

Cloud Operating Systems

Virtual Memory and Structural Setup

Fabian Rauscher, Daniel Gruss, Andreas Kogler

2023-03-13

Setup

Free pages 994 F9 MemInfo F10 Locks F11 Stacktrace F12 Threads

This is on term 0, you should see me now
Kernel end address is 0xffffffff80165000
Now enabling Interrupts...

SWEB-Pseudo-Shell starting...

SWEB: />

- Upstream: <https://github.com/IAIK/sweb>

- Upstream: <https://github.com/IAIK/sweb>
- SWEB Tutorials: <https://www.iaik.tugraz.at/teaching/materials/os/tutorials/>

- Upstream: <https://github.com/IAIK/sweb>
- SWEB Tutorials: <https://www.iaik.tugraz.at/teaching/materials/os/tutorials/>
- Useful links in the Discord `cloudos-announcements` channel

- Upstream: <https://github.com/IAIK/sweb>
- SWEB Tutorials: <https://www.iaik.tugraz.at/teaching/materials/os/tutorials/>
- Useful links in the Discord `cloudos-announcements` channel
- We recommend you work on Linux





Building and running SWEB



Building and running SWEB

- `mkdir -p /tmp/sweb`



Building and running SWEB

- `mkdir -p /tmp/swab`
- `cd /tmp/swab`



Building and running SWEB

- `mkdir -p /tmp/swab`
- `cd /tmp/swab`
- `cmake /path/to/sourcecode/of/swab`



Building and running SWEB

- `mkdir -p /tmp/swab`
- `cd /tmp/swab`
- `cmake /path/to/sourcecode/of/swab`
- `make`



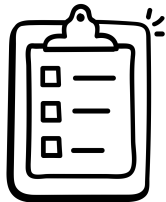
Building and running SWEB

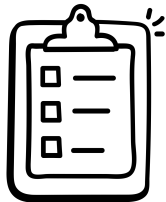
- `mkdir -p /tmp/swab`
- `cd /tmp/swab`
- `cmake /path/to/sourcecode/of/swab`
- `make`
- `make kvm`



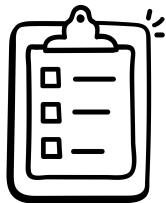


- Condition (condition variable)

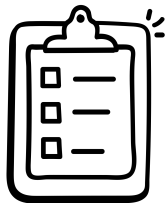




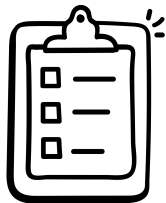
- Condition (condition variable)
- Mutex



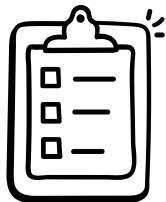
- Condition (condition variable)
- Mutex
- PageManager (allocate and free pages)



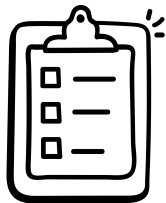
- Condition (condition variable)
- Mutex
- PageManager (allocate and free pages)
- ArchMemory (manages virtual address space)



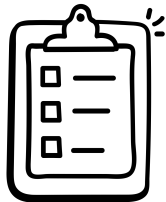
- Condition (condition variable)
- Mutex
- PageManager (allocate and free pages)
- ArchMemory (manages virtual address space)
- main (boot code)



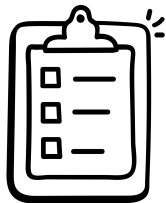
- Condition (condition variable)
- Mutex
- PageManager (allocate and free pages)
- ArchMemory (manages virtual address space)
- main (boot code)
- Syscall (syscall handling)



- Condition (condition variable)
- Mutex
- PageManager (allocate and free pages)
- ArchMemory (manages virtual address space)
- main (boot code)
- Syscall (syscall handling)
- uSTL (STL, provides vector, string, map,...)



- Condition (condition variable)
- Mutex
- PageManager (allocate and free pages)
- ArchMemory (manages virtual address space)
- main (boot code)
- Syscall (syscall handling)
- uSTL (STL, provides vector, string, map,...)
- Loader (user space binary loader)

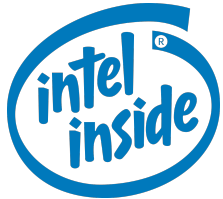


- Condition (condition variable)
- Mutex
- PageManager (allocate and free pages)
- ArchMemory (manages virtual address space)
- main (boot code)
- Syscall (syscall handling)
- uSTL (STL, provides vector, string, map,...)
- Loader (user space binary loader)

→ For more info:

<https://www.iaik.tugraz.at/teaching/materials/os/tutorials/>

VMX







- Intels Virtual-Machine Extension



- Intel's Virtual-Machine Extension
- Provides hardware support for virtualization



- Intel's Virtual-Machine Extension
- Provides hardware support for virtualization
- Two different classes of software:



- Intel's Virtual-Machine Extension
- Provides hardware support for virtualization
- Two different classes of software:
 - Virtual-machine monitors (VMM)



- Intel's Virtual-Machine Extension
- Provides hardware support for virtualization
- Two different classes of software:
 - Virtual-machine monitors (VMM)
 - Guest software







- DO NOT READ IT LIKE A BOOK



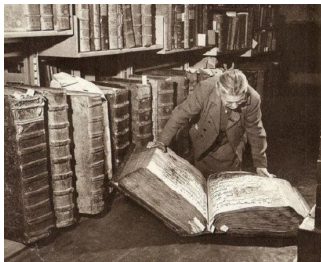
- DO NOT READ IT LIKE A BOOK
- Look things up that you need



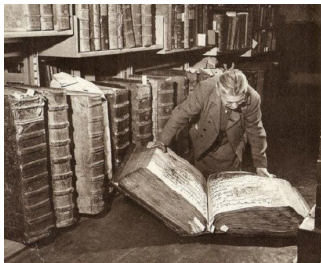
- DO NOT READ IT LIKE A BOOK
- Look things up that you need
- Important volumes for this lecture:



- DO NOT READ IT LIKE A BOOK
- Look things up that you need
- Important volumes for this lecture:
 - Volume 3C: virtual machine extensions (vmx)



- DO NOT READ IT LIKE A BOOK
- Look things up that you need
- Important volumes for this lecture:
 - Volume 3C: virtual machine extensions (vmx)
 - Volume 3D: appendices with values for constants



- DO NOT READ IT LIKE A BOOK
- Look things up that you need
- Important volumes for this lecture:
 - Volume 3C: virtual machine extensions (vmx)
 - Volume 3D: appendices with values for constants
 - Volume 3A: operating-system support environment







- Gives the guest the illusion of running on real hardware



- Gives the guest the illusion of running on real hardware
- Has full control of the processors resources



- Gives the guest the illusion of running on real hardware
- Has full control of the processors resources
- Manages the ability of the guest to access



- Gives the guest the illusion of running on real hardware
- Has full control of the processors resources
- Manages the ability of the guest to access
 - Processor Resources



- Gives the guest the illusion of running on real hardware
- Has full control of the processors resources
- Manages the ability of the guest to access
 - Processor Resources
 - Physical Memory



- Gives the guest the illusion of running on real hardware
- Has full control of the processors resources
- Manages the ability of the guest to access
 - Processor Resources
 - Physical Memory
 - Interrupts



- Gives the guest the illusion of running on real hardware
- Has full control of the processors resources
- Manages the ability of the guest to access
 - Processor Resources
 - Physical Memory
 - Interrupts
 - I/O



- Gives the guest the illusion of running on real hardware
- Has full control of the processors resources
- Manages the ability of the guest to access
 - Processor Resources
 - Physical Memory
 - Interrupts
 - I/O
 - ...

Getting Started

¹Intel SDM Volume 3, 24.7



¹Intel SDM Volume 3, 24.7



¹Intel SDM Volume 3, 24.7



- Check for VMX support: `CPUID.1:ECX.VMX[bit 5] = 1`

¹Intel SDM Volume 3, 24.7



- Check for VMX support: $\text{CPUID.1:ECX.VMX}[\text{bit } 5] = 1$
- Enable VMX by setting bit 13 in CR4 ¹

¹Intel SDM Volume 3, 24.7





- VMX root operation



- VMX root operation
 - new instructions (VMX instructions)



- VMX root operation
 - new instructions (VMX instructions)
 - restrictions on certain control register values





- VMX root operation
 - new instructions (VMX instructions)
 - restrictions on certain control register values
 - used for the VMM



- VMX root operation
 - new instructions (VMX instructions)
 - restrictions on certain control register values
 - used for the VMM
- VMX non-root operation



- VMX root operation
 - new instructions (VMX instructions)
 - restrictions on certain control register values
 - used for the VMM
- VMX non-root operation
 - more restricted than normal operation



- VMX root operation
 - new instructions (VMX instructions)
 - restrictions on certain control register values
 - used for the VMM
- VMX non-root operation
 - more restricted than normal operation
 - certain instructions and events cause VM exits



- VMX root operation
 - new instructions (VMX instructions)
 - restrictions on certain control register values
 - used for the VMM
- VMX non-root operation
 - more restricted than normal operation
 - certain instructions and events cause VM exits
 - used for the guest



- VMX root operation
 - new instructions (VMX instructions)
 - restrictions on certain control register values
 - used for the VMM
- VMX non-root operation
 - more restricted than normal operation
 - certain instructions and events cause VM exits
 - used for the guest
- Transition between the two



- VMX root operation
 - new instructions (VMX instructions)
 - restrictions on certain control register values
 - used for the VMM
- VMX non-root operation
 - more restricted than normal operation
 - certain instructions and events cause VM exits
 - used for the guest
- Transition between the two
 - VM entries: VMX root operation → VMX non-root operation



- VMX root operation
 - new instructions (VMX instructions)
 - restrictions on certain control register values
 - used for the VMM
- VMX non-root operation
 - more restricted than normal operation
 - certain instructions and events cause VM exits
 - used for the guest
- Transition between the two
 - VM entries: VMX root operation \rightarrow VMX non-root operation
 - VM exits: VMX non-root operation \rightarrow VMX root operation

¹Intel SDM Volume 3, 24.8

²Intel SDM Volume 3, 24.7



¹Intel SDM Volume 3, 24.8

²Intel SDM Volume 3, 24.7



Preparations

¹Intel SDM Volume 3, 24.8

²Intel SDM Volume 3, 24.7



Preparations

- Setup CR0¹

¹Intel SDM Volume 3, 24.8

²Intel SDM Volume 3, 24.7



Preparations

- Setup CR0¹
 - Set bits specified by IA32_VMX_CR0_FIXED0 (0x486)

¹Intel SDM Volume 3, 24.8

²Intel SDM Volume 3, 24.7



Preparations

- Setup CR0¹
 - Set bits specified by IA32_VMX_CR0_FIXED0 (0x486)
 - Clear bits specified by IA32_VMX_CR0_FIXED1 (0x487)

¹Intel SDM Volume 3, 24.8

²Intel SDM Volume 3, 24.7



Preparations

- Setup CR0¹
 - Set bits specified by IA32_VMX_CR0_FIXED0 (0x486)
 - Clear bits specified by IA32_VMX_CR0_FIXED1 (0x487)
- Setup CR4¹

¹Intel SDM Volume 3, 24.8

²Intel SDM Volume 3, 24.7



Preparations

- Setup CR0¹
 - Set bits specified by IA32_VMX_CR0_FIXED0 (0x486)
 - Clear bits specified by IA32_VMX_CR0_FIXED1 (0x487)
- Setup CR4¹
 - Set bits specified by IA32_VMX_CR4_FIXED0 (0x488)

¹Intel SDM Volume 3, 24.8

²Intel SDM Volume 3, 24.7



Preparations

- Setup CR0¹
 - Set bits specified by IA32_VMX_CR0_FIXED0 (0x486)
 - Clear bits specified by IA32_VMX_CR0_FIXED1 (0x487)
- Setup CR4¹
 - Set bits specified by IA32_VMX_CR4_FIXED0 (0x488)
 - Clear bits specified by IA32_VMX_CR4_FIXED1 (0x489)

¹Intel SDM Volume 3, 24.8

²Intel SDM Volume 3, 24.7



Preparations

- Setup CR0¹
 - Set bits specified by IA32_VMX_CR0_FIXED0 (0x486)
 - Clear bits specified by IA32_VMX_CR0_FIXED1 (0x487)
- Setup CR4¹
 - Set bits specified by IA32_VMX_CR4_FIXED0 (0x488)
 - Clear bits specified by IA32_VMX_CR4_FIXED1 (0x489)
- **Ensure** that bit 2 of IA32_FEATURE_CONTROL (0x3A) is set²

¹Intel SDM Volume 3, 24.8

²Intel SDM Volume 3, 24.7

¹Intel SDM Volume 3, 31



¹Intel SDM Volume 3, 31

VMXON region ²



¹Intel SDM Volume 3, 31



VMXON region ²

- Used by the processor to support VMX operation

¹Intel SDM Volume 3, 31



VMXON region ²

- Used by the processor to support VMX operation
- Up to 4KB in size

¹Intel SDM Volume 3, 31



VMXON region ²

- Used by the processor to support VMX operation
- Up to 4KB in size
- VMXON pointer needs to be a 4KB aligned valid physical address

¹Intel SDM Volume 3, 31



VMXON region ²

- Used by the processor to support VMX operation
- Up to 4KB in size
- VMXON pointer needs to be a 4KB aligned valid physical address
- Bits 30:0 must contain the VMCS revision identifier (IA32_VMX_BASIC, 0x480)

¹Intel SDM Volume 3, 31



VMXON region ²

- Used by the processor to support VMX operation
- Up to 4KB in size
- VMXON pointer needs to be a 4KB aligned valid physical address
- Bits 30:0 must contain the VMCS revision identifier (IA32_VMX_BASIC, 0x480)
- Can be loaded using the VMXON instruction to enter VMX operation

¹Intel SDM Volume 3, 31





- VMXON



- VMXON
- VMXOFF



- VMXON
- VMXOFF
- INVEPT





- VMXON
- VMXOFF
- INVEPT
- INVVPID



- VMXON
- VMXOFF
- INVEPT
- INVVPID
- VMCLEAR



- VMXON
- VMXOFF
- INVEPT
- INVVPID
- VMCLEAR
- VMPTRLD



- VMXON
- VMXOFF
- INVEPT
- INVVPID
- VMCLEAR
- VMPTRLD
- VMPTRSTR



- VMXON
- VMXOFF
- INVEPT
- INVVPID
- VMCLEAR
- VMPTRLD
- VMPTRSTR
- VMLAUNCH/VMRESUME



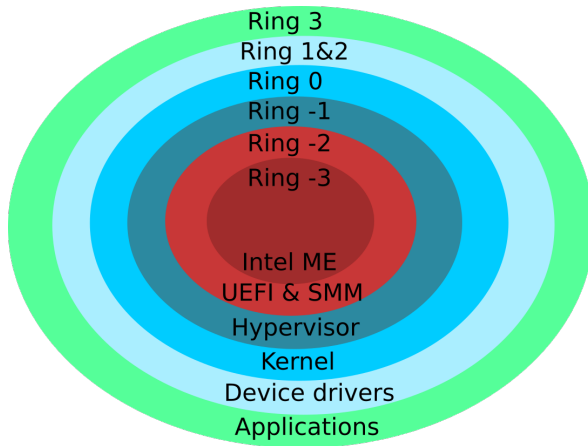
- VMXON
- VMXOFF
- INVEPT
- INVVPID
- VMCLEAR
- VMPTRLD
- VMPTRSTR
- VMLAUNCH/VMRESUME
- VMREAD

