# Cloud Operating Systems
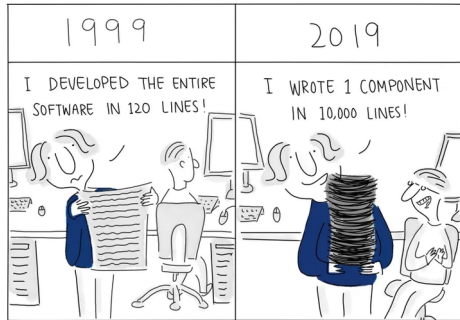
**Daniel Gruss**
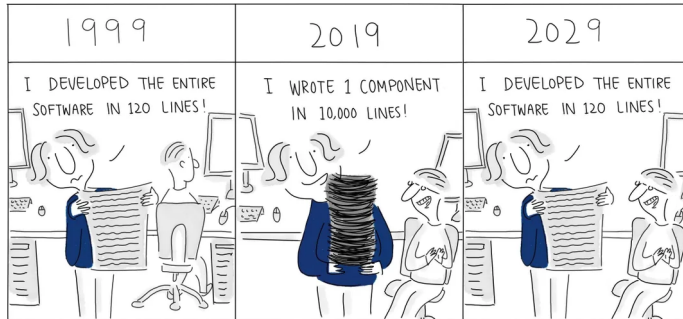
2023-03-01
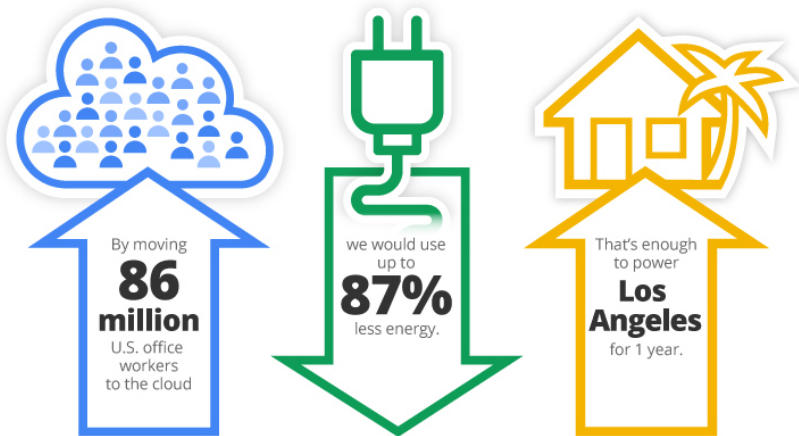
YESS...

MOAR SOFTWARE COMPLEXITY

# Moving to the cloud can save up to 87% of IT energy



By moving **86 million** U.S. office workers to the cloud

we would use up to **87%** less energy.

That's enough to power **Los Angeles** for 1 year.

# Cloud means Efficiency

Apply OS techniques - Example?

- Processes used to have access to all physical memory $\rightarrow$ that's not efficient!

Daniel Gruss

Apply OS techniques - Example?

- Processes used to have access to all physical memory $\rightarrow$ that's not efficient!
- $\rightarrow$ Virtualize memory $\rightarrow$ processes can share resources of one machine and utilize it better

Daniel Gruss

Apply OS techniques - Example?

- Processes used to have access to all physical memory $\rightarrow$ that's not efficient!
- $\rightarrow$ Virtualize memory $\rightarrow$ processes can share resources of one machine and utilize it better
- Processes need all the same pages $\rightarrow$ that's not efficient!

Apply OS techniques - Example?

- Processes used to have access to all physical memory $\rightarrow$ that's not efficient!
- $\rightarrow$ Virtualize memory $\rightarrow$ processes can share resources of one machine and utilize it better
- Processes need all the same pages $\rightarrow$ that's not efficient!
- $\rightarrow$ Let them share memory, using COW, page deduplication, etc.

Apply OS techniques - Example?

- Processes used to have access to all physical memory → that's not efficient!
- → Virtualize memory → processes can share resources of one machine and utilize it better
- Processes need all the same pages → that's not efficient!
- → Let them share memory, using COW, page deduplication, etc.
- Processes often cannot do anything but wait → that's not efficient!

Apply OS techniques - Example?

- Processes used to have access to all physical memory $\rightarrow$ that's not efficient!
- $\rightarrow$ Virtualize memory $\rightarrow$ processes can share resources of one machine and utilize it better
- Processes need all the same pages $\rightarrow$ that's not efficient!
- $\rightarrow$ Let them share memory, using COW, page deduplication, etc.
- Processes often cannot do anything but wait $\rightarrow$ that's not efficient!
- $\rightarrow$ Let other processes run in between

- Efficiency

Daniel Gruss

- Efficiency
- Isolation of tenants (security, reliability, availability)

- Efficiency
- Isolation of tenants (security, reliability, availability)
- Abstraction of hardware

Daniel Gruss

**Virtualization** allows to represent resources in a computer in a way they can be used easily and without the need to know details of their properties

Daniel Gruss

- Decouple operating system from hardware

- Decouple operating system from hardware
  - "computer in computer" - implemented in software

Daniel Gruss

- Decouple operating system from hardware
  - "computer in computer" - implemented in software
  - includes devices (network, keyboard, sound…)

- Decouple operating system from hardware
  - "computer in computer" - implemented in software
  - includes devices (network, keyboard, sound…)
- OS in VM "sees" its hardware, irrespective from the actual hardware in use

- Decouple operating system from hardware
  - "computer in computer" - implemented in software
  - includes devices (network, keyboard, sound…)
- OS in VM "sees" its hardware, irrespective from the actual hardware in use
- OS does not know if HW is concurrently used by other VMS

Daniel Gruss

- Cheaper hardware: one server for one task was common

- Cheaper hardware: one server for one task was common
- most of these servers: idle time 90%

- Cheaper hardware: one server for one task was common
- most of these servers: idle time 90%
- cost issue:

- Cheaper hardware: one server for one task was common
- most of these servers: idle time 90%
- cost issue:
  - support, maintenance

- Cheaper hardware: one server for one task was common
- most of these servers: idle time 90%
- cost issue:
  - support, maintenance
  - power consumption (operation, cooling)
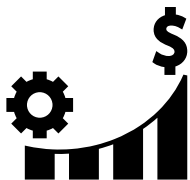
- Cheaper hardware: one server for one task was common
- most of these servers: idle time 90%
- cost issue:
    - support, maintenance
    - power consumption (operation, cooling)
    - space

