

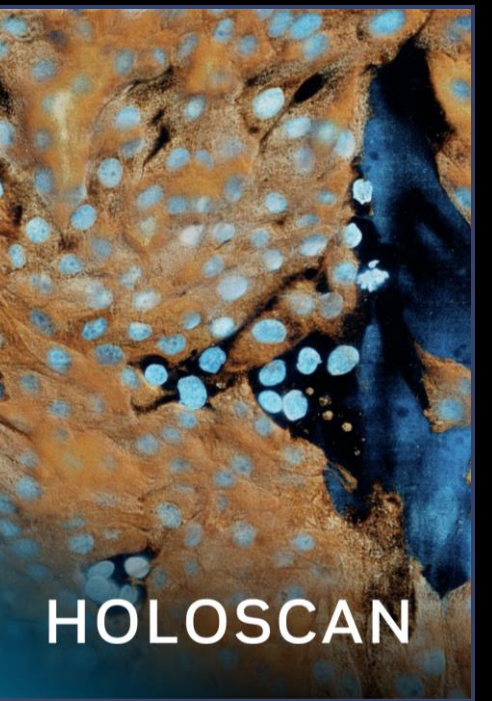
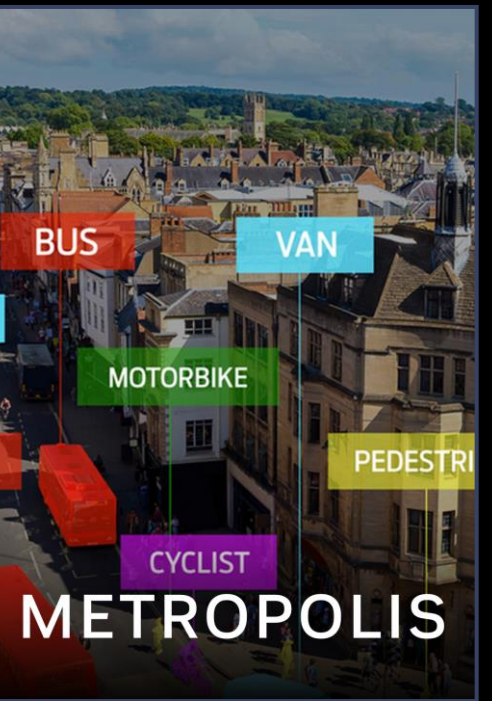
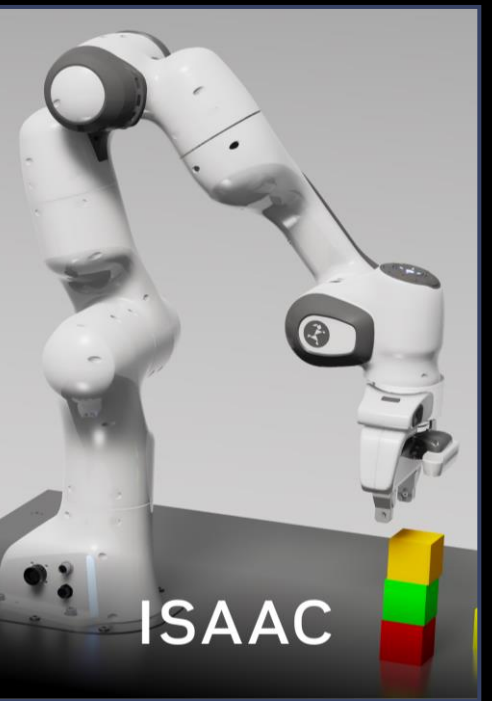
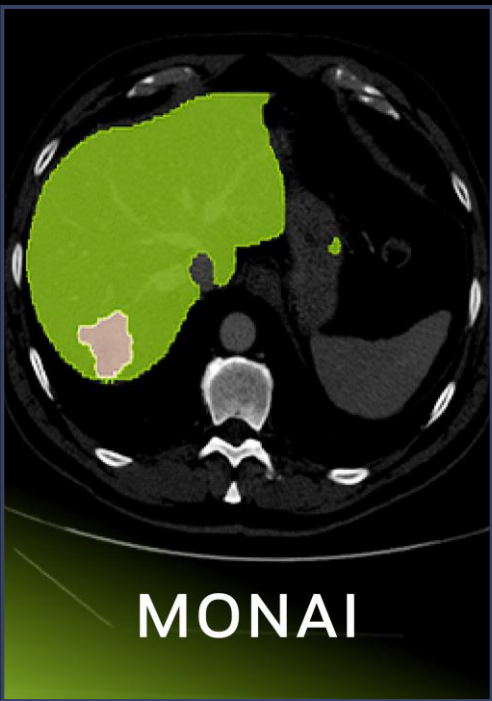
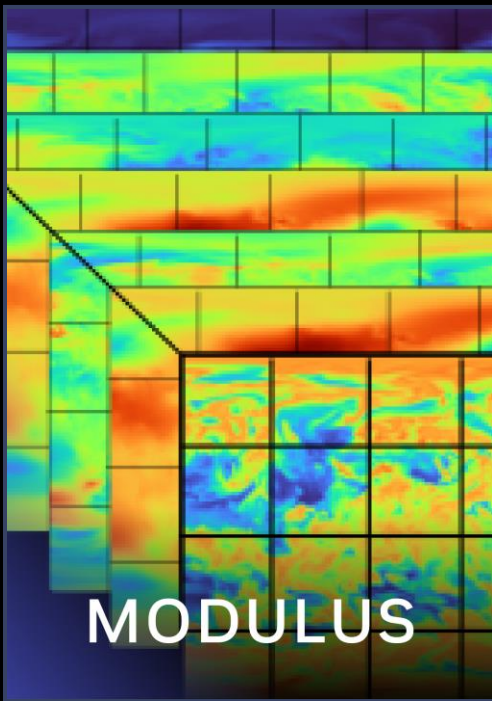


