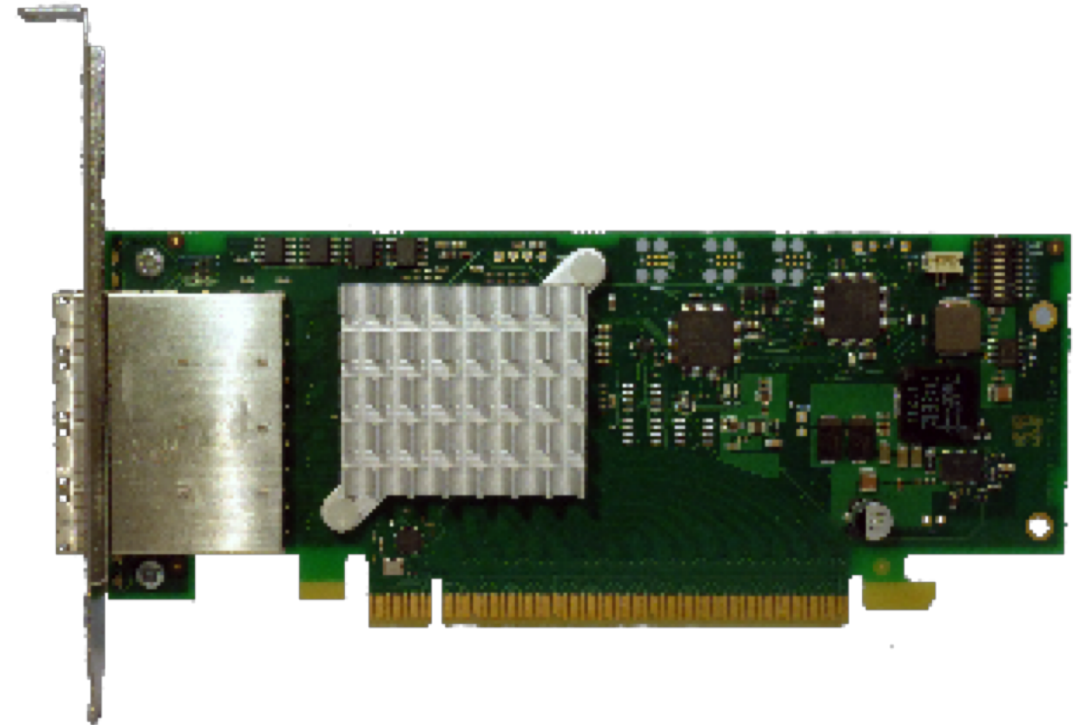# S9709
# *Dynamic Sharing of GPUs and IO in a PCIe Network*

**Håkon Kvale Stensland**

Senior Research Scientist / Associate Professor
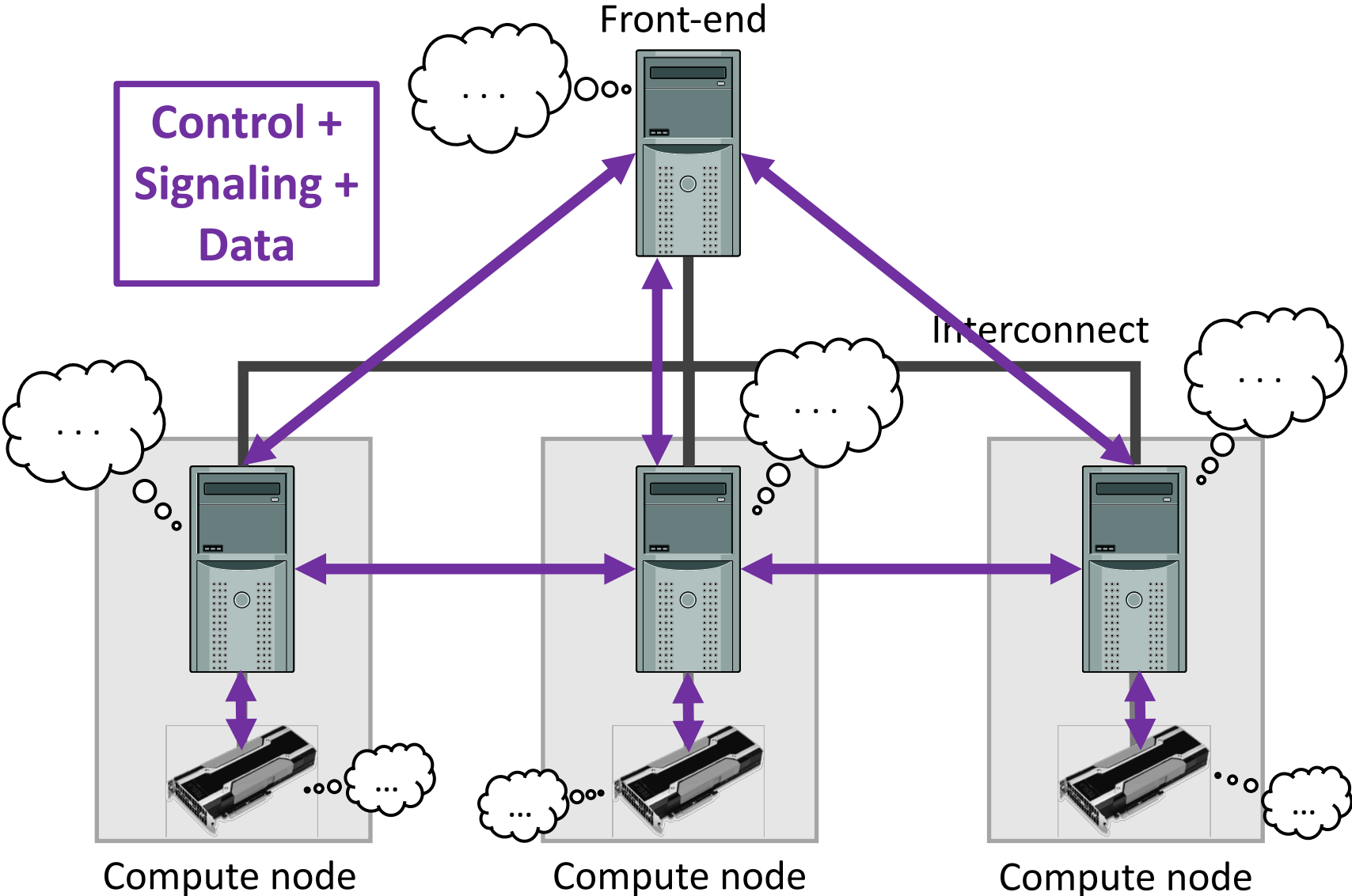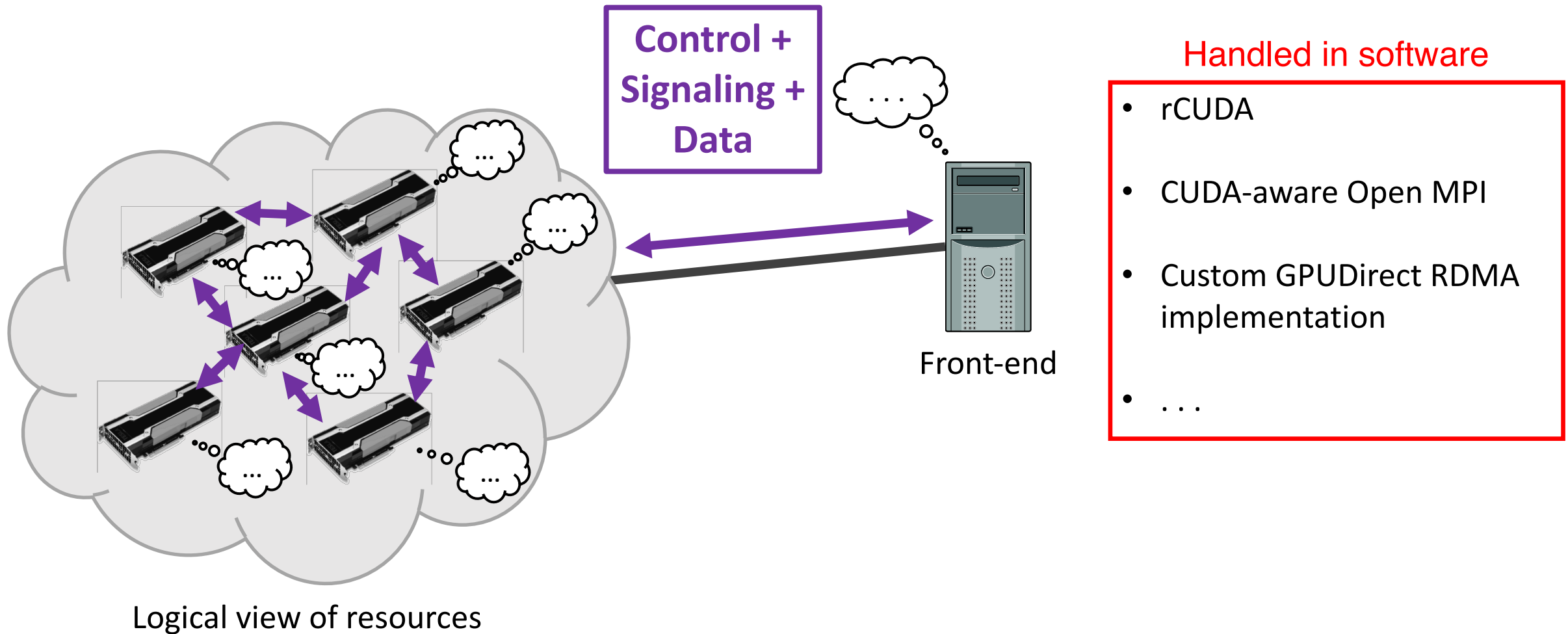Simula Research Laboratory / University of Oslo

# Outline

- Motivation

- PCIe Overview

- Non-Transparent Bridges

- Dolphin SmartIO
  - Example Application
  - NVMe sharing
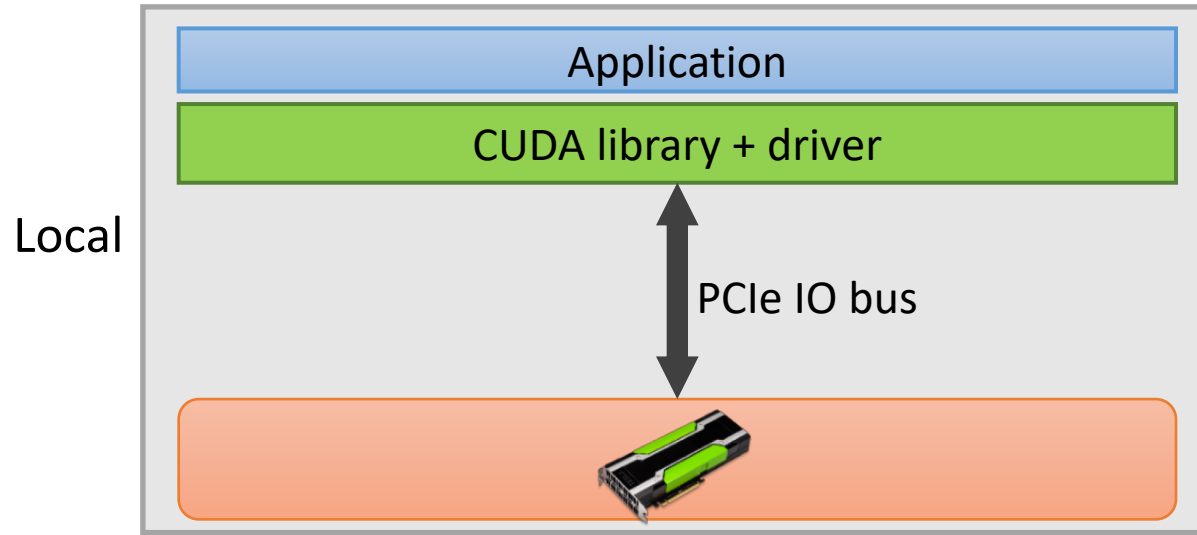
- SmartIO in Virtual Machines

# Distributed applications may need to access and use IO resources that are physically located inside remote hosts
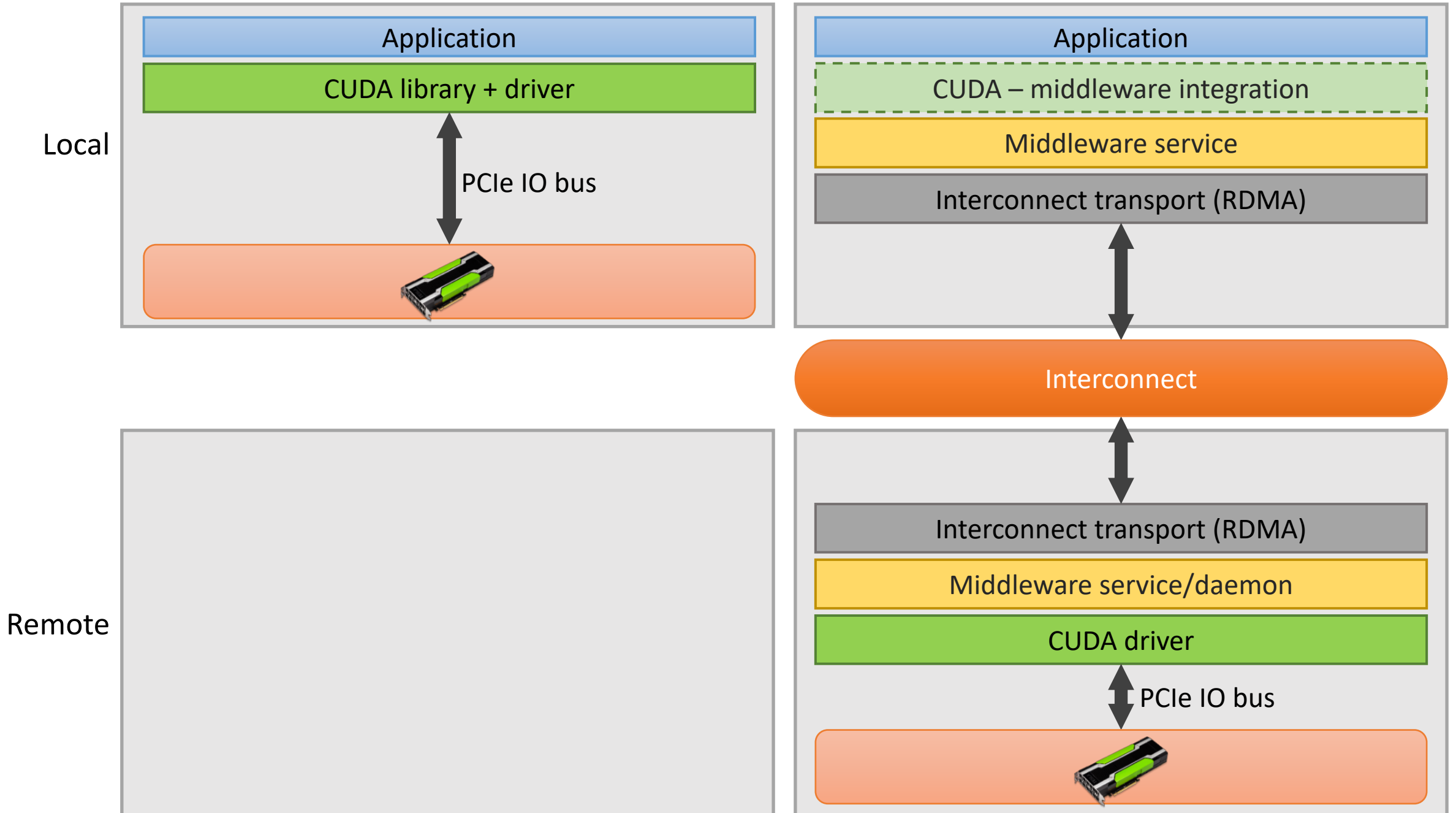
Control + Signaling + Data