- Cheaper hardware: one server for one task was common
- most of these servers: idle time 90%
- cost issue:
  - support, maintenance
  - power consumption (operation, cooling)
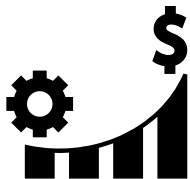  - space
- Virtualization allows consolidation

- Cheaper hardware: one server for one task was common
- most of these servers: idle time 90%
- cost issue:
  - support, maintenance
  - power consumption (operation, cooling)
  - space
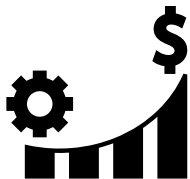- Virtualization allows consolidation
  - multiple servers on one box

Daniel Gruss

- Better hardware utilization
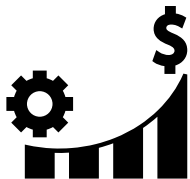
　　　　　　　　　　　　　　　　　　　　　　　　　　　　　　　　Daniel Gruss

- Better hardware utilization
- Lower administration cost

Daniel Gruss

- Better hardware utilization
- Lower administration cost
- long-term operations of older applications

Daniel Gruss

- Better hardware utilization
- Lower administration cost
- long-term operations of older applications
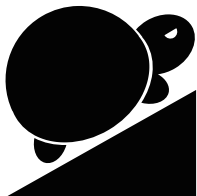- lower down-times

Daniel Gruss

- Better hardware utilization
- Lower administration cost
- long-term operations of older applications
- lower down-times
- simple migration to more powerful hardware

Daniel Gruss

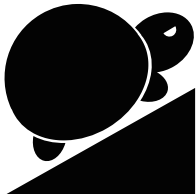- Performance cost: slower I/O operation

- Performance cost: slower I/O operation
- single point of failure: requires better hardware reliability

Daniel Gruss

- Performance cost: slower I/O operation
- single point of failure: requires better hardware reliability
- security gets more complex

- Virtualization no significant role in internet hosting

Daniel Gruss

- Virtualization no significant role in internet hosting
- often PaaS

- Virtualization no significant role in internet hosting
- often PaaS
- Web hosts (FTP access, HTTP website)

- Virtualization no significant role in internet hosting
- often PaaS
- Web hosts (FTP access, HTTP website)
- Isolation on the OS level (tenants as users)

- Virtualization no significant role in internet hosting
- often PaaS
- Web hosts (FTP access, HTTP website)
- Isolation on the OS level (tenants as users)
- no hardware support $\rightarrow$ expensive $+$ many problems

- OS-level Virtualization

Daniel Gruss

- OS-level Virtualization
- Para-Virtualization

Daniel Gruss

- OS-level Virtualization
- Para-Virtualization
- Full Virtualization

Daniel Gruss

- OS-level Virtualization
- Para-Virtualization
- Full Virtualization
- Hardware-Assisted Virtualization

Daniel Gruss

- integrated into kernel

Daniel Gruss

- integrated into kernel
- all application software intended to run in a virtual environment get strictly separated virtual runtime environments (container, jail)

- integrated into kernel
- all application software intended to run in a virtual environment get strictly separated virtual runtime environments (container, jail)
- no separate kernels - only process level virtualization

- integrated into kernel
- all application software intended to run in a virtual environment get strictly separated virtual runtime environments (container, jail)
- no separate kernels - only process level virtualization
- can't run other OSes - only for applications

- integrated into kernel
- all application software intended to run in a virtual environment get strictly separated virtual runtime environments (container, jail)
- no separate kernels - only process level virtualization
- can't run other OSes - only for applications
- examples: OpenVZ, Docker, (s)chroot

Daniel Gruss

Daniel Gruss

- Cooperation with OS: OS is aware of virtualization

Daniel Gruss

- Cooperation with OS: OS is aware of virtualization
- needs to modify guest

- Cooperation with OS: OS is aware of virtualization
- needs to modify guest
- not usable for closed source OSes

- OS not aware of being virtualized

Daniel Gruss

- OS not aware of being virtualized
- no need to adapt guest

Daniel Gruss

- OS not aware of being virtualized
- no need to adapt guest
- reduced performance

Daniel Gruss

- OS not aware of being virtualized
- no need to adapt guest
- reduced performance
    - up to 25%

Daniel Gruss

- OS not aware of being virtualized
- no need to adapt guest
- reduced performance
    - up to 25%
- full virtualization of HW required (e.g., emulation via qemu)

Daniel Gruss

- OS not aware of being virtualized
- no need to adapt guest
- reduced performance
    - up to 25%
- full virtualization of HW required (e.g., emulation via qemu)
    - virtual machines not allowed to access physical components

- OS not aware of being virtualized
- no need to adapt guest
- reduced performance
  - up to 25%
- full virtualization of HW required (e.g., emulation via qemu)
  - virtual machines not allowed to access physical components
  - every physical component has to be virtualized and requires drivers in OS

- Guest no longer runs in kernel mode (Ring 0)

Daniel Gruss

- Guest no longer runs in kernel mode (Ring 0)
  - parts that require kernel privileges won't run

- Guest no longer runs in kernel mode (Ring 0)
    - parts that require kernel privileges won't run
- hypervisor (VMM) changes binaries of guest-OS on the fly

- Guest no longer runs in kernel mode (Ring 0)
  - parts that require kernel privileges won't run
- hypervisor (VMM) changes binaries of guest-OS on the fly
- allows supporting any OS

- Guest no longer runs in kernel mode (Ring 0)
  - parts that require kernel privileges won't run
- hypervisor (VMM) changes binaries of guest-OS on the fly
- allows supporting any OS
  - no need to change source

- Guest no longer runs in kernel mode (Ring 0)
  - parts that require kernel privileges won't run
- hypervisor (VMM) changes binaries of guest-OS on the fly
- allows supporting any OS
  - no need to change source
- high performance penalty

- First full x86 virtualization

- First full x86 virtualization
- hypervisor continuously reads program code before it is executed (pre-scan)

- First full x86 virtualization
- hypervisor continuously reads program code before it is executed (pre-scan)
- looking for relevant commands

Daniel Gruss

- First full x86 virtualization
- hypervisor continuously reads program code before it is executed (pre-scan)
- looking for relevant commands
  - change of system state

