

The Challenges of Building AI Infra on Virtualization

Xin He & Hao Hong

System Technologies and Engineering (STE) team, ByteDance



Agenda

- Background
- Design & Implementation
- Future Work

Background



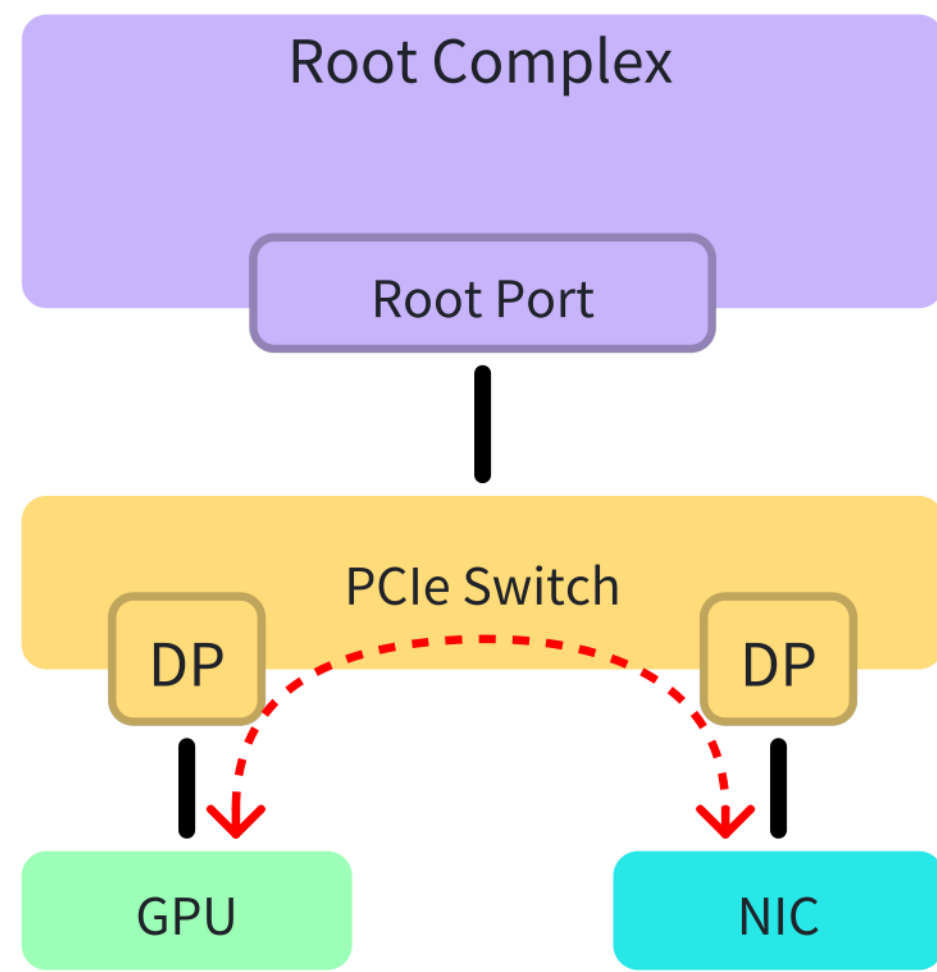
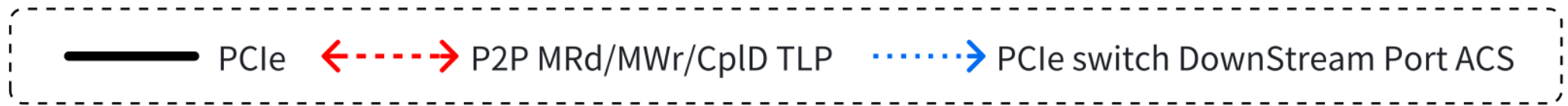
Background

Main issue:

- In virtualization scenarios, there will be a serious performance degradation for PCIe peer-to-peer (p2p) communication due to the enablement of IOMMU
- Various high-precision (millisecond-level) monitoring agents can result in a high number of VMEXITs due to frequent PIO and RDPMC operations

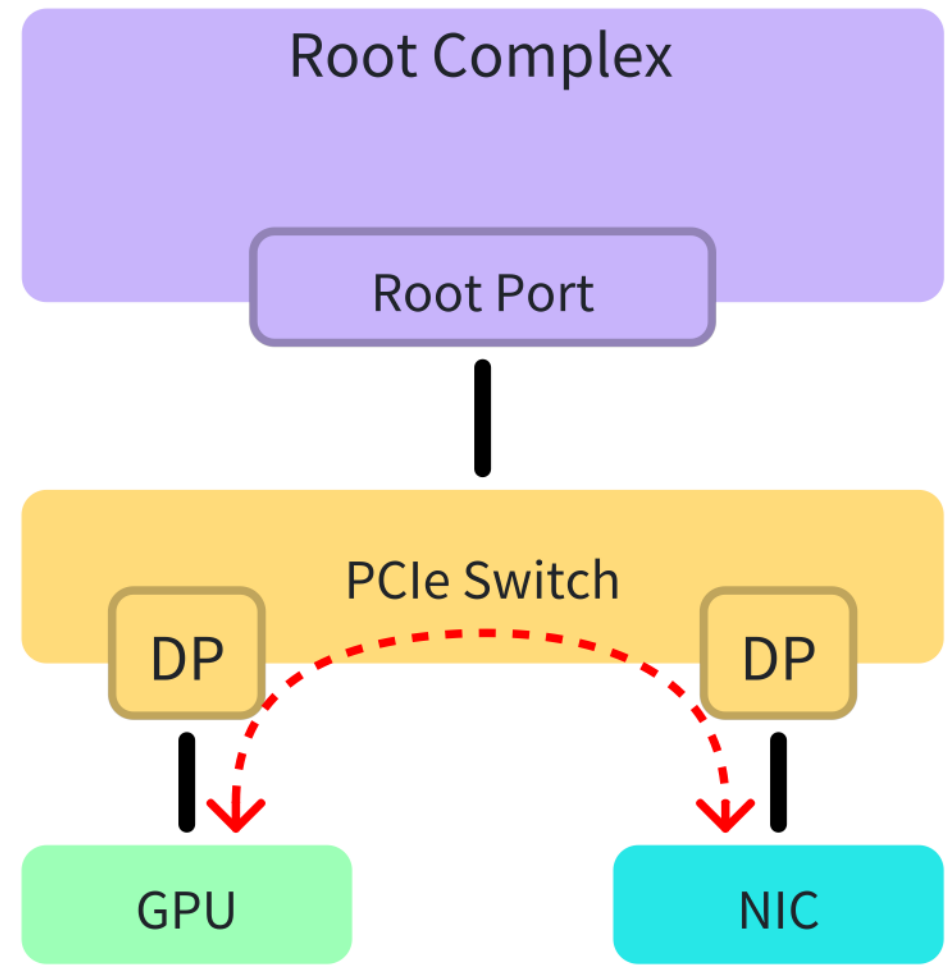
Design & Implementation — Direct P2P

Direct P2P

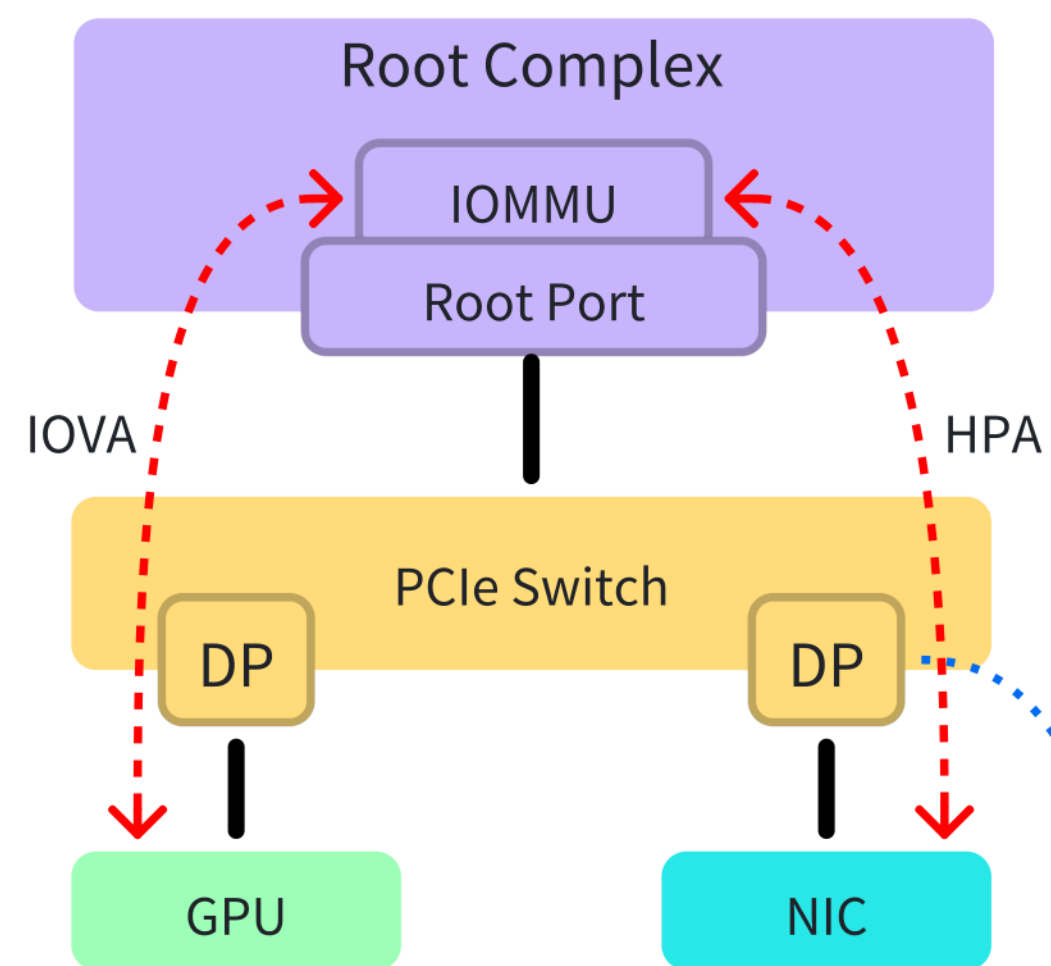


Host(Disable IOMMU)

Direct P2P



Host(Disable IOMMU)

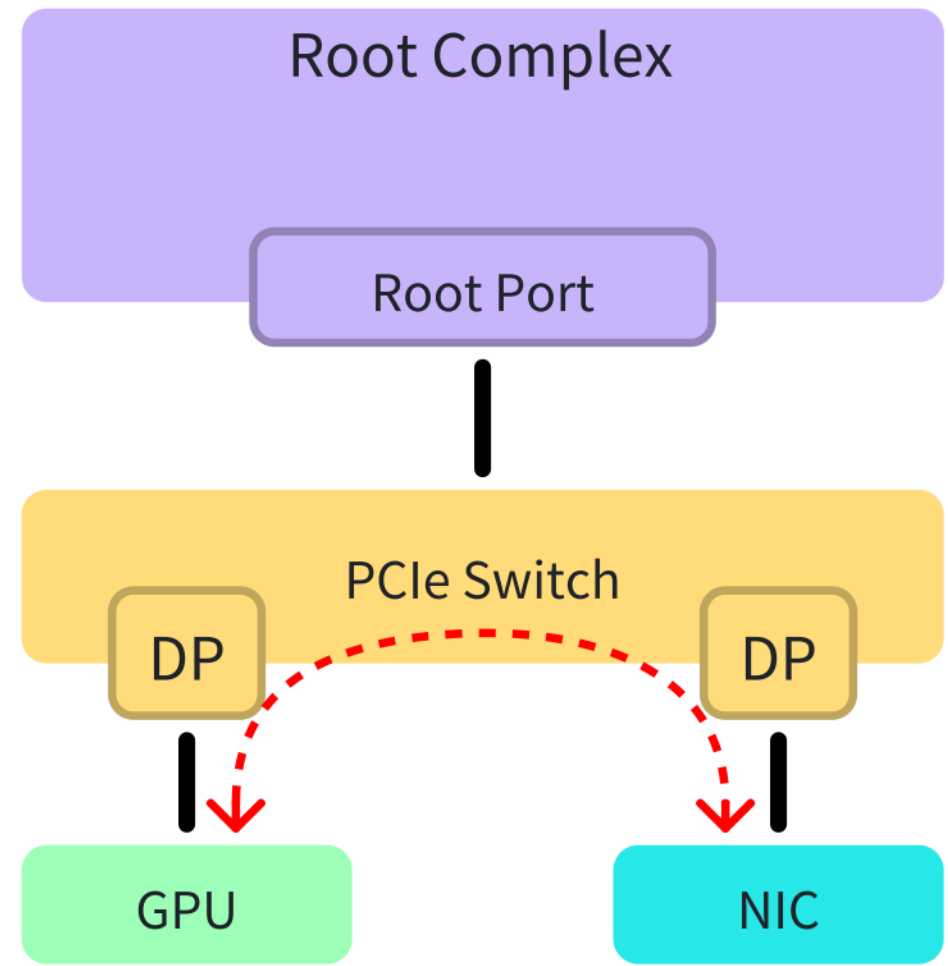
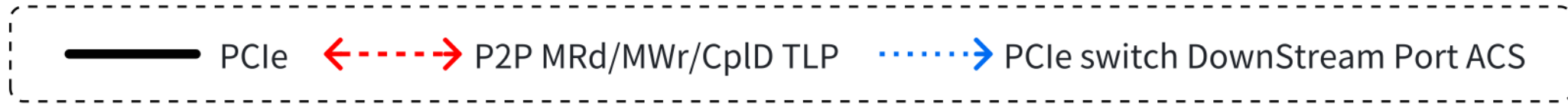


VM(Enable IOMMU & ACS)

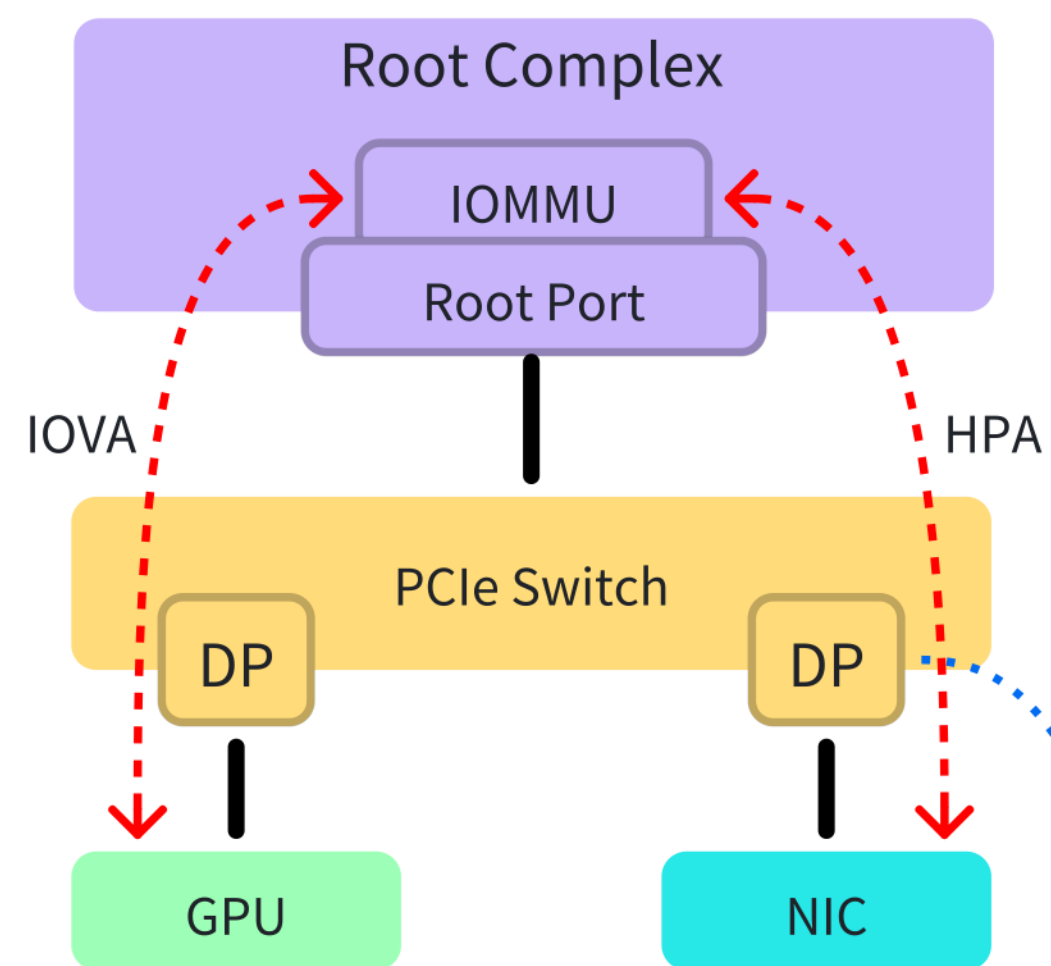
- DP ---- Downstream Port
- RR ---- Request Redirect
- CR ---- Completion Redirect
- DT ---- Direct Translated

