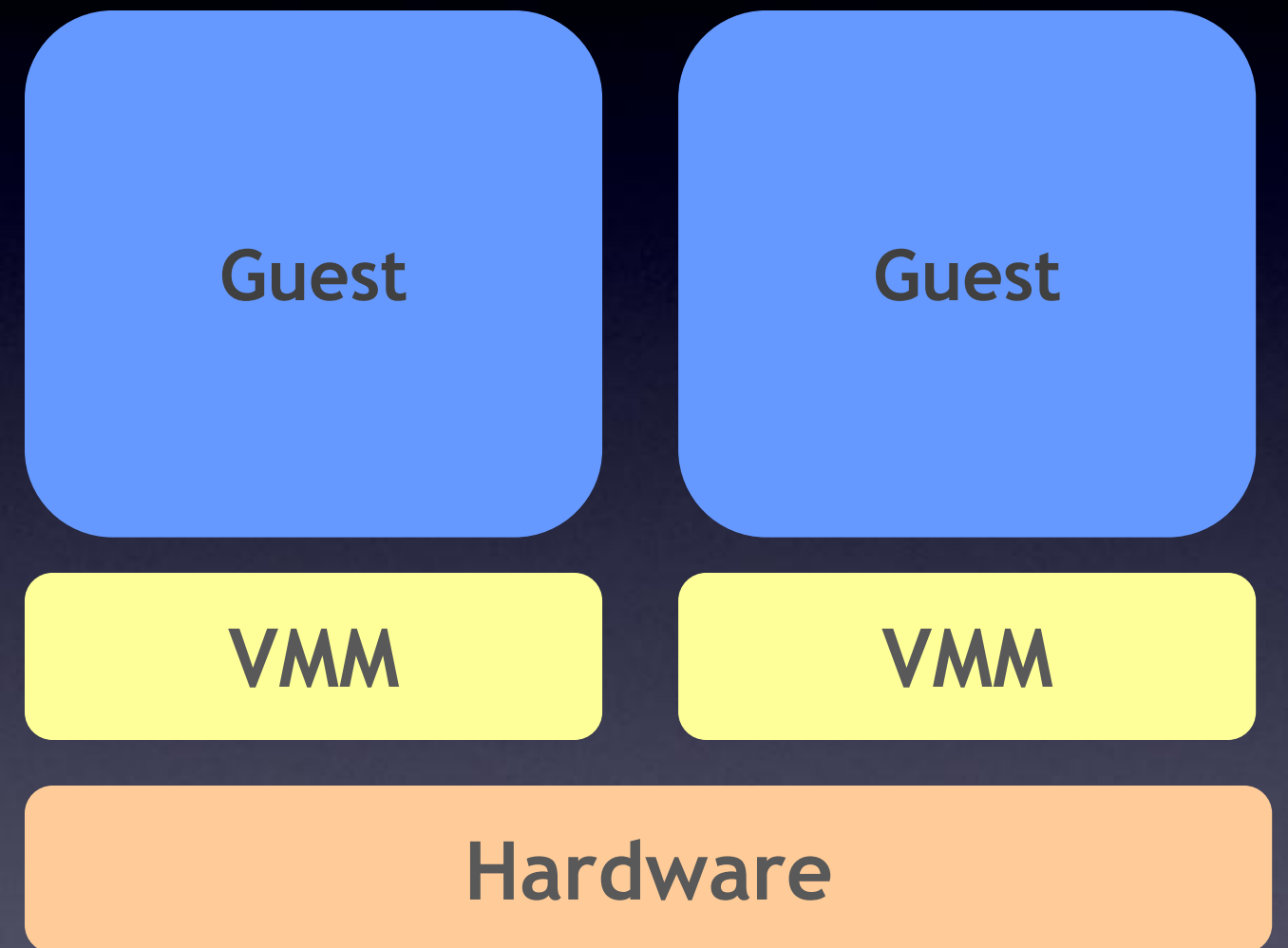
% virtualized workloads world wide



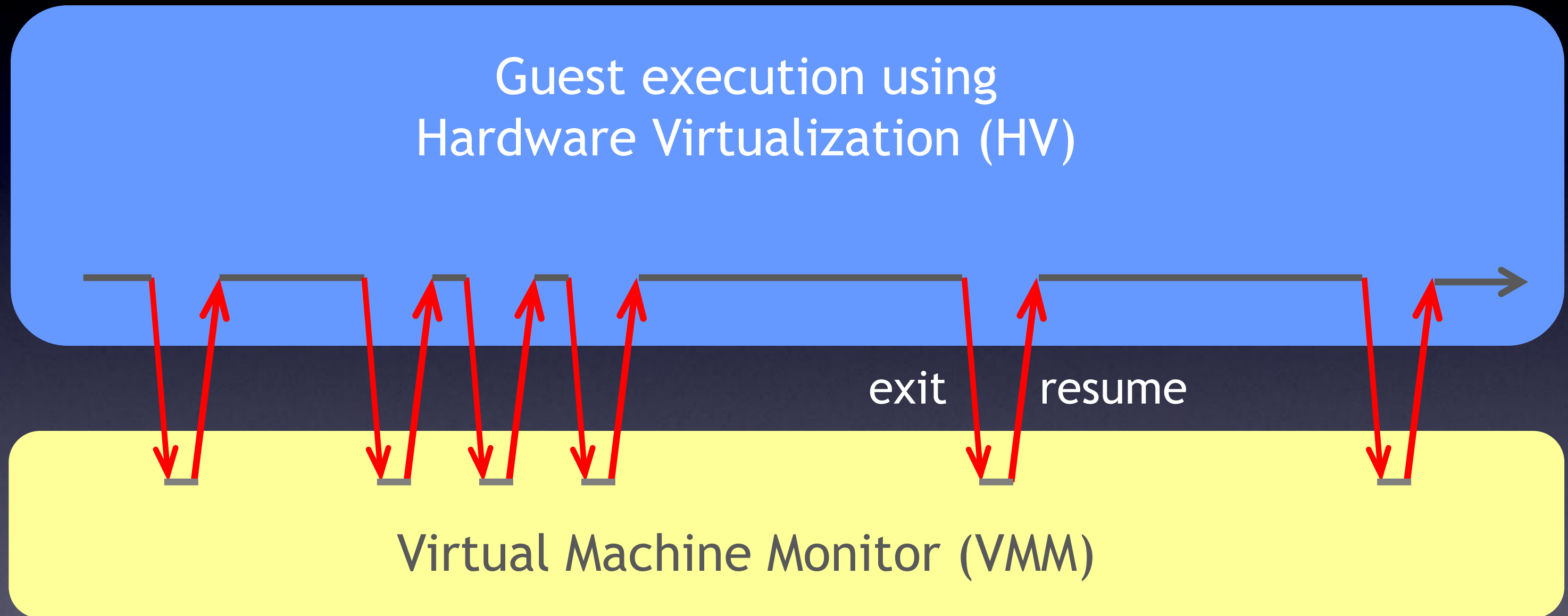
Source: IDC Worldwide Virtualization Tracker

x86 Virtualization

- Virtual Machine Monitor (VMM) abstracts physical hardware
- Two approaches:
 - Binary Translation (BT)
 - Hardware-Assisted (HV)
(Intel VT-x and AMD-V)