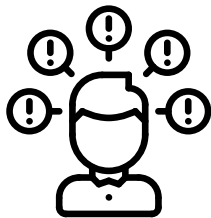- First full x86 virtualization
- hypervisor continuously reads program code before it is executed (pre-scan)
- looking for relevant commands
  - change of system state
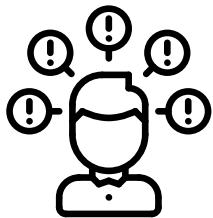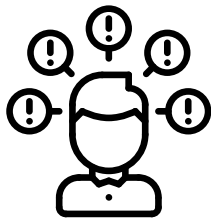  - commands depending on CPU state

- First full x86 virtualization
- hypervisor continuously reads program code before it is executed (pre-scan)
- looking for relevant commands
  - change of system state
  - commands depending on CPU state
- sets breakpoint and lets OS run

Daniel Gruss

Daniel Gruss

- Diverse problems were to be solved when virtualizing on IA-32:
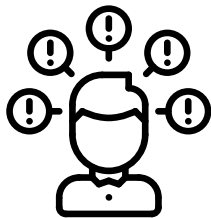
Daniel Gruss

- Diverse problems were to be solved when virtualizing on IA-32:
  - Ring Problems

- Diverse problems were to be solved when virtualizing on IA-32:
  - Ring Problems
  - Address Space Compression

- Diverse problems were to be solved when virtualizing on IA-32:
    - Ring Problems
    - Address Space Compression
    - Non-Faulting Access to Priv. State

- Diverse problems were to be solved when virtualizing on IA-32:
    - Ring Problems
    - Address Space Compression
    - Non-Faulting Access to Priv. State
    - SYSENTER / SYSEXIT

- Diverse problems were to be solved when virtualizing on IA-32:
  - Ring Problems
  - Address Space Compression
  - Non-Faulting Access to Priv. State
  - SYSENTER / SYSEXIT
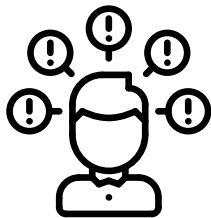  - Interrupt Virtualization

- Diverse problems were to be solved when virtualizing on IA-32:
  - Ring Problems
  - Address Space Compression
  - Non-Faulting Access to Priv. State
  - SYSENTER / SYSEXIT
  - Interrupt Virtualization
  - Hidden States

- usually: application run in ring 3, kernel in ring 0

Daniel Gruss

- usually: application run in ring 3, kernel in ring 0
- guest may not run in ring 0

- usually: application run in ring 3, kernel in ring 0
- guest may not run in ring 0
- ring de-privileging needed: guest must run in ring $> 0$

- usually: application run in ring 3, kernel in ring 0
- guest may not run in ring 0
- ring de-privileging needed: guest must run in ring $> 0$
  - most often 1 or 3

- guest has to run in a ring it has not been developed for

- guest has to run in a ring it has not been developed for
- certain instructions contain privilege level in result (e.g. `PUSH CS`)

Daniel Gruss

- guest has to run in a ring it has not been developed for
- certain instructions contain privilege level in result (e.g. `PUSH CS`)
- guest OS can find out ring it is running in

- guest has to run in a ring it has not been developed for
- certain instructions contain privilege level in result (e.g. PUSH CS)
- guest OS can find out ring it is running in
- may result in diverse problems

- Guest expects to have full address space available

Daniel Gruss

- Guest expects to have full address space available
- hypervisor requires part of address space

- Guest expects to have full address space available
- hypervisor requires part of address space
  - control structures for switching between guest and hypervisor

- Guest expects to have full address space available
- hypervisor requires part of address space
    - control structures for switching between guest and hypervisor
- Access to these areas not allowed for guest. Invokes switch to hypervisor who has to emulate these accesses

- unprivileged software may not access certain elements of the CPU state

Daniel Gruss

- unprivileged software may not access certain elements of the CPU state
- access by guest results in fault: hypervisor can emulate instructions

Daniel Gruss

- unprivileged software may not access certain elements of the CPU state
- access by guest results in fault: hypervisor can emulate instructions
- IA-32 possesses instructions that do not induce a fault:

- unprivileged software may not access certain elements of the CPU state
- access by guest results in fault: hypervisor can emulate instructions
- IA-32 possesses instructions that do not induce a fault:
    - Registers GDTR, IDTR, LDTR and TR are only modifiable in ring 0

- unprivileged software may not access certain elements of the CPU state
- access by guest results in fault: hypervisor can emulate instructions
- IA-32 possesses instructions that do not induce a fault:
    - Registers GDTR, IDTR, LDTR and TR are only modifiable in ring 0
    - can be executed in any ring without fault (without function)

- special commands for fast syscalls

- special commands for fast syscalls
- SYSENTER always switches to ring 0

Daniel Gruss

- special commands for fast syscalls
- SYSENTER always switches to ring 0
- SYSEXIT can only be executed in ring 0

- special commands for fast syscalls
- SYSENTER always switches to ring 0
- SYSEXIT can only be executed in ring 0
- ring 1 thus is problematic

- special commands for fast syscalls
- SYSENTER always switches to ring 0
- SYSEXIT can only be executed in ring 0
- ring 1 thus is problematic
  - SYSENTER switches to hypervisor $\rightarrow$ has to emulate

- special commands for fast syscalls
- SYSENTER always switches to ring 0
- SYSEXIT can only be executed in ring 0
- ring 1 thus is problematic
  - SYSENTER switches to hypervisor $\rightarrow$ has to emulate
  - SYSEXIT switches to hypervisor $\rightarrow$ has to emulate

- interrupts can be masked (so they do not occur if not welcome)

- interrupts can be masked (so they do not occur if not welcome)
- controlled by IF-flag in EFLAGS-Register

- interrupts can be masked (so they do not occur if not welcome)
- controlled by IF-flag in EFLAGS-Register
- Interrupts managed by VM though

- interrupts can be masked (so they do not occur if not welcome)
- controlled by IF-flag in EFLAGS-Register
- Interrupts managed by VM though
- change of IF $\rightarrow$ fault to hypervisor

- interrupts can be masked (so they do not occur if not welcome)
- controlled by IF-flag in EFLAGS-Register
- Interrupts managed by VM though
- change of IF $\rightarrow$ fault to hypervisor
- OS do this quite often $\rightarrow$ performance problem

- interrupts can be masked (so they do not occur if not welcome)
- controlled by IF-flag in EFLAGS-Register
- Interrupts managed by VM though
- change of IF $\rightarrow$ fault to hypervisor
- OS do this quite often $\rightarrow$ performance problem
- forwarding of virtual interrupts must consider IF

Daniel Gruss

- Not all state-information accessible via registers

Daniel Gruss

- Not all state-information accessible via registers
- cannot be saved and restored when switching between VMs