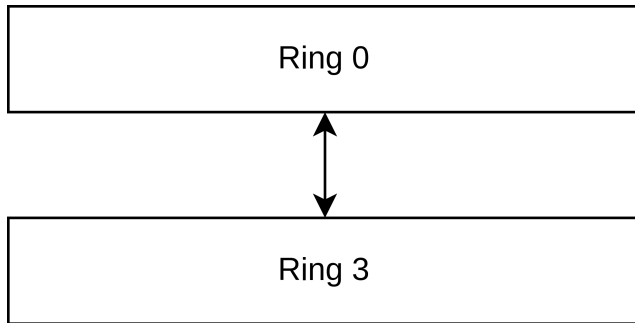


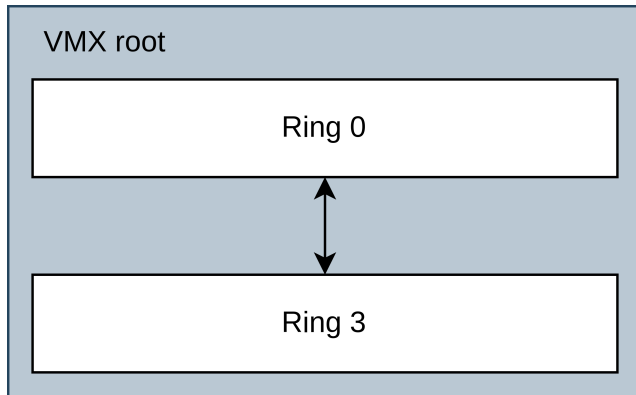
- VMXON
- VMXOFF
- INVEPT
- INVVPID
- VMCLEAR
- VMPTRLD
- VMPTRSTR
- VMLAUNCH/VMRESUME
- VMREAD
- VMWRITE

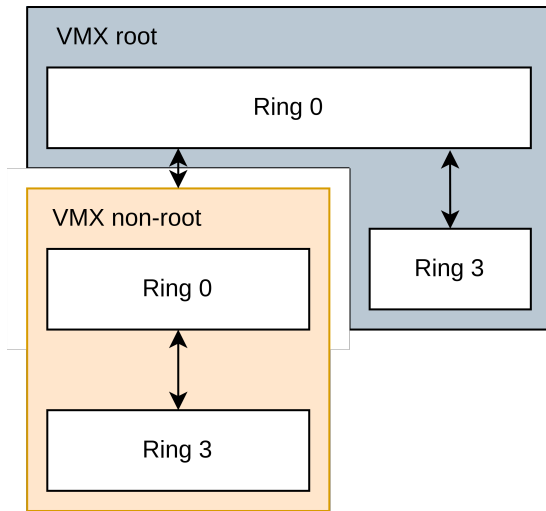


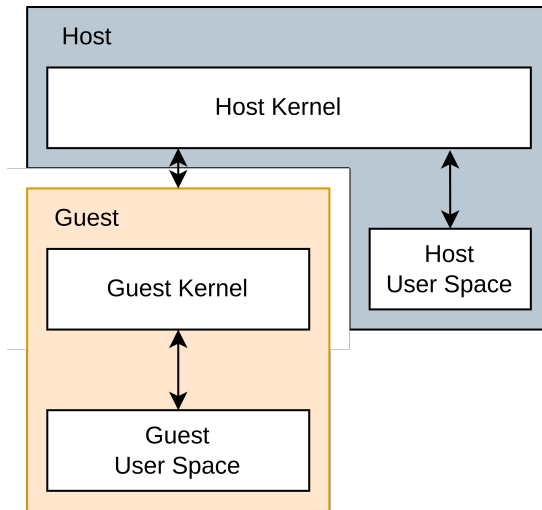
- VMXON
- VMXOFF
- INVEPT
- INVVPID
- VMCLEAR
- VMPTRLD
- VMPTRSTR
- VMLAUNCH/VMRESUME
- VMREAD
- VMWRITE
- VMCALL

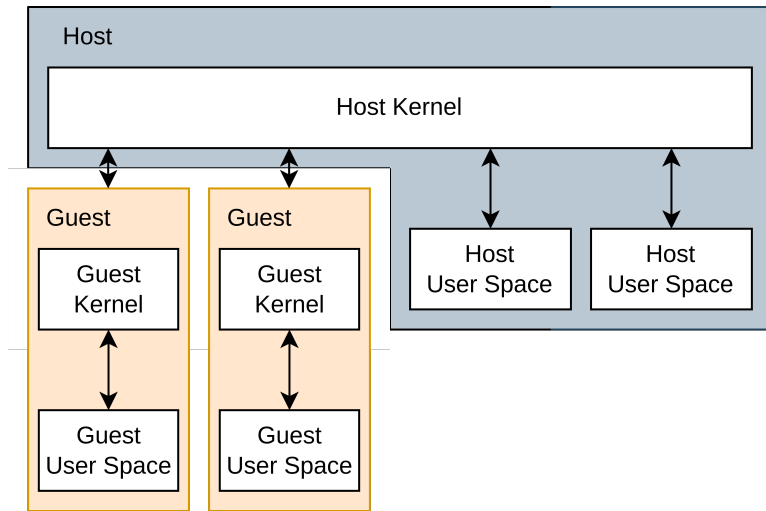


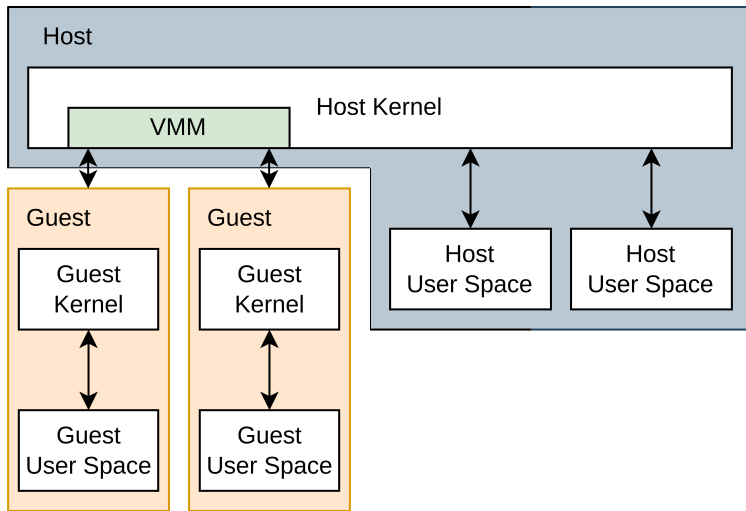


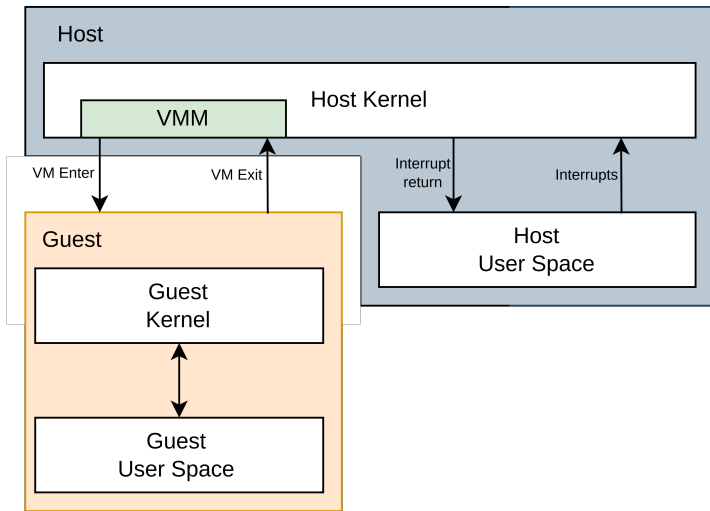






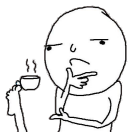






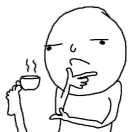
VMCS

²Intel SDM Volume 3, 25.2



²Intel SDM Volume 3, 25.2

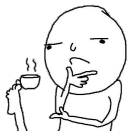
Virtual-Machine Control Data Structure (VMCS)



²Intel SDM Volume 3, 25.2

Virtual-Machine Control Data Structure (VMCS)

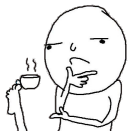
- Manages VM entries and VM exits



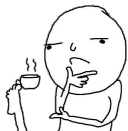
²Intel SDM Volume 3, 25.2

Virtual-Machine Control Data Structure (VMCS)

- Manages VM entries and VM exits
- Used to setup processor behavior in VMX non-root operation



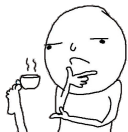
²Intel SDM Volume 3, 25.2



Virtual-Machine Control Data Structure (VMCS)

- Manages VM entries and VM exits
- Used to setup processor behavior in VMX non-root operation
- Can be manipulated using VMCLEAR, VMPTRLD, VMREAD, VMWRITE

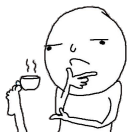
²Intel SDM Volume 3, 25.2



Virtual-Machine Control Data Structure (VMCS)

- Manages VM entries and VM exits
- Used to setup processor behavior in VMX non-root operation
- Can be manipulated using VMCLEAR, VMPTRLD, VMREAD, VMWRITE
- One VMCS per virtual processor

²Intel SDM Volume 3, 25.2



Virtual-Machine Control Data Structure (VMCS)

- Manages VM entries and VM exits
- Used to setup processor behavior in VMX non-root operation
- Can be manipulated using VMCLEAR, VMPTRLD, VMREAD, VMWRITE
- One VMCS per virtual processor
- The current VMCS can be accessed using the VMWRITE and VMREAD instructions

²Intel SDM Volume 3, 25.2

¹Intel SDM Volume 3, 25.2



¹Intel SDM Volume 3, 25.2



¹Intel SDM Volume 3, 25.2



- Up to 4KB in size

¹Intel SDM Volume 3, 25.2



- Up to 4KB in size
- VMCS pointer needs to be a 4KB aligned valid physical address

¹Intel SDM Volume 3, 25.2



- Up to 4KB in size
- VMCS pointer needs to be a 4KB aligned valid physical address
- Bits 30:0 must contain the VMCS revision identifier (IA32_VMX_BASIC, 0x480)¹

¹Intel SDM Volume 3, 25.2







- Guest configuration



- Guest configuration
- Some registers are not easy to restore



- Guest configuration
- Some registers are not easy to restore
 - CR4/CR0



- Guest configuration
- Some registers are not easy to restore
 - CR4/CR0
 - Segment bases and access rights



- Guest configuration
- Some registers are not easy to restore
 - CR4/CR0
 - Segment bases and access rights
 - CR3



- Guest configuration
- Some registers are not easy to restore
 - CR4/CR0
 - Segment bases and access rights
 - CR3
 - ...

- Natural-Width fields.
- 16-bits fields.
- 32-bits fields.
- 64-bits fields.

Copyright 2017, [@Noteworthy](#) (Intel Manual of July 2017)

CONTROL FIELDS

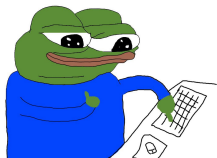
Pin-Based VM-Execution Controls	External-interrupt exiting		NMI exiting		Virtual NMIs	
	Activate VMX-preemption timer			Process posted interrupts		
Primary processor-based VM-execution controls	Interrupt-window exiting			Use TSC offsetting		
	HLT exiting		INVLPG exiting		MWAIT exiting	
	RDTSC exiting		CR3-load exiting		CR3-store exiting	
	CR8-store exiting		Use TPR shadow		NMI-window exiting	
	Unconditional I/O exiting		Use I/O bitmaps		Monitor trap flag	
	MONITOR exiting		PAUSE exiting		Activate secondary controls	
Secondary processor-based VM-execution controls	Virtualize APIC accesses		Enable EPT		Descriptor-table exiting	
	Virtualize x2APIC mode		Enable VPID		WBINVD exiting	
	APIC-register virtualization		Virtual-interrupt delivery		PAUSE-loop exiting	
	RDRAND exiting		Enable INVPCID		Enable VM functions	
	Enable ENCLS exiting		RDSEED exiting		Enable PML	
	Conceal VMX non-root operation from Intel PT			Enable XSAVES/XRSTORS		
	Mode-based execute control for EPT			Use TSC scaling		
Exception Bitmap			I/O-Bitmap Addresses		TSC-offset	
Guest/Host Masks for CR0		Guest/Host Masks for CR4		Read Shadows for CR0		
Read Shadows for CR4						
CR3-target value 0		CR3-target value 1		CR3-target value 2		
CR3-target value 3		CR3-target count				
APIC Virtualization	APIC-access address		Virtual-APIC address		TPR threshold	
	EOI-exit bitmap 0		EOI-exit bitmap 1		EOI-exit bitmap 2	
	EOI-exit bitmap 3					
Posted-interrupt notification vector			Posted-interrupt descriptor address			
Read bitmap for low MSRs		Read bitmap for high MSRs		Write bitmap for low MSRs		
Write bitmap for high MSRs						
Executive-VMCS Pointer			Extended-Page-Table Pointer		Virtual-Processor Identifier	
PLE_Gap		PLE_Window		VM-function controls		
VMREAD bitmap		VMWRITE bitmap				
ENCLS-exiting bitmap			PML address			
Virtualization-exception information address		EPTP index		XSS-exiting bitmap		

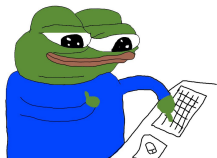
GUEST STATE AREA

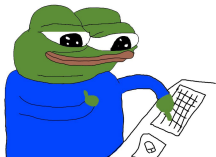
CR0	CR3			CR4	
DR7					
RSP	RIP			RFLAGS	
CS	Selector	Base Address	Segment Limit	Access Right	
SS	Selector	Base Address	Segment Limit	Access Right	
DS	Selector	Base Address	Segment Limit	Access Right	
ES	Selector	Base Address	Segment Limit	Access Right	
FS	Selector	Base Address	Segment Limit	Access Right	
GS	Selector	Base Address	Segment Limit	Access Right	
LDTR	Selector	Base Address	Segment Limit	Access Right	
TR	Selector	Base Address	Segment Limit	Access Right	
GDTR	Selector	Base Address	Limit	Access Right	
IDTR	Selector	Base Address	Limit	Access Right	
IA32_DEBUGCTL	IA32_SYSENTER_CS	IA32_SYSENTER_ESP		IA32_SYSENTER_EIP	
IA32_PERF_GLOBAL_CTRL	IA32_PAT	IA32_EFER		IA32_BNDCFGS	
SMBASE					
Activity state	Interruptibility state				
Pending debug exceptions					
VMCS link pointer					
VMX-preemption timer value					
Page-directory-pointer-table entries	PDPTE0	PDPTE1	PDPTE2	PDPTE3	
Guest interrupt status					
PML index					

HOST STATE AREA

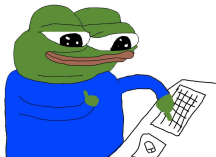
CRO		CR3		CR4	
RSP			RIP		
CS	Selector				
SS	Selector				
DS	Selector				
ES	Selector				
FS	Selector	Base Address			
GS	Selector	Base Address			
TR	Selector	Base Address			
GDTR	Base Address				
IDTR	Base Address				
IA32_SYSENTER_CS		IA32_SYSENTER_ESP		IA32_SYSENTER_EIP	
IA32_PERF_GLOBAL_CTRL		IA32_PAT		IA32_EFER	







- What about things that are not in the guest state area, like general-purpose registers?



- What about things that are not in the guest state area, like general-purpose registers?
- We have to save and restore them by hand!



