# Virtual Machines - history

4381 Model Group 24 (1988): 20 MHz dual CPU, 128KB cache, 64MB RAM, 890K USD (2M USD@2020)

ARCHITECTURAL PRINCIPLES FOR VIRTUAL COMPUTER SYSTEMS

Robert P. Goldberg

Harvard University

AD-772 809

Prepared for:

Electronic Systems Division

February 1973

DISTRIBUTED BY:

NTIS

National Technical Information Service
U. S. DEPARTMENT OF COMMERCE
5285 Port Royal Road, Springfield Va. 22151

- ▸ **The first era of VM**
  - ▸ 1972 – IBM VM for S/370
    - ▪ Coexistence of different OSes
    - ▪ Time-sharing and virtual memory for OSes not implementing them
    - ▪ Debugging of OSes
      - ▪ Including a VM under a VM
  - ▸ Every tenth S/370 used a VM
  - ▸ 1980… – Gradual decline
    - ▪ Mainframes overcome by cheaper architectures (minicomputers, PC)
      - ▪ New hardware did not support VM
    - ▪ The growing dominance of Unix
      - ▪ VM is a complication for inter-process communication

▸ **Formal Requirements
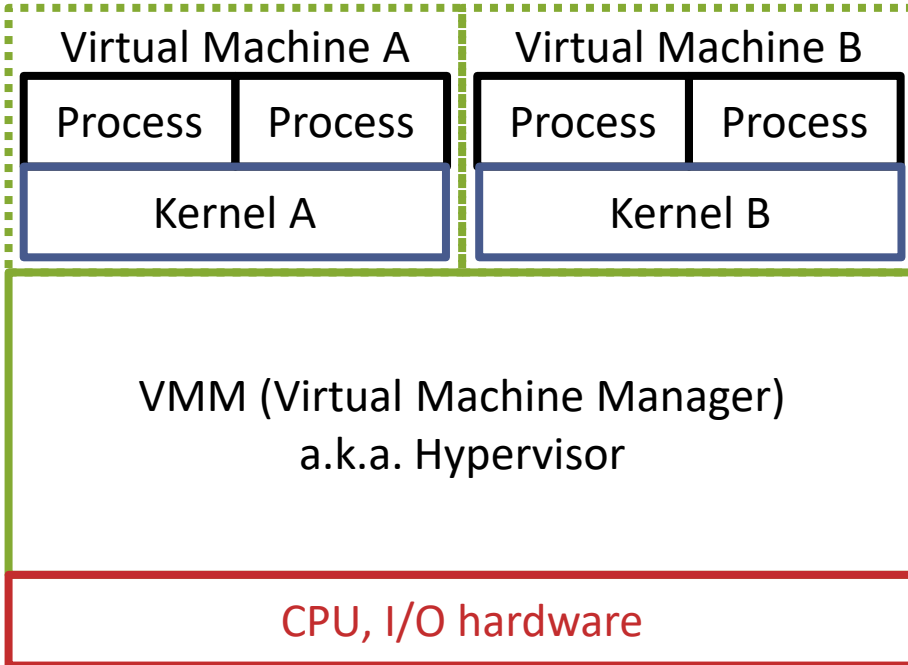for Virtualizable Third Generation Architectures**

 ▸ Gerald J. Popek and Robert P. Goldberg, 1974

 ▸ Equivalence / Fidelity

   ▪ A program running under the VMM should exhibit a behavior essentially identical to that demonstrated when running on an equivalent machine directly.

 ▸ Resource control / Safety

   ▪ The VMM must be in complete control of the virtualized resources

 ▸ Efficiency / Performance

   ▪ A statistically dominant fraction of machine instructions must be executed without VMM intervention

# CPU virtualization

# Types of Virtual Machine Systems
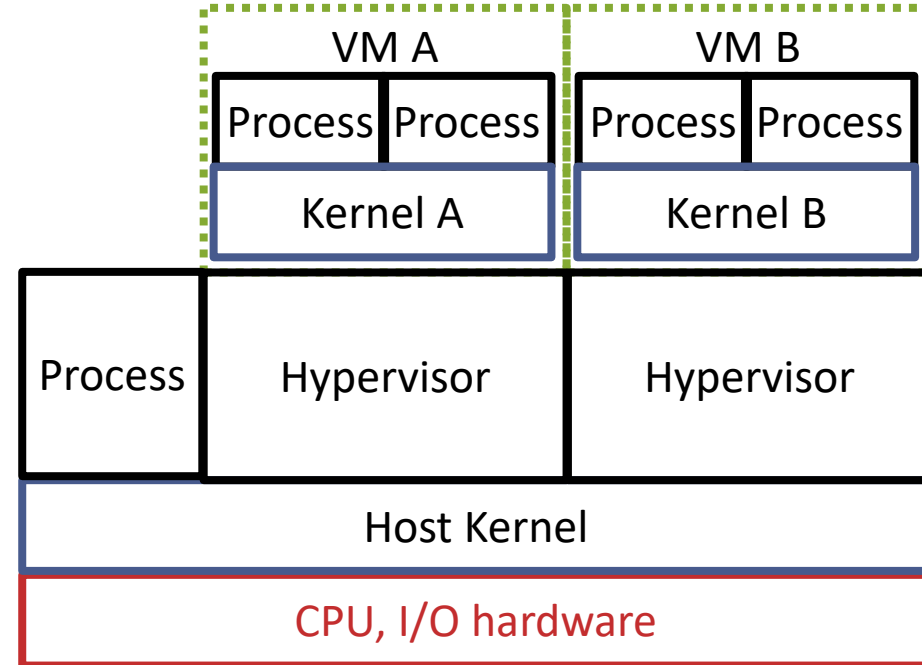
## Type 1 (Bare Metal) Hypervisor

**Example: VMWare ESXi**

| Virtual Machine A | Virtual Machine B |
|---|---|
| Process | Process | Process | Process |
| Kernel A | Kernel B |

**VMM (Virtual Machine Manager)
a.k.a. Hypervisor**

CPU, I/O hardware

▸ Hypervisor on bare metal

  ▸ Hypervisor directly performs all hardware access (CPU configuration, I/O)

    ▪ Requires device drivers
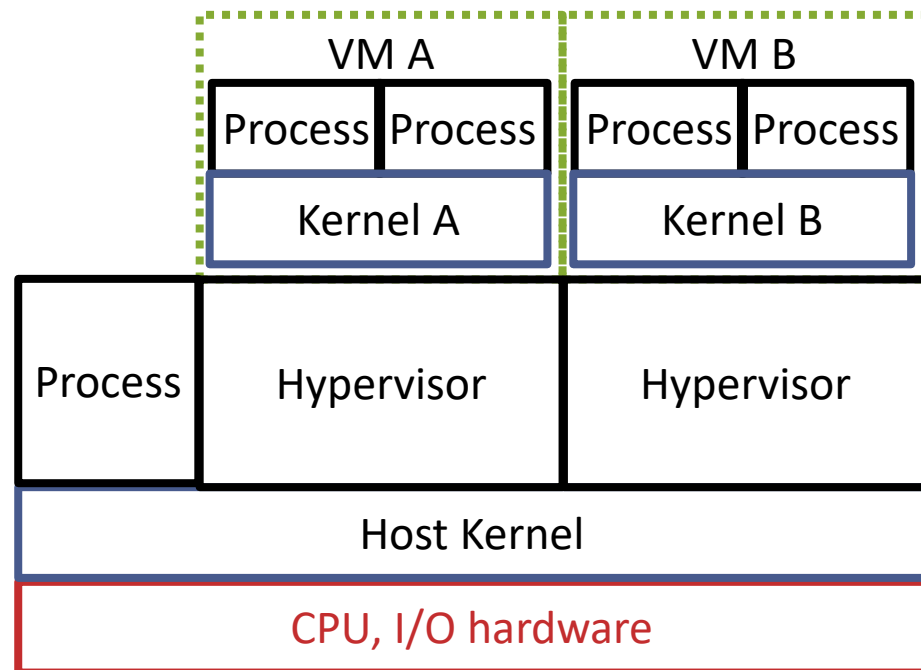    ▪ Complex but fast

## Type 2 (Hosted) Hypervisor
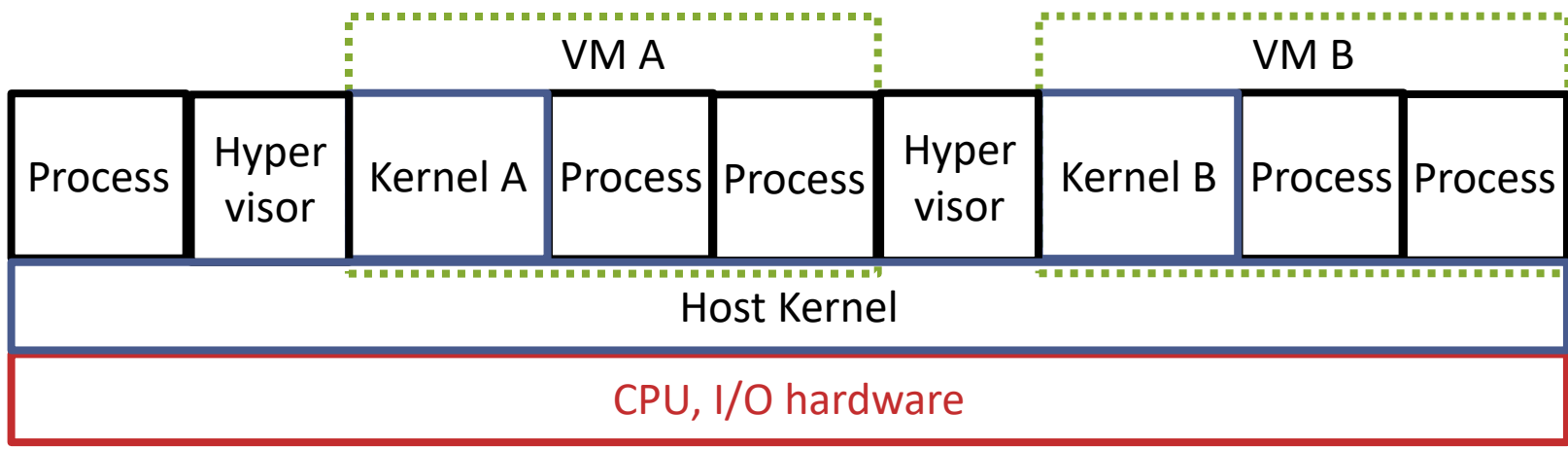
**Example: VMWare Workstation Player**

| VM A | VM B |
|---|---|
| Process | Process | Process | Process |
| Kernel A | Kernel B |

| Process | Hypervisor | Hypervisor |

Host Kernel

CPU, I/O hardware

▸ Hypervisor above an host OS

  ▸ Hypervisor is a (privileged) process

    ▪ Often one per VM
    ▪ I/O access performed by host kernel
    ▪ CPU control requires support from the host kernel (debugging services)
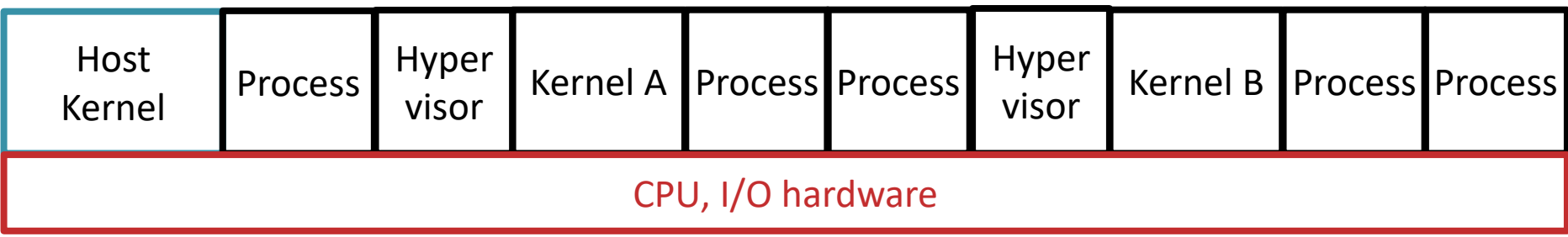
# Beware

▸ Pictures like this are misleading

| | VM A | | VM B | |
|---|---|---|---|---|
| | Process | Process | Process | Process |
| | Kernel A | | Kernel B | |

| Process | Hypervisor | Hypervisor |
|---|---|---|

| Host Kernel |
|---|

| CPU, I/O hardware |
|---|

# Beware

▶ The host kernel actually sees this:

| | | VM A | | | | VM B | | |
|---|---|---|---|---|---|---|---|---|
| Process | Hyper visor | Kernel A | Process | Process | Hyper visor | Kernel B | Process | Process |