Daniel Gruss

- Two new operating modes:

Daniel Gruss

- Two new operating modes:
  - VMX root operation

- Two new operating modes:
  - VMX root operation
    - for hypervisor

Daniel Gruss

- Two new operating modes:
  - VMX root operation
    - for hypervisor
  - VMX non-root operation

- Two new operating modes:
  - VMX root operation
    - for hypervisor
  - VMX non-root operation
    - controlled by hypervisor

Daniel Gruss

- Two new operating modes:
  - VMX root operation
    - for hypervisor
  - VMX non-root operation
    - controlled by hypervisor
    - supports VMs

- Two new operating modes:
  - VMX root operation
    - for hypervisor
  - VMX non-root operation
    - controlled by hypervisor
    - supports VMs
- Both modes have ring 0-3

- Two new operating modes:
  - VMX root operation
    - for hypervisor
  - VMX non-root operation
    - controlled by hypervisor
    - supports VMs
- Both modes have ring 0-3
- guest can run in ring 0

- Two new operating modes:
  - VMX root operation
    - for hypervisor
  - VMX non-root operation
    - controlled by hypervisor
    - supports VMs
- Both modes have ring 0-3
- guest can run in ring 0
- hypervisor said to be running in "ring -1"

Ring 3
Ring 1&2
Ring 0
Ring -1
Ring -2
Ring -3

Intel ME
UEFI & SMM
Hypervisor
Kernel
Device drivers
Applications

- VM entry: root operation $\rightarrow$ non-root operation

- VM entry: root operation → non-root operation
- VM exit: non-root operation → root operation

- VM entry: root operation → non-root operation
- VM exit: non-root operation → root operation
- VMCS: Virtual Machine Control Structure

- VM entry: root operation $\rightarrow$ non-root operation
- VM exit: non-root operation $\rightarrow$ root operation
- VMCS: Virtual Machine Control Structure
  - Guest-state-area

- VM entry: root operation $\rightarrow$ non-root operation
- VM exit: non-root operation $\rightarrow$ root operation
- VMCS: Virtual Machine Control Structure
  - Guest-state-area
  - Host-state-area

- VM entry: root operation → non-root operation
- VM exit: non-root operation → root operation
- VMCS: Virtual Machine Control Structure
  - Guest-state-area
  - Host-state-area
- Entry/Exit loads/safes information using the proper area

- Contains elements comprising the state of the virtual CPU of a VMCS

Daniel Gruss

- Contains elements comprising the state of the virtual CPU of a VMCS
- VM-exit requires loading certain registers (like segment registers, CR3, IRTR...)

- Contains elements comprising the state of the virtual CPU of a VMCS
- VM-exit requires loading certain registers (like segment registers, CR3, IRTR...)
- GSA contains fields for these registers

- Contains elements comprising the state of the virtual CPU of a VMCS
- VM-exit requires loading certain registers (like segment registers, CR3, IRTR...)
- GSA contains fields for these registers
- GSA contains fields for other information not readable via registers

- Contains elements comprising the state of the virtual CPU of a VMCS
- VM-exit requires loading certain registers (like segment registers, CR3, IRTR...)
- GSA contains fields for these registers
- GSA contains fields for other information not readable via registers
    - e.g. interruptability state

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

CopyLeft 2017, @Noteworthy (Intel Manuel of July 2017)

# 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 | Segment Limit | Access Right |
| IDTR | Selector | Base Address | Segment 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 | | |
|---|---|---|
| CR0 | 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 |

- Addressed using physical addresses

Daniel Gruss

- Addressed using physical addresses
- not part of guest address space

Daniel Gruss

- Addressed using physical addresses
- not part of guest address space
- hypervisor may run in different address space as guest (CR3 part of state)

- Addressed using physical addresses
- not part of guest address space
- hypervisor may run in different address space as guest (CR3 part of state)
- VM-exits leave detailed information on reason for exit in VMCS

- Addressed using physical addresses
- not part of guest address space
- hypervisor may run in different address space as guest (CR3 part of state)
- VM-exits leave detailed information on reason for exit in VMCS
    - exit reason

# VMCS

- Addressed using physical addresses
- not part of guest address space
- hypervisor may run in different address space as guest (CR3 part of state)
- VM-exits leave detailed information on reason for exit in VMCS
  - exit reason
  - exit qualification

Daniel Gruss

# VM-EXIT CONTROL FIELDS

| VM-Exit Controls | Save debug controls | | Host address space size | | | Load IA32_PERF_GLOBAL_CTRL | |
|---|---|---|---|---|---|---|---|
| | Acknowledge interrupt on exit | | Save IA32_PAT | Load IA32_PAT | | Save IA32_EFER | Load IA32_EFER |
| | Save VMX preemption timer value | | Clear IA32_BNDCFGS | | | Conceal VM exits from Intel PT | |
| VM-Exit Controls for MSRs | VM-exit MSR-store count | | VM-exit MSR-store address | | | | |
| | VM-exit MSR-load count | | VM-exit MSR-load address | | | | |

# VM-EXIT INFORMATION FIELDS

| Basic VM-Exit Information | Exit reason | | Exit qualification | |
|---|---|---|---|---|
| | Guest-linear address | | Guest-physical address | |
| VM Exits Due to Vectored Events | VM-exit interruption information | | VM-exit interruption error code | |
| VM Exits That Occur During Event Delivery | IDT-vectoring information | | IDT-vectoring error code | |
| VM Exits Due to Instruction Execution | VM-exit instruction length | | VM-exit instruction information | |
| | I/O RCX | I/O RSI | I/O RDI | I/O RIP |
| VM-instruction error field | | | | |

- Example: MOV CR

Daniel Gruss

- Example: MOV CR
- Exit reason: "control register access"

- Example: MOV CR
- Exit reason: "control register access"
- Exit qualification:

Daniel Gruss

- Example: MOV CR
- Exit reason: "control register access"
- Exit qualification:
  - which CR

Daniel Gruss

- Example: MOV CR
- Exit reason: "control register access"
- Exit qualification:
  - which CR
  - direction (Rx→CR or CR→Rx)

- Example: MOV CR
- Exit reason: "control register access"
- Exit qualification:
  - which CR
  - direction (Rx→CR or CR→Rx)
  - register used

