# S8511
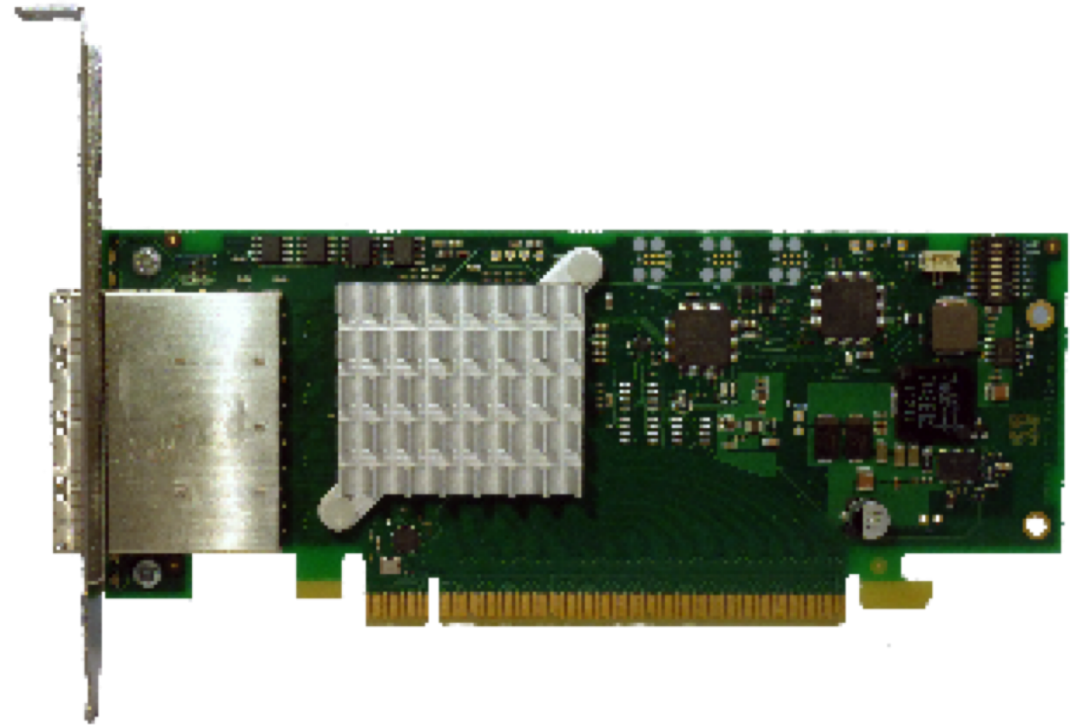# *SmartIO: Dynamic Sharing of GPUs and IO in a PCIe Cluster*

**Håkon Kvale Stensland**
Research Scientist / Associate Professor
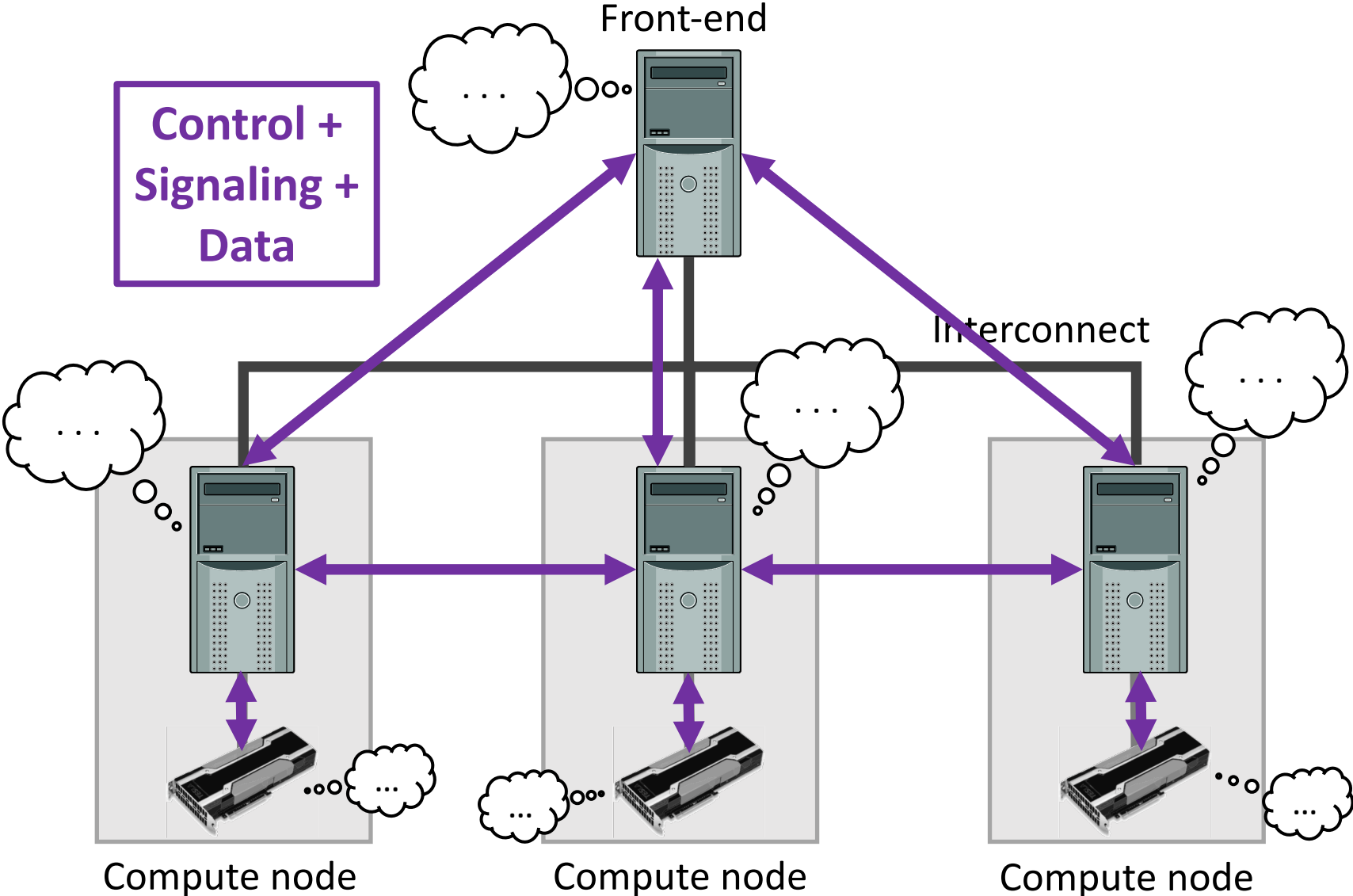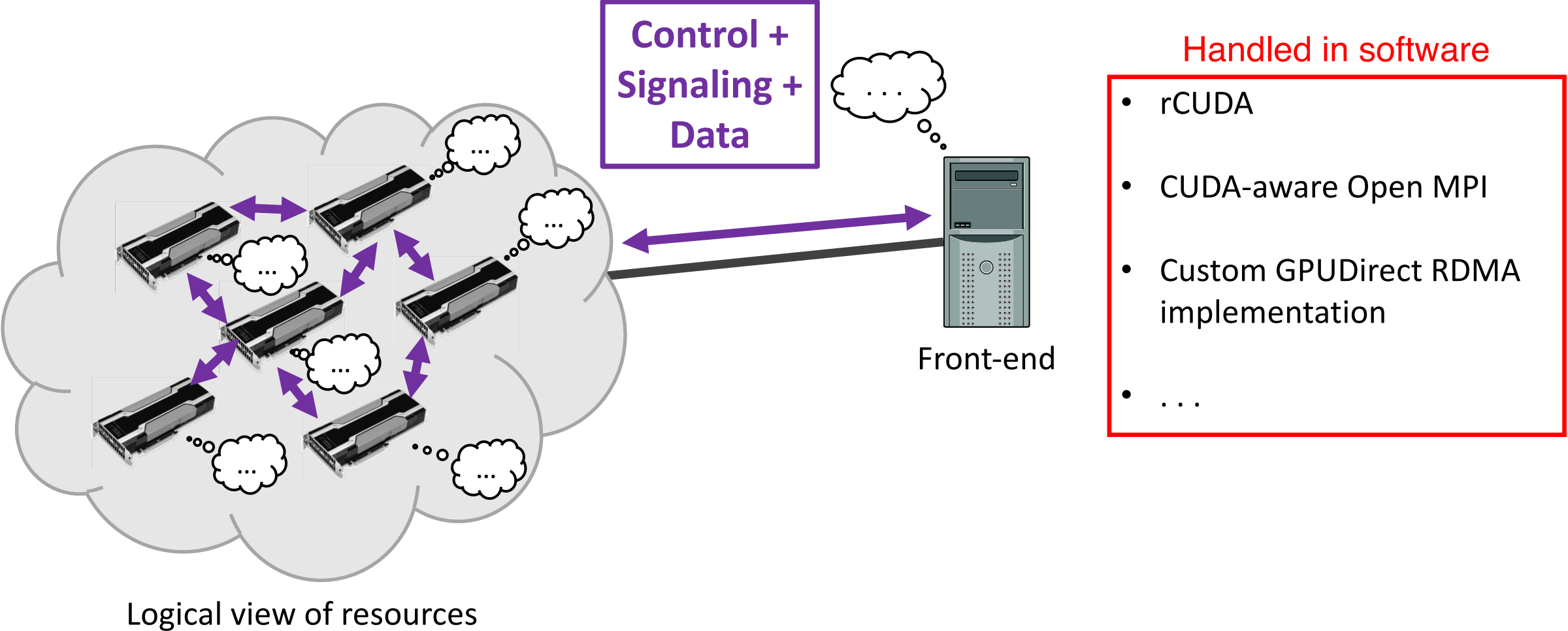Simula Research Laboratory / University of Oslo