- What about things that are not in the guest state area, like general-purpose registers?
 - We have to save and restore them by hand!
- We need a context switch

- Host EFER

- Host EFER
- Host PAT

- Host EFER
- Host PAT
- Host CRx

- Host EFER
- Host PAT
- Host CRx
- Host RSP

- Host EFER
- Host PAT
- Host CRx
- Host RSP
- Host RIP

- Host EFER
- Host PAT
- Host CRx
- Host RSP
- Host RIP
- Host segment selectors

- Host EFER
- Host PAT
- Host CRx
- Host RSP
- Host RIP
- Host segment selectors
- TR selector and base address

- Host EFER
- Host PAT
- Host CRx
- Host RSP
- Host RIP
- Host segment selectors
- TR selector and base address
- GDTR and IDTR base address

- Guest EFER

- Guest EFER
- Guest PAT

- Guest EFER
- Guest PAT
- Guest CRx

- Guest EFER
- Guest PAT
- Guest CRx
- Guest RSP

- Guest EFER
- Guest PAT
- Guest CRx
- Guest RSP
- Guest RIP

- Guest EFER
- Guest PAT
- Guest CRx
- Guest RSP
- Guest RIP
- Guest RFLAGS

- Guest EFER
- Guest PAT
- Guest CRx
- Guest RSP
- Guest RIP
- Guest RFLAGS
- Guest segment selectors, base addresses, limits, access rights

- Guest EFER
- Guest PAT
- Guest CRx
- Guest RSP
- Guest RIP
- Guest RFLAGS
- Guest segment selectors, base addresses, limits, access rights
- GDTR and IDTR base address, limit

- Guest EFER
- Guest PAT
- Guest CRx
- Guest RSP
- Guest RIP
- Guest RFLAGS
- Guest segment selectors, base addresses, limits, access rights
- GDTR and IDTR base address, limit
- VMCS link pointer (set it to -1)

- Pin-Based VM-Execution Controls

- Pin-Based VM-Execution Controls
- Primary Processor-Based VM-Execution Controls

- Pin-Based VM-Execution Controls
- Primary Processor-Based VM-Execution Controls
- Secondary Processor-Based VM-Execution Controls

- Pin-Based VM-Execution Controls
- Primary Processor-Based VM-Execution Controls
- Secondary Processor-Based VM-Execution Controls
- VM entry controls

- Pin-Based VM-Execution Controls
- Primary Processor-Based VM-Execution Controls
- Secondary Processor-Based VM-Execution Controls
- VM entry controls
- VM exit controls

- Pin-Based VM-Execution Controls
- Primary Processor-Based VM-Execution Controls
- Secondary Processor-Based VM-Execution Controls
- VM entry controls
- VM exit controls
- EPTP

Running VMs

¹Intel SDM Volume 3, 25.1



¹Intel SDM Volume 3, 25.1



¹Intel SDM Volume 3, 25.1



- VMLAUNCH

¹Intel SDM Volume 3, 25.1



- VMLAUNCH
→ VMCS not launched yet

¹Intel SDM Volume 3, 25.1



- VMLAUNCH
 - VMCS not launched yet
- VMRESUME

¹Intel SDM Volume 3, 25.1



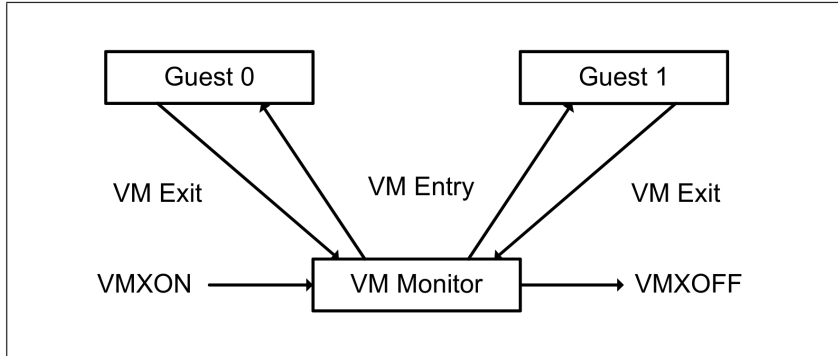
- VMLAUNCH
 - VMCS not launched yet
- VMRESUME
 - Continue launched VMCS

¹Intel SDM Volume 3, 25.1

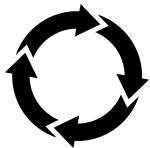


- VMLAUNCH
 - VMCS not launched yet
- VMRESUME
 - Continue launched VMCS
- VMCLEAR resets the launch state!

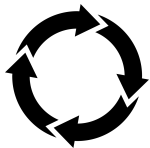
¹Intel SDM Volume 3, 25.1



¹Intel SDM Volume 3, 25.1

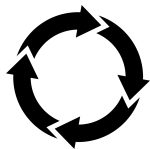


¹Intel SDM Volume 3, 25.1

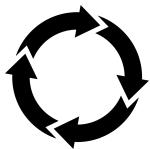


¹Intel SDM Volume 3, 25.1

- Current

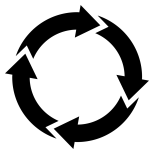


¹Intel SDM Volume 3, 25.1



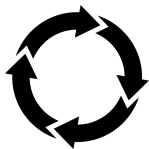
- Current
 - only one VMCS can is current at a time

¹Intel SDM Volume 3, 25.1



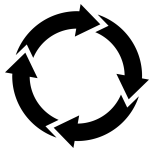
- Current
 - only one VMCS can be current at a time
 - VMWRITE, VMREAD, VMLAUNCH, VMRESUME operate on the current VMCS

¹Intel SDM Volume 3, 25.1



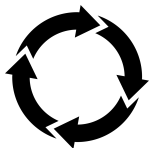
- Current
 - only one VMCS can be current at a time
 - VMWRITE, VMREAD, VMLAUNCH, VMRESUME operate on the current VMCS
- Active

¹Intel SDM Volume 3, 25.1



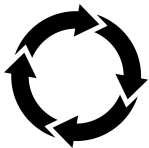
- Current
 - only one VMCS can be current at a time
 - VMWRITE, VMREAD, VMLAUNCH, VMRESUME operate on the current VMCS
- Active
 - VMCS is managed by the CPU

¹Intel SDM Volume 3, 25.1



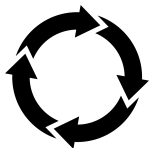
- Current
 - only one VMCS can be current at a time
 - VMWRITE, VMREAD, VMLAUNCH, VMRESUME operate on the current VMCS
- Active
 - VMCS is managed by the CPU
 - multiple VMCS can be active at a time

¹Intel SDM Volume 3, 25.1



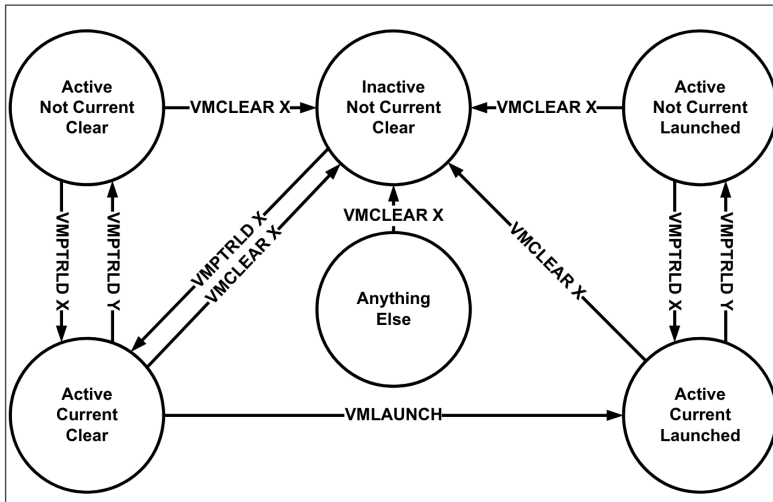
- Current
 - only one VMCS can be current at a time
 - VMWRITE, VMREAD, VMLAUNCH, VMRESUME operate on the current VMCS
- Active
 - VMCS is managed by the CPU
 - multiple VMCS can be active at a time
- Launched

¹Intel SDM Volume 3, 25.1



- Current
 - only one VMCS can be current at a time
 - VMWRITE, VMREAD, VMLAUNCH, VMRESUME operate on the current VMCS
- Active
 - VMCS is managed by the CPU
 - multiple VMCS can be active at a time
- Launched
 - VMCS is launched and can only be continued with VMRESUME

¹Intel SDM Volume 3, 25.1



¹Intel SDM Volume 3, 25.1



¹Intel SDM Volume 3, 25.1



¹Intel SDM Volume 3, 25.1

- VMCLEAR



¹Intel SDM Volume 3, 25.1

- VMCLEAR
 - initializes a new VMCS



¹Intel SDM Volume 3, 25.1

- VMCLEAR
 - initializes a new VMCS
 - sets active VMCS to inactive, not current, and not launched



¹Intel SDM Volume 3, 25.1



- VMCLEAR
 - initializes a new VMCS
 - sets active VMCS to inactive, not current, and not launched
- VMPTRLD

¹Intel SDM Volume 3, 25.1



- VMCLEAR
 - initializes a new VMCS
 - sets active VMCS to inactive, not current, and not launched
- VMPTRLD
 - sets VMCS to active and current

¹Intel SDM Volume 3, 25.1



- VMCLEAR
 - initializes a new VMCS
 - sets active VMCS to inactive, not current, and not launched
- VMPTRLD
 - sets VMCS to active and current
 - sets previous current VMCS to not current

¹Intel SDM Volume 3, 25.1



- VMCLEAR
 - initializes a new VMCS
 - sets active VMCS to inactive, not current, and not launched
- VMPTRLD
 - sets VMCS to active and current
 - sets previous current VMCS to not current
- VMLAUNCH

¹Intel SDM Volume 3, 25.1



- VMCLEAR
 - initializes a new VMCS
 - sets active VMCS to inactive, not current, and not launched
- VMPTRLD
 - sets VMCS to active and current
 - sets previous current VMCS to not current
- VMLAUNCH
 - sets VMCS to launched

¹Intel SDM Volume 3, 25.1



- VMCLEAR
 - initializes a new VMCS
 - sets active VMCS to inactive, not current, and not launched
- VMPTRLD
 - sets VMCS to active and current
 - sets previous current VMCS to not current
- VMLAUNCH
 - sets VMCS to launched
- VMRESMUE

¹Intel SDM Volume 3, 25.1

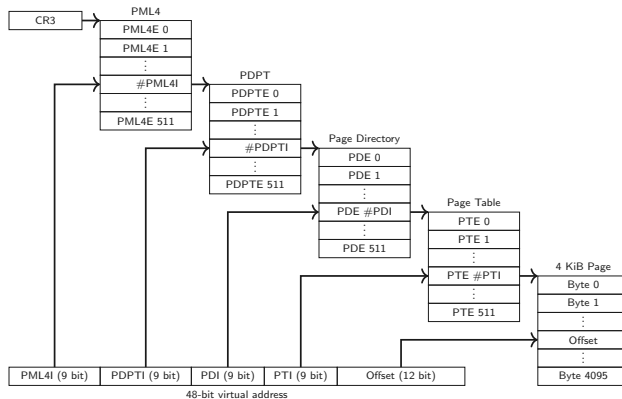


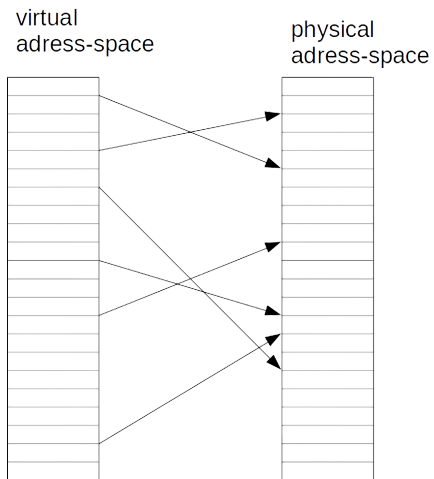
- VMCLEAR
 - initializes a new VMCS
 - sets active VMCS to inactive, not current, and not launched
- VMPTRLD
 - sets VMCS to active and current
 - sets previous current VMCS to not current
- VMLAUNCH
 - sets VMCS to launched
- VMRESMUE
 - can only be used on launched VMCS

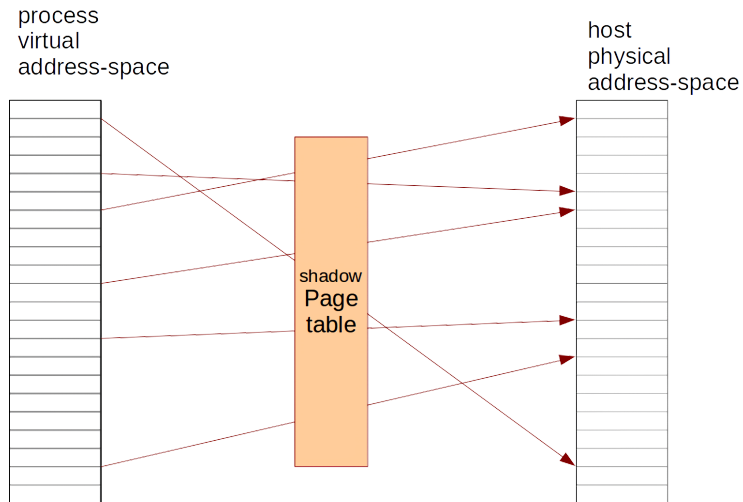
¹Intel SDM Volume 3, 25.1

EPT

How can we give a guest access to a physical address space without giving it access to the hosts physical address space?











- guest runs on a normal page table





- guest runs on a normal page table
- guest can not modify their real page table



- guest runs on a normal page table
- guest can not modify their real page table
- instead the guest has a copy of the real page table



- guest runs on a normal page table
- guest can not modify their real page table
- instead the guest has a copy of the real page table
- when the guest changes page table copy



- guest runs on a normal page table
- guest can not modify their real page table
- instead the guest has a copy of the real page table
- when the guest changes page table copy
 - Hypervisor has to catch access



- guest runs on a normal page table
- guest can not modify their real page table
- instead the guest has a copy of the real page table
- when the guest changes page table copy
 - Hypervisor has to catch access
 - update shadow (real) page table



- guest runs on a normal page table
 - guest can not modify their real page table
 - instead the guest has a copy of the real page table
 - when the guest changes page table copy
 - Hypervisor has to catch access
 - update shadow (real) page table
- guest can manage virtual memory without access to arbitrary physical memory or real PPNs





Advantages:





Advantages:

- + address translations are fast



Advantages:

- + address translations are fast
- + requires only minimal hardware support



Advantages:

- + address translations are fast
- + requires only minimal hardware support

Disadvantages:



Advantages:

- + address translations are fast
- + requires only minimal hardware support

Disadvantages:

- page table modifications are slow



Advantages:

- + address translations are fast
- + requires only minimal hardware support

Disadvantages:

- page table modifications are slow
- only works with paging enabled



Advantages:

- + address translations are fast
- + requires only minimal hardware support

Disadvantages:

- page table modifications are slow
- only works with paging enabled
- x2 memory overhead for page tables in the guest



Advantages:

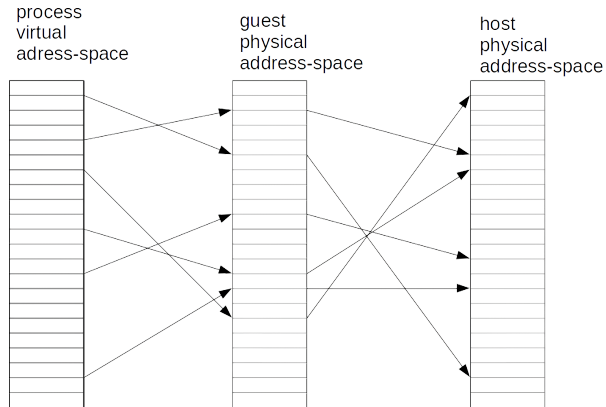
- + address translations are fast
- + requires only minimal hardware support

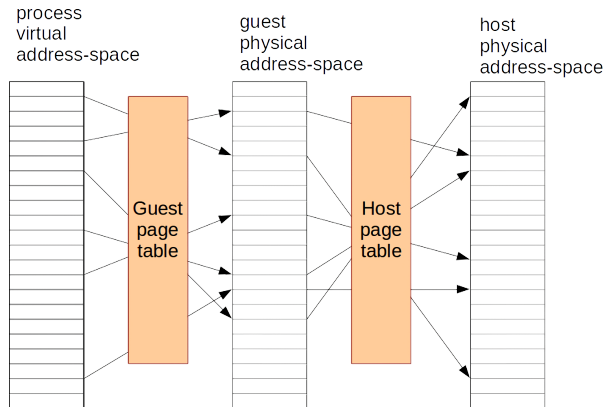
Disadvantages:

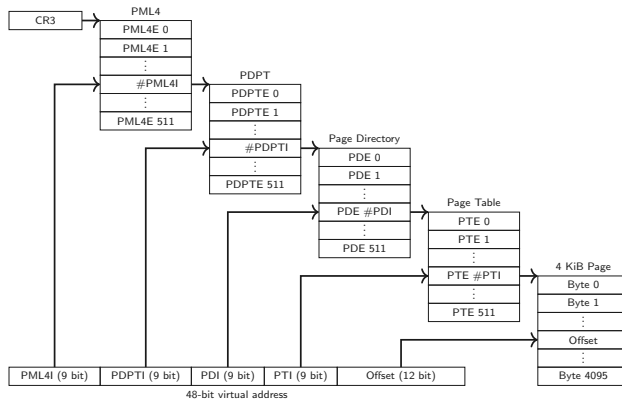
- page table modifications are slow
- only works with paging enabled
- x2 memory overhead for page tables in the guest
- complicated to implement

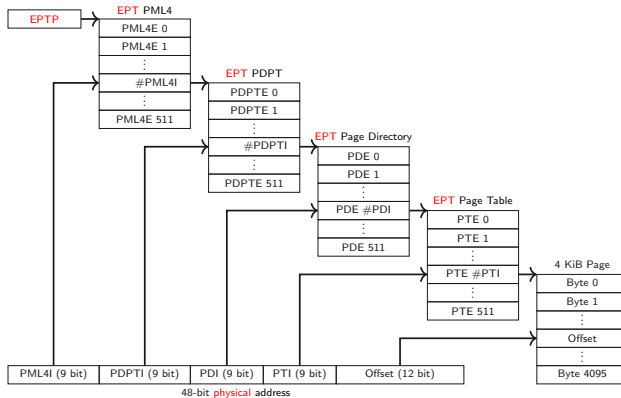


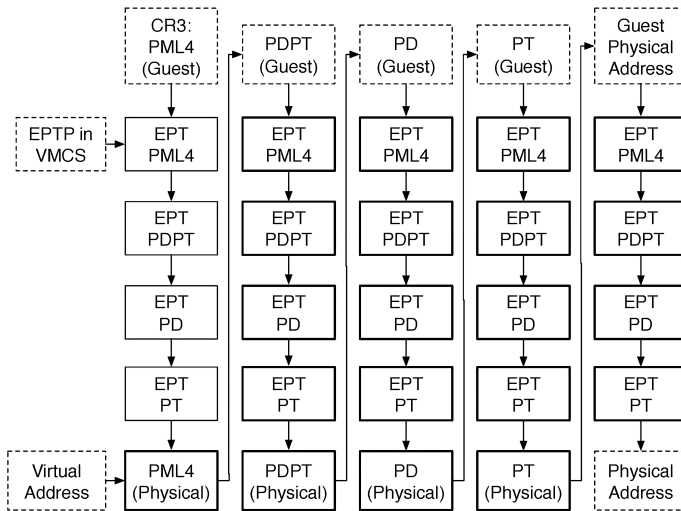
WE NEED TO GO DEEPER

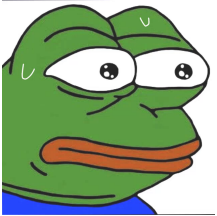


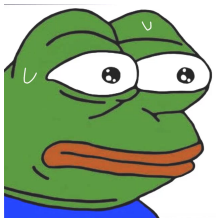


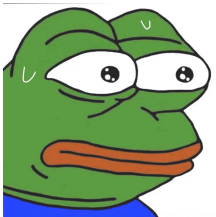




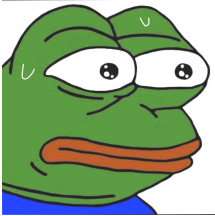




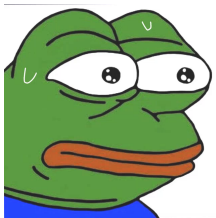




- EPT translations are cached in the TLB



- EPT translations are cached in the TLB
- Translations are tagged with the EPTP



- EPT translations are cached in the TLB
- Translations are tagged with the EPTP
- We can invalidate all TLB entries for a given EPTP or all EPT TLB entries using the INVEPT instruction

VM Exit

¹Intel SDM Volume 3, 28.1



¹Intel SDM Volume 3, 28.1



¹Intel SDM Volume 3, 28.1

- An event that causes a VM exit directly **does not** update the architectural state



¹Intel SDM Volume 3, 28.1

- An event that causes a VM exit directly **does not** update the architectural state
- An event that causes a VM exit indirectly **does** update the architectural state (example: PF causes a write to CR2 and CR2 exiting is enabled)



¹Intel SDM Volume 3, 28.1



- An event that causes a VM exit directly **does not** update the architectural state
- An event that causes a VM exit indirectly **does** update the architectural state (example: PF causes a write to CR2 and CR2 exiting is enabled)
- If a VM exit is triggered by executing specified instructions we can gain further information

¹Intel SDM Volume 3, 28.1



- An event that causes a VM exit directly **does not** update the architectural state
- An event that causes a VM exit indirectly **does** update the architectural state (example: PF causes a write to CR2 and CR2 exiting is enabled)
- If a VM exit is triggered by executing specified instructions we can gain further information
 - VM exit qualification → VM-exit reason specific information

¹Intel SDM Volume 3, 28.1



- An event that causes a VM exit directly **does not** update the architectural state
- An event that causes a VM exit indirectly **does** update the architectural state (example: PF causes a write to CR2 and CR2 exiting is enabled)
- If a VM exit is triggered by executing specified instructions we can gain further information
 - VM exit qualification → VM-exit reason specific information
 - VM exit instruction length → length of the instruction that caused the VM exit

¹Intel SDM Volume 3, 28.1



- An event that causes a VM exit directly **does not** update the architectural state
- An event that causes a VM exit indirectly **does** update the architectural state (example: PF causes a write to CR2 and CR2 exiting is enabled)
- If a VM exit is triggered by executing specified instructions we can gain further information
 - VM exit qualification → VM-exit reason specific information
 - VM exit instruction length → length of the instruction that caused the VM exit
 - VM exit instruction information → detailed information on the instruction that caused the VM exit

¹Intel SDM Volume 3, 28.1



- An event that causes a VM exit directly **does not** update the architectural state
- An event that causes a VM exit indirectly **does** update the architectural state (example: PF causes a write to CR2 and CR2 exiting is enabled)
- If a VM exit is triggered by executing specified instructions we can gain further information
 - VM exit qualification → VM-exit reason specific information
 - VM exit instruction length → length of the instruction that caused the VM exit
 - VM exit instruction information → detailed information on the instruction that caused the VM exit
- Host RFLAGS is cleared (except bit 1)

¹Intel SDM Volume 3, 28.1

Table 27-5. Exit Qualification for I/O Instructions

Bit Position(s)	Contents
2:0	Size of access: 0 = 1-byte 1 = 2-byte 3 = 4-byte Other values not used
3	Direction of the attempted access (0 = OUT, 1 = IN)
4	String instruction (0 = not string; 1 = string)
5	REP prefixed (0 = not REP; 1 = REP)
6	Operand encoding (0 = DX, 1 = immediate)
15:7	Not currently defined
31:16	Port number (as specified in DX or in an immediate operand)
63:32	Not currently defined. These bits exist only on processors that support Intel 64 architecture.

- Instructions that cause VM exits unconditionally

- Instructions that cause VM exits unconditionally
 - CPUID

- Instructions that cause VM exits unconditionally
 - CPUID
 - GETSEC

- Instructions that cause VM exits unconditionally
 - CPUID
 - GETSEC
 - INVD

- Instructions that cause VM exits unconditionally
 - CPUID
 - GETSEC
 - INVD
 - XSETBV

- Instructions that cause VM exits unconditionally
 - CPUID
 - GETSEC
 - INVD
 - XSETBV
 - VMX instructions (excluding VMFUNC)

- Instructions that cause VM exits unconditionally
 - CPUID
 - GETSEC
 - INVD
 - XSETBV
 - VMX instructions (excluding VMFUNC)
- Instructions that cause VM exits conditionally

- Instructions that cause VM exits unconditionally
 - CPUID
 - GETSEC
 - INVD
 - XSETBV
 - VMX instructions (excluding VMFUNC)
- Instructions that cause VM exits conditionally
 - RDMSR/WRMSR

- Instructions that cause VM exits unconditionally
 - CPUID
 - GETSEC
 - INVD
 - XSETBV
 - VMX instructions (excluding VMFUNC)
- Instructions that cause VM exits conditionally
 - RDMSR/WRMSR
 - IN/OUT

- Instructions that cause VM exits unconditionally
 - CPUID
 - GETSEC
 - INVD
 - XSETBV
 - VMX instructions (excluding VMFUNC)
- Instructions that cause VM exits conditionally
 - RDMSR/WRMSR
 - IN/OUT
 - MOV to/from CRx

- Instructions that cause VM exits unconditionally
 - CPUID
 - GETSEC
 - INVD
 - XSETBV
 - VMX instructions (excluding VMFUNC)
- Instructions that cause VM exits conditionally
 - RDMSR/WRMSR
 - IN/OUT
 - MOV to/from CRx
 - PAUSE

- Instructions that cause VM exits unconditionally
 - CPUID
 - GETSEC
 - INVD
 - XSETBV
 - VMX instructions (excluding VMFUNC)
- Instructions that cause VM exits conditionally
 - RDMSR/WRMSR
 - IN/OUT
 - MOV to/from CRx
 - PAUSE
 - ...

- Instructions that cause VM exits unconditionally
 - CPUID
 - GETSEC
 - INVD
 - XSETBV
 - VMX instructions (excluding VMFUNC)
- Instructions that cause VM exits conditionally
 - RDMSR/WRMSR
 - IN/OUT
 - MOV to/from CRx
 - PAUSE
 - ...
- Exceptions

- Instructions that cause VM exits unconditionally
 - CPUID
 - GETSEC
 - INVD
 - XSETBV
 - VMX instructions (excluding VMFUNC)
- Instructions that cause VM exits conditionally
 - RDMSR/WRMSR
 - IN/OUT
 - MOV to/from CRx
 - PAUSE
 - ...
- Exceptions
- VMX-preemption timer

- Instructions that cause VM exits unconditionally
 - CPUID
 - GETSEC
 - INVD
 - XSETBV
 - VMX instructions (excluding VMFUNC)
- Instructions that cause VM exits conditionally
 - RDMSR/WRMSR
 - IN/OUT
 - MOV to/from CRx
 - PAUSE
 - ...
- Exceptions
- VMX-preemption timer
- NMIs

- Instructions that cause VM exits unconditionally
 - CPUID
 - GETSEC
 - INVD
 - XSETBV
 - VMX instructions (excluding VMFUNC)
- Instructions that cause VM exits conditionally
 - RDMSR/WRMSR
 - IN/OUT
 - MOV to/from CRx
 - PAUSE
 - ...
- Exceptions
- VMX-preemption timer
- NMIs
- ...

```
main:
  push  rbx
  xor   eax, eax
  sub   rsp, 0x10
  cpuid
  mov   rdi, rsp
  mov   DWORD PTR [rsp], ebx
  mov   DWORD PTR [rsp+0x4], ecx
  mov   DWORD PTR [rsp+0x8], edx
  mov   DWORD PTR [rsp+0x10], 0x0
  call  401030 <puts@plt>
  add   rsp, 0x10
  xor   eax, eax
  pop   rbx
  ret
```

Registers

rax	<rax>
rbx	<rbx>
rcx	<rcx>
rdx	<rdx>
rdi	<rdi>
rsp	0x7fff7830478
rip	main+0


```
main:
```

```
push rbx  
xor eax, eax  
sub rsp, 0x10  
cpuid  
mov rdi, rsp  
mov DWORD PTR [rsp], ebx  
mov DWORD PTR [rsp+0x4], ecx  
mov DWORD PTR [rsp+0x8], edx  
mov DWORD PTR [rsp+0x10], 0x0  
call 401030 <puts@plt>  
add rsp, 0x10  
xor eax, eax  
pop rbx  
ret
```

Registers

rax	<rax>
rbx	<rbx>
rcx	<rcx>
rdx	<rdx>
rdi	<rdi>
rsp	0x7fff7830478
rip	main+0

```
main:
  push  rbx
  xor   eax, eax
  sub   rsp, 0x10
  cpuid
  mov   rdi, rsp
  mov   DWORD PTR [rsp], ebx
  mov   DWORD PTR [rsp+0x4], ecx
  mov   DWORD PTR [rsp+0x8], edx
  mov   DWORD PTR [rsp+0x10], 0x0
  call  401030 <puts@plt>
  add   rsp, 0x10
  xor   eax, eax
  pop   rbx
  ret
```

Registers

rax	<rax>
rbx	<rbx>
rcx	<rcx>
rdx	<rdx>
rdi	<rdi>
rsp	0x7fff7830470
rip	main+1

```
main:
  push  rbx
  xor   eax, eax
  sub   rsp, 0x10
  cpuid
  mov   rdi, rsp
  mov   DWORD PTR [rsp], ebx
  mov   DWORD PTR [rsp+0x4], ecx
  mov   DWORD PTR [rsp+0x8], edx
  mov   DWORD PTR [rsp+0x10], 0x0
  call  401030 <puts@plt>
  add   rsp, 0x10
  xor   eax, eax
  pop   rbx
  ret
```

Registers

rax	0x0
rbx	<rbx>
rcx	<rcx>
rdx	<rdx>
rdi	<rdi>
rsp	0x7ffff7830470
rip	main+3

```
main:
  push  rbx
  xor   eax, eax
  sub   rsp, 0x10
  cpuid
  mov   rdi, rsp
  mov   DWORD PTR [rsp], ebx
  mov   DWORD PTR [rsp+0x4], ecx
  mov   DWORD PTR [rsp+0x8], edx
  mov   DWORD PTR [rsp+0x10], 0x0
  call  401030 <puts@plt>
  add   rsp, 0x10
  xor   eax, eax
  pop   rbx
  ret
```

Registers

rax	0x0
rbx	<rbx>
rcx	<rcx>
rdx	<rdx>
rdi	<rdi>
rsp	0x7ffff7830460
rip	main+7

```
main:
  push  rbx
  xor    eax, eax
  sub   rsp, 0x10
  cpuid
  mov   rdi, rsp
  mov   DWORD PTR [rsp], ebx
  mov   DWORD PTR [rsp+0x4], ecx
  mov   DWORD PTR [rsp+0x8], edx
  mov   DWORD PTR [rsp+0x10], 0x0
  call  401030 <puts@plt>
  add   rsp, 0x10
  xor   eax, eax
  pop   rbx
  ret
```

Registers

rax	0xd
rbx	0x756e6547
rcx	0x6c65746e
rdx	0x49656e69
rdi	<rdi>
rsp	0x7fff7830460
rip	main+9

```
main:
  push  rbx
  xor   eax, eax
  sub   rsp, 0x10
  cpuid
  mov   rdi, rsp
  mov   DWORD PTR [rsp], ebx
  mov   DWORD PTR [rsp+0x4], ecx
  mov   DWORD PTR [rsp+0x8], edx
  mov   DWORD PTR [rsp+0x10], 0x0
  call  401030 <puts@plt>
  add   rsp, 0x10
  xor   eax, eax
  pop   rbx
  ret
```

Registers

rax	0xd
rbx	"Genu"
rcx	"ntel"
rdx	"inel"
rdi	<rdi>
rsp	0x7fff7830460
rip	main+9

```
main:
  push  rbx
  xor   eax, eax
  sub   rsp, 0x10
  cpuid
  mov   rdi, rsp
  mov   DWORD PTR [rsp], ebx
  mov   DWORD PTR [rsp+0x4], ecx
  mov   DWORD PTR [rsp+0x8], edx
  mov   DWORD PTR [rsp+0x10], 0x0
  call  401030 <puts@plt>
  add   rsp, 0x10
  xor   eax, eax
  pop   rbx
  ret
```