# 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 | RDPMC exiting |
| | RDTSC exiting | CR3-load exiting | CR3-store exiting | CR8-load exiting |
| | CR8-store exiting | Use TPR shadow | NMI-window exiting | MOV-DR exiting |
| | Unconditional I/O exiting | Use I/O bitmaps | Monitor trap flag | Use MSR bitmaps |
| | MONITOR exiting | PAUSE exiting | Activate secondary controls |
| **Secondary processor-based VM-execution controls** | Virtualize APIC accesses | Enable EPT | Descriptor-table exiting | Enable RDTSCP |
| | Virtualize x2APIC mode | Enable VPID | WBINVD exiting | Unrestricted guest |
| | APIC-register virtualization | Virtual-interrupt delivery | PAUSE-loop exiting |
| | RDRAND exiting | Enable INVPCID | Enable VM functions | VMCS shadowing |
| | Enable ENCLS exiting | RDSEED exiting | Enable PML | EPT-violation #VE |
| | 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 low 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 | |

The next step ($\approx$ 2005):

- Virtualization Hardware Extensions for Intel and AMD

The next step ($\approx$ 2005):

- Virtualization Hardware Extensions for Intel and AMD
- $\rightarrow$ substantially lower overheads for VMs

The next step ($\approx$ 2005):

- Virtualization Hardware Extensions for Intel and AMD
- $\rightarrow$ substantially lower overheads for VMs
- $\rightarrow$ better isolation

Daniel Gruss

The next step ($\approx$ 2005):

- Virtualization Hardware Extensions for Intel and AMD
- → substantially lower overheads for VMs
- → better isolation
- → IaaS VMs become widely used

- Support for interrupt-virtualization

Daniel Gruss

- Support for interrupt-virtualization
  - VM-exit with every external interrupt (cannot be masked by guest)

- Support for interrupt-virtualization
  - VM-exit with every external interrupt (cannot be masked by guest)
  - VM-exit when guest-OS ready to accept interrupts (EFLAGS.IF==1)

- Support for interrupt-virtualization
  - VM-exit with every external interrupt (cannot be masked by guest)
  - VM-exit when guest-OS ready to accept interrupts (EFLAGS.IF==1)
- Support for CR0 and CR4-virtualization

- Support for interrupt-virtualization
  - VM-exit with every external interrupt (cannot be masked by guest)
  - VM-exit when guest-OS ready to accept interrupts (EFLAGS.IF==1)
- Support for CR0 and CR4-virtualization
  - VM-exit with any change of these registers

- Support for interrupt-virtualization
  - VM-exit with every external interrupt (cannot be masked by guest)
  - VM-exit when guest-OS ready to accept interrupts (EFLAGS.IF==1)
- Support for CR0 and CR4-virtualization
  - VM-exit with any change of these registers
  - can be set on which bits this shall happen

- Address Space Compression

Daniel Gruss

- Address Space Compression
  - change of address space with any switch guest/hypervisor

- Address Space Compression
  - change of address space with any switch guest/hypervisor
  - guest owns full virtual address space

- Address Space Compression
  - change of address space with any switch guest/hypervisor
  - guest owns full virtual address space
- Ring Problems, SYSENTER/SYSEXIT

Daniel Gruss

- Address Space Compression
  - change of address space with any switch guest/hypervisor
  - guest owns full virtual address space
- Ring Problems, SYSENTER/SYSEXIT
  - Guest can now run in ring 0

- Non-faulting Access to Privileged State

Daniel Gruss

- Non-faulting Access to Privileged State
  - access raise fault into hypervisor

- Non-faulting Access to Privileged State
  - access raise fault into hypervisor
- Hidden State

Daniel Gruss

- Non-faulting Access to Privileged State
  - access raise fault into hypervisor
- Hidden State
  - Saved into VMCS
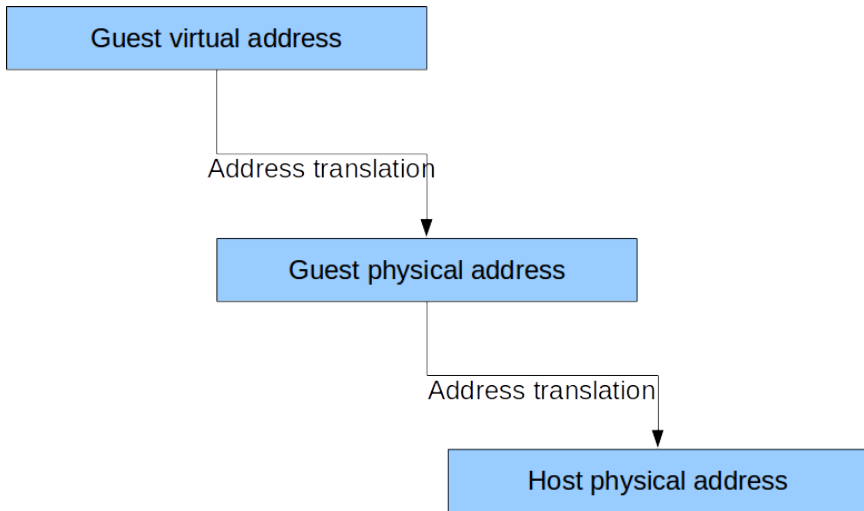
- Hypervisor uses virtual memory

Daniel Gruss

- Hypervisor uses virtual memory
- guest OS uses virtual memory

- Hypervisor uses virtual memory
- guest OS uses virtual memory
- hardware supports page tables

Daniel Gruss

- Hypervisor uses virtual memory
- guest OS uses virtual memory
- hardware supports page tables
- how does this work?

Daniel Gruss

- Hypervisor uses virtual memory
- guest OS uses virtual memory
- hardware supports page tables
- how does this work?
  - shadow page tables

- Hypervisor uses virtual memory
- guest OS uses virtual memory
- hardware supports page tables
- how does this work?
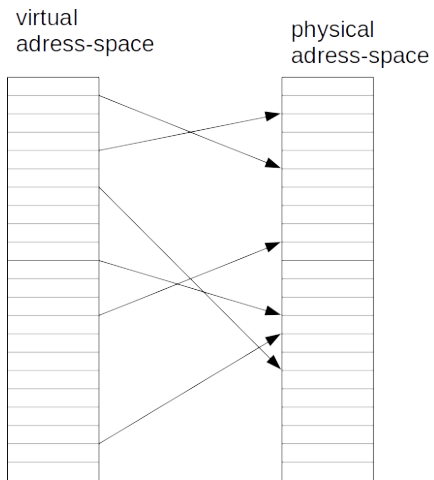    - shadow page tables
    - hardware support

