



KVM/ARM: The Design and Implementation of the Linux ARM Hypervisor

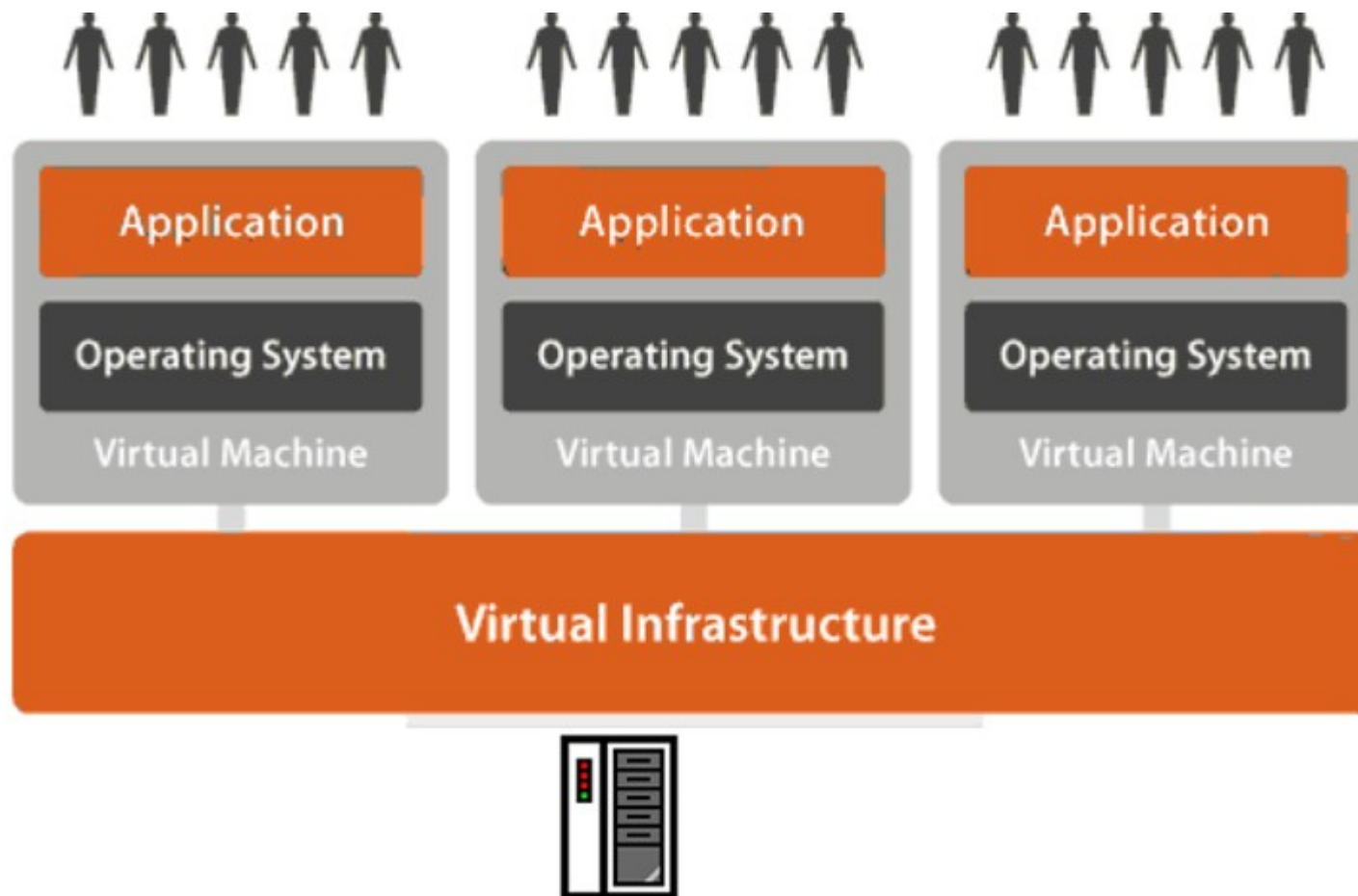
Fall 2014

Presented By: Probir Roy



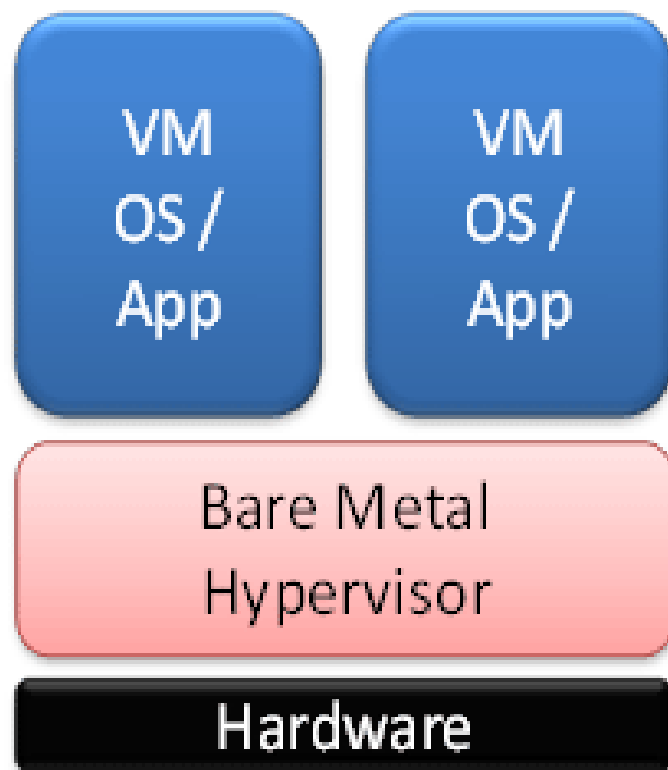
Virtualization & Hypervisor

WILLIAM
& MARY



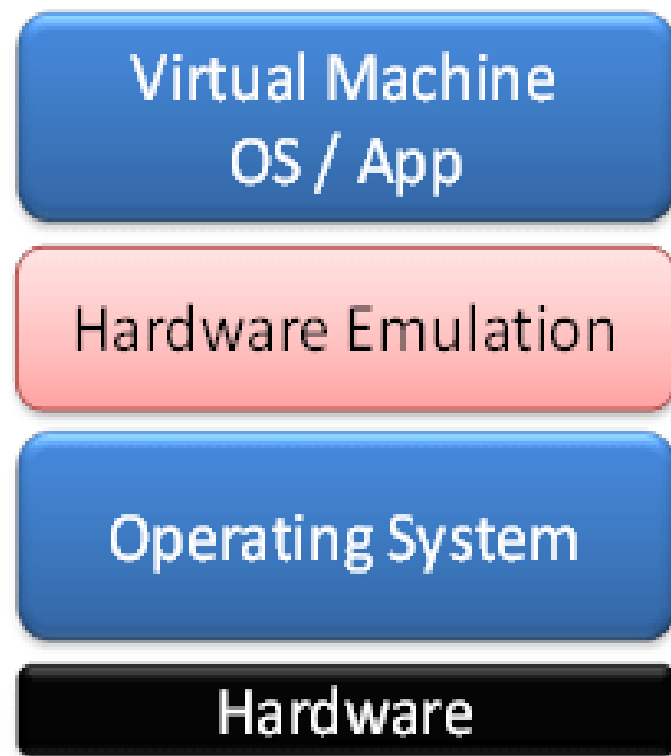


Which type of Hypervisor is better?



Bare Metal / Native Hypervisor

Xen, VMWare ESX

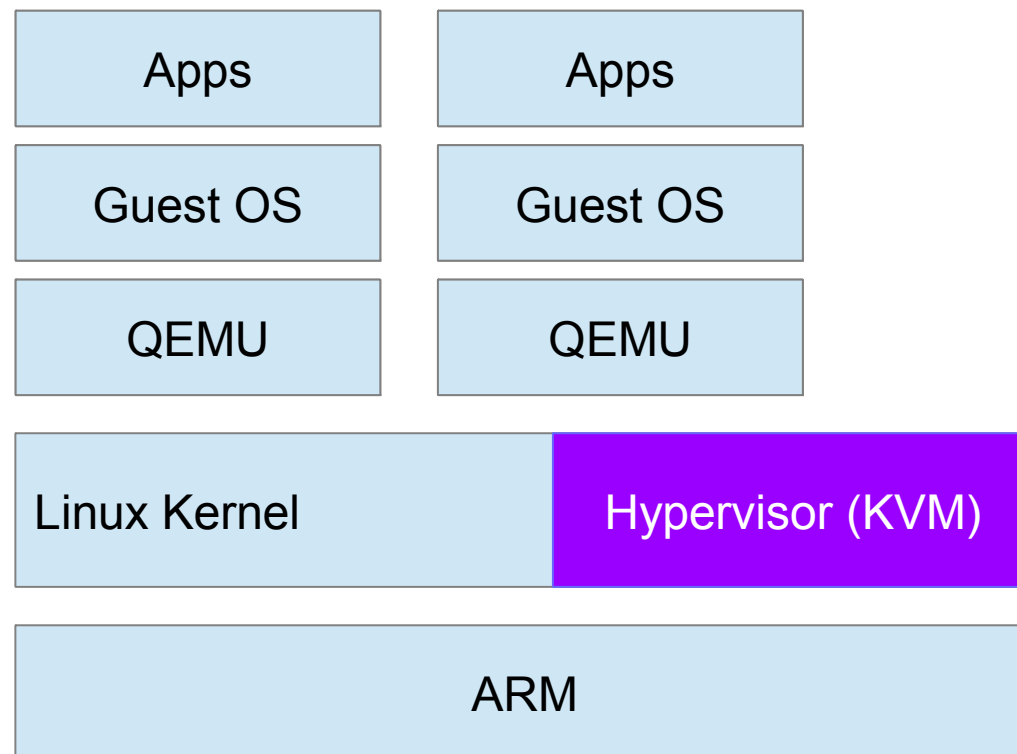


Hosted Hypervisor

VMWare Workstation



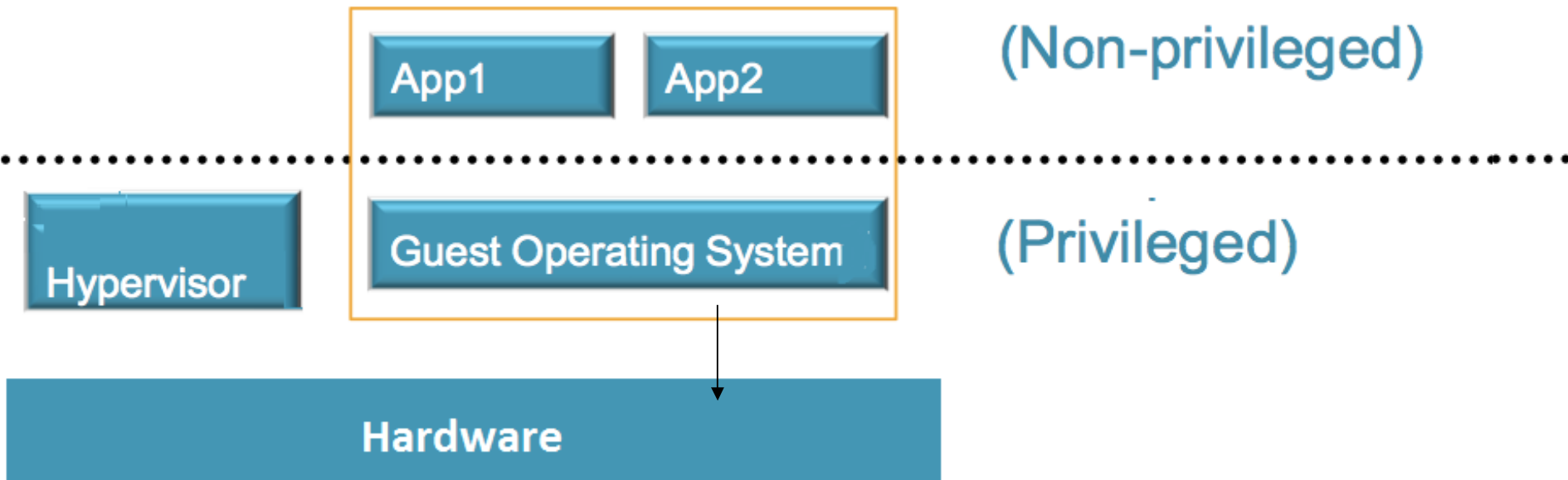
KVM/ARM: The Design and Implementation of the Linux ARM Hypervisor