**Host Kernel**

**CPU, I/O hardware**

▶ The CPU sees this:

| Host Kernel | Process | Hyper visor | Kernel A | Process | Process | Hyper visor | Kernel B | Process | Process |
|---|---|---|---|---|---|---|---|---|---|

**CPU, I/O hardware**

# Flavors of Type 1 Hypervisors

## Traditional

**Example: VMWare ESXi**

| Virtual Machine A | Virtual Machine B |
|---|---|
| Process \| Process | Process \| Process |
| Kernel A | Kernel B |

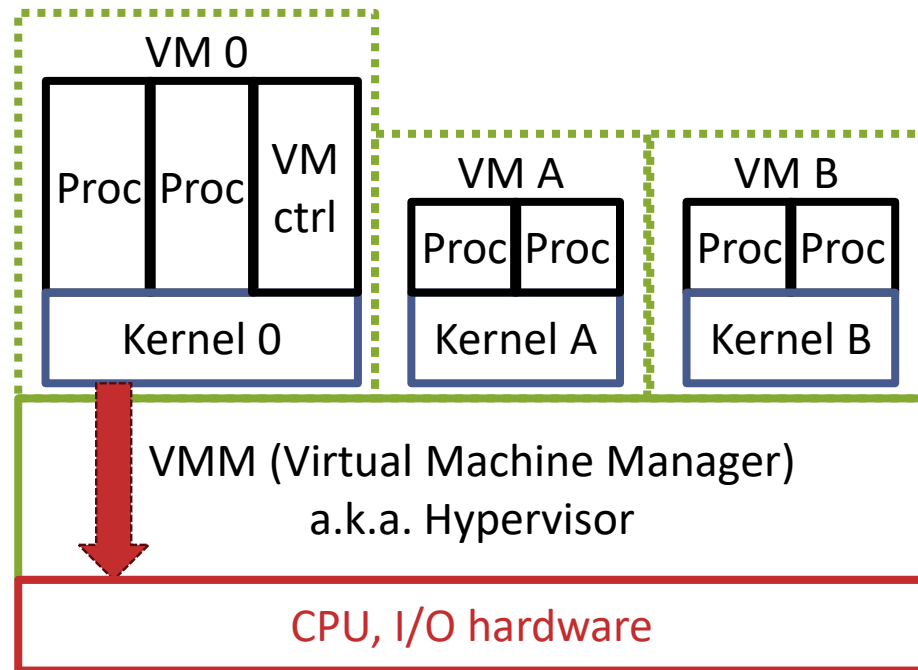**VMM (Virtual Machine Manager)
a.k.a. Hypervisor**

**CPU, I/O hardware**

▸ Hypervisor performs I/O

  ▪ Requires device drivers tailored for the hypervisor

  ▪ Too costly development
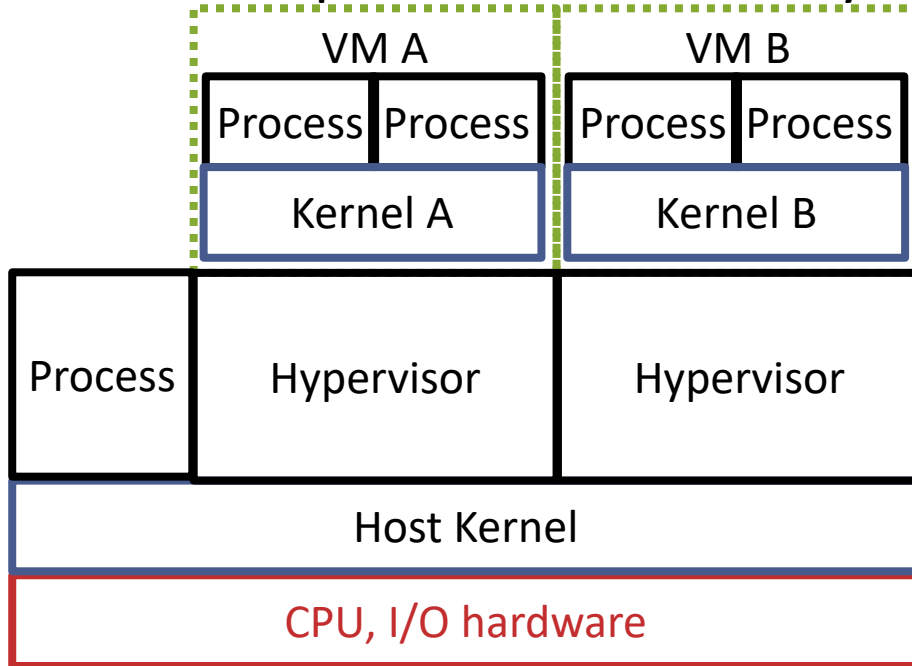
## With root partition (Microsoft terminology)

**Example: Microsoft Windows + Hyper-V**

**VM 0**

| Proc | Proc | VM ctrl |
|---|---|---|
| Kernel 0 | | |

**VM A**

| Proc | Proc |
|---|---|
| Kernel A | |

**VM B**

| Proc | Proc |
|---|---|
| Kernel B | |

**VMM (Virtual Machine Manager)
a.k.a. Hypervisor**

**CPU, I/O hardware**

▸ Hypervisor only controls CPU

  ▸ VM 0 aka Root partition

    ▪ Allowed to directly access I/O hardware

    ▪ Standard OS with device drivers

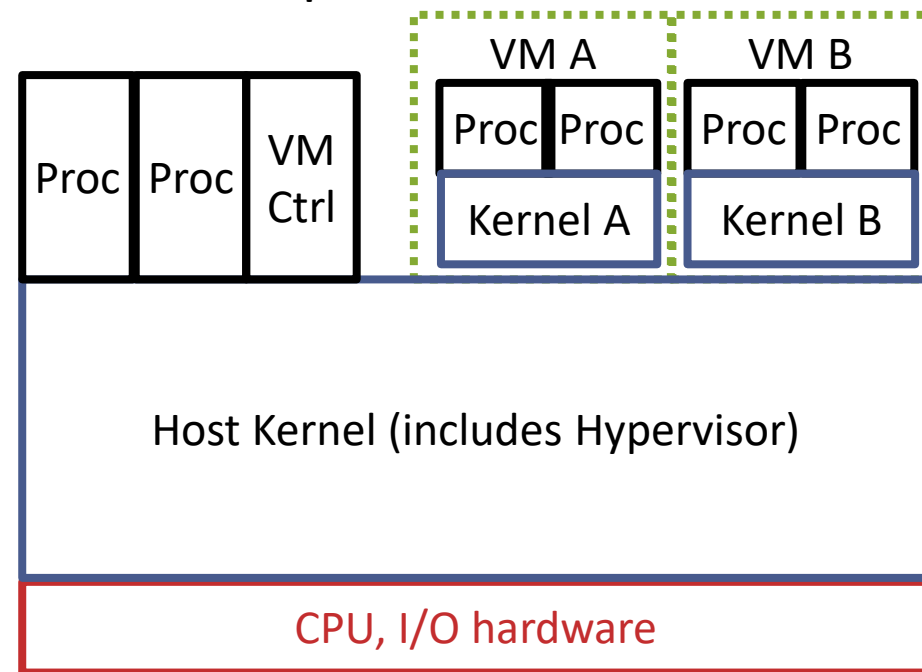  ▸ Hypervisor forwards I/O requests

## Implemented in user-space

**Example: VMWare Workstation Player**



## Implemented in a kernel

**Example: Linux KVM**



▸ ## Hypervisor above an host OS

- ▸ Hypervisor is a (privileged) process
  - ▪ Often one per VM
  - ▪ I/O access performed by host kernel
  - ▪ CPU control requires support from the host kernel (debugging services)
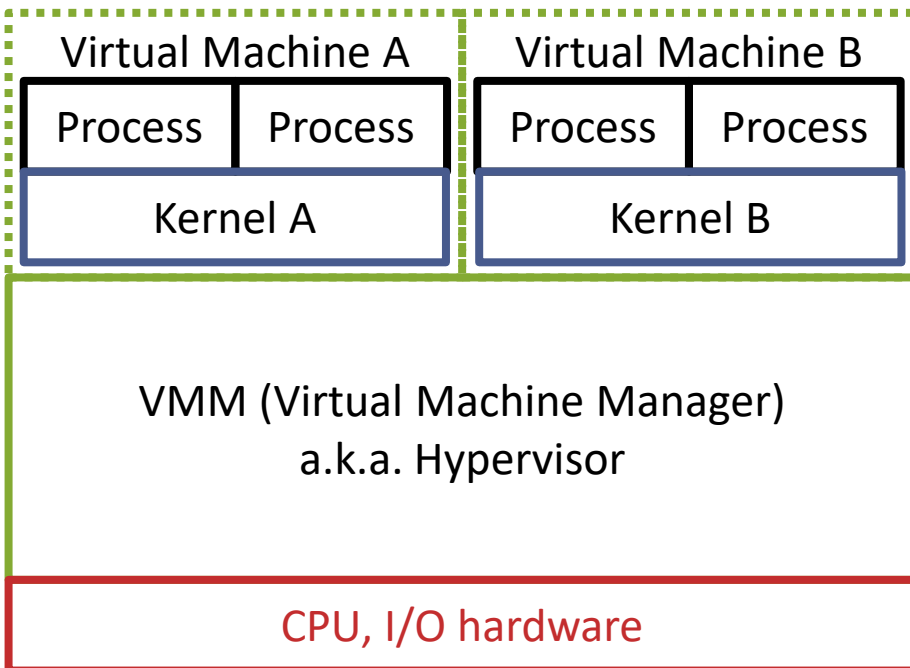
▸ ## Hypervisor integrated in kernel

- ▸ Fast
  - ▪ No need to indirect CPU control via kernel service
- ▸ Complex and dangerous
  - ▪ Kernels were not designed for this

## Traditional type 1 hypervisor

**Example: VMWare ESXi**

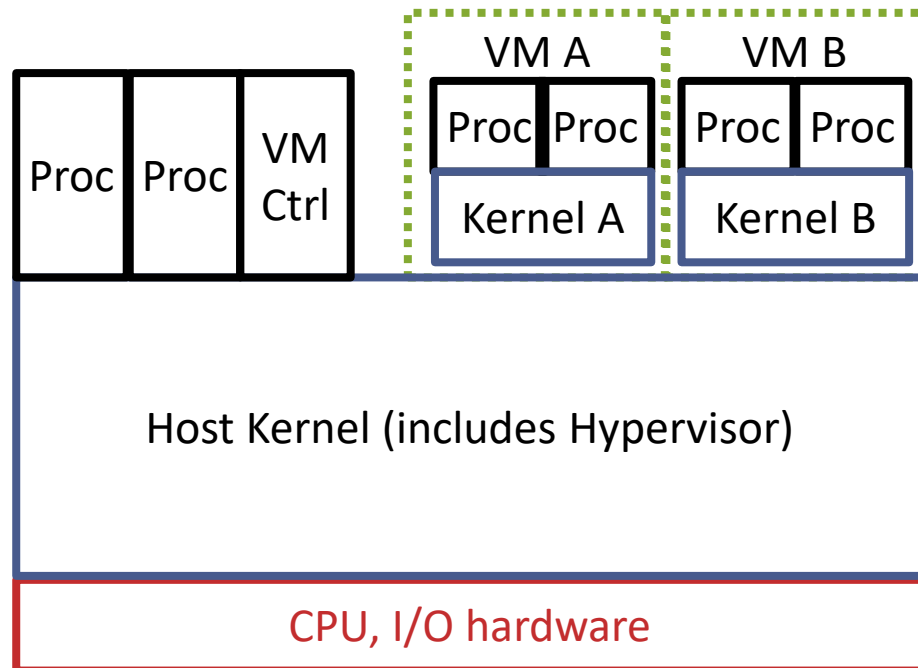| Virtual Machine A | Virtual Machine B |
|---|---|
| Process / Process | Process / Process |
| Kernel A | Kernel B |