# Outline

- Motivation

- PCIe Overview

- Non-Transparent Bridges

- Dolphin SmartIO

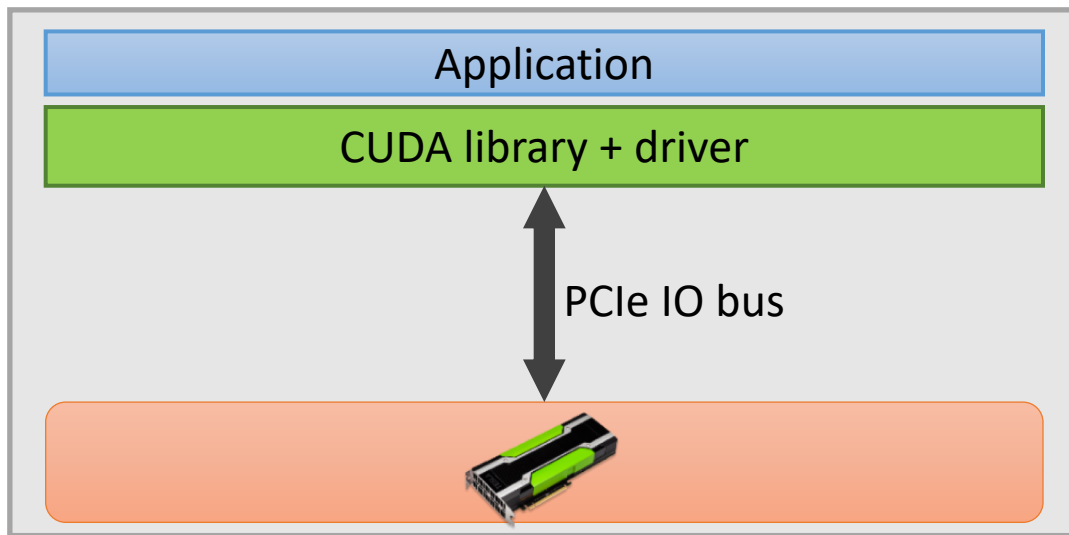# Distributed applications may need to access and use IO resources that are physically located inside remote hosts



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



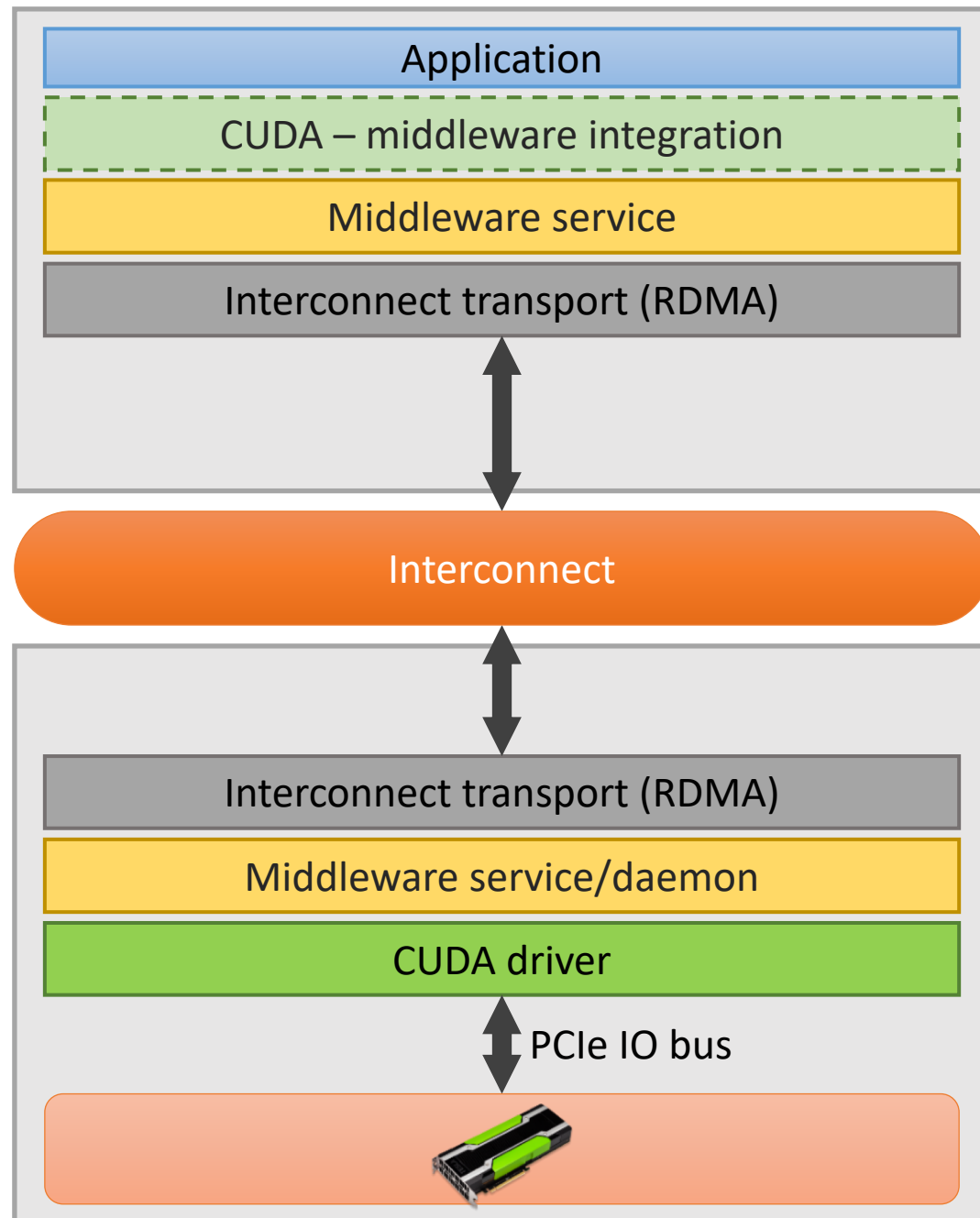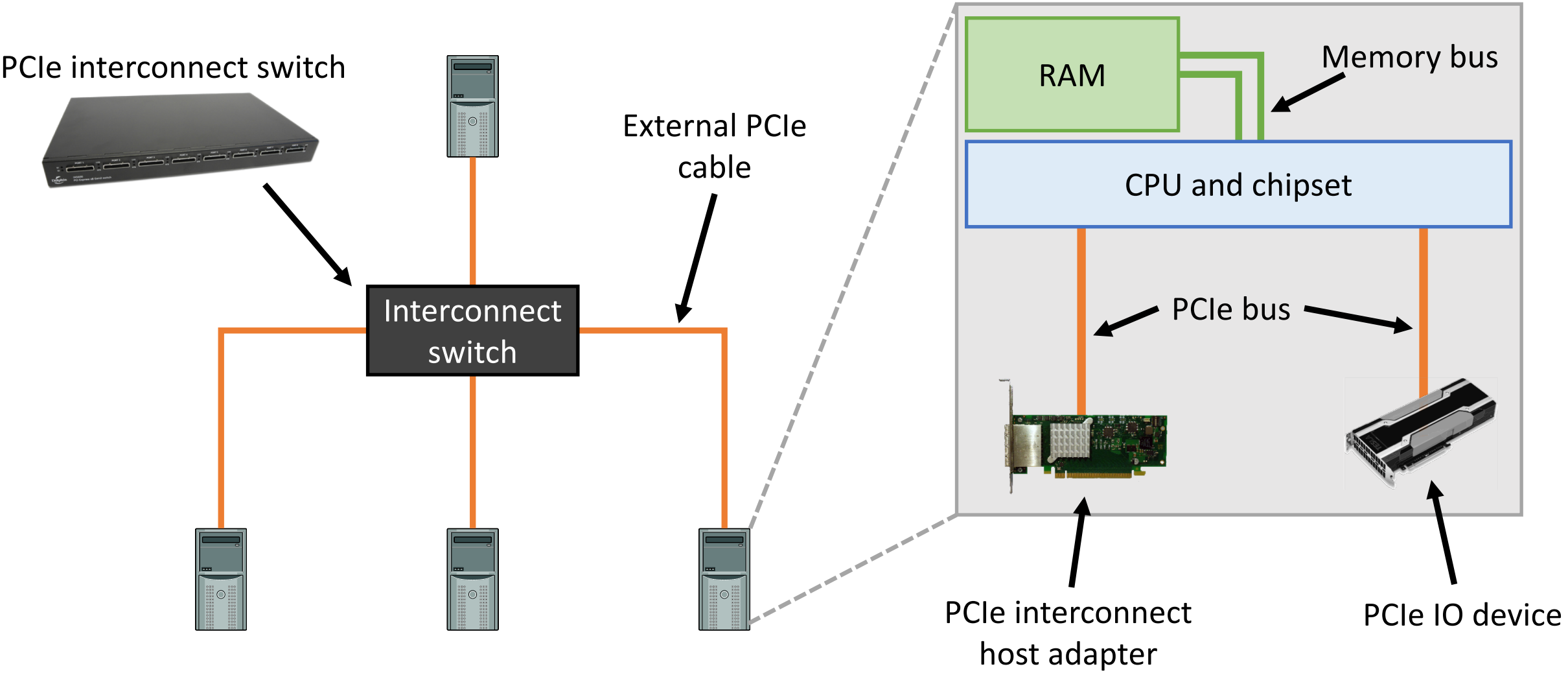**Control + Signaling + Data**

Front-end

Logical view of resources

Handled in software

- rCUDA
- CUDA-aware Open MPI
- Custom GPUDirect RDMA implementation
- . . .