Access Control Services (ACS)			
RR	CR	DT
on	on	off

Direct P2P



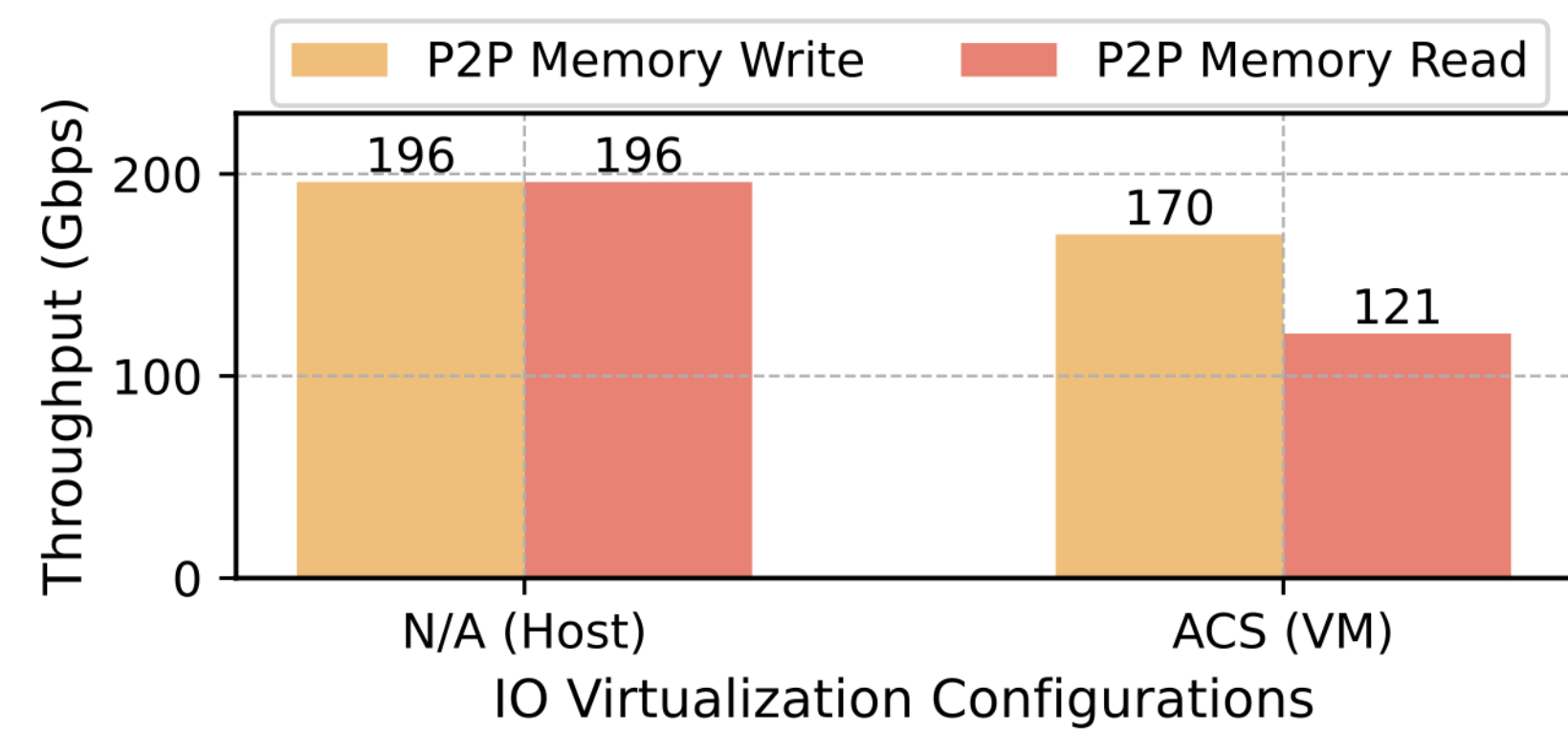
Host(Disable IOMMU)



VM(Enable IOMMU & ACS)

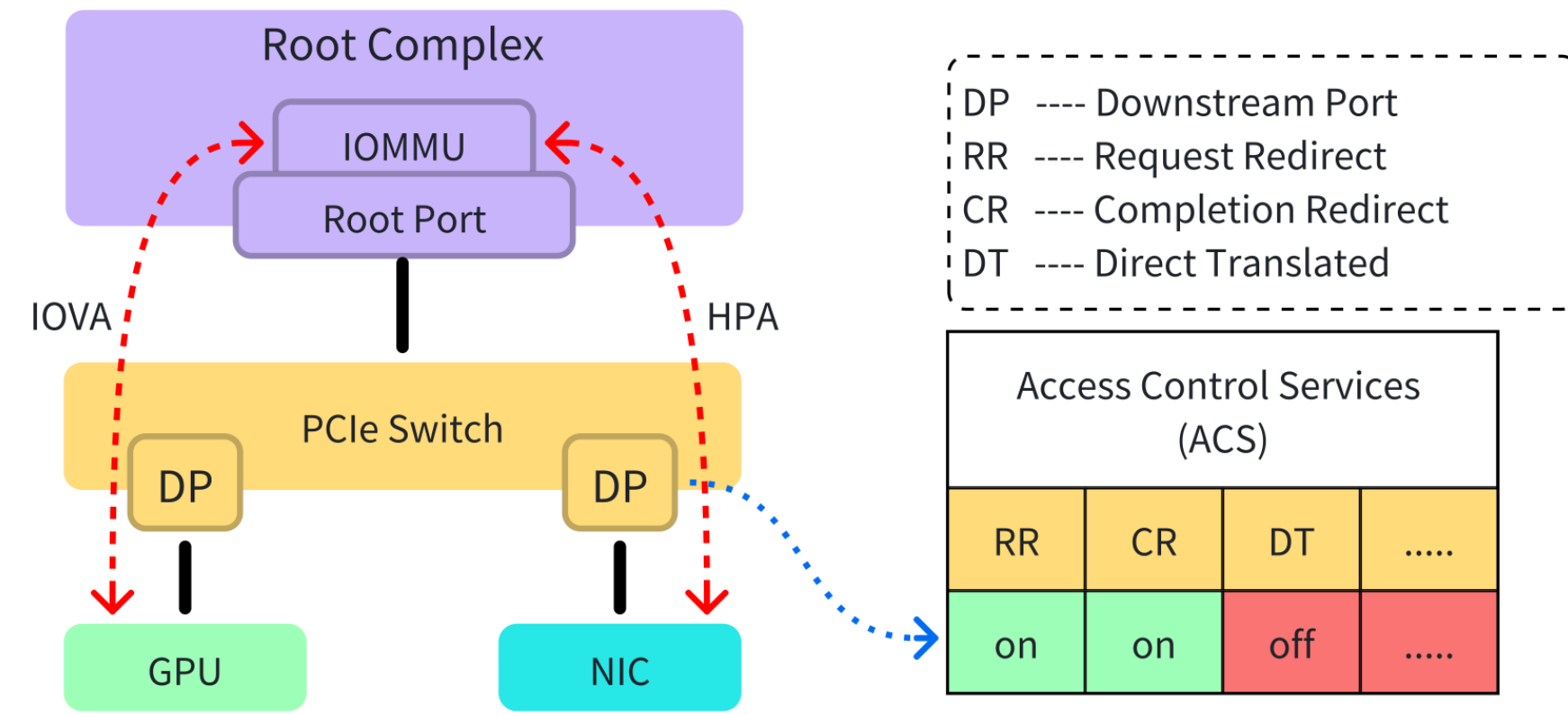
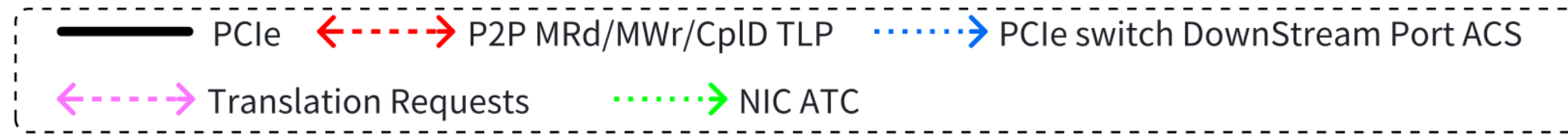
- DP ---- Downstream Port
- RR ---- Request Redirect
- CR ---- Completion Redirect
- DT ---- Direct Translated

Access Control Services (ACS)			
RR	CR	DT
on	on	off



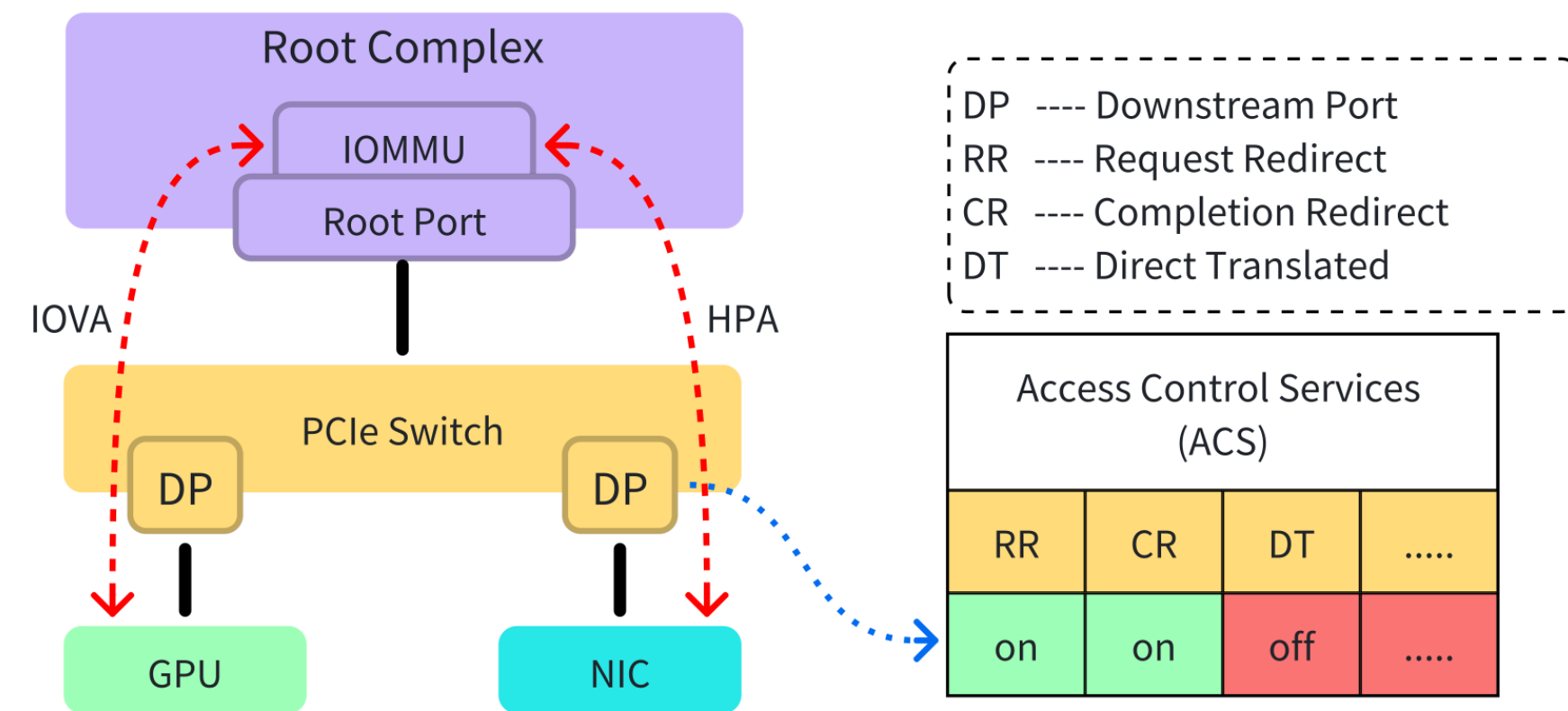
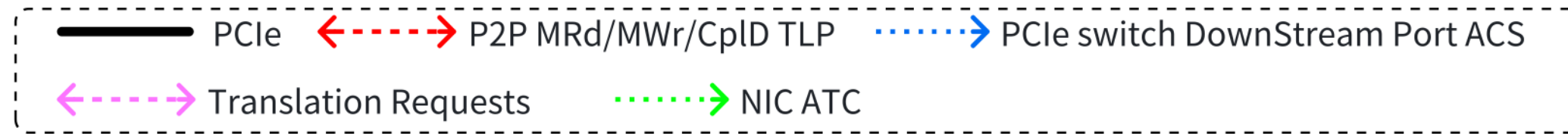
CPU: Intel iceLake
 GPU: A800
 NIC: CX6

Direct P2P

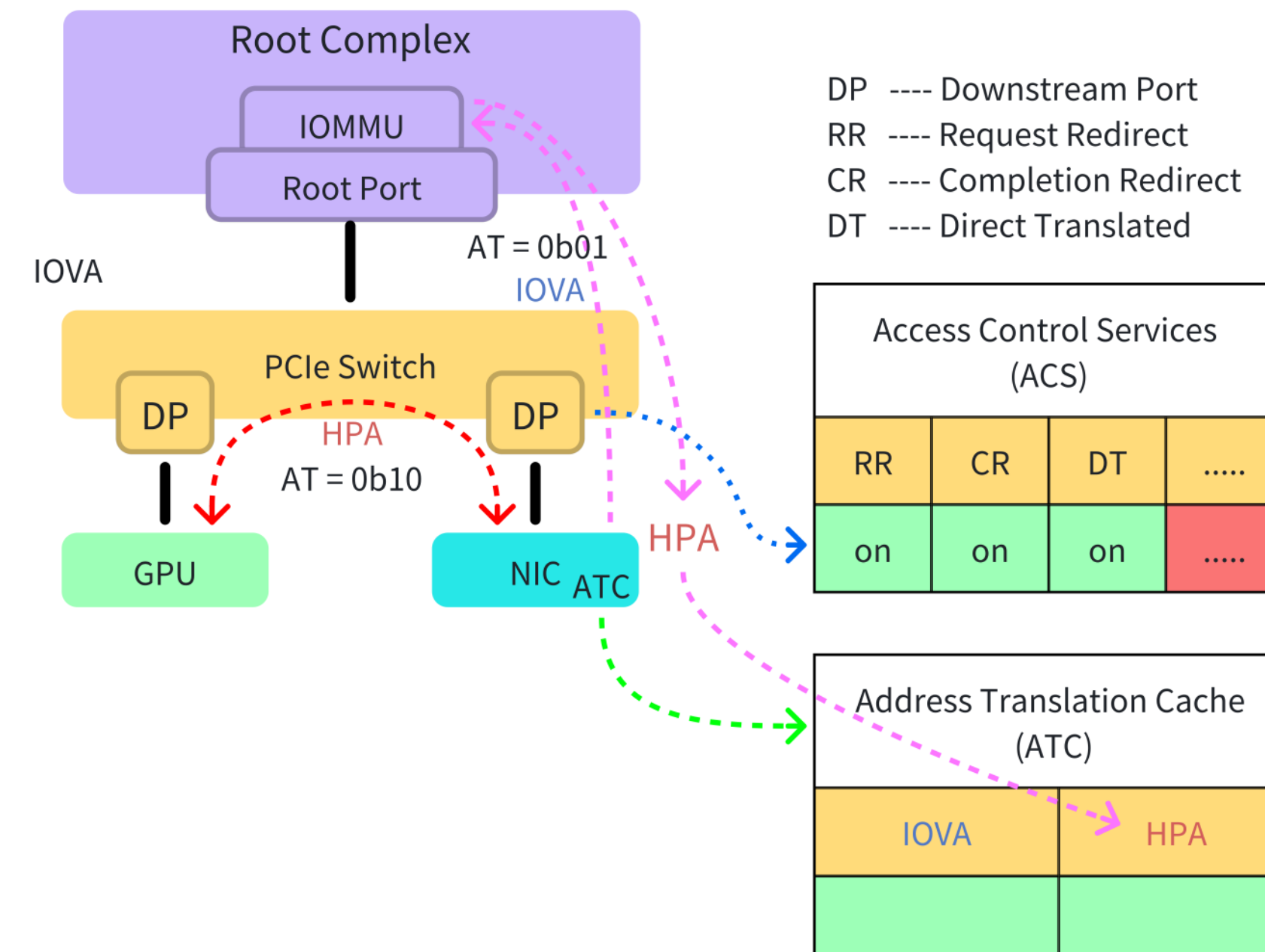


VM(Enable IOMMU & ACS)