Full Virtualization



Prior ARMV7



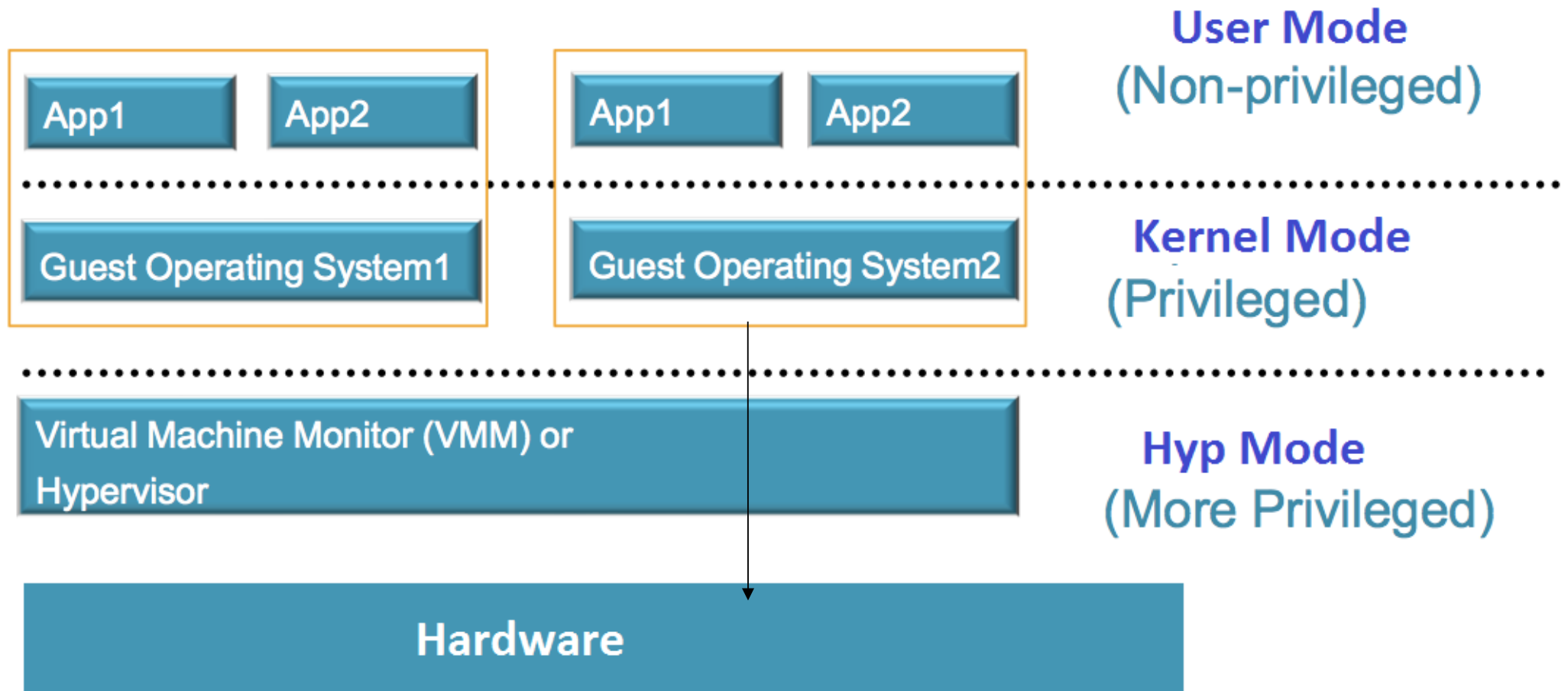
Not Classically Virtualizable



Trap & Emulate



ARM V7 Virtualization Extension

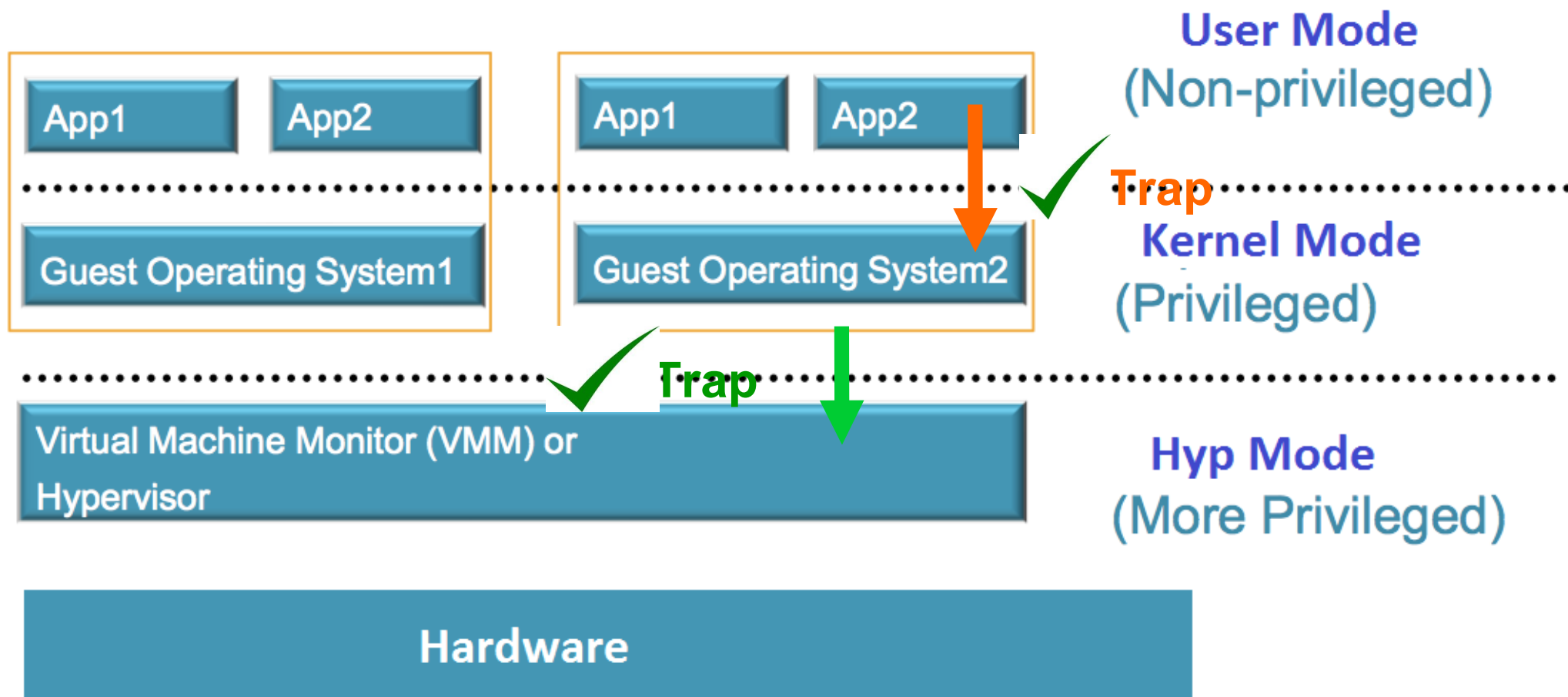


Trap & Emulate



CPU Virtualization

Why trapping at Kernel Mode is useful?



From Hyp Mode, hardware Configurable to trap sensitive instructions and interrupt to Hyp Mode

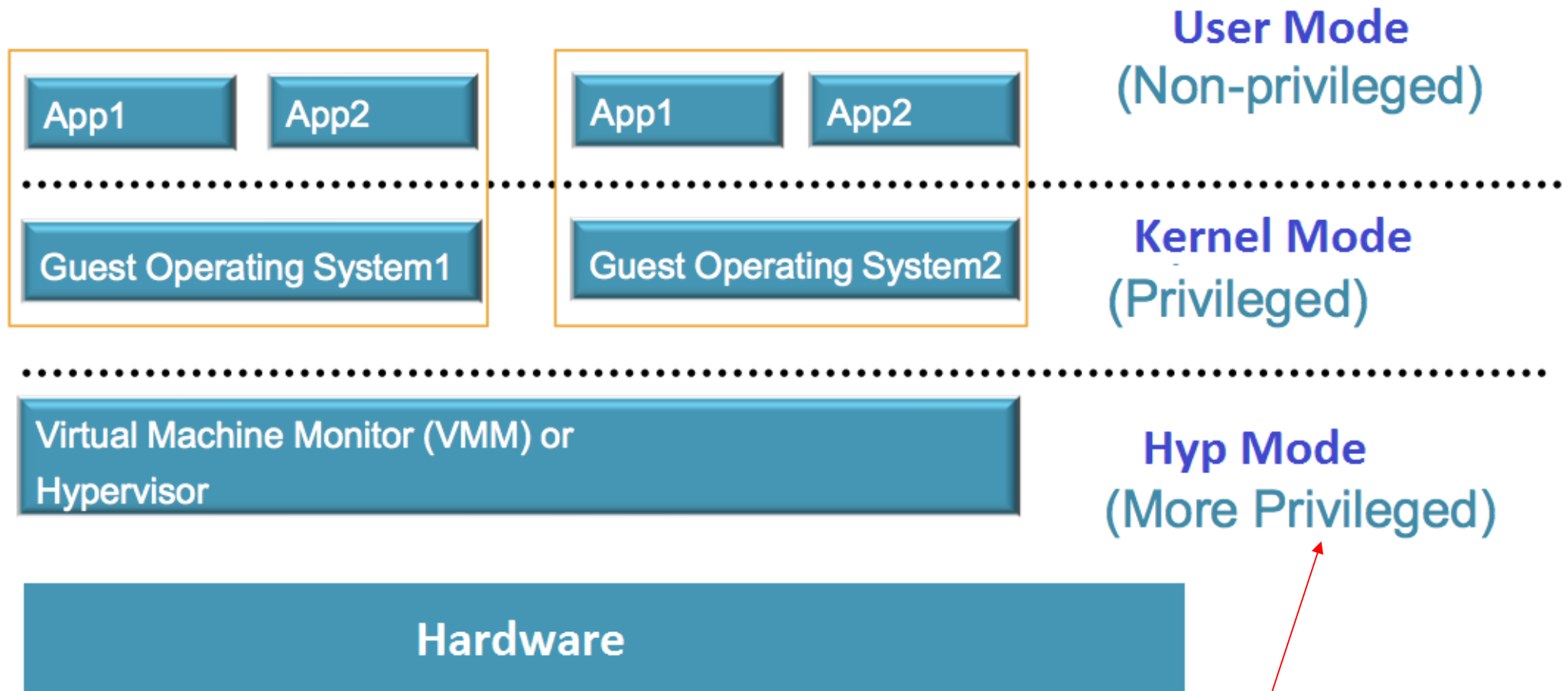
Hardware Configurable to trap sensitive instructions and interrupt directly to VM's Kernel mode



CPU Virtualization



Hypervisor should be lightweight and simple to program



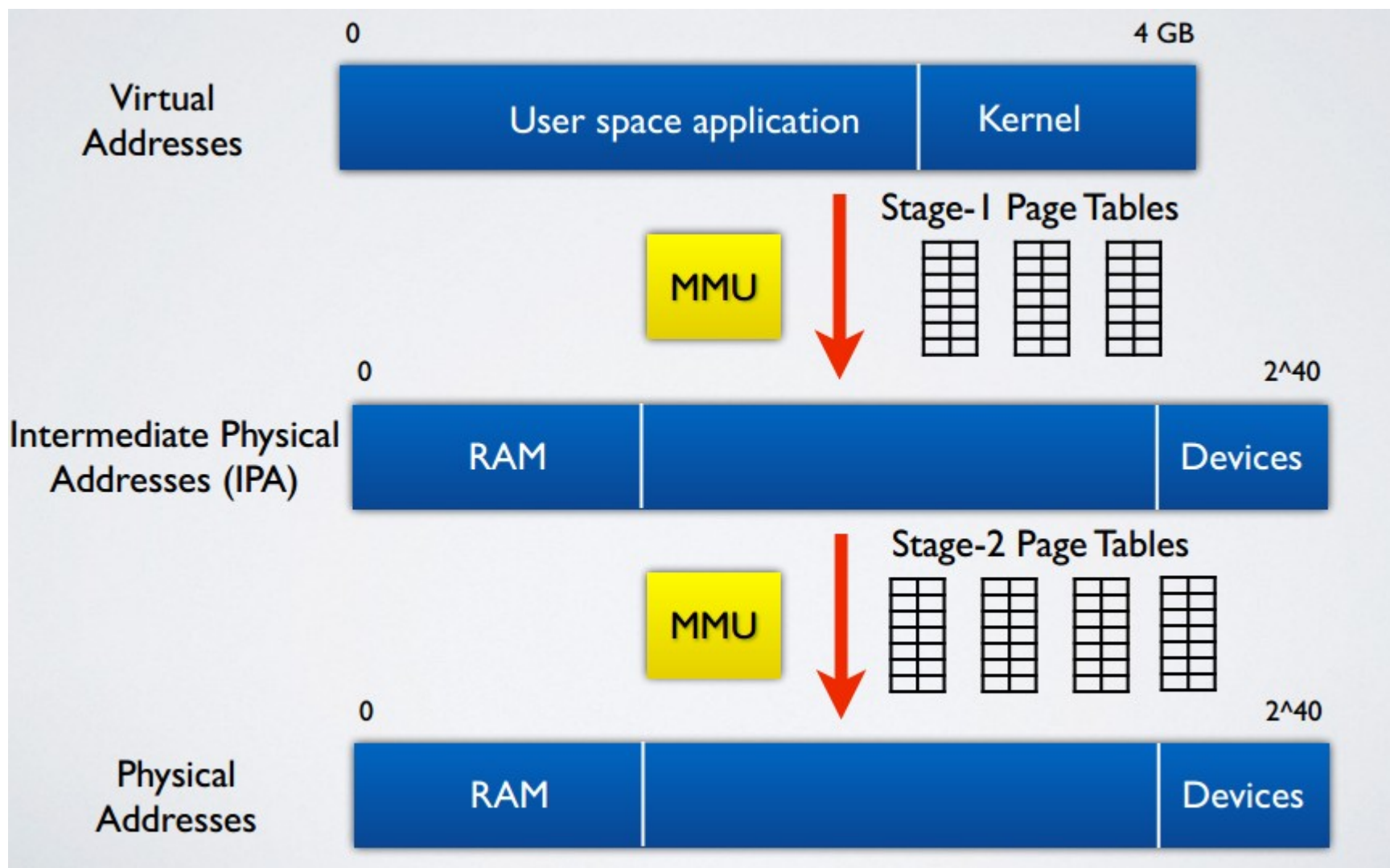
Reduced Number of Control Registers

Hyp mode has its own separate address space

Page tables entries is protected from user mode, as they should not be shared with user mode



Memory Virtualization

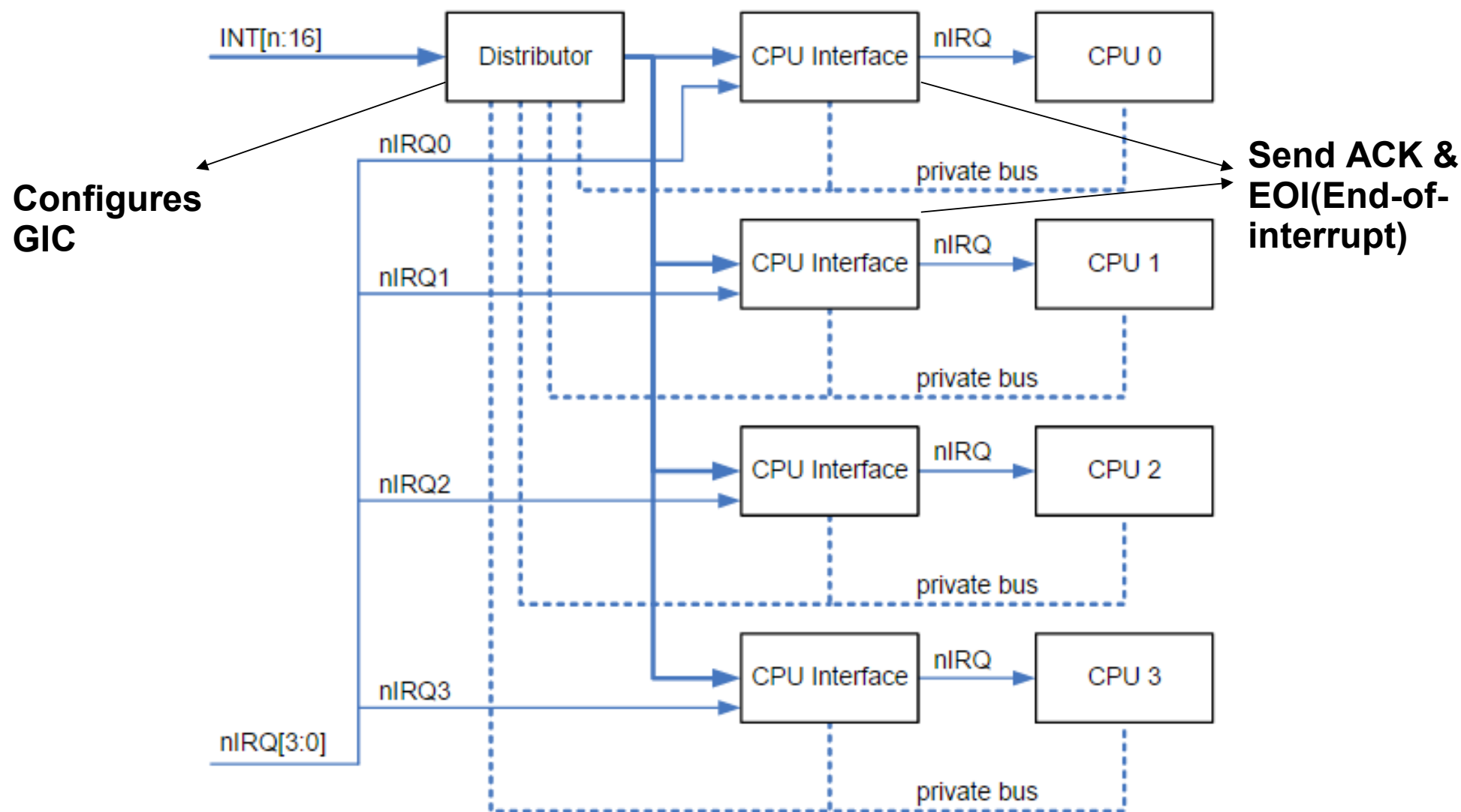


Hardware support to virtualize physical memory: Stage 2 Page Tables
Enabled/Disabled from Hyp Mode