Hardware-Assisted x86 Virtualization



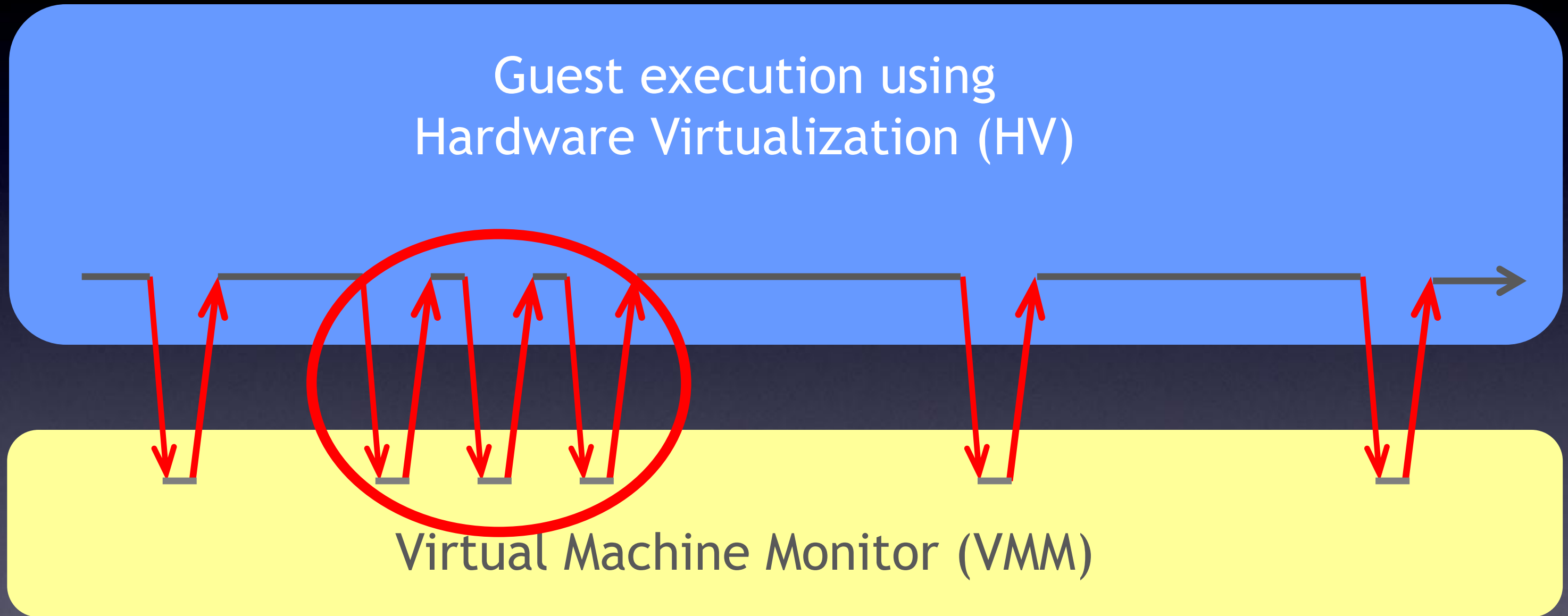
Virtualization Exits Are Expensive

Microarchitecture	Launch Date	Hardware Exit + Resume (cycles)
Prescott	2005	3963
Merom	2006	1579
Penryn	2008	1266
Nehalem	2009	1009
Westmere	2010	761
Sandy Bridge	2011	784

This Talk: Cluster HV Exits

Guest execution using
Hardware Virtualization (HV)

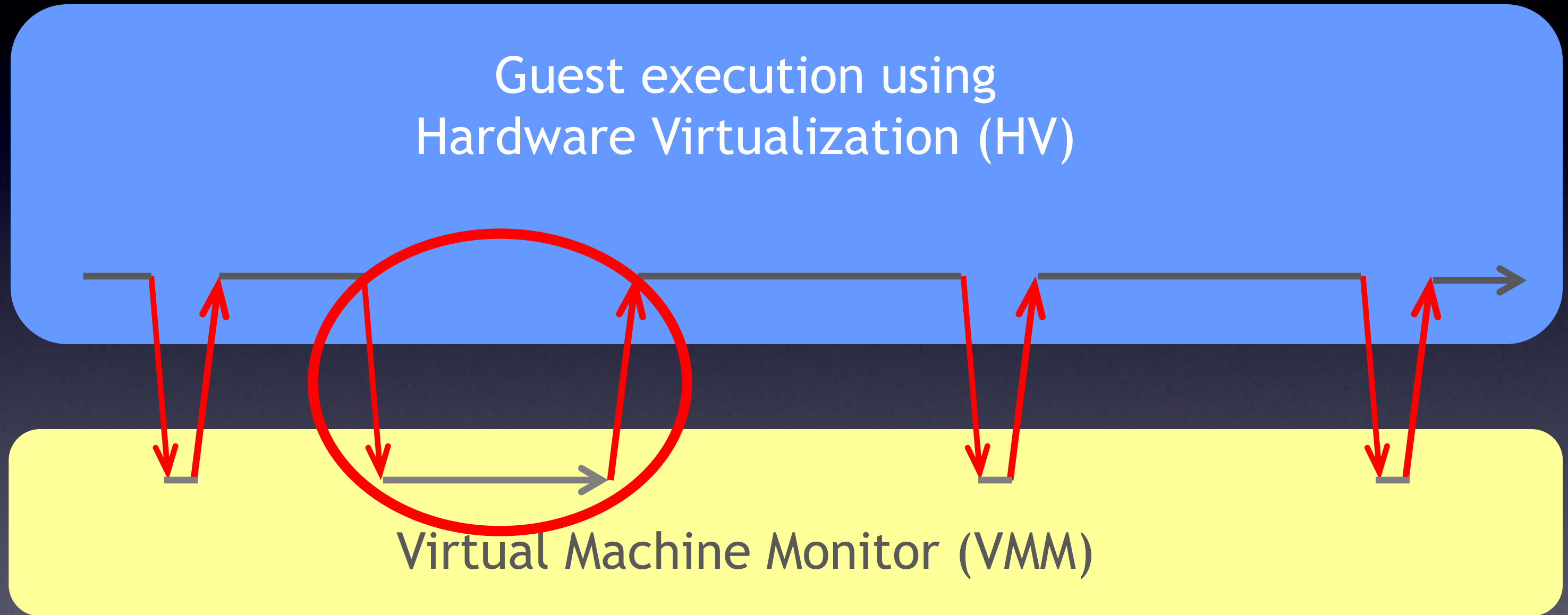
Virtual Machine Monitor (VMM)



This Talk: Cluster HV Exits

Guest execution using
Hardware Virtualization (HV)