VMM (Virtual Machine Manager)
a.k.a. Hypervisor

CPU, I/O hardware

▸ Hypervisor does everything

 ▸ CPU control, time sharing, and I/O in the same project

 ▸ Complex and dangerous

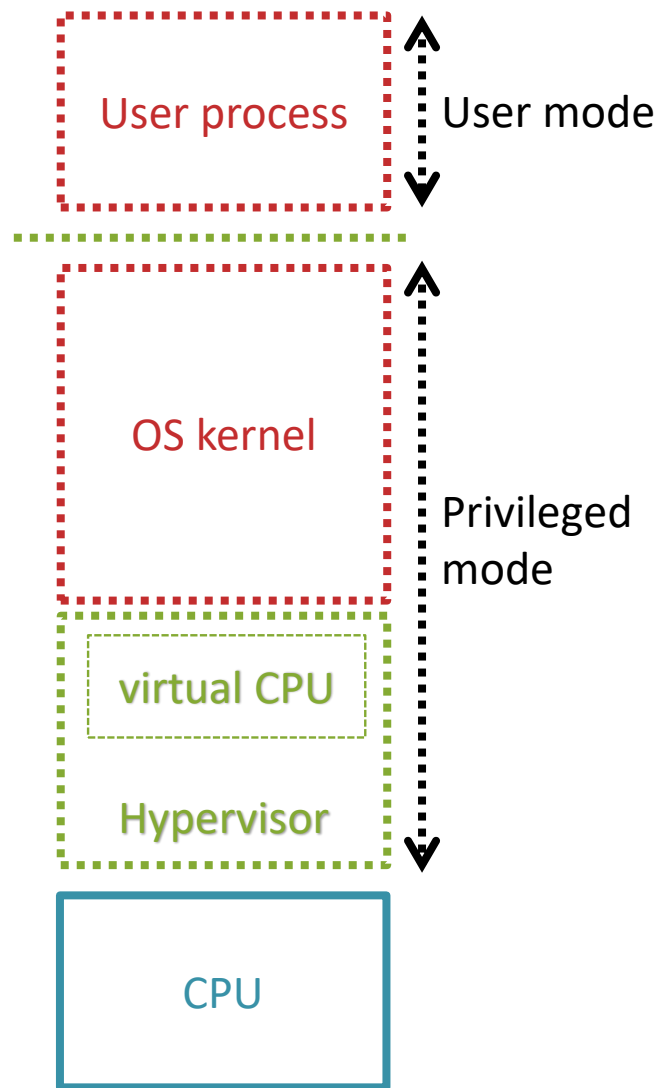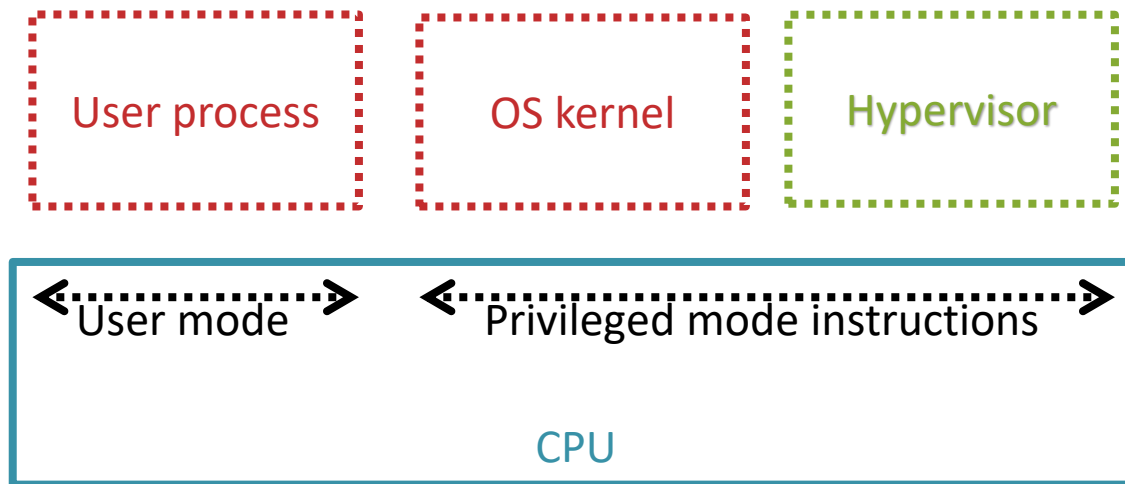## Type 2 implemented in a kernel

**Example: Linux KVM**

| Proc | Proc | VM Ctrl | | VM A | | VM B |
|---|---|---|---|---|---|---|
| | | | | Proc / Proc | | Proc / Proc |
| | | | | Kernel A | | Kernel B |

Host Kernel (includes Hypervisor)

CPU, I/O hardware

▸ Hypervisor implanted in kernel

 ▸ CPU control, time sharing, and I/O in the same project

 ▸ Complex and dangerous

User process — User mode

OS kernel

virtual CPU

Hypervisor
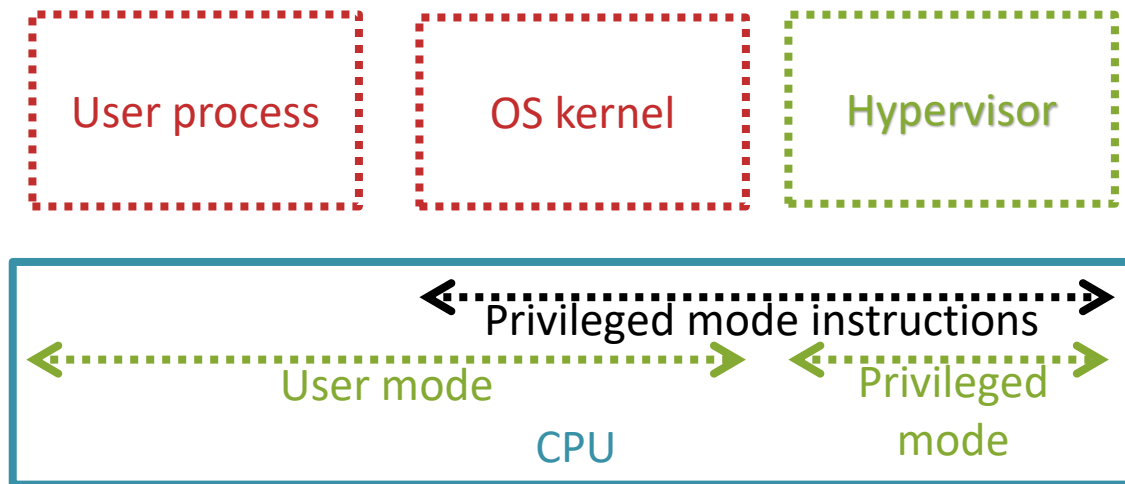
Privileged mode

CPU

- ▶ True virtualization
- ▶ Three layers of software
  - ▶ Guest user processes
    - ▪ Only user-mode CPU instructions
  - ▶ Guest OS kernel
    - ▪ All CPU instructions
    - ▪ Privileged mode expected
      - ▪ But shall not be granted
  - ▶ Hypervisor
    - ▪ All CPU instructions
    - ▪ Exclusive control over the hardware

- ▶ This picture is misleading

User process    OS kernel    Hypervisor
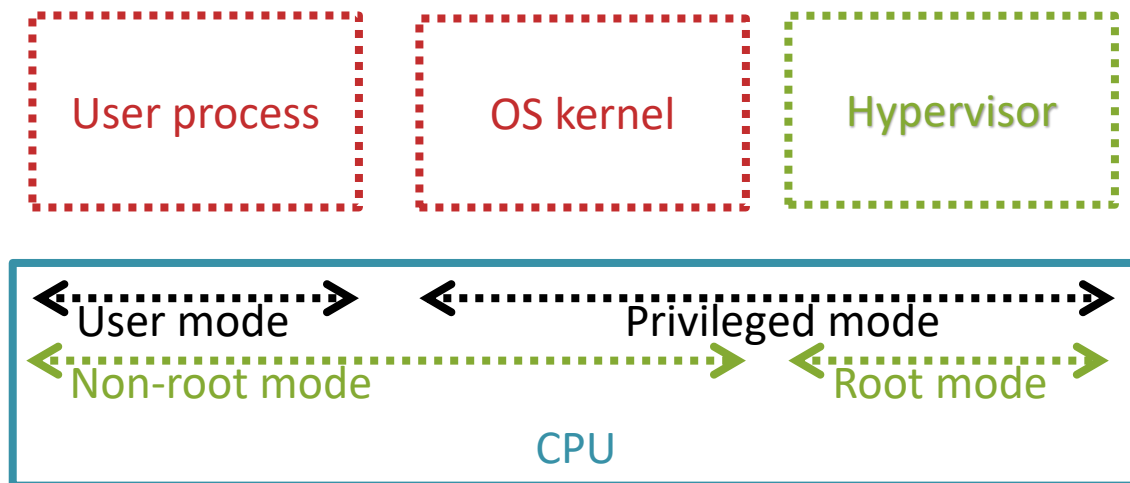
User mode    Privileged mode instructions

CPU

▸ ## The correct picture

▸ All levels of the software directly interact with CPU by executing instructions

▸ OS kernels use privileged instructions

- We must not allow their direct execution (Popek-Goldberg: Safety)

- We must allow direct execution of the other instructions (Performance)

- OS kernel must run with a different privilege setting than the Hypervisor

▸ *Compression of privileges*

- Mapping of 3 privilege levels onto the 2 levels available in typical CPU

| User process | OS kernel | Hypervisor |
|---|---|---|

Privileged mode instructions

User mode ← → Privileged mode

CPU

## Trap-and-emulate (IBM 1972)

- OS kernels use privileged instructions but run in the user mode
  - Every privileged instruction in the kernel causes a *trap* (synchronous interrupt)
  - The hypervisor emulates the instruction
  - The emulation allows verification of access rights, virtualization etc.
- Performance considerations
  - Every syscall goes through hypervisor
  - Every I/O instruction in kernel is emulated
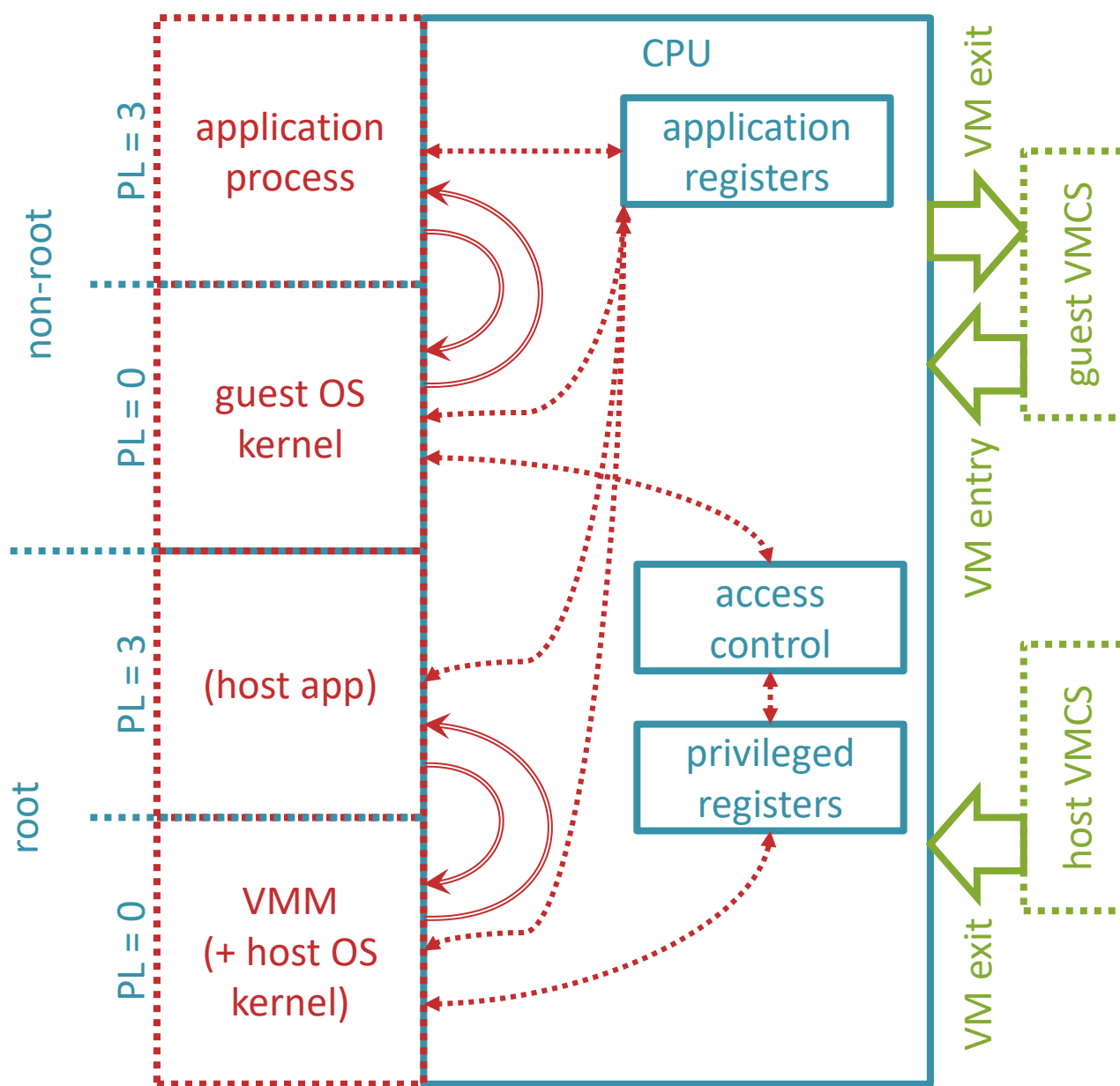    - S/370 had "channel programs" = single I/O instruction started the whole I/O operation