Interrupt Virtualization

Generic Interrupt Controller

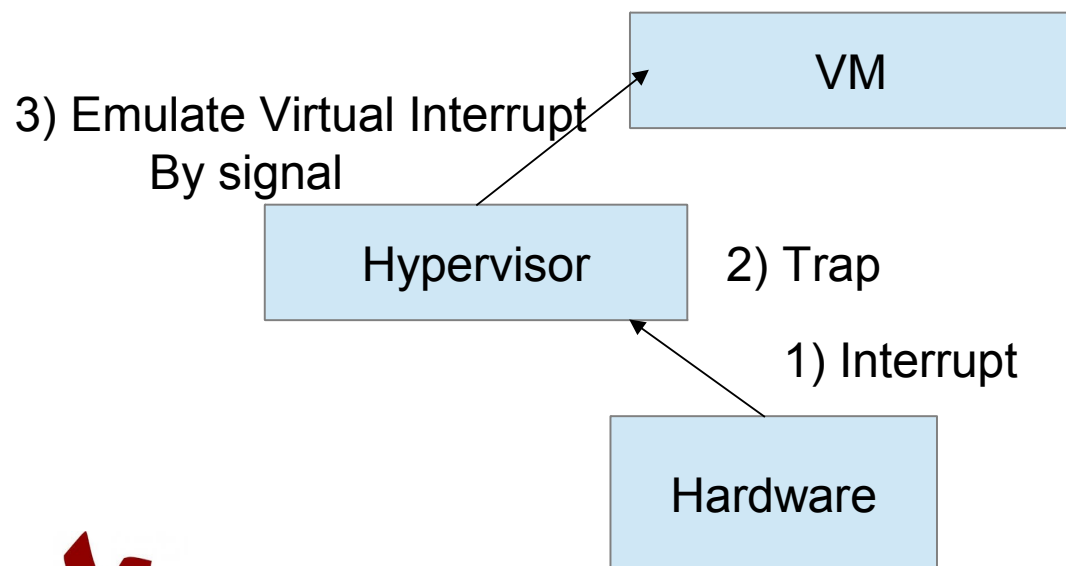




Interrupt Virtualization

Generic Interrupt Controller

Trapping Interrupt in Hyp Mode



**Cumbersome &
Expensive**



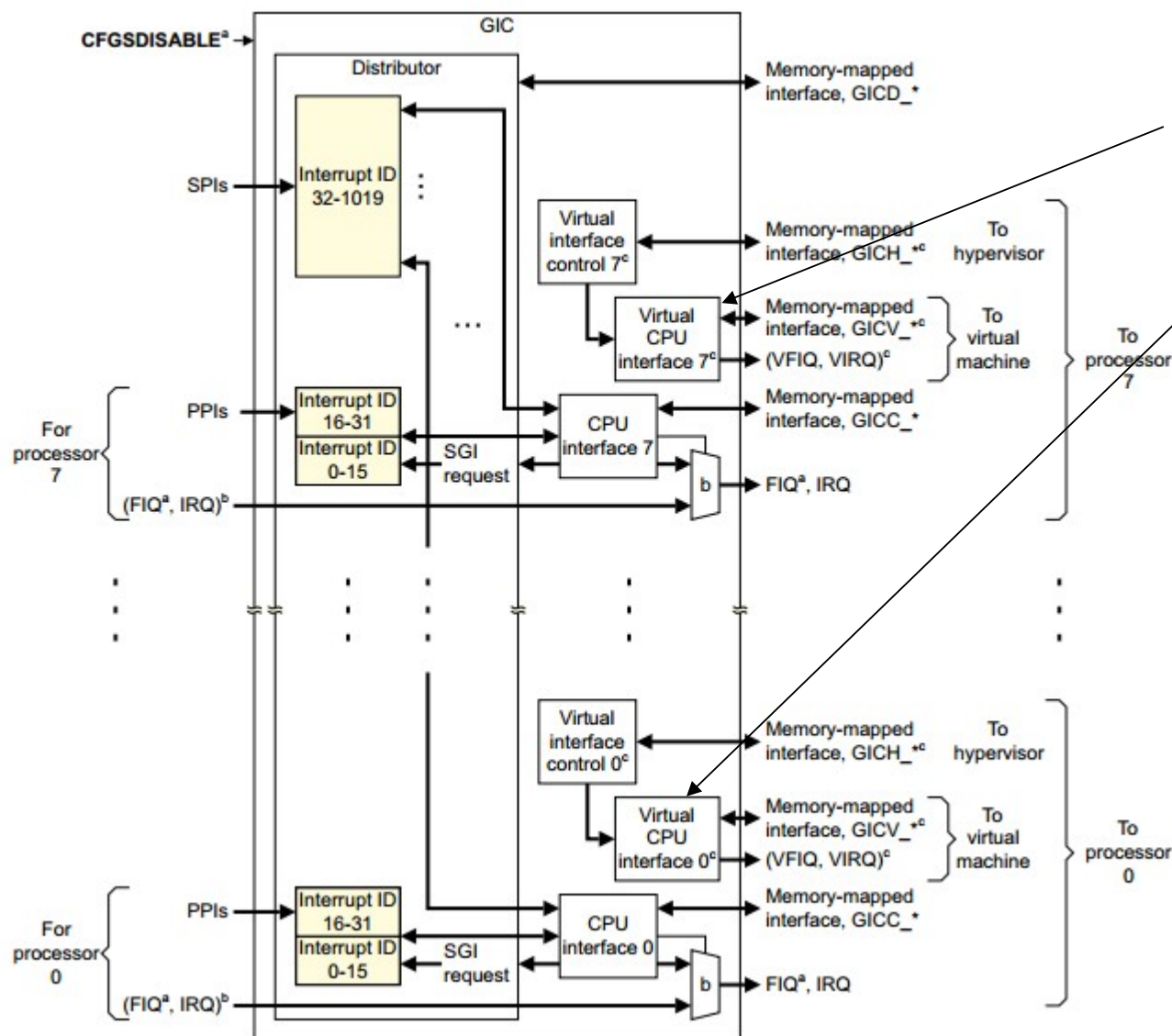
Interrupt Virtualization

Generic Interrupt Controller (V2.0)

Virtual GIC

Why Virtual Distributor is required?

**But No
Virtual
Distributor**



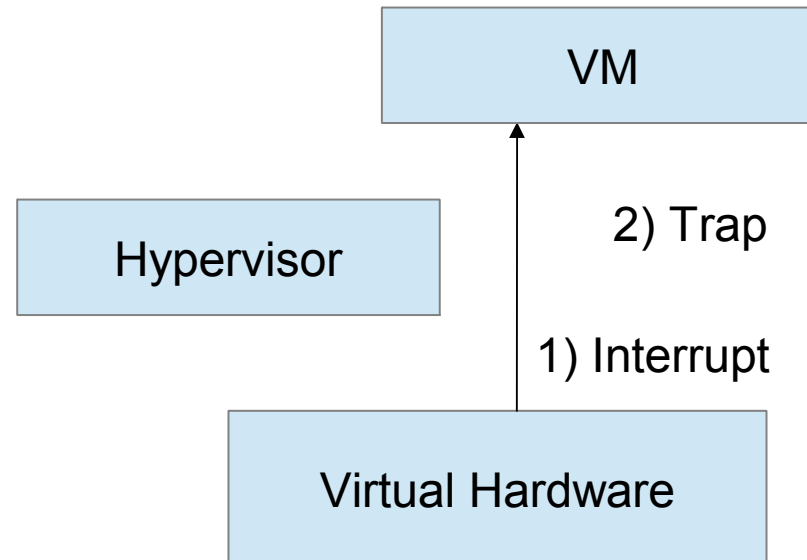


Interrupt Virtualization

Generic Interrupt Controller (V2.0)

Virtual GIC

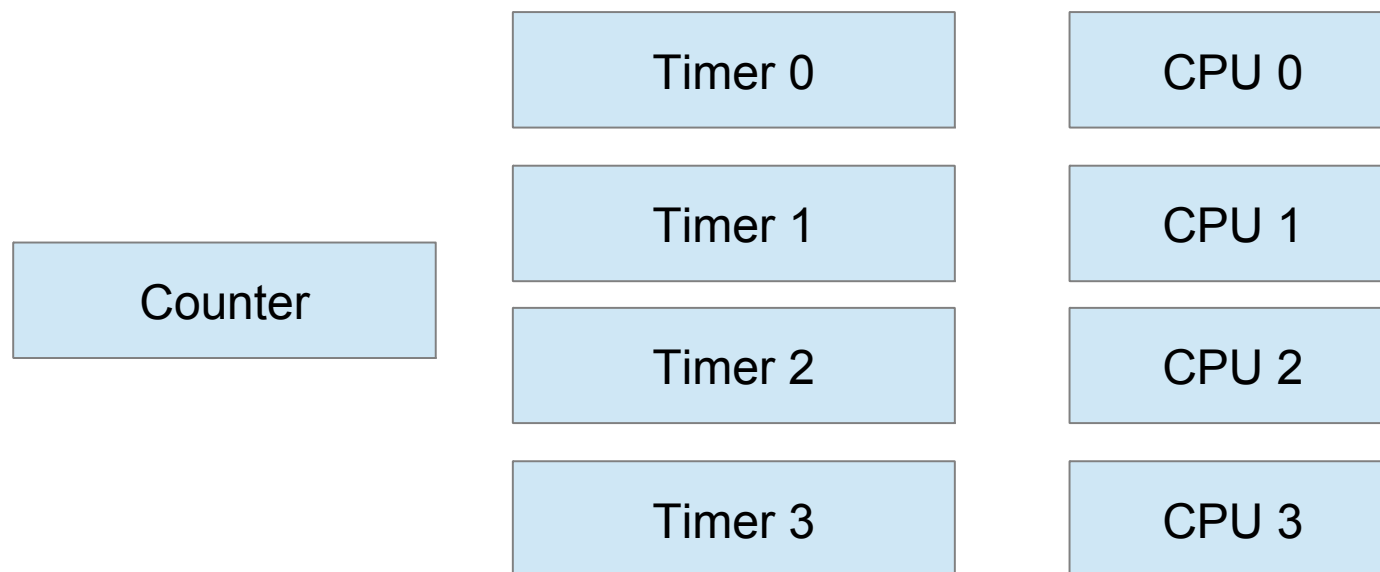
Trapping Interrupt in Kernel Mode





Timer Virtualization

Generic Timer support



Can be configured from Hyp Mode

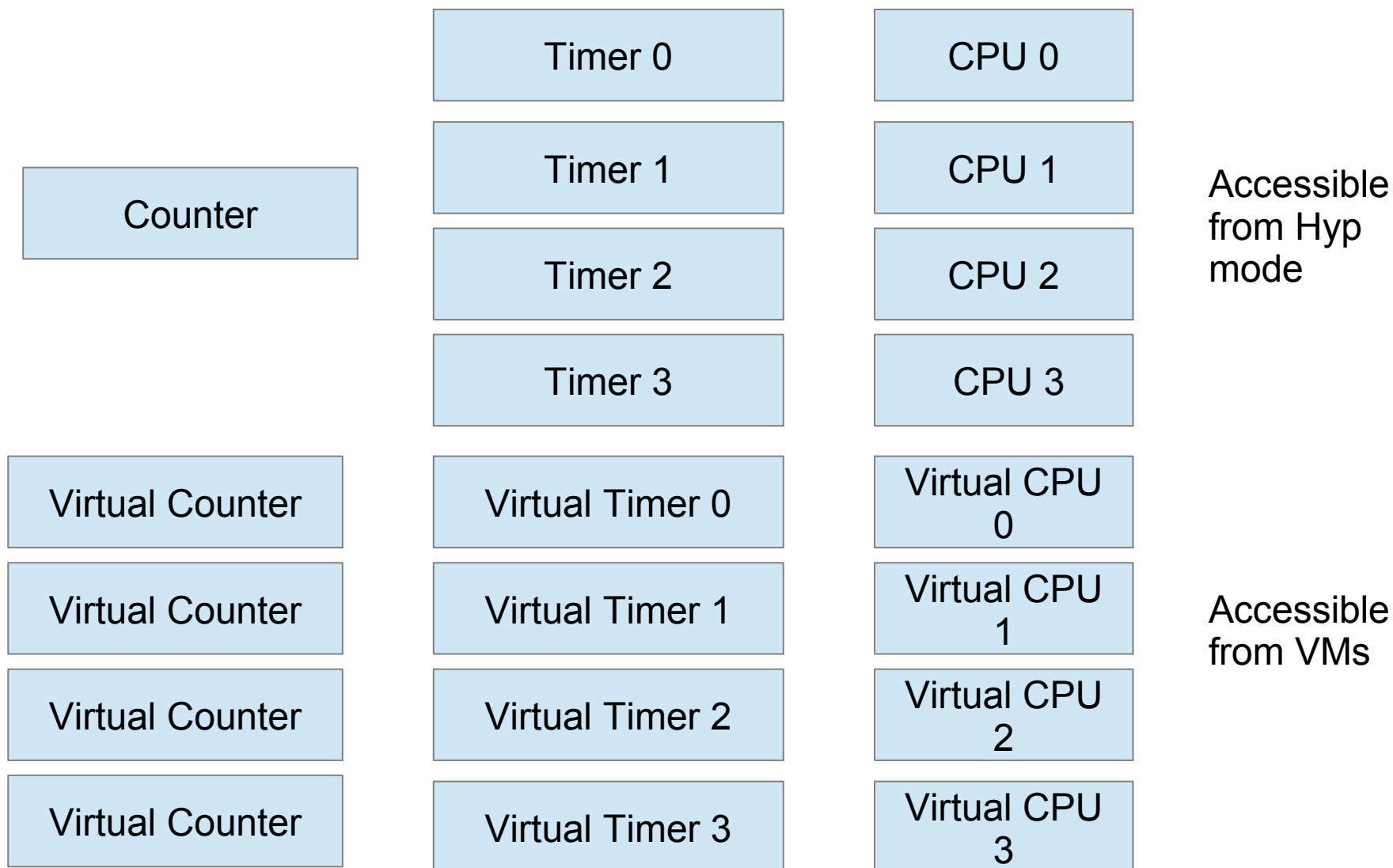
But requires Timer operations to be performed by Hyp

ARM introduces Virtual Timer and Counter



Timer Virtualization

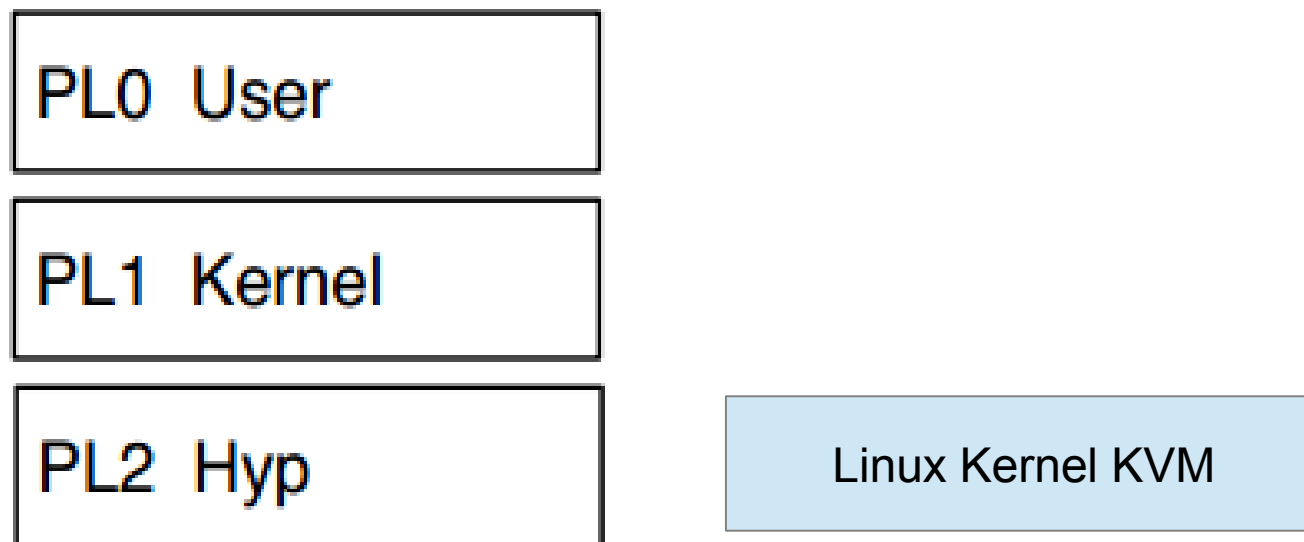
Generic and Virtualized Timer support



Why Virtual Timer support is required?