Registers

rax	0xd
rbx	"Genu"
rcx	"ntel"
rdx	"inel"
rdi	0x7ffff7830460
rsp	0x7ffff7830460
rip	main+12

```
main:
  push  rbx
  xor   eax, eax
  sub   rsp, 0x10
  cpuid
  mov   rdi, rsp
  mov   DWORD PTR [rsp], ebx
  mov   DWORD PTR [rsp+0x4], ecx
  mov   DWORD PTR [rsp+0x8], edx
  mov   DWORD PTR [rsp+0x10], 0x0
  call  401030 <puts@plt>
  add   rsp, 0x10
  xor   eax, eax
  pop   rbx
  ret
```

Registers

rax	0xd
rbx	"Genu"
rcx	"ntel"
rdx	"inel"
rdi	0x7ffff7830460
rsp	0x7ffff7830460
rip	main+16


```
main:
  push  rbx
  xor   eax, eax
  sub   rsp, 0x10
  cpuid
  mov   rdi, rsp
  mov   DWORD PTR [rsp], ebx
  mov   DWORD PTR [rsp+0x4], ecx
  mov   DWORD PTR [rsp+0x8], edx
  mov   DWORD PTR [rsp+0x10], 0x0
  call  401030 <puts@plt>
  add   rsp, 0x10
  xor   eax, eax
  pop   rbx
  ret
```

Registers

rax	0xd
rbx	"Genu"
rcx	"ntel"
rdx	"inel"
rdi	0x7ffff7830460
rsp	0x7ffff7830460
rip	main+21

```
main:
  push  rbx
  xor    eax, eax
  sub   rsp, 0x10
  cpuid
  mov   rdi, rsp
  mov   DWORD PTR [rsp], ebx
  mov   DWORD PTR [rsp+0x4], ecx
  mov   DWORD PTR [rsp+0x8], edx
  mov   DWORD PTR [rsp+0x10], 0x0
  call  401030 <puts@plt>
  add   rsp, 0x10
  xor   eax, eax
  pop   rbx
  ret
```

Registers

rax	0xd
rbx	"Genu"
rcx	"ntel"
rdx	"inel"
rdi	0x7ffff7830460
rsp	0x7ffff7830460
rip	main+26

```
main:
  push  rbx
  xor   eax, eax
  sub   rsp, 0x10
  cpuid
  mov   rdi, rsp
  mov   DWORD PTR [rsp], ebx
  mov   DWORD PTR [rsp+0x4], ecx
  mov   DWORD PTR [rsp+0x8], edx
  mov   DWORD PTR [rsp+0x10], 0x0
  call  401030 <puts@plt>
  add   rsp, 0x10
  xor   eax, eax
  pop   rbx
  ret
```

Registers

rax	0xd
rbx	"Genu"
rcx	"ntel"
rdx	"inel"
rdi	0x7ffff7830460
rsp	0x7ffff7830460
rip	main+34

```
main:
  push  rbx
  xor   eax, eax
  sub   rsp, 0x10
  cpuid
  mov   rdi, rsp
  mov   DWORD PTR [rsp], ebx
  mov   DWORD PTR [rsp+0x4], ecx
  mov   DWORD PTR [rsp+0x8], edx
  mov   DWORD PTR [rsp+0x10], 0x0
  call  401030 <puts@plt>
  add   rsp, 0x10
  xor   eax, eax
  pop   rbx
  ret
```

Registers

rax	0xc
rbx	"Genu"
rcx	<rcx>
rdx	<rdx>
rdi	<rdi>
rsp	0x7ffff7830460
rip	main+39

```
main:
  push  rbx
  xor   eax, eax
  sub   rsp, 0x10
  cpuid
  mov   rdi, rsp
  mov   DWORD PTR [rsp], ebx
  mov   DWORD PTR [rsp+0x4], ecx
  mov   DWORD PTR [rsp+0x8], edx
  mov   DWORD PTR [rsp+0x10], 0x0
  call  401030 <puts@plt>
  add   rsp, 0x10
  xor   eax, eax
  pop   rbx
  ret
```

Registers

rax	0xc
rbx	"Genu"
rcx	<rcx>
rdx	<rdx>
rdi	<rdi>
rsp	0x7ffff7830470
rip	main+43

```
main:
  push  rbx
  xor    eax, eax
  sub   rsp, 0x10
  cpuid
  mov   rdi, rsp
  mov   DWORD PTR [rsp], ebx
  mov   DWORD PTR [rsp+0x4], ecx
  mov   DWORD PTR [rsp+0x8], edx
  mov   DWORD PTR [rsp+0x10], 0x0
  call  401030 <puts@plt>
  add   rsp, 0x10
  xor   eax, eax
  pop   rbx
  ret
```

Registers

rax	0x0
rbx	"Genu"
rcx	<rcx>
rdx	<rdx>
rdi	<rdi>
rsp	0x7ffff7830470
rip	main+45

```
main:
  push  rbx
  xor   eax, eax
  sub   rsp, 0x10
  cpuid
  mov   rdi, rsp
  mov   DWORD PTR [rsp], ebx
  mov   DWORD PTR [rsp+0x4], ecx
  mov   DWORD PTR [rsp+0x8], edx
  mov   DWORD PTR [rsp+0x10], 0x0
  call  401030 <puts@plt>
  add   rsp, 0x10
  xor   eax, eax
  pop   rbx
  ret
```

Registers

rax	0x0
rbx	<rbx>
rcx	<rcx>
rdx	<rdx>
rdi	<rdi>
rsp	0x7ffff7830460
rip	main+46

```

main:
  push  rbx
  xor    eax, eax
  sub    rsp, 0x10
  cpuid
  mov    rdi, rsp
  mov    DWORD PTR [rsp], ebx
  mov    DWORD PTR [rsp+0x4], ecx
  mov    DWORD PTR [rsp+0x8], edx
  mov    DWORD PTR [rsp+0x10], 0x0
  call   401030 <puts@plt>
  add    rsp, 0x10
  xor    eax, eax
  pop    rbx
  ret
    
```

Registers

rax	<rax>
rbx	<rbx>
rcx	<rcx>
rdx	<rdx>
rdi	<rdi>
rsp	0x7ffff7830478
rip	main+0

main:

```

push    rbx
xor     eax, eax
sub     rsp, 0x10
cpuid
mov     rdi, rsp
mov     DWORD PTR [rsp], ebx
mov     DWORD PTR [rsp+0x4], ecx
mov     DWORD PTR [rsp+0x8], edx
mov     DWORD PTR [rsp+0x10], 0x0
call   401030 <puts@plt>
add     rsp, 0x10
xor     eax, eax
pop     rbx
ret
    
```

Registers

rax	<rax>
rbx	<rbx>
rcx	<rcx>
rdx	<rdx>
rdi	<rdi>
rsp	0x7ffff7830478
rip	main+0

```

main:
  push  rbx
  xor   eax, eax
  sub   rsp, 0x10
  cpuid
  mov   rdi, rsp
  mov   DWORD PTR [rsp], ebx
  mov   DWORD PTR [rsp+0x4], ecx
  mov   DWORD PTR [rsp+0x8], edx
  mov   DWORD PTR [rsp+0x10], 0x0
  call  401030 <puts@plt>
  add   rsp, 0x10
  xor   eax, eax
  pop   rbx
  ret
    
```

Registers

rax	<rax>
rbx	<rbx>
rcx	<rcx>
rdx	<rdx>
rdi	<rdi>
rsp	0x7ffff7830470
rip	main+1

```

main:
  push  rbx
  xor   eax, eax
  sub   rsp, 0x10
  cpuid
  mov   rdi, rsp
  mov   DWORD PTR [rsp], ebx
  mov   DWORD PTR [rsp+0x4], ecx
  mov   DWORD PTR [rsp+0x8], edx
  mov   DWORD PTR [rsp+0x10], 0x0
  call  401030 <puts@plt>
  add   rsp, 0x10
  xor   eax, eax
  pop   rbx
  ret
    
```

Registers

rax	0x0
rbx	<rbx>
rcx	<rcx>
rdx	<rdx>
rdi	<rdi>
rsp	0x7ffff7830470
rip	main+3

```

main:
  push  rbx
  xor   eax, eax
  sub   rsp, 0x10
  cpuid
  mov   rdi, rsp
  mov   DWORD PTR [rsp], ebx
  mov   DWORD PTR [rsp+0x4], ecx
  mov   DWORD PTR [rsp+0x8], edx
  mov   DWORD PTR [rsp+0x10], 0x0
  call  401030 <puts@plt>
  add   rsp, 0x10
  xor   eax, eax
  pop   rbx
  ret
    
```

Registers

rax	0x0
rbx	<rbx>
rcx	<rcx>
rdx	<rdx>
rdi	<rdi>
rsp	0x7ffff7830460
rip	main+7

```

main:
  push  rbx
  xor   eax, eax
  sub   rsp, 0x10
  cpuid
  mov   rdi, rsp
  mov   DWORD PTR [rsp], ebx
  mov   DWORD PTR [rsp+0x4], ecx
  mov   DWORD PTR [rsp+0x8], edx
  mov   DWORD PTR [rsp+0x10], 0x0
  call  401030 <puts@plt>
  add   rsp, 0x10
  xor   eax, eax
  pop   rbx
  ret
    
```



Registers

rax	0x0
rbx	<rbx>
rcx	<rcx>
rdx	<rdx>
rdi	<rdi>
rsp	0x7ffff7830460
rip	main+7

1. VM Exit → VMM takes control

Registers

rax	0x0
rbx	<rbx>
rcx	<rcx>
rdx	<rdx>
rdi	<rdi>
rsp	0x7ffff7830460
rip	main+7

1. VM Exit → VMM takes control
2. Check VM Exit reason → CPUID

Registers

rax	0x0
rbx	<rbx>
rcx	<rcx>
rdx	<rdx>
rdi	<rdi>
rsp	0x7ffff7830460
rip	main+7

1. VM Exit → VMM takes control
2. Check VM Exit reason → CPUID
3. Emulate CPUID

Registers

rax	0x0
rbx	<rbx>
rcx	<rcx>
rdx	<rdx>
rdi	<rdi>
rsp	0x7fff7830460
rip	main+7

1. VM Exit → VMM takes control
2. Check VM Exit reason → CPUID
3. Emulate CPUID
 - 3.1 Check rax → vendor string (rax = 0)

Registers

rax	0x0
rbx	<rbx>
rcx	<rcx>
rdx	<rdx>
rdi	<rdi>
rsp	0x7ffff7830460
rip	main+7

1. VM Exit → VMM takes control
2. Check VM Exit reason → CPUID
3. Emulate CPUID
 - 3.1 Check rax → vendor string (rax = 0)
 - 3.2 Set rbx, rcx, rdx to the vendor string

Registers

rax	0x0
rbx	<rbx>
rcx	<rcx>
rdx	<rdx>
rdi	<rdi>
rsp	0x7ffff7830460
rip	main+7

1. VM Exit → VMM takes control
2. Check VM Exit reason → CPUID
3. Emulate CPUID
 - 3.1 Check rax → vendor string (rax = 0)
 - 3.2 Set rbx, rcx, rdx to the vendor string

Registers

rax	0x0
rbx	"SWEB"
rcx	"SWEB"
rdx	"SWEB"
rdi	<rdi>
rsp	0x7fff7830460
rip	main+7

1. VM Exit → VMM takes control
2. Check VM Exit reason → CPUID
3. Emulate CPUID
 - 3.1 Check rax → vendor string (rax = 0)
 - 3.2 Set rbx, rcx, rdx to the vendor string
 - 3.3 Set rax to the maximum input value

Registers

rax	0x0
rbx	"SWEB"
rcx	"SWEB"
rdx	"SWEB"
rdi	<rdi>
rsp	0x7fff7830460
rip	main+7

1. VM Exit → VMM takes control
2. Check VM Exit reason → CPUID
3. Emulate CPUID
 - 3.1 Check rax → vendor string (rax = 0)
 - 3.2 Set rbx, rcx, rdx to the vendor string
 - 3.3 Set rax to the maximum input value

Registers

rax	0xd
rbx	"SWEB"
rcx	"SWEB"
rdx	"SWEB"
rdi	<rdi>
rsp	0x7fff7830460
rip	main+7

1. VM Exit → VMM takes control
2. Check VM Exit reason → CPUID
3. Emulate CPUID
 - 3.1 Check rax → vendor string (rax = 0)
 - 3.2 Set rbx, rcx, rdx to the vendor string
 - 3.3 Set rax to the maximum input value
4. VM Resume

Registers

rax	0xd
rbx	"SWEB"
rcx	"SWEB"
rdx	"SWEB"
rdi	<rdi>
rsp	0x7fff7830460
rip	main+7

```

main:
  push  rbx
  xor   eax, eax
  sub   rsp, 0x10
  cpuid
  mov   rdi, rsp
  mov   DWORD PTR [rsp], ebx
  mov   DWORD PTR [rsp+0x4], ecx
  mov   DWORD PTR [rsp+0x8], edx
  mov   DWORD PTR [rsp+0x10], 0x0
  call  401030 <puts@plt>
  add   rsp, 0x10
  xor   eax, eax
  pop   rbx
  ret
    
```

Registers

rax	0xd
rbx	"SWEB"
rcx	"SWEB"
rdx	"SWEB"
rdi	<rdi>
rsp	0x7ffff7830460
rip	main+7

```

main:
  push  rbx
  xor   eax, eax
  sub   rsp, 0x10
  cpuid
  mov   rdi, rsp
  mov   DWORD PTR [rsp], ebx
  mov   DWORD PTR [rsp+0x4], ecx
  mov   DWORD PTR [rsp+0x8], edx
  mov   DWORD PTR [rsp+0x10], 0x0
  call  401030 <puts@plt>
  add   rsp, 0x10
  xor   eax, eax
  pop   rbx
  ret
    
```



Registers

rax	0xd
rbx	"SWEB"
rcx	"SWEB"
rdx	"SWEB"
rdi	<rdi>
rsp	0x7ffff7830460
rip	main+7

1. VM Exit → VMM takes control
2. Check VM Exit reason → CPUID
3. Emulate CPUID
 - 3.1 Check rax → vendor string (rax = 0)
 - 3.2 Set rbx, rcx, rdx to the vendor string
 - 3.3 Set rax to the maximum input value
 - 3.4 Advance rip by the instruction length (field VM-exit instruction length)
4. VM Resume

Registers

rax	0xd
rbx	"SWEB"
rcx	"SWEB"
rdx	"SWEB"
rdi	<rdi>
rsp	0x7ffff7830460
rip	main+7

1. VM Exit → VMM takes control
2. Check VM Exit reason → CPUID
3. Emulate CPUID
 - 3.1 Check rax → vendor string (rax = 0)
 - 3.2 Set rbx, rcx, rdx to the vendor string
 - 3.3 Set rax to the maximum input value
 - 3.4 Advance rip by the instruction length (field VM-exit instruction length)
4. VM Resume

Registers

rax	0xd
rbx	"SWEB"
rcx	"SWEB"
rdx	"SWEB"
rdi	<rdi>
rsp	0x7ffff7830460
rip	main+9

```

main:
  push  rbx
  xor    eax, eax
  sub    rsp, 0x10
  cpuid
  mov    rdi, rsp
  mov    DWORD PTR [rsp], ebx
  mov    DWORD PTR [rsp+0x4], ecx
  mov    DWORD PTR [rsp+0x8], edx
  mov    DWORD PTR [rsp+0x10], 0x0
  call   401030 <puts@plt>
  add    rsp, 0x10
  xor    eax, eax
  pop    rbx
  ret
    
```