Direct P2P

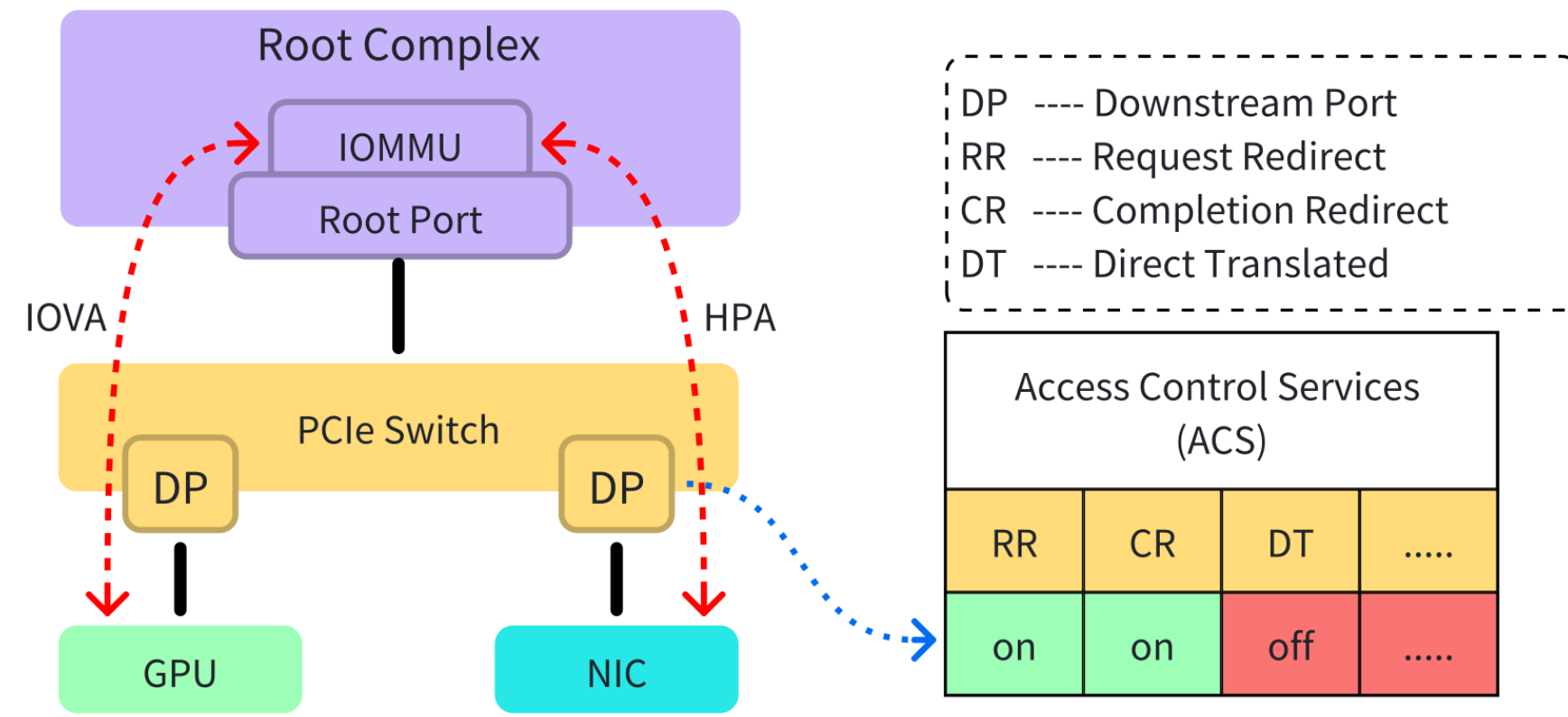
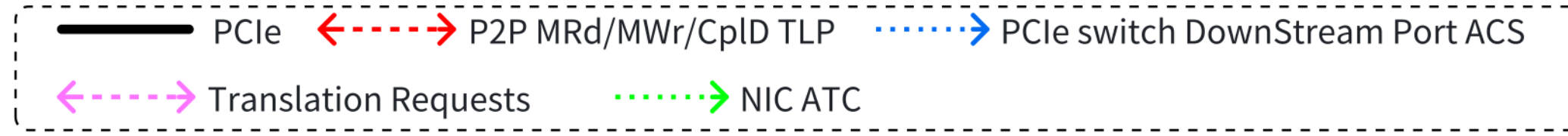


VM(Enable IOMMU & ACS)

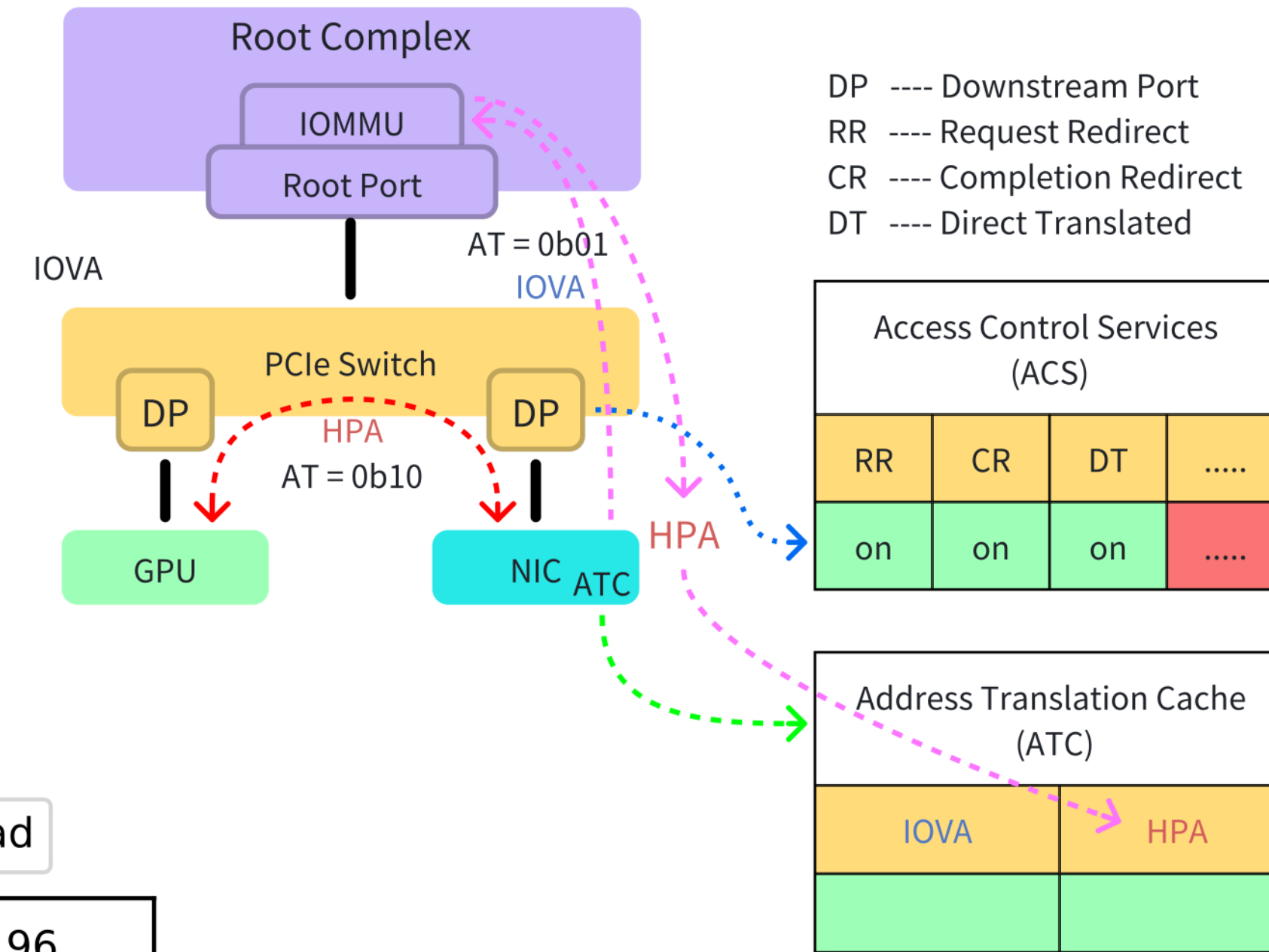


VM(Enable IOMMU & ACS & ATS)