Front-end

Interconnect

Compute node

Compute node

Compute node

# Software abstractions simplify the use and allocation of resources in a cluster and facilitate development of distributed applications
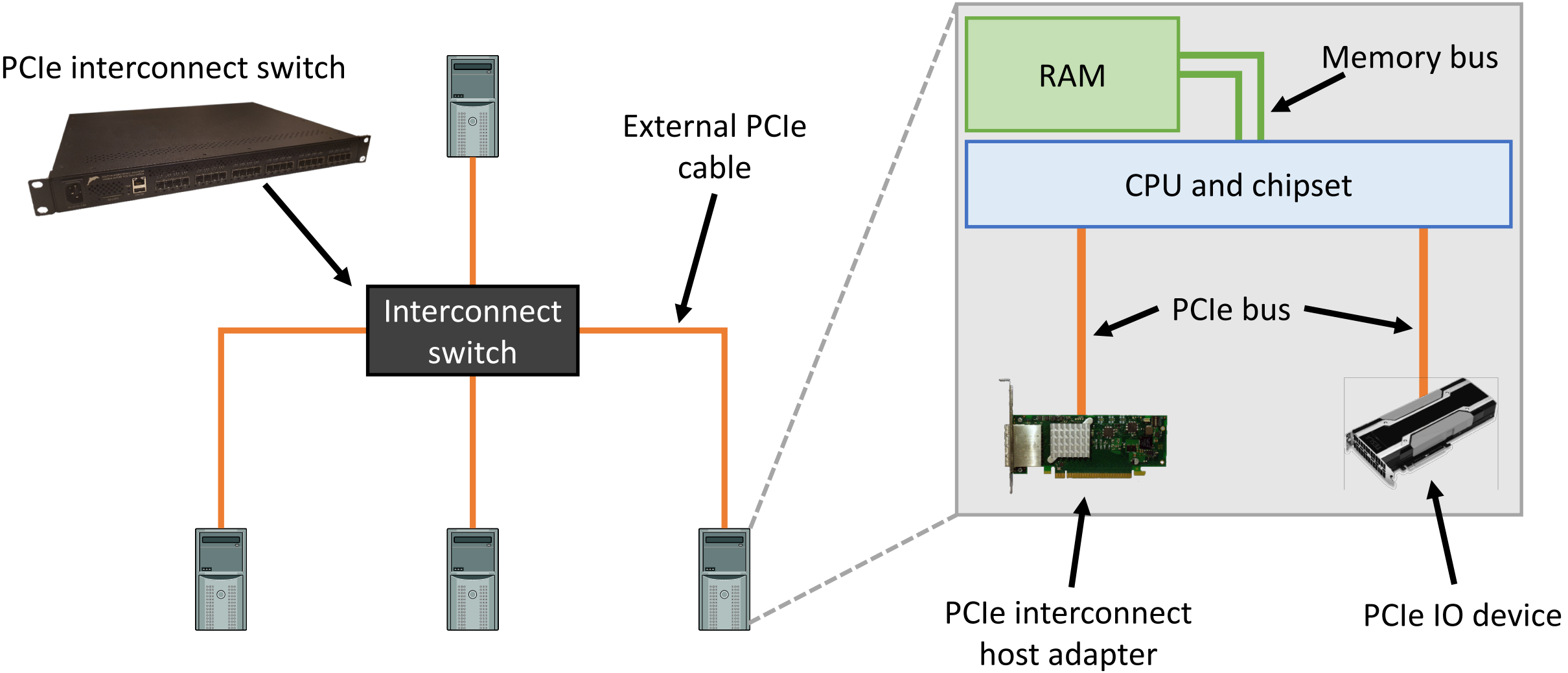


**Control + Signaling + Data**
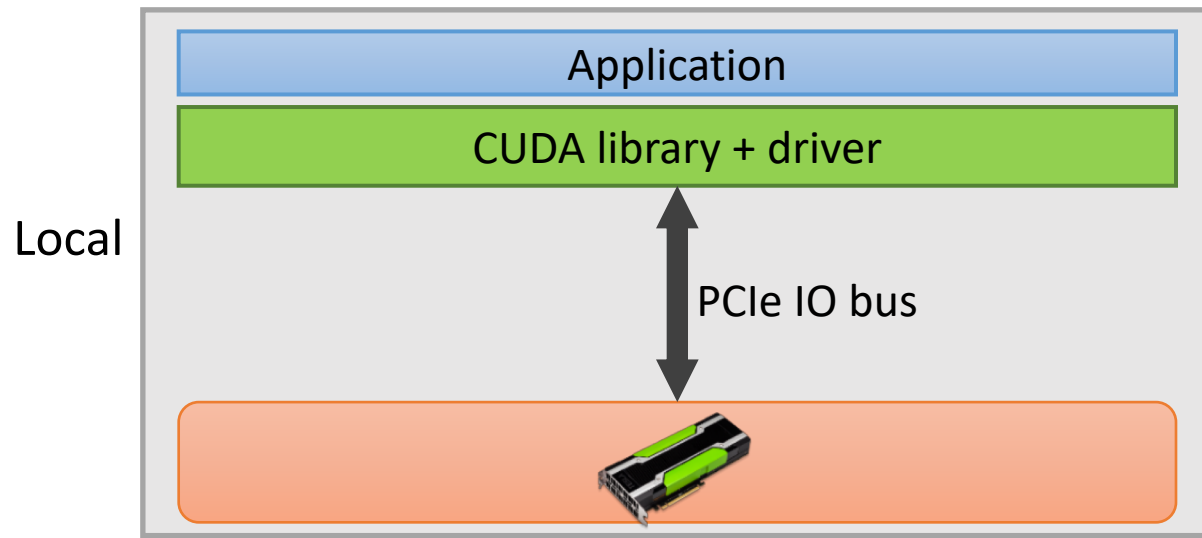
Front-end

Logical view of resources

Handled in software

- rCUDA

- CUDA-aware Open MPI

- Custom GPUDirect RDMA implementation

- . . .

## Local resource

### Local

| Application |
| CUDA library + driver |

PCIe IO bus

## Remote resource using **middleware**

| Application |
| CUDA – middleware integration |
| Middleware service |
| Interconnect transport (RDMA) |

Interconnect

### Remote

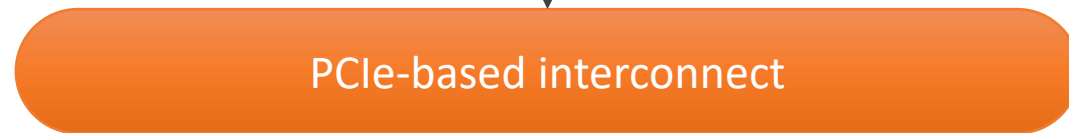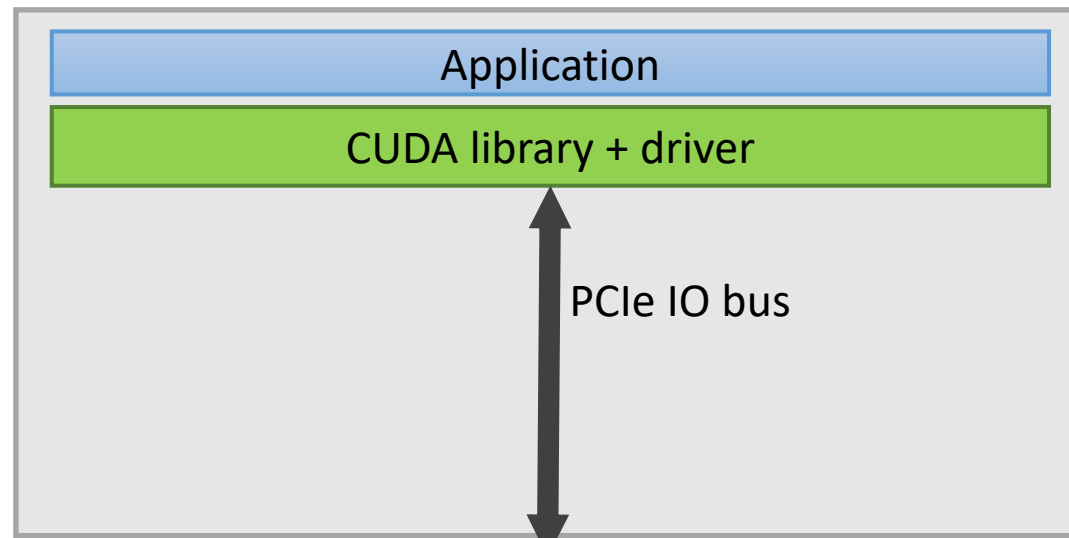| Interconnect transport (RDMA) |
| Middleware service/daemon |
| CUDA driver |

PCIe IO bus

# In PCIe clusters, the same fabric is used both as local IO bus within a single node and as the interconnect between separate nodes
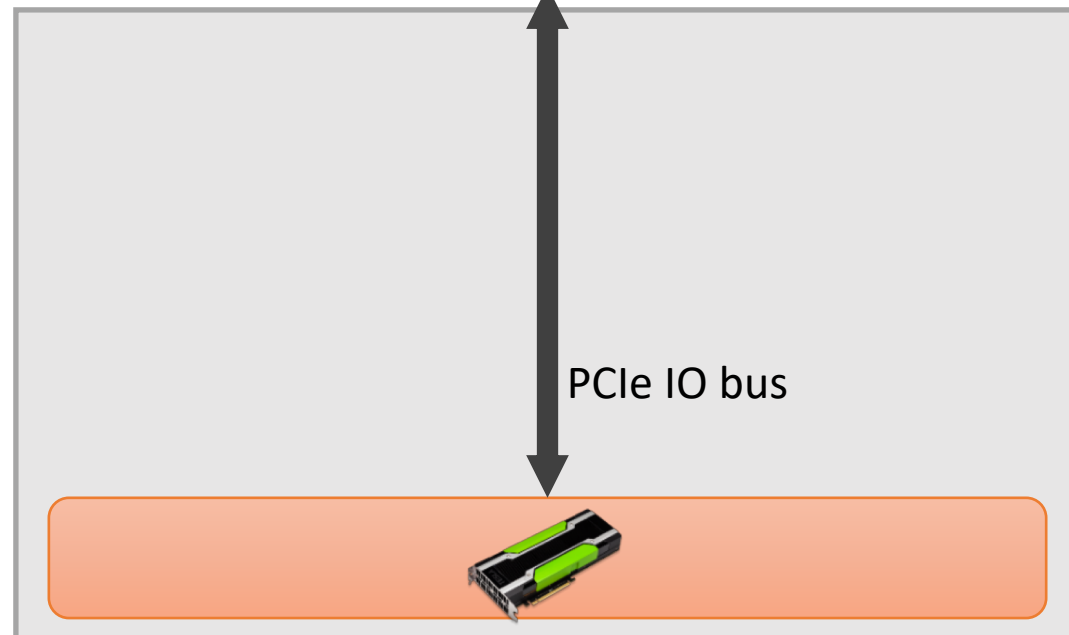