Registers

rax	0xd
rbx	"SWEB"
rcx	"SWEB"
rdx	"SWEB"
rdi	<rdi>
rsp	0x7ffff7830460
rip	main+9

```

main:
  push  rbx
  xor   eax, eax
  sub   rsp, 0x10
  cpuid
  mov   rdi, rsp
  mov   DWORD PTR [rsp], ebx
  mov   DWORD PTR [rsp+0x4], ecx
  mov   DWORD PTR [rsp+0x8], edx
  mov   DWORD PTR [rsp+0x10], 0x0
  call  401030 <puts@plt>
  add   rsp, 0x10
  xor   eax, eax
  pop   rbx
  ret
    
```

Registers

rax	0xd
rbx	"SWEB"
rcx	"SWEB"
rdx	"SWEB"
rdi	0x7ffff7830460
rsp	0x7ffff7830460
rip	main+12

```

main:
  push  rbx
  xor   eax, eax
  sub   rsp, 0x10
  cpuid
  mov   rdi, rsp
  mov   DWORD PTR [rsp], ebx
  mov   DWORD PTR [rsp+0x4], ecx
  mov   DWORD PTR [rsp+0x8], edx
  mov   DWORD PTR [rsp+0x10], 0x0
  call  401030 <puts@plt>
  add   rsp, 0x10
  xor   eax, eax
  pop   rbx
  ret
    
```

Registers

rax	0xd
rbx	"SWEB"
rcx	"SWEB"
rdx	"SWEB"
rdi	0x7ffff7830460
rsp	0x7ffff7830460
rip	main+16

```

main:
  push  rbx
  xor   eax, eax
  sub   rsp, 0x10
  cpuid
  mov   rdi, rsp
  mov   DWORD PTR [rsp], ebx
  mov   DWORD PTR [rsp+0x4], ecx
  mov   DWORD PTR [rsp+0x8], edx
  mov   DWORD PTR [rsp+0x10], 0x0
  call  401030 <puts@plt>
  add   rsp, 0x10
  xor   eax, eax
  pop   rbx
  ret
    
```

Registers

rax	0xd
rbx	"SWEB"
rcx	"SWEB"
rdx	"SWEB"
rdi	0x7ffff7830460
rsp	0x7ffff7830460
rip	main+21

```

main:
  push  rbx
  xor    eax, eax
  sub   rsp, 0x10
  cpuid
  mov   rdi, rsp
  mov   DWORD PTR [rsp], ebx
  mov   DWORD PTR [rsp+0x4], ecx
  mov   DWORD PTR [rsp+0x8], edx
  mov   DWORD PTR [rsp+0x10], 0x0
  call  401030 <puts@plt>
  add   rsp, 0x10
  xor   eax, eax
  pop   rbx
  ret
    
```

Registers

rax	0xd
rbx	"SWEB"
rcx	"SWEB"
rdx	"SWEB"
rdi	0x7ffff7830460
rsp	0x7ffff7830460
rip	main+26

```

main:
  push  rbx
  xor   eax, eax
  sub   rsp, 0x10
  cpuid
  mov   rdi, rsp
  mov   DWORD PTR [rsp], ebx
  mov   DWORD PTR [rsp+0x4], ecx
  mov   DWORD PTR [rsp+0x8], edx
  mov   DWORD PTR [rsp+0x10], 0x0
  call  401030 <puts@plt>
  add   rsp, 0x10
  xor   eax, eax
  pop   rbx
  ret
    
```

Registers

rax	0xd
rbx	"SWEB"
rcx	"SWEB"
rdx	"SWEB"
rdi	0x7ffff7830460
rsp	0x7ffff7830460
rip	main+34


```

main:
  push  rbx
  xor   eax, eax
  sub   rsp, 0x10
  cpuid
  mov   rdi, rsp
  mov   DWORD PTR [rsp], ebx
  mov   DWORD PTR [rsp+0x4], ecx
  mov   DWORD PTR [rsp+0x8], edx
  mov   DWORD PTR [rsp+0x10], 0x0
  call  401030 <puts@plt>
  add   rsp, 0x10
  xor   eax, eax
  pop   rbx
  ret
    
```

Registers

rax	0xc
rbx	"SWEB"
rcx	<rcx>
rdx	<rdx>
rdi	<rdi>
rsp	0x7ffff7830460
rip	main+39

```

main:
  push  rbx
  xor   eax, eax
  sub   rsp, 0x10
  cpuid
  mov   rdi, rsp
  mov   DWORD PTR [rsp], ebx
  mov   DWORD PTR [rsp+0x4], ecx
  mov   DWORD PTR [rsp+0x8], edx
  mov   DWORD PTR [rsp+0x10], 0x0
  call  401030 <puts@plt>
  add   rsp, 0x10
  xor   eax, eax
  pop   rbx
  ret
    
```

Registers

rax	0xc
rbx	"SWEB"
rcx	<rcx>
rdx	<rdx>
rdi	<rdi>
rsp	0x7ffff7830470
rip	main+43

```

main:
  push  rbx
  xor    eax, eax
  sub   rsp, 0x10
  cpuid
  mov   rdi, rsp
  mov   DWORD PTR [rsp], ebx
  mov   DWORD PTR [rsp+0x4], ecx
  mov   DWORD PTR [rsp+0x8], edx
  mov   DWORD PTR [rsp+0x10], 0x0
  call  401030 <puts@plt>
  add   rsp, 0x10
  xor   eax, eax
  pop   rbx
  ret
    
```

Registers

rax	0x0
rbx	"SWEB"
rcx	<rcx>
rdx	<rdx>
rdi	<rdi>
rsp	0x7ffff7830470
rip	main+45

```

main:
  push  rbx
  xor   eax, eax
  sub   rsp, 0x10
  cpuid
  mov   rdi, rsp
  mov   DWORD PTR [rsp], ebx
  mov   DWORD PTR [rsp+0x4], ecx
  mov   DWORD PTR [rsp+0x8], edx
  mov   DWORD PTR [rsp+0x10], 0x0
  call  401030 <puts@plt>
  add   rsp, 0x10
  xor   eax, eax
  pop   rbx
  ret
    
```

Registers

rax	0x0
rbx	<rbx>
rcx	<rcx>
rdx	<rdx>
rdi	<rdi>
rsp	0x7ffff7830460
rip	main+46







- The only purpose of this instruction is to trigger a VM exit



- The only purpose of this instruction is to trigger a VM exit
- Can be used by the guest to request hypervisor assistance



- The only purpose of this instruction is to trigger a VM exit
- Can be used by the guest to request hypervisor assistance
- Calling convention is hypervisor specific

```
main:
  push  rbx
  xor   eax, eax
  sub   rsp, 0x10
  vmcall
  mov   rdi, rsp
  mov   DWORD PTR [rsp], ebx
  mov   DWORD PTR [rsp+0x4], ecx
  mov   DWORD PTR [rsp+0x8], edx
  mov   DWORD PTR [rsp+0x10], 0x0
  call  401030 <puts@plt>
  add   rsp, 0x10
  xor   eax, eax
  pop   rbx
  ret
```

Registers

rax	<rax>
rbx	<rbx>
rcx	<rcx>
rdx	<rdx>
rdi	<rdi>
rsp	0x7ffff7830478
rip	main+0

```
main:
```

```
push rbx  
xor eax, eax  
sub rsp, 0x10  
vmcall  
mov rdi, rsp  
mov DWORD PTR [rsp], ebx  
mov DWORD PTR [rsp+0x4], ecx  
mov DWORD PTR [rsp+0x8], edx  
mov DWORD PTR [rsp+0x10], 0x0  
call 401030 <puts@plt>  
add rsp, 0x10  
xor eax, eax  
pop rbx  
ret
```

Registers

rax	<rax>
rbx	<rbx>
rcx	<rcx>
rdx	<rdx>
rdi	<rdi>
rsp	0x7ffff7830478
rip	main+0

```
main:
  push  rbx
  xor   eax, eax
  sub   rsp, 0x10
  vmcall
  mov   rdi, rsp
  mov   DWORD PTR [rsp], ebx
  mov   DWORD PTR [rsp+0x4], ecx
  mov   DWORD PTR [rsp+0x8], edx
  mov   DWORD PTR [rsp+0x10], 0x0
  call  401030 <puts@plt>
  add   rsp, 0x10
  xor   eax, eax
  pop   rbx
  ret
```

Registers

rax	<rax>
rbx	<rbx>
rcx	<rcx>
rdx	<rdx>
rdi	<rdi>
rsp	0x7ffff7830470
rip	main+1

```
main:
  push  rbx
  xor   eax, eax
  sub   rsp, 0x10
  vmcall
  mov   rdi, rsp
  mov   DWORD PTR [rsp], ebx
  mov   DWORD PTR [rsp+0x4], ecx
  mov   DWORD PTR [rsp+0x8], edx
  mov   DWORD PTR [rsp+0x10], 0x0
  call  401030 <puts@plt>
  add   rsp, 0x10
  xor   eax, eax
  pop   rbx
  ret
```

Registers

rax	0x0
rbx	<rbx>
rcx	<rcx>
rdx	<rdx>
rdi	<rdi>
rsp	0x7ffff7830470
rip	main+3

```
main:
  push  rbx
  xor    eax, eax
  sub   rsp, 0x10
  vmcall
  mov   rdi, rsp
  mov   DWORD PTR [rsp], ebx
  mov   DWORD PTR [rsp+0x4], ecx
  mov   DWORD PTR [rsp+0x8], edx
  mov   DWORD PTR [rsp+0x10], 0x0
  call  401030 <puts@plt>
  add   rsp, 0x10
  xor   eax, eax
  pop   rbx
  ret
```

Registers

rax	0x0
rbx	<rbx>
rcx	<rcx>
rdx	<rdx>
rdi	<rdi>
rsp	0x7ffff7830460
rip	main+7

```
main:
  push  rbx
  xor   eax, eax
  sub   rsp, 0x10
  vmcall
  mov   rdi, rsp
  mov   DWORD PTR [rsp], ebx
  mov   DWORD PTR [rsp+0x4], ecx
  mov   DWORD PTR [rsp+0x8], edx
  mov   DWORD PTR [rsp+0x10], 0x0
  call  401030 <puts@plt>
  add   rsp, 0x10
  xor   eax, eax
  pop   rbx
  ret
```



Registers

rax	0x0
rbx	<rbx>
rcx	<rcx>
rdx	<rdx>
rdi	<rdi>
rsp	0x7ffff7830460
rip	main+7

1. VM Exit → VMM takes control

Registers

rax	0x0
rbx	<rbx>
rcx	<rcx>
rdx	<rdx>
rdi	<rdi>
rsp	0x7ffff7830460
rip	main+7

1. VM Exit → VMM takes control
2. Check VM Exit reason → VM call

Registers

rax	0x0
rbx	<rbx>
rcx	<rcx>
rdx	<rdx>
rdi	<rdi>
rsp	0x7fff7830460
rip	main+7

1. VM Exit → VMM takes control
2. Check VM Exit reason → VM call
3. Handle VM call

Registers

rax	0x0
rbx	<rbx>
rcx	<rcx>
rdx	<rdx>
rdi	<rdi>
rsp	0x7fff7830460
rip	main+7

1. VM Exit → VMM takes control
2. Check VM Exit reason → VM call
3. Handle VM call
 - 3.1 Check rax → vendor string (rax = 0)

Registers

rax	0x0
rbx	<rbx>
rcx	<rcx>
rdx	<rdx>
rdi	<rdi>
rsp	0x7fff7830460
rip	main+7

1. VM Exit → VMM takes control
2. Check VM Exit reason → VM call
3. Handle VM call
 - 3.1 Check rax → vendor string (rax = 0)
 - 3.2 Set rbx, rcx, rdx to the vendor string

Registers

rax	0x0
rbx	<rbx>
rcx	<rcx>
rdx	<rdx>
rdi	<rdi>
rsp	0x7ffff7830460
rip	main+7

1. VM Exit → VMM takes control
2. Check VM Exit reason → VM call
3. Handle VM call
 - 3.1 Check rax → vendor string (rax = 0)
 - 3.2 Set rbx, rcx, rdx to the vendor string

Registers

rax	0x0
rbx	"SWEB"
rcx	"SWEB"
rdx	"SWEB"
rdi	<rdi>
rsp	0x7fff7830460
rip	main+7

1. VM Exit → VMM takes control
2. Check VM Exit reason → VM call
3. Handle VM call
 - 3.1 Check rax → vendor string (rax = 0)
 - 3.2 Set rbx, rcx, rdx to the vendor string
 - 3.3 Set rax to the maximum input value

Registers

rax	0x0
rbx	"SWEB"
rcx	"SWEB"
rdx	"SWEB"
rdi	<rdi>
rsp	0x7ffff7830460
rip	main+7

1. VM Exit → VMM takes control
2. Check VM Exit reason → VM call
3. Handle VM call
 - 3.1 Check rax → vendor string (rax = 0)
 - 3.2 Set rbx, rcx, rdx to the vendor string
 - 3.3 Set rax to the maximum input value

Registers

rax	0xd
rbx	"SWEB"
rcx	"SWEB"
rdx	"SWEB"
rdi	<rdi>
rsp	0x7ffff7830460
rip	main+7

1. VM Exit → VMM takes control
2. Check VM Exit reason → VM call
3. Handle VM call
 - 3.1 Check rax → vendor string (rax = 0)
 - 3.2 Set rbx, rcx, rdx to the vendor string
 - 3.3 Set rax to the maximum input value
 - 3.4 Advance rip by the instruction length (field VM-exit instruction length)

Registers

rax	0xd
rbx	"SWEB"
rcx	"SWEB"
rdx	"SWEB"
rdi	<rdi>
rsp	0x7ffff7830460
rip	main+7

1. VM Exit → VMM takes control
2. Check VM Exit reason → VM call
3. Handle VM call
 - 3.1 Check rax → vendor string (rax = 0)
 - 3.2 Set rbx, rcx, rdx to the vendor string
 - 3.3 Set rax to the maximum input value
 - 3.4 Advance rip by the instruction length (field VM-exit instruction length)

Registers

rax	0xd
rbx	"SWEB"
rcx	"SWEB"
rdx	"SWEB"
rdi	<rdi>
rsp	0x7ffff7830460
rip	main+10

1. VM Exit → VMM takes control
2. Check VM Exit reason → VM call
3. Handle VM call
 - 3.1 Check rax → vendor string (rax = 0)
 - 3.2 Set rbx, rcx, rdx to the vendor string
 - 3.3 Set rax to the maximum input value
 - 3.4 Advance rip by the instruction length (field VM-exit instruction length)
4. VM Resume

Registers

rax	0xd
rbx	"SWEB"
rcx	"SWEB"
rdx	"SWEB"
rdi	<rdi>
rsp	0x7ffff7830460
rip	main+10

```
main:
  push  rbx
  xor    eax, eax
  sub   rsp, 0x10
  vmcall
  mov   rdi, rsp
  mov   DWORD PTR [rsp], ebx
  mov   DWORD PTR [rsp+0x4], ecx
  mov   DWORD PTR [rsp+0x8], edx
  mov   DWORD PTR [rsp+0x10], 0x0
  call  401030 <puts@plt>
  add   rsp, 0x10
  xor   eax, eax
  pop   rbx
  ret
```

Registers

rax	0xd
rbx	"SWEB"
rcx	"SWEB"
rdx	"SWEB"
rdi	<rdi>
rsp	0x7ffff7830460
rip	main+10

```
main:
  push  rbx
  xor   eax, eax
  sub   rsp, 0x10
  vmcall
  mov   rdi, rsp
  mov   DWORD PTR [rsp], ebx
  mov   DWORD PTR [rsp+0x4], ecx
  mov   DWORD PTR [rsp+0x8], edx
  mov   DWORD PTR [rsp+0x10], 0x0
  call  401030 <puts@plt>
  add   rsp, 0x10
  xor   eax, eax
  pop   rbx
  ret
```