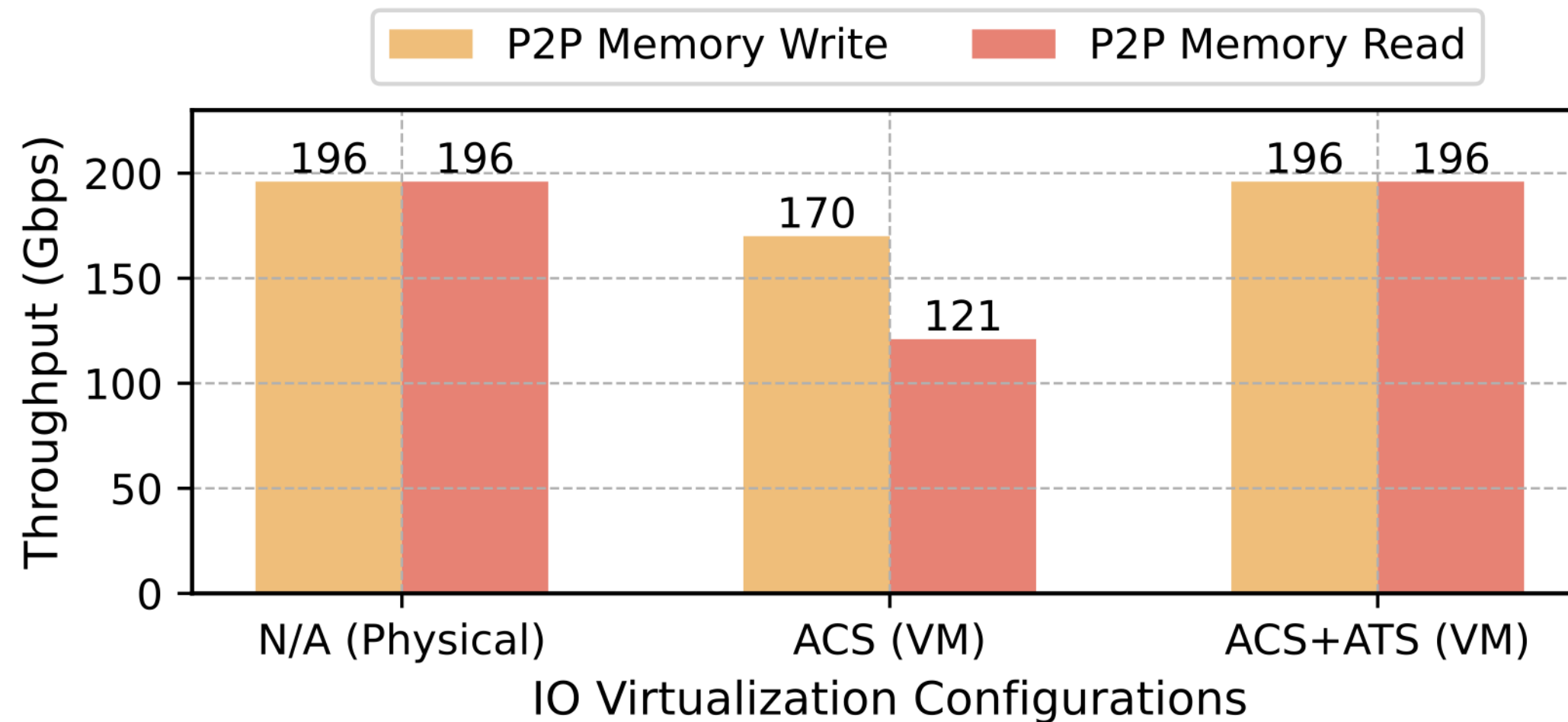
Direct P2P



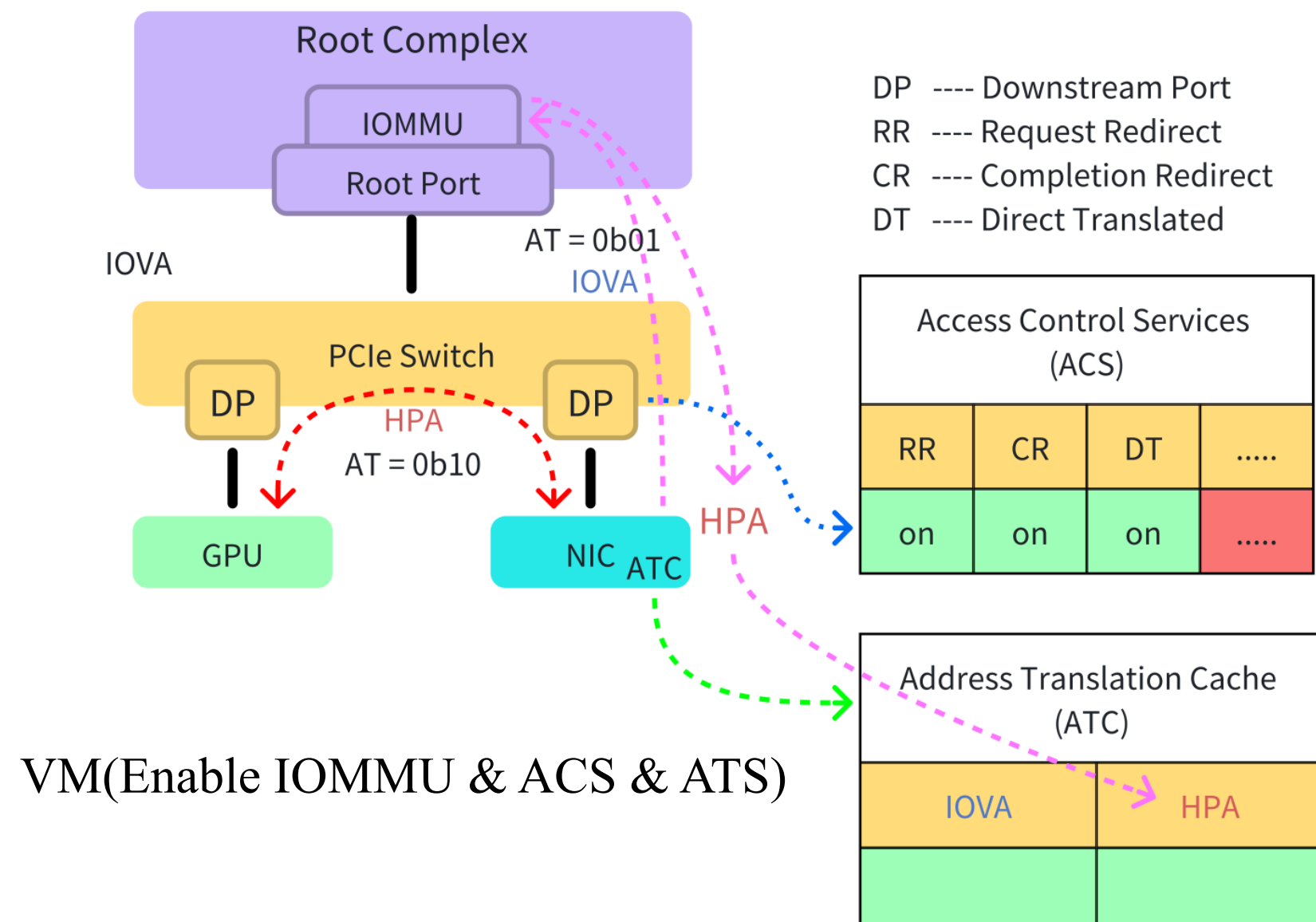
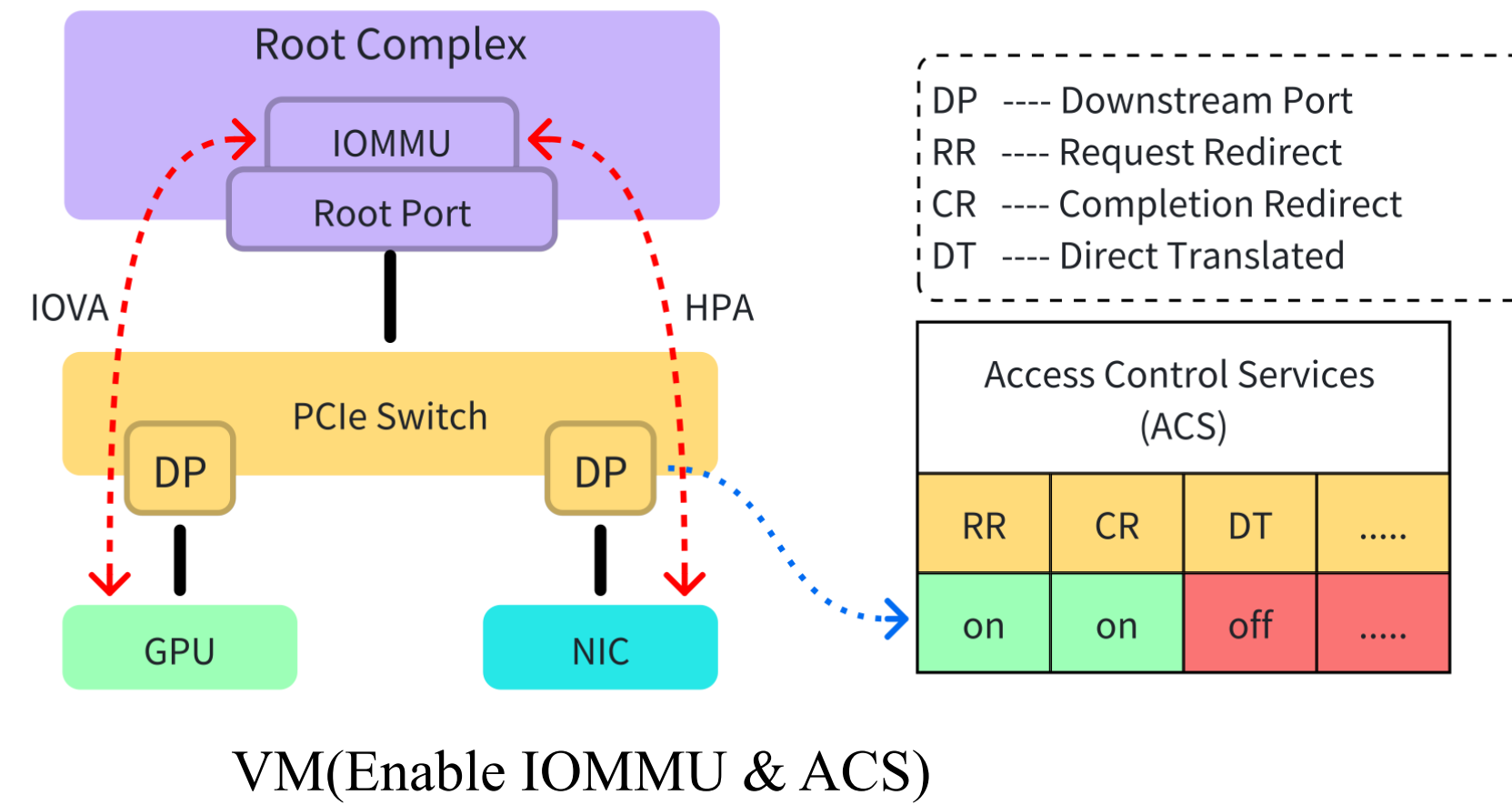
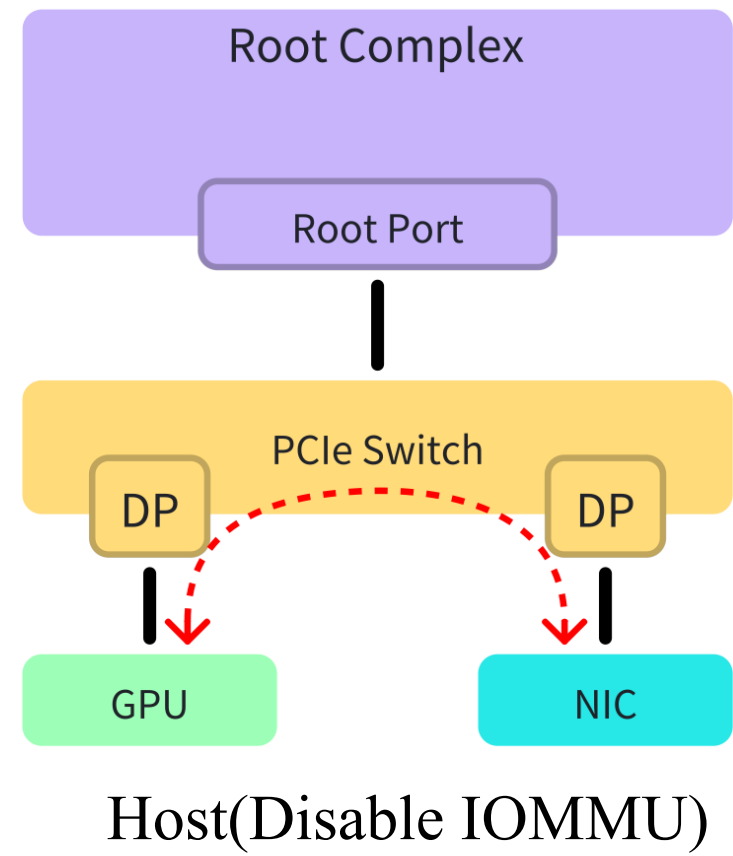
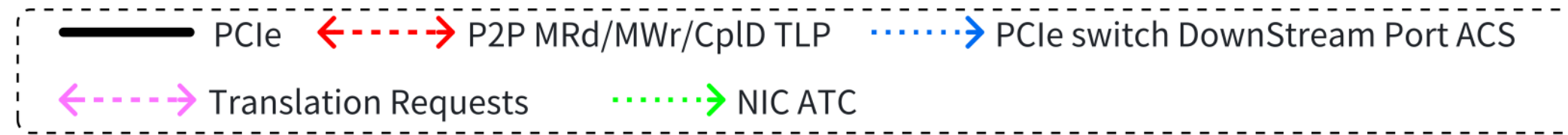
VM(Enable IOMMU & ACS)



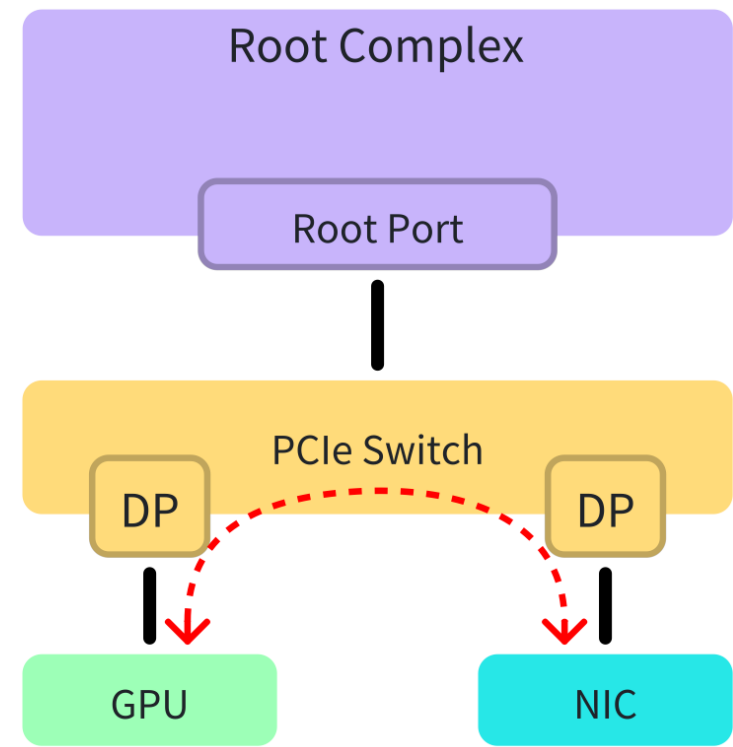
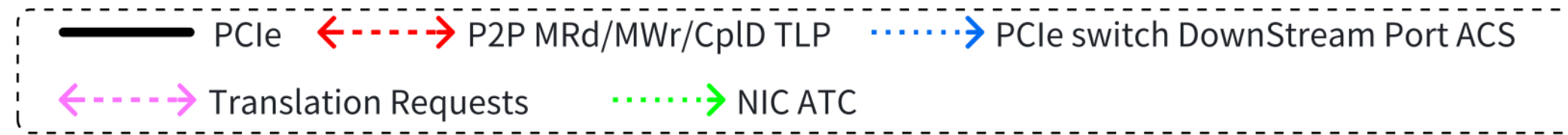
VM(Enable IOMMU & ACS & ATS)



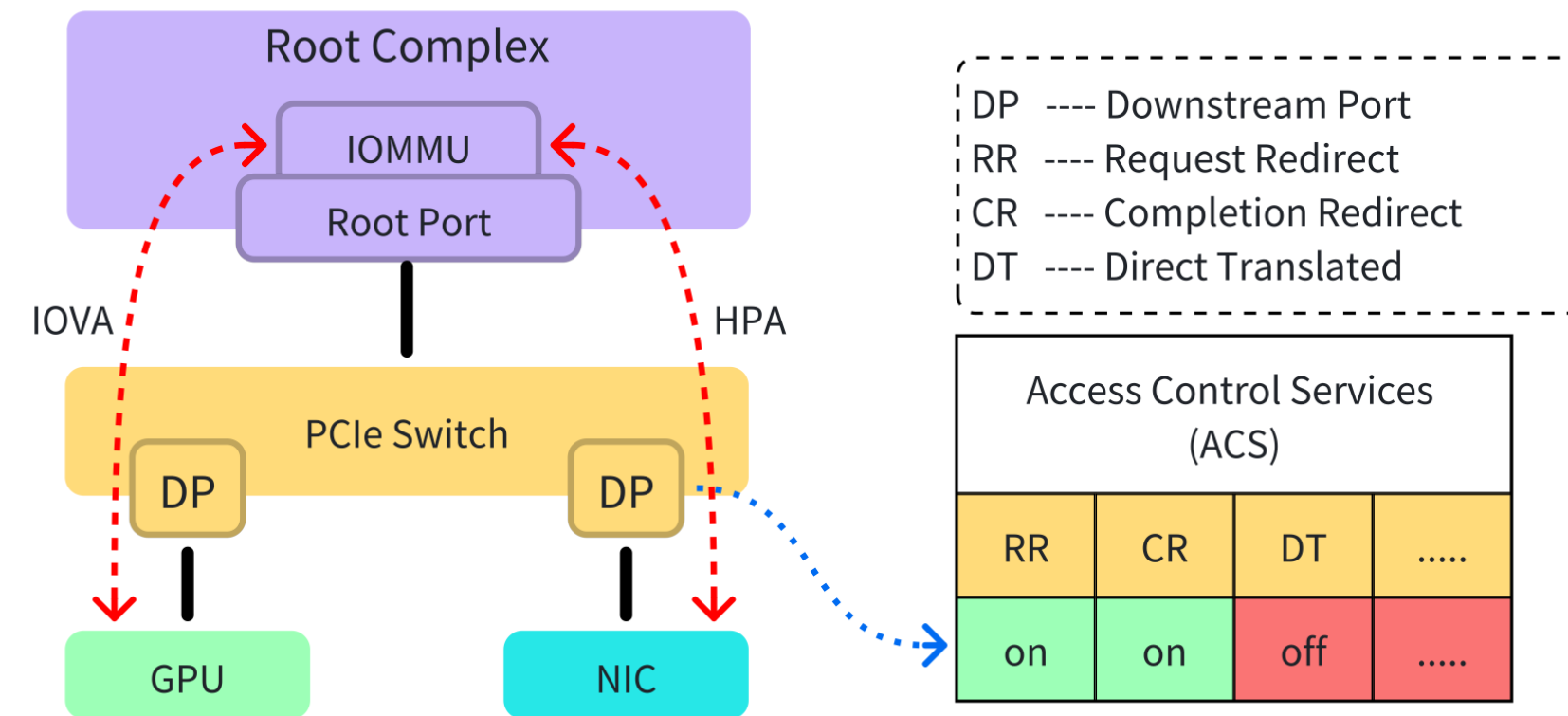
Direct P2P



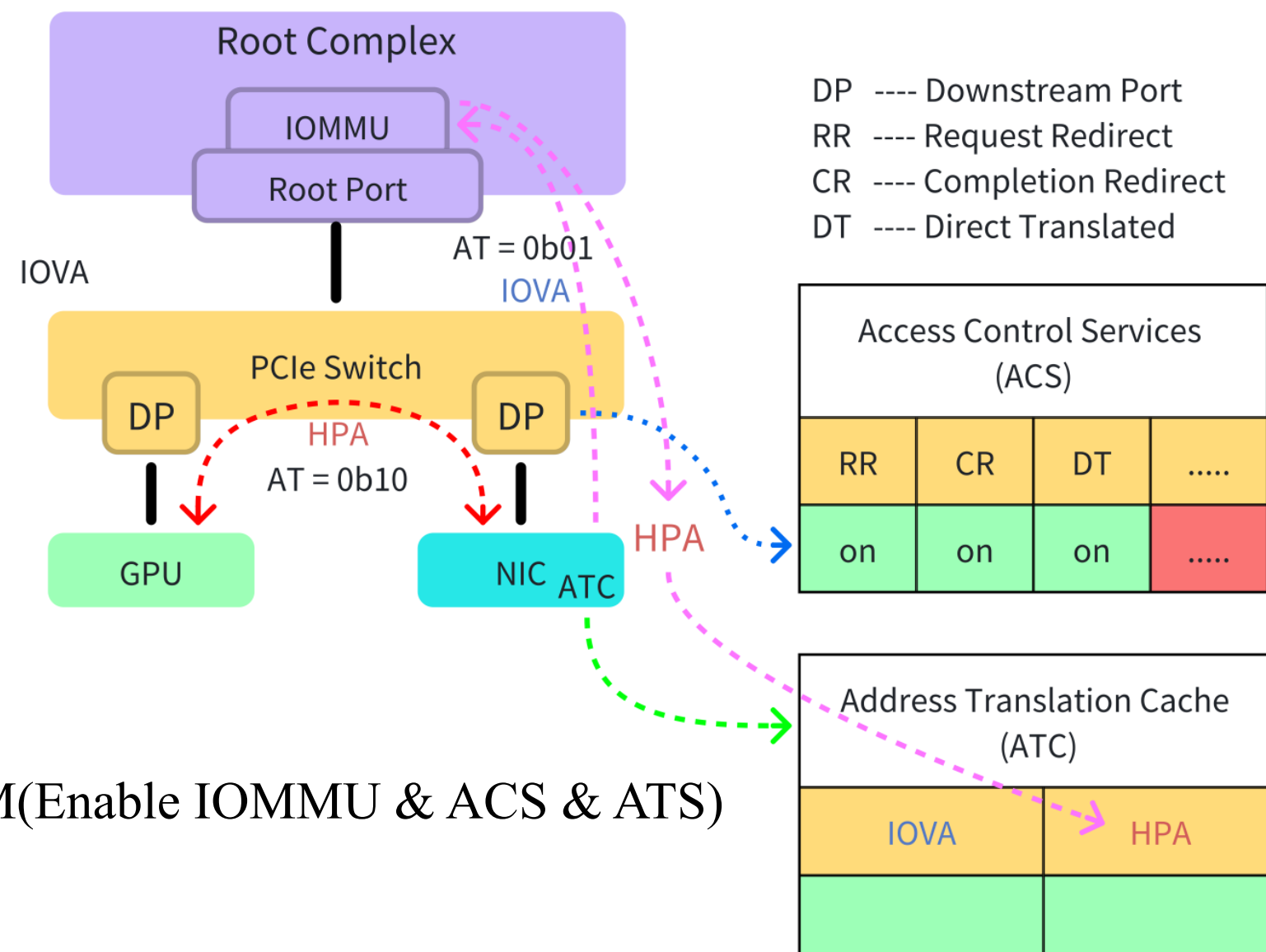
Direct P2P



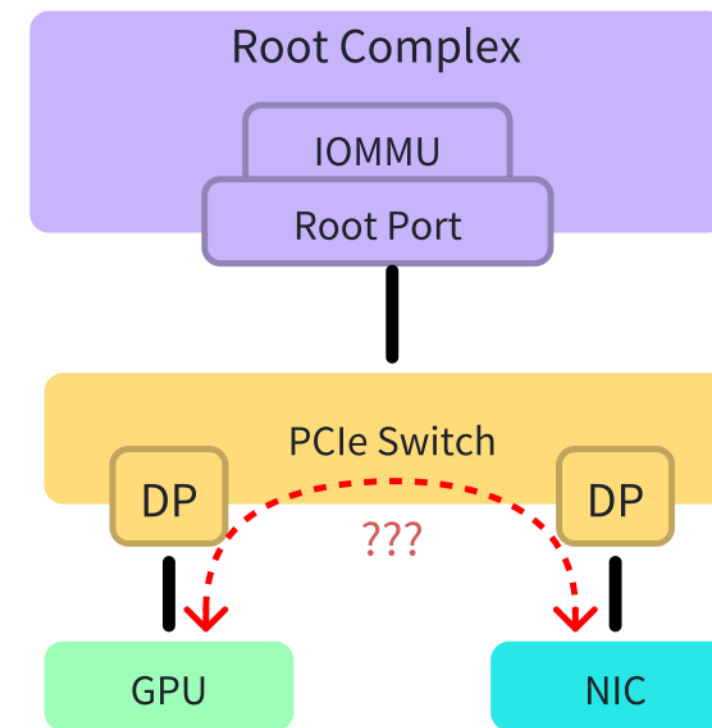
Host(Disable IOMMU)