PCIe interconnect switch

Interconnect switch

External PCIe cable

RAM

Memory bus

CPU and chipset

PCIe bus

PCIe interconnect host adapter

PCIe IO device

Local resource

Remote resource over **native fabric**

Local

Remote

Application

CUDA library + driver

PCIe IO bus

Application

CUDA library + driver

PCIe IO bus

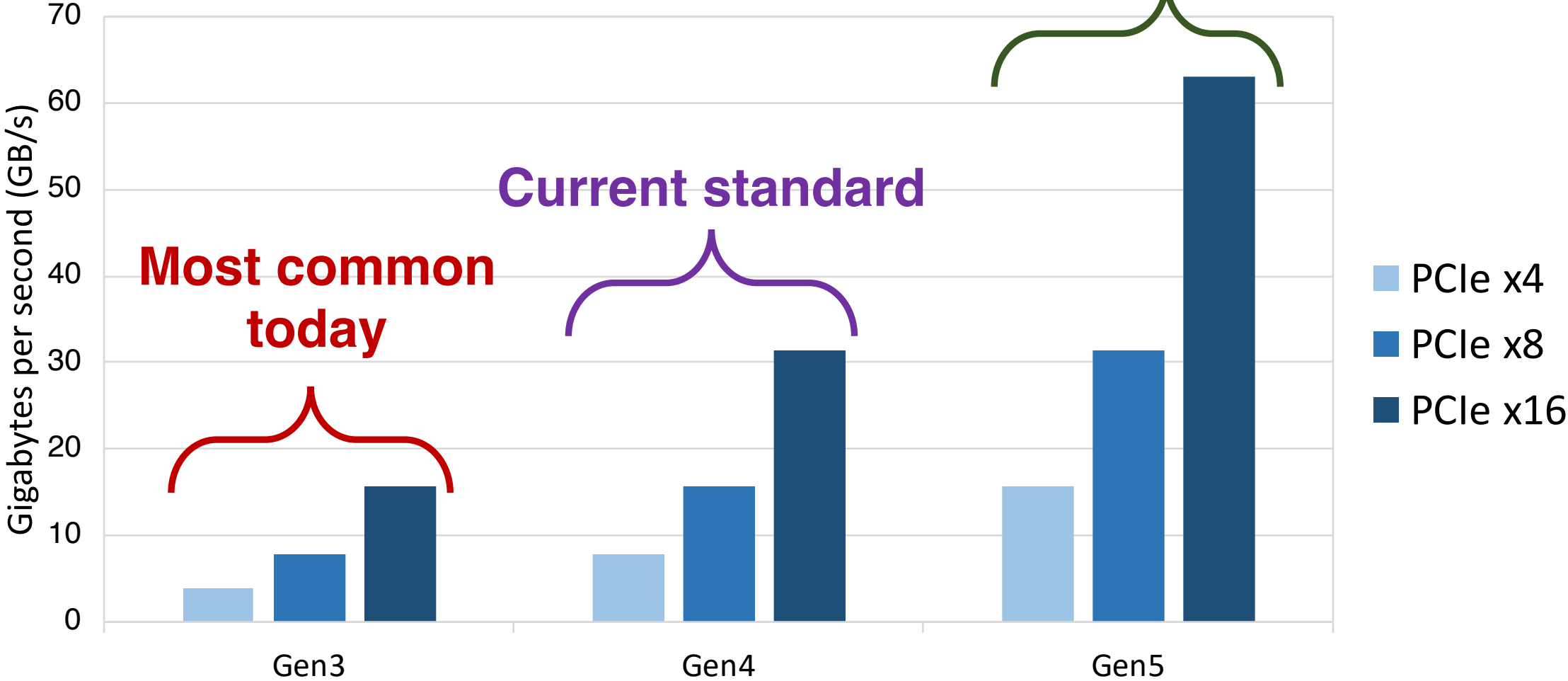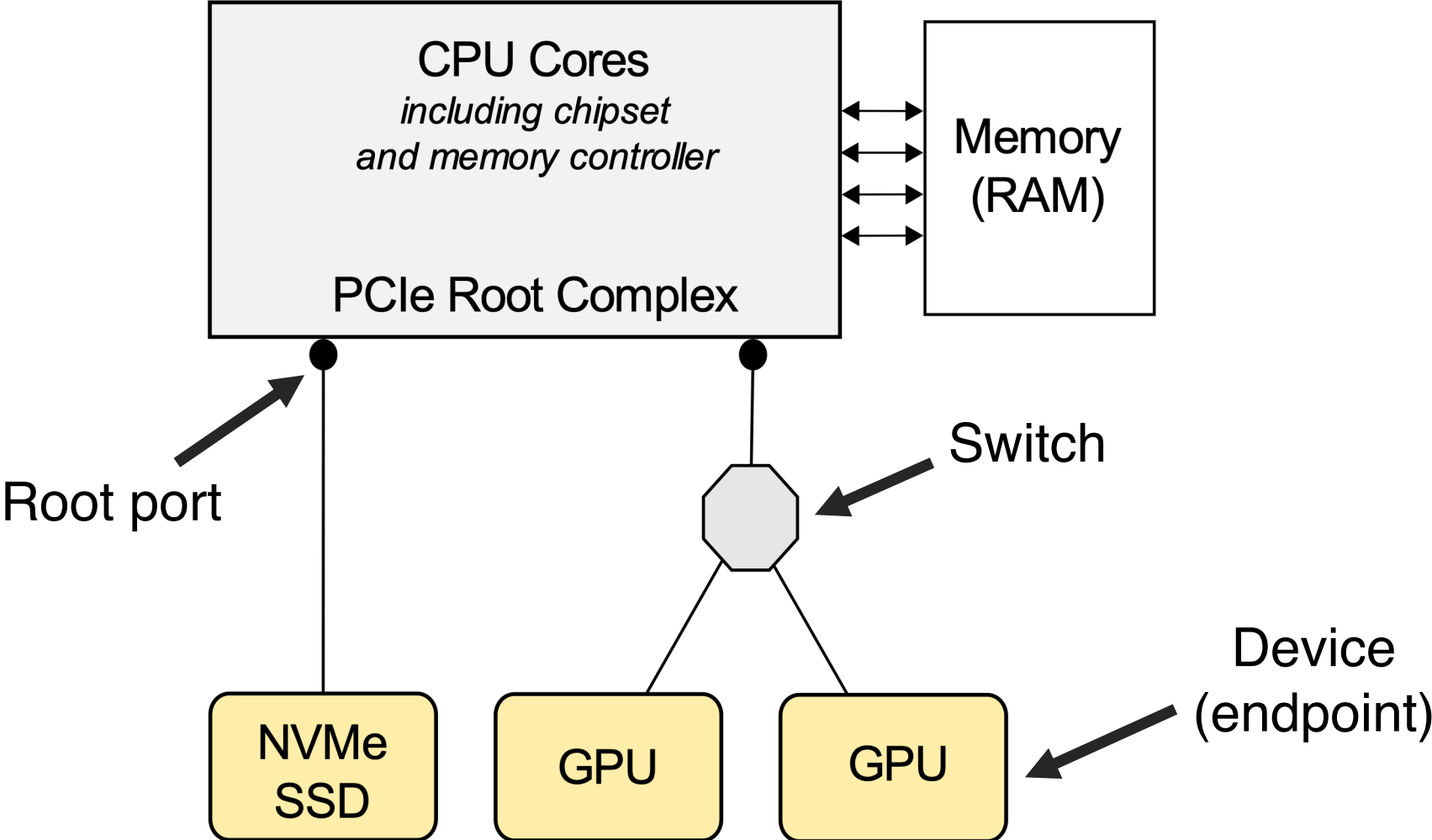PCIe-based interconnect

PCIe IO bus

# PCIe Overview

# PCI Express (PCIe) is the most widely adopted I/O interconnection technology used in computer systems today



Bar chart titled with y-axis "Gigabytes per second (GB/s)" ranging from 0 to 70, showing three generations (Gen3, Gen4, Gen5) each with three bars: PCIe x4, PCIe x8, and PCIe x16.

- **Most common today** — labeled over Gen3
- **Current standard** — labeled over Gen4
- **Near future-ish** — labeled over Gen5

Legend:
- PCIe x4
- PCIe x8
- PCIe x16

# The PCIe fabric is structured as a tree, where devices form the leaf nodes (endpoints) and the CPU is on top of the root



CPU Cores
*including chipset and memory controller*

PCIe Root Complex

Memory (RAM)
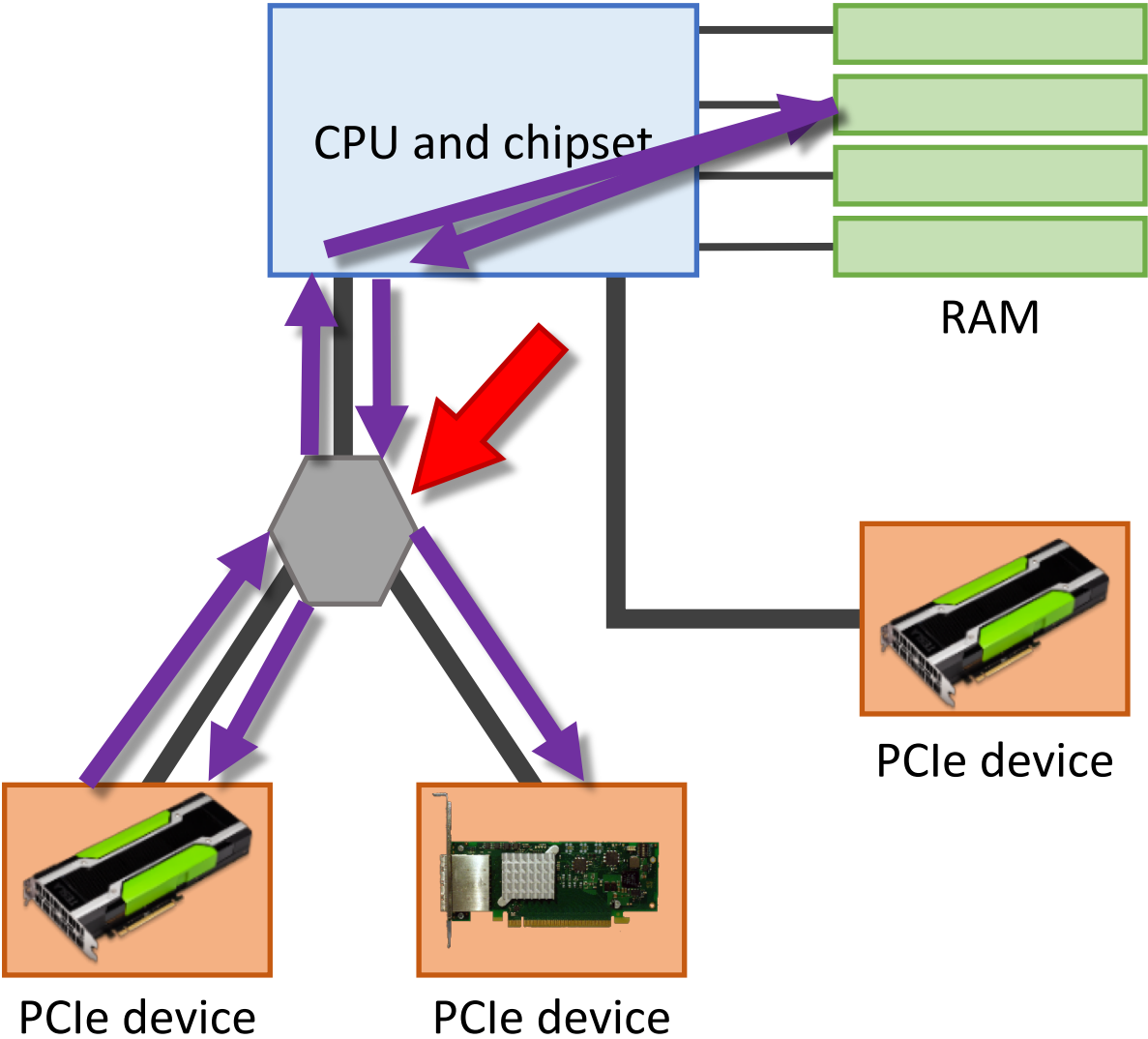
Root port

Switch

Device (endpoint)

NVMe SSD

GPU

GPU

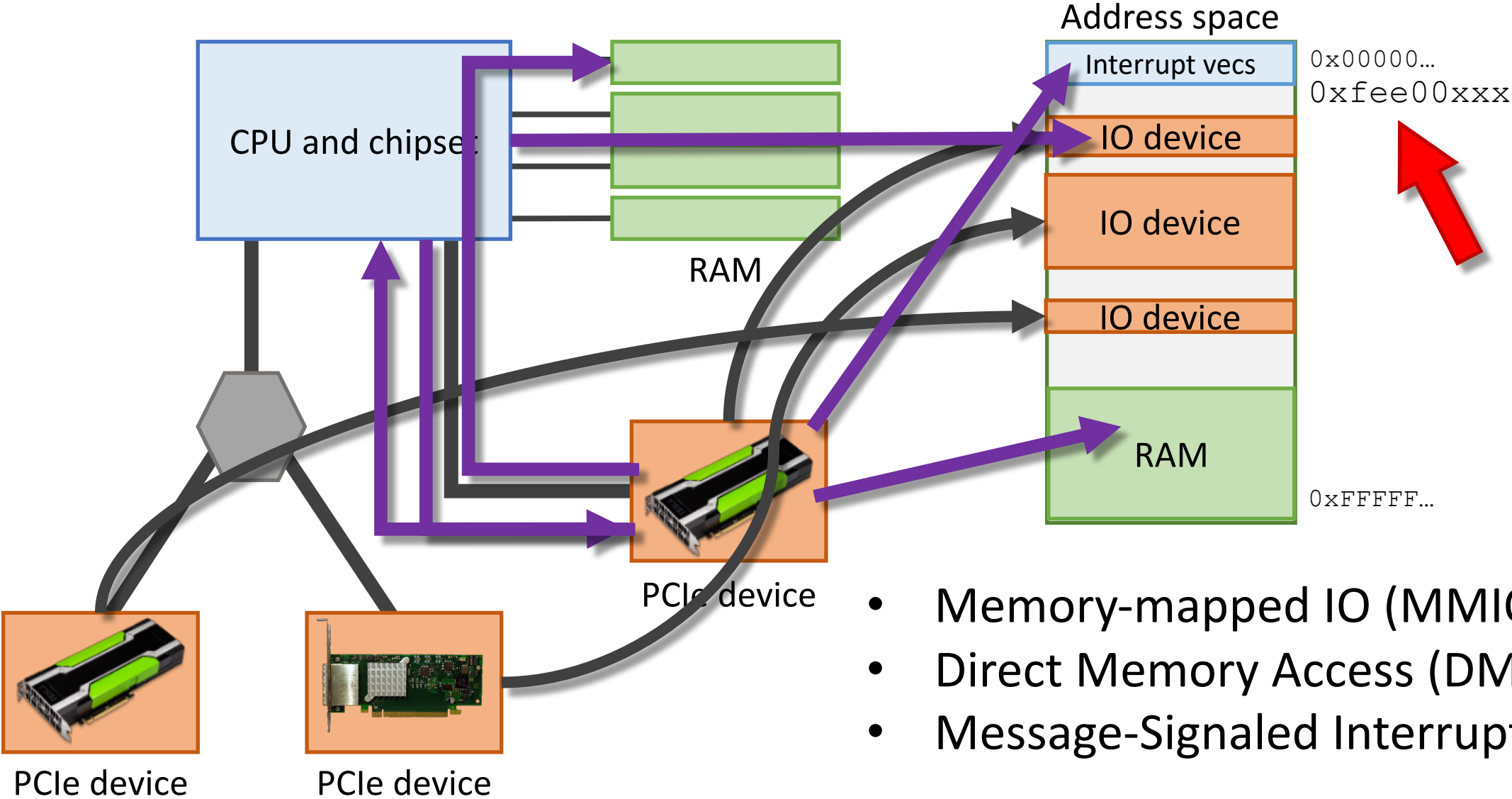# The PCIe fabric is structured as a tree, where devices form the leaf nodes (endpoints) and the CPU is on top of the root

# Memory reads and writes are handled by PCIe as transactions that are packet-switched through the fabric depending on the address



CPU and chipset

RAM

PCIe device

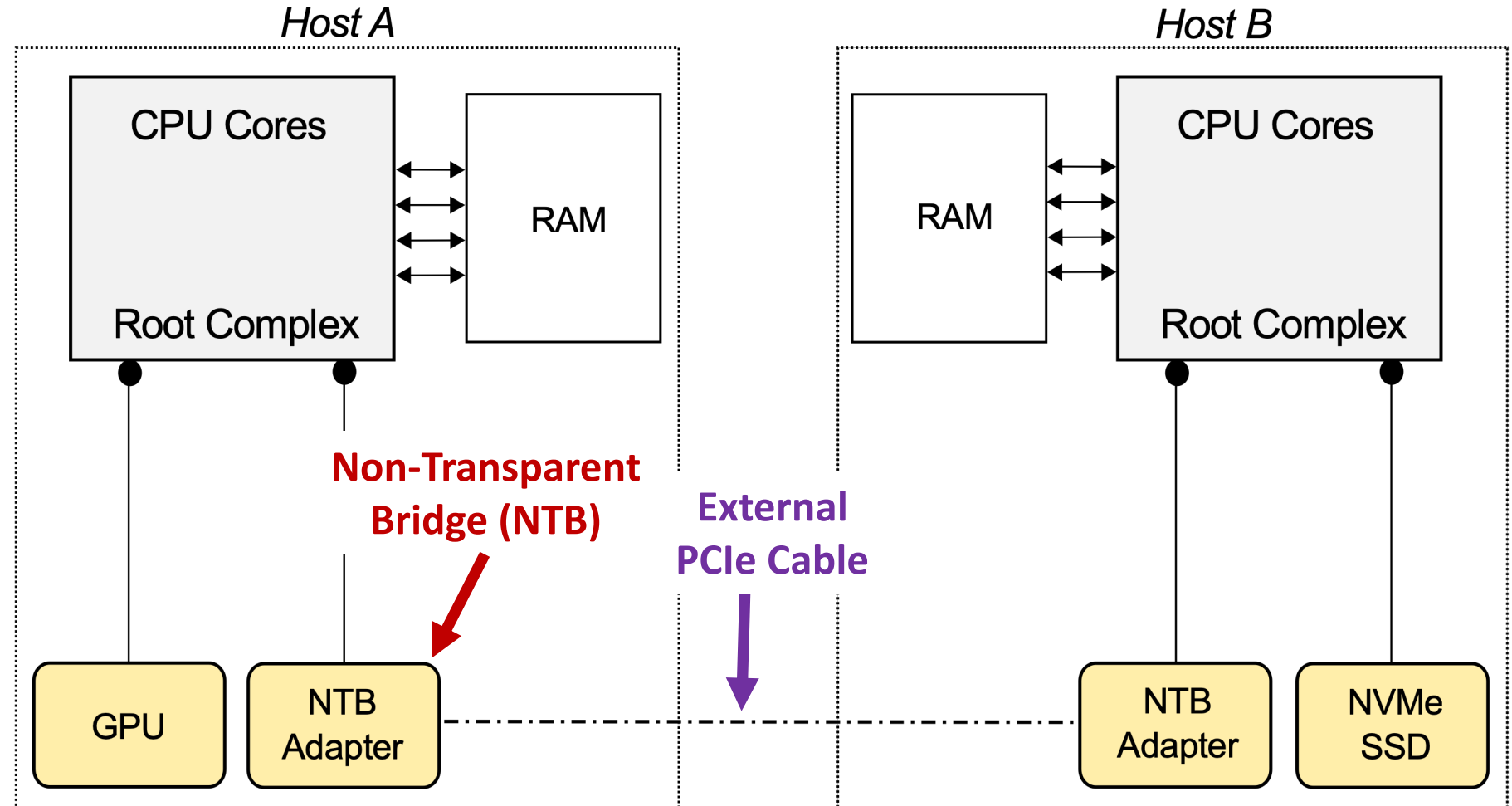PCIe device

PCIe device

- Upstream
- Downstream
- Peer-to-peer (shortest path)