Virtual Machine Monitor (VMM)



Outline

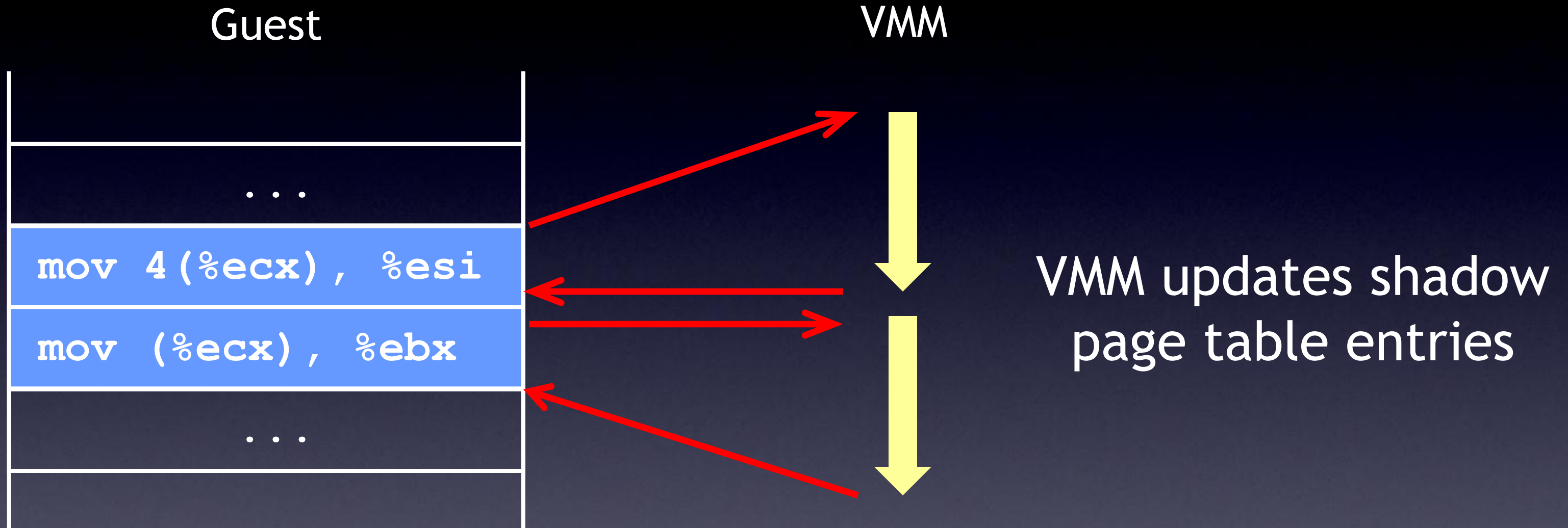
- Exit Pairs
- Exit Clusters
- Nested Virtual Machines
- Results
- Conclusions

Exit Pairs

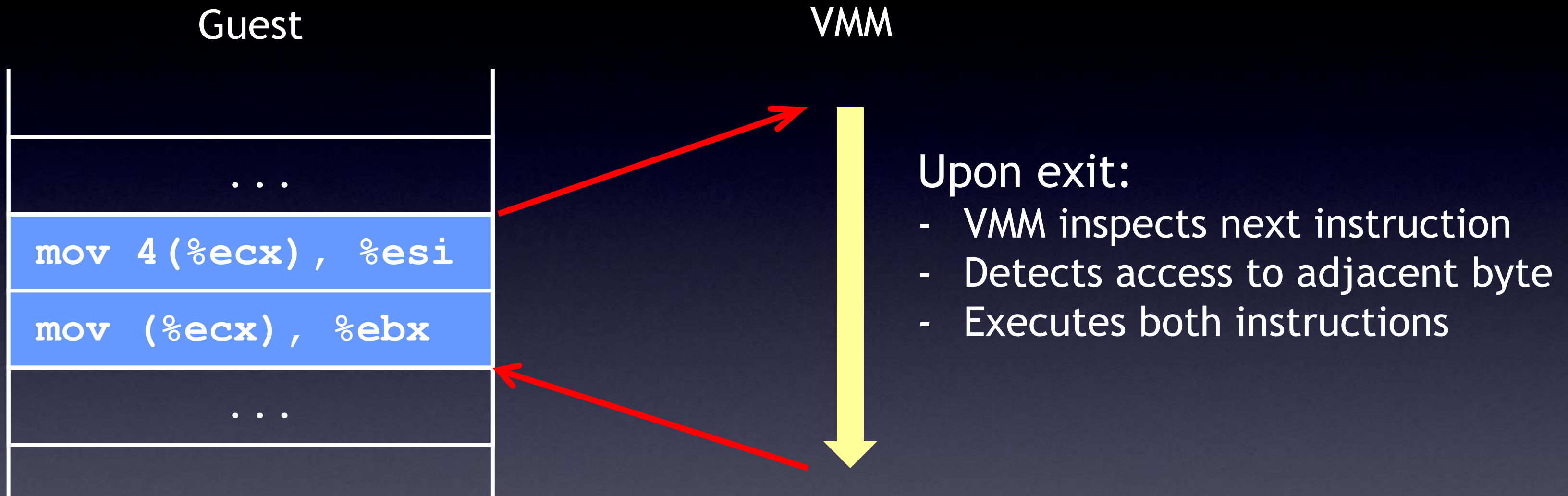


- 32-bit Guest OS using Physical Address Extension (PAE)
- Shadow paging
- Page table entry updates use two 32-bit writes
- Each write causes an exit

Exit Pairs



Exit Pairs



Exit Pairs



VMM

Upon exit:

- VMM inspects next instruction
- Detects access to adjacent byte
- Executes both instructions
- Optimizes execution

Exit Clusters

- **Upon exit:**
 - Scan a few (up to 16) instructions downstream
 - Identify exiting instructions
 - Form a cluster of instructions, executed all at once
- **Challenges:**
 - Cluster formation
 - Efficient execution

Cluster Formation



Cluster Formation

Guest

<code>in %al, %dx</code>
<code>out 0x80, %al</code>
<code>mov %al, %cl</code>
<code>mov %dl, 0xc0</code>
<code>out %al, %dx</code>
<code>out 0x80, %al</code>
<code>xchg %ah, %al</code>
<code>xor %cl, %cl</code>

VMM



HV exit occurs
Decode instructions

Cluster Formation

Guest

<code>in %al, %dx</code>
<code>out 0x80, %al</code>
<code>mov %al, %cl</code>
<code>mov %dl, 0xc0</code>
<code>out %al, %dx</code>
<code>out 0x80, %al</code>
<code>xchg %ah, %al</code>
<code>xor %cl, %cl</code>

VMM



HV exit occurs
Decode instructions
Identify exiting instructions

Cluster Formation

Guest

<code>in %a1, %dx</code>
<code>out 0x80, %a1</code>
<code>mov %a1, %c1</code>
<code>mov %d1, 0xc0</code>
<code>out %a1, %dx</code>
<code>out 0x80, %a1</code>