User process     OS kernel     Hypervisor

User mode     Privileged mode

Non-root mode     Root mode

CPU

▸ **Intel/AMD root/non-root modes**

  ▸ **Control transferred between levels**

    ▪ User -> Kernel: Simultaneously with switching CPU to the privileged mode

      ▪ SYSCALL, some synchronous (software) interrupts

    ▪ User/Kernel -> Hypervisor: VM Exit event = switch to hypervisor mode

      ▪ Asynchronous (hardware) interrupts

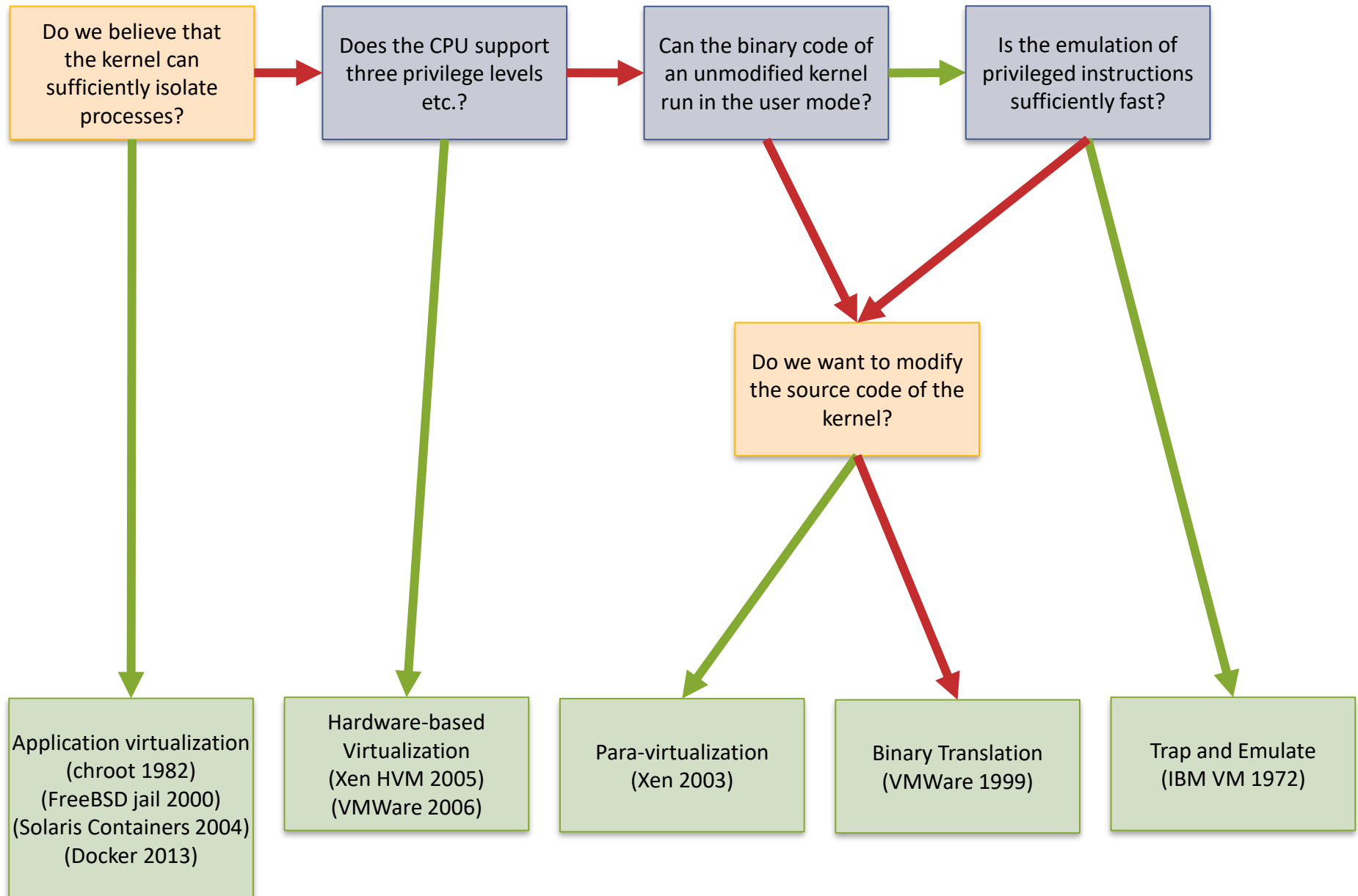      ▪ Some synchronous (software) interrupts (e.g. page faults)

# Hardware support for virtualization – a new dimension of privilege



- ▸ Intel VT-x / AMD-V
  - ▸ Details differ
- ▸ „Root" mode
  - ▸ Like a CPU without virtualization
  - ▸ Usable to run the host OS
- ▸ „Non-root" mode
  - ▸ Limited access to the privileged of the CPU state
  - ▸ Unwanted actions cause „VM exit"
- ▸ Mode switch
  - ▸ A part of the CPU state is read from/stored to memory
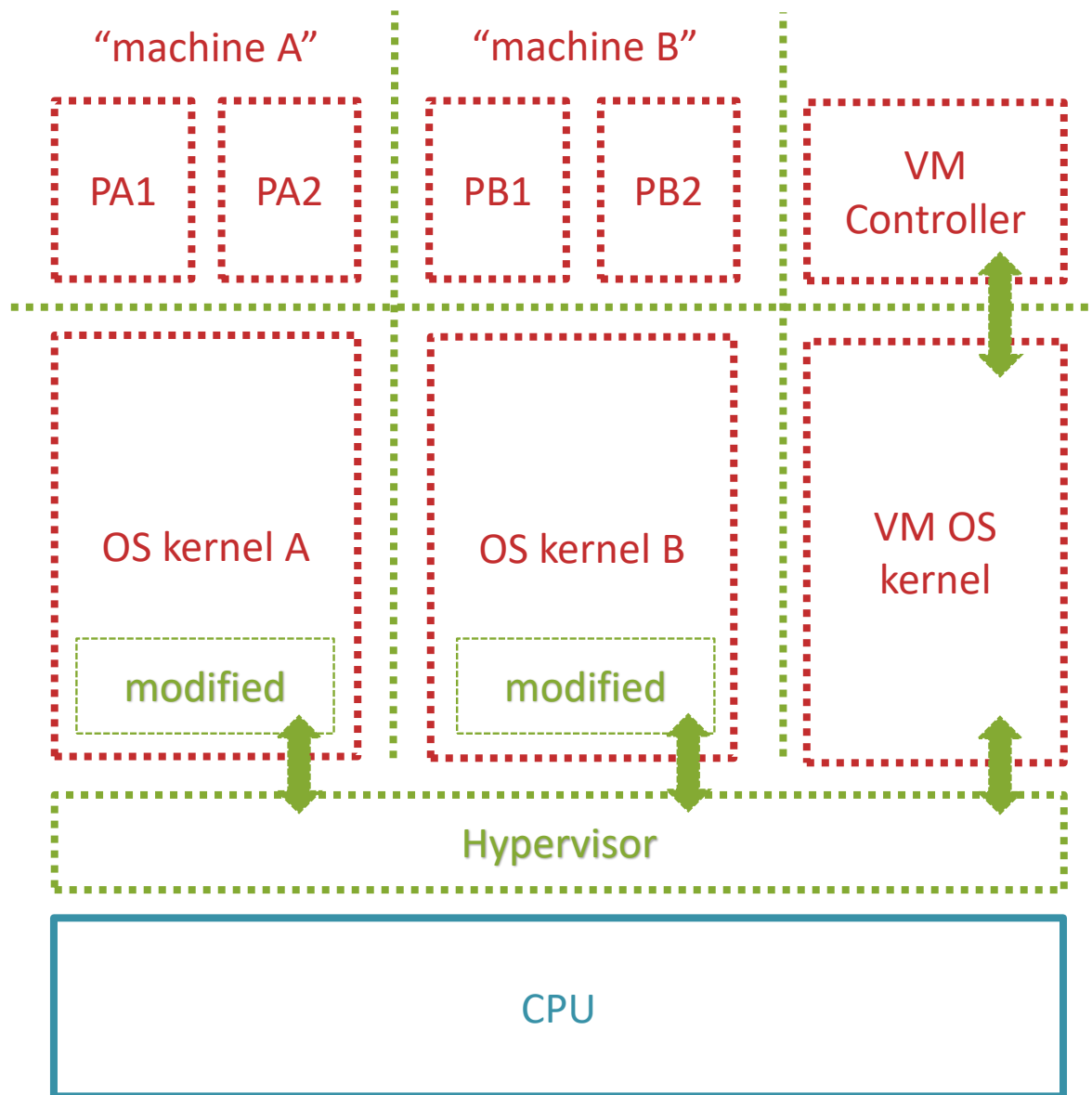  - ▸ Address-space switch included

# Approach to virtualization

Do we believe that the kernel can sufficiently isolate processes?

Does the CPU support three privilege levels etc.?

Can the binary code of an unmodified kernel run in the user mode?

Is the emulation of privileged instructions sufficiently fast?

Do we want to modify the source code of the kernel?

Application virtualization
(chroot 1982)
(FreeBSD jail 2000)
(Solaris Containers 2004)
(Docker 2013)

Hardware-based Virtualization
(Xen HVM 2005)
(VMWare 2006)

Para-virtualization
(Xen 2003)

Binary Translation
(VMWare 1999)

Trap and Emulate
(IBM VM 1972)

## The second era of VM

- 1999 – VMWare Workstation
  - Software virtualization (Binary Translation)
  - VMM as an app in Windows NT
- 2002 – VMWare ESX Server
  - VMM replaces the host OS
- 2003 – Xen
  - Para-virtualization
    - Host OS modified (in source code)
- 2007 – Linux KVM
  - VMM integrated into the OS kernel
- 2008 – Microsoft Hyper-V
  - VMM cooperates with the host OS
    - Non-cooperating guest OS possible