Hypervisor Architecture



Firstly, Linux is written to work in kernel mode and **would not run unmodified** in Hyp Mode

Some registers and table formats are different in Hyp mode than in kernel mode

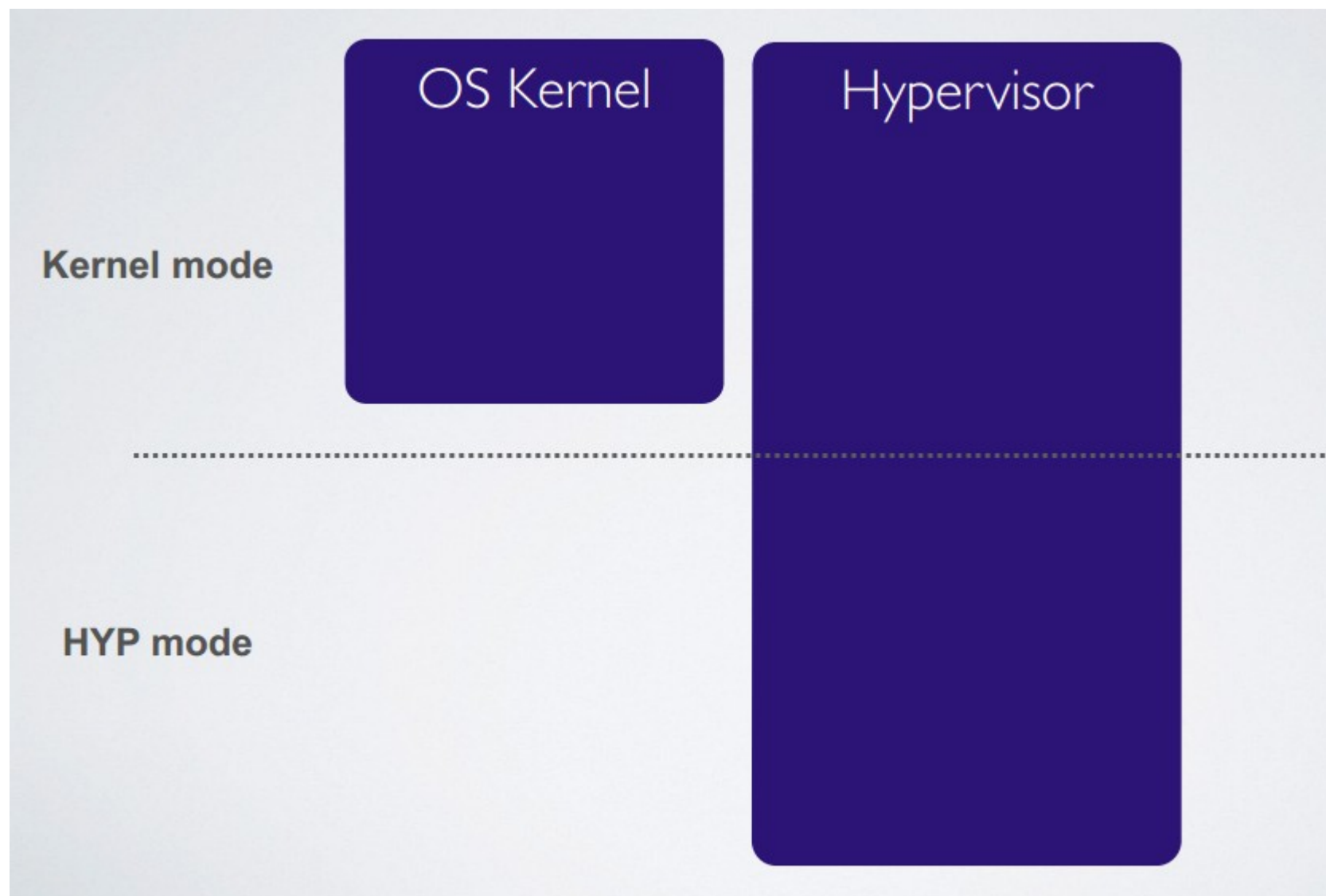
Secondly, Running entire kernel in Hyp Mode would adversely **affect native performance**

Hyp mode has its own separate address space. This requires explicitly map user space data while accessing user space memory