Local resource

Remote resource using **middleware**

Local

| Application |
| CUDA library + driver |

PCIe IO bus

Remote

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

Interconnect

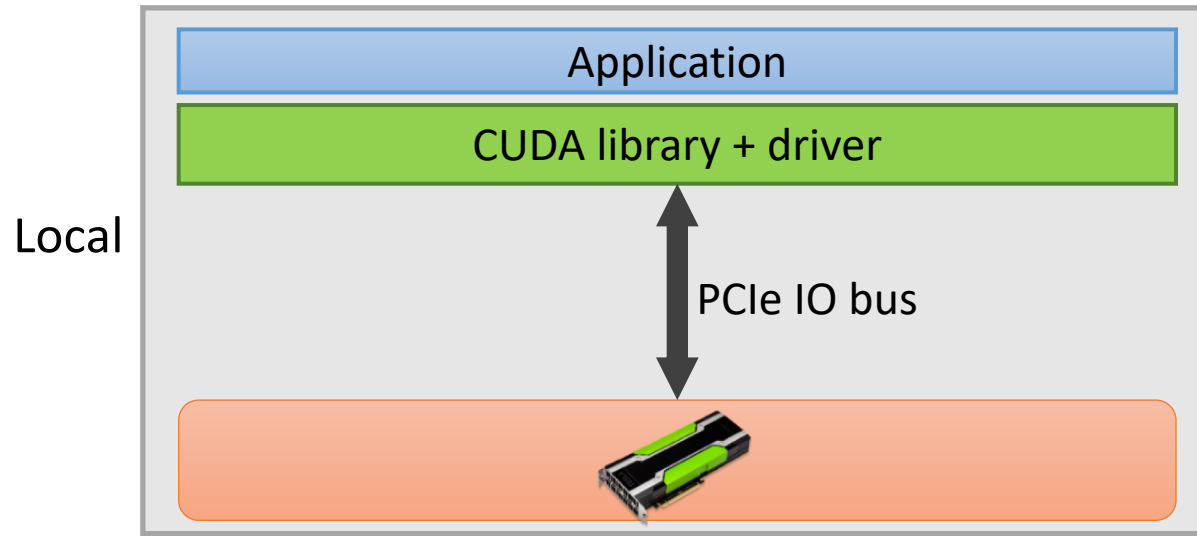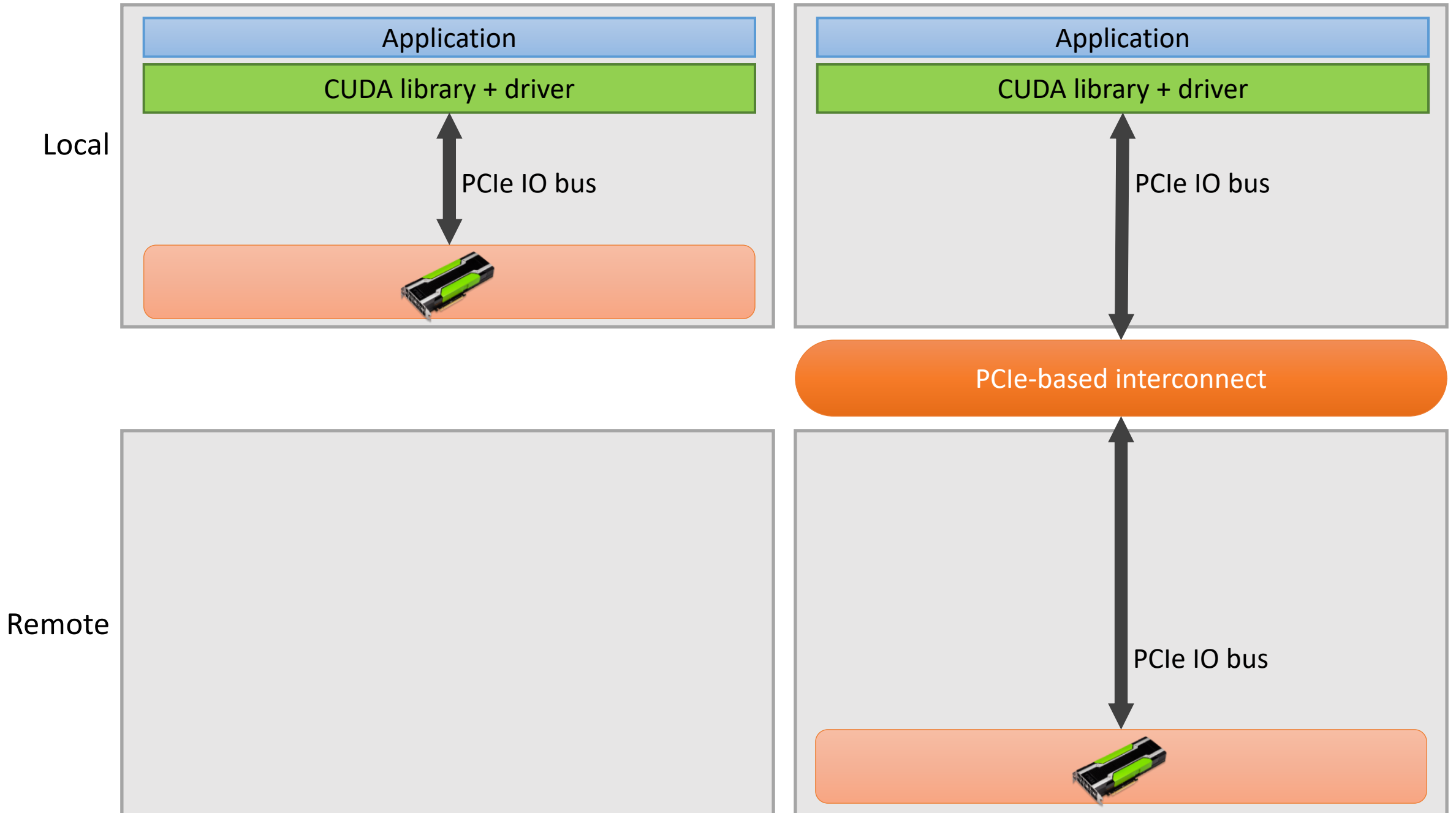| Interconnect transport (RDMA) |
| Middleware service/daemon |
| CUDA driver |

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
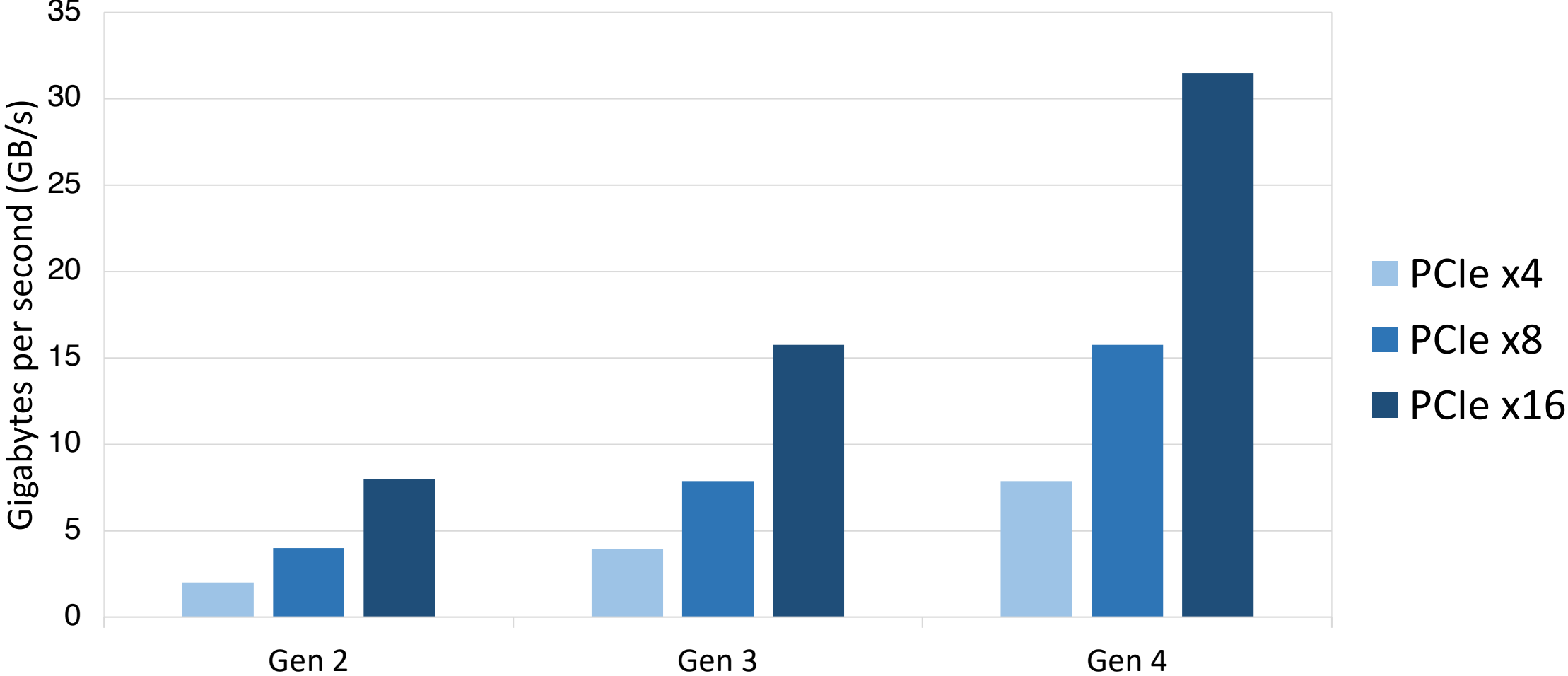
PCIe interconnect switch