# IO devices and the CPU share the same physical address space, allowing devices to access system memory and other devices



Address space

Interrupt vecs — 0x00000…
0xfee00xxx

IO device

IO device

IO device

RAM — 0xFFFFF…

CPU and chipset

RAM

PCIe device
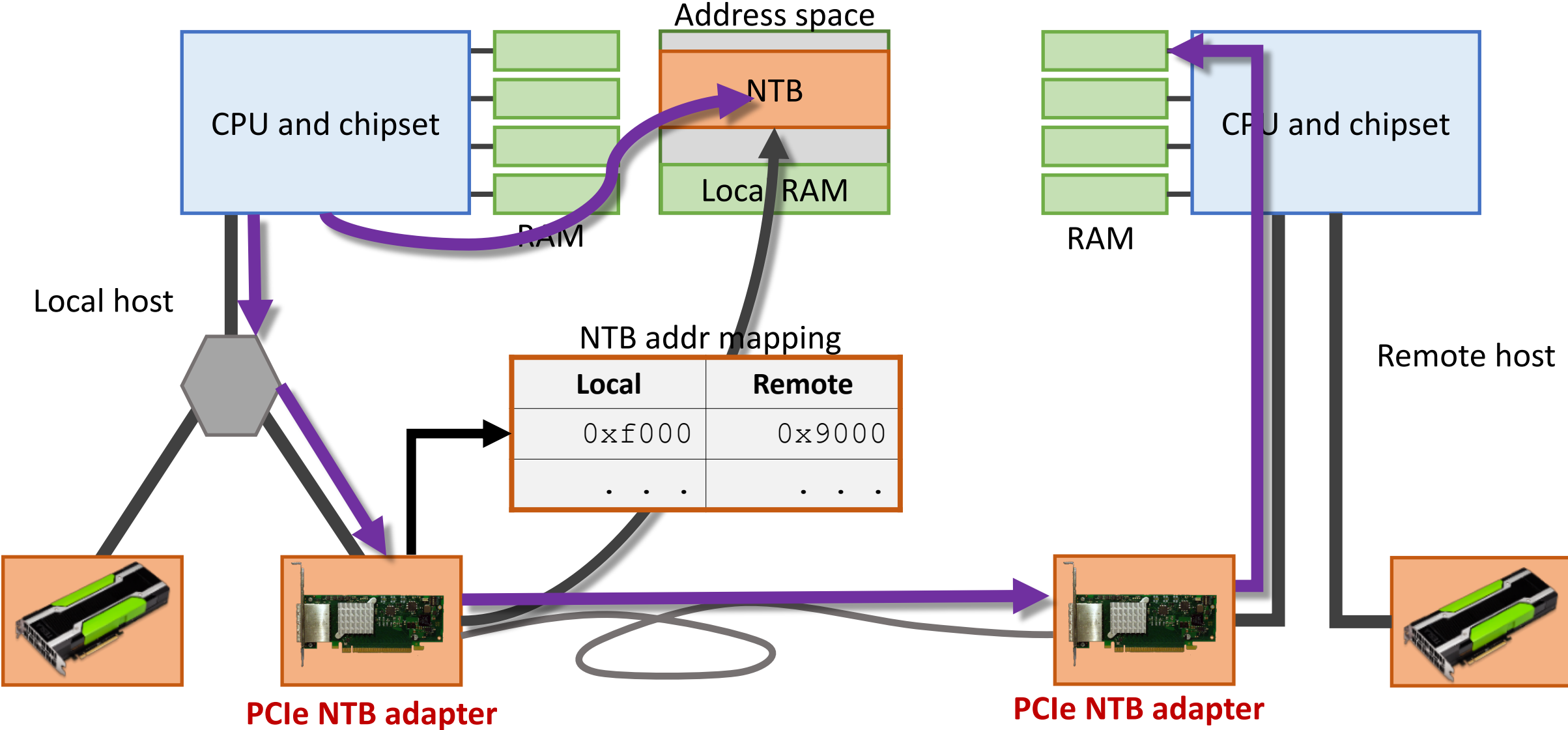
PCIe device

PCIe device

PCIe device

- Memory-mapped IO (MMIO / PIO)
- Direct Memory Access (DMA)
- Message-Signaled Interrupts (MSI-X)

# Non-Transparent Bridges

# We can interconnect separate PCIe root complexes and translate addresses between them using a non-transparent bridge (NTB)
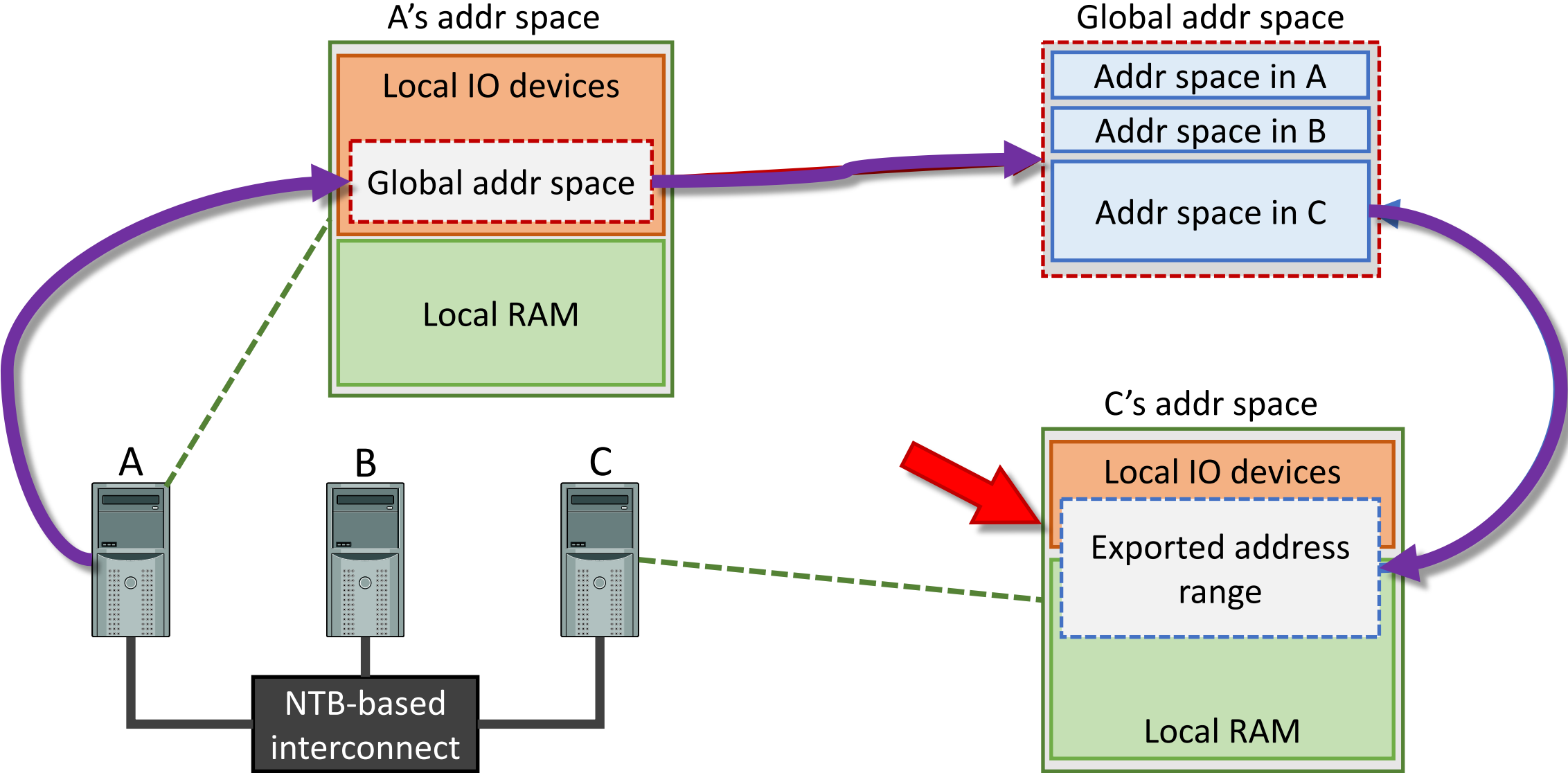
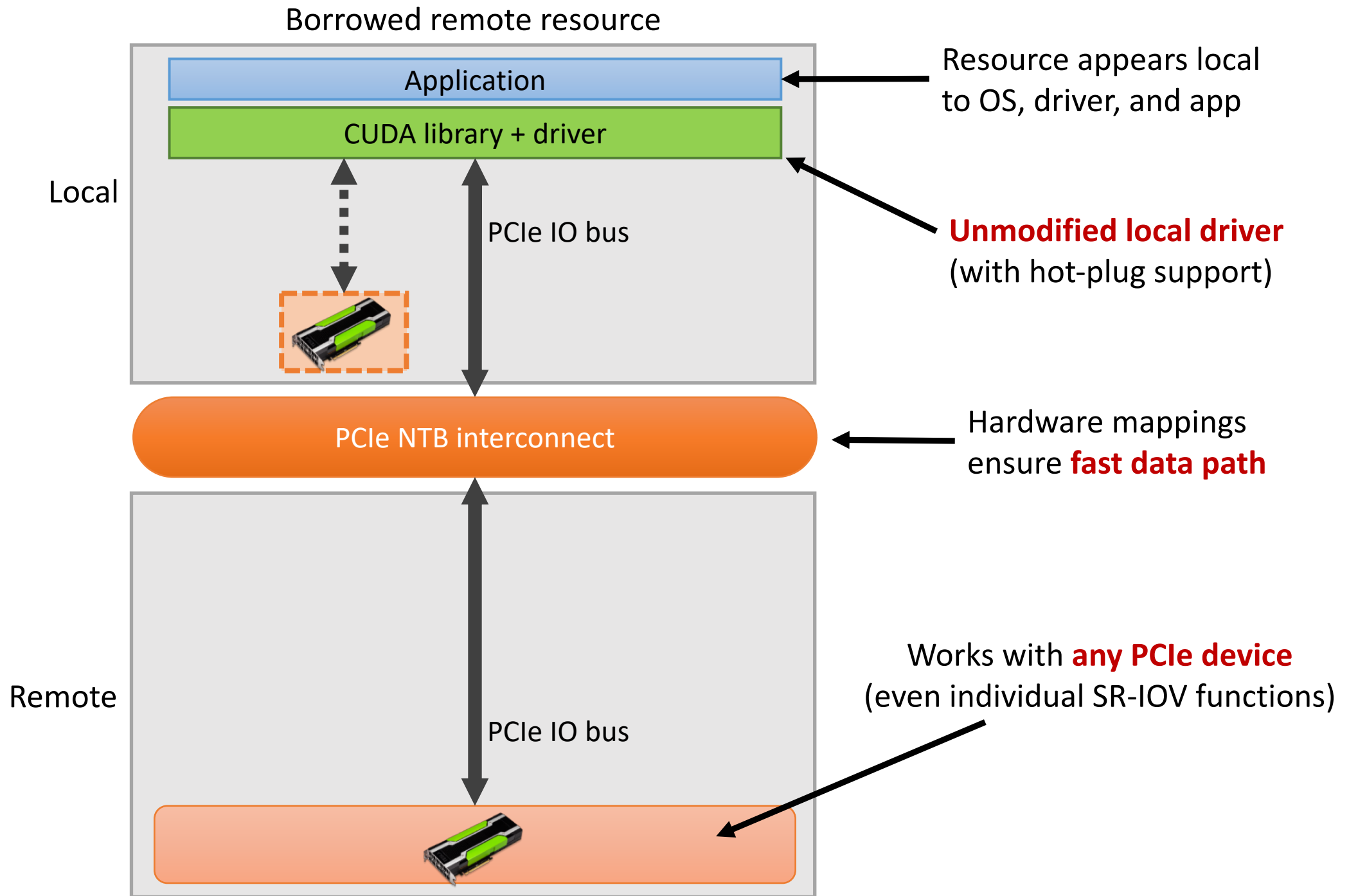# Remote address space can be mapped into local address space by using PCIe Non-Transparent Bridges (NTBs)
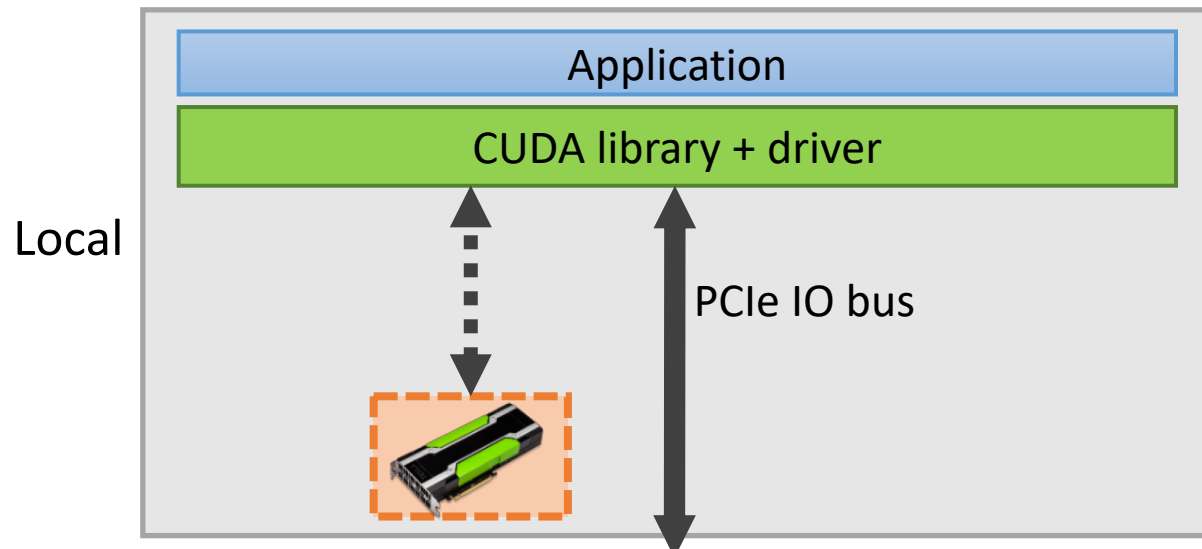
# Using NTBs, each node in the cluster take part in a shared address space and have their own "window" into the global address space
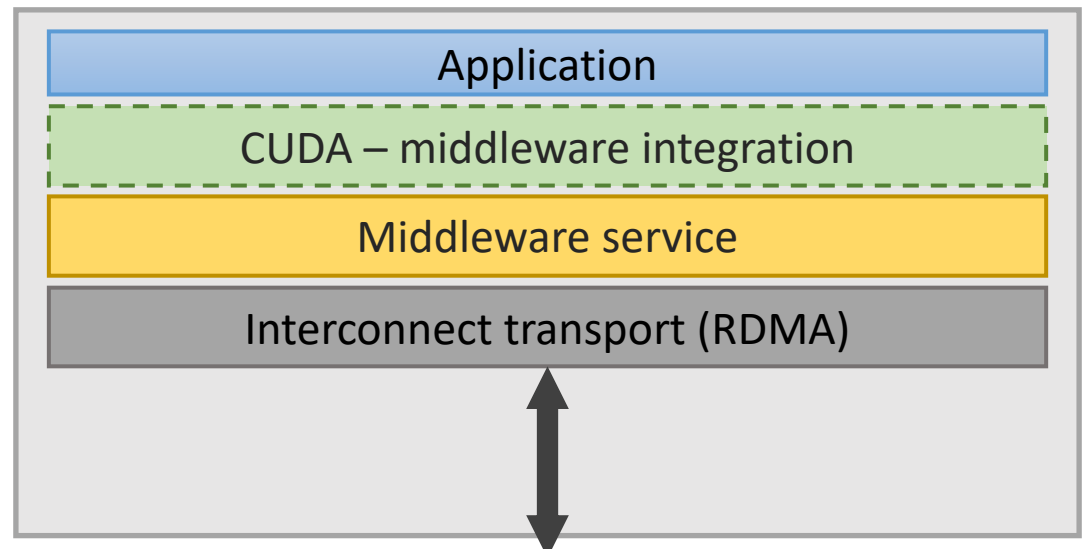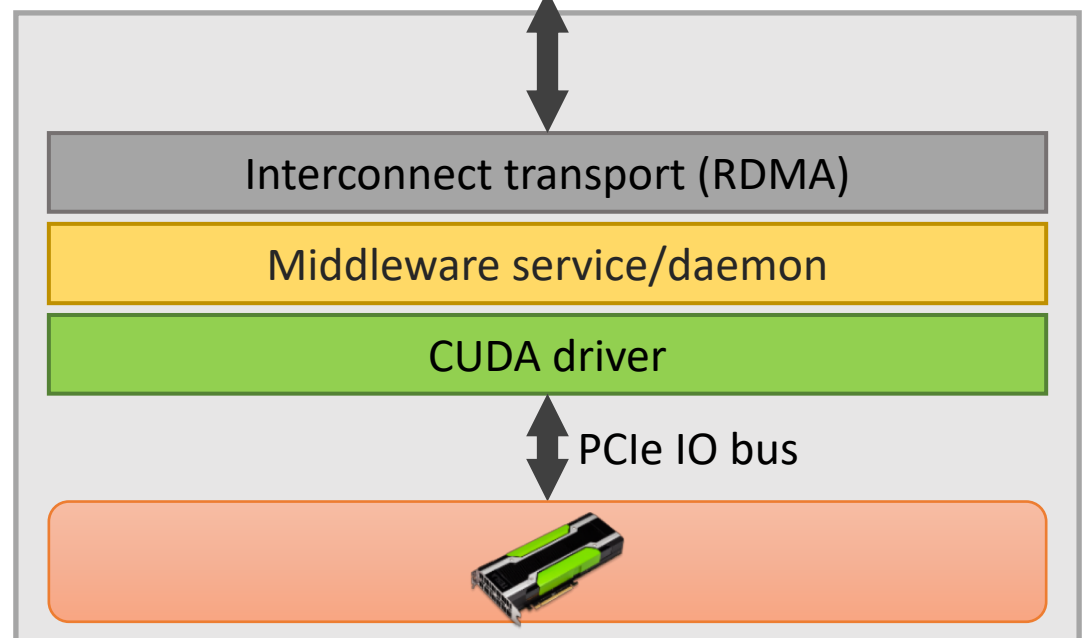
# SmartIO

