



# **CSE 511: Operating Systems Design**

## **Lecture 12** Virtualization

# Virtualization

- Allows multiple instances of an operating system to run on a single computer
  - Originally introduced for IBM VM/370
  - Has later been revitalized for modern platforms
- A **hypervisor** is a software layer that
  - Separates the **virtual** hardware an OS sees from the **actual** hardware
  - Arbitrates access to physical resources such as CPU or memory

# Virtualization

- Widely known hypervisors
  - Xen, KVM, VMware ESX, Hyper V, VirtualBox, qemu etc
- Virtualization
  - Improves **isolation** and **reliability** because each OS runs independently from the others on its own **virtual** processors so that OS failures are **contained**
  - Increases **resource utilization** as the same hardware can be used for multiple purpose
  - Leads to better **productivity** as large pieces of software can be preconfigured and installed very easily



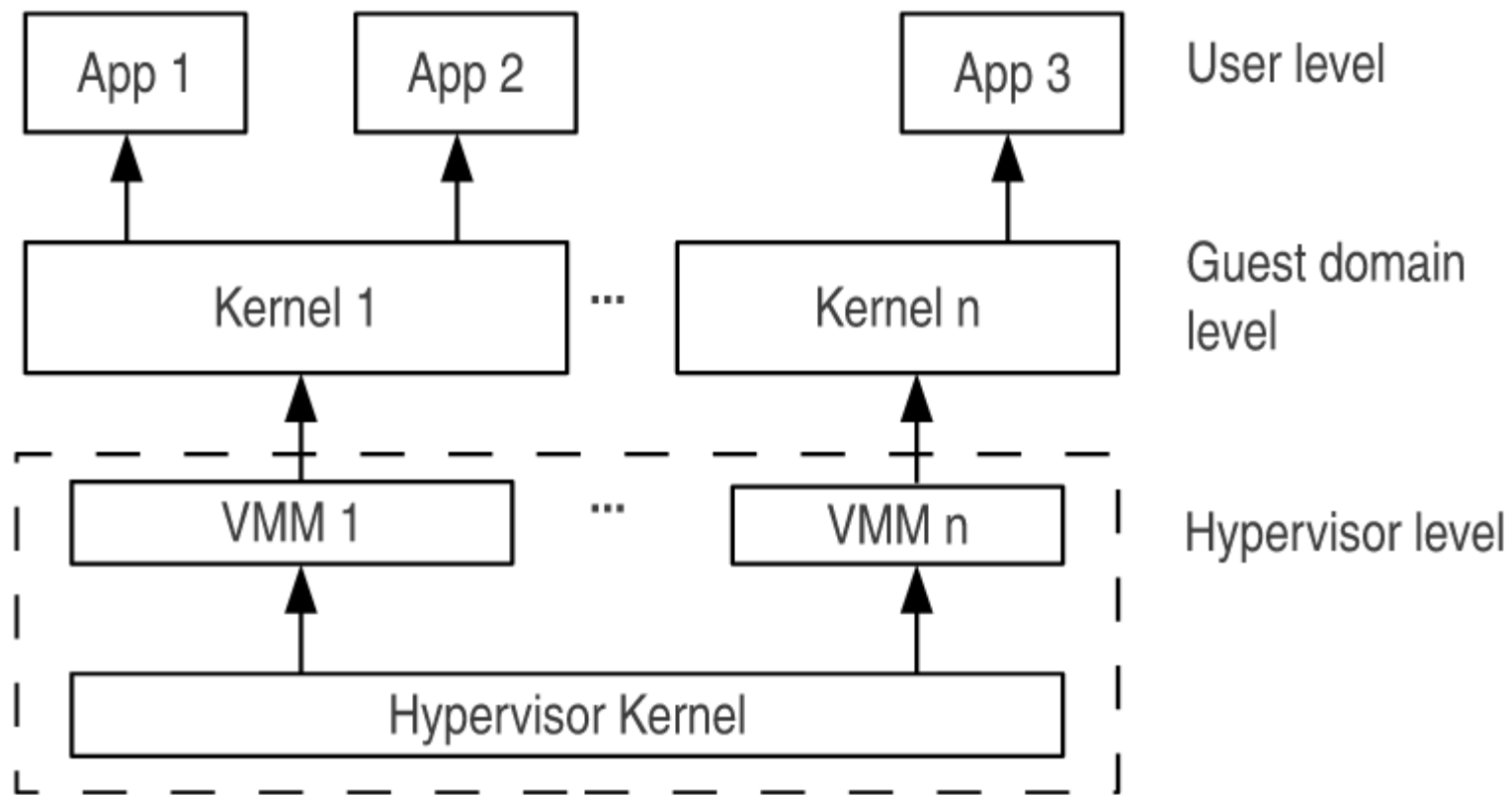
# Virtualization

- Virtual appliances
  - Software bundles containing preconfigured packages with their own OS aimed at simplified distribution and deployment
  - They run along with other virtual appliances and general purpose OS on a single machine
- Infrastructure as a service (IaaS) solutions
  - Rely on virtualization to provide business solutions for server consolidation
- Hypervisors can be utilized to provide better security and robustness for operating systems

# Hypervisor Types

- Type I
  - Allows execution of virtual machines in **guest** domains
- Type II
  - Runs on top of some OS, known as a **host** OS
- Type I hypervisors avoid unnecessary overheads, latencies; they are also not subject to host OS vulnerabilities
  - The hypervisor is at the lowest layer: the hypervisor kernel and Virtual Machine Monitors (VMMs)
  - The hypervisor kernel has direct access to hardware, resource allocation
  - The VMM layer is responsible for virtualizing resources

# Hypervisor Types



# CPU Virtualization

- A hypervisor needs to provide guests with ***virtual*** CPUs
  - Entities that allow guests to have access to CPU resources while allowing the hypervisor to schedule different operating systems on the available ***physical*** CPUs or cores