External PCIe cable

Interconnect switch

RAM

Memory bus

CPU and chipset

PCIe bus

PCIe interconnect host adapter

PCIe IO device

# Local resource

# Remote resource over **native fabric**

| Application |
| --- |
| CUDA library + driver |

**Local**

PCIe IO bus

| Application |
| --- |
| CUDA library + driver |

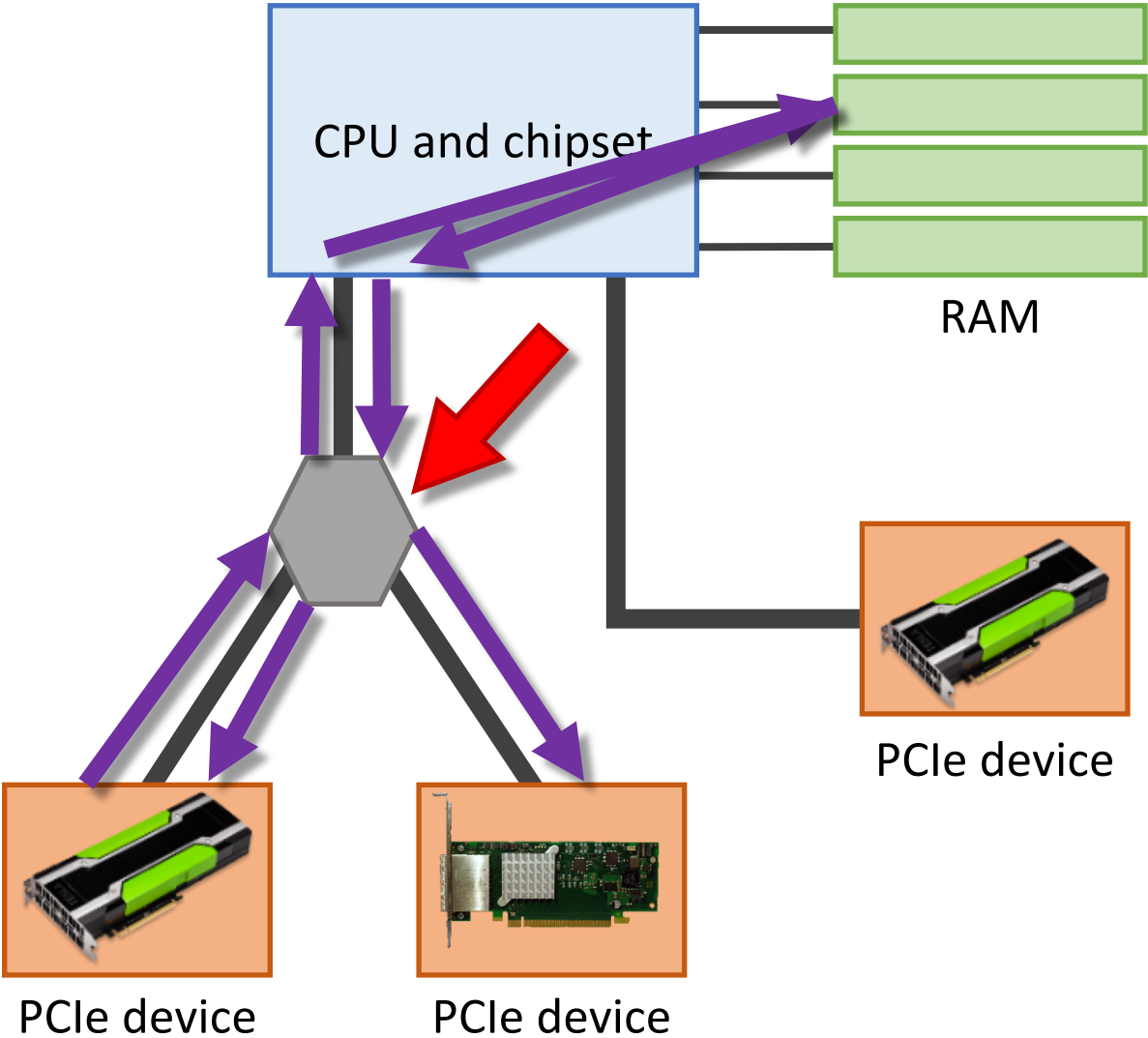PCIe IO bus

PCIe-based interconnect

**Remote**

PCIe IO bus

# PCIe Overview

# PCIe is the dominant IO bus technology in computers today, and can also be used as a high-bandwidth low-latency interconnect
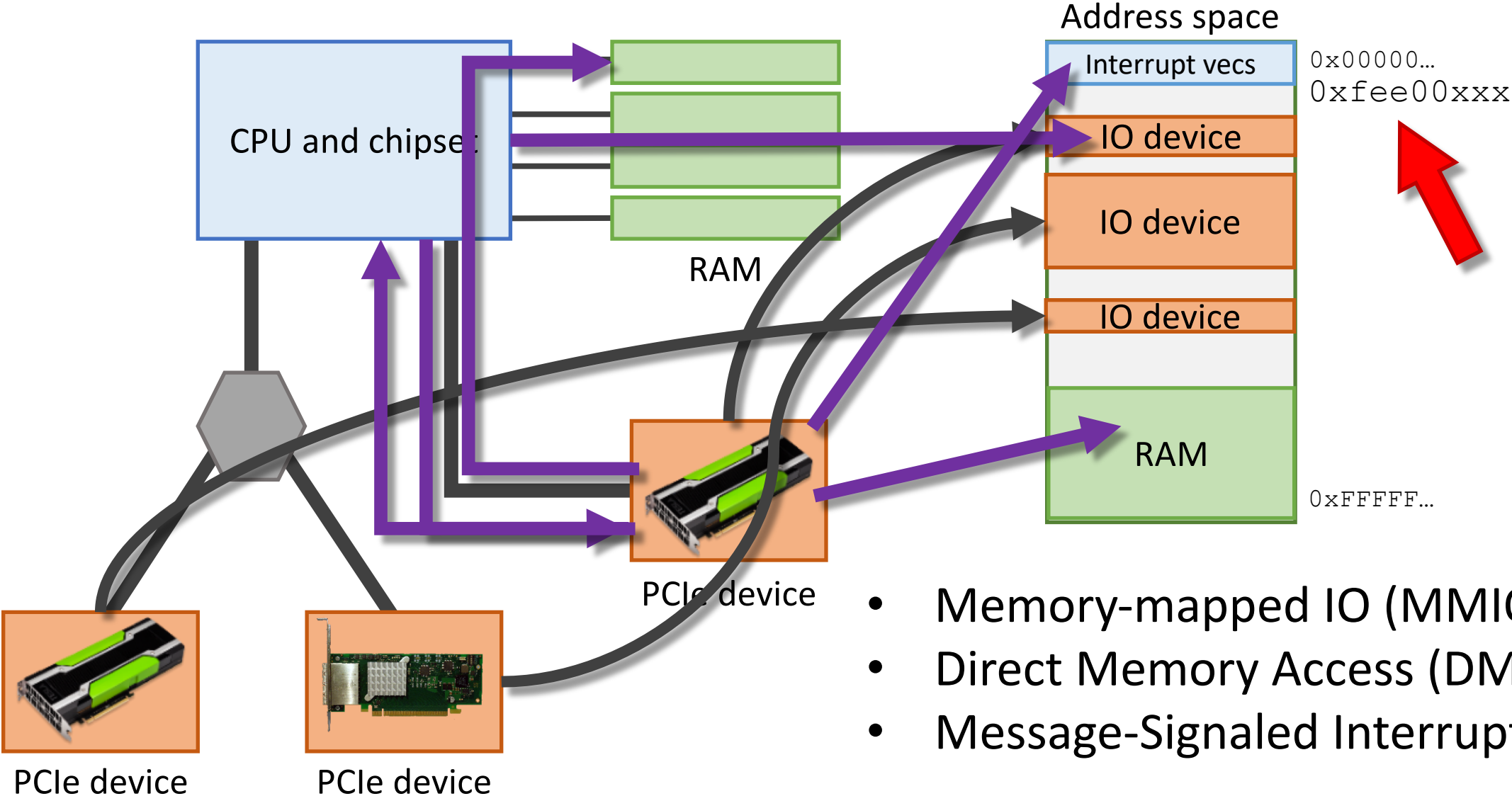
# Memory reads and writes are handled by PCIe as transactions that are packet-switched through the fabric depending on the address



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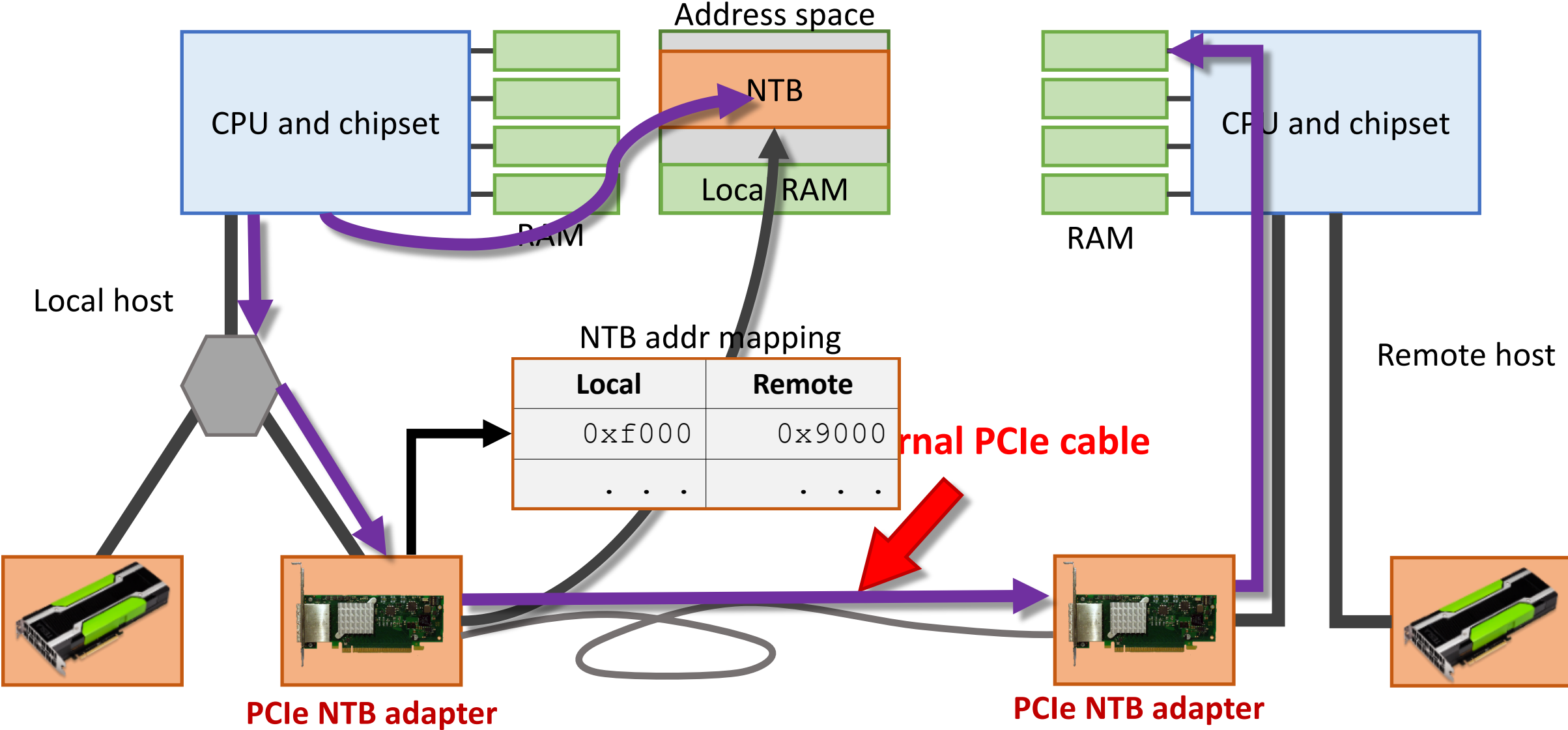- Peer-to-peer (shortest path)