Registers

rax	0xd
rbx	"SWEB"
rcx	"SWEB"
rdx	"SWEB"
rdi	<rdi>
rsp	0x7ffff7830460
rip	main+10

```
main:
  push  rbx
  xor    eax, eax
  sub   rsp, 0x10
  vmcall
  mov   rdi, rsp
  mov   DWORD PTR [rsp], ebx
  mov   DWORD PTR [rsp+0x4], ecx
  mov   DWORD PTR [rsp+0x8], edx
  mov   DWORD PTR [rsp+0x10], 0x0
  call  401030 <puts@plt>
  add   rsp, 0x10
  xor   eax, eax
  pop   rbx
  ret
```

Registers

rax	0xd
rbx	"SWEB"
rcx	"SWEB"
rdx	"SWEB"
rdi	0x7ffff7830460
rsp	0x7ffff7830460
rip	main+13

```
main:
  push  rbx
  xor   eax, eax
  sub   rsp, 0x10
  vmcall
  mov   rdi, rsp
  mov   DWORD PTR [rsp], ebx
  mov   DWORD PTR [rsp+0x4], ecx
  mov   DWORD PTR [rsp+0x8], edx
  mov   DWORD PTR [rsp+0x10], 0x0
  call  401030 <puts@plt>
  add   rsp, 0x10
  xor   eax, eax
  pop   rbx
  ret
```

Registers

rax	0xd
rbx	"SWEB"
rcx	"SWEB"
rdx	"SWEB"
rdi	0x7ffff7830460
rsp	0x7ffff7830460
rip	main+17

```
main:
  push  rbx
  xor   eax, eax
  sub   rsp, 0x10
  vmcall
  mov   rdi, rsp
  mov   DWORD PTR [rsp], ebx
  mov   DWORD PTR [rsp+0x4], ecx
  mov   DWORD PTR [rsp+0x8], edx
  mov   DWORD PTR [rsp+0x10], 0x0
  call  401030 <puts@plt>
  add   rsp, 0x10
  xor   eax, eax
  pop   rbx
  ret
```

Registers

rax	0xd
rbx	"SWEB"
rcx	"SWEB"
rdx	"SWEB"
rdi	0x7ffff7830460
rsp	0x7ffff7830460
rip	main+22


```
main:
  push  rbx
  xor    eax, eax
  sub   rsp, 0x10
  vmcall
  mov   rdi, rsp
  mov   DWORD PTR [rsp], ebx
  mov   DWORD PTR [rsp+0x4], ecx
  mov   DWORD PTR [rsp+0x8], edx
  mov   DWORD PTR [rsp+0x10], 0x0
  call  401030 <puts@plt>
  add   rsp, 0x10
  xor   eax, eax
  pop   rbx
  ret
```

Registers

rax	0xd
rbx	"SWEB"
rcx	"SWEB"
rdx	"SWEB"
rdi	0x7ffff7830460
rsp	0x7ffff7830460
rip	main+27

```
main:
  push  rbx
  xor    eax, eax
  sub   rsp, 0x10
  vmcall
  mov   rdi, rsp
  mov   DWORD PTR [rsp], ebx
  mov   DWORD PTR [rsp+0x4], ecx
  mov   DWORD PTR [rsp+0x8], edx
  mov   DWORD PTR [rsp+0x10], 0x0
  call  401030 <puts@plt>
  add   rsp, 0x10
  xor   eax, eax
  pop   rbx
  ret
```

Registers

rax	0xd
rbx	"SWEB"
rcx	"SWEB"
rdx	"SWEB"
rdi	0x7ffff7830460
rsp	0x7ffff7830460
rip	main+35

```
main:
  push  rbx
  xor   eax, eax
  sub   rsp, 0x10
  vmcall
  mov   rdi, rsp
  mov   DWORD PTR [rsp], ebx
  mov   DWORD PTR [rsp+0x4], ecx
  mov   DWORD PTR [rsp+0x8], edx
  mov   DWORD PTR [rsp+0x10], 0x0
  call  401030 <puts@plt>
  add   rsp, 0x10
  xor   eax, eax
  pop   rbx
  ret
```

Registers

rax	0xc
rbx	"SWEB"
rcx	<rcx>
rdx	<rdx>
rdi	<rdi>
rsp	0x7ffff7830460
rip	main+40

```
main:
  push  rbx
  xor   eax, eax
  sub   rsp, 0x10
  vmcall
  mov   rdi, rsp
  mov   DWORD PTR [rsp], ebx
  mov   DWORD PTR [rsp+0x4], ecx
  mov   DWORD PTR [rsp+0x8], edx
  mov   DWORD PTR [rsp+0x10], 0x0
  call  401030 <puts@plt>
  add   rsp, 0x10
  xor   eax, eax
  pop   rbx
  ret
```

Registers

rax	0xc
rbx	"SWEB"
rcx	<rcx>
rdx	<rdx>
rdi	<rdi>
rsp	0x7ffff7830470
rip	main+44

```
main:
  push  rbx
  xor    eax, eax
  sub   rsp, 0x10
  vmcall
  mov   rdi, rsp
  mov   DWORD PTR [rsp], ebx
  mov   DWORD PTR [rsp+0x4], ecx
  mov   DWORD PTR [rsp+0x8], edx
  mov   DWORD PTR [rsp+0x10], 0x0
  call  401030 <puts@plt>
  add   rsp, 0x10
  xor   eax, eax
  pop   rbx
  ret
```

Registers

rax	0x0
rbx	"SWEB"
rcx	<rcx>
rdx	<rdx>
rdi	<rdi>
rsp	0x7ffff7830470
rip	main+46

```
main:
  push  rbx
  xor    eax, eax
  sub   rsp, 0x10
  vmcall
  mov   rdi, rsp
  mov   DWORD PTR [rsp], ebx
  mov   DWORD PTR [rsp+0x4], ecx
  mov   DWORD PTR [rsp+0x8], edx
  mov   DWORD PTR [rsp+0x10], 0x0
  call  401030 <puts@plt>
  add   rsp, 0x10
  xor   eax, eax
  pop   rbx
  ret
```

Registers

rax	0x0
rbx	<rbx>
rcx	<rcx>
rdx	<rdx>
rdi	<rdi>
rsp	0x7ffff7830460
rip	main+47

- We can control I/O and MSR accesses

- We can control I/O and MSR accesses
- Bitmaps define which I/O ports and MSRs cause VM exits

- We can control I/O and MSR accesses
- Bitmaps define which I/O ports and MSRs cause VM exits
- Giving the guest direct access to external hardware can be dangerous

- We can control I/O and MSR accesses
- Bitmaps define which I/O ports and MSRs cause VM exits
- Giving the guest direct access to external hardware can be dangerous
 - Guests could mess up the host state

- We can control I/O and MSR accesses
- Bitmaps define which I/O ports and MSRs cause VM exits
- Giving the guest direct access to external hardware can be dangerous
 - Guests could mess up the host state
 - DMAs don't care about our EPT

- We can control I/O and MSR accesses
 - Bitmaps define which I/O ports and MSRs cause VM exits
 - Giving the guest direct access to external hardware can be dangerous
 - Guests could mess up the host state
 - DMAs don't care about our EPT
- Disable access or use bitmaps to limit the ports/MSRS

VIOLATION

VIOLATION

- Same as page fault just for EPT translations

VIOLATION

- Same as page fault just for EPT translations
- Guest memory accesses that violate the permissions set in the EPT cause VM exits

VIOLATION

VIOLATION

- Same as page fault just for EPT translations
- Guest memory accesses that violate the permissions set in the EPT cause VM exits
- The VMM can handle EPT violations accordingly

VIOLATION

- Same as page fault just for EPT translations
- Guest memory accesses that violate the permissions set in the EPT cause VM exits
- The VMM can handle EPT violations accordingly
- This gives us a lot of room to play around with:

VIOLATION

- Same as page fault just for EPT translations
- Guest memory accesses that violate the permissions set in the EPT cause VM exits
- The VMM can handle EPT violations accordingly
- This gives us a lot of room to play around with:
 - On-demand paging

VIOLATION

- Same as page fault just for EPT translations
- Guest memory accesses that violate the permissions set in the EPT cause VM exits
- The VMM can handle EPT violations accordingly
- This gives us a lot of room to play around with:
 - On-demand paging
 - Copy-on-Write

VIOLATION

- Same as page fault just for EPT translations
- Guest memory accesses that violate the permissions set in the EPT cause VM exits
- The VMM can handle EPT violations accordingly
- This gives us a lot of room to play around with:
 - On-demand paging
 - Copy-on-Write
 - Deduplication

VIOLATION

- Same as page fault just for EPT translations
- Guest memory accesses that violate the permissions set in the EPT cause VM exits
- The VMM can handle EPT violations accordingly
- This gives us a lot of room to play around with:
 - On-demand paging
 - Copy-on-Write
 - Deduplication
 - EPT Hooking

VIOLATION

- Same as page fault just for EPT translations
- Guest memory accesses that violate the permissions set in the EPT cause VM exits
- The VMM can handle EPT violations accordingly
- This gives us a lot of room to play around with:
 - On-demand paging
 - Copy-on-Write
 - Deduplication
 - EPT Hooking
 - ...







- Raised by external devices (PIT, disk, network card, ...)



- Raised by external devices (PIT, disk, network card, ...)
- External devices continue operating while a guest is scheduled



- Raised by external devices (PIT, disk, network card, ...)
 - External devices continue operating while a guest is scheduled
- Interrupts for the host might occur while the guest is scheduled



- Raised by external devices (PIT, disk, network card, ...)
 - External devices continue operating while a guest is scheduled
- Interrupts for the host might occur while the guest is scheduled
- VMCS can be configured for a VM exit on external interrupts



- Raised by external devices (PIT, disk, network card, ...)
 - External devices continue operating while a guest is scheduled
- Interrupts for the host might occur while the guest is scheduled
- VMCS can be configured for a VM exit on external interrupts
 - Once the host sets $IF=1$ the external interrupt is delivered normally



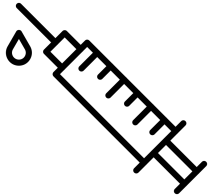


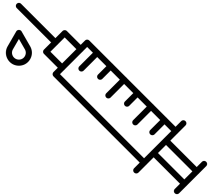


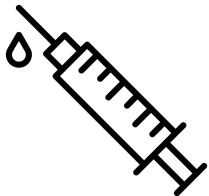
- We can force a VM exit as soon as the guest is able to receive interrupts (Guest RFLAGS.IF = 1)



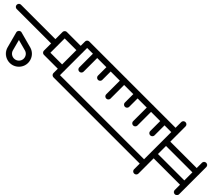
- We can force a VM exit as soon as the guest is able to receive interrupts (Guest RFLAGS.IF = 1)
- Very useful for injecting interrupts



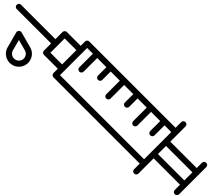




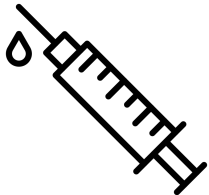
- The host can inject events into the guest



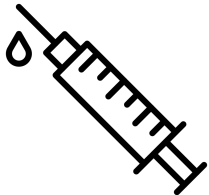
- The host can inject events into the guest
 - External interrupts



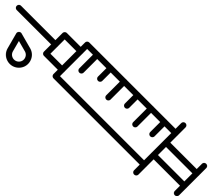
- The host can inject events into the guest
 - External interrupts
 - Non-maskable interrupts



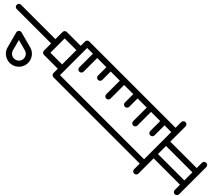
- The host can inject events into the guest
 - External interrupts
 - Non-maskable interrupts
 - Exceptions



- The host can inject events into the guest
 - External interrupts
 - Non-maskable interrupts
 - Exceptions
 - Software interrupts



- The host can inject events into the guest
 - External interrupts
 - Non-maskable interrupts
 - Exceptions
 - Software interrupts
 - ...



- The host can inject events into the guest
 - External interrupts
 - Non-maskable interrupts
 - Exceptions
 - Software interrupts
 - ...
- Very useful for device emulation

Table 25-16. Format of the VM-Entry Interruption-Information Field

Bit Position(s)	Content
7:0	Vector of interrupt or exception
10:8	Interruption type: 0: External interrupt 1: Reserved 2: Non-maskable interrupt (NMI) 3: Hardware exception (e.g., #PF) 4: Software interrupt (INT <i>n</i>) 5: Privileged software exception (INT1) 6: Software exception (INT3 or INTO) 7: Other event
11	Deliver error code (0 = do not deliver; 1 = deliver)
30:12	Reserved
31	Valid







- Set interrupt information in VM-entry interruption-information



- Set interrupt information in VM-entry interruption-information
- Optional: set error code in VM-entry exception error code



- Set interrupt information in VM-entry interruption-information
- Optional: set error code in VM-entry exception error code
- For software interrupts: set instruction length in VM-entry instruction length



- Set interrupt information in VM-entry interruption-information
- Optional: set error code in VM-entry exception error code
- For software interrupts: set instruction length in VM-entry instruction length
- Injection of external interrupts is only allowed if guest RFLAGS.IF = 1



- Set interrupt information in VM-entry interruption-information
- Optional: set error code in VM-entry exception error code
- For software interrupts: set instruction length in VM-entry instruction length
- Injection of external interrupts is only allowed if guest $RFLAGS.IF = 1$
 - Interrupt-window exiting

1. Check pending interrupts
→ Timer, Keyboard

Int. Window Exiting	0
Int. Info.	0
rflags.IF	0
pending IRQs	Timer, Keyboard

1. Check pending interrupts
→ Timer, Keyboard
2. Check guest IF flag → 0

Int. Window Exiting	0
Int. Info.	0
rflags.IF	0
pending IRQs	Timer, Keyboard

1. Check pending interrupts
→ Timer, Keyboard
2. Check guest IF flag → 0
3. Set interrupt-window exiting

Int. Window Exiting	1
Int. Info.	0
rflags.IF	0
pending IRQs	Timer, Keyboard

1. Check pending interrupts
→ Timer, Keyboard
2. Check guest IF flag → 0
3. Set interrupt-window exiting
4. Continue guest

Int. Window Exiting	1
Int. Info.	0
rflags.IF	0
pending IRQs	Timer, Keyboard

```
push rbp
mov rbp, rsp
sti
nop
pop rbp
ret
```

Int. Window Exiting	1
Int. Info.	0
rflags.IF	0
pending IRQs	Timer, Keyboard

```
push    rbp
mov     rbp, rsp
sti
nop
pop     rbp
ret
```

Int. Window Exiting	1
Int. Info.	0
rflags.IF	0
pending IRQs	Timer, Keyboard

```
push    rbp
mov     rbp, rsp
sti
nop
pop     rbp
ret
```

Int. Window Exiting	1
Int. Info.	0
rflags.IF	0
pending IRQs	Timer, Keyboard

```
push    rbp
mov     rbp, rsp
sti
nop
pop     rbp
ret
```

Int. Window Exiting	1
Int. Info.	0
rflags.IF	1
pending IRQs	Timer, Keyboard


```
push rbp
mov rbp, rsp
sti
nop
pop rbp
ret
```



Int. Window Exiting	1
Int. Info.	0
rflags.IF	1
pending IRQs	Timer, Keyboard

1. VM Exit → VMM takes control

Int. Window Exiting	1
Int. Info.	0
rflags.IF	1
pending IRQs	Timer, Keyboard

1. VM Exit → VMM takes control
2. Check VM Exit reason
→ interrupt window

Int. Window Exiting	1
Int. Info.	0
rflags.IF	1
pending IRQs	Timer, Keyboard

1. VM Exit → VMM takes control
2. Check VM Exit reason
 - interrupt window
 - 2.1 Check pending interrupts
 - Timer, Keyboard

Int. Window Exiting	1
Int. Info.	0
rflags.IF	1
pending IRQs	Timer, Keyboard

1. VM Exit → VMM takes control
2. Check VM Exit reason
 - interrupt window
 - 2.1 Check pending interrupts
 - Timer, Keyboard
 - 2.2 Set interruption information