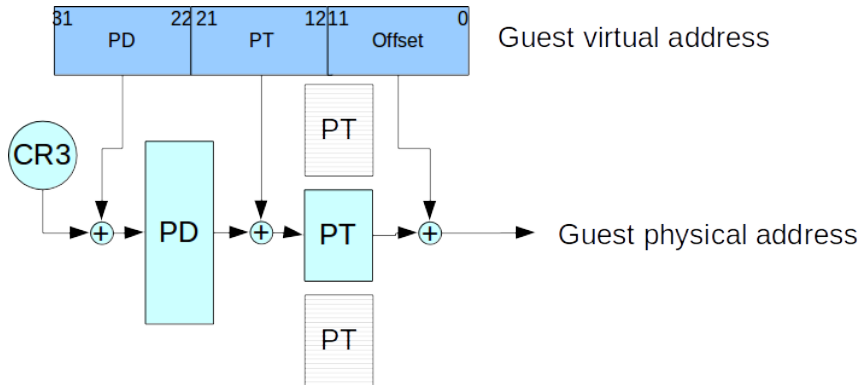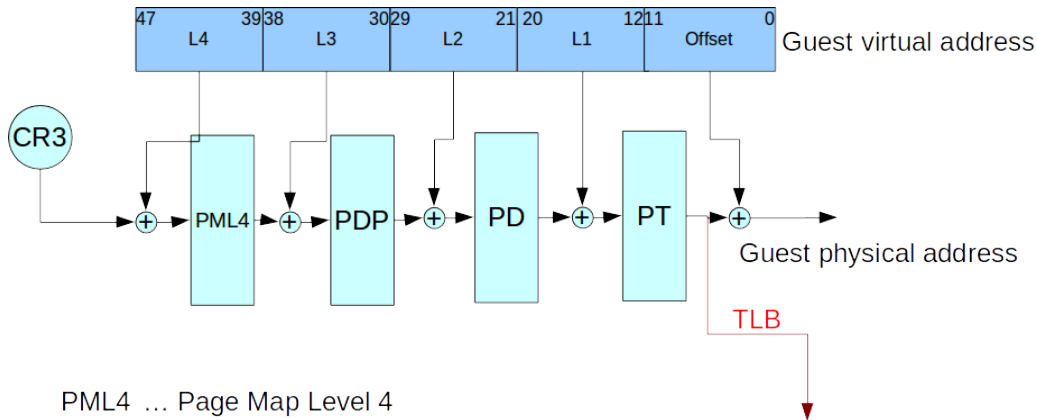All problems in computer science can be solved by another level of indirection.

All problems in computer science can be solved by another level of indirection.

But that usually will create another problem.

David Wheeler

virtual
adress-space

physical
adress-space

Guest virtual address

Guest physical address
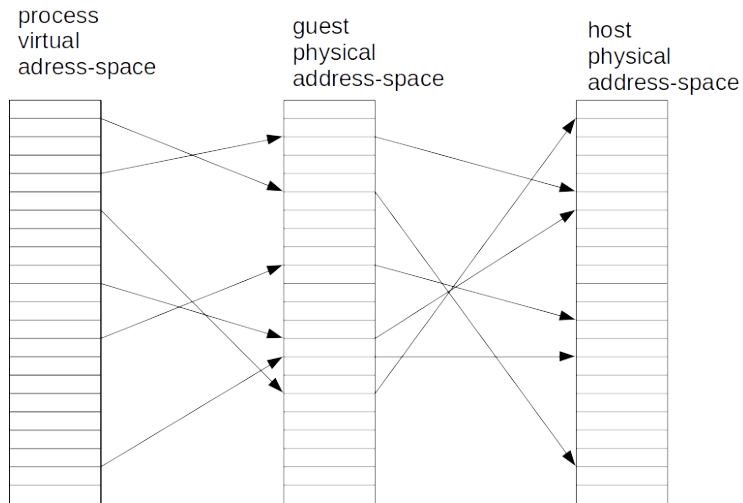
PML4 ... Page Map Level 4
PDP ... Page Directory Pointer
PD ... Page Directory
PT ... Page Table

process
virtual
adress-space

guest
physical
address-space

host
physical
address-space

process
virtual
adress-space

guest
physical
adress-space

host
physical
adress-space

process virtual address-space

guest physical address-space

host physical address-space

Guest page table

Host page table

- **merges** both page tables into one that the HW uses

- **merges** both page tables into one that the HW uses
- when guest changes own page table

- **merges** both page tables into one that the HW uses
- when guest changes own page table
  - Hypervisor has to catch access

Daniel Gruss

- **merges** both page tables into one that the HW uses
- when guest changes own page table
  - Hypervisor has to catch access
  - update shadow page table

- when HW changes shadow page table

- when HW changes shadow page table
- update guest PT

- when HW changes shadow page table
- update guest PT
  - expensive!

- when HW changes shadow page table
- update guest PT
  - expensive!
  - page faults caught by hypervisor

Daniel Gruss

- when HW changes shadow page table
- update guest PT
  - expensive!
  - page faults caught by hypervisor
  - must run through guest PTs

- when HW changes shadow page table
- update guest PT
  - expensive!
  - page faults caught by hypervisor
  - must run through guest PTs
  - must emulate accessed and modified bits for guest

"guest page walk"

- lots of memory accesses....

- lots of memory accesses....
- but how many exactly?

Daniel Gruss

max. number of memory accesses per address translation

- 5 on guest level

max. number of memory accesses per address translation

- 5 on guest level
- each induces 5 on host level

max. number of memory accesses per address translation

- 5 on guest level
- each induces 5 on host level
- makes 25!

Daniel Gruss

- depending on application: 3.9-4.6 times slower

- depending on application: 3.9-4.6 times slower
- but: TLB

Daniel Gruss

- EPT only used if VM active

- EPT only used if VM active
- Translations tagged in TLB with EPT base pointer

- EPT only used if VM active
- Translations tagged in TLB with EPT base pointer
  - differentiate TLB-entries of different VMs

- EPT only used if VM active
- Translations tagged in TLB with EPT base pointer
  - differentiate TLB-entries of different VMs
  - TLB-flush per guest possible

- EPT only used if VM active
- Translations tagged in TLB with EPT base pointer
  - differentiate TLB-entries of different VMs
  - TLB-flush per guest possible
- VPID: virtual processor ID

- EPT only used if VM active
- Translations tagged in TLB with EPT base pointer
    - differentiate TLB-entries of different VMs
    - TLB-flush per guest possible
- VPID: virtual processor ID
    - unique value for each VM