CPU and chipset

RAM

PCIe device

PCIe device

PCIe device

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

Address space

CPU and chipset

RAM

PCIe device

PCIe device

PCIe device

PCIe device

Interrupt vecs    0x00000…
                  0xfee00xxx

IO device

IO device
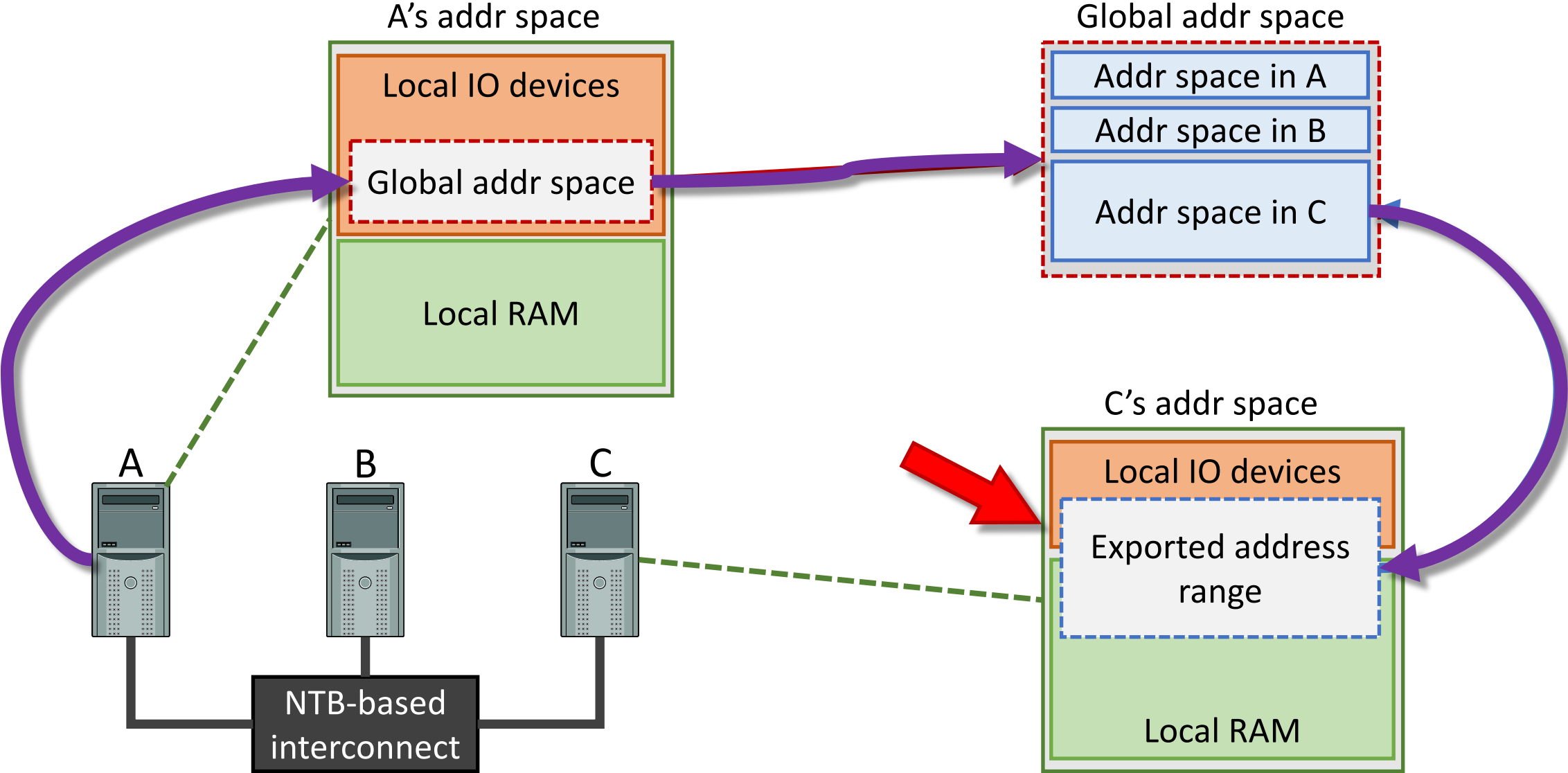
IO device

RAM    0xFFFFF…

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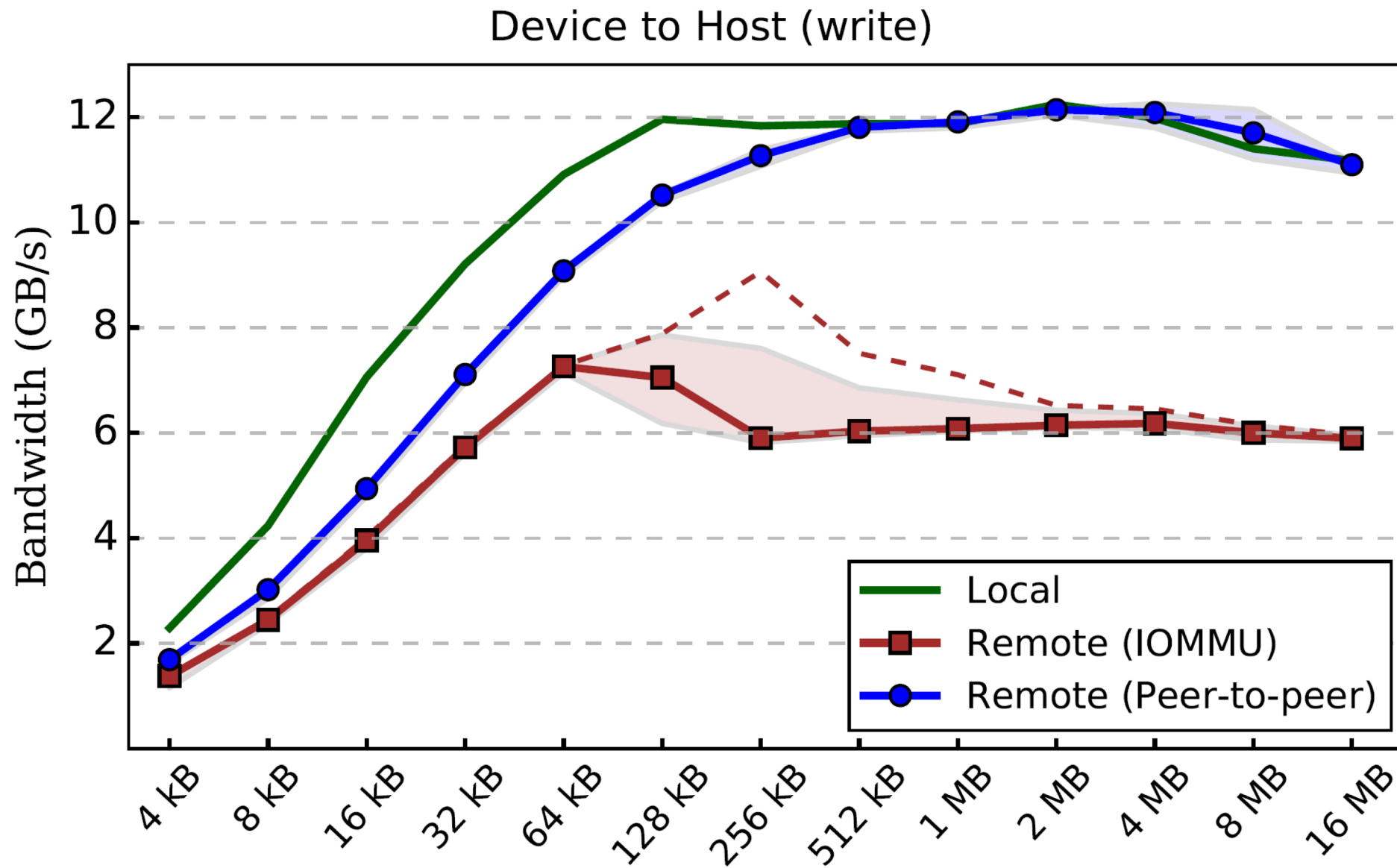- Message-Signaled Interrupts (MSI-X)