VMM

HV exit occurs
Decode instructions
Identify exiting instructions
Form cluster

Cluster Formation

Guest

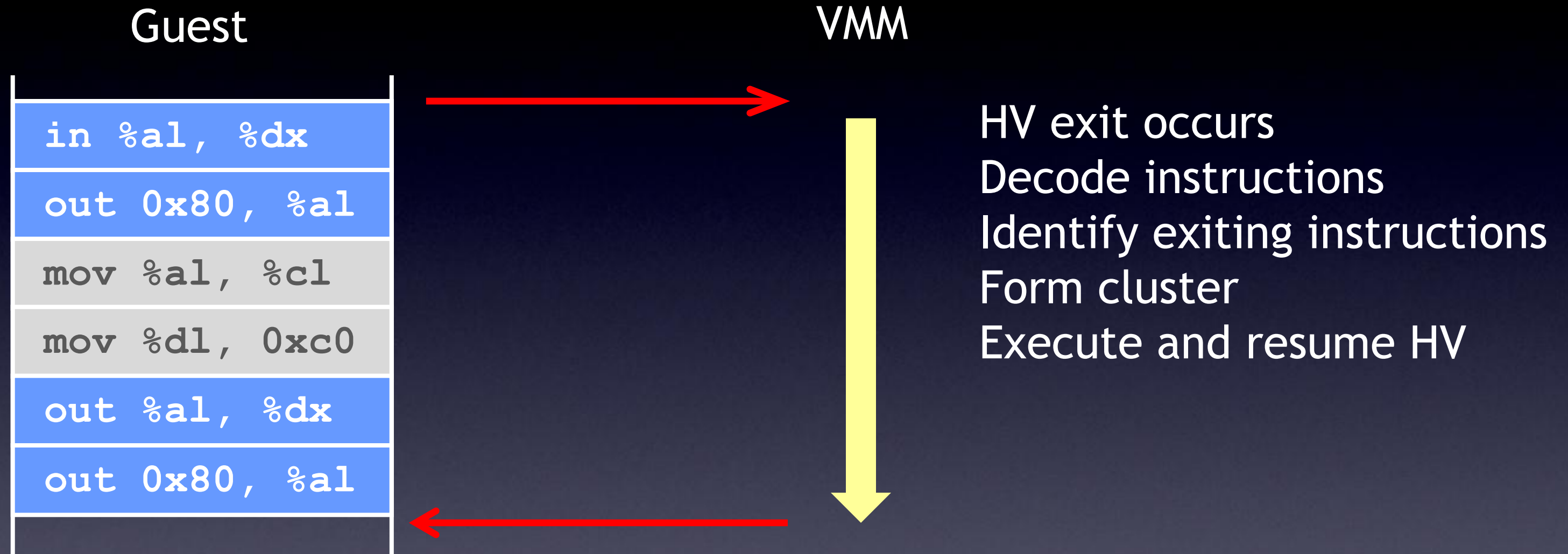
<code>in %a1, %dx</code>
<code>out 0x80, %a1</code>
<code>mov %a1, %c1</code>
<code>mov %d1, 0xc0</code>
<code>out %a1, %dx</code>
<code>out 0x80, %a1</code>

} gap fillers

VMM

HV exit occurs
Decode instructions
Identify exiting instructions
Form cluster

Cluster Formation



Efficient Exit Handling

- **Cluster translation:**
 - Generate a translation for the cluster
 - Insert it in a Translation Cache (TC)
 - Specialize translations (e.g., on addressing mode)
- **Translation reuse:**
 - Compile once
 - Reuse for all subsequent exits

Cluster Translation

Guest

<code>in %a1, %dx</code>
<code>out 0x80, %a1</code>
<code>mov %a1, %c1</code>
<code>mov %d1, 0xc0</code>
<code>out %a1, %dx</code>
<code>out 0x80, %a1</code>

Coherency
Checks

Code
Translation



Dynamic Cluster Formation

Guest

<code>out %al, %dx</code>
<code>mov %cx, -0x12(%bp)</code>
<code>sub %si, %cx</code>
<code>mov %dl, 0x3c5</code>
<code>mov %al, 0x1</code>
<code>out %al, %dx</code>
<code>shr %bl</code>
<code>sbb %ah, %ah</code>
<code>shl %al</code>
<code>mov es:(%di), %ah</code>

Dynamic Cluster Formation

Guest

<code>out %al, %dx</code>
<code>mov %cx, -0x12(%bp)</code>
<code>sub %si, %cx</code>
<code>mov %d1, 0x3c5</code>
<code>mov %a1, 0x1</code>
<code>out %al, %dx</code>
<code>shr %bl</code>
<code>sbb %ah, %ah</code>
<code>shl %a1</code>
<code>mov es:(%di), %ah</code>

Do memory accesses
cause exits?

Dynamic Cluster Formation

Guest

<code>out %al, %dx</code>
<code>mov %cx, -0x12(%bp)</code>
<code>sub %si, %cx</code>
<code>mov %d1, 0x3c5</code>
<code>mov %a1, 0x1</code>
<code>out %a1, %dx</code>
<code>shr %b1</code>
<code>sbb %ah, %ah</code>
<code>shl %a1</code>
<code>mov es:(%di), %ah</code>

Do memory accesses
cause exits?

Sometimes.

Instruction Classification

- **Strongly exiting:**
 - Always cause exits
 - Examples: `in`, `out`, `cpuid`
- **Weakly exiting :**
 - Dynamic exiting behavior
 - Example: memory accesses
 - Runtime VMM support for detecting such cases

Dynamic Cluster Formation

Guest

3	out %al, %dx
0	mov %cx, -0x12(%bp)
0	sub %si, %cx
0	mov %d1, 0x3c5
0	mov %a1, 0x1
2	out %a1, %dx
0	shr %b1
0	sbb %ah, %ah
0	shl %a1
2	mov es:(%di), %ah