- Complete and transparent CPU virtualization traditionally required ***architectural support***

# CPU Virtualization

- Popek and Goldberg formulated necessary and sufficient conditions for virtualization
  - ***Fidelity***: requires that the virtual environment in which programs run should be indistinguishable from the real hardware
  - ***Performance***: a substantial amount of instructions need to be executed directly by hardware and without any additional hypervisor handling
  - ***Safety***: requires that the hypervisor has full control of all system resources



# CPU Virtualization

- ***Sensitive instructions*** used by guests need to be trapped and emulated if they could affect the correctness of the hypervisor's behavior
  - ***Trap-and-emulate*** is a technique for intercepting CPU instructions and executing them using special (hypervisor) handlers
- The original x86 architecture does not meet described criteria, as it does not allow to intercept ***all*** necessary instructions when executed in ***unprivileged*** mode

# CPU Virtualization

- Techniques to overcome this problem
  - **Binary translation:** machine code is being processed by a translator program, which allows prefetching, inspecting and special treatment of all relevant instructions (e.g., VMware)
  - **Paravirtualization:** adaptations of the guest OS kernel code that avoids the use of untrappable instructions, reducing emulation and management costs and yielding better performance (e.g., Xen)

# CPU Virtualization

- Xen's original paravirtualization relied on the 4-ring protection hierarchy present in 32-bit x86 processors
  - Ring 0: privileged hypervisor code
  - Ring 1: guest kernel
  - Ring 3: user processes
- Segmentation protected guest kernel code and data from user processes, but guest kernels had to be adapted to successfully run in this mode
  - x86-64 effectively removed segmentation

# CPU Virtualization

- The original paravirtualization approach is not so practical for x86-64
  - Need to run the kernel in Ring 3 and switch page tables every single system call
- But CPU vendors introduced **hardware extensions** (Intel VT-x, AMD-V) that allowed the execution of unchanged guest kernels
  - Safe virtualization of the CPU by introducing a VMM mode distinct from the mode in which privileged guest kernel code executes