Hypervisor Architecture

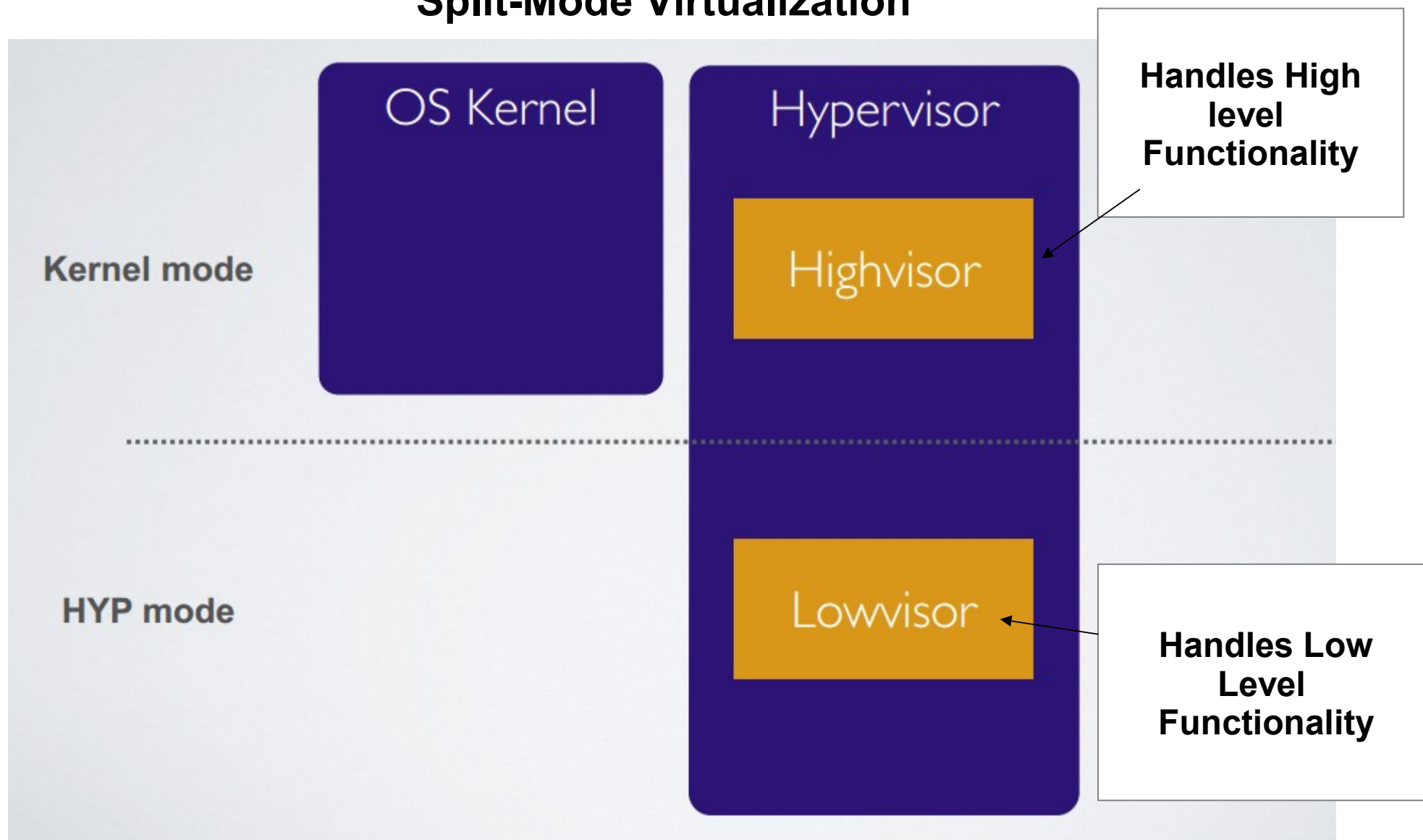
Split-Mode Virtualization





Hypervisor Architecture

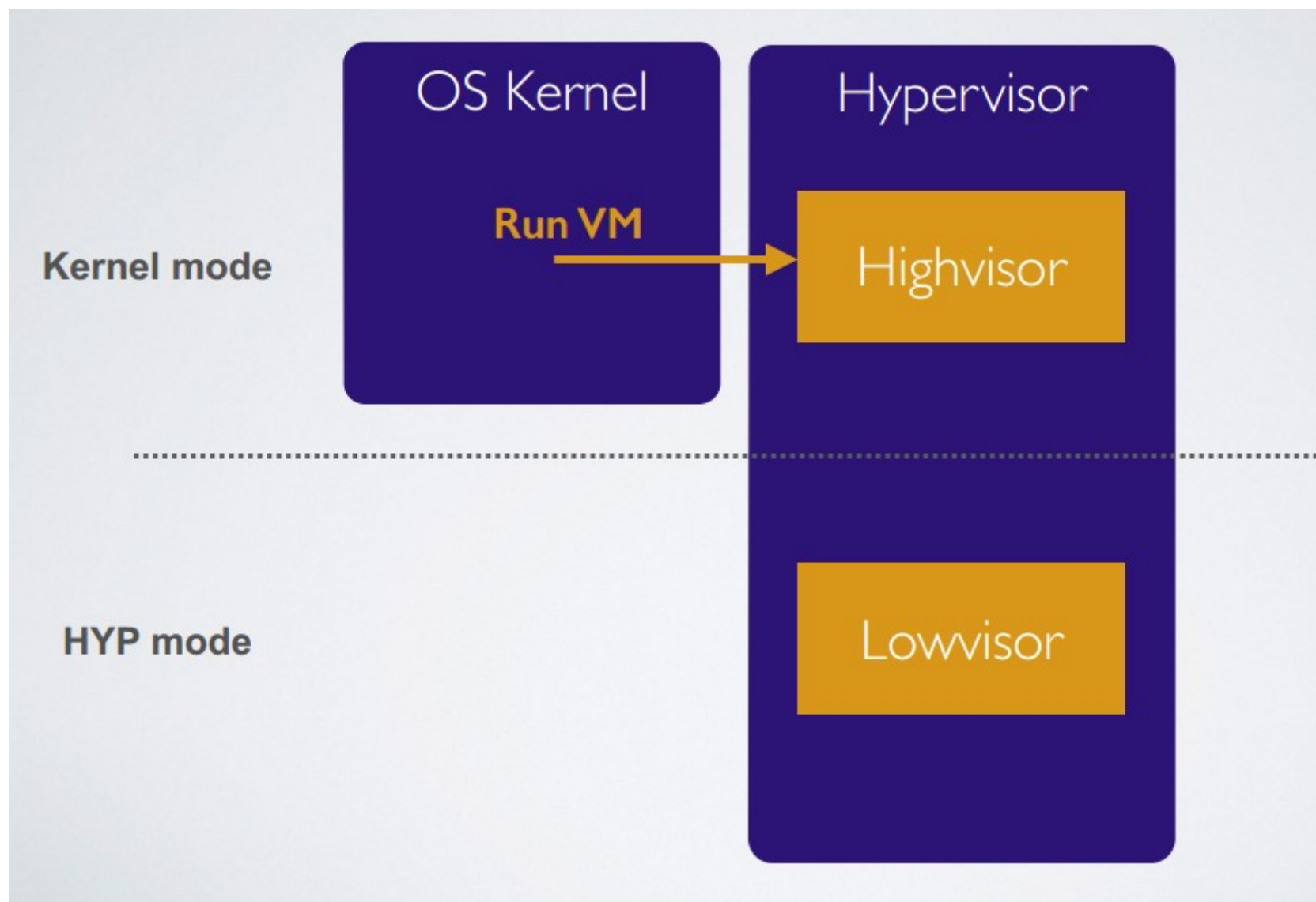
Split-Mode Virtualization





Hypervisor Architecture

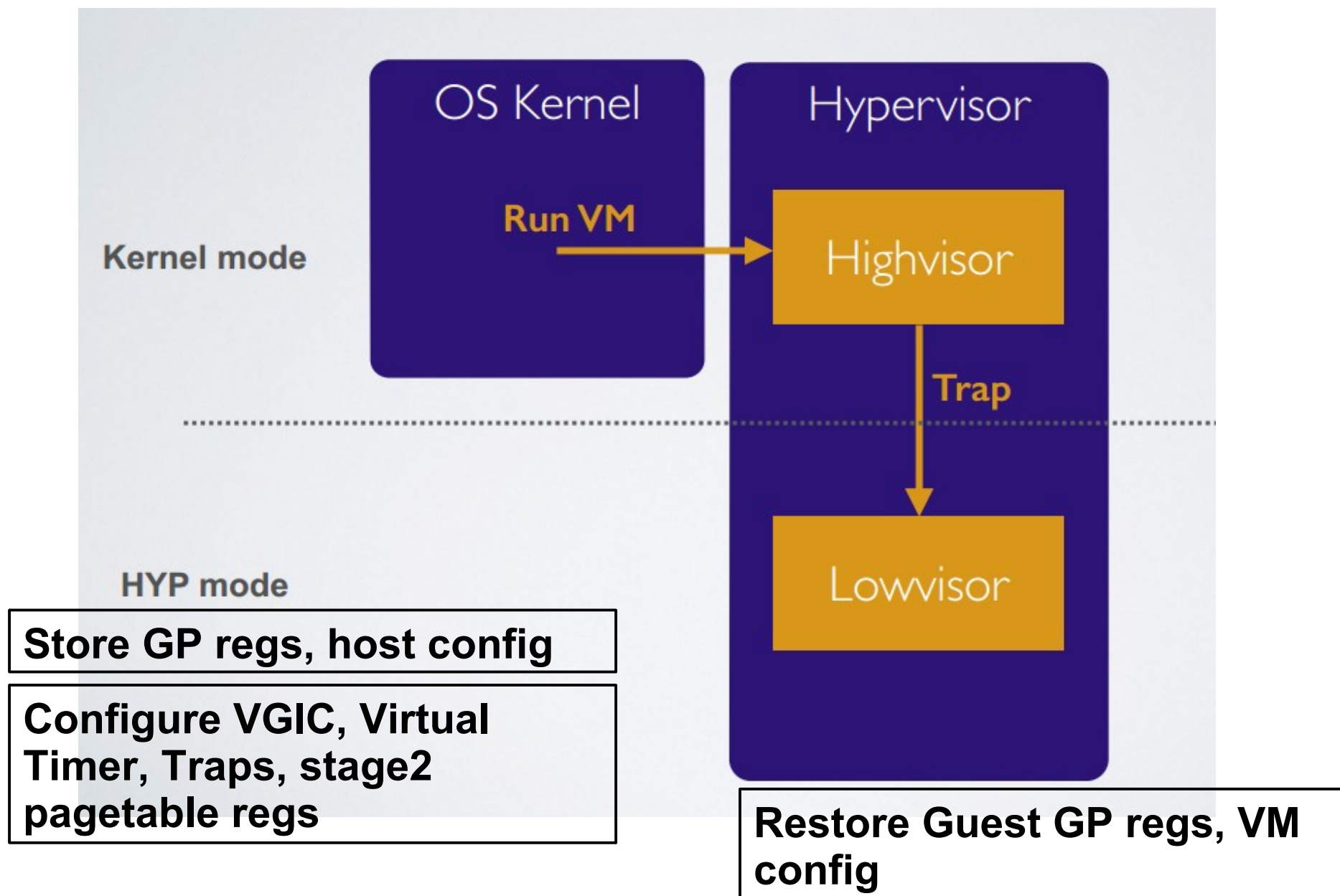
Split-Mode Virtualization





Hypervisor Architecture

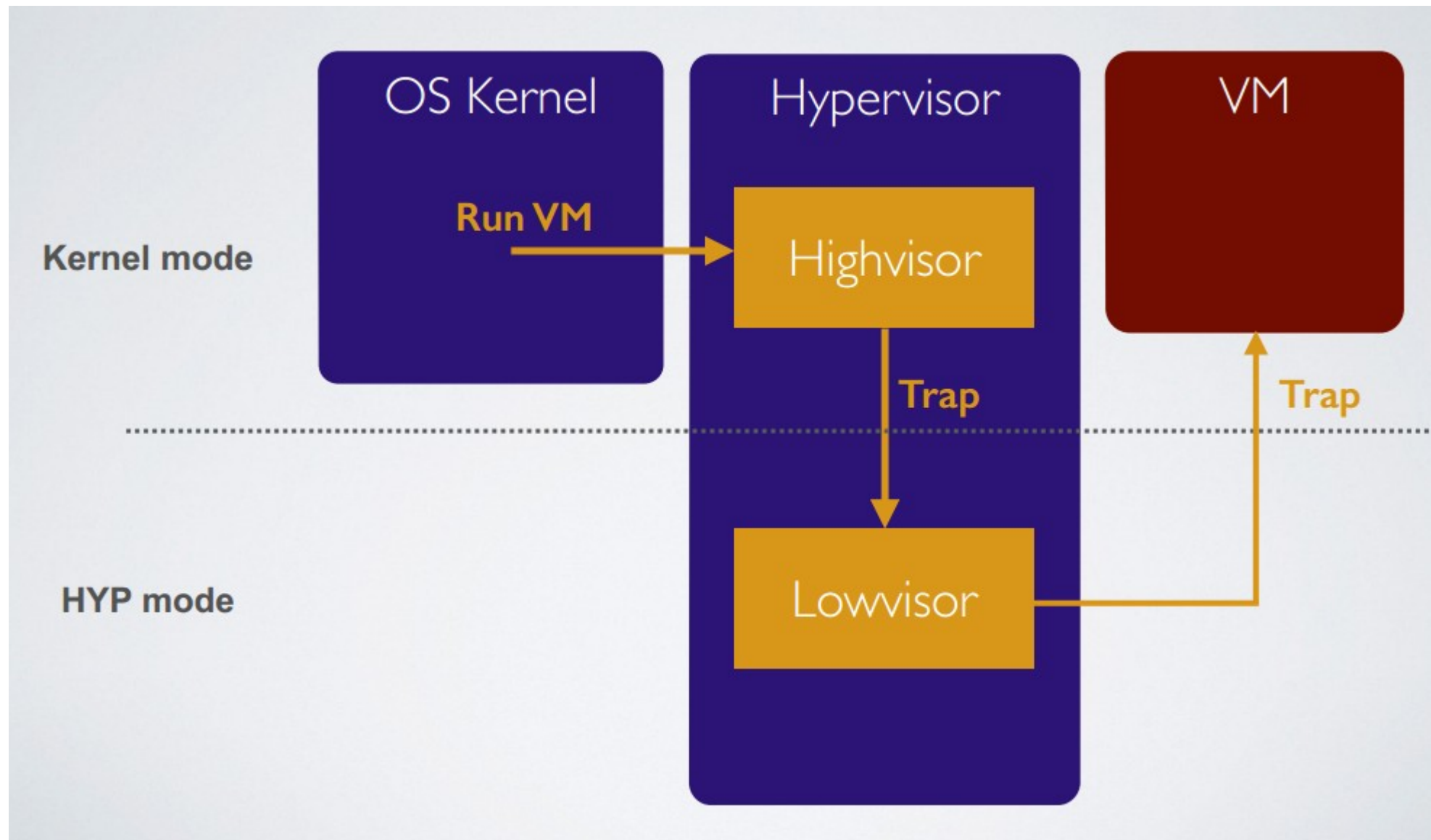
Split-Mode Virtualization





Hypervisor Architecture

Split-Mode Virtualization

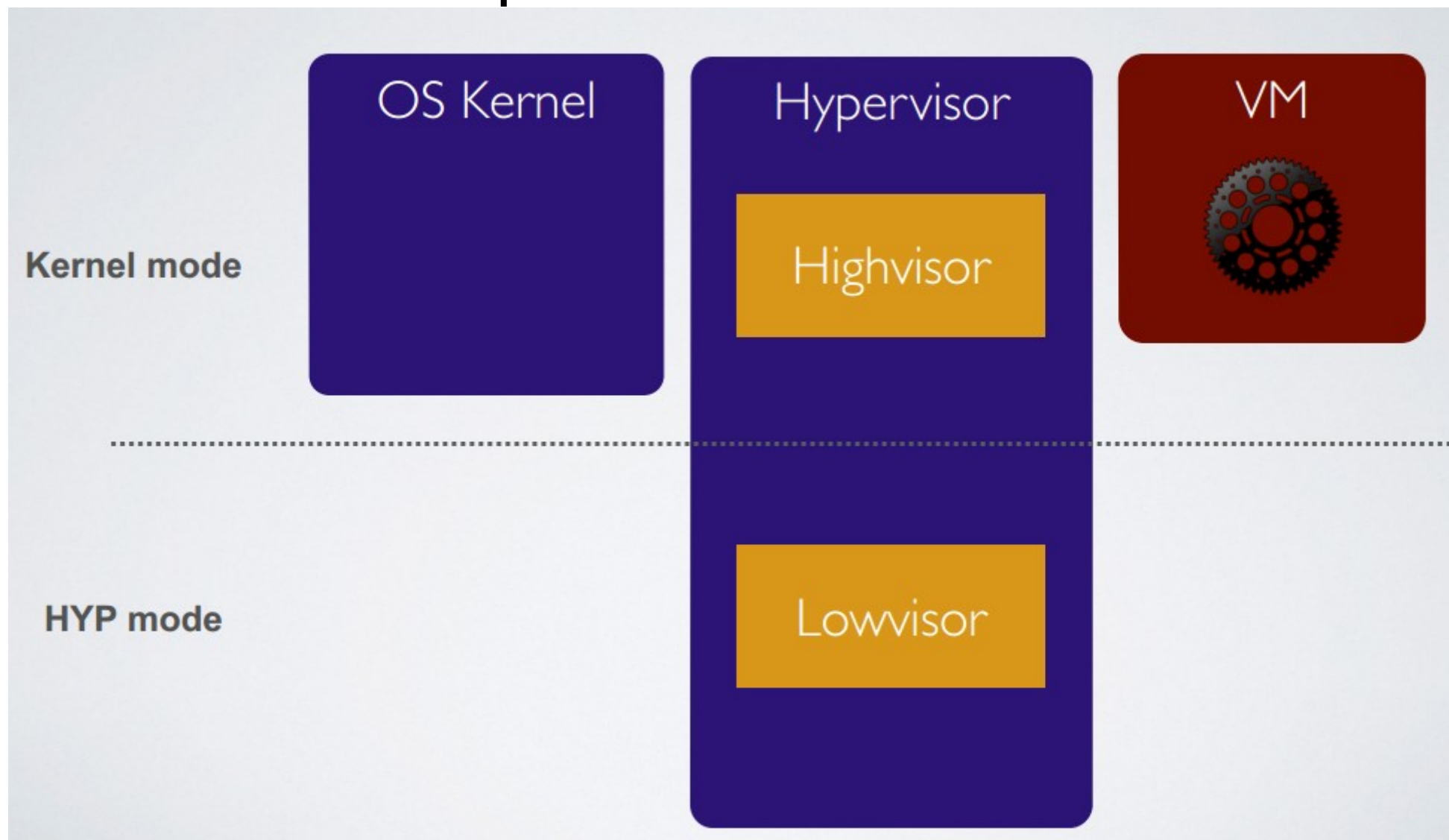


Lowvisor creates correct execution context by configuring hardware



Hypervisor Architecture

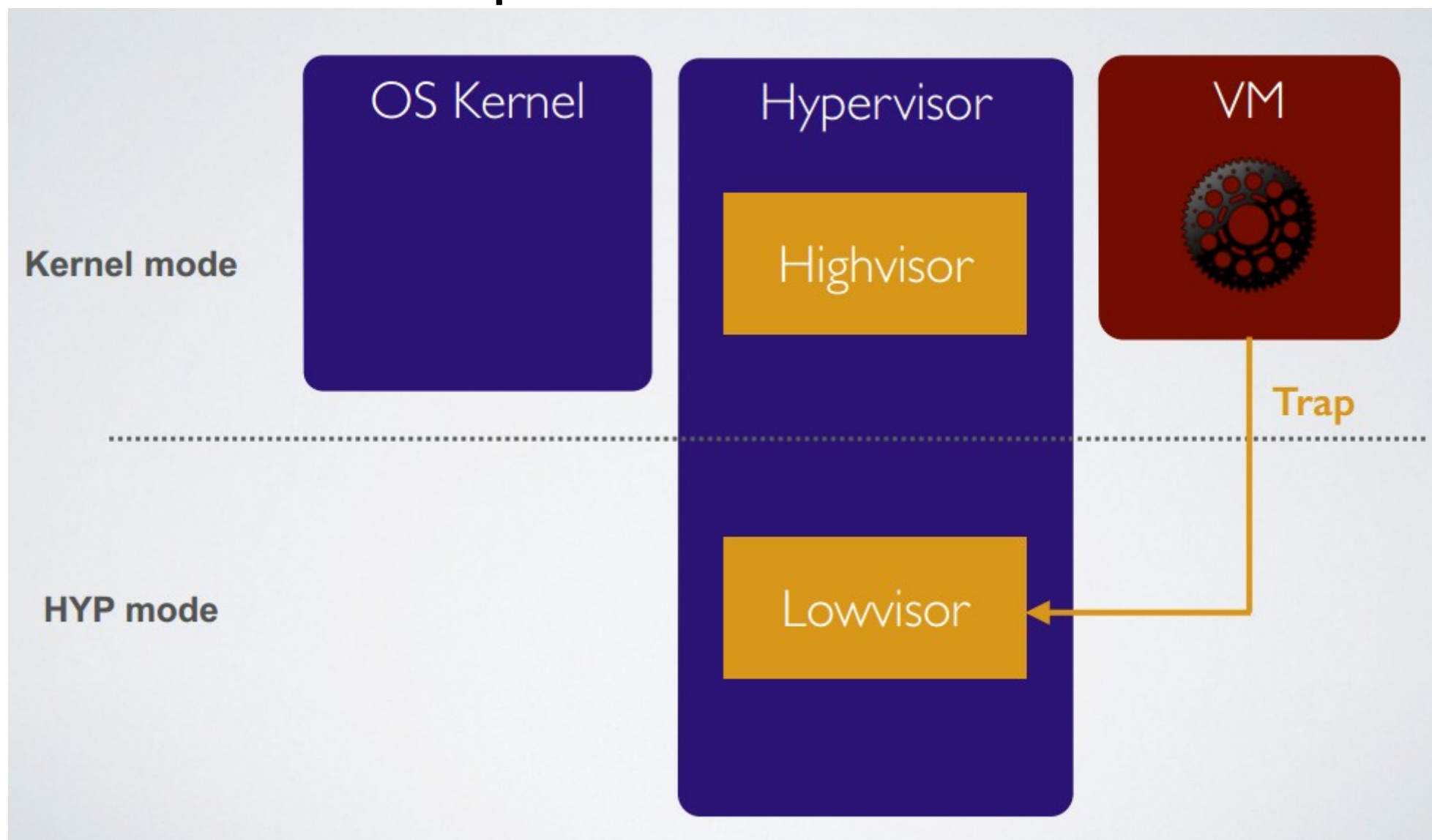
Split-Mode Virtualization





Hypervisor Architecture

Split-Mode Virtualization

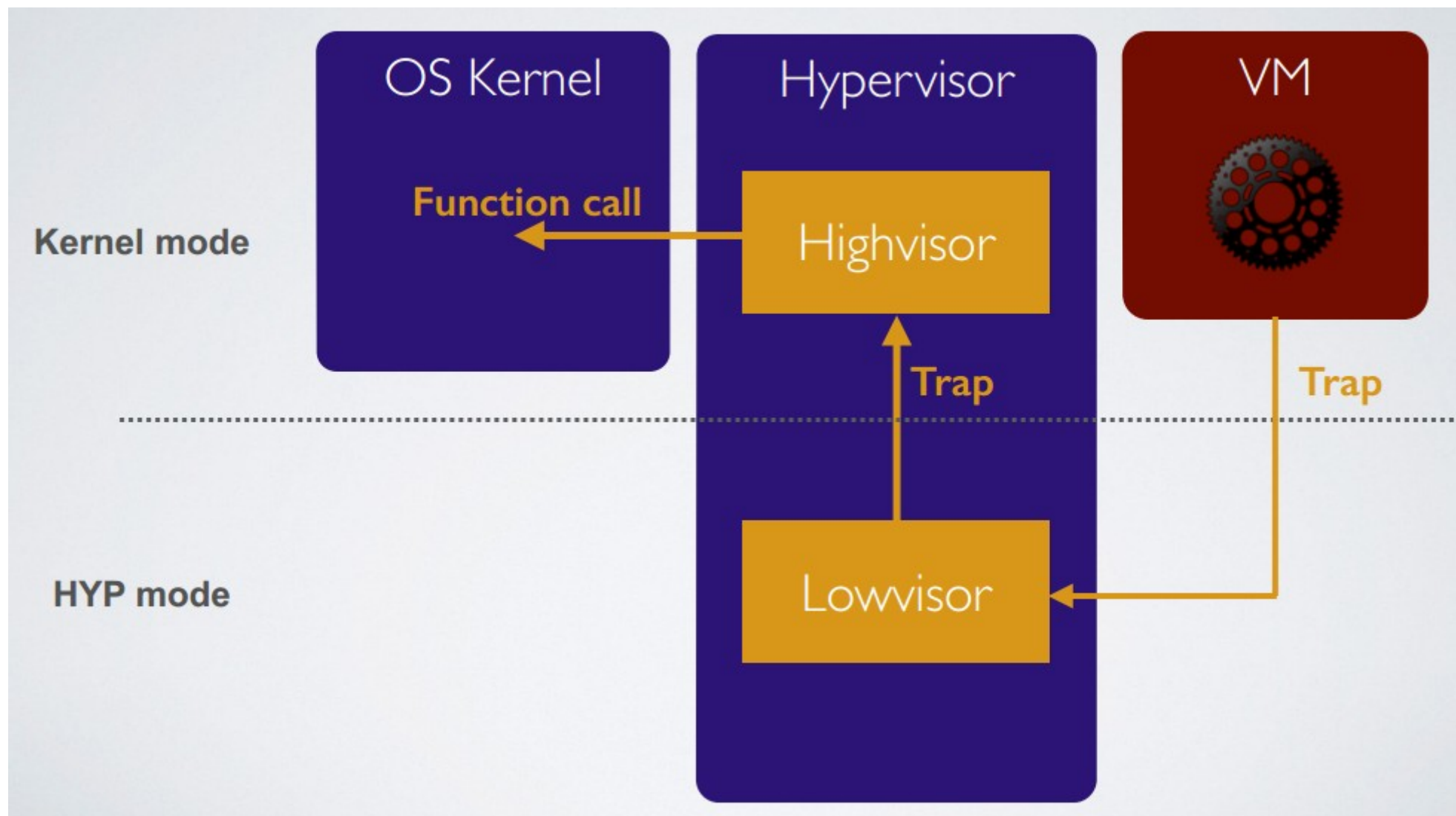


Lowvisor: Handles Interrupts and Exceptions



Hypervisor Architecture

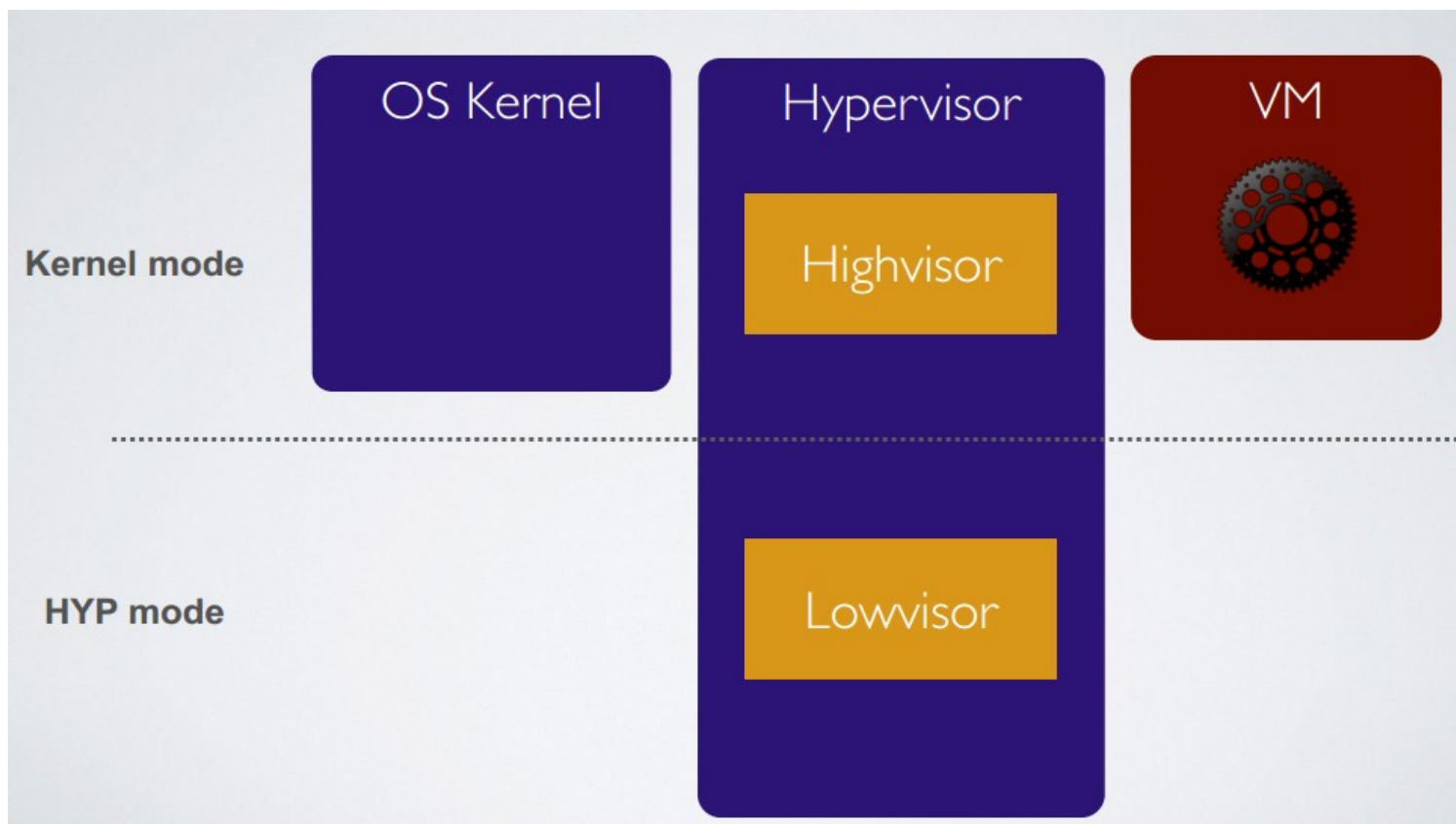
Split-Mode Virtualization





Hypervisor Architecture

What are the Functionalities of LowVisor?