VM(Enable IOMMU & ACS)

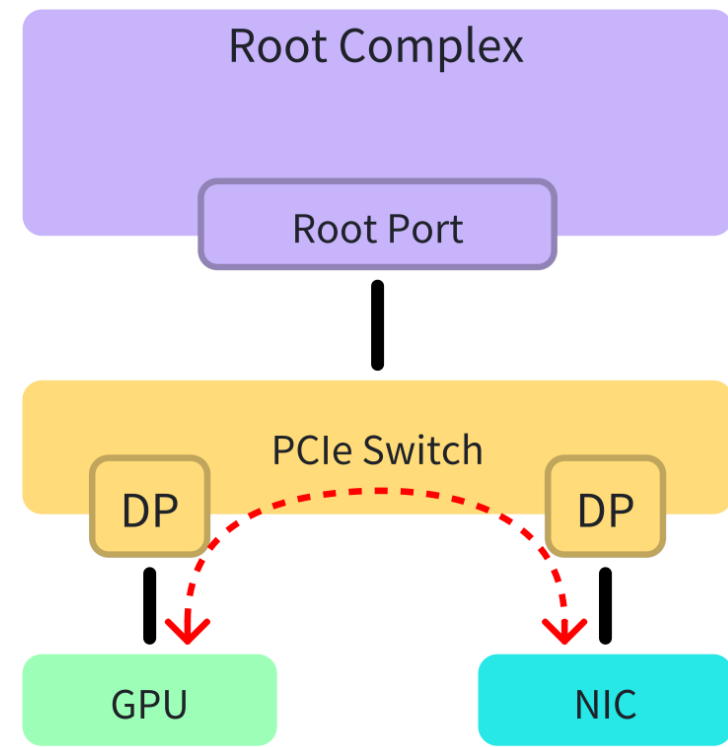
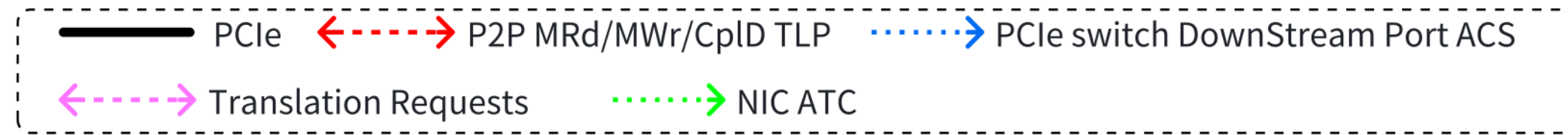


VM(Enable IOMMU & ACS & ATs)

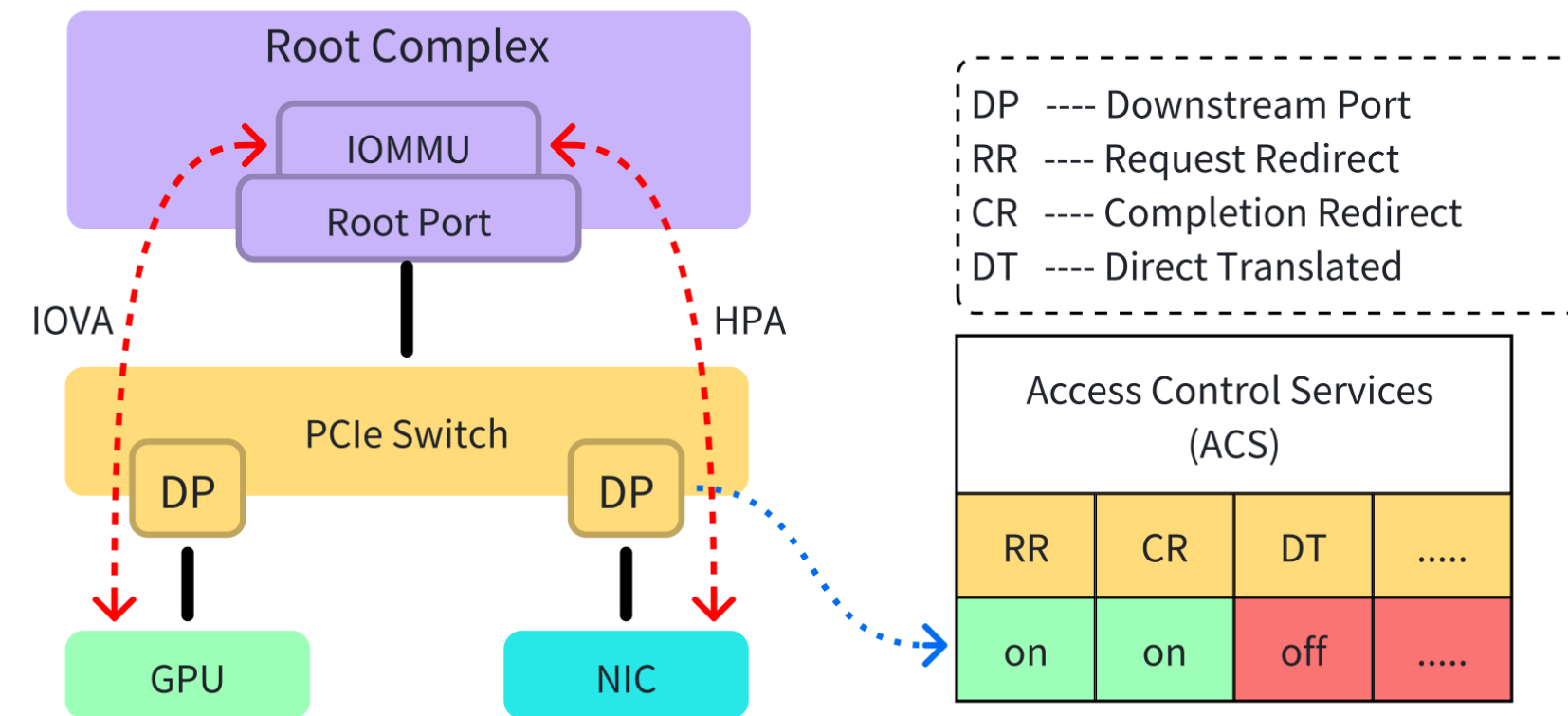


VM(Enable IOMMU & ACS & NO_ATs)

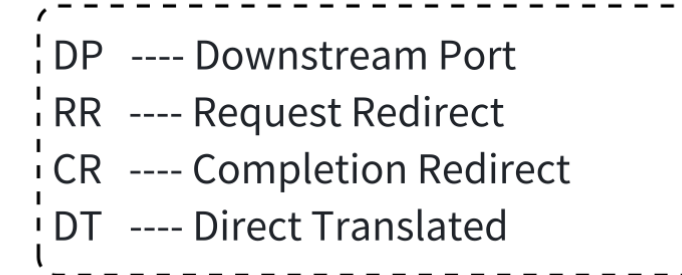
Direct P2P



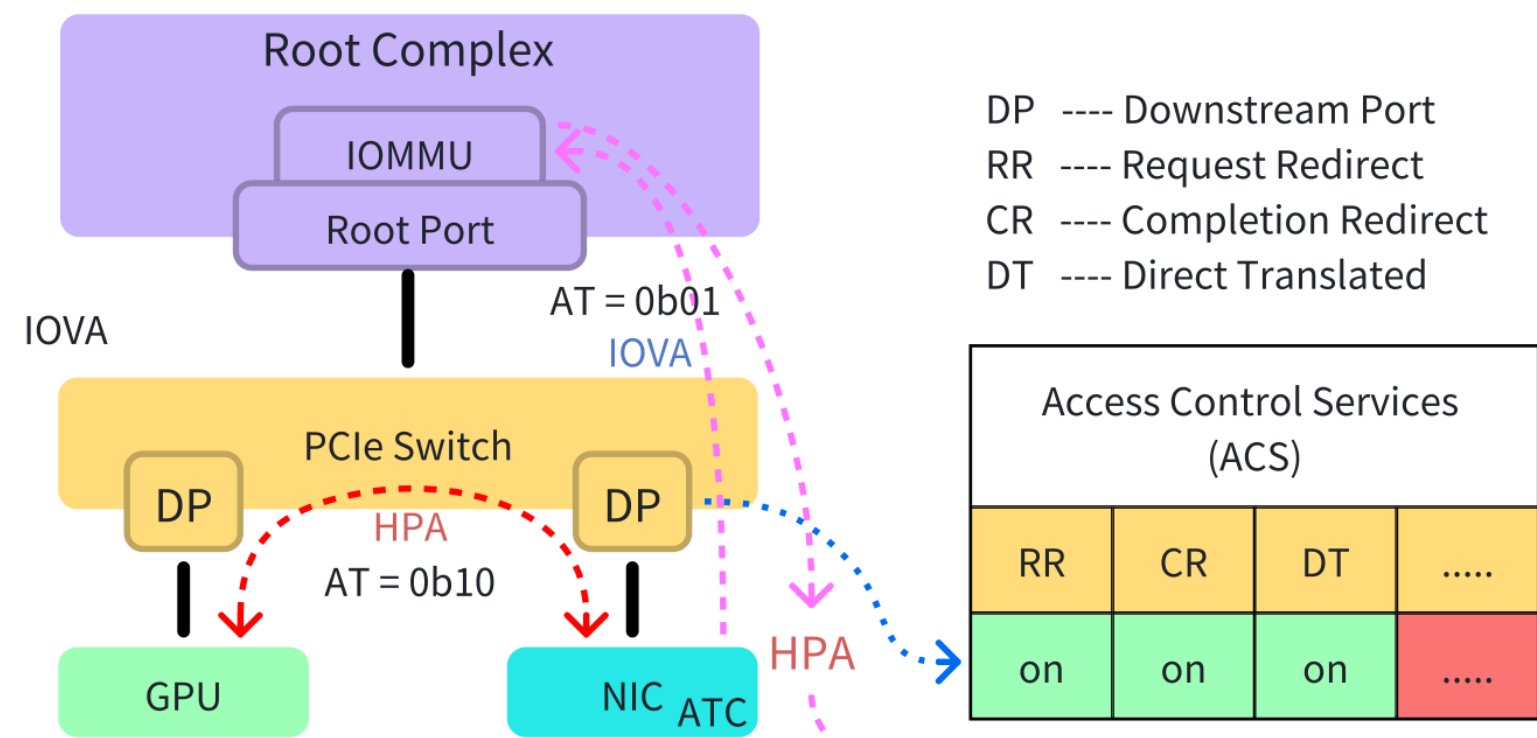
Host(Disable IOMMU)



VM(Enable IOMMU & ACS)



Access Control Services (ACS)			
RR	CR	DT
on	on	off

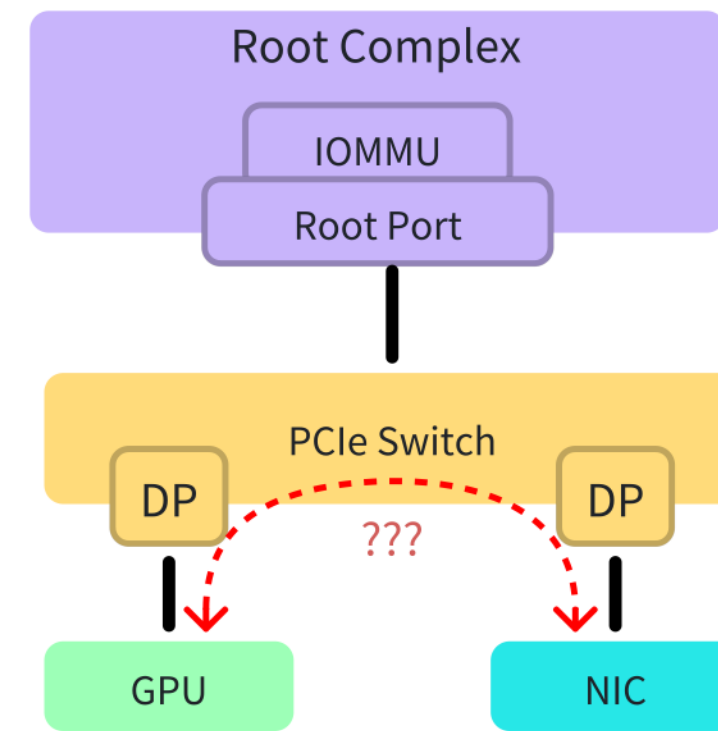


VM(Enable IOMMU & ACS & ATs)

DP ---- Downstream Port
 RR ---- Request Redirect
 CR ---- Completion Redirect
 DT ---- Direct Translated

Access Control Services (ACS)			
RR	CR	DT
on	on	on

Address Translation Cache (ATC)	
IOVA	HPA

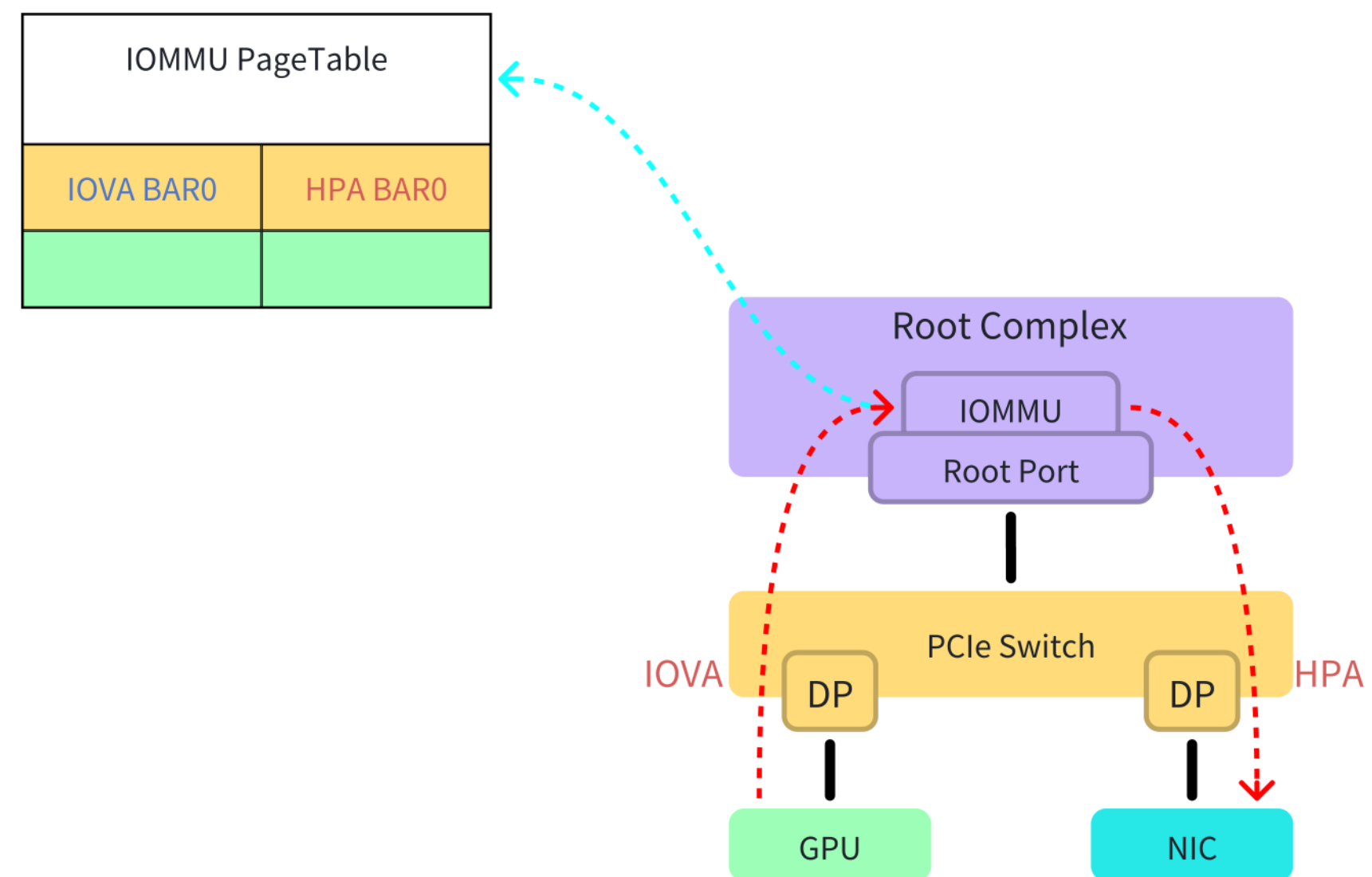
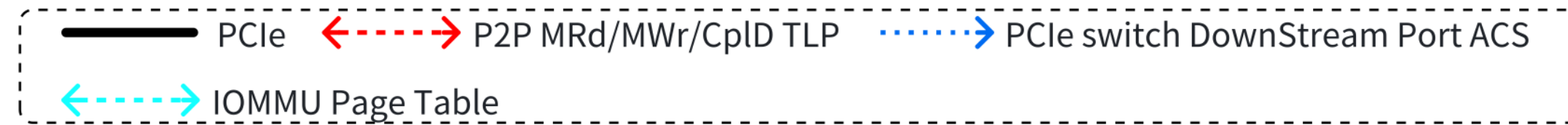