Translation postponed
Count exits
Translate on 3rd exit

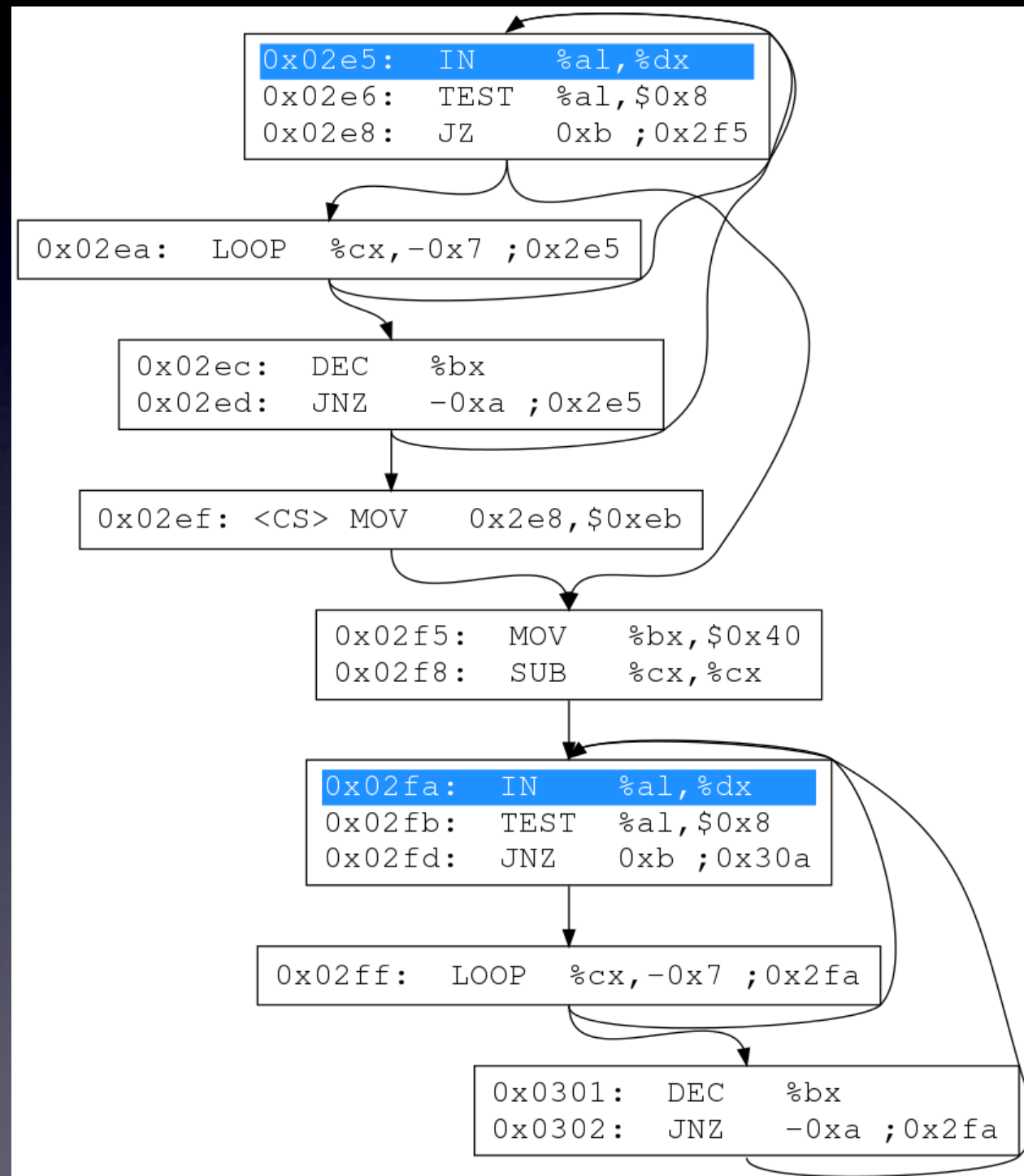
Dynamic Cluster Formation

Guest

3	<code>out %al, %dx</code>
0	<code>mov %cx, -0x12(%bp)</code>
0	<code>sub %si, %cx</code>
0	<code>mov %d1, 0x3c5</code>
0	<code>mov %al, 0x1</code>
2	<code>out %al, %dx</code>
0	<code>shr %bl</code>
0	<code>sbb %ah, %ah</code>
0	<code>shl %al</code>
2	<code>mov es:(%di), %ah</code>