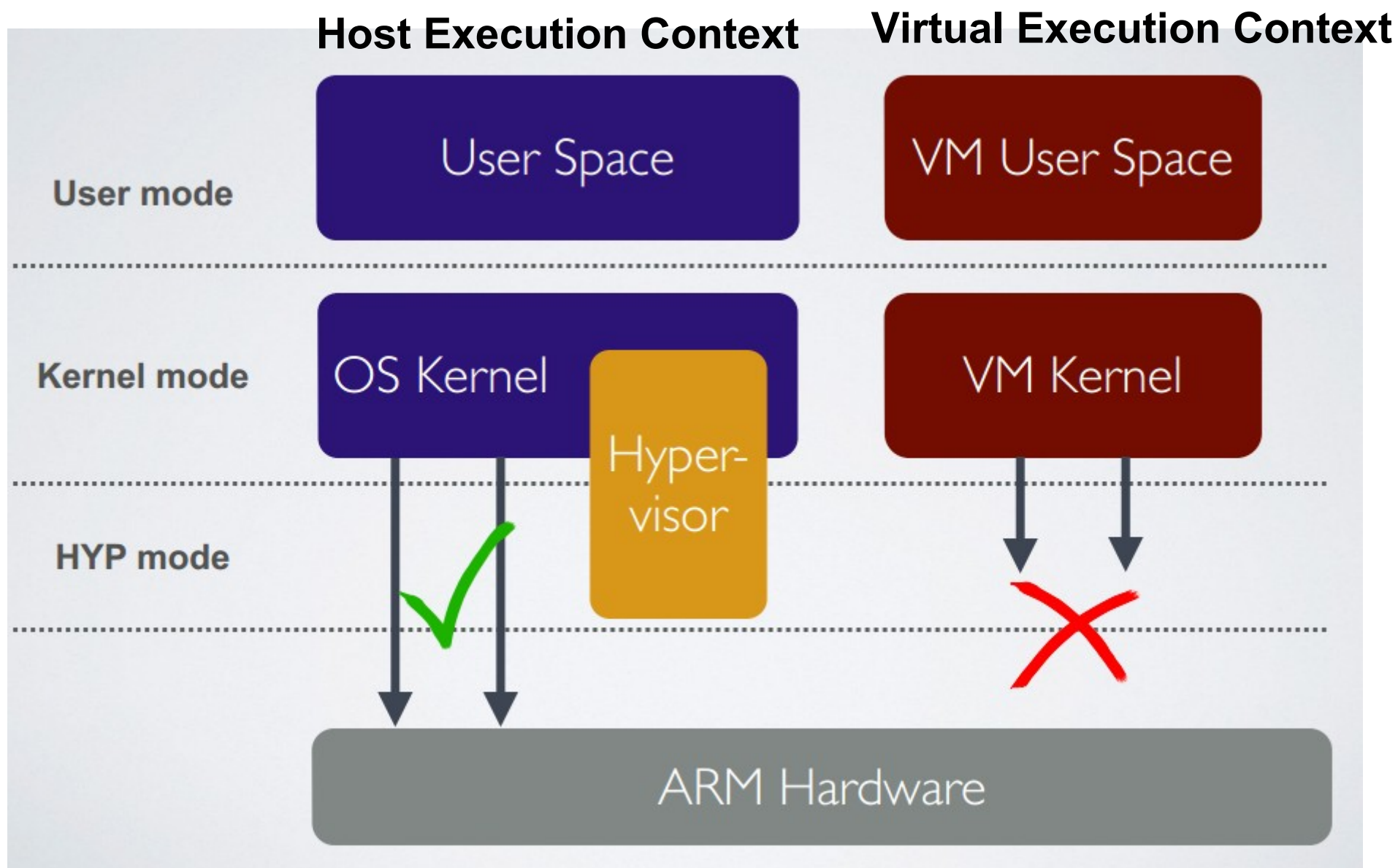
Trap handler: Handles Interrupt and Exception