VM(Enable IOMMU & ACS & NO_ATs)

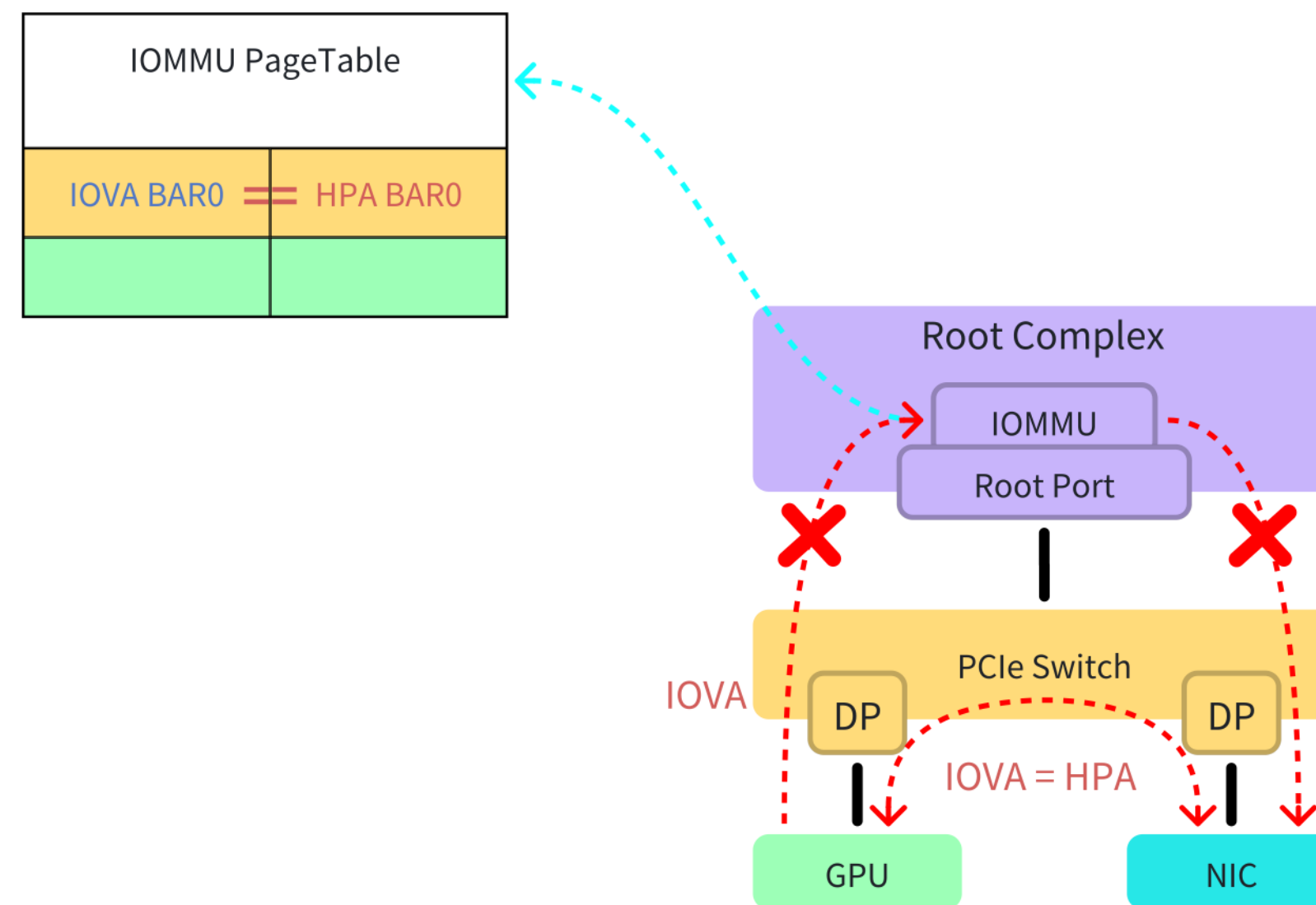
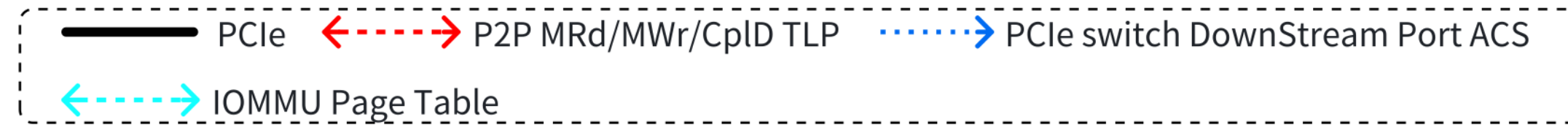
Q1: P2P without IOMMU?

Q2: ACS Configuration

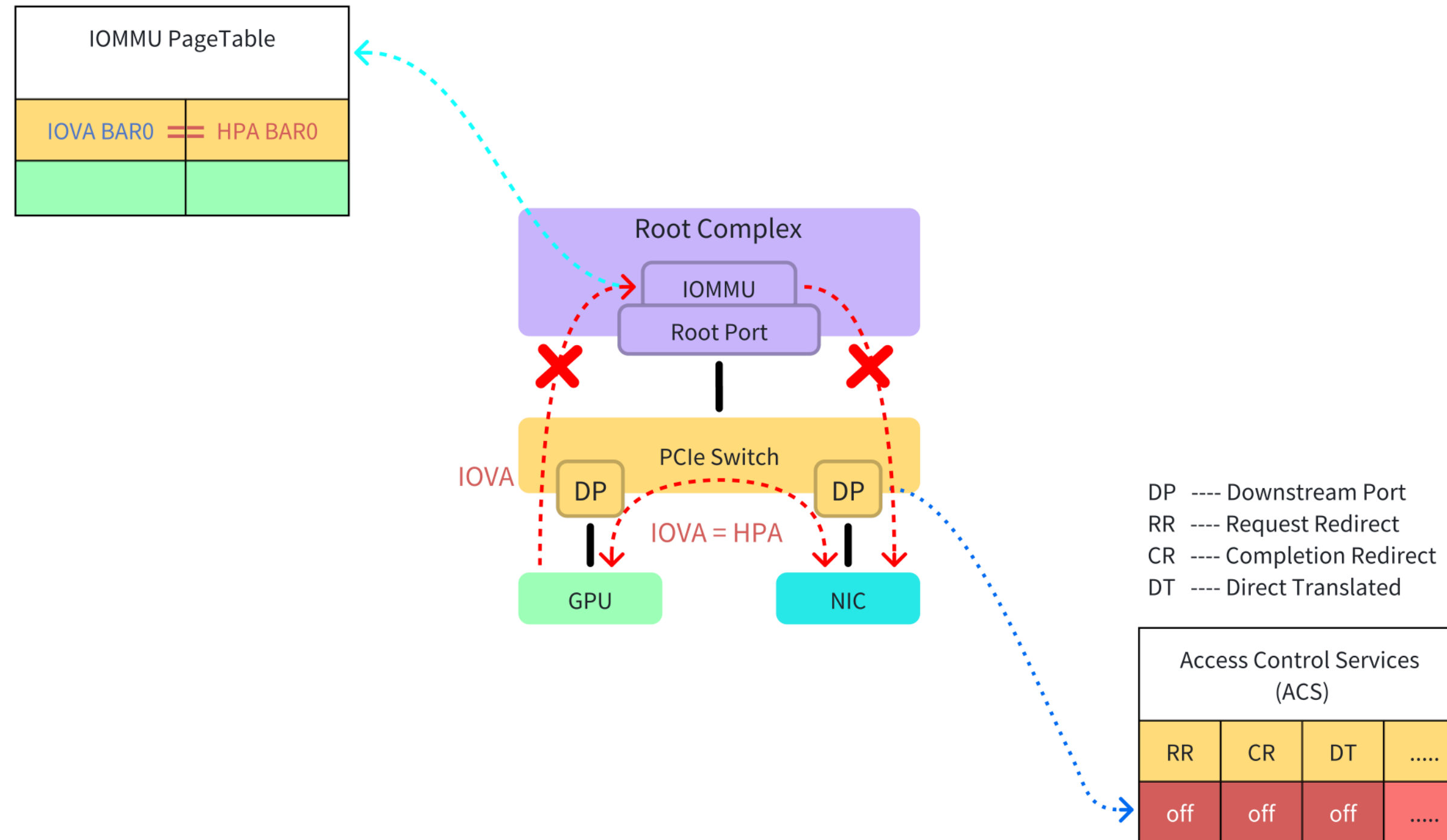
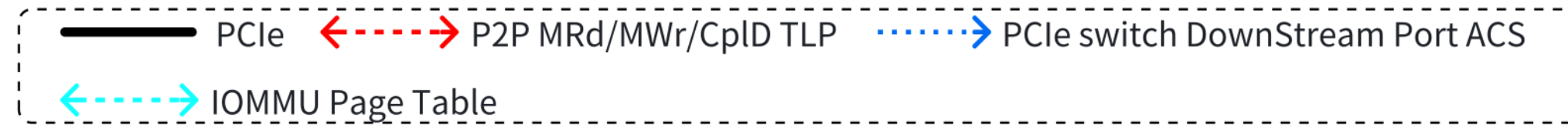
Direct P2P



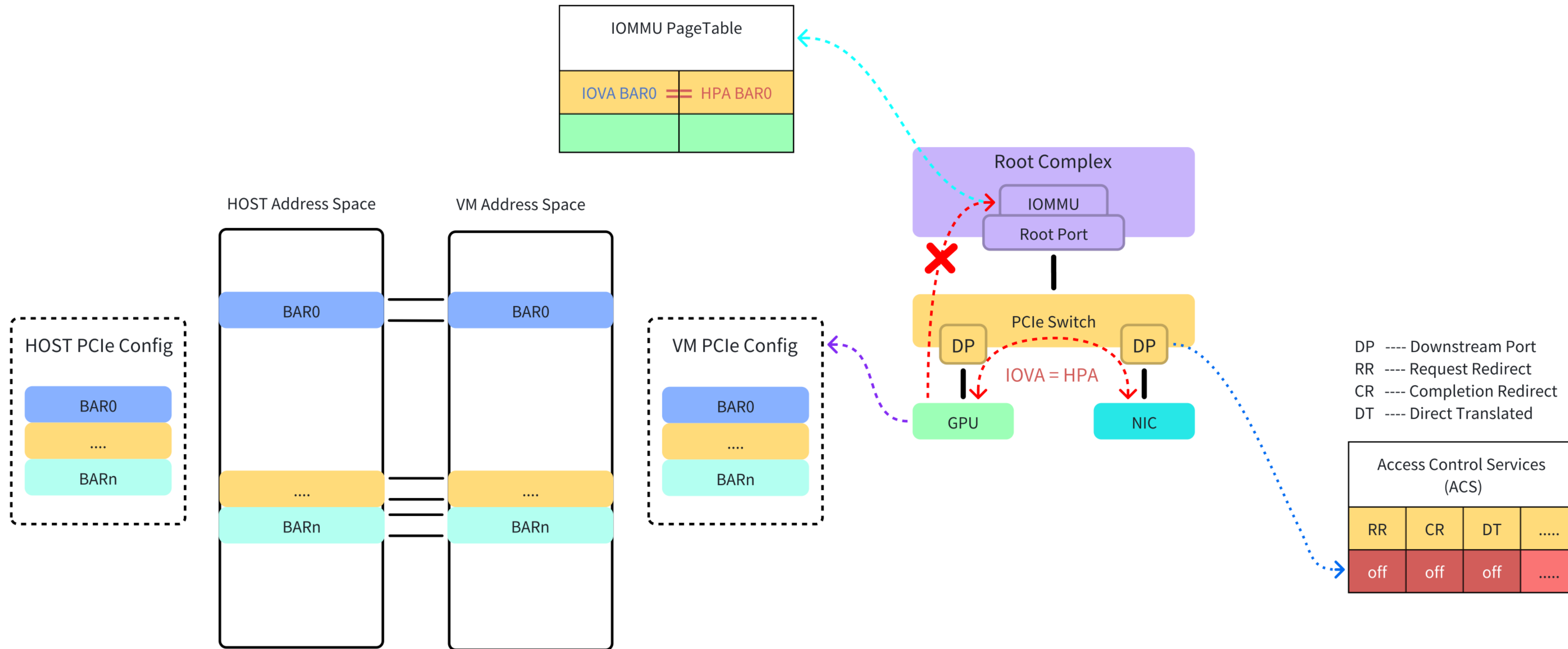
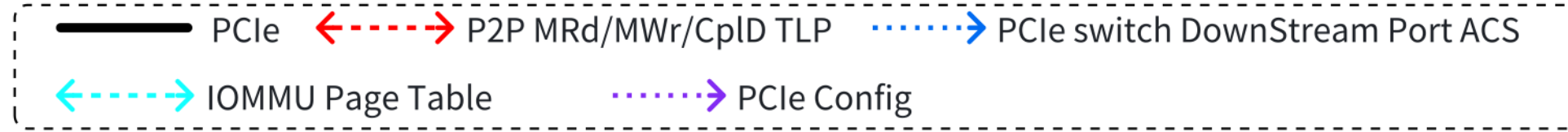
Direct P2P