## The x86 is unsuitable for VM

- Legacy of the Intel 80286 CPU
  - 1982 – still in the first era of VM
- The first mitigation attempts
  - 2005 – Intel VT-x
  - 2006 – AMD-V
- Gradually improved performance
  - Improved HW support
  - Para-virtualization in critical OS parts
- Performance loss is now insignificant for most applications
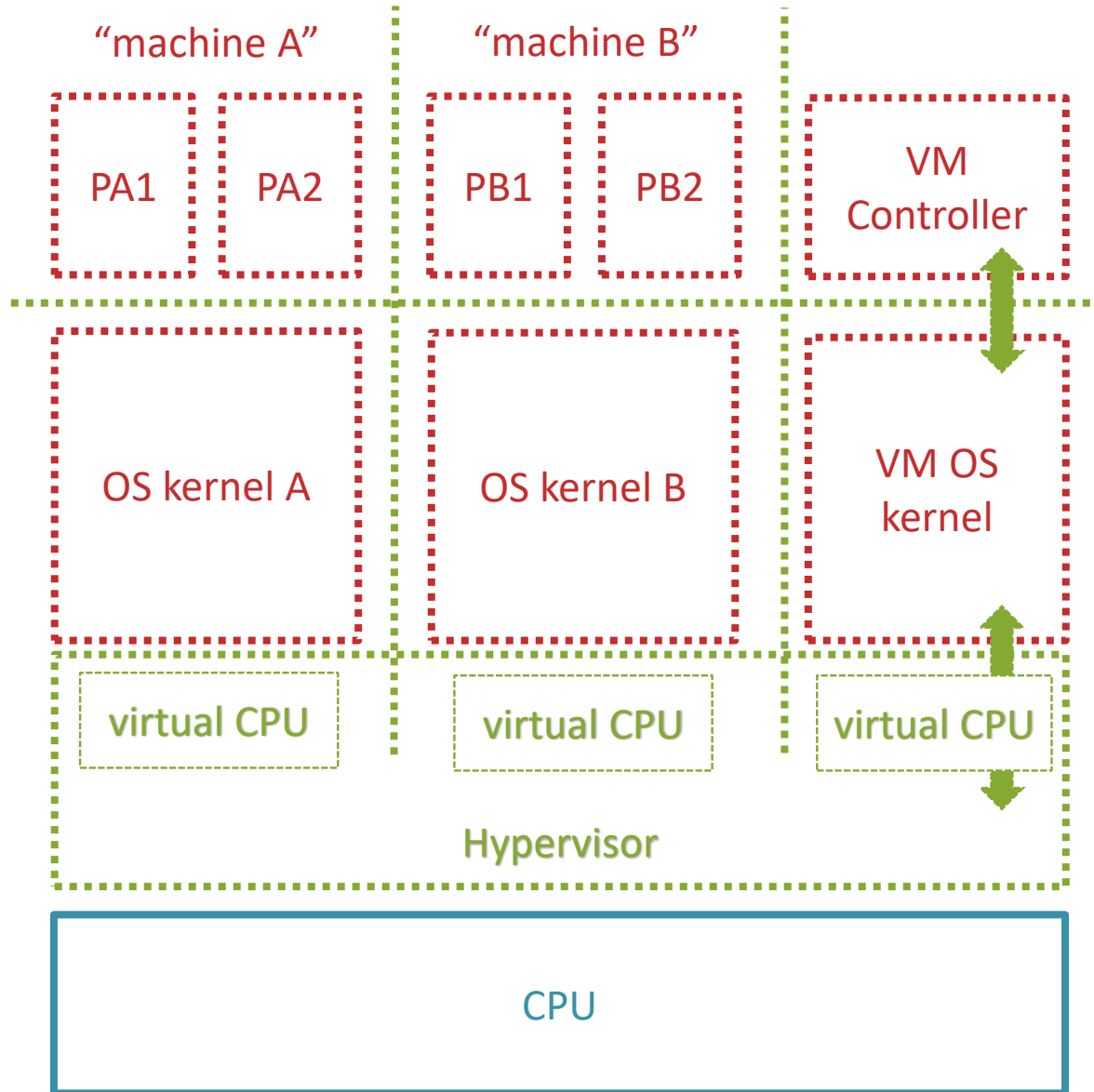  - Variations in performance too big for Real-Time applications and performance measurement

# VMM implementation

# Para-virtualization

| "machine A" | | "machine B" | | |
|---|---|---|---|---|
| PA1 | PA2 | PB1 | PB2 | VM Controller |
| OS kernel A | | OS kernel B | | VM OS kernel |
| modified | | modified | | |

Hypervisor

CPU

- ▸ Para-virtualization
  - ▸ Lower layers of OS kernels are modified
  - ▸ Instead of controlling hardware, these layers call the hypervisor
- ▸ Hypervisor (VMM)
  - ▸ Creates an illusion of a machine dedicated for each kernel
  - ▸ The illusion is not perfect; difficult parts replaced by cooperation of the modified kernel
- ▸ VM controller
  - ▸ Provides administrator control
- ▸ VM OS kernel
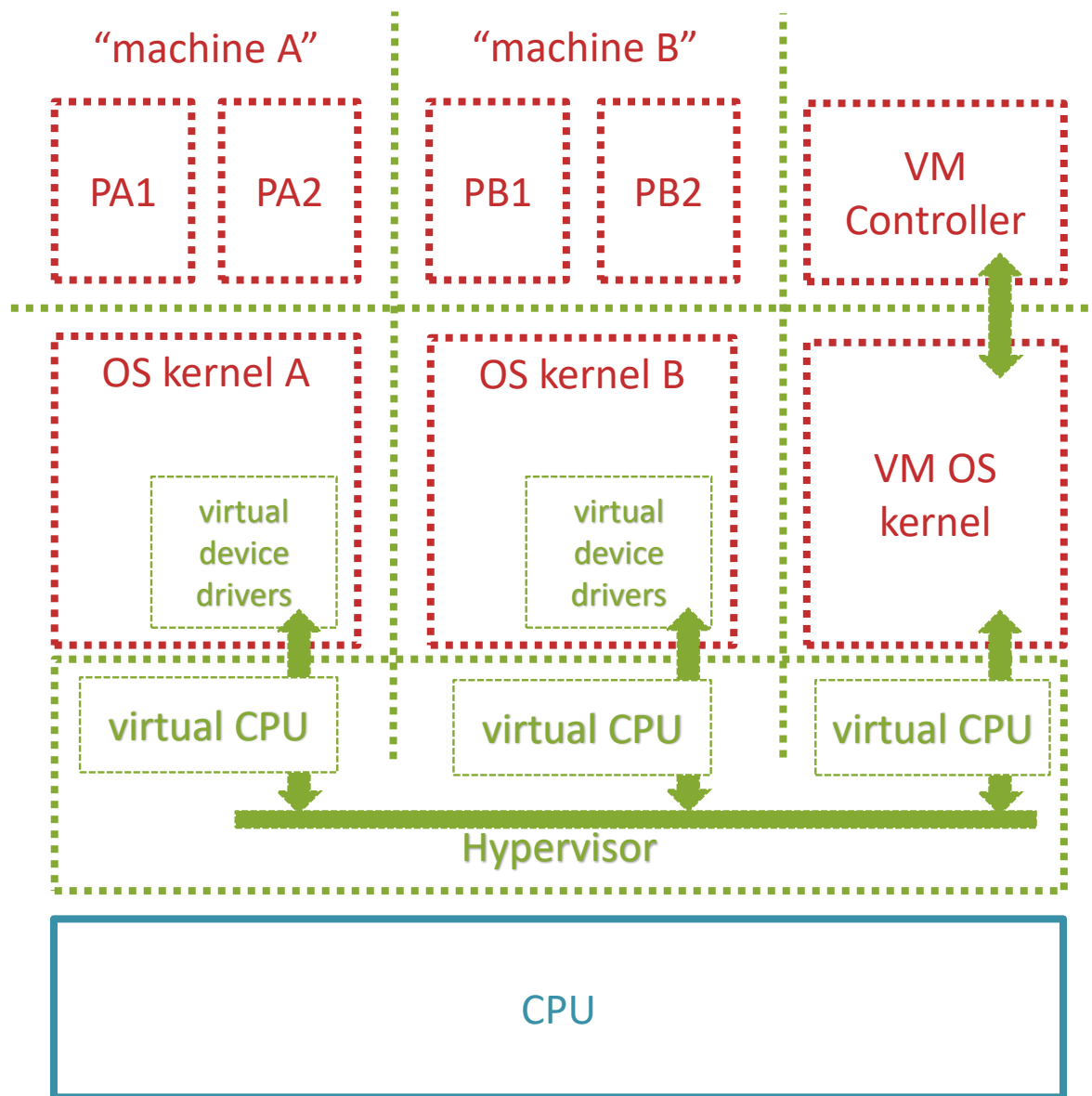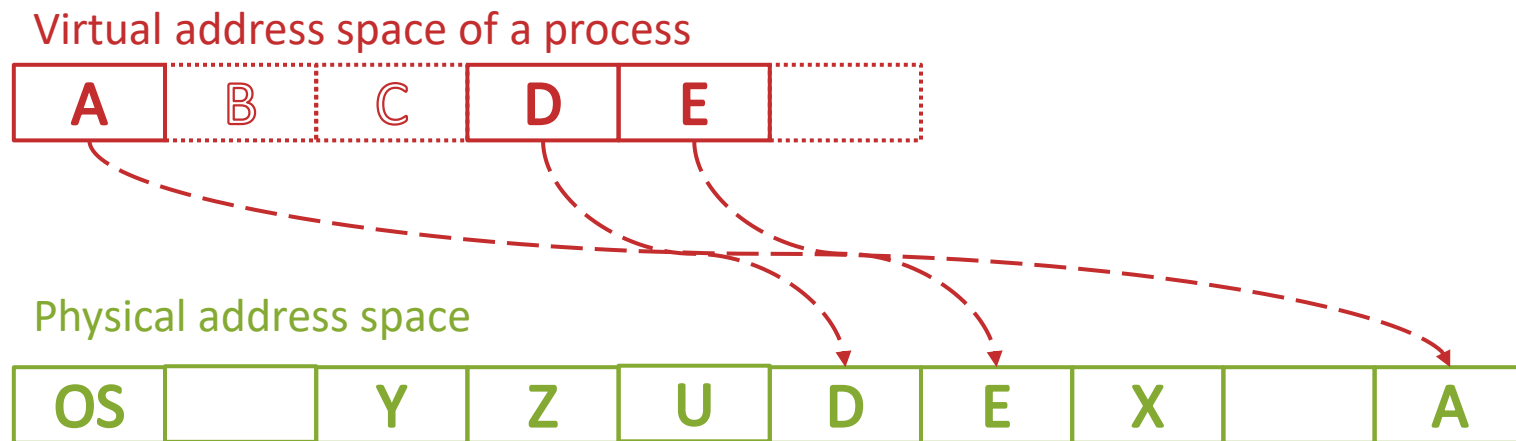  - ▸ Provides file and network services for the controller and hypervisor

# True virtualization



- ▶ True virtualization
  - ▶ OS kernels directly work with virtual CPUs and other HW
- ▶ Hypervisor (VMM)
  - ▶ Creates an illusion of a machine dedicated for each kernel
  - ▶ The illusion is perfect, emulating every bit of CPU and other HW
  - ▶ Modern physical CPUs help creating this illusion
- ▶ VM controller
  - ▶ Provides administrator control
- ▶ VM OS kernel
  - ▶ Provides file and network services for the controller and hypervisor

- ▸ **Hardware support makes CPU virtualization "easy"**
  - ▸ Negligible overhead
  - ▸ Implementing a hypervisor is still a tremendous task
- ▸ **This does not apply to most other HW**
  - ▸ OS kernels "modified" by the means of device drivers
  - ▸ Actions forwarded to the VM OS
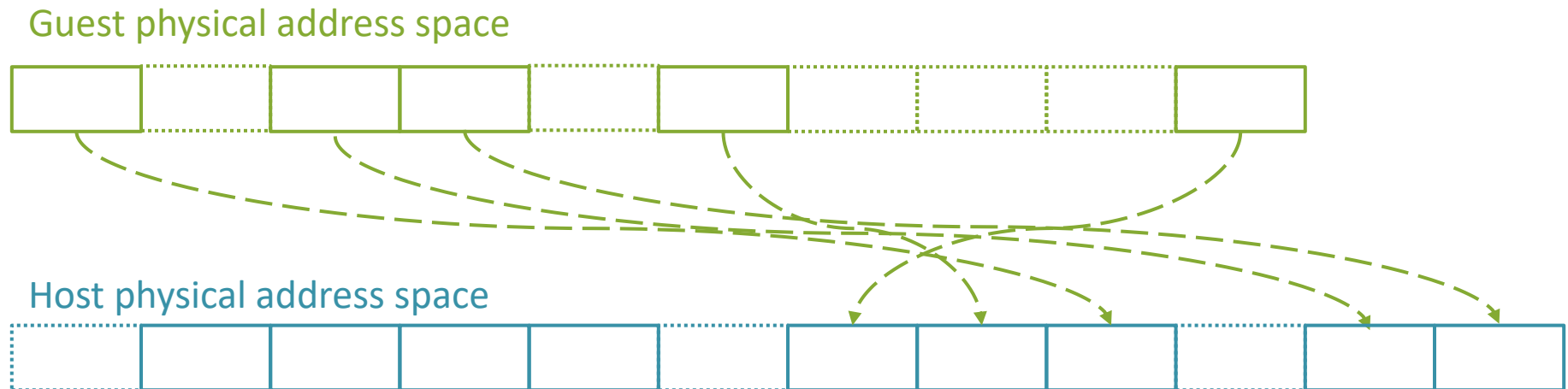- ▸ **Fast communication infrastructure**
  - ▸ Implemented in the Hypervisor