# Borrowed remote resource

## Local

| Application | ← Resource appears local to OS, driver, and app |

| CUDA library + driver | ← **Unmodified local driver** (with hot-plug support) |

PCIe IO bus

**PCIe NTB interconnect** ← Hardware mappings ensure **fast data path**

## Remote

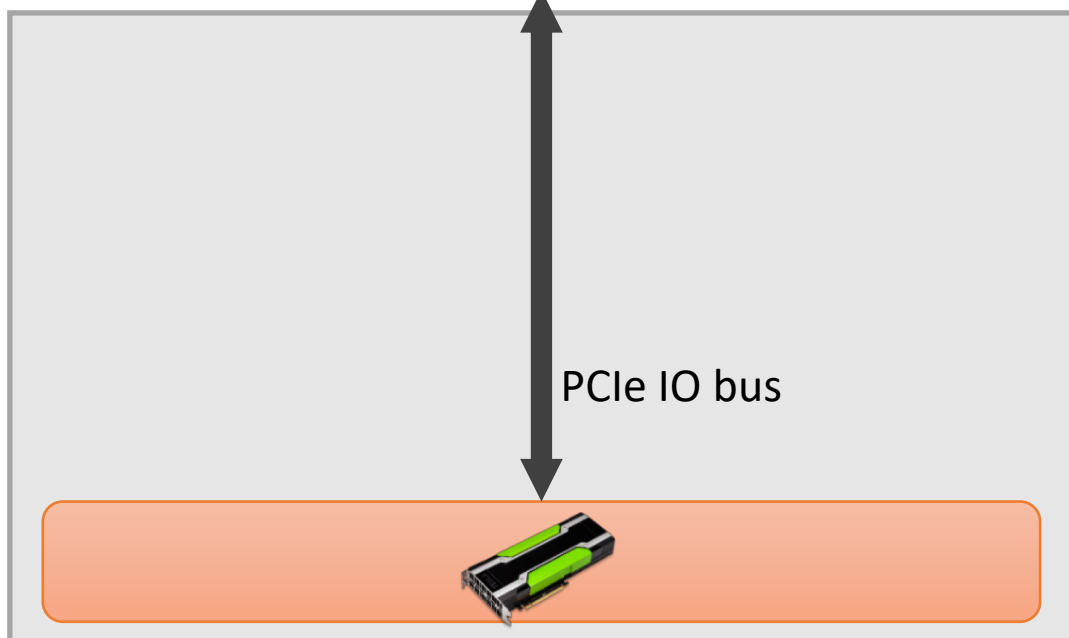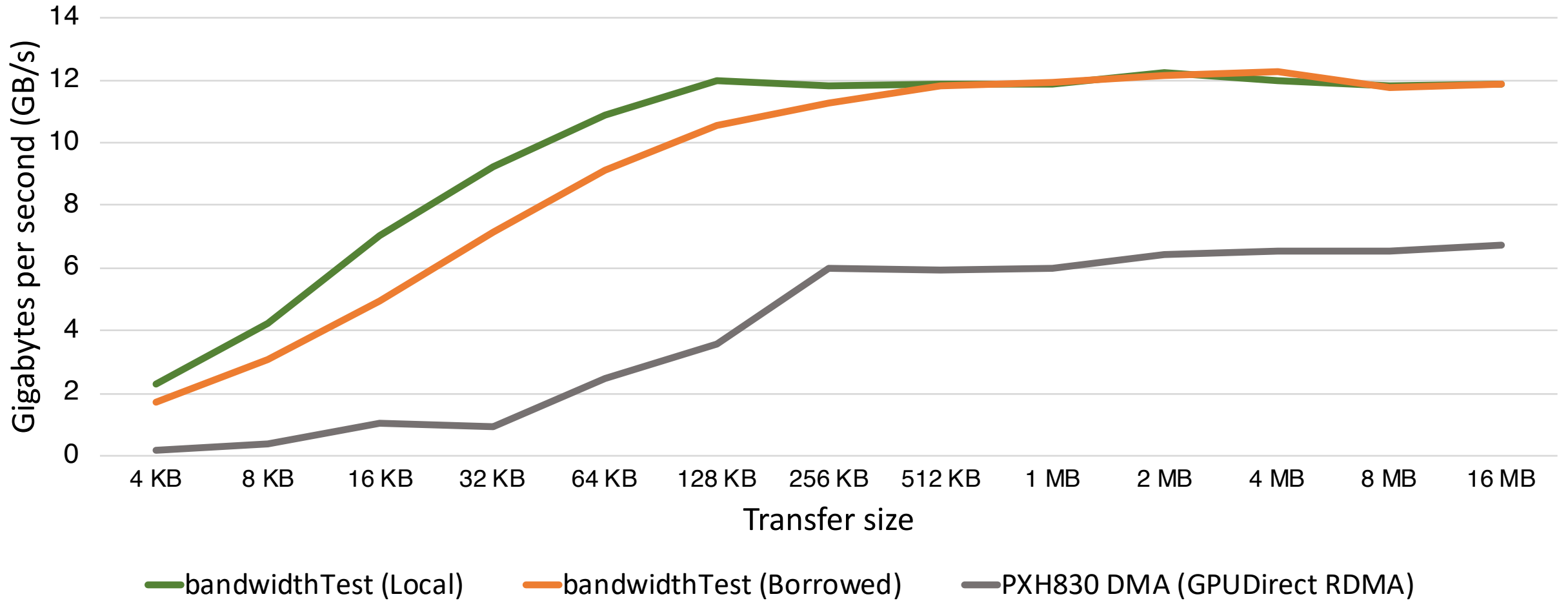PCIe IO bus

Works with **any PCIe device** (even individual SR-IOV functions)

**Borrowed** remote resource

Remote resource using middleware

Local

Remote

Application

CUDA library + driver

PCIe IO bus

PCIe NTB interconnect

PCIe IO bus

Application

CUDA – middleware integration

Middleware service

Interconnect transport (RDMA)