Translation postponed
Count exits
Translate on 3rd exit

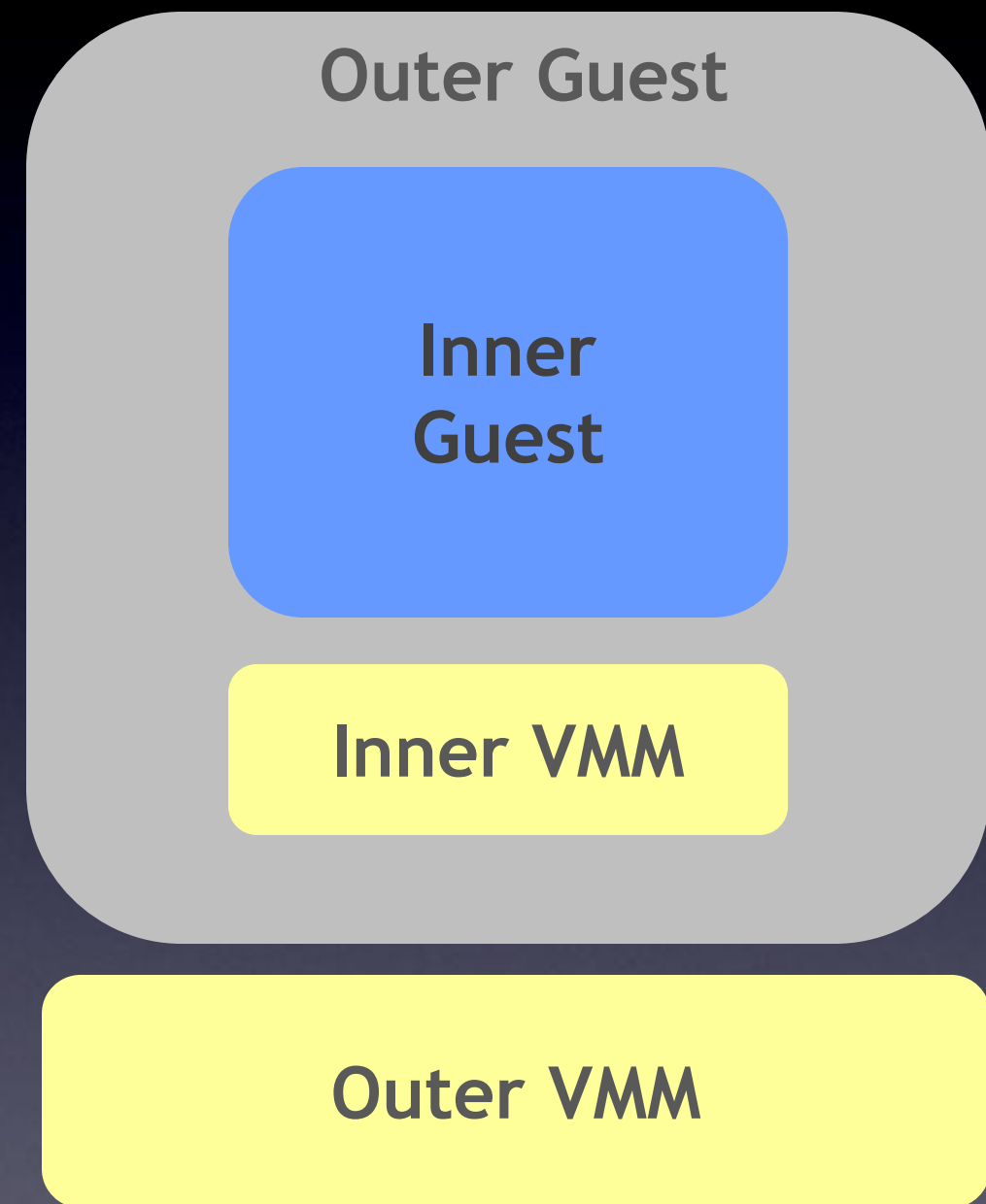
Cluster With Complex Control-Flow



Nested Virtual Machines

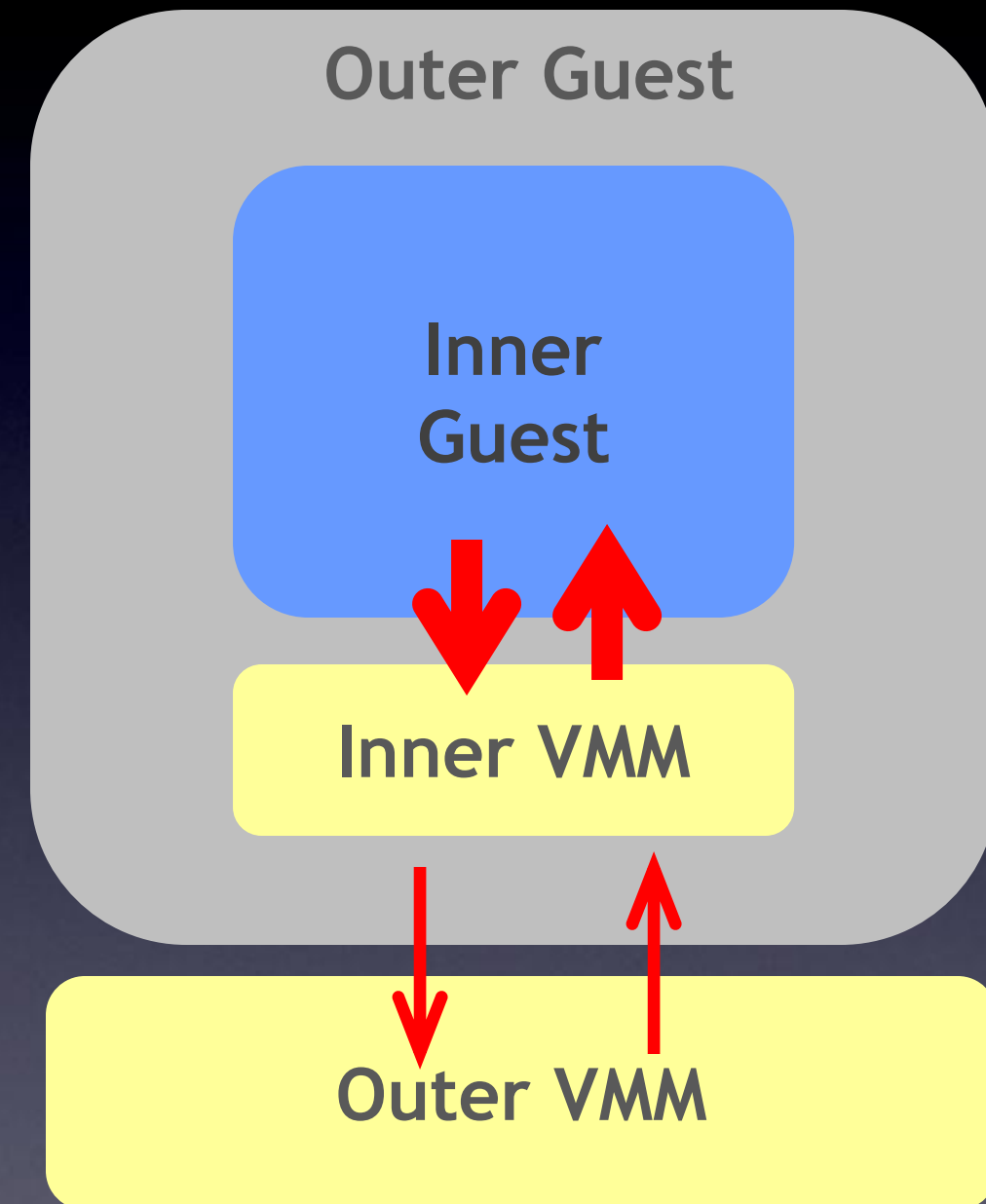
Nested VMs

- Run hypervisor as a guest [Ben-Yehuda et al., OSDI'10]
- Simulate large-scale virtualized environments with fewer hosts
- Training, testing, debugging
- “Windows XP mode” in Win 7