Set Correct Execution context by configuring hardware

Perform Context Switch between Host and VM execution context



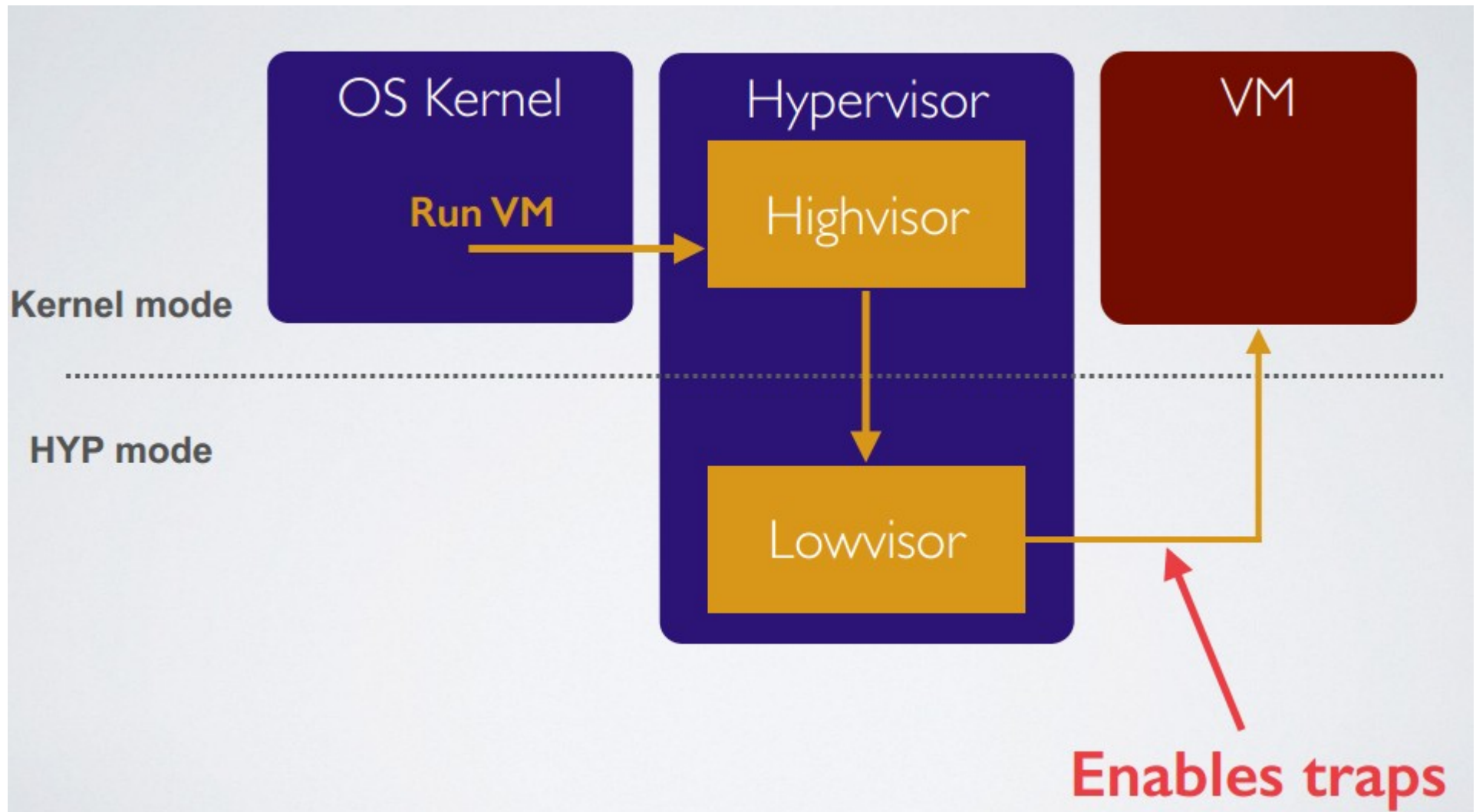
Hypervisor Architecture



Perform Context Switch between Host and VM execution context



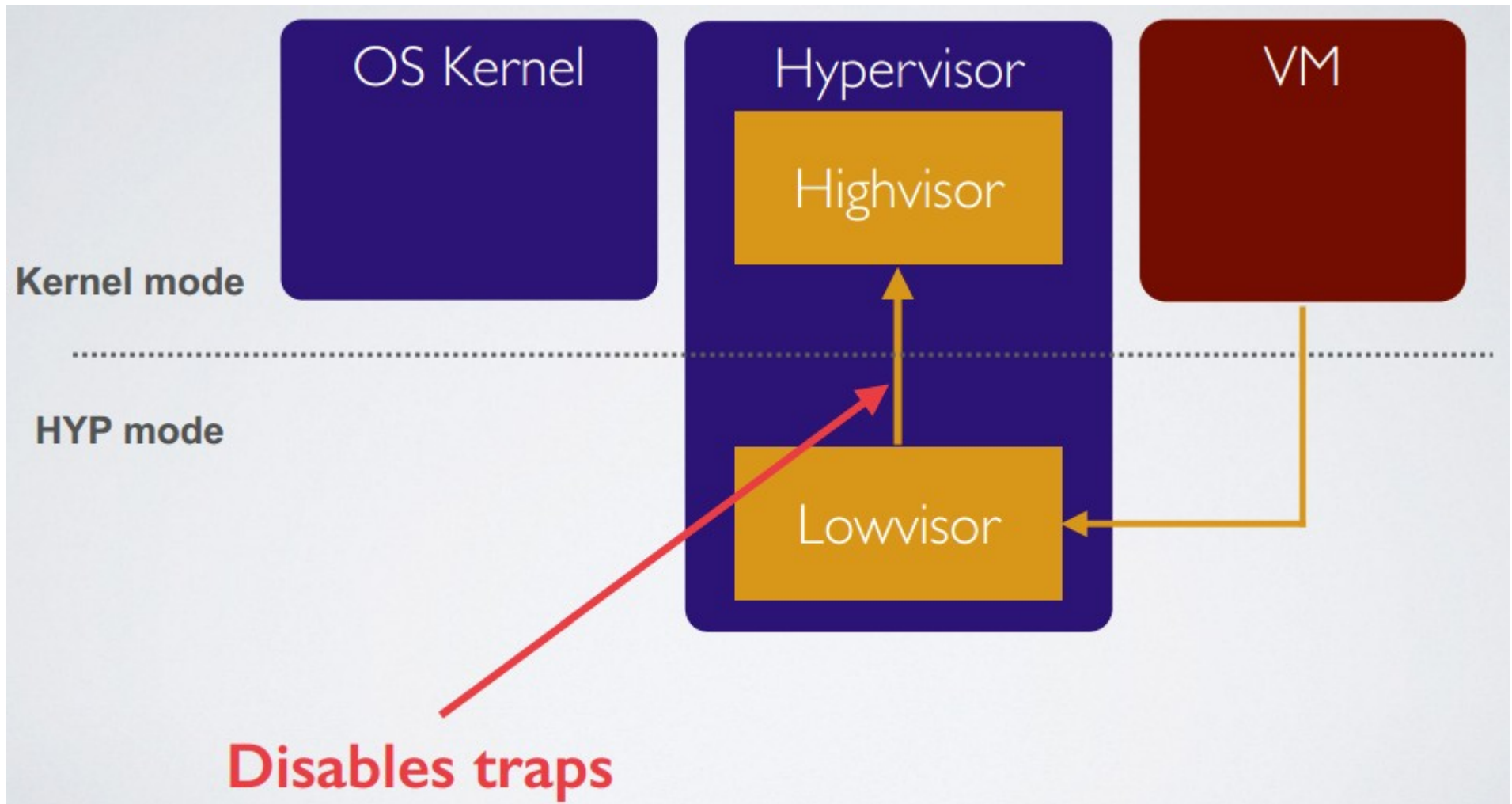
CPU Virtualization



Perform Context Switch between Host and VM execution context



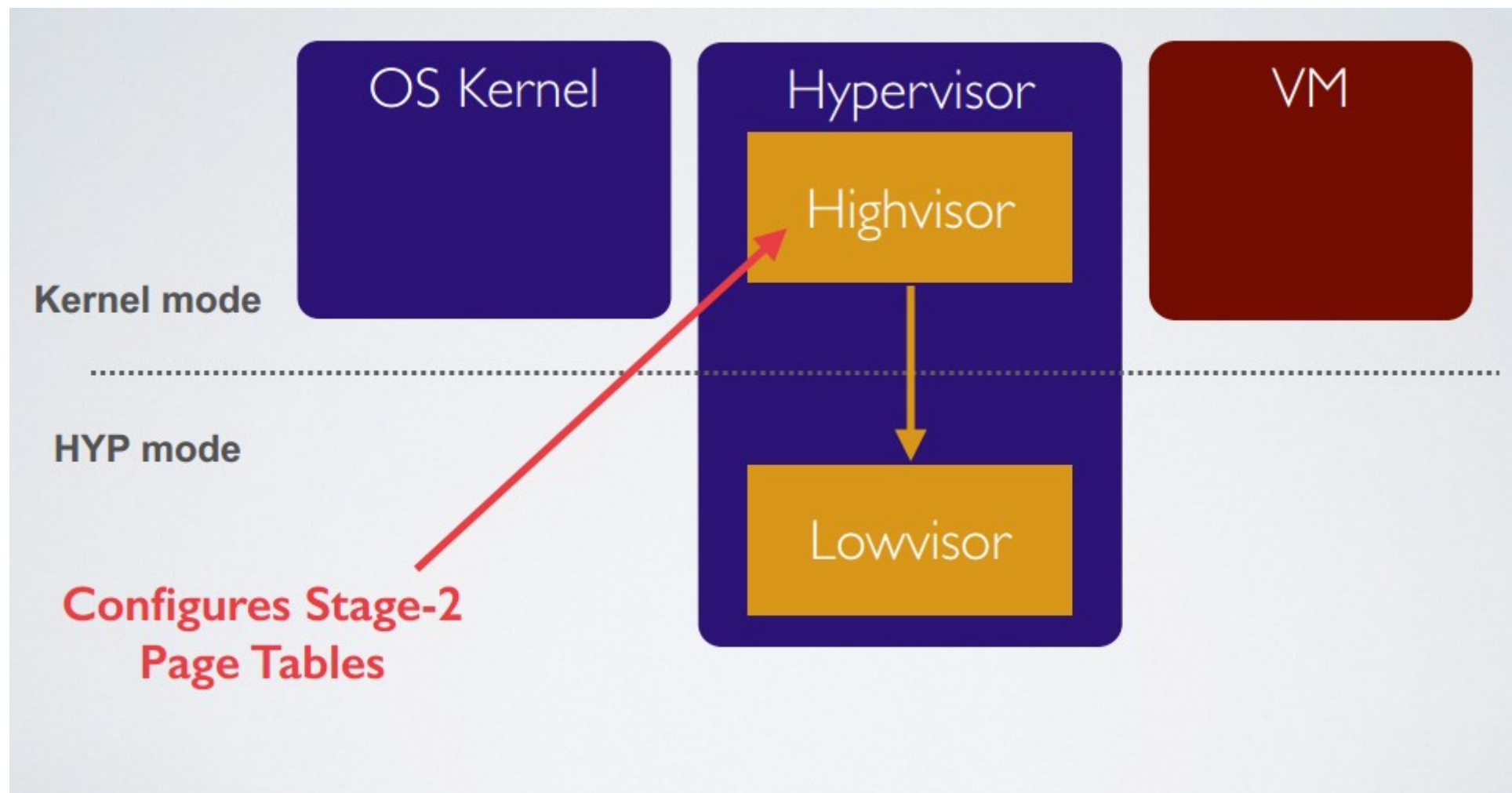
CPU Virtualization



Perform Context Switch between Host and VM execution context



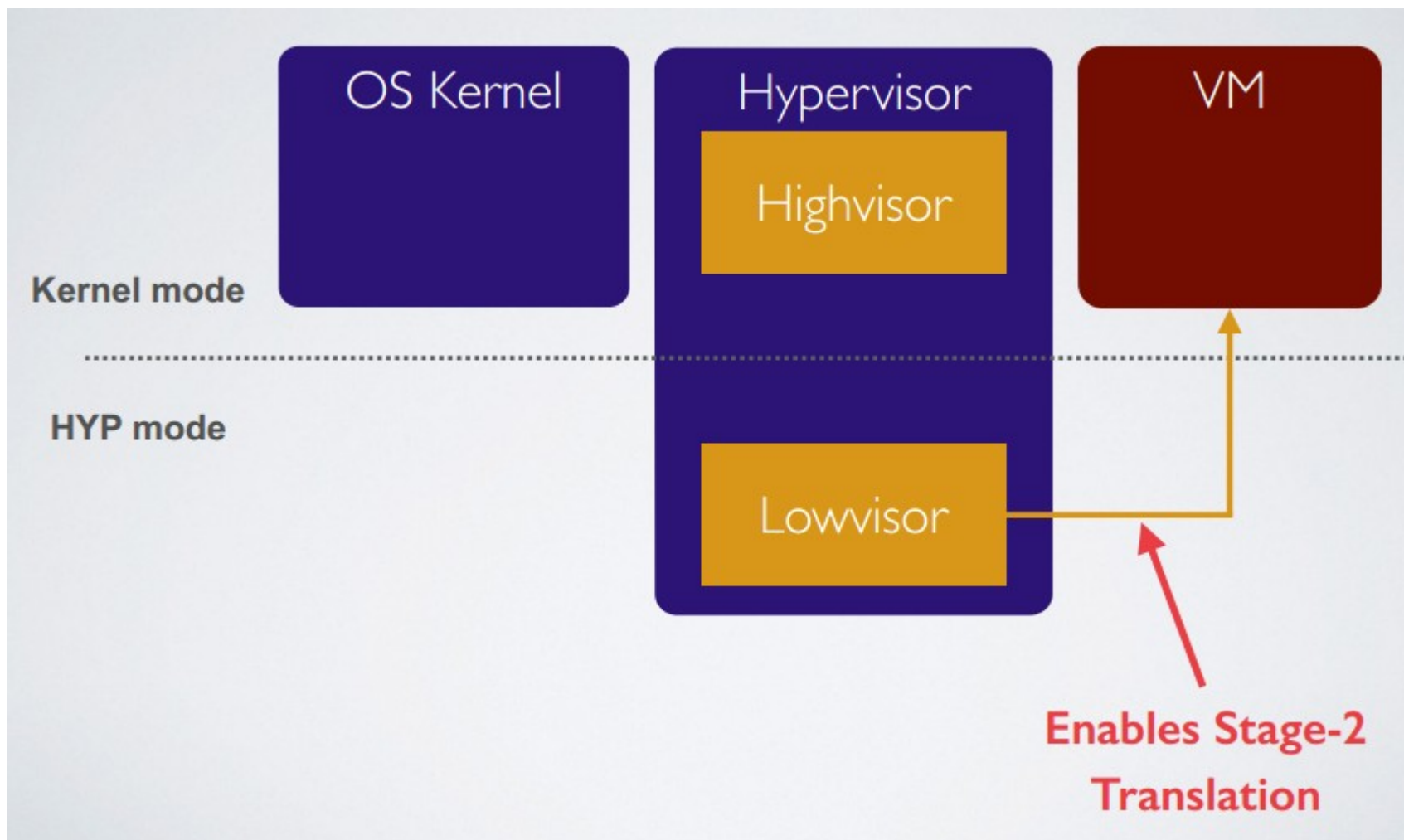
Memory Virtualization



Configuring page tables is a high level Functionality



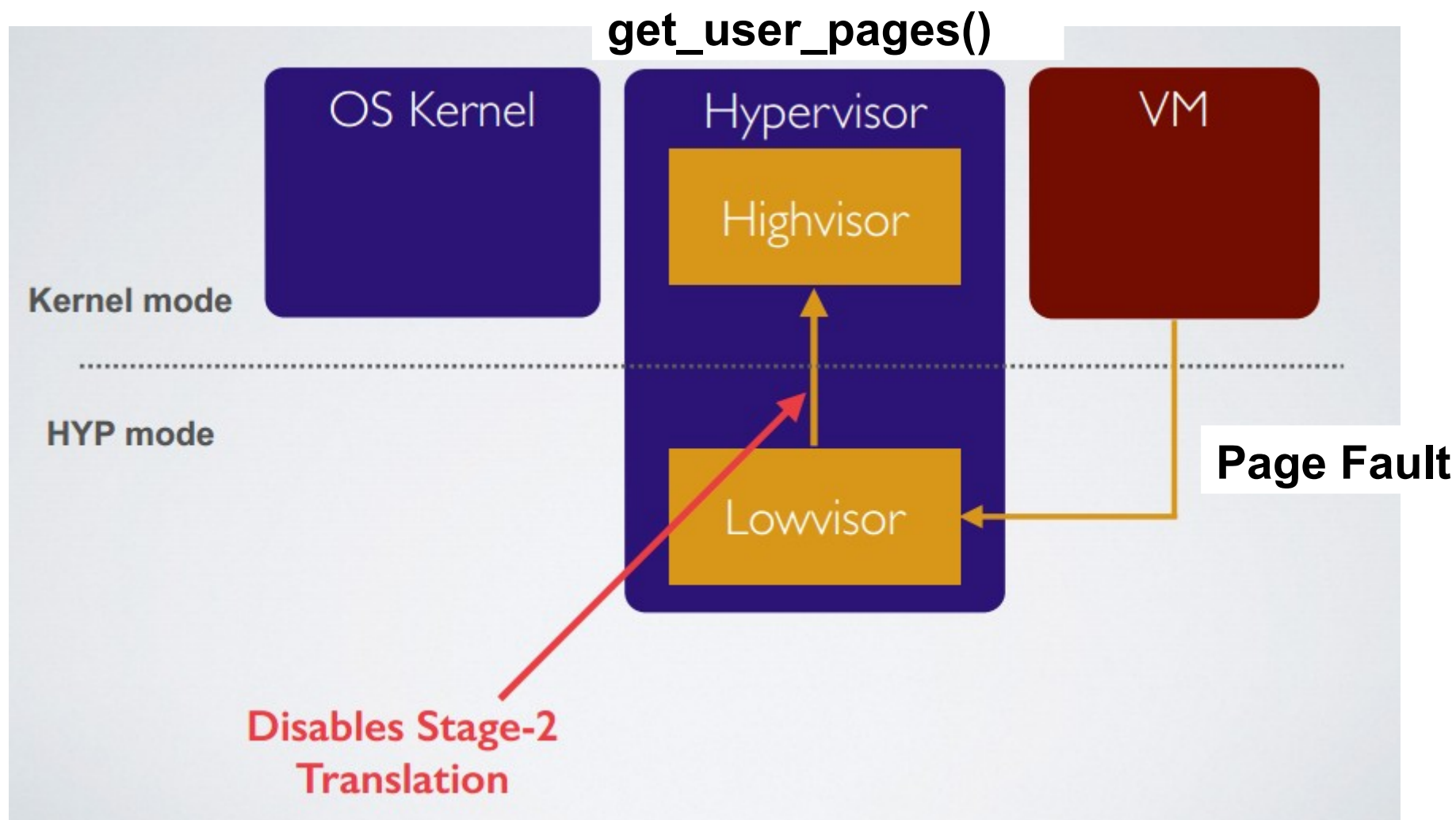
Memory Virtualization



But LowVisor has hardware access as it runs in Hyp Mode

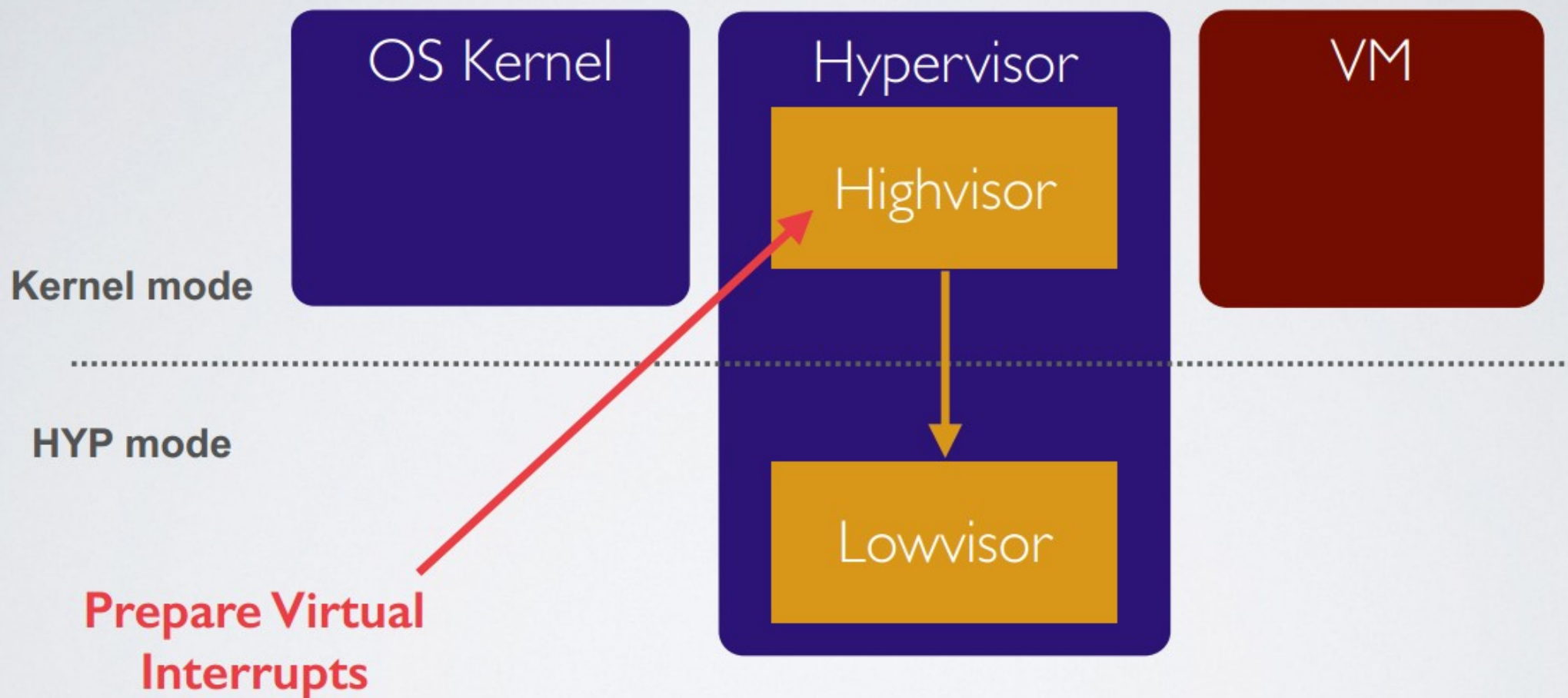


Memory Virtualization



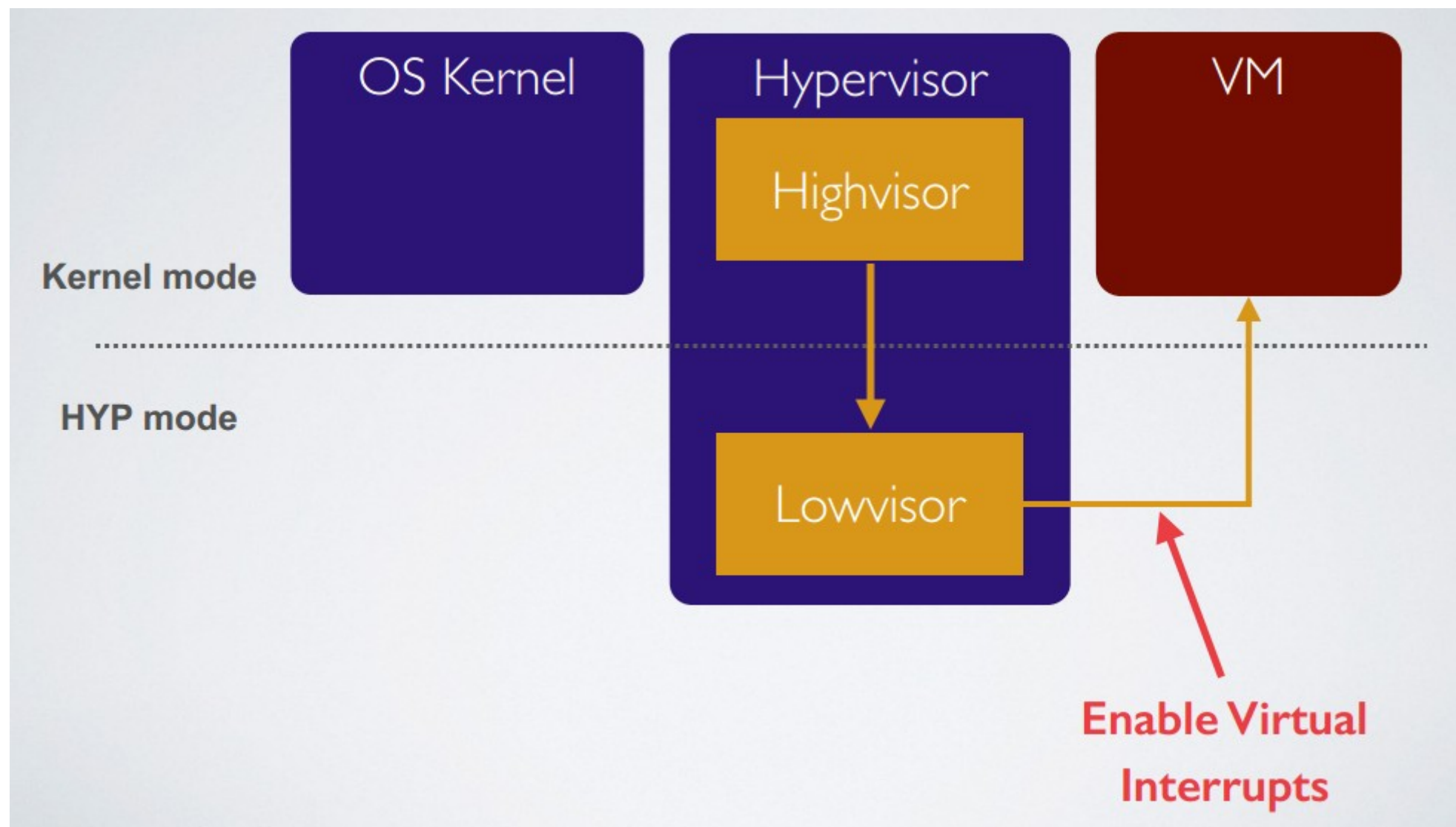


Interrupt Virtualization



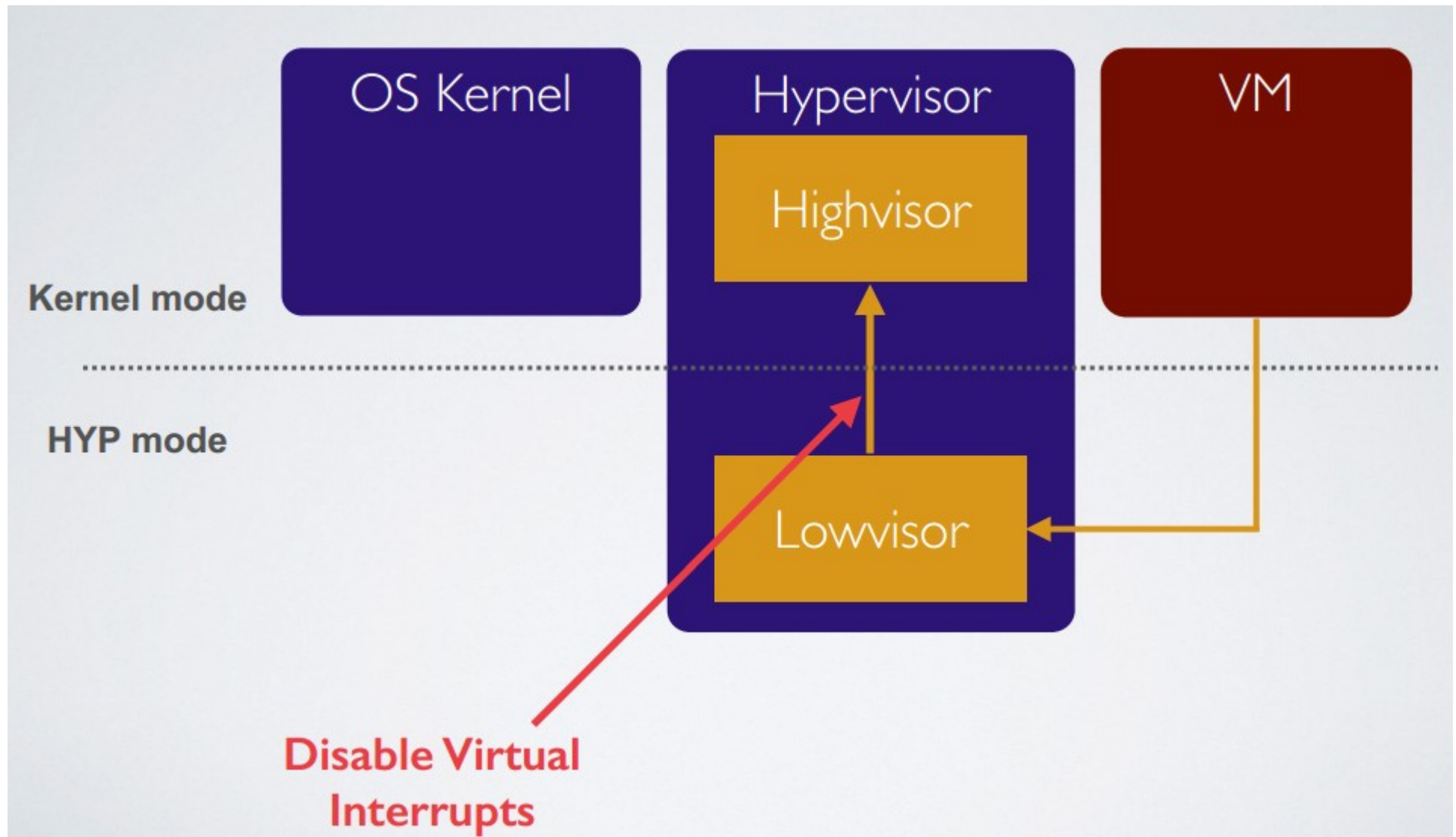


Interrupt Virtualization





Interrupt Virtualization

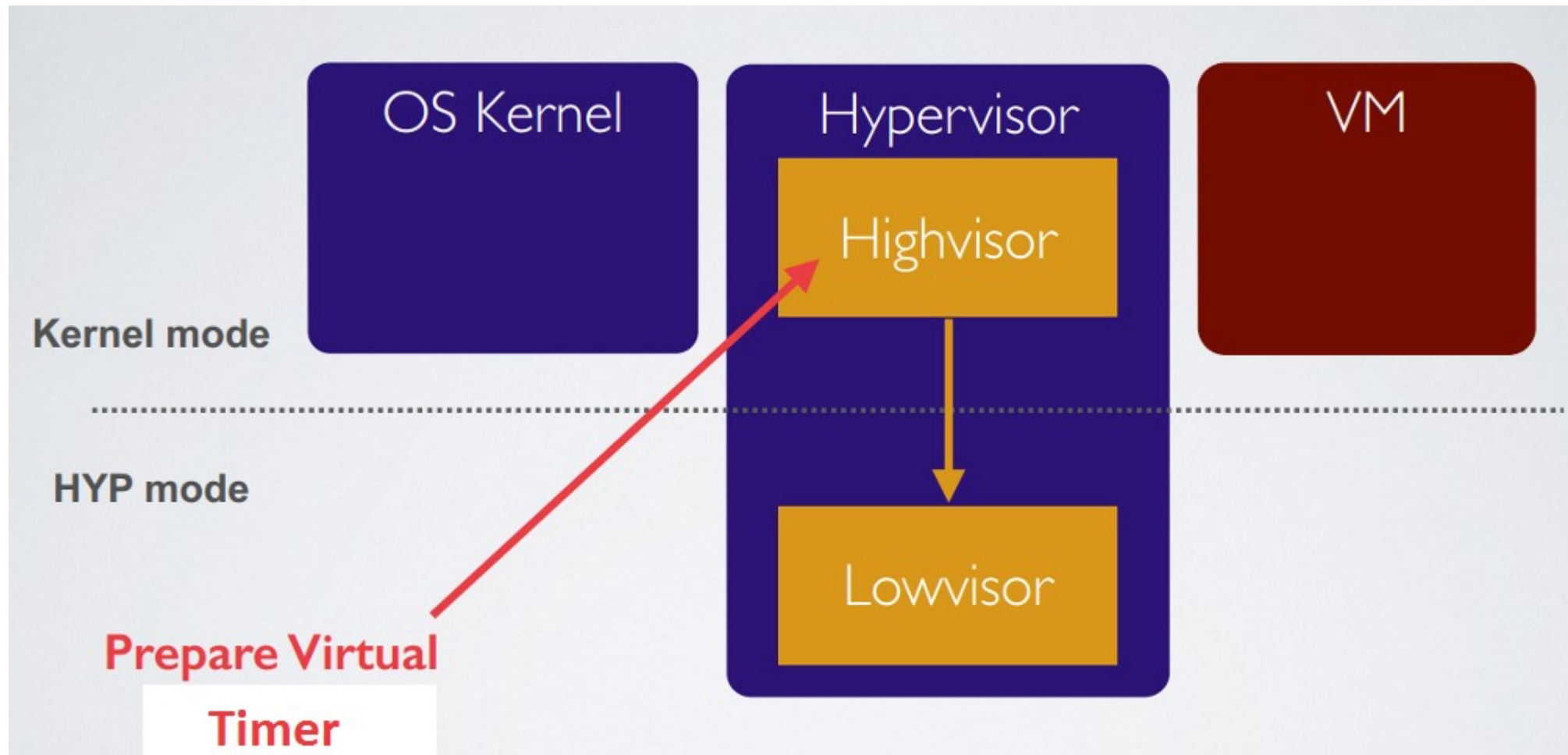


KVM/ARM implements virtual GIC Distributor in software

What happens when a virtual interrupt occurs?

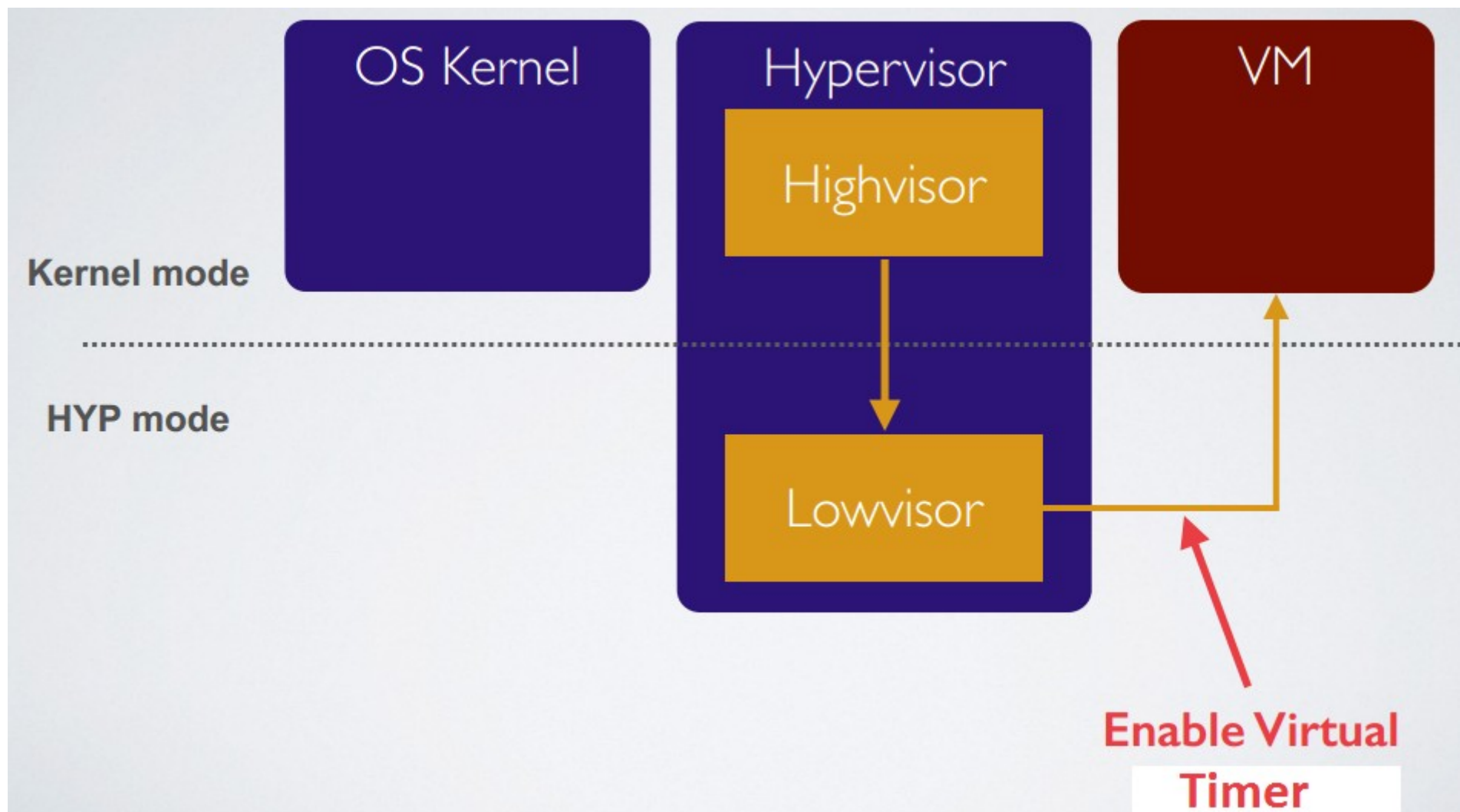


Timer Virtualization



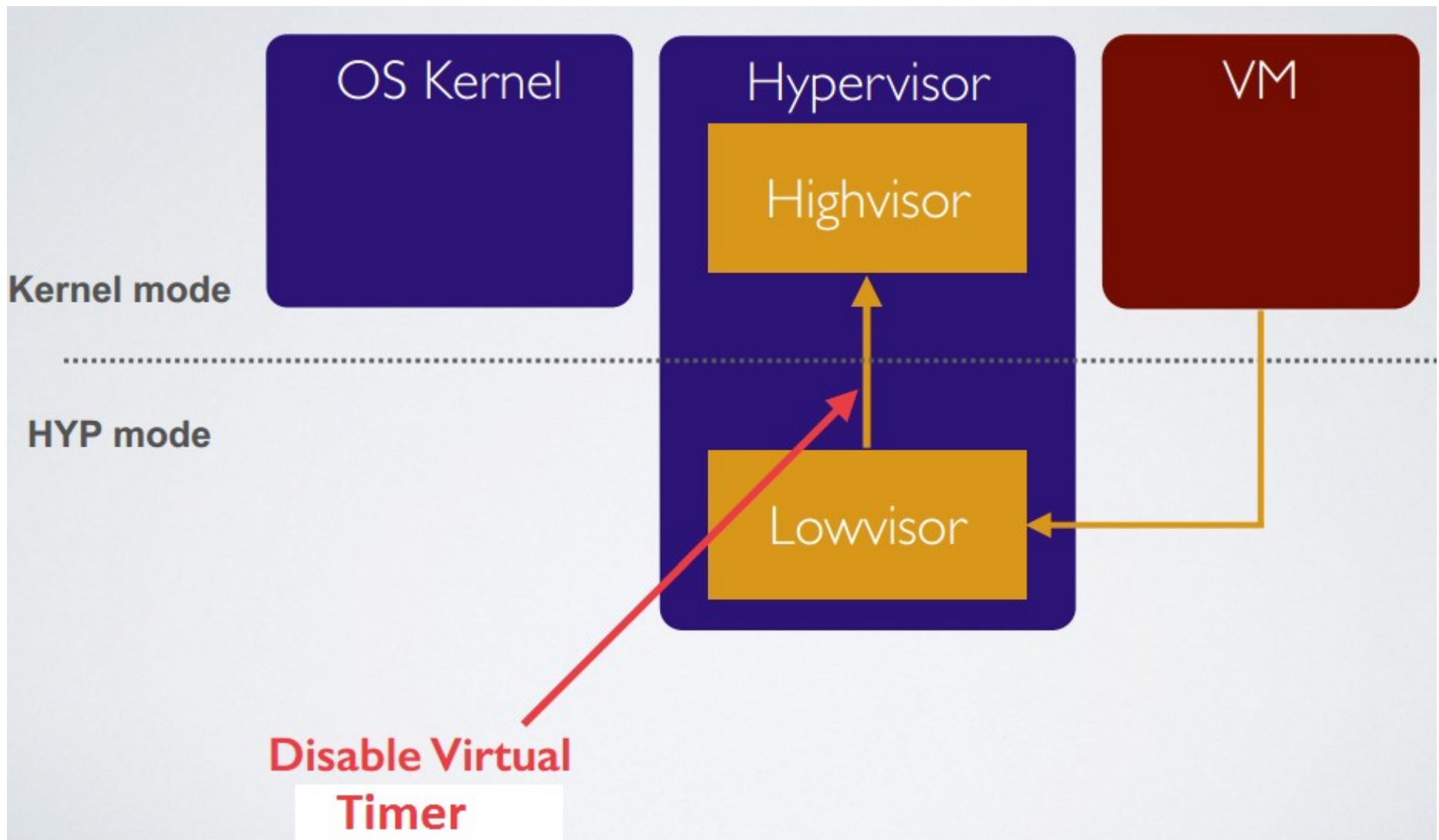


Timer Virtualization





Timer Virtualization



Virtual timers cannot directly raise virtual interrupts, but always raise hardware interrupts



Implementation & Experimental Setup



Arndale Board - Exynos 5250
Dual-Core Cortex A-15 @ 1.4 GHz
eSATA SSD



MacBook Air
Dual-Core Intel i7 @ 1.8 GHz
SATA SSD



Implementation & Experimental Setup

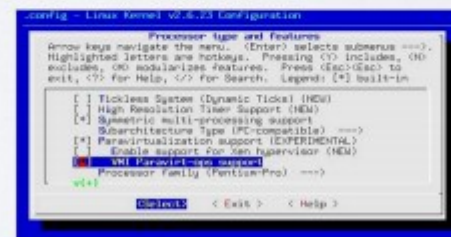
WILLIAM
& MARY



Same amount of memory
across all runs



Kernel version 3.10



Same kernel config



Performance on Micro Benchmark

Two world switches

Micro Test	ARM	ARM no VGIC/vtimers	x86 laptop	x86 server
Hypercall	5,326	2,270	1,336	1,638
Trap	27	27	632	821
I/O Kernel	5,990	2,850	3,19	3,291
I/O User	10,119	6,704	10,98	12,218
...	14,366	32,951	17,1	21,177
EOI+ACK	427	13,726	2	2,305

Cost of saving and
restoring VGIC states