# CPU Virtualization

- Later generations added support for MMU virtualization via ***nested paging*** without resorting to ***shadow page tables***
- Well known hypervisors such as Xen and VMware now provide support for hardware virtualization
- Xen modes
  - ***PV***: paravirtualization
  - ***HVM, PVHVM***: using hardware extensions
  - ***PVH***: in between

# CPU Virtualization

Poor Performance  
 Scope for Improvement  
 Optimal Performance

PV = Paravirtualized  
 VS = Software Virtualized (QEMU)  
 VH = Hardware Virtualized  
 HA = Hardware Accelerated



x86 Shortcut	Mode	With					
HVM / Fully Virtualized	HVM		VS	VS <sup>1</sup>	VS	VH	Yes
HVM + PV drivers	HVM	PV Drivers Installed	PV	VS <sup>1</sup>	VS	VH	Yes
PVHVM	HVM	PVHVM Capable Guest	PV	PV <sup>2</sup>	VS	VH	Yes
PVH	PVH	PVH Capable Guest	PV	HA <sup>3</sup>	PV <sup>4</sup>	VH	No
PV	PV		PV	PV	PV <sup>5</sup>	PV	No
<b>ARM</b>							
N/A	N/A		PV	VH	PV <sup>6</sup>	VH	No

\* The picture is taken from  
[https://wiki.xenproject.org/wiki/Xen\\_Project\\_Software\\_Overview](https://wiki.xenproject.org/wiki/Xen_Project_Software_Overview)

# Memory Virtualization

- Conventional OS manage memory assuming that they have unrestricted control over physical memory
- Hypervisors introduce special adaptations to guest OS allowing to run multiple guest domains at the same time
  - *Virtual*, *physical*, and *machine* memory address space
  - *Machine* memory is the actual memory exposed by hardware



# Memory Virtualization

- Physical and virtual memory are guest specific and have their conventional meaning
  - Physical addresses correspond to memory exposed by the VMM rather than hardware
- Hypervisors and CPUs provide support for virtual-to-physical, physical-to-machine and virtual-to-machine address translation
  - Page granularity: hypervisors and guests refer to the physical and machine addresses using their physical and machine frame numbers





# Memory Virtualization

- Fully-virtualized guests are completely unaware of the underlying hypervisor
  - They know nothing about machine frames and treat physical frames as if they represented actual hardware exposed memory
- Guests create page tables containing virtual-to-physical memory mappings
  - Cannot be used by the hardware since physical frames do not correspond to actual memory location

# Memory Virtualization

- The hypervisor traps any attempt to load page tables and replace them with ***shadow page tables***
  - ***Shadow page tables*** are created and maintained by the hypervisor
  - Entries are created on demand
  - Contain virtual-to-machine memory mappings generated by the hypervisor from guest page tables and the physical-to-machine (P2M) hypervisor translation tables
  - Transparent (not exposed to guests)

# Memory Virtualization

- CPU vendors also introduced *nested page tables*
  - Eliminates shadow page tables by providing an additional hardware translation layer
  - The first layer is for regular OS page tables containing virtual-to-physical mappings
  - The second layer is for physical-to-machine mappings and is managed by the hypervisor

# Memory Virtualization

- Paravirtualization
  - Guests are adapted to assist the hypervisor with memory management
  - They have direct access to the P2M translation table, allowing them to perform physical-to-machine translation
  - No shadow page tables: guests create page tables with virtual-to-machine mappings
    - The hypervisor will verify the consistency of page tables