# Virtualization of Virtual Memory

# Virtual memory in a physical computer

Virtual address space of a process

| A | B | C | D | E | |
|---|---|---|---|---|---|

Physical address space

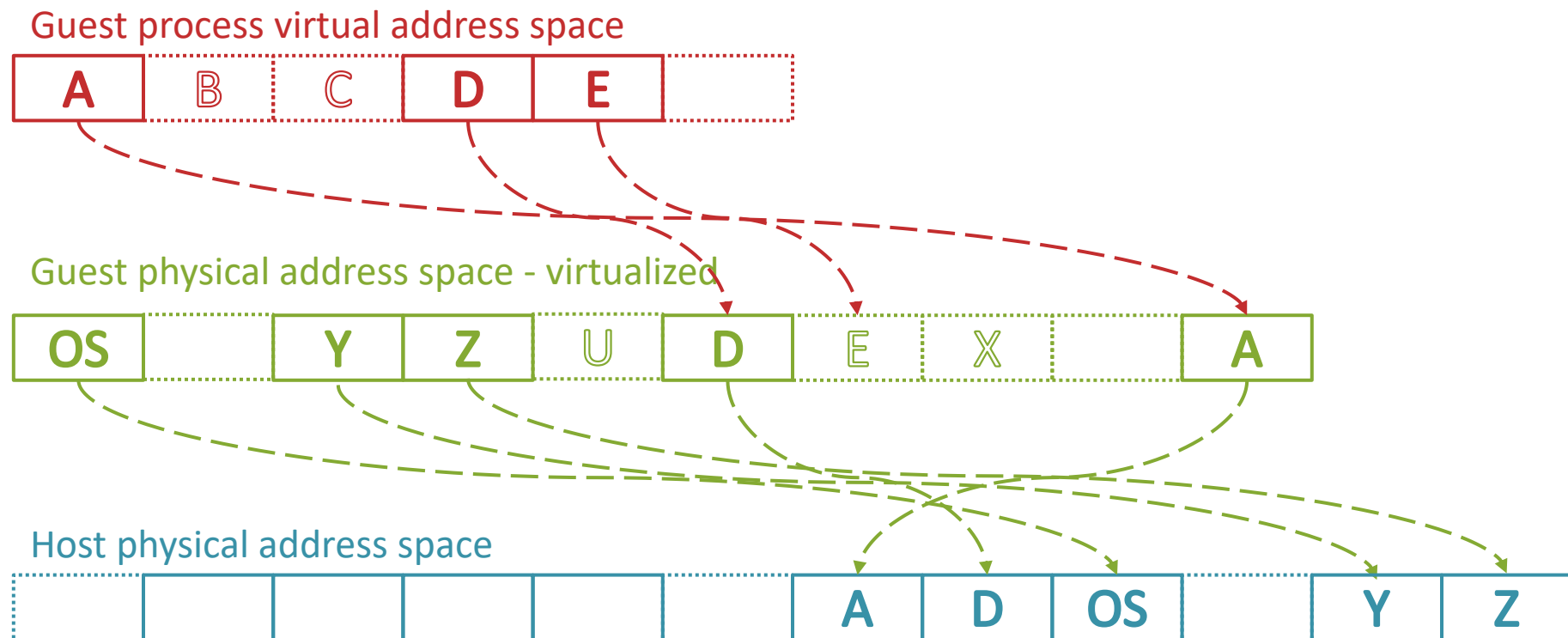| OS | | Y | Z | U | D | E | X | | A |
|----|--|---|---|---|---|---|---|--|---|

- ▸ All code works in a virtual address space, including the OS kernel
- ▸ OS defines the mapping of virtual to physical addresses (including for itself)
  - ▸ Intel/AMD: 2 to 5 levels of page tables, stored in the physical memory
    - ▪ CPU translates addresses using the TLB in most cases
    - ▪ On a TLB miss, the CPU will read the page tables to fill the TLB
    - ▪ On a page-table miss, the CPU will wake-up the OS by executing a synchronous interrupt

Guest physical address space

Host physical address space

▸ **The hypervisor must allow co-existence of several VMs**

▸ **The physical address space of each VM is virtualized**

  ▸ The mapping is defined by the hypervisor

  ▸ An equivalent of page mapping by an OS

Guest process virtual address space
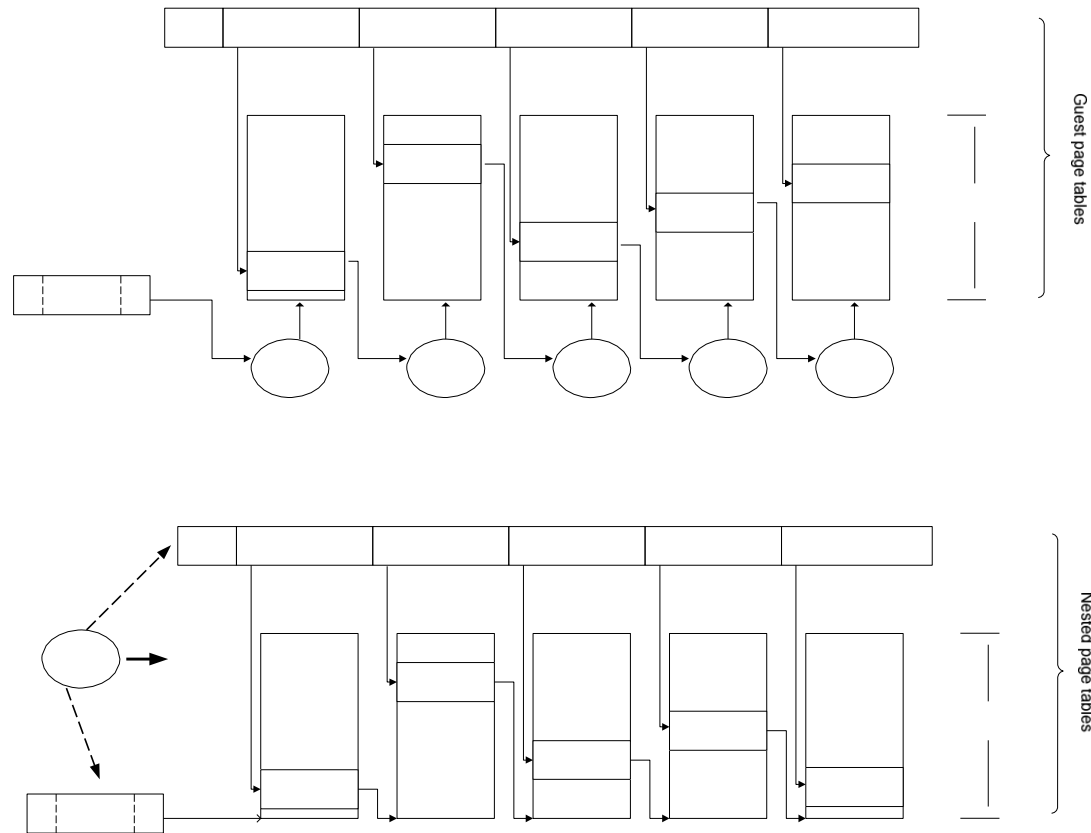
| A | B | C | D | E | |
|---|---|---|---|---|---|

Guest physical address space - virtualized

| OS | | Y | Z | U | D | E | X | | A |
|----|--|---|---|---|---|---|---|--|---|

Host physical address space

| | | | | | | A | D | OS | | Y | Z |
|--|--|--|--|--|--|---|---|----|--|---|---|

▶ A composition of two mappings

  ▶ The mapping defined by the guest OS for a process
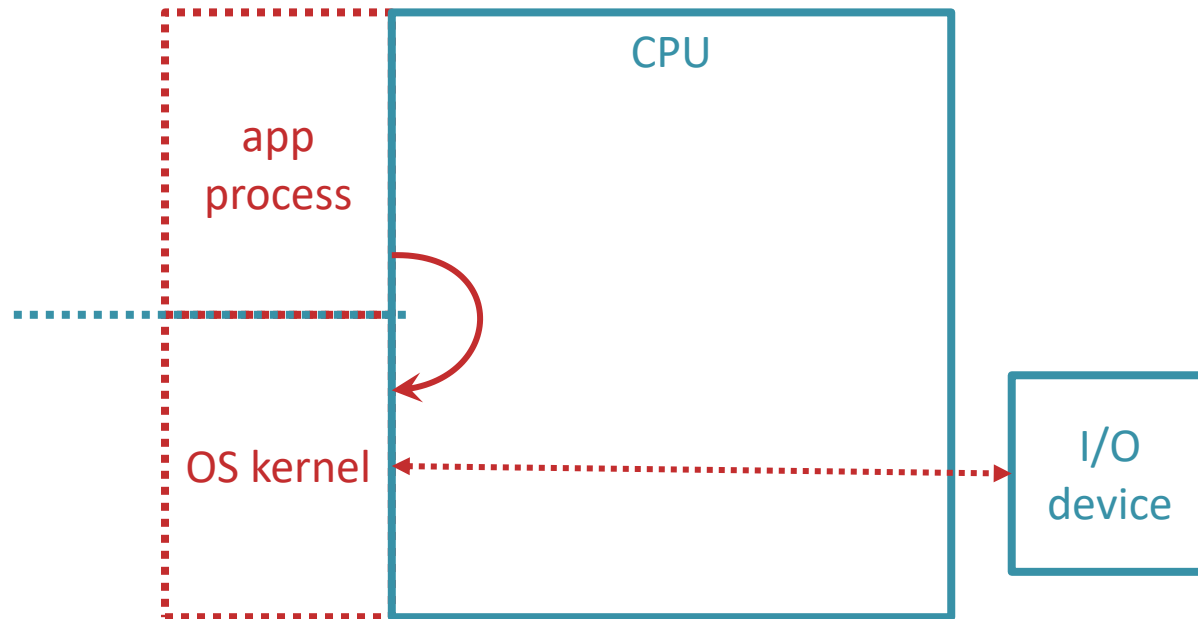
  ▶ The guest-host mapping defined by the hypervizor