# Non-Transparent Bridges

# Remote address space can be mapped into local address space by using PCIe Non-Transparent Bridges (NTBs)
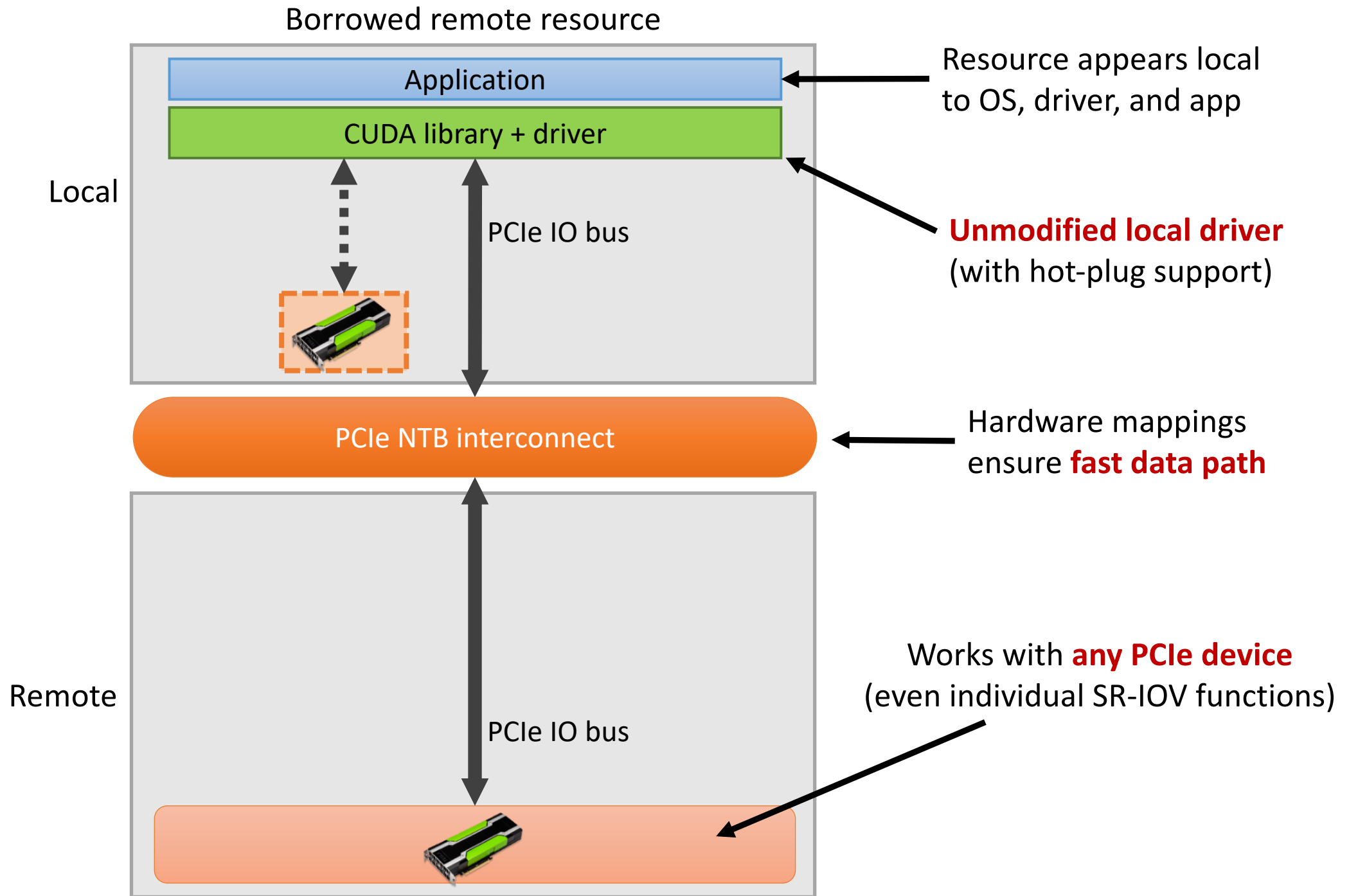
# Using NTBs, each node in the cluster take part in a shared address space and have their own "window" into the global address space
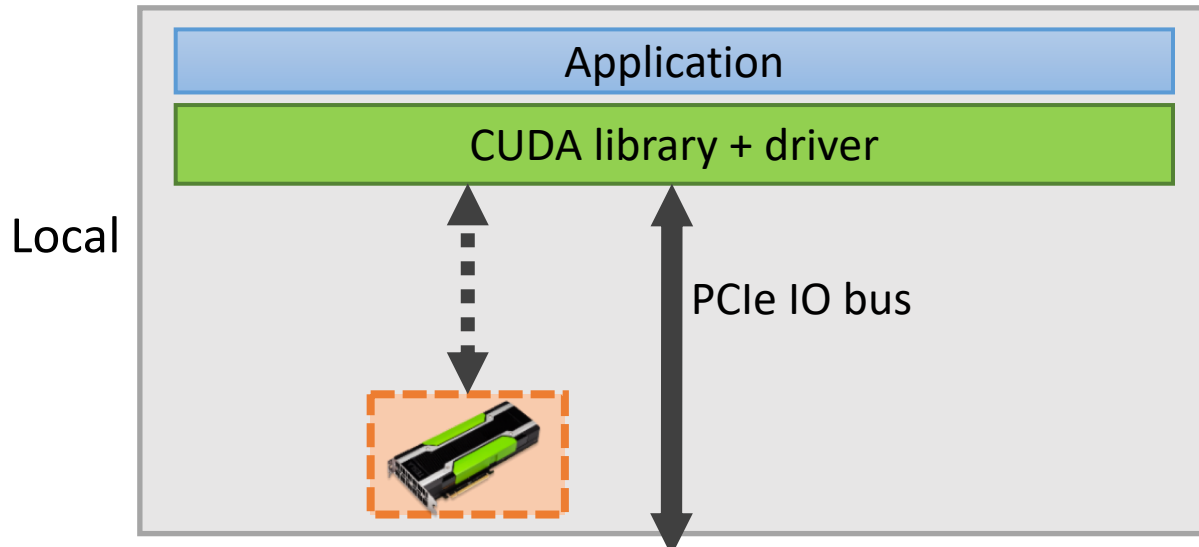


A's addr space

Local IO devices

Global addr space

Local RAM

Global addr space

Addr space in A

Addr space in B

Addr space in C

A    B    C

NTB-based interconnect

C's addr space

Local IO devices

Exported address range

Local RAM
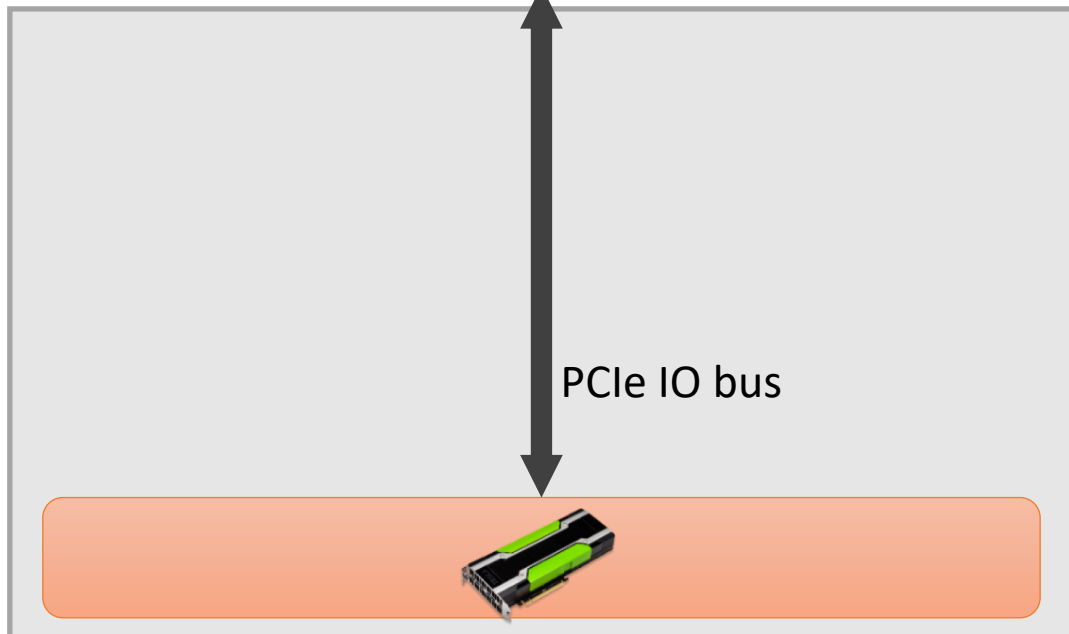
# Device to host transfers: Comparing local to borrowed GPU



Device to Host (write)

# 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**
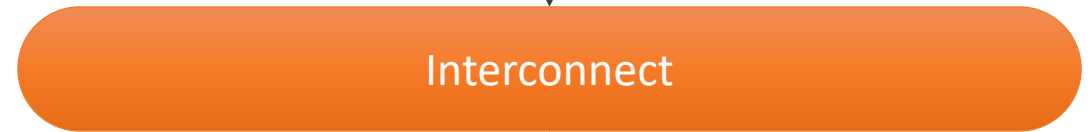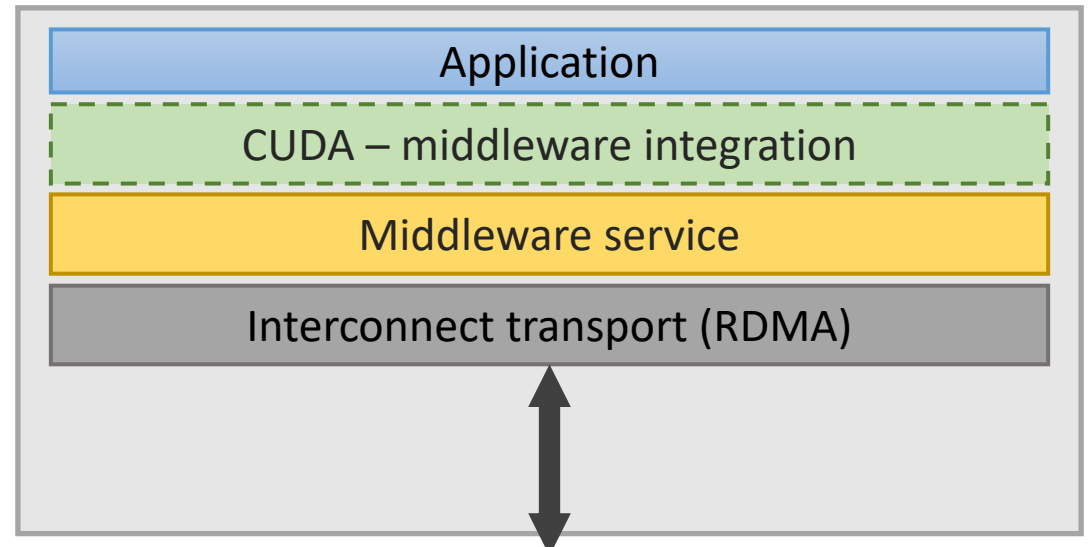
Remote