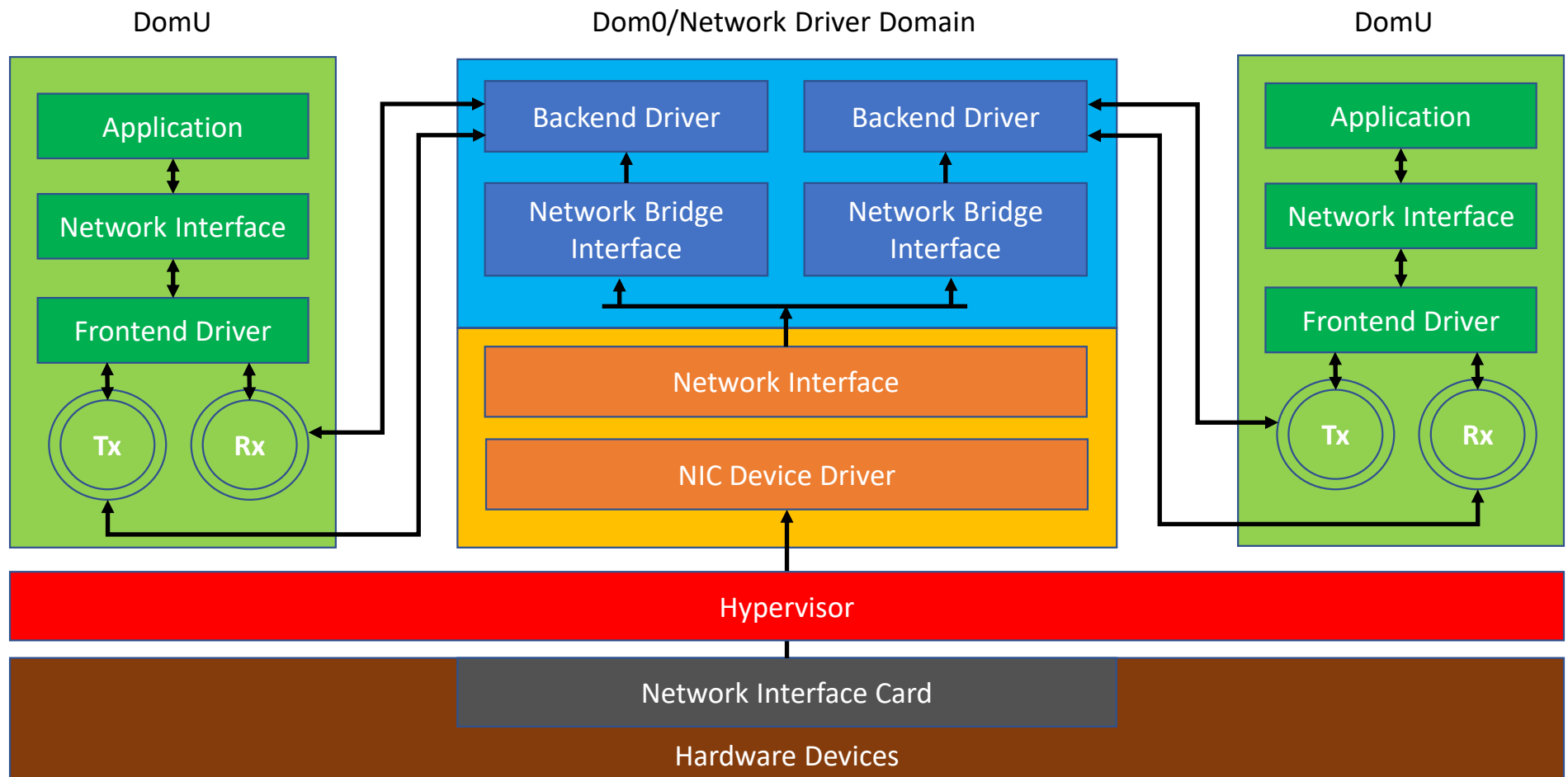
# I/O Virtualization

- A hypervisor needs to organize device sharing
- Device drivers are usually executed in
  - a host OS (Type II hypervisors)
  - specialized driver domains
  - a dedicated privileged domain (e.g., Dom0 in case of Xen)
- The privileged domain has unrestricted access to hardware