Interconnect

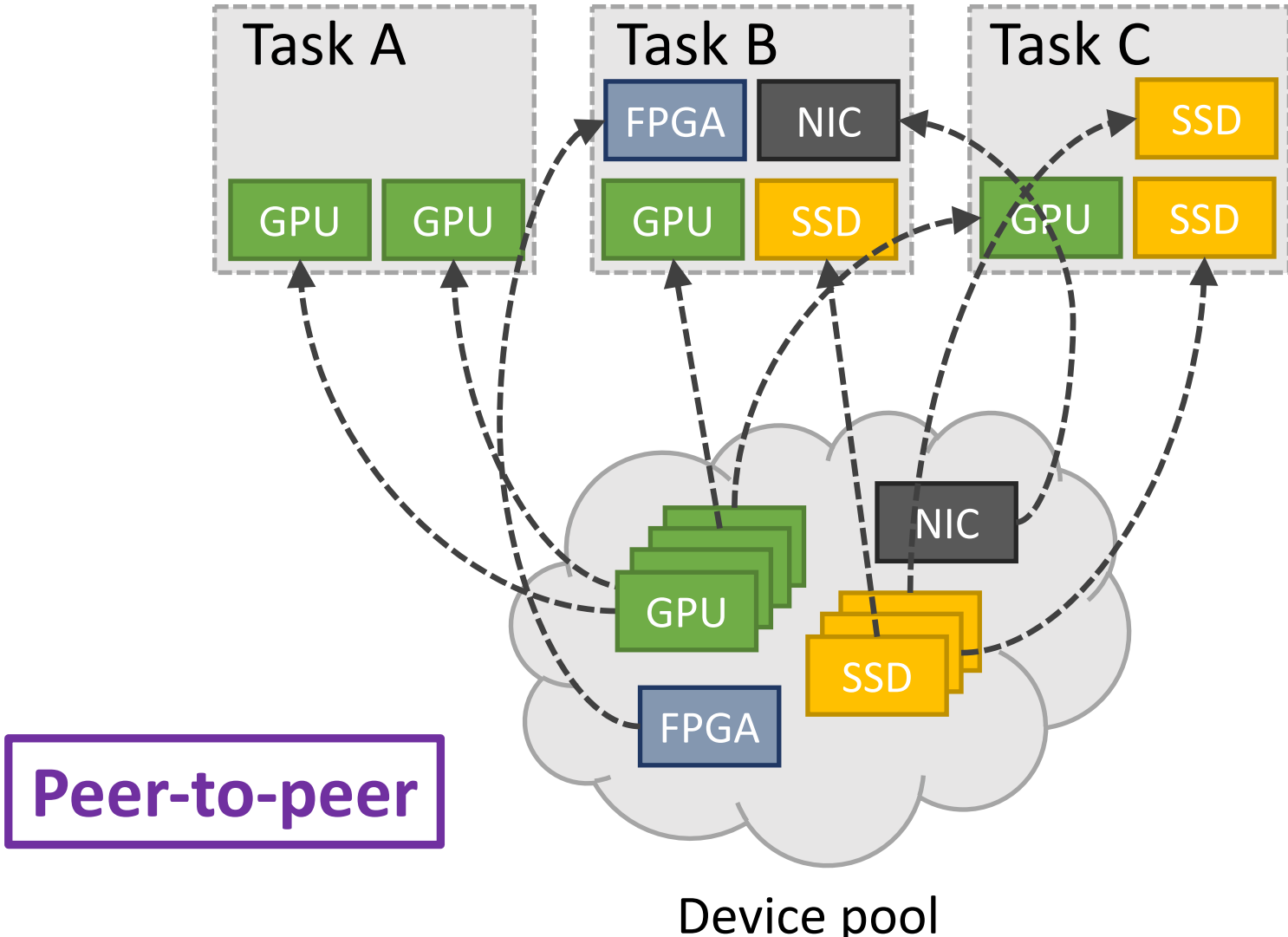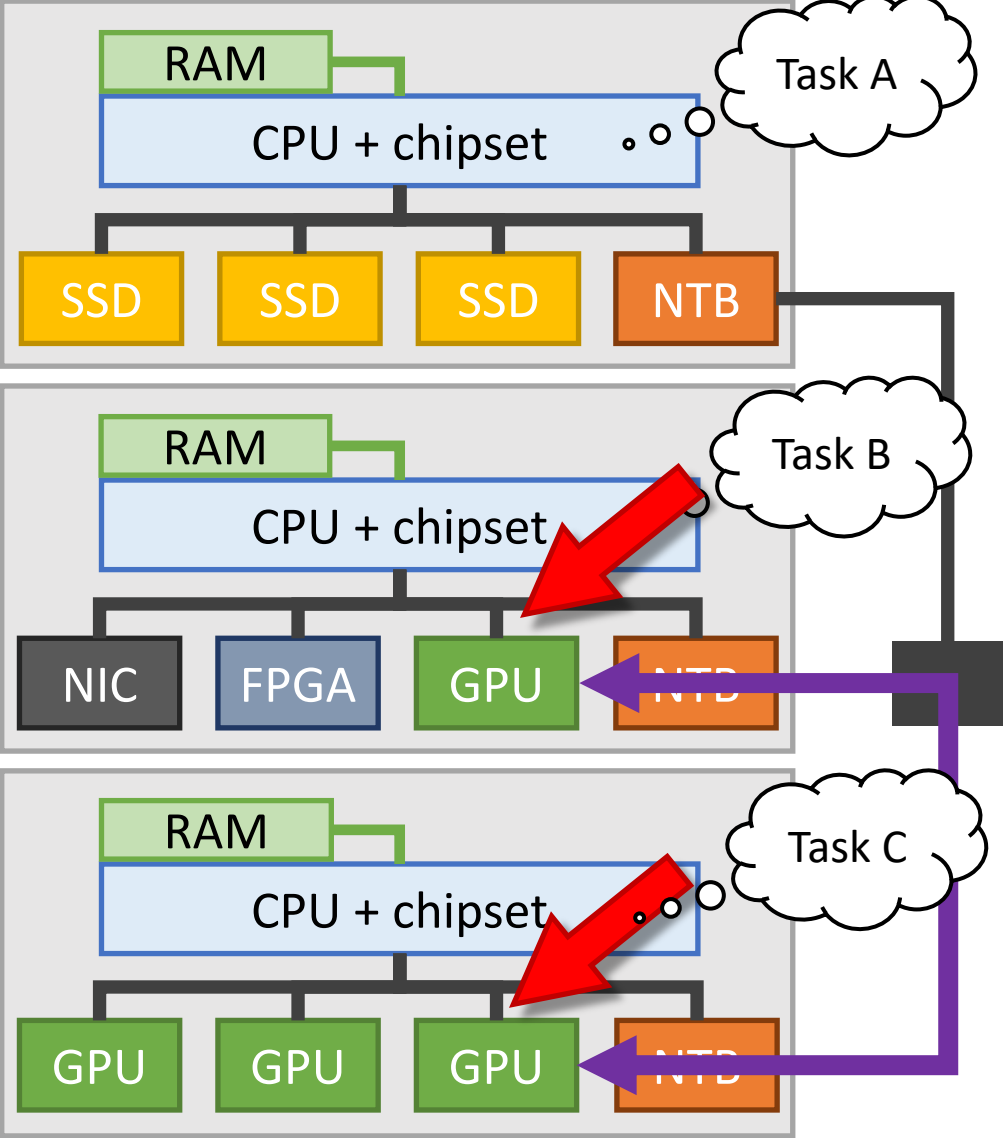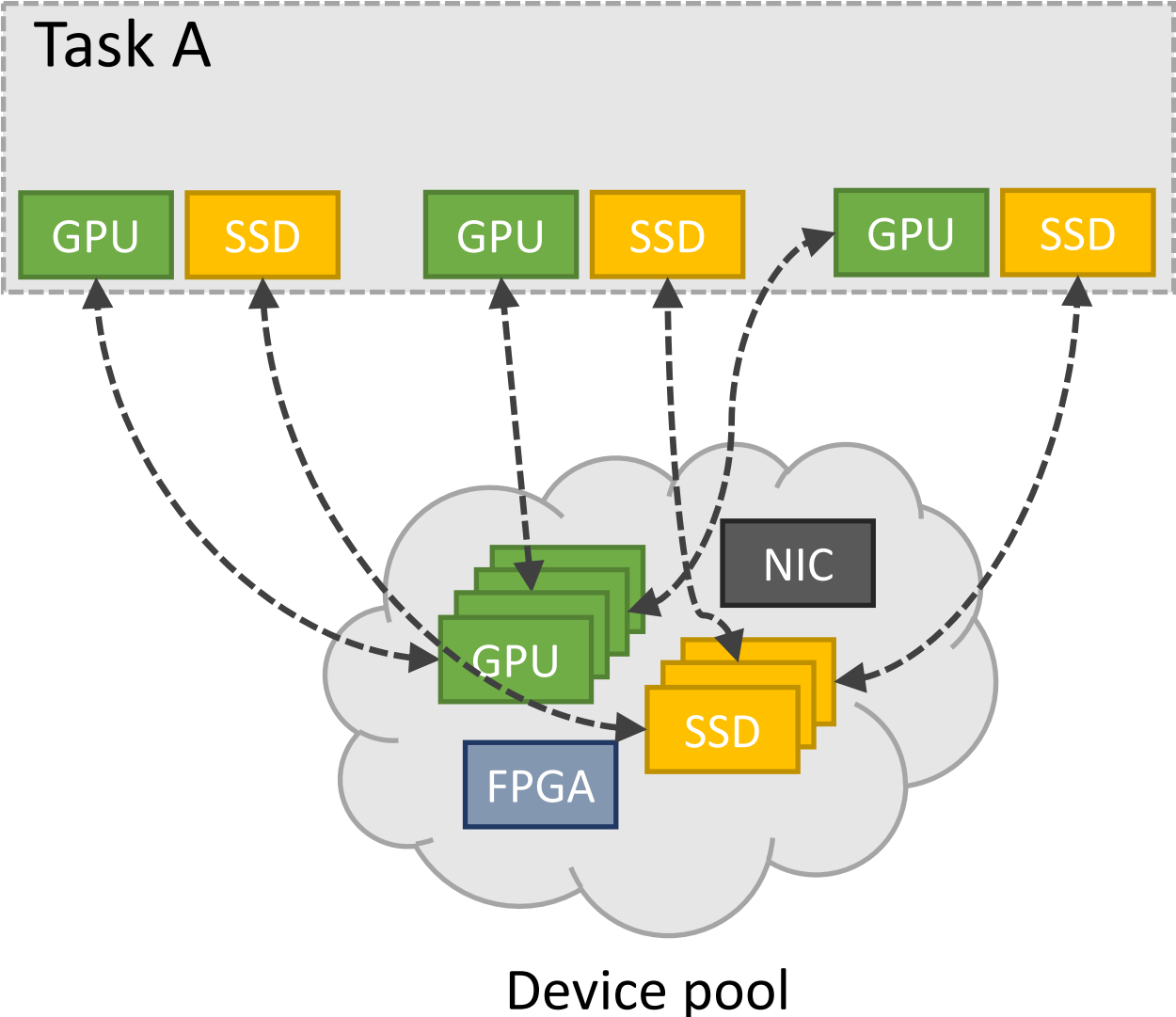Interconnect transport (RDMA)

Middleware service/daemon

CUDA driver

PCIe IO bus

**Device to host transfers: Comparing local to borrowed GPU**

_Gigabytes per second (GB/s)_ vs _Transfer size_

X-axis: 4 KB, 8 KB, 16 KB, 32 KB, 64 KB, 128 KB, 256 KB, 512 KB, 1 MB, 2 MB, 4 MB, 8 MB, 16 MB

Legend:
— bandwidthTest (Local)
— bandwidthTest (Borrowed)
— PXH830 DMA (GPUDirect RDMA)

# Using Device Lending, nodes in a PCIe cluster can share resources through a process of borrowing and giving back devices



Peer-to-peer

Device pool

# Using Device Lending, nodes in a PCIe cluster can share resources through a process of borrowing and giving back devices
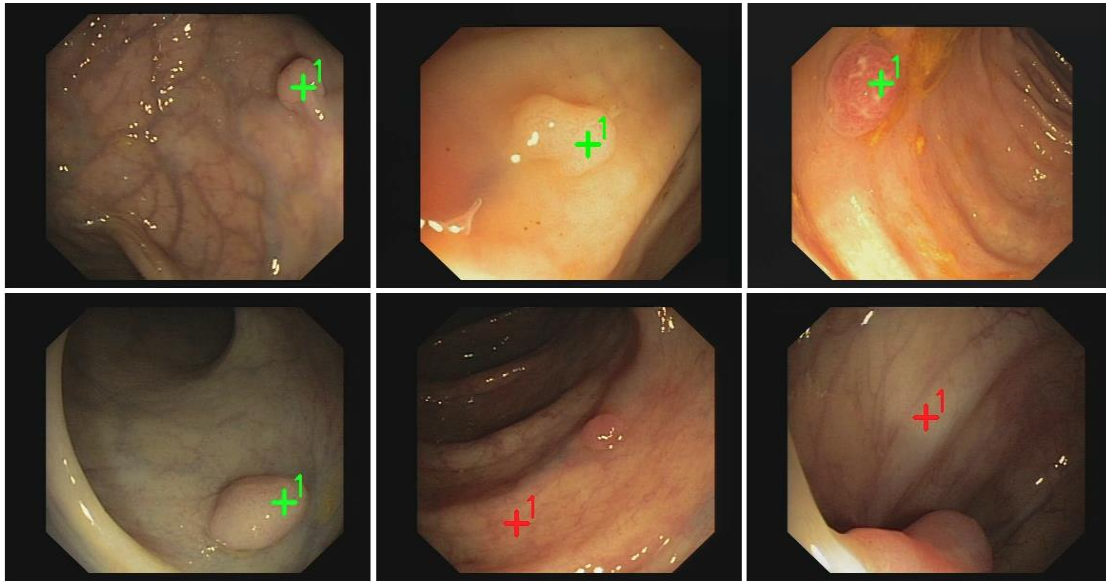


Device pool

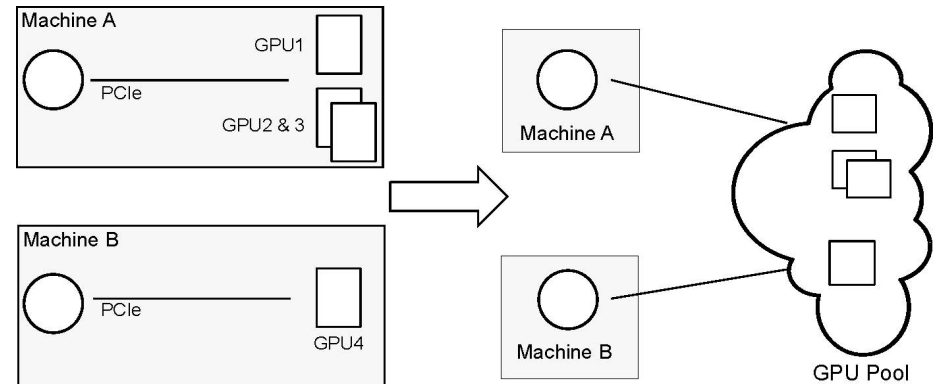# Example Application

## Processing of Medical Videos

P9258 - Efficient Processing of Medical Videos in a Multi-auditory Environment Using GPU Lending
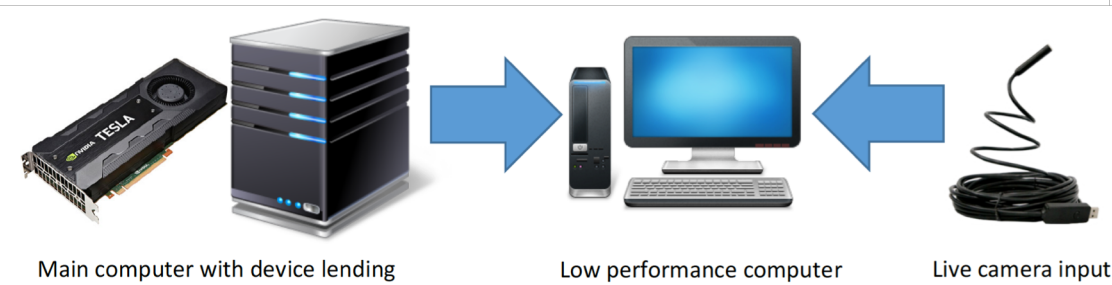
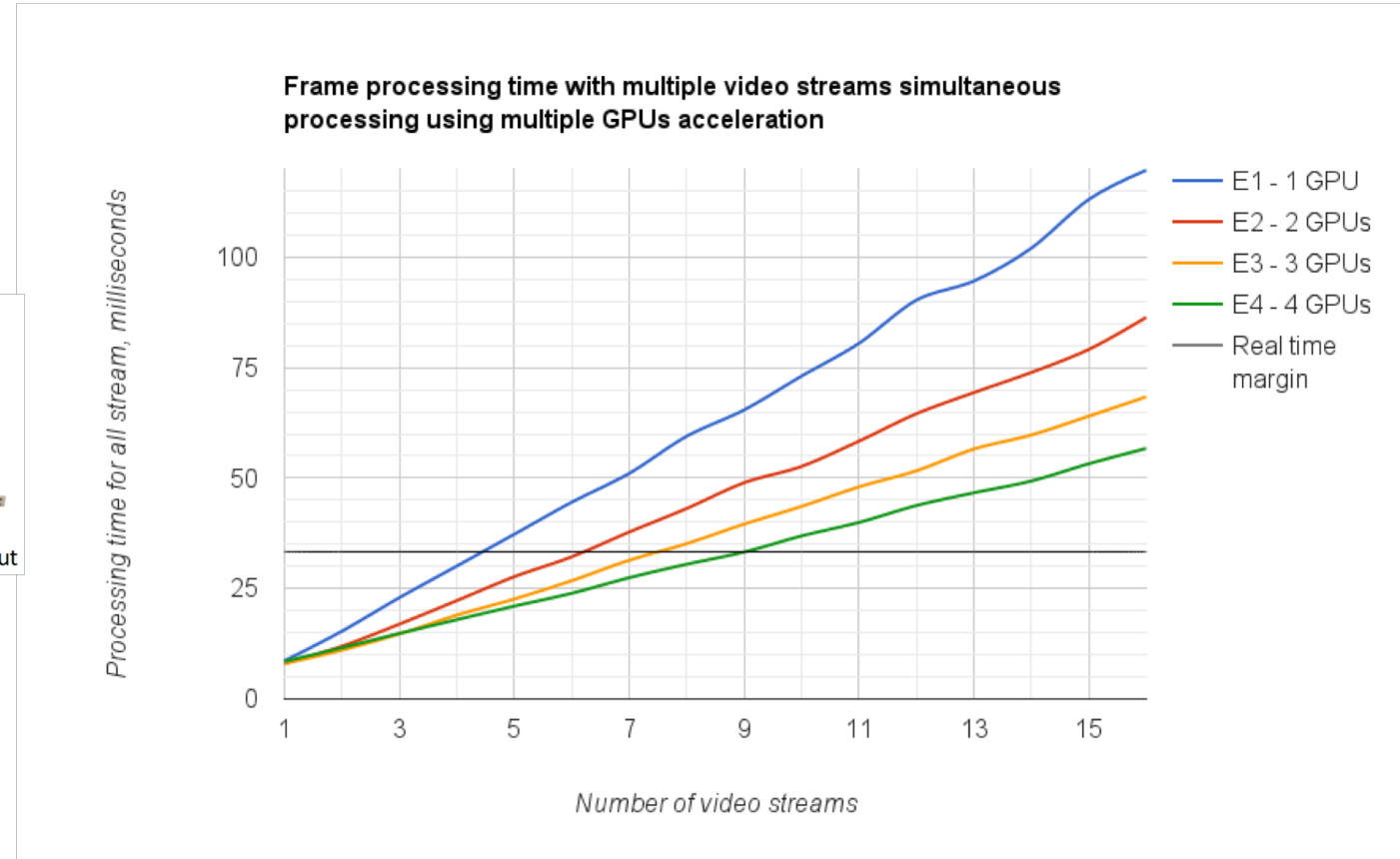# Scenario: Real-time computer-aided polyp detection