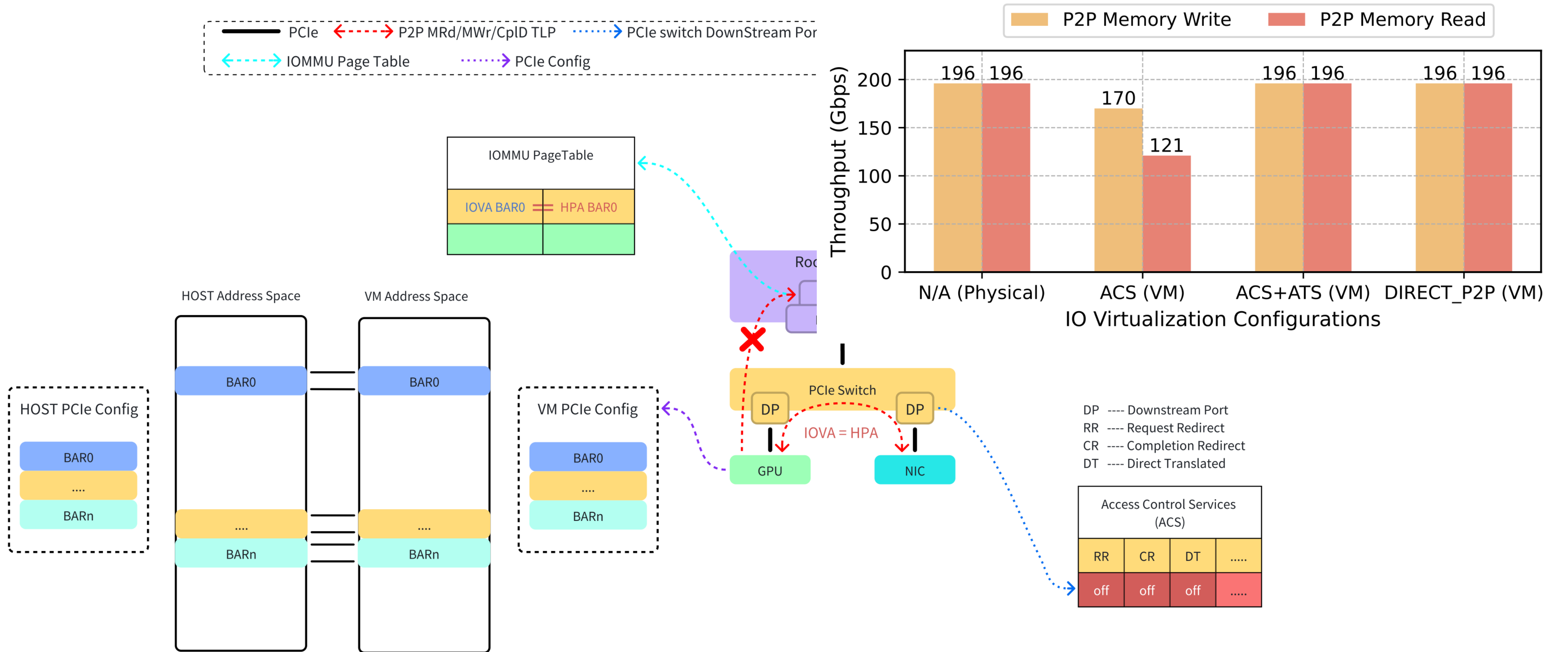
Direct P2P



Direct P2P



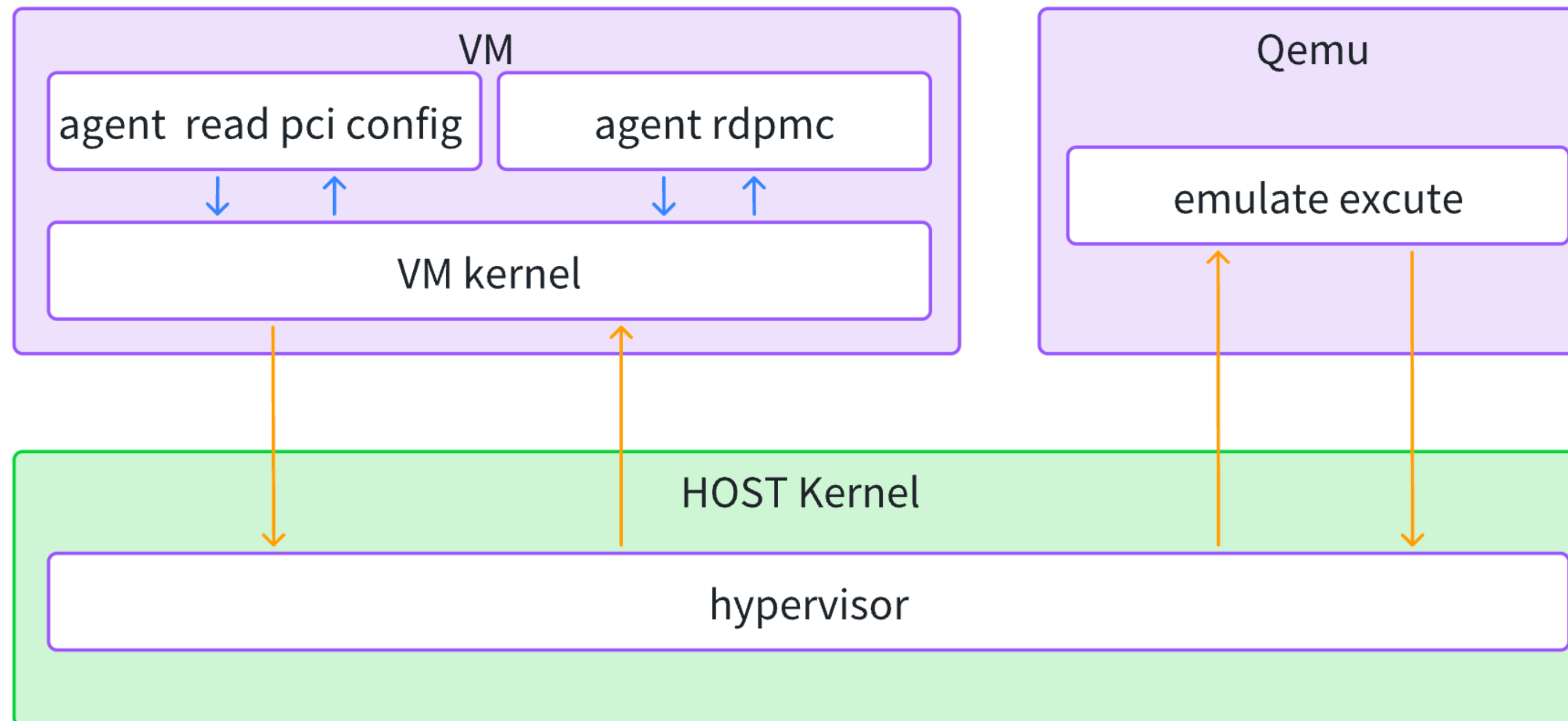
Direct P2P



High-precision monitoring agents

High-precision monitoring agent

PCI config read/rdpmc Path



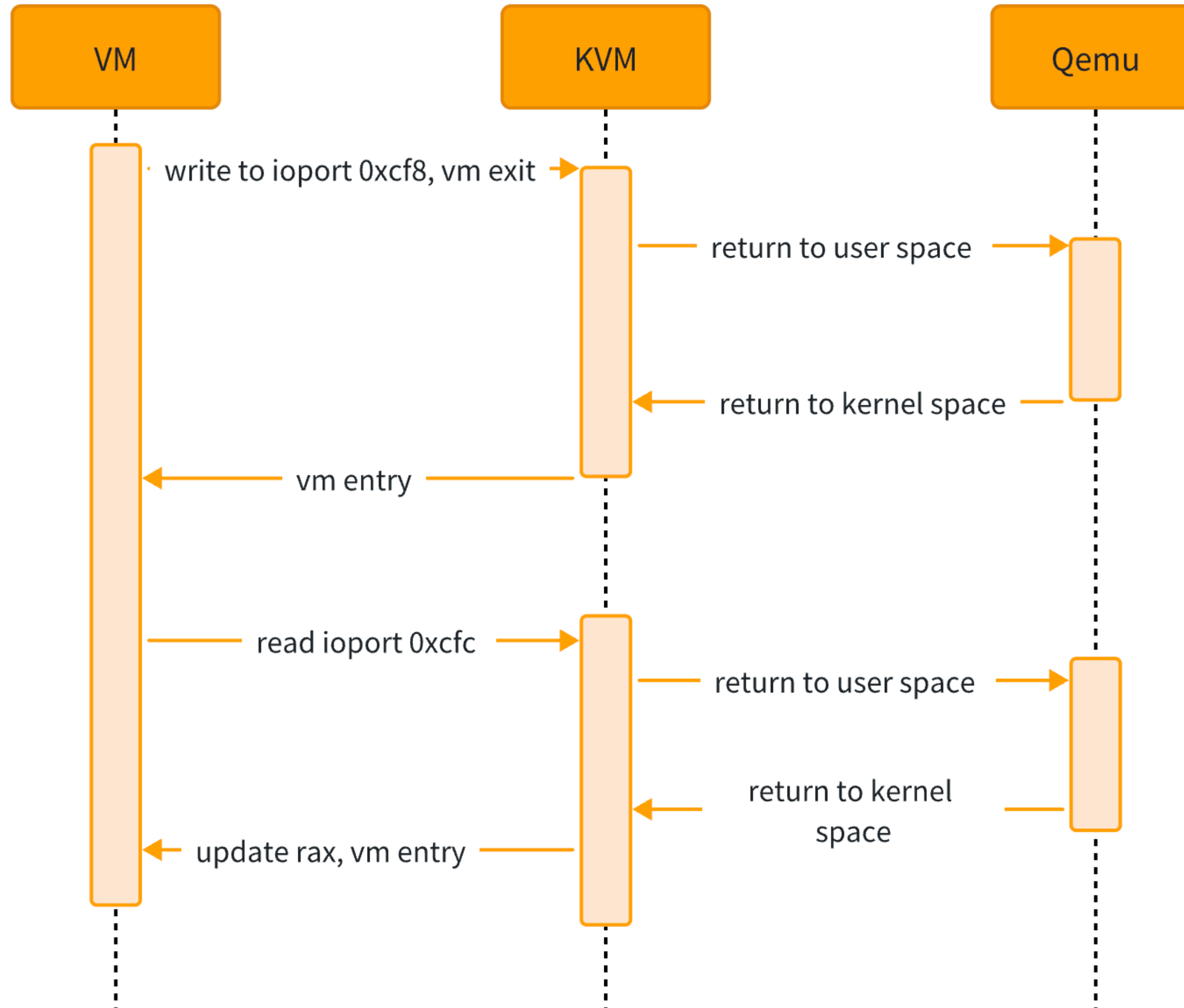
Context switch cost same as the Host

Context switch cost because of virtualization

Design & Implementation — PCI PortIO

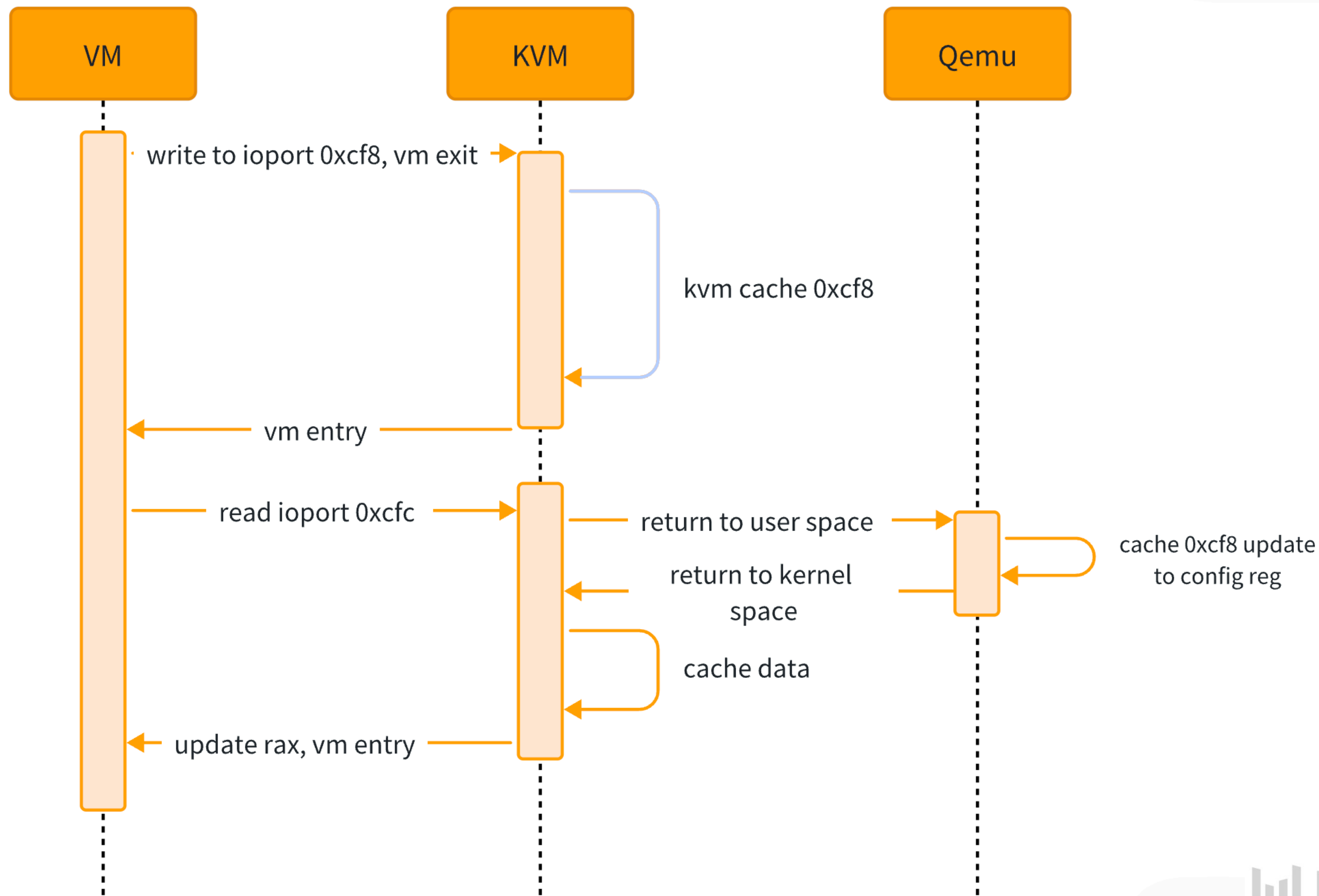
PCI PortIO

Original read path



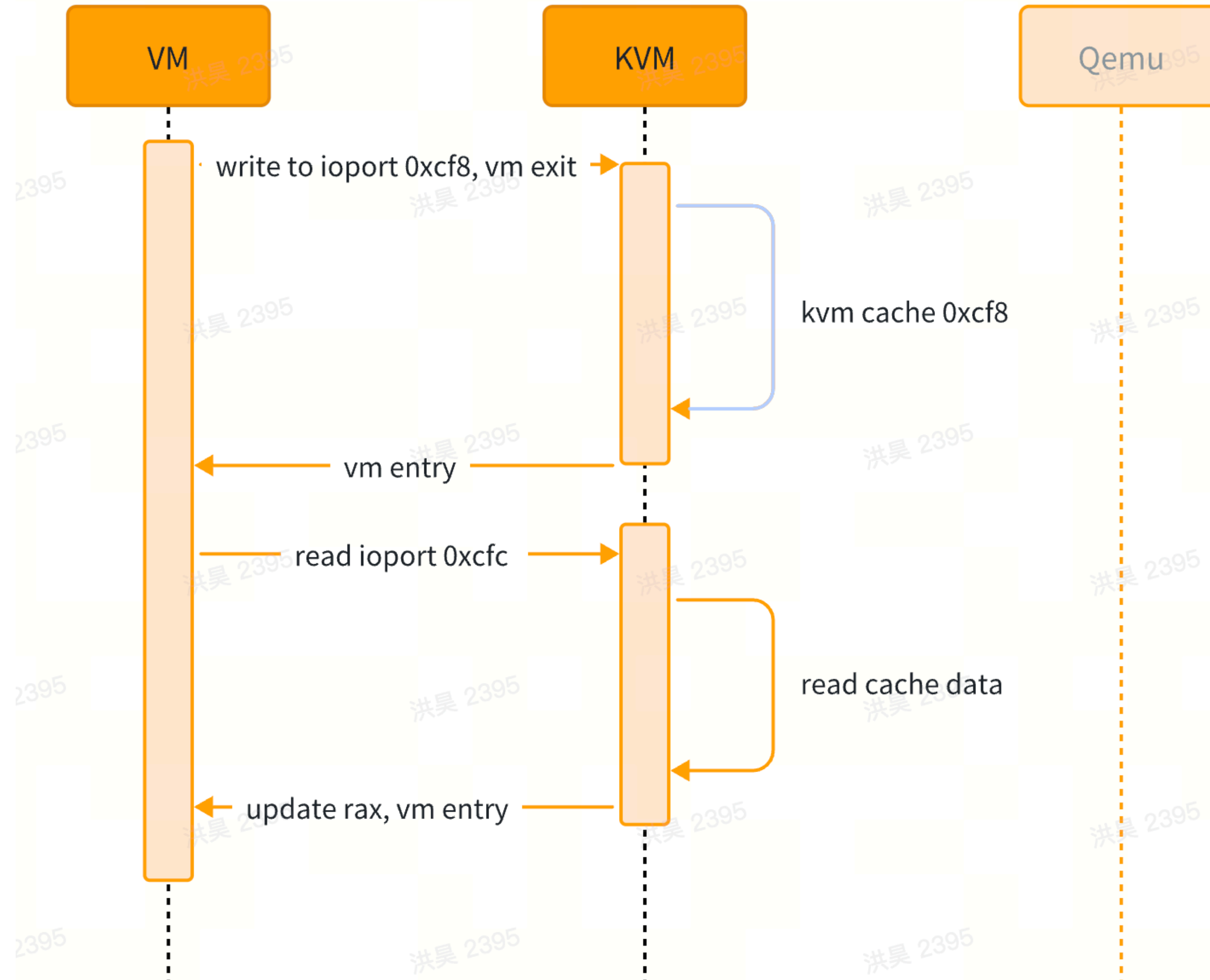
PCI PortIO

First read path



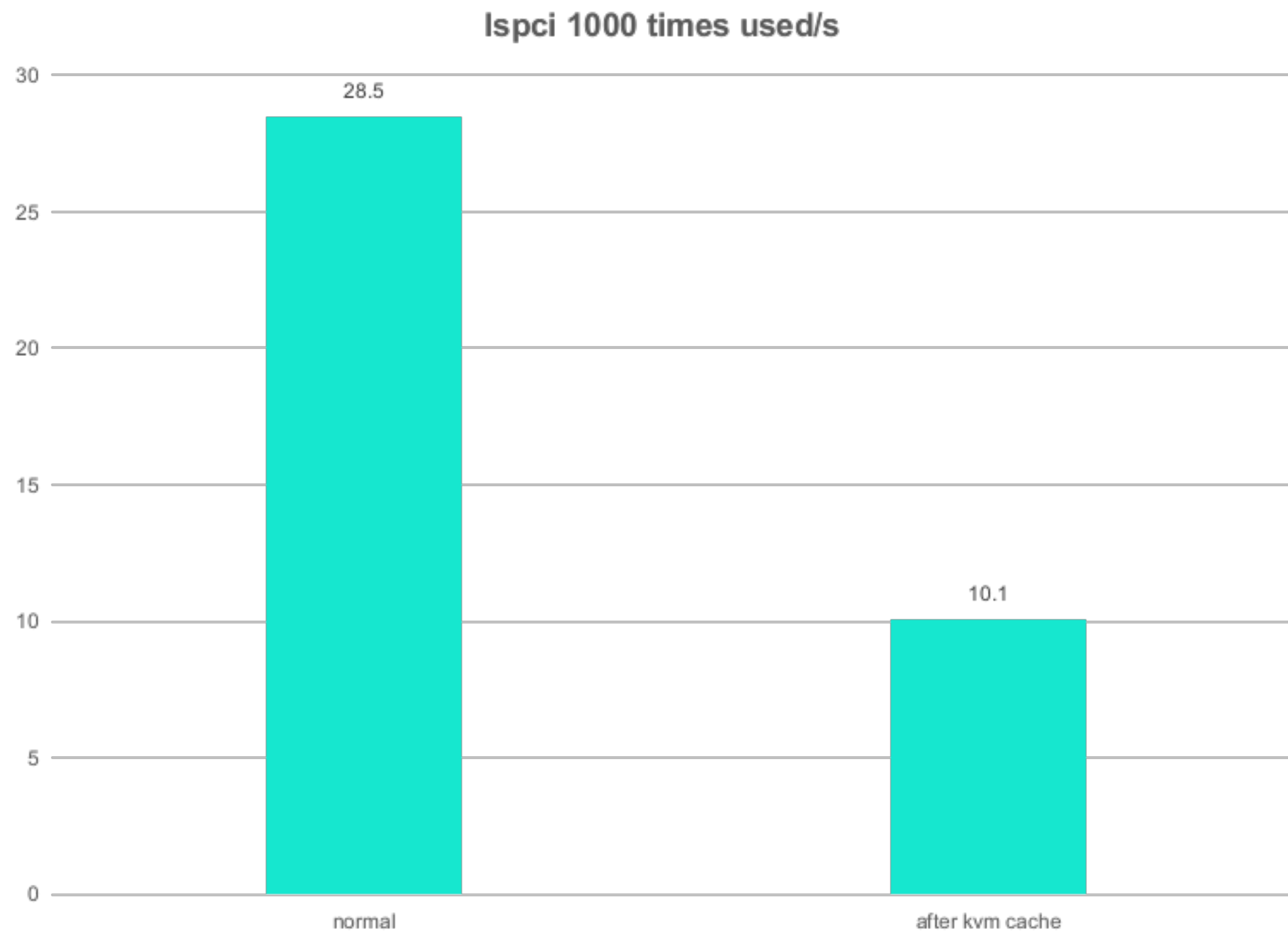
PCI PortIO

Read cache path



PCI PortIO

Performance comparison

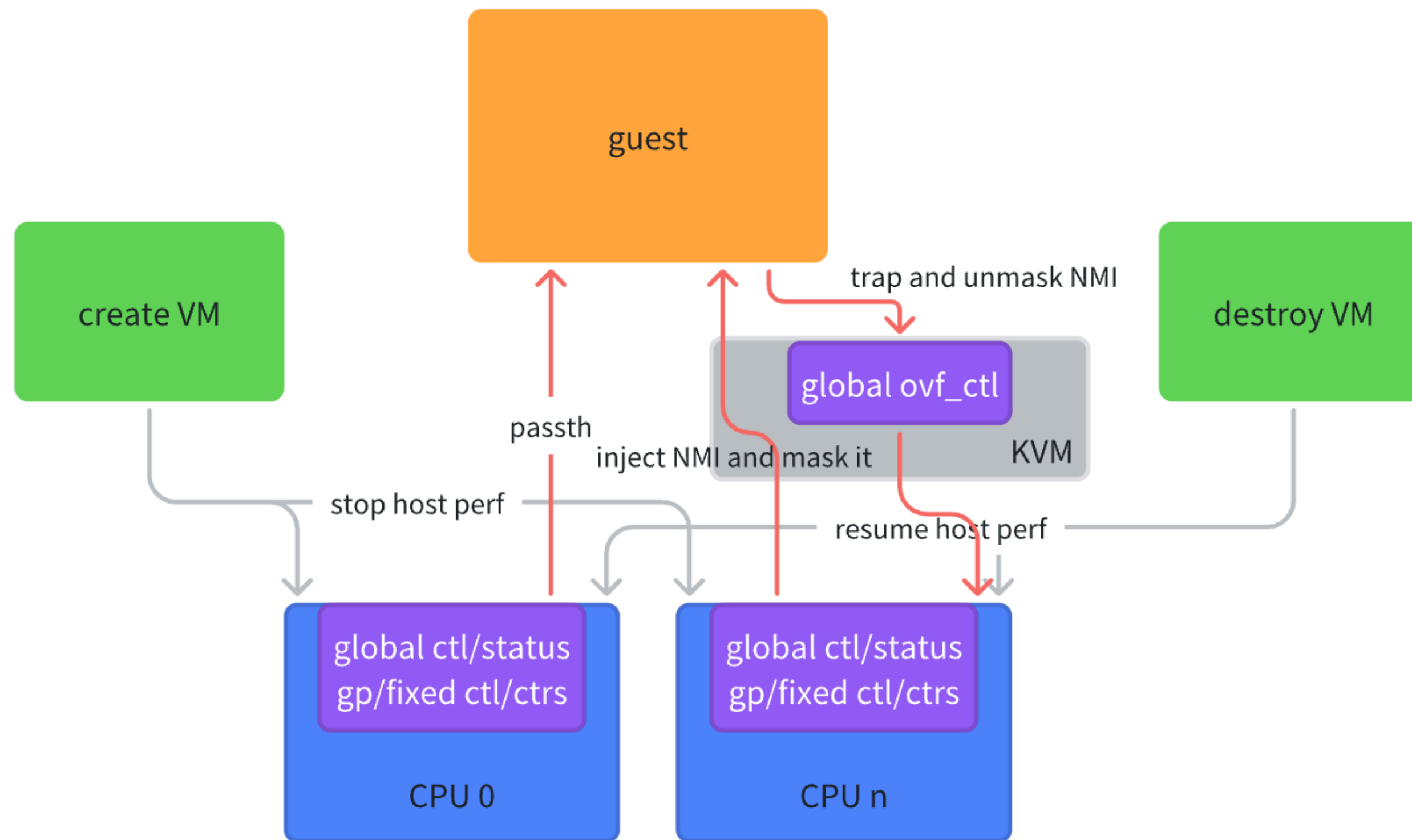


The implementation of `lspci` is consistent with the high-precision monitoring agents. They are both read pci config space.

Design & Implementation — Passthrough PMU