# I/O Virtualization

- All other guest domains
  - Complete emulation of devices
    - the privileged domain will emulate devices and guests will access them as if devices were some real piece of hardware
  - A split driver model
    - drivers consist of two parts that interact with each other using inter-domain communication
  - Direct I/O: device is disabled in the privileged domain and reassigned to a specified guest

# I/O Virtualization



# I/O Virtualization

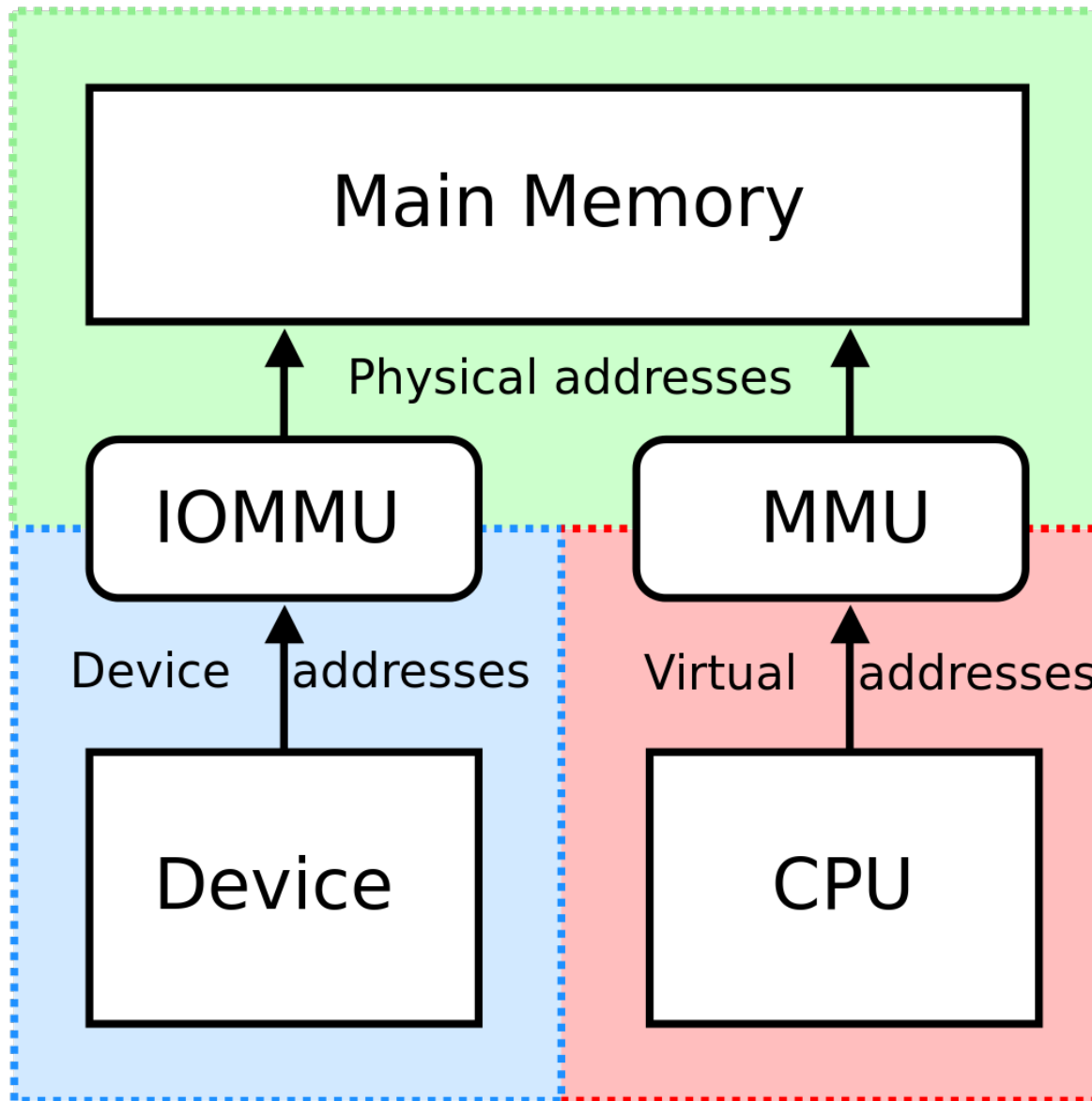
- Drivers can run in the privileged domain (Dom0) or in dedicated driver domains
  - Backend drivers run in Dom0 or driver domains
  - Frontend drivers run in other (guest) domains
- Direct I/O: giving a direct access to a device, “PCI passthrough”
  - Xen/KVM can safely assign access to devices on a PCI bus to guests other than the privileged domain
  - To avoid collisions, the privileged domain excludes those PCI addresses from the I/O space it manages