- PCIe fiber cables can be up to 100 meters.
- Enable "thin clients" to use GPUs in remote machine room

# Flexible sharing of GPU resources between multiple examination rooms



Main computer with device lending — Low performance computer — Live camera input

Frame processing time with multiple video streams simultaneous processing using multiple GPUs acceleration

- E1 - 1 GPU
- E2 - 2 GPUs
- E3 - 3 GPUs
- E4 - 4 GPUs
- Real time margin

Processing time for all stream, milliseconds

Number of video streams

- System uses a combination of classic computer vision algorithms and machine learning.
- Research prototype since 2016.
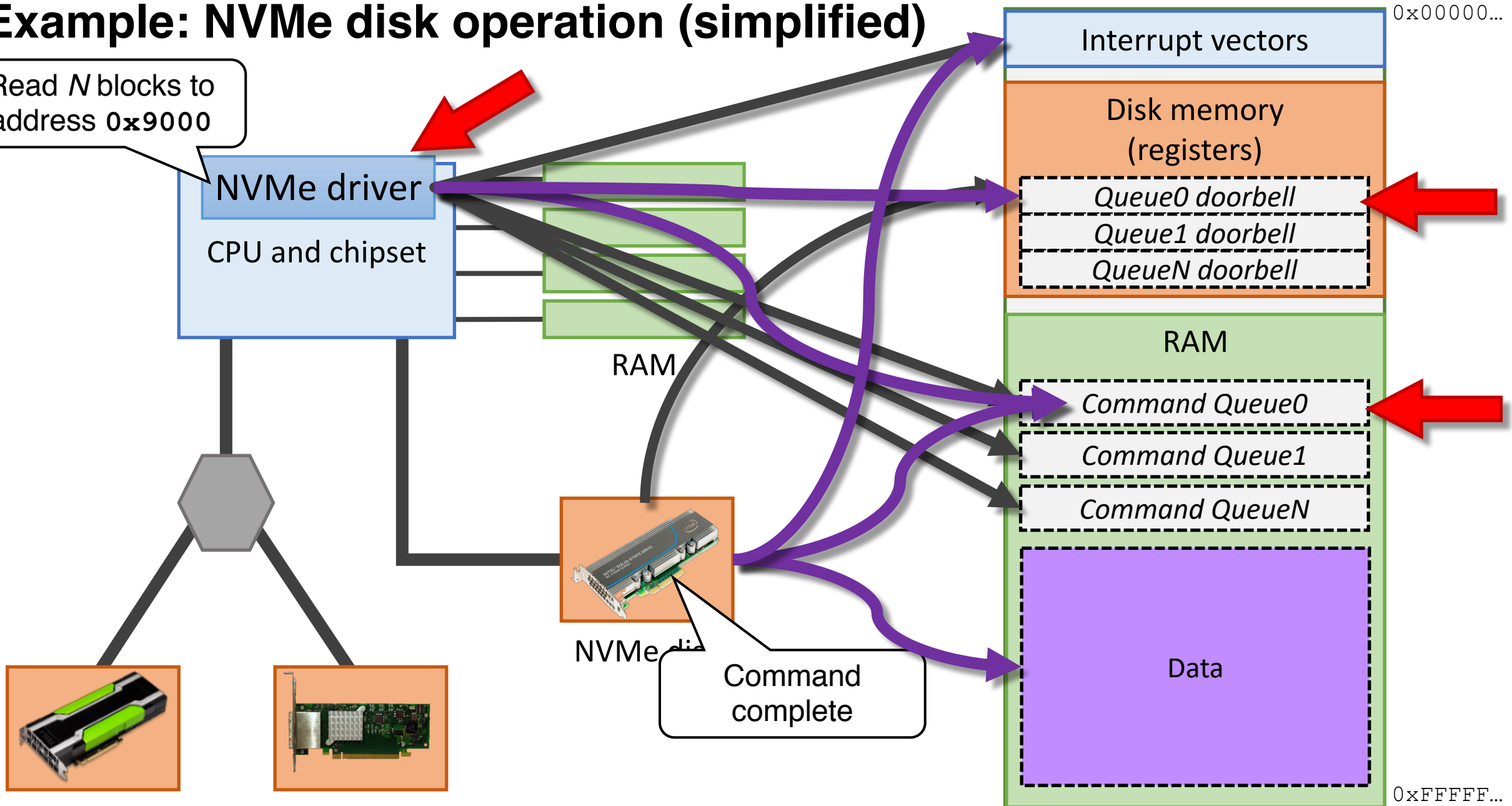
# Sharing of NVMe drives

For more details:
S9563 - Efficient Distributed Storage I/O using NVMe and GPU Direct in a PCIe Network
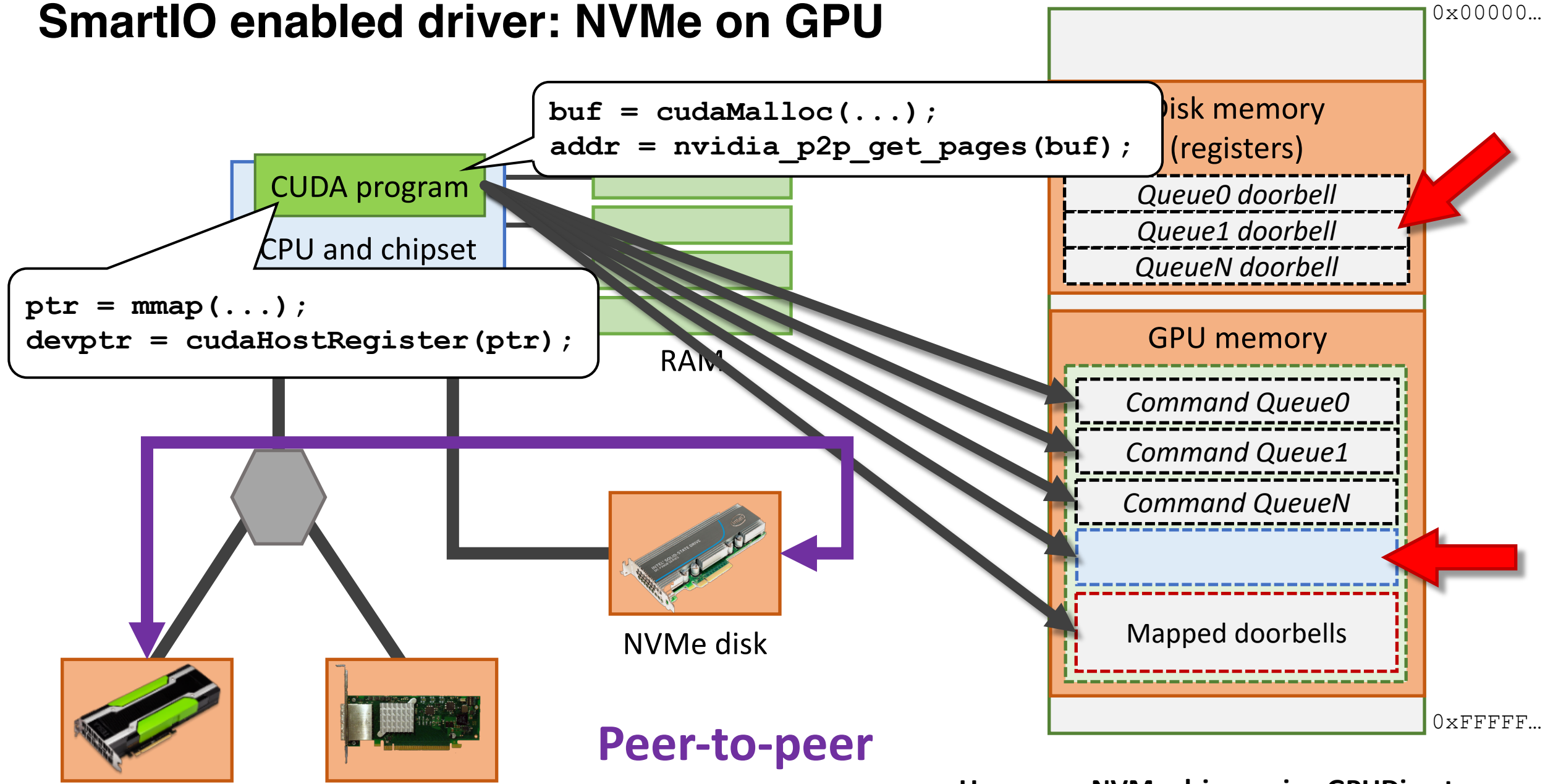or
*Visit Dolphin Interconnect Solutions in booth 1520*

# Example: NVMe disk operation (simplified)

Read *N* blocks to address `0x9000`

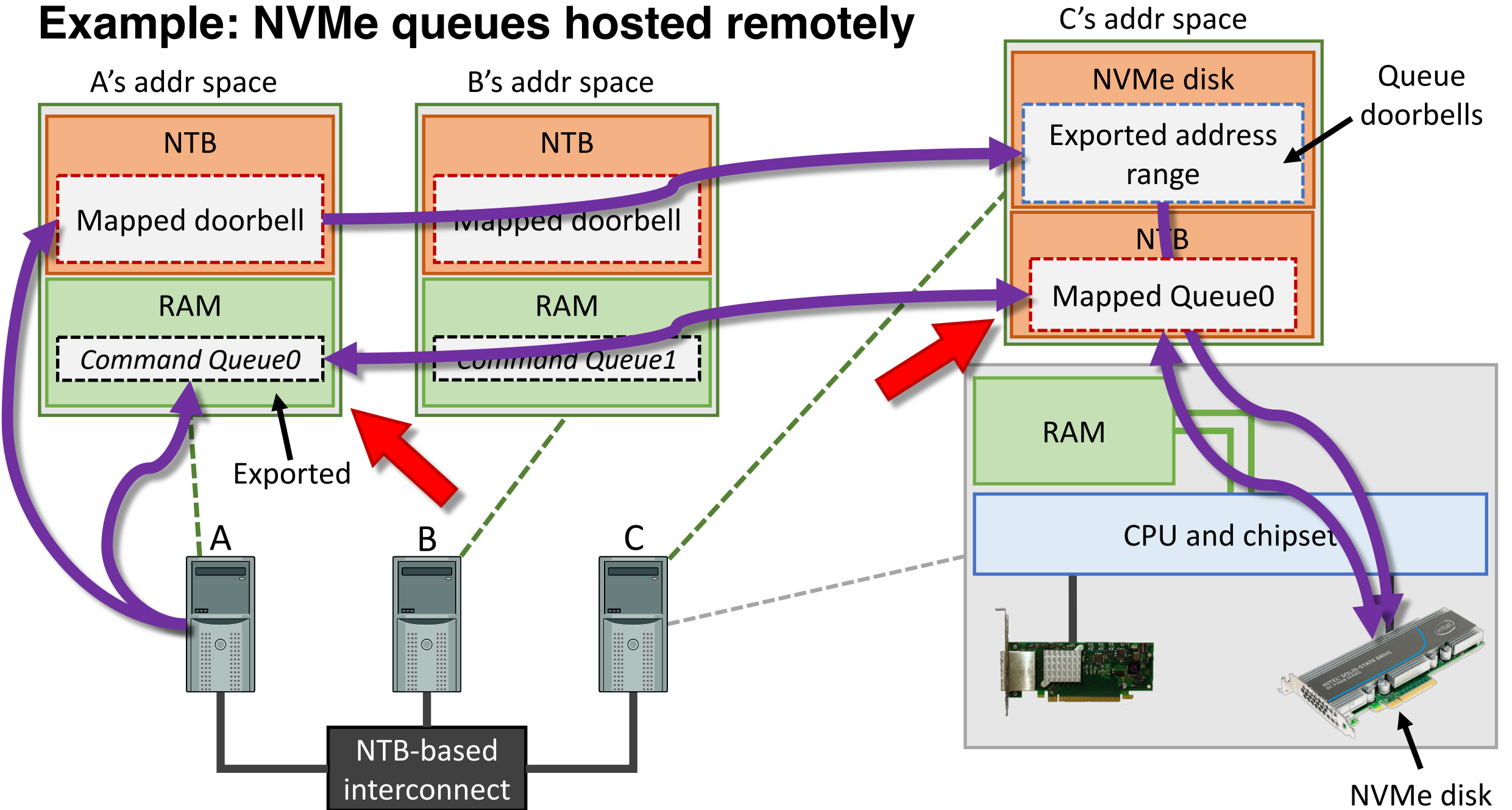NVMe driver

CPU and chipset

RAM

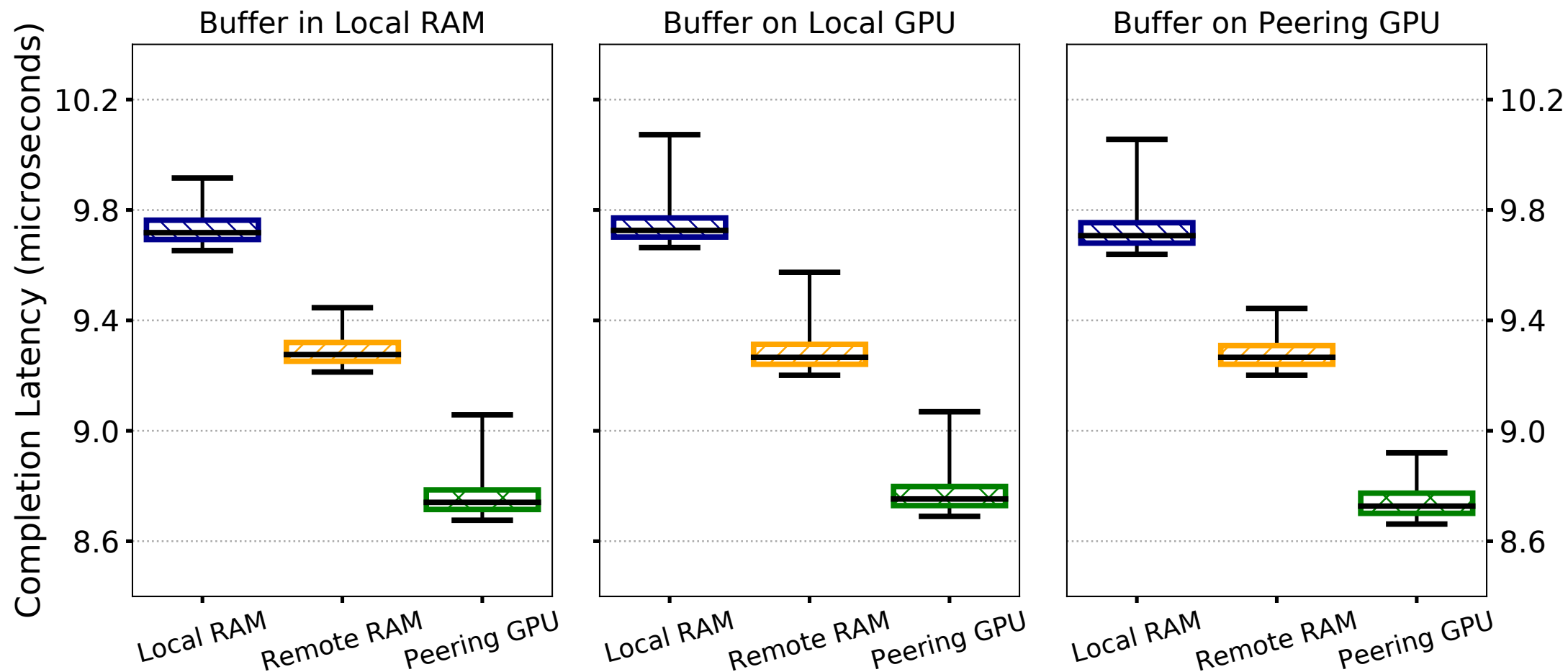NVMe disk

Command complete

Interrupt vectors

Disk memory (registers)