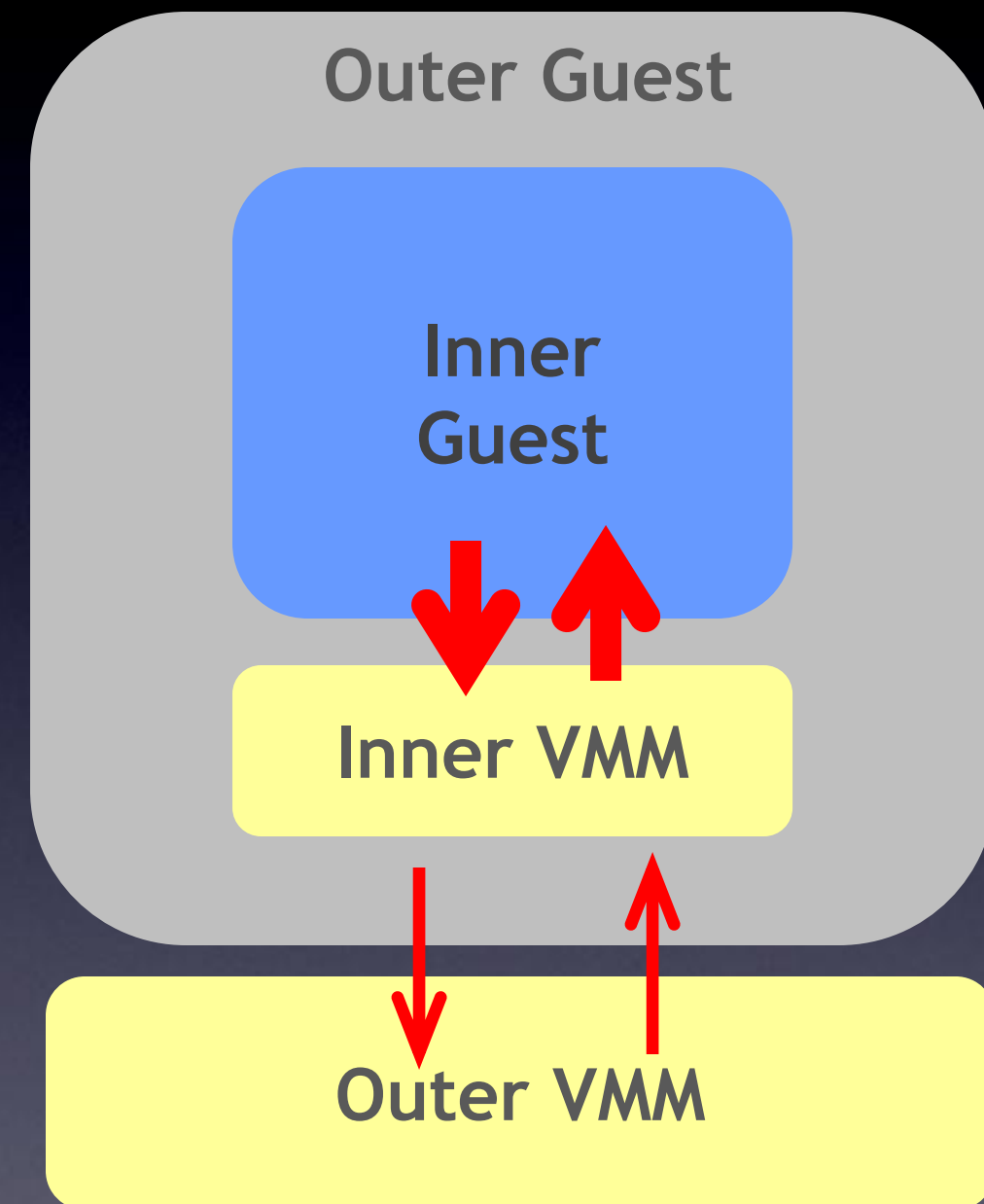
HV for Nested VMs

- Outer VMM transitions
 - Handled in hardware
- Inner VMM transitions
 - Virtual hardware exits are emulated in software



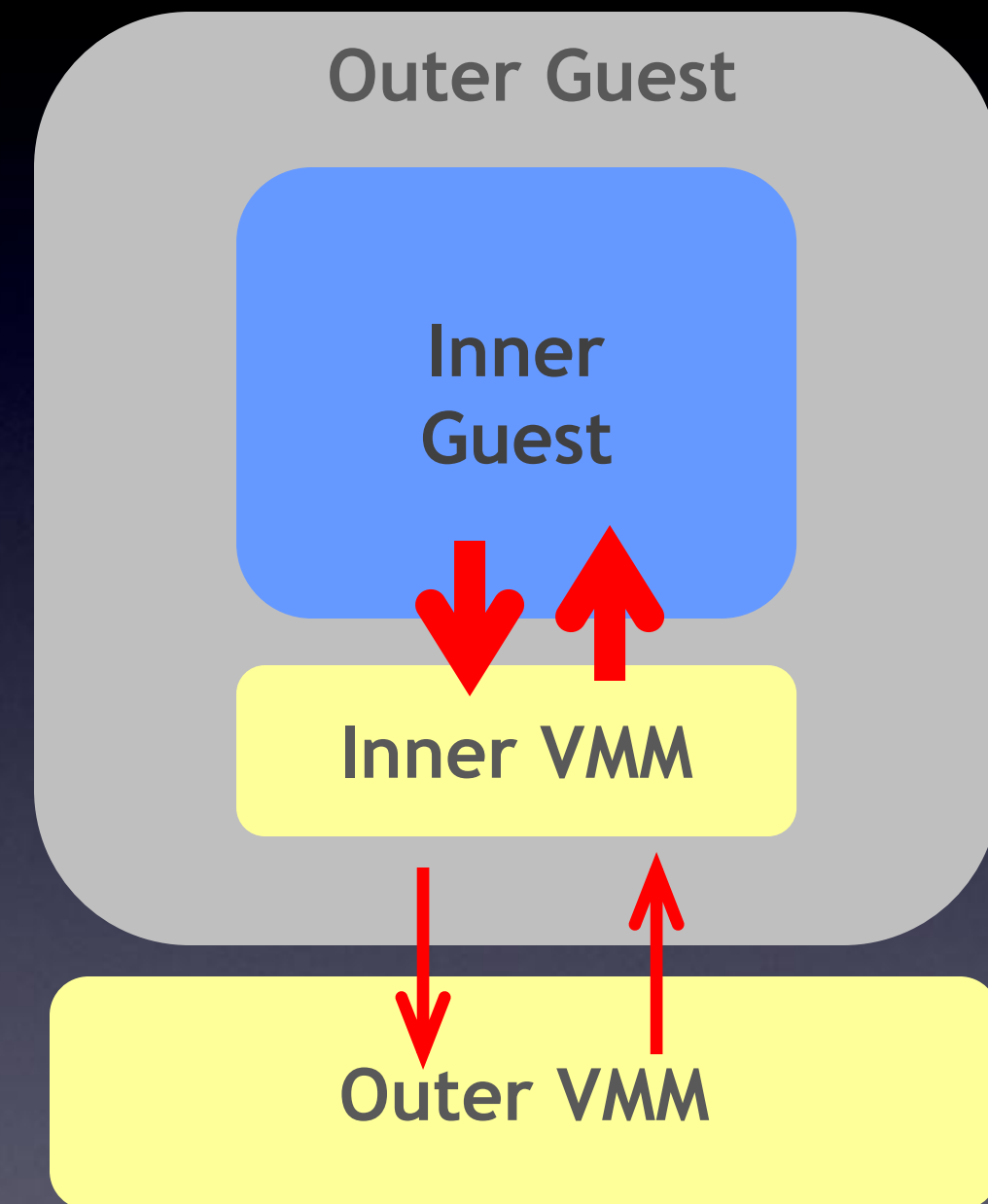
HV for Nested VMs

- Outer VMM transitions
 - Handled in hardware
- Inner VMM transitions
 - Virtual hardware exits are emulated in software
 - Virtual software exit path has lots of exiting instructions (e.g., `vmread`, `vmwrite`)