# I/O Virtualization

- To make PCI passthrough safe, the physical presence of an Input/Output memory management unit (IOMMU) is required
  - Also the hypervisor should support that
  - Specific implementation of IOMMU are Intel VT-d, AMD-Vi
- An IOMMU remaps and protects addresses and interrupts used by memory-mapped I/O devices
  - It thus protects from devices and drivers that might make improper use of DMA or interrupts

# I/O Virtualization



\* The picture is taken from  
[https://en.wikipedia.org/wiki/Input%E2%80%93output\\_memory\\_management\\_unit](https://en.wikipedia.org/wiki/Input%E2%80%93output_memory_management_unit)

# Why Do We Need Virtualization?

- Aside from running legacy OSs or applications, a typical argument would include “convenience” and “isolation”
- **Convenience**
  - Not so clear because of contending solutions, e.g., containers (LxC, Docker, etc)
- **Isolation**
  - A very valid argument due to multitude of system calls, device drivers, etc
  - When one layer breaks (e.g., the Meltdown CPU bug), another layer can still protect
  - The argument is contended by microkernels

# Sharing vs. Isolation

- **Threads** provide virtually no isolation
  - They share almost all crucial program resources such as virtual memory space and file descriptors
- **Processes** provide good isolation with respect to memory accesses and file descriptor management
  - But they share the same file system, process ID namespace, system devices, etc
- **Containers** (LxC) provide different user environments through process filtering, changing root file system and networking isolation
  - Do not allow to run unrelated OSs, share same device drivers, memory management, CPU scheduler and other OS critical components

# Sharing vs. Isolation

- Containers lack complete isolation of core components
  - They are limiting for people who want complete isolation of user environments
- **Virtual Machines** solve this problem by isolating **entire** OSs from each other and leaving only fundamental mechanisms (e.g., network sockets) for communication between various OSs running on the same host
- **Distributed Systems** run their components on different physical hosts, and each host may use any of these methods to organize isolation or sharing within a host
  - Cloud systems provide specific services to remote users such as IaaS (Infrastructure as a Service) or SaaS (Storage as a Service)

# Summary: What Do Hypervisors Manage?

- **Scheduling:** Each guest OS (domain) has its own virtual CPUs (vCPUs)
  - Similar to scheduling of tasks/threads in operating systems
- **Memory:** Each guest OS has to get a slice of physical memory
- **Devices:** Full Emulation, Special Drivers, or Direct I/O
- In summary, a hypervisor can be considered as a special “OS” that schedules guests, allocates memory, gives access to devices
  - The difference is that the level of abstraction is much low-level

# Hypervisor API

- It depends on the hypervisor
  - Xen supports “hypercalls”, which are similar to system calls but used in the hypervisor context
  - A guest kernel can issue hypercalls
    - Avoids trapping and emulation, requests can be batched together to reduce costs, etc
- CPU support
  - Intel VT-x: vmcall, AMD-V: vmxcall
- VT-x and AMD-V diverge but provide comparable capabilities
  - Xen supports both (known as VMX: Intel and SVM: AMD)