hardware support
to save and restore state

Table 3: Micro-Architectural Cycle Counts



Performance on Micro Benchmark



Switching the hardware mode

Micro Test	ARM	ARM no VGIC/vtimers	x86 laptop	x86 server
Hypercall	5,326	2,270	1,336	1,638
Trap	27	27	632	821
I/O Kernel	5,990	2,850	3,190	3,291
I/O User	10,119	6,704	10,955	12,218
IPI	14,366	32,951	17,153	21,177
IOI+ACK	427	13,726	2,153	2,305

Manipulate two registers to trap

Table 3: Micro-Architectural Cycle Counts

hardware supported save and restore state



Performance on Micro Benchmark

cost of an I/O
operation
from the VM to
a device

Micro Test	ARM	ARM no VGIC/vtimers	x86 laptop	x86 server
Hypercall	5,326	2,270	1,336	1,638
Trap	27	27	632	821
I/O Kernel	5,990	2,850	3,190	3,291
I/O User	10,119	6,704	10,985	12,218
IPI	14,366	32,951	17,138	21,177
IO+ACK	427	13,526	2,043	2,305

Table 3: Micro-Architectural Cycle Count

saving and restoring
VGIC

x86 KVM
saves and restores
additional state lazily

hardware supported
save and restore state



Performance on Micro Benchmark

cost of
Inter-processor
Interrupt

Micro Test	ARM	ARM no VGIC/vtimers	x86 laptop	x86 server
Hypercall	5,326	2,270	1,336	1,638
Trap	27	27	632	821
I/O Kernel	5,990	2,850	3,190	3,291
I/O User	10,119	6,704	10,985	12,218
IPI	14,366	32,951	17,138	21,177
EOI+ACK	427	13,726	2,043	2,305

Hardware supported
VGIC

Table 3: Micro-Architectural Cycle counts

EOI/ACK must be handled
in Root mode in x86



Recommendation to Hardware Designers

Share Kernel Mode Memory Model in Hyp Mode

- Current implementation is different and require extra effort for memory management**

Make VGIC state access fast, or at least infrequent

- Cause overhead because of slow MMIO access. Summary registers describing virtual interrupt state could be an improvement**

Completely avoid IPI traps

- As Virtual Distributors are implemented in software, this cause some overhead.**



Summary



KVM/ARM was the first system to run unmodified guests on hardware

Facilitated by split-mode virtualization

Comparison with KVM



Questions?