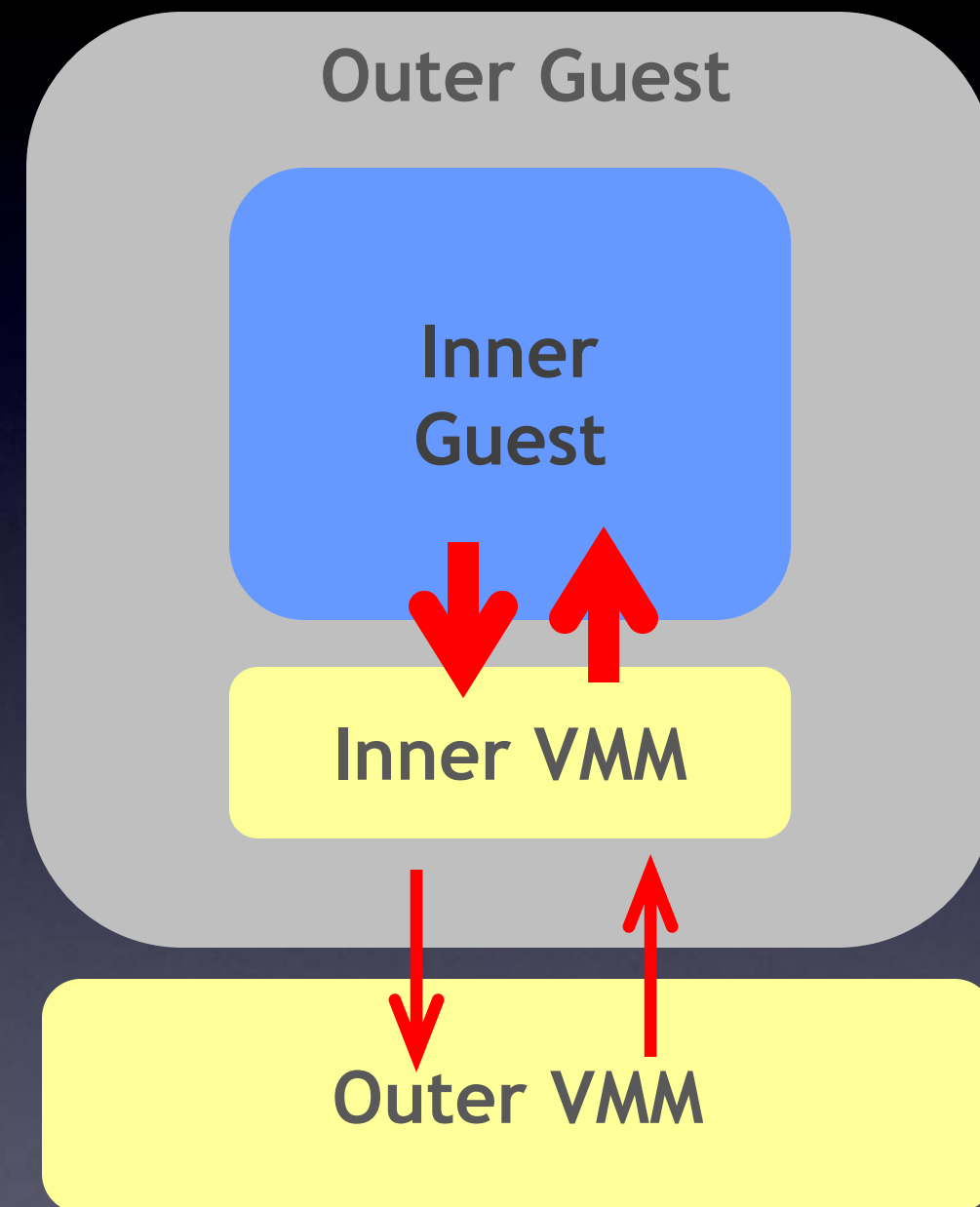
HV for Nested VMs

- Outer VMM transitions
 - Handled in hardware
- Inner VMM transitions
 - Virtual hardware exits are emulated in software
 - Virtual software exit path has lots of exiting instructions (e.g., `vmread`, `vmwrite`)
 - 10x slowdown!



Exit Avoidance

- **Cluster inner exits**
 - Big impact due to high cost of inner exits
- **Cluster outer exits:**
 - Big impact due to the high frequency of outer exits
 - Opportunity: cluster-friendly VMM instruction scheduling



Experimental Results

Implementation

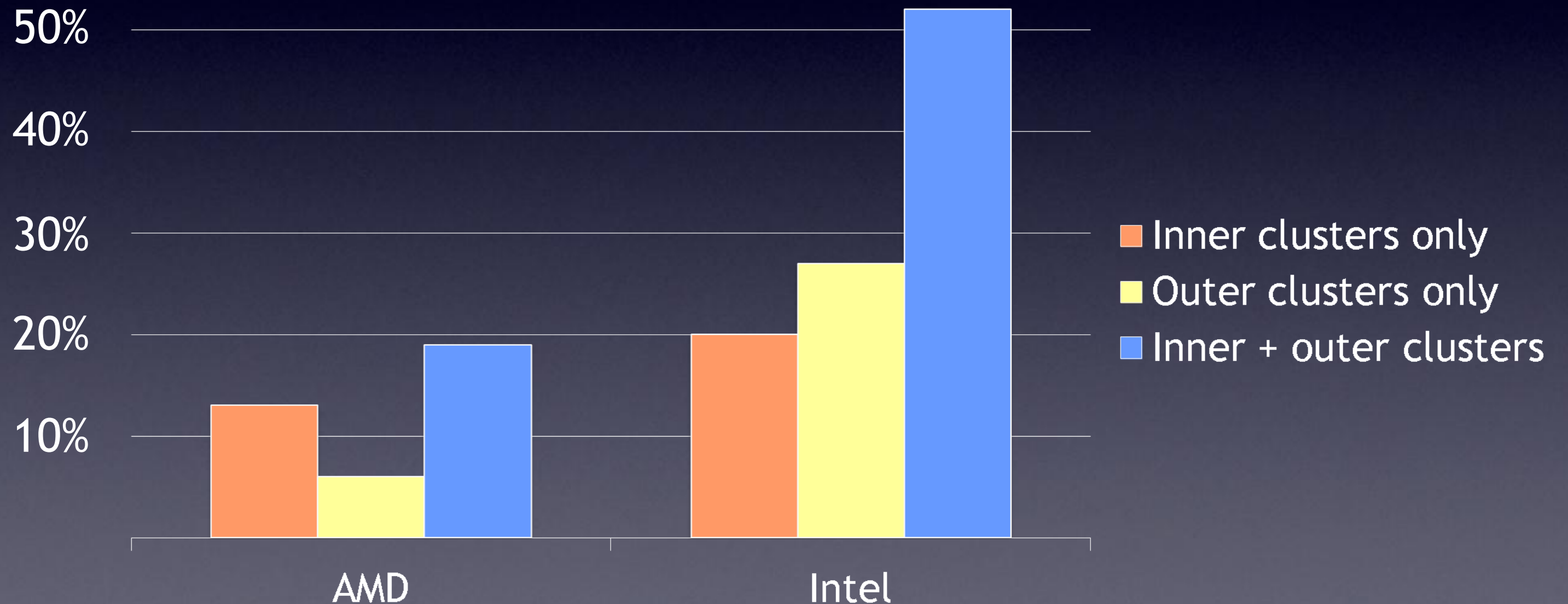
- Fully implemented in VMware products
 - Workstation, Fusion, ESX
- Evolved over many years from exit pairs to complex clusters
- Validated by use in the field

Non-Nested VMs

- **PassMark** (2D graphics benchmark)
 - Clustering improves score by 50% - 80%
- **VMmark** (virtualization benchmark)
 - Consistent exit rate reduction
 - No measurable runtime improvement
- Netperf against a VM with **virtual e1000 NIC**
 - 24% reduction of exits per packet roundtrip

Nested VM Speedup

Kernel-compile workload in the inner VM



Future Directions

- Support more complex clusters
 - Non-contiguous clusters
 - Exits in the middle of loops
- Optimize memory accesses in clusters
 - Cache and reuse work for accesses on the same page
- Applications to virtualized low-latency workloads

Conclusions

- HV exits are expensive, can sap performance
- Software exit clustering complements hardware optimizations for HV exits
- Clustering is key to enabling reasonable performance for nested VMs