*Queue0 doorbell*
*Queue1 doorbell*
*QueueN doorbell*

RAM

*Command Queue0*
*Command Queue1*
*Command QueueN*

Data

`0x00000...`

`0xFFFFF...`

# SmartIO enabled driver: NVMe on GPU

`0x00000…`

```
buf = cudaMalloc(...);
addr = nvidia_p2p_get_pages(buf);
```

```
ptr = mmap(...);
devptr = cudaHostRegister(ptr);
```

CUDA program

CPU and chipset

RAM

Disk memory (registers)

*Queue0 doorbell*
*Queue1 doorbell*
*QueueN doorbell*

GPU memory

*Command Queue0*
*Command Queue1*
*Command QueueN*
Mapped doorbells

NVMe disk

GPU

**Peer-to-peer**

**Userspace NVMe driver using GPUDirect**
*https://github.com/enfiskutensykkel/ssd-gpu-dma*

`0xFFFFF…`

# Example: NVMe queues hosted remotely



A's addr space

B's addr space

C's addr space

NVMe disk

Queue doorbells

NTB — Mapped doorbell

RAM — Command Queue0

Exported

NTB — Mapped doorbell

RAM — Command Queue1

Exported address range

NTB — Mapped Queue0

RAM

CPU and chipset

NVMe disk
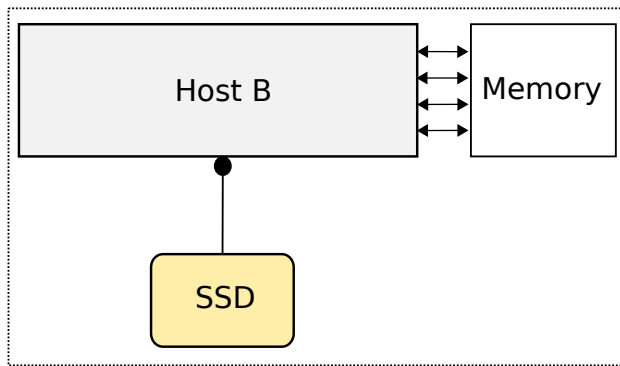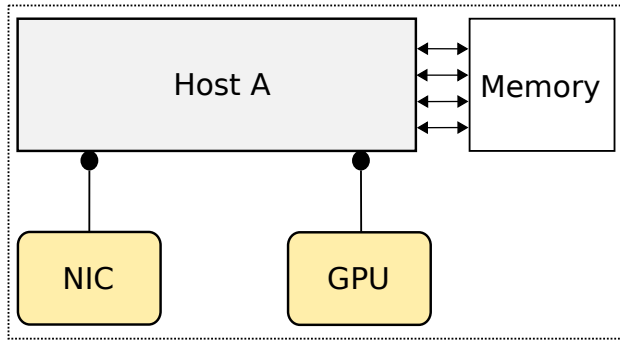
A    B    C

NTB-based interconnect
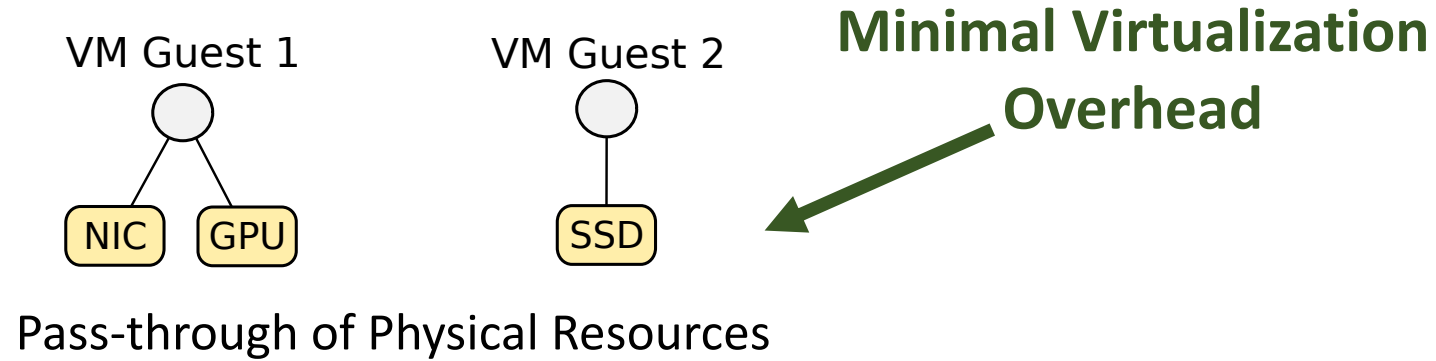
# SmartIO in Virtual Machines

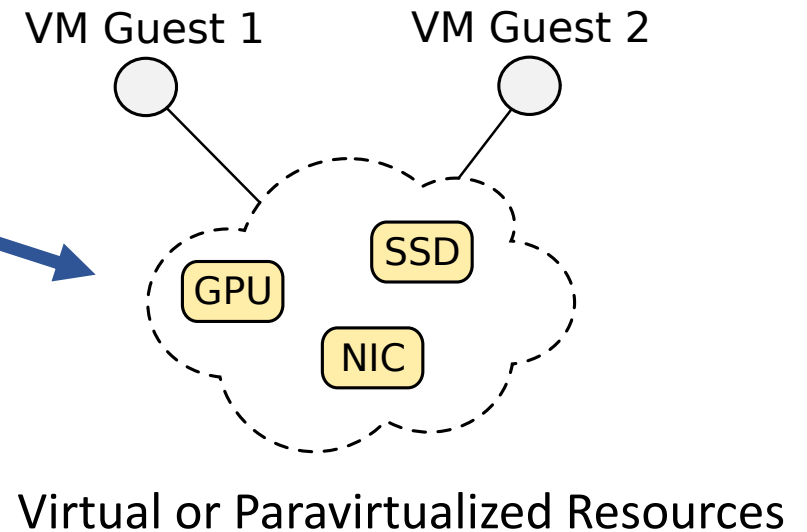# SmartIO fully supports to lend devices to virtual machines running in Linux KVM uning Virtual Function IO API (VFIO)

# Pass-through allows physical devices to be used by VMs with minimal overhead, but is not as flexible as resource virtualization
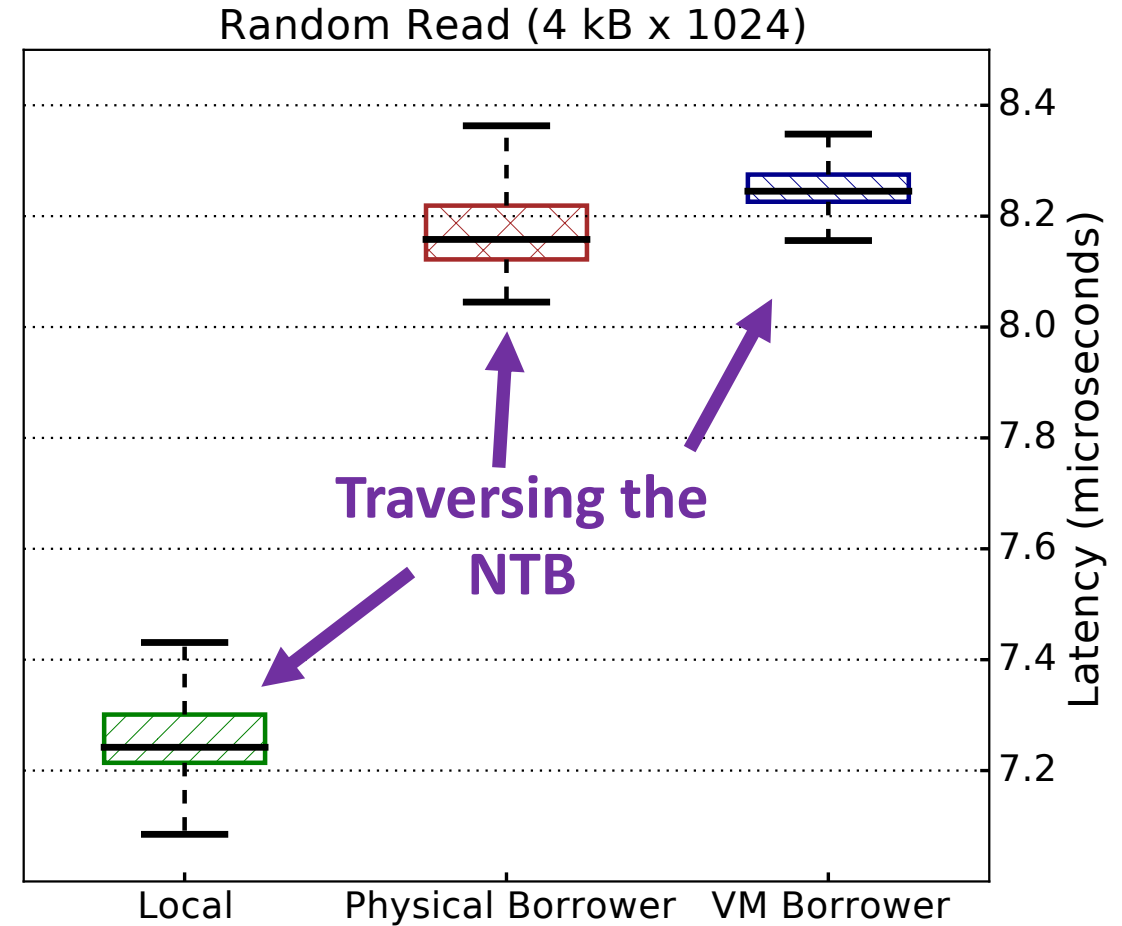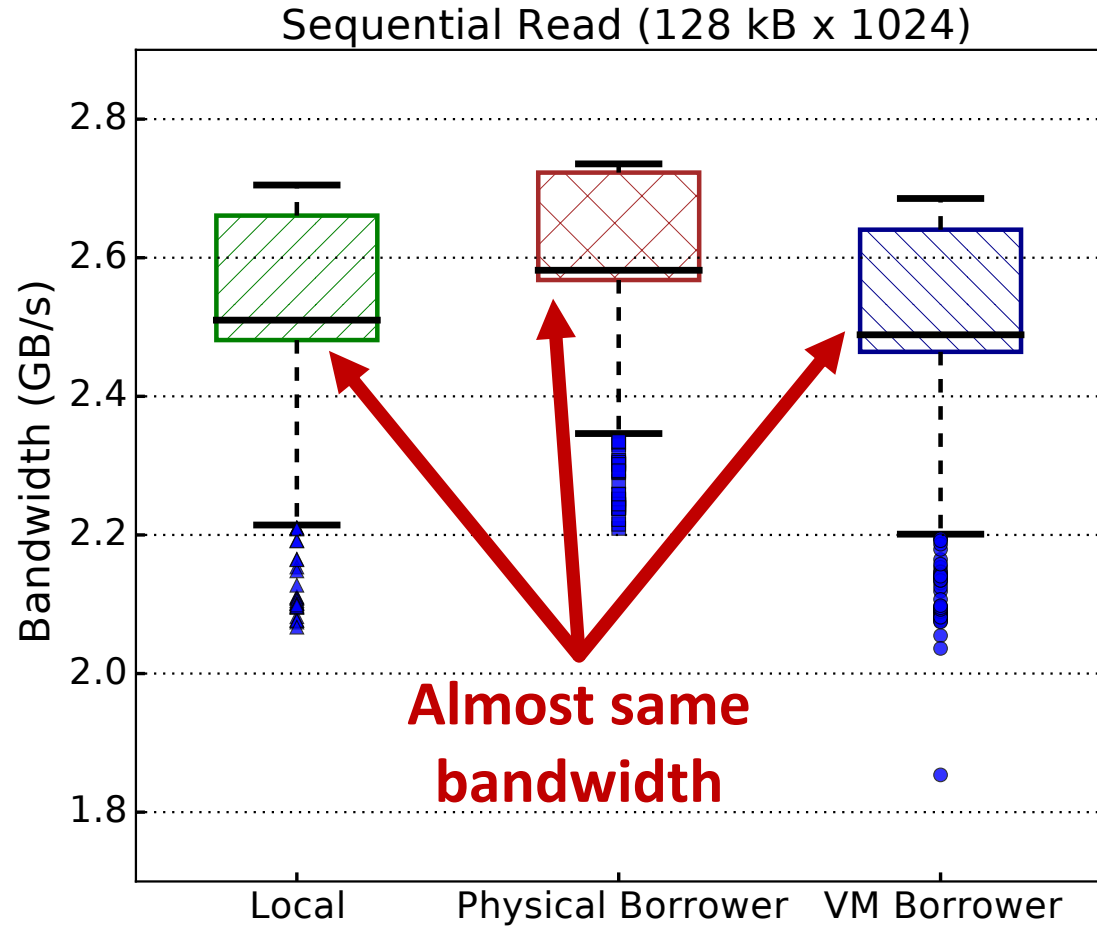


**Minimal Virtualization Overhead**

Pass-through of Physical Resources

**Dynamic Provisioning & Flexible Composition**

Physical View

Virtual or Paravirtualized Resources

# Passing through a remote NVMe disk to a VM only adds the latency of traversing the NTB and is comparable to a physical borrower



Sequential Read (128 kB x 1024)

Random Read (4 kB x 1024)

Almost same bandwidth

Traversing the NTB

Local  Physical Borrower  VM Borrower

Guest OS: Ubuntu 17.04, Host OS: CentOS 7
VM: Qemu 2.17 using KVM
NVMe Disk: Intel 900P Optane (PCIe x4 Gen3)

# Thank you!

Selected
publications

*"Device Lending in PCI Express Networks"*
ACM NOSSDAV 2016

*"Efficient Processing of Video in a Multi Auditory Environment using Device Lending of GPUs"*
ACM Multimedia Systems 2016 (MMSys'16)

*"Flexible Device Sharing in PCIe Clusters using Device Lending",* International Conference on Parallel Processing Companion (ICPP'18 Comp)

haakonks@simula.no

**SmartIO & Device Lending demo with GPUs, NVMe and more**
Visit Dolphin in the exhibition area (booth 1520)

simula  Dolphin  UNIVERSITAS OSLOENSIS MDCCCXI