- PCIe IO bus

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

## Borrowed remote resource

**Local**

| Application |
| CUDA library + driver |

PCIe IO bus

PCIe NTB interconnect

**Remote**

PCIe IO bus

## Remote resource using middleware

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

Interconnect

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

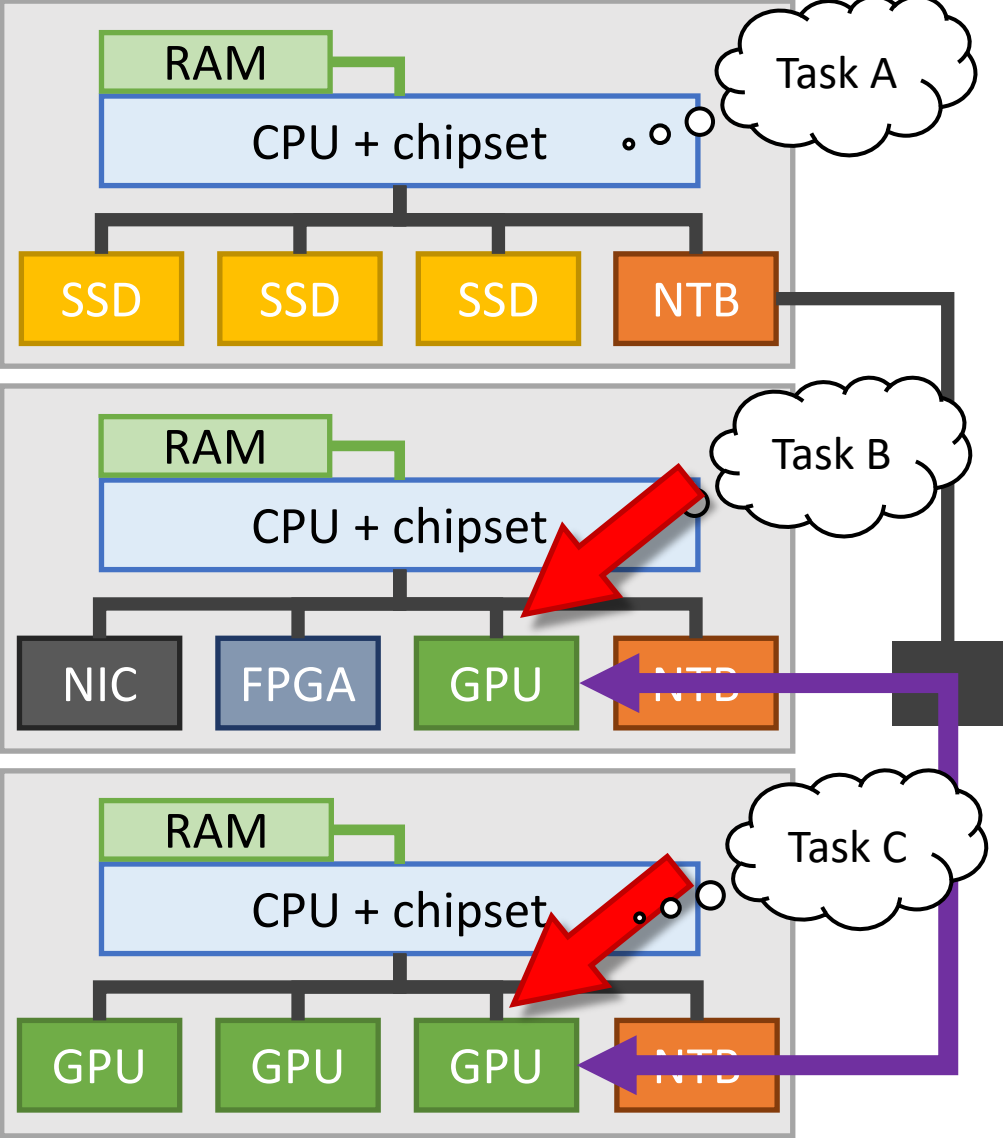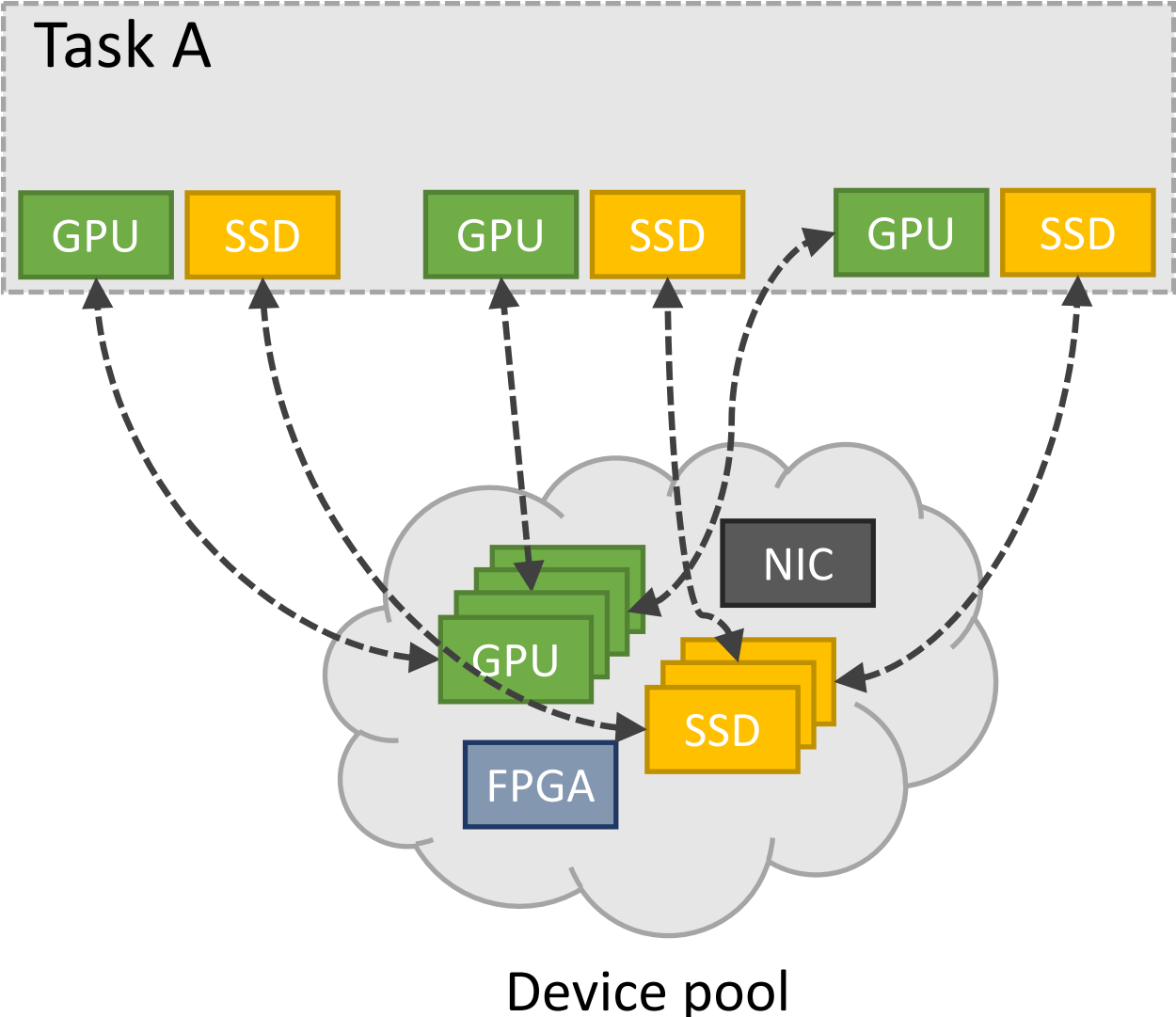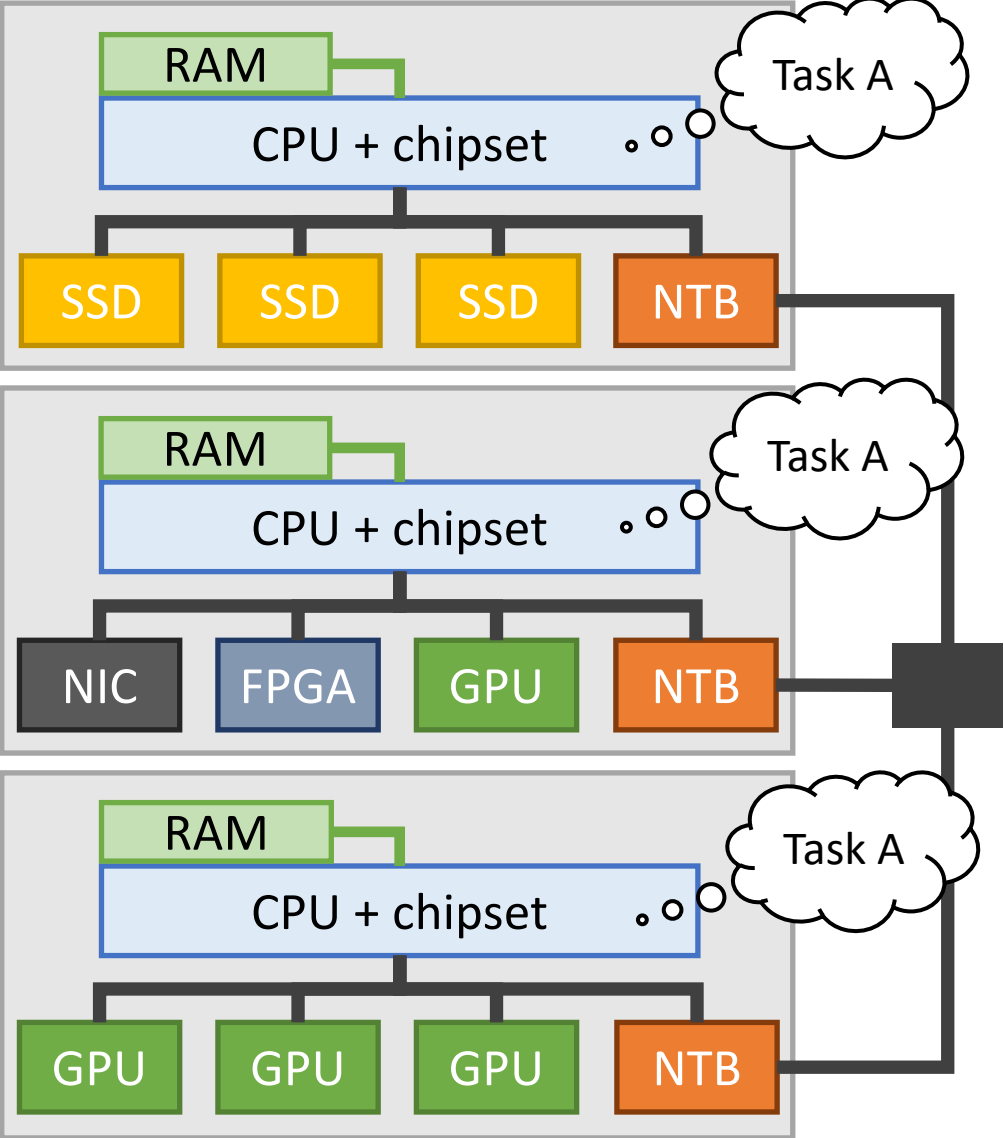PCIe IO bus