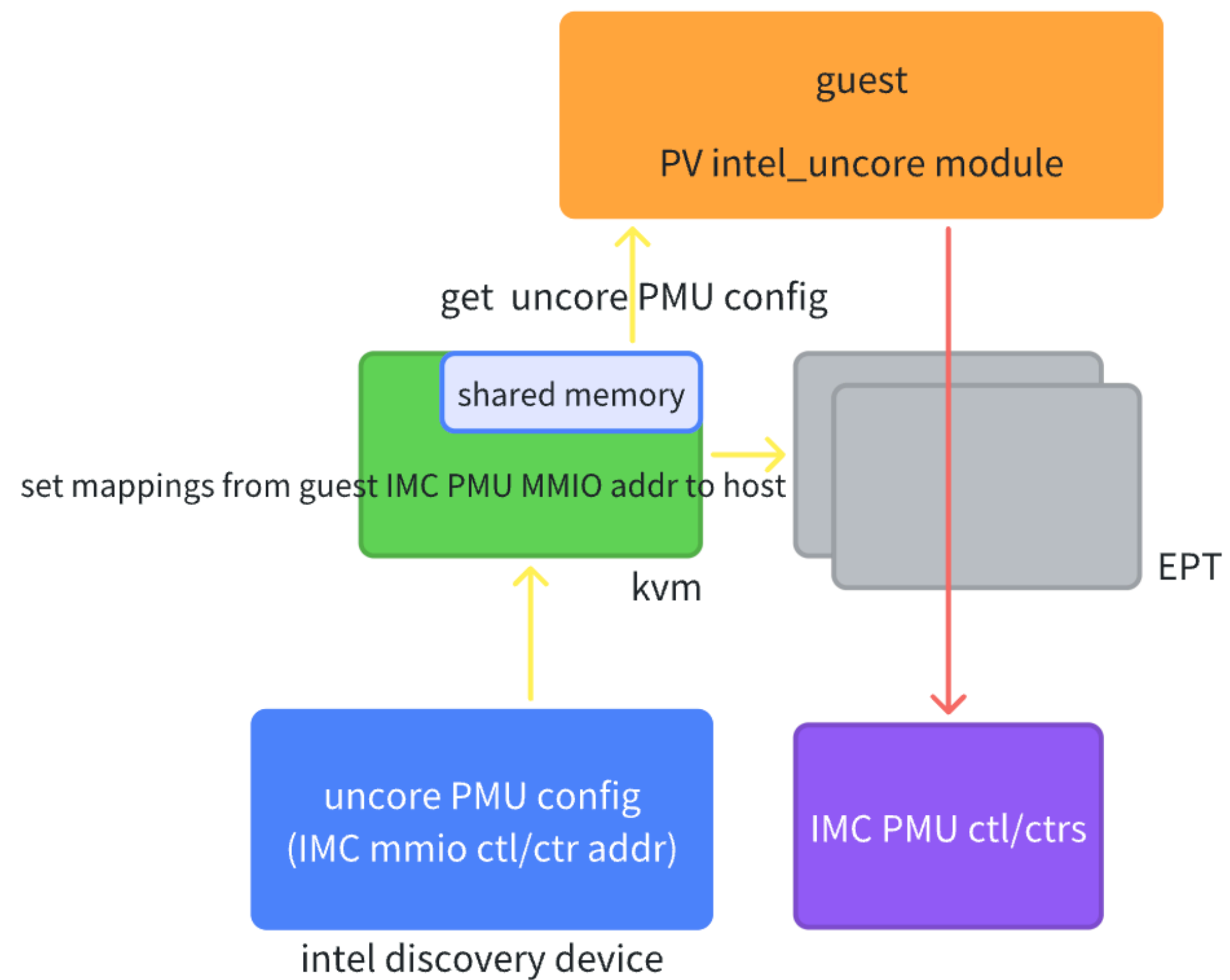
Passthrough Core PMU

eliminate core PMU MSR VM-Exits overhead



Uncore PMU Virtualization

need to deploy high-precision monitoring agents in guest, take IMC PMU for example



Future Work



Future Work

- upstream
- solve H2D & D2H virtualization degradation by DMA devirt

Thank You

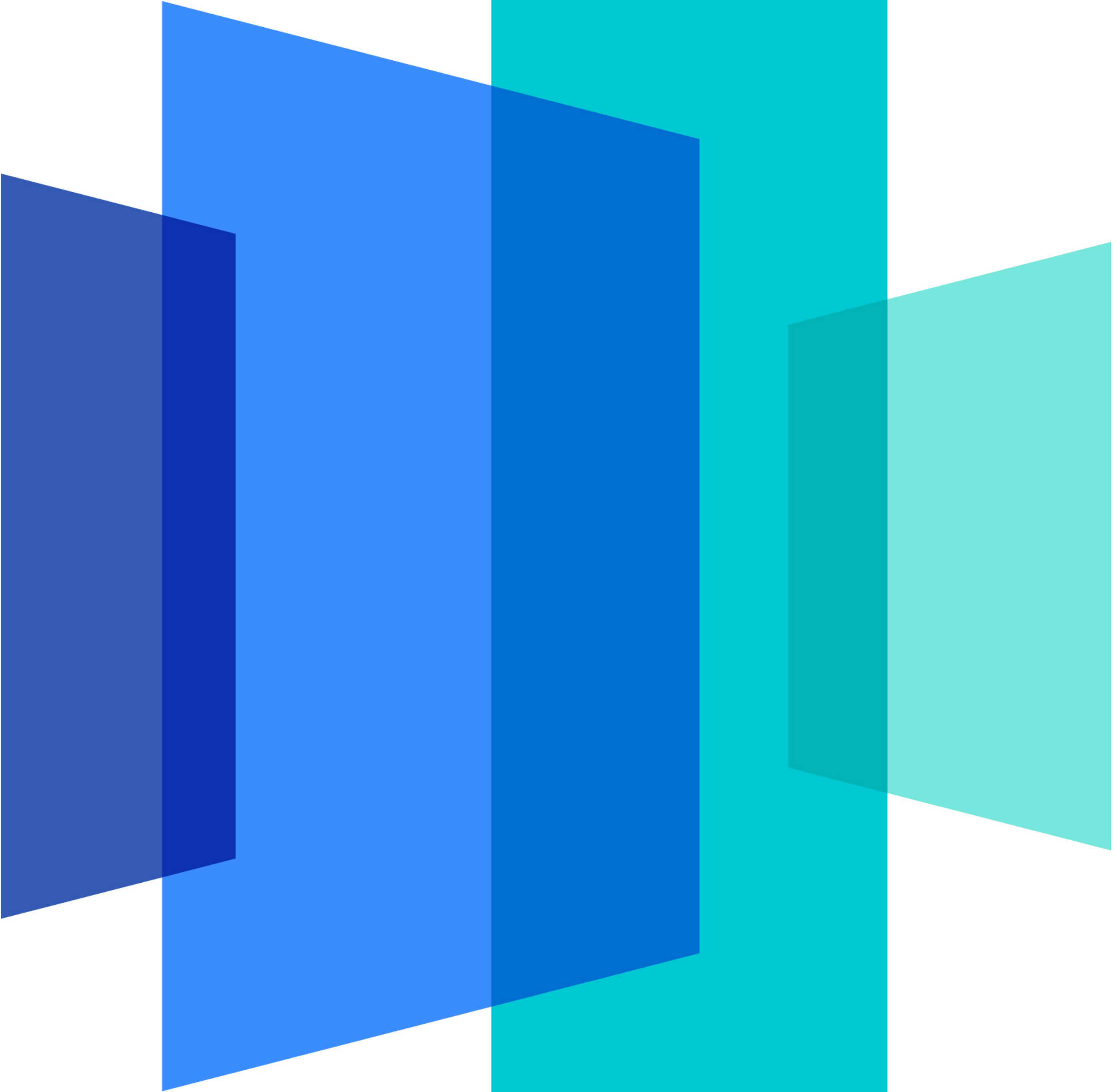
Q & A

Contact Info: hexin.op@bytedance.com
honghao.dante@bytedance.com



Can't overcommit

- Core PMU: Host cannot use perf when VM running
- Uncore PMU: Can't used in multiple VM



 ByteDance

The ByteDance logo consists of a stylized icon on the left, which is a series of four vertical bars of increasing height from left to right, colored in shades of blue and teal. To the right of the icon, the word "ByteDance" is written in a bold, sans-serif font in a dark blue color.