▶ **The CPU holds two mappings and performs the composition**

- ▶ Each (guest-physical) address in the GPT (starting with the CR3) must be translated by the NPT (to a host-physical address)
  - ▪ For 5 levels, 25 memory-accesses required!
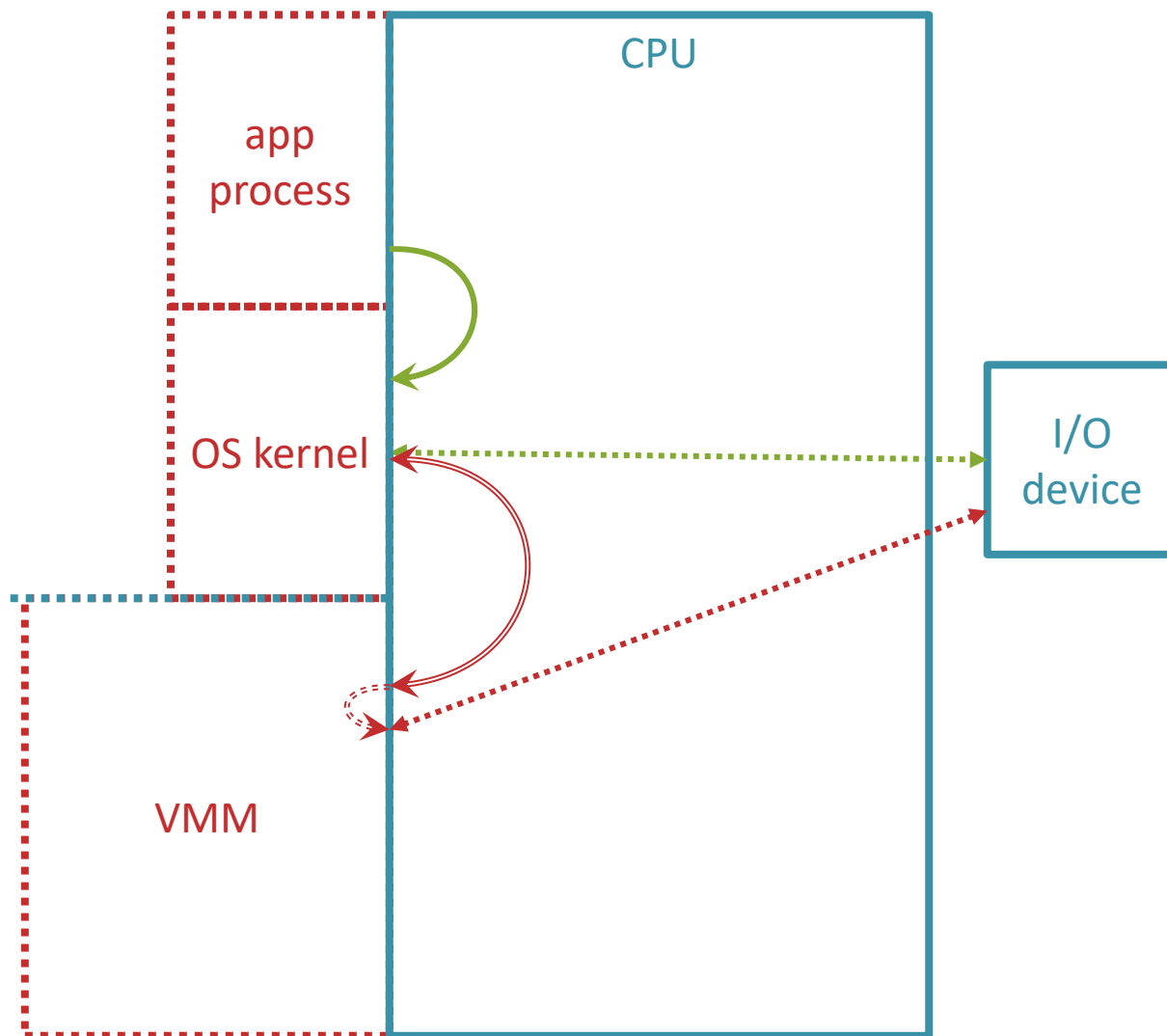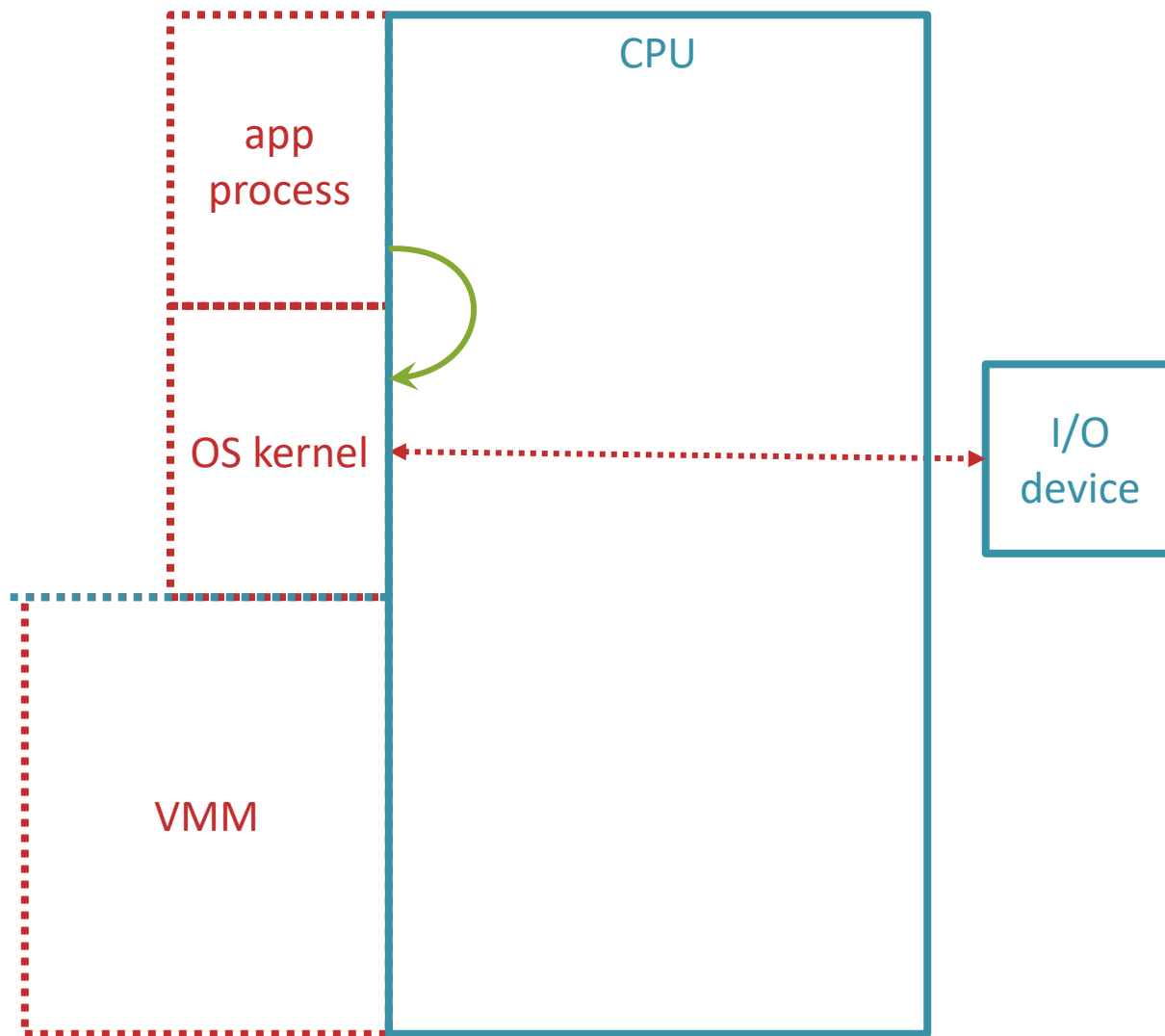- ▶ The TLB stores the composed mapping

# Virtualization of I/O

- App processes must perform I/O by invoking an kernel syscall
- The OS kernel communicates with the I/O device
  - I/O instructions (privileged), or
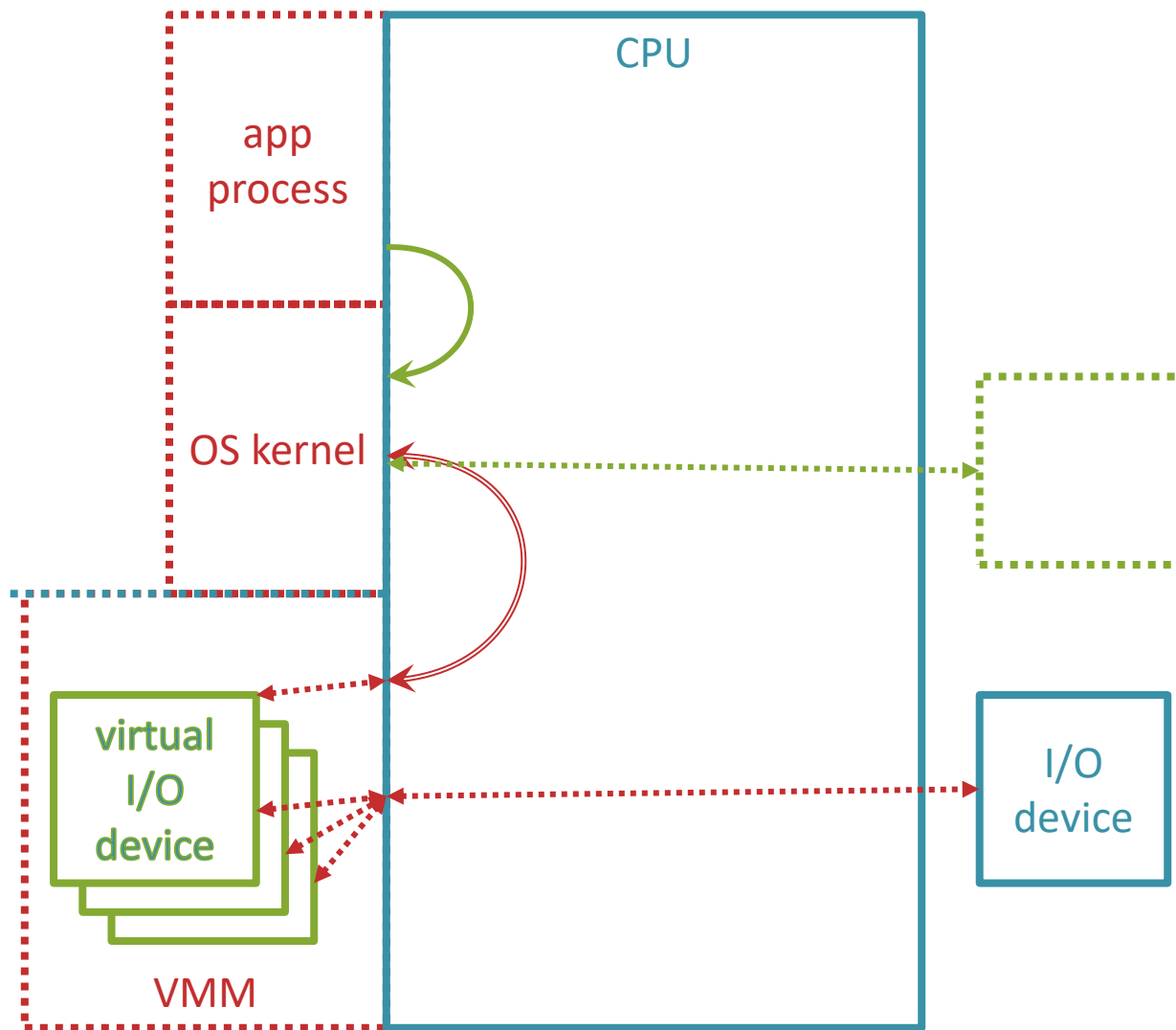  - Memory-mapped I/O device (protected by virtual memory mapping)

- ‣ Privileged I/O instructions cause synchronous interrupts
- ‣ Executed by an instruction emulator in the VMM
- ‣ Besides the I/O device, the related interrupt system and/or DMA controller must also be virtualized

- ‣ Exclusive mode
  - ‣ Only one VM can access the device

# I/O access in a VM – exclusive mode
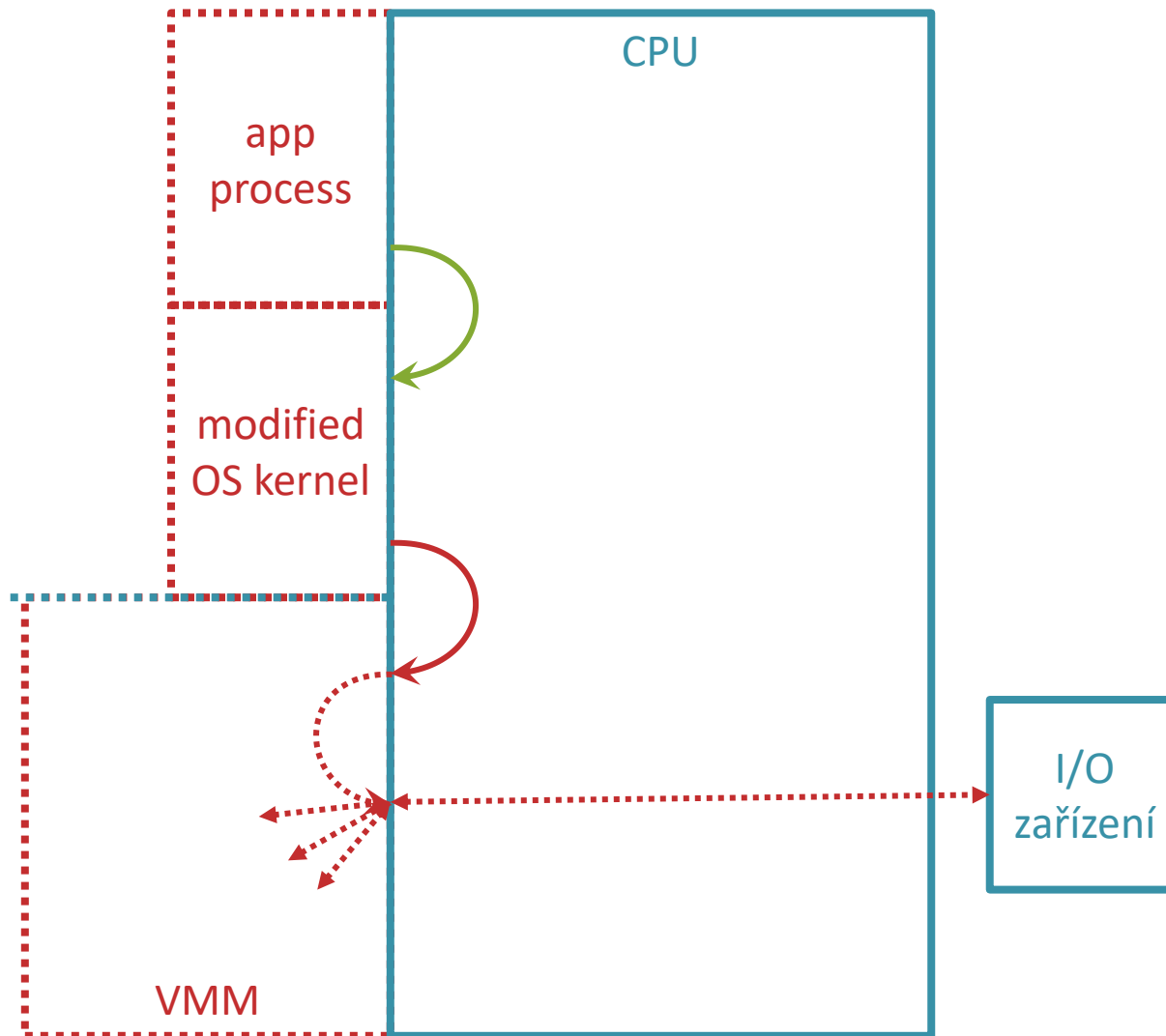


CPU

app
process

OS kernel

VMM

I/O
device

▸ When the communication with the device (including the DMA etc) is possible using non-privileged instructions