- EPT only used if VM active
- Translations tagged in TLB with EPT base pointer
  - differentiate TLB-entries of different VMs
  - TLB-flush per guest possible
- VPID: virtual processor ID
  - unique value for each VM
  - translations tagged in TLB using VPID

| 63 62 61 60 59 58 57 56 55 54 53 52 51 | M | M-1 | 32 31 30 29 28 27 26 25 24 23 22 21 20 19 18 17 16 15 14 13 12 11 10 9 8 7 6 5 4 3 2 1 0 | | | | | | EPTP[3] |
|---|---|---|---|---|---|---|---|---|---|
| Reserved | | | Address of EPT PML4 table | Rsvd. | S S 2 | A/D | EPT PWL–1 | EPT PS MT | **EPTP[3]** |
| Ignored | Rsvd. | | Address of EPT page-directory-pointer table | Ig n. | X U | Ig4 n. | A | Reserved | X W R | **PML4E: present[6]** |
| SVE[7] | | | Ignored | | | | | | 0 0 0 | **PML4E: not present** |
| SVE | Ign. | S S 8 | Ignored | Rsvd. | Physical address of 1GB page | Reserved | Ig n. X U | D A 1 | I P A T | EPT MT | X W R | **PDPTE: 1GB page** |
| Ignored | Rsvd. | | Address of EPT page directory | Ig n. X U | Ig n. | A 0 | Rsvd. | X W R | **PDPTE: page directory** |
| SVE | | | Ignored | | | | | | 0 0 0 | **PDTPE: not present** |
| SVE | Ign. | S S S | Ignored | Rsvd. | Physical address of 2MB page | Reserved | Ig n. X U | D A 1 | I P A T | EPT MT | X W R | **PDE: 2MB page** |
| Ignored | Rsvd. | | Address of EPT page table | Ig n. X U | Ig n. | A 0 | Rsvd. | X W R | **PDE: page table** |
| SVE | | | Ignored | | | | | | 0 0 0 | **PDE: not present** |
| SVE | Ig n. P P 9 | S S S | Ignored | Rsvd. | Physical address of 4KB page | Ig n. X U | D A | I g n | I P A T | EPT MT | X W R | **PTE: 4KB page** |
| SVE | | | Ignored | | | | | | 0 0 0 | **PTE: not present** |

Figure 28-1.  Formats of EPTP and EPT Paging-Structure Entries

1. Enable VMX via `CR4`

1. Enable VMX via CR4
2. Allocate a VMXON region and use the VMXON instruction

Daniel Gruss

1. Enable VMX via CR4
2. Allocate a VMXON region and use the VMXON instruction
3. Allocate an MSR Bitmap region (we don't want a trap for all MSRs)

Daniel Gruss

1. Enable VMX via `CR4`
2. Allocate a `VMXON` region and use the `VMXON` instruction
3. Allocate an MSR Bitmap region (we don't want a trap for all MSRs)
4. Use `VMCLEAR` instruction

1. Enable VMX via CR4
2. Allocate a VMXON region and use the VMXON instruction
3. Allocate an MSR Bitmap region (we don't want a trap for all MSRs)
4. Use VMCLEAR instruction
5. Execute VMPTRLD to make a VMCS the "current VMCS"

1. Enable VMX via CR4
2. Allocate a VMXON region and use the VMXON instruction
3. Allocate an MSR Bitmap region (we don't want a trap for all MSRs)
4. Use VMCLEAR instruction
5. Execute VMPTRLD to make a VMCS the "current VMCS"
6. Allocate a VMCS region and set up the VMCS (using VMWRITEs)

1. Enable VMX via CR4
2. Allocate a VMXON region and use the VMXON instruction
3. Allocate an MSR Bitmap region (we don't want a trap for all MSRs)
4. Use VMCLEAR instruction
5. Execute VMPTRLD to make a VMCS the "current VMCS"
6. Allocate a VMCS region and set up the VMCS (using VMWRITEs)
7. Use the VMLAUNCH

Similar problem as with Userspace-Kernelspace Isolation:

1. user needs help for some operations (e.g., HW interaction)

Similar problem as with Userspace-Kernelspace Isolation:

1. user needs help for some operations (e.g., HW interaction)
$\rightarrow$ can use a syscall!

Daniel Gruss

Similar problem as with Userspace-Kernelspace Isolation:

1. user needs help for some operations (e.g., HW interaction)
$\rightarrow$ can use a syscall!
2. What about VMs?

Similar problem as with Userspace-Kernelspace Isolation:

1. user needs help for some operations (e.g., HW interaction)
$\rightarrow$ can use a syscall!
2. What about VMs?
3. Same concept different level:

Similar problem as with Userspace-Kernelspace Isolation:

1. user needs help for some operations (e.g., HW interaction)

$\rightarrow$ can use a syscall!

2. What about VMs?

3. Same concept different level:

$\rightarrow$ Hypercalls!

Similar problem as with Userspace-Kernelspace Isolation:

1. user needs help for some operations (e.g., HW interaction)
$\rightarrow$ can use a syscall!
2. What about VMs?
3. Same concept different level:
$\rightarrow$ Hypercalls!

Similar problem as with Userspace-Kernelspace Isolation:

1. user needs help for some operations (e.g., HW interaction)
$\rightarrow$ can use a syscall!
2. What about VMs?
3. Same concept different level:
$\rightarrow$ Hypercalls! via the `vmcall` instruction

Optimization

- Full virtualization often not needed

Optimization

- Full virtualization often not needed
- Serverless / Edge Computing (it's still a form of cloud computing)

Optimization

- Full virtualization often not needed
- Serverless / Edge Computing (it's still a form of cloud computing)
- Virtualization is not for free → why not skip it and just use OS level isolation?
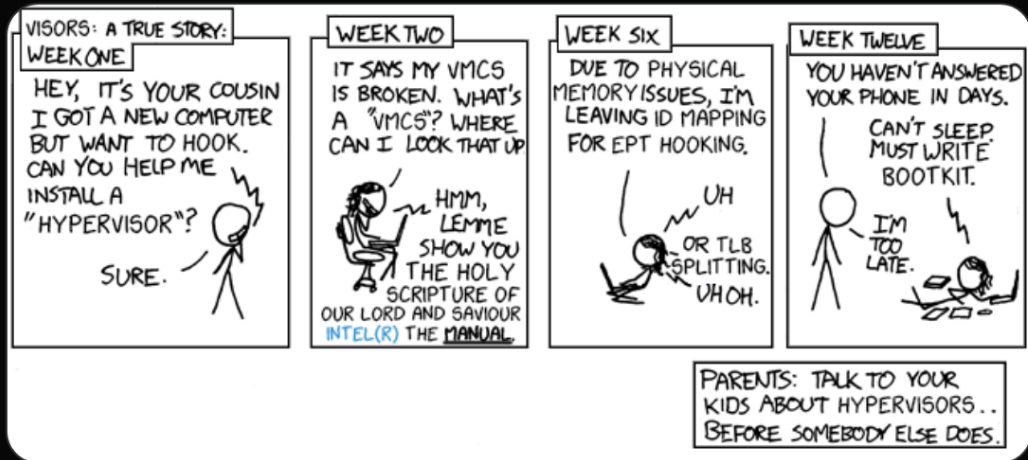
Optimization

- Full virtualization often not needed
- Serverless / Edge Computing (it's still a form of cloud computing)
- Virtualization is not for free → why not skip it and just use OS level isolation?
- Context switches between processes are expensive → why not skip process isolation and just use language-level isolation?