Int. Window Exiting	1
Int. Info.	0x80000020
rflags.IF	1
pending IRQs	Timer, Keyboard

1. VM Exit → VMM takes control
2. Check VM Exit reason
 - interrupt window
 - 2.1 Check pending interrupts
 - Timer, Keyboard
 - 2.2 Set interruption information
 - 2.3 Continue guest

Int. Window Exiting	1
Int. Info.	0x80000020
rflags.IF	1
pending IRQs	Keyboard

```
push  rbp
mov   rbp, rsp
sti
nop
pop   rbp
ret
```

Int. Window Exiting	1
Int. Info.	0x80000020
rflags.IF	1
pending IRQs	Keyboard

```
push rbp
mov rbp, rsp
sti
nop
pop rbp
ret
```



Int. Window Exiting	1
Int. Info.	0x80000020
rflags.IF	1
pending IRQs	Keyboard


```
<irq handler>
```

```
iret
```

Int. Window Exiting	1
Int. Info.	0x80000020
rflags.IF	0
pending IRQs	Keyboard

<irq handler>

iret

Int. Window Exiting	1
Int. Info.	0x80000020
rflags.IF	0
pending IRQs	Keyboard

```
push  rbp
mov   rbp, rsp
sti
nop
pop   rbp
ret
```

Int. Window Exiting	1
Int. Info.	0x80000020
rflags.IF	1
pending IRQs	Keyboard

```
push rbp
mov rbp, rsp
sti
nop
pop rbp
ret
```



Int. Window Exiting	1
Int. Info.	0x80000020
rflags.IF	1
pending IRQs	Keyboard

1. VM Exit → VMM takes control

Int. Window Exiting	1
Int. Info.	0x20
rflags.IF	1
pending IRQs	Keyboard

1. VM Exit → VMM takes control
2. Check VM Exit reason
→ interrupt window

Int. Window Exiting	1
Int. Info.	0x20
rflags.IF	1
pending IRQs	Keyboard

1. VM Exit → VMM takes control
2. Check VM Exit reason
 - interrupt window
 - 2.1 Check pending interrupts
 - Keyboard

Int. Window Exiting	1
Int. Info.	0x20
rflags.IF	1
pending IRQs	Keyboard

1. VM Exit → VMM takes control
2. Check VM Exit reason
 - interrupt window
 - 2.1 Check pending interrupts
 - Keyboard
 - 2.2 Set interruption information

Int. Window Exiting	1
Int. Info.	0x80000021
rflags.IF	1
pending IRQs	Keyboard

1. VM Exit → VMM takes control
2. Check VM Exit reason
 - interrupt window
 - 2.1 Check pending interrupts
 - Keyboard
 - 2.2 Set interruption information
 - 2.3 No other pending interrupts
 - clear interrupt-window exiting

Int. Window Exiting	0
Int. Info.	0x80000021
rflags.IF	1
pending IRQs	

1. VM Exit → VMM takes control
2. Check VM Exit reason
 - interrupt window
 - 2.1 Check pending interrupts
 - Keyboard
 - 2.2 Set interruption information
 - 2.3 No other pending interrupts
 - clear interrupt-window exiting
 - 2.4 Continue guest

Int. Window Exiting	0
Int. Info.	0x80000021
rflags.IF	1
pending IRQs	

```
push  rbp
mov   rbp, rsp
sti
nop
pop   rbp
ret
```

Int. Window Exiting	0
Int. Info.	0x80000021
rflags.IF	1
pending IRQs	

```
push  rbp
mov   rbp, rsp
sti
nop
pop   rbp
ret
```



Int. Window Exiting	0
Int. Info.	0x80000021
rflags.IF	1
pending IRQs	

```
<irq handler>
```

```
iret
```

Int. Window Exiting	0
Int. Info.	0x80000021
rflags.IF	0
pending IRQs	

<irq handler>

iret

Int. Window Exiting	0
Int. Info.	0x80000021
rflags.IF	0
pending IRQs	

```
push  rbp
mov   rbp, rsp
sti
nop
pop   rbp
ret
```

Int. Window Exiting	0
Int. Info.	0x80000021
rflags.IF	1
pending IRQs	

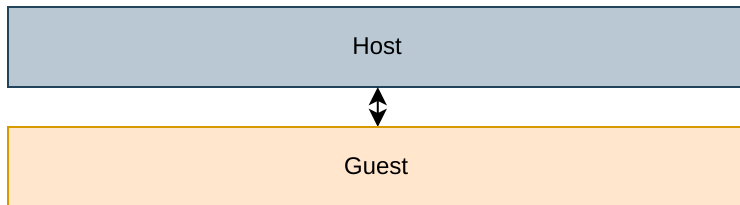
```
push  rbp
mov   rbp, rsp
sti
nop
pop   rbp
ret
```

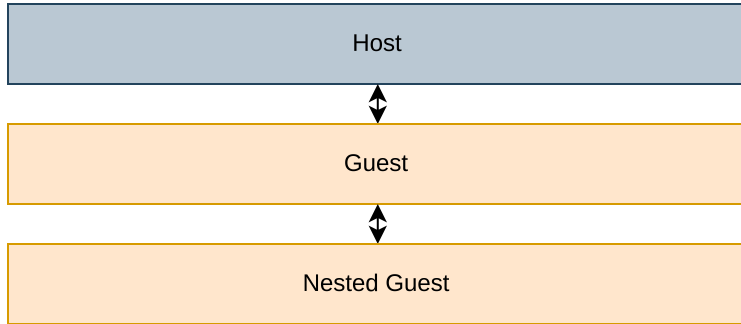
Int. Window Exiting	0
Int. Info.	0x80000021
rflags.IF	1
pending IRQs	

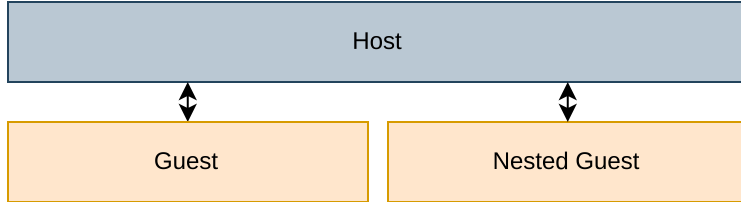

```
push rbp
mov rbp, rsp
sti
nop
pop rbp
ret
```

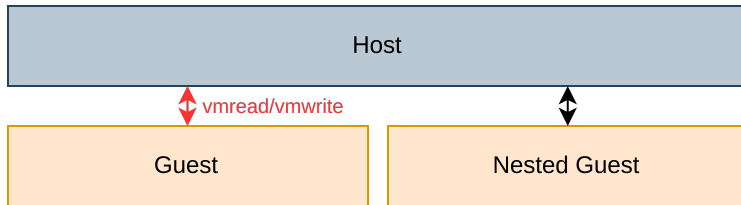
Int. Window Exiting	0
Int. Info.	0x80000021
rflags.IF	1
pending IRQs	

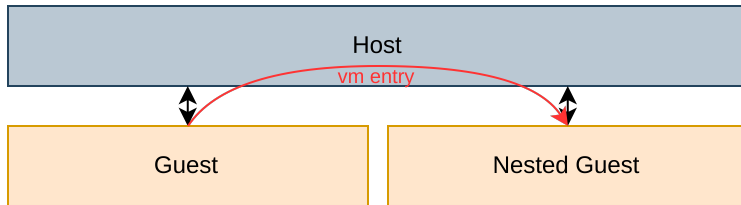
Nested Virtualization

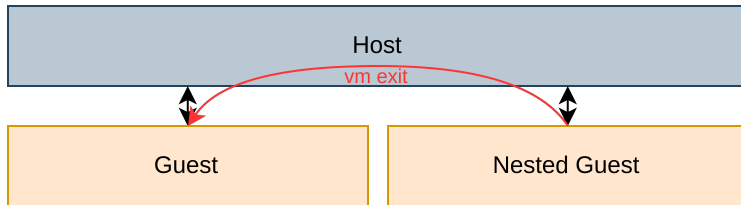


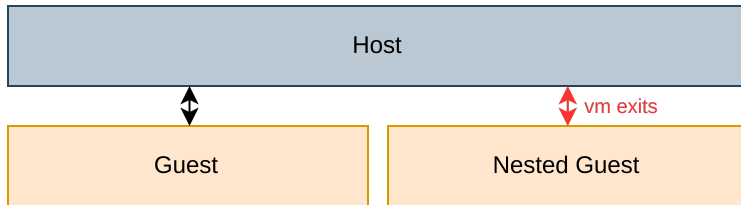
















- "Real" nested VMX would be too expensive





- "Real" nested VMX would be too expensive
- Run the nested guest as a normal guest and emulate VMX



- "Real" nested VMX would be too expensive
- Run the nested guest as a normal guest and emulate VMX
- vmread and vmwrite cause a VM exit



- "Real" nested VMX would be too expensive
- Run the nested guest as a normal guest and emulate VMX
- vmread and vmwrite cause a VM exit
 - Host maintains a real VMCS and configurations for the nested guests



- "Real" nested VMX would be too expensive
- Run the nested guest as a normal guest and emulate VMX
- vmread and vmwrite cause a VM exit
 - Host maintains a real VMCS and configurations for the nested guests
 - VM entry causes a VM exit to the host



- "Real" nested VMX would be too expensive
- Run the nested guest as a normal guest and emulate VMX
- vmread and vmwrite cause a VM exit
 - Host maintains a real VMCS and configurations for the nested guests
 - VM entry causes a VM exit to the host
- Host schedules nested guest



- "Real" nested VMX would be too expensive
- Run the nested guest as a normal guest and emulate VMX
- vmread and vmwrite cause a VM exit
- Host maintains a real VMCS and configurations for the nested guests
- VM entry causes a VM exit to the host
- Host schedules nested guest
- Nested guest VM exits can be forwarded to the guest or handled by the host





- Can be enabled in the VMCS





- Can be enabled in the VMCS
- VMCS link pointer contains the address of a shadow VMCS



- Can be enabled in the VMCS
- VMCS link pointer contains the address of a shadow VMCS
- Guest vmwrite and vmread modify the shadow VMCS



- Can be enabled in the VMCS
 - VMCS link pointer contains the address of a shadow VMCS
 - Guest vmwrite and vmread modify the shadow VMCS
- less VM exits



- Can be enabled in the VMCS
 - VMCS link pointer contains the address of a shadow VMCS
 - Guest vmwrite and vmread modify the shadow VMCS
- less VM exits
- Host can use vmread and vmwrite on the shadow VMCS



- Can be enabled in the VMCS
 - VMCS link pointer contains the address of a shadow VMCS
 - Guest vmwrite and vmread modify the shadow VMCS
- less VM exits
- Host can use vmread and vmwrite on the shadow VMCS
 - the shadow VMCS can not be launched



- Can be enabled in the VMCS
 - VMCS link pointer contains the address of a shadow VMCS
 - Guest vmwrite and vmread modify the shadow VMCS
- less VM exits
- Host can use vmread and vmwrite on the shadow VMCS
 - the shadow VMCS can not be launched
 - Host still has to maintain a normal VMCS for the nested



- Can be enabled in the VMCS
 - VMCS link pointer contains the address of a shadow VMCS
 - Guest vmwrite and vmread modify the shadow VMCS
- less VM exits
- Host can use vmread and vmwrite on the shadow VMCS
 - the shadow VMCS can not be launched
 - Host still has to maintain a normal VMCS for the nested
 - Bitmaps can be used to restrict the guest from accessing VMCS fields



- Can be enabled in the VMCS
 - VMCS link pointer contains the address of a shadow VMCS
 - Guest vmwrite and vmread modify the shadow VMCS
- less VM exits
- Host can use vmread and vmwrite on the shadow VMCS
 - the shadow VMCS can not be launched
 - Host still has to maintain a normal VMCS for the nested
 - Bitmaps can be used to restrict the guest from accessing VMCS fields
- Host does not always need to copy the whole shadow VMCS







- Up to 25 memory accesses per address dereference with EPT



- Up to 25 memory accesses per address dereference with EPT
- 125 if we add a third layer



- Up to 25 memory accesses per address dereference with EPT
 - 125 if we add a third layer
- use normal EPT for nested guests



- Up to 25 memory accesses per address dereference with EPT
 - 125 if we add a third layer
- use normal EPT for nested guests
- Guest maintains a copy of the real EPT



- Up to 25 memory accesses per address dereference with EPT
 - 125 if we add a third layer
- use normal EPT for nested guests
- Guest maintains a copy of the real EPT
 - Host maintains the real nested guest EPT





- Host has to...



- Host has to...
 - Update the real EPT if the shadow changes





- Host has to...
 - Update the real EPT if the shadow changes
 - Copy accessed and dirty bits from the real EPT to the shadow



- Host has to...
 - Update the real EPT if the shadow changes
 - Copy accessed and dirty bits from the real EPT to the shadow
- updating real EPT can be done on demand:



- Host has to...
 - Update the real EPT if the shadow changes
 - Copy accessed and dirty bits from the real EPT to the shadow
- updating real EPT can be done on demand:
 1. schedule nested guest with empty EPT



- Host has to...
 - Update the real EPT if the shadow changes
 - Copy accessed and dirty bits from the real EPT to the shadow
- updating real EPT can be done on demand:
 1. schedule nested guest with empty EPT
 2. on EPT violation check the shadow EPT



- Host has to...
 - Update the real EPT if the shadow changes
 - Copy accessed and dirty bits from the real EPT to the shadow
- updating real EPT can be done on demand:
 1. schedule nested guest with empty EPT
 2. on EPT violation check the shadow EPT
 3. map page if it exists, forward VM exit to guest if it does not



- Host has to...
 - Update the real EPT if the shadow changes
 - Copy accessed and dirty bits from the real EPT to the shadow
- updating real EPT can be done on demand:
 1. schedule nested guest with empty EPT
 2. on EPT violation check the shadow EPT
 3. map page if it exists, forward VM exit to guest if it does not
- set shadow EPT pages to read only to detect updates to existing mappings

Exercise







- Goal: write a hypervisor that runs SWEB inside of SWEB



- Goal: write a hypervisor that runs SWEB inside of SWEB
- Base requirements:



- Goal: write a hypervisor that runs SWEB inside of SWEB
- Base requirements:
 - The guest has to eventually get to the SWEB entry and boot normally
<Boot Code>



- Goal: write a hypervisor that runs SWEB inside of SWEB
- Base requirements:
 - The guest has to eventually get to the SWEB entry and boot normally <Boot Code>
 - Emulate video output, PIC, PIT and the keyboard



- Goal: write a hypervisor that runs SWEB inside of SWEB
- Base requirements:
 - The guest has to eventually get to the SWEB entry and boot normally <Boot Code>
 - Emulate video output, PIC, PIT and the keyboard
 - Forward the host VFSSyscalls to the guest via hypercalls



- Goal: write a hypervisor that runs SWEB inside of SWEB
- Base requirements:
 - The guest has to eventually get to the SWEB entry and boot normally <Boot Code>
 - Emulate video output, PIC, PIT and the keyboard
 - Forward the host VFSSyscalls to the guest via hypercalls
 - The guest should run normally and receive regular timer interrupts



- Goal: write a hypervisor that runs SWEB inside of SWEB
- Base requirements:
 - The guest has to eventually get to the SWEB entry and boot normally <Boot Code>
 - Emulate video output, PIC, PIT and the keyboard
 - Forward the host VFSSyscalls to the guest via hypercalls
 - The guest should run normally and receive regular timer interrupts
 - Running a guest should not block the host from running other processes







- Start the guest in protected mode



- Start the guest in protected mode
- Copy the currently running SWEB



- Start the guest in protected mode
- Copy the currently running SWEB
 - reset multiboot remainder and fd_num_



- Start the guest in protected mode
- Copy the currently running SWEB
 - reset multiboot remainder and fd_num_
- Checkout `parseMultibootHeader` for creating a new multiboot header