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



Peer-to-peer

Task A

Task B

Task C

Device pool

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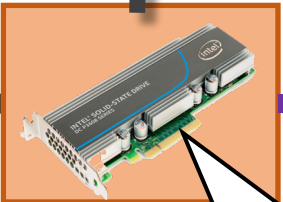


Device pool

# Example: NVMe disk operation (simplified)

Read *N* blocks to address `0x9000`

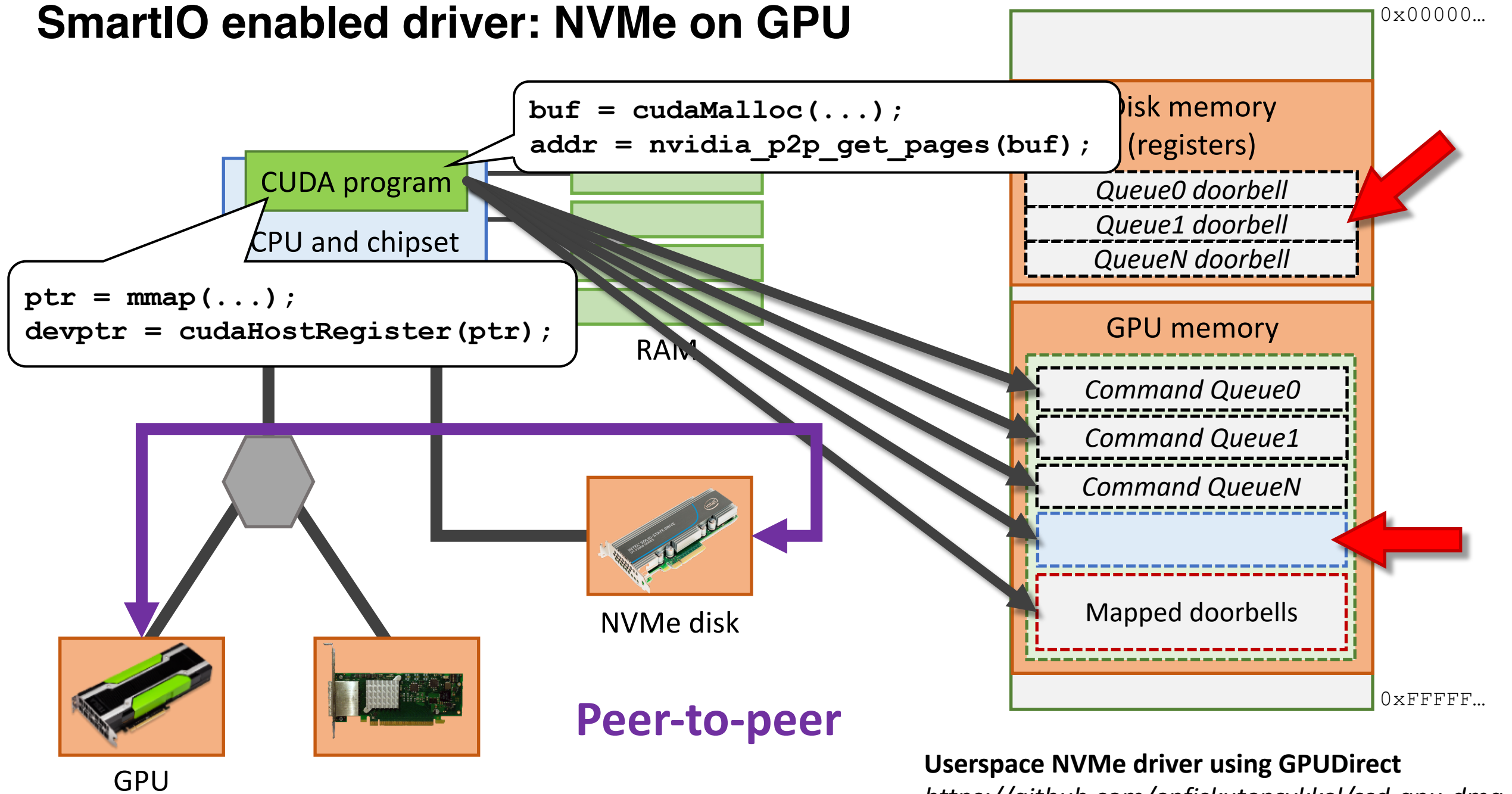NVMe driver

CPU and chipset

RAM

NVMe disk

Command complete

Interrupt vectors

0x00000...

Disk memory (registers)

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

RAM

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

Data

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

NVMe disk

GPU

**Peer-to-peer**

Disk memory (registers)

*Queue0 doorbell*

*Queue1 doorbell*

*QueueN doorbell*

GPU memory

*Command Queue0*

*Command Queue1*

*Command QueueN*

Mapped doorbells

0xFFFFF...

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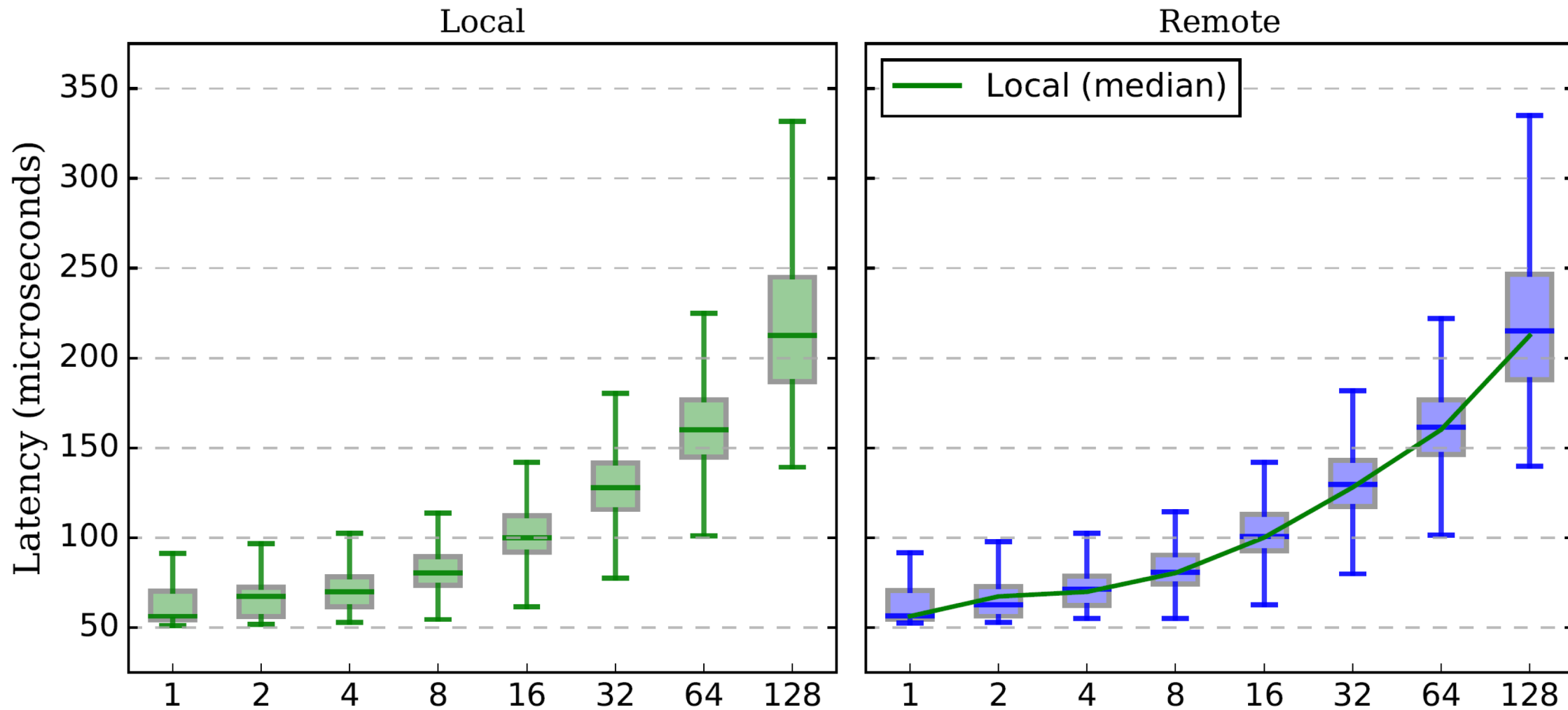*https://github.com/enfiskutensykkel/ssd-gpu-dma*

# Example: NVMe queues hosted remotely

# Read latency for reading blocks from a NVMe disk into a GPU: Local versus borrowed disk



Local vs. Remote NVMe Read Performance

# 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)

*"PCIe Device Lending"*
University of Oslo 2015

My email address

haakonks@simula.no

**SmartIO & Device Lending demo and more**
Visit Dolphin in the exhibition area (booth 523)