Cloud Operating Systems $\rightarrow$ Hardware-assisted virtualization

Talk to your kids about hypervisors...before someone else does



9  ↻ 19  ♡ 60  ↑

- Seminar-style

- Seminar-style
- You code

- Seminar-style
- You code
- You plan

- Seminar-style
- You code
- You plan
- You present

Daniel Gruss, Andreas Kogler, Fabian Rauscher, Sudheendra Neela

Daniel Gruss

- 100 P. = 100%

- 100 P. $= 100\%$
- 87.5 P. $\rightarrow 1$

Daniel Gruss

- 100 P. $= 100\%$
- 87.5 P. $\rightarrow$ 1
- 75 P. $\rightarrow$ 2

- 100 P. $= 100\%$
- 87.5 P. $\rightarrow$ 1
- 75 P. $\rightarrow$ 2
- 62.5 P. $\rightarrow$ 3

- 100 P. = 100%
- 87.5 P. $\rightarrow$ 1
- 75 P. $\rightarrow$ 2
- 62.5 P. $\rightarrow$ 3
- 50 P. $\rightarrow$ 4

- 10 participants $\rightarrow$ 2-3 teams with each 4 participants (default)

- 10 participants $\rightarrow$ 2-3 teams with each 4 participants (default)
- 5 ECTS = 500h with 125h per team member

- 10 participants $\rightarrow$ 2-3 teams with each 4 participants (default)
- 5 ECTS = 500h with 125h per team member
- Team of 3? Same effort but +5 points

- 10 participants $\rightarrow$ 2-3 teams with each 4 participants (default)
- 5 ECTS $=$ 500h with 125h per team member
- Team of 3? Same effort but $+5$ points
- Team of 2? Same effort but $+10$ points

- 10 participants $\rightarrow$ 2-3 teams with each 4 participants (default)
- 5 ECTS = 500h with 125h per team member
- Team of 3? Same effort but +5 points
- Team of 2? Same effort but +10 points
$\rightarrow$ send us your registration until Friday March 10

- Deadlines: Friday 23:59

- Deadlines: Friday 23:59
- Grace Period: 48 hours **but no support**

Daniel Gruss

- 24.3. Structure Setup
  Estimated Team Effort: 125h, Points: 5P.

Daniel Gruss

- 24.3. Structure Setup
  Estimated Team Effort: 125h, Points: 5P.
- 21.4. Executing Guest Code + Video Output
  Estimated Team Effort: 125h, Points: 15P. $\rightarrow$ AG1

- 24.3. Structure Setup
  Estimated Team Effort: 125h, Points: 5P.
- 21.4. Executing Guest Code + Video Output
  Estimated Team Effort: 125h, Points: 15P. → AG1
- 5.5. Interrupt + Emulate PIC + Public Feature Bidding
  Estimated Team Effort: 100h, Points: 5P.

- 24.3. Structure Setup
  Estimated Team Effort: 125h, Points: 5P.
- 21.4. Executing Guest Code + Video Output
  Estimated Team Effort: 125h, Points: 15P. → AG1
- 5.5. Interrupt + Emulate PIC + Public Feature Bidding
  Estimated Team Effort: 100h, Points: 5P.
- 26.5. Boot Guest SWEB Shell + Virtualize Disk + Private Feature Bidding
  Estimated Team Effort: 75h, Points: 35P. → AG2

- 24.3. Structure Setup
  Estimated Team Effort: 125h, Points: 5P.
- 21.4. Executing Guest Code + Video Output
  Estimated Team Effort: 125h, Points: 15P. → AG1
- 5.5. Interrupt + Emulate PIC + Public Feature Bidding
  Estimated Team Effort: 100h, Points: 5P.
- 26.5. Boot Guest SWEB Shell + Virtualize Disk + Private Feature Bidding
  Estimated Team Effort: 75h, Points: 35P. → AG2
- 2.6. Feature PoC in Booted Guest SWEB
  Estimated Team Effort: 75h, Points: 10P.

- 24.3. Structure Setup
  Estimated Team Effort: 125h, Points: 5P.
- 21.4. Executing Guest Code + Video Output
  Estimated Team Effort: 125h, Points: 15P. → AG1
- 5.5. Interrupt + Emulate PIC + Public Feature Bidding
  Estimated Team Effort: 100h, Points: 5P.
- 26.5. Boot Guest SWEB Shell + Virtualize Disk + Private Feature Bidding
  Estimated Team Effort: 75h, Points: 35P. → AG2
- 2.6. Feature PoC in Booted Guest SWEB
  Estimated Team Effort: 75h, Points: 10P.
- 16.6. Feature Implementation Done + Final Presentation and Demo in Booted Guest SWEB
  Estimated Team Effort: 75h, Points: 30P. → AG3

- 24.3. Structure Setup
  Estimated Team Effort: 125h, Points: 5P.
- 21.4. Executing Guest Code + Video Output
  Estimated Team Effort: 125h, Points: 15P. → AG1
- 5.5. Interrupt + Emulate PIC + Public Feature Bidding
  Estimated Team Effort: 100h, Points: 5P.
- 26.5. Boot Guest SWEB Shell + Virtualize Disk + Private Feature Bidding
  Estimated Team Effort: 75h, Points: 35P. → AG2
- 2.6. Feature PoC in Booted Guest SWEB
  Estimated Team Effort: 75h, Points: 10P.
- 16.6. Feature Implementation Done + Final Presentation and Demo in Booted Guest SWEB
  Estimated Team Effort: 75h, Points: 30P. → AG3
- 16.6. Successful Live Presentation at 21:00, Bonus Points: 5P.