  ▪ Memory-mapped devices, or

  ▪ configurable access into I/O address space

▸ Exclusive mode

▸ Only one VM can access the device
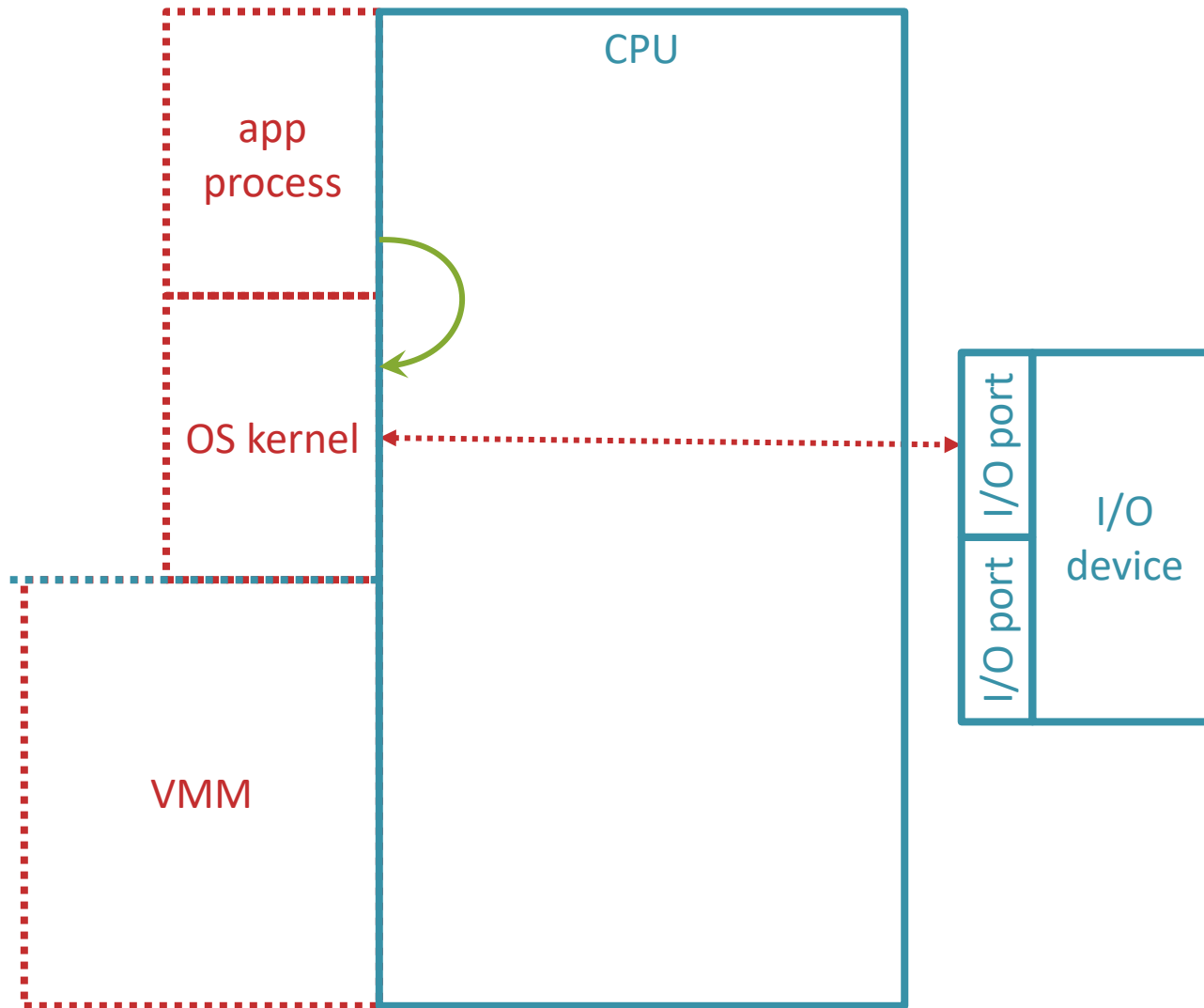
▸ Suitable for the host-OS running in a privileged VM

- Access to the I/O device caught and emulated by the VMM
- The emulation state is independent for each VM
- The emulated type of hardware does not have to exactly match the physical device

▶ Shared mode

- VMM extracts logical actions from the emulated virtual devices
- The logical actions are performed by the physical device

Diagram labels: app process, CPU, OS kernel, virtual I/O device, VMM, I/O device

- ▸ Guest OS modified
  - ▸ Modified source code, or
  - ▸ a device driver for a non-existent device
- ▸ Advantages
  - ▸ The modified guest OS sends logical commands instead of physical I/O
  - ▸ Emulation of I/O instructions not needed
  - ▸ Single logical command instead of a sequence of I/O instructions
  - ▸ Synchronization of logical commands from different VMs is simpler

Diagram labels: CPU, app process, modified OS kernel, VMM, I/O zařízení

# I/O access in a VM – multi-port I/O device



app
process

CPU

OS kernel

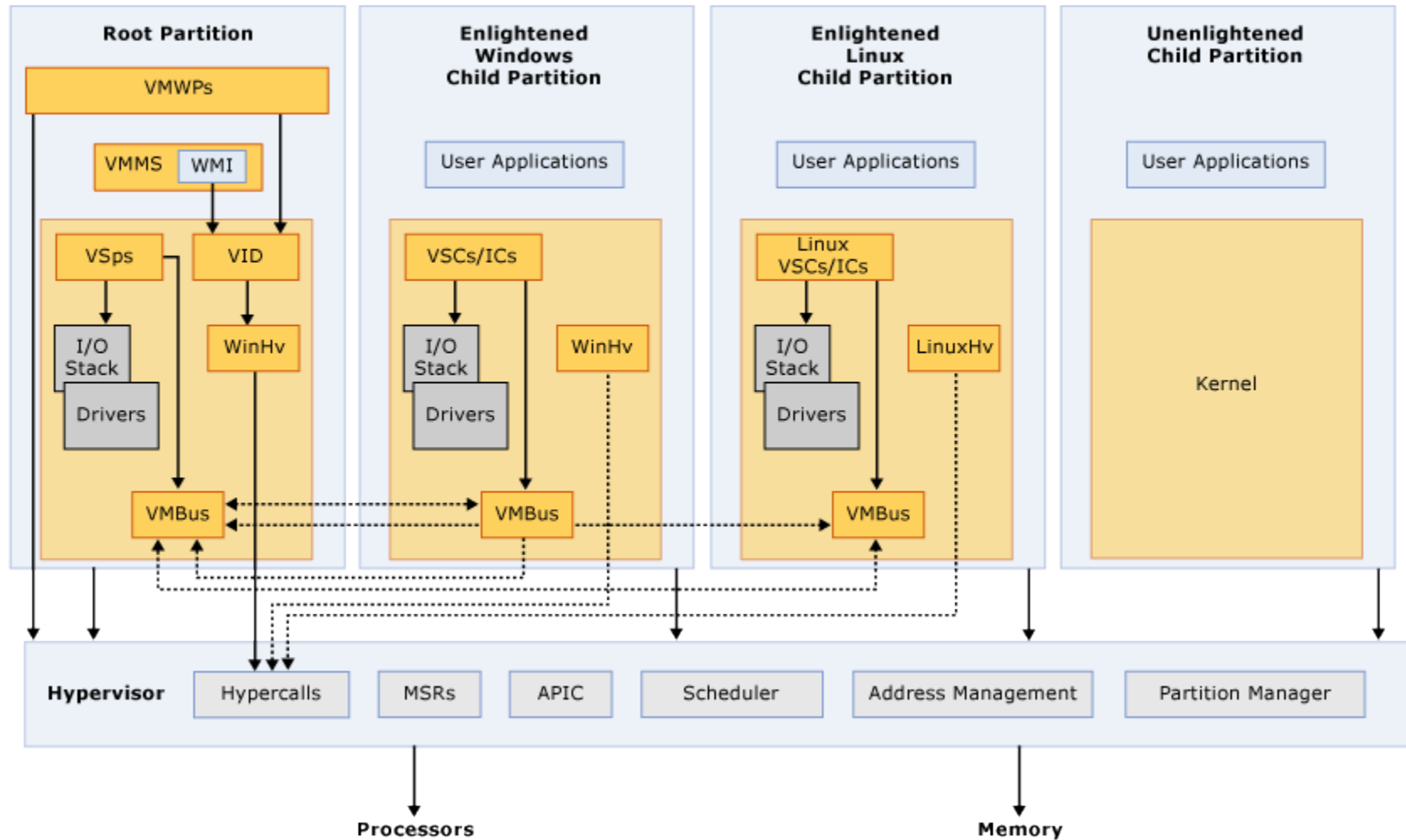I/O port

I/O port

I/O
device

VMM

- ▸ Non-privileged access to I/O
- ▸ Configurable I/O space protection required
- ▸ **Shared mode**
  - ▸ The I/O device presents itself more than once in the I/O space - ports
  - ▸ The I/O device maintains an independent state for each port
  - ▸ The I/O device synchronizes logical commands from the ports on the shared physical device

  - ▸ Expensive hardware
    - ▪ Mostly NICs

# VM-VMM Communication
## (Example: Microsoft Hyper-V)

Hyper-V High Level Architecture

## Partition

- A set of virtual processors and other hardware, plus its configuration
- Root partition – typically used to run the Host OS and VM Management

## Inter-partition messaging

- The hypervisor supports a simple message-based inter-partition communication mechanism.
- Messages can be sent by the hypervisor to a partition or can be sent from one partition to another.

## Guest Physical Address Space

- The GPA mappings are defined by the partition's parent.
  - At the time they are mapped, they are specified in terms of the parent's GPA space.

## Guest Virtual Address Space

- The hypervisor exposes operations to flush the TLB (on one virtual processor).

# Guest-to-Hypervisor Interface (Microsoft Hyper-V)

- ▸ Virtual MSRs
  - ▪ Physical MSRs used by Kernels to read/alter CPU configuration
  - ▸ VMM emulates additional Machine Status Registers (MSR) not present in HW
    - ▪ VMM-aware VM Kernel can read/write virtual MSRs to exchange configuration information with VMM
  - ▸ Emulation too slow for real communication

- ▸ Hypercall
  - ▸ Call Hypervisor from Guest (privileged mode)
  - ▸ Exposed as procedure call to a special guest-physical page
    - ▪ Provided by Hypervisor on request from Guest (via a virtual MSR)
    - ▪ VM Kernel must map the guest-physical page to a guest-virtual page
    - ▪ The page contains either special instructions or nothing – both cases cause VM exit
  - ▸ Arguments passed/returned in registers or VPAP

- ▸ Virtual Processor Assist Page (VPAP)
  - ▸ Special guest-physical page per virtual processor (core/logical thread)
    - ▪ Both Hypervisor and Guest can read/write

## ▶ Hypercall

- ▶ Call Hypervisor from Guest (privileged mode)
- ▶ Exposed as procedure call to a special guest-physical page
- ▶ Arguments passed/returned in registers or VPAP

- ▶ One Hypercall may serve several logical requests
  - ▪ Chained into an array of arguments

- ▶ All Hypercalls return within 50 microseconds
  - ▪ Avoids blocking in the Hypervisor (giant lock?)
  - ▪ Longer requests serviced in continuation-style
    - ▪ The Hypercall return address is set before the instruction that invoked it
    - ▪ Arguments adjusted to indicate that part of the request is already done
    - ▪ On the next VM Entry, the Hypercall is entered again