# NVIDIA HPC Networking Platform

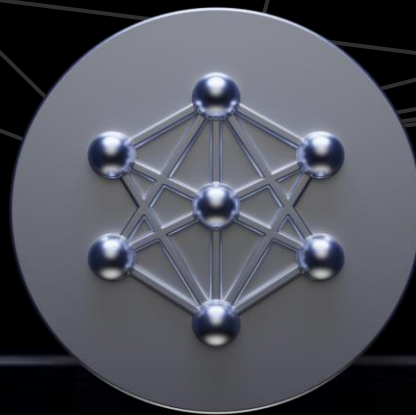
MUG 2023



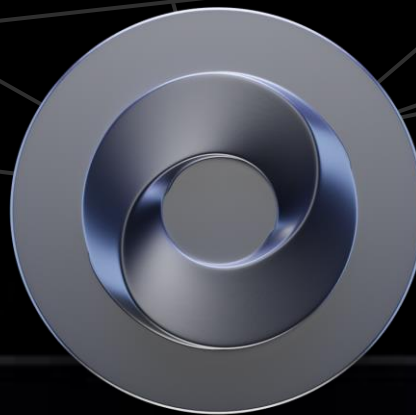
APPLICATION  
FRAMEWORKS



PLATFORM

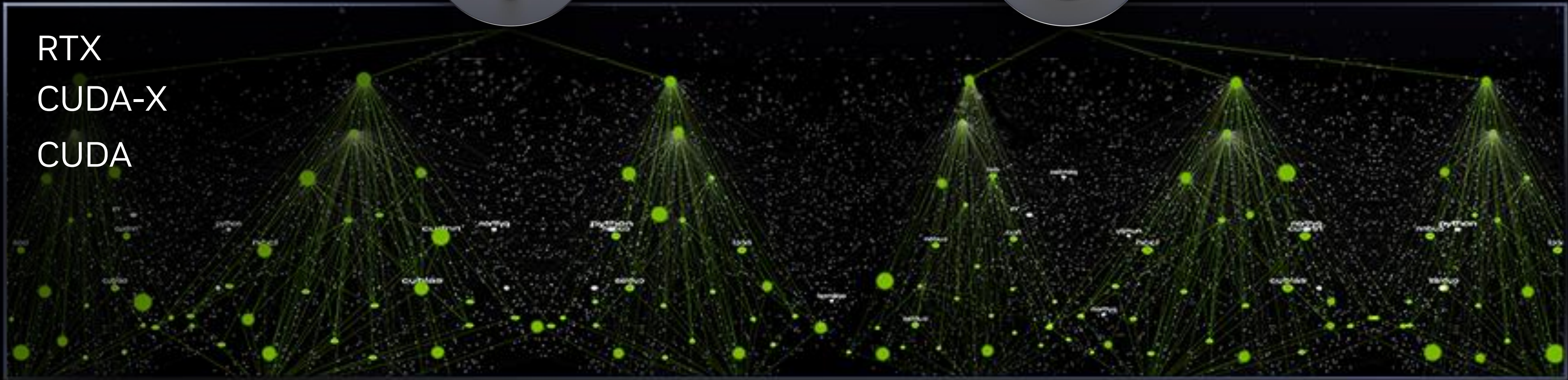


NVIDIA AI

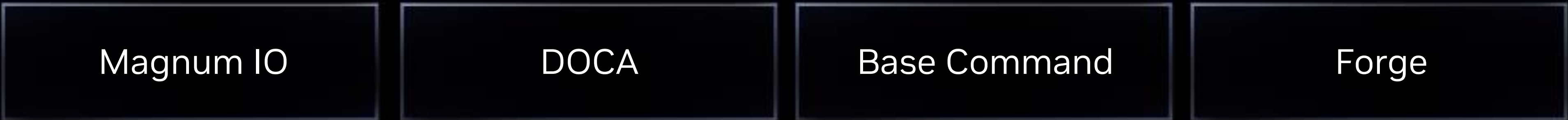


NVIDIA OMNIVERSE

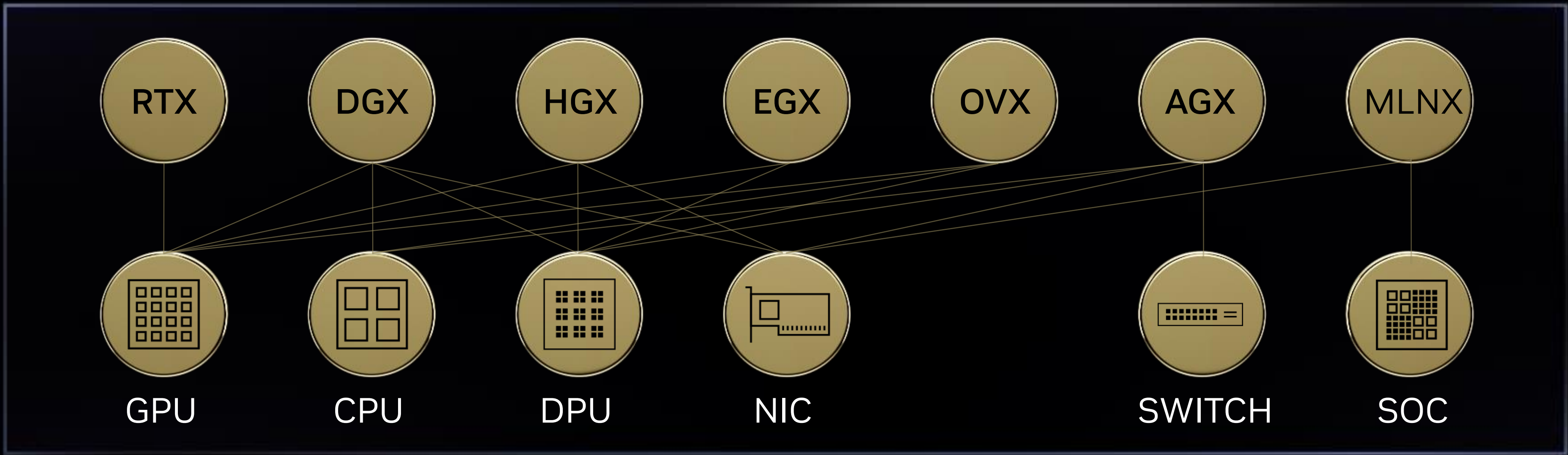
ACCELERATION  
LIBRARIES



SYSTEM  
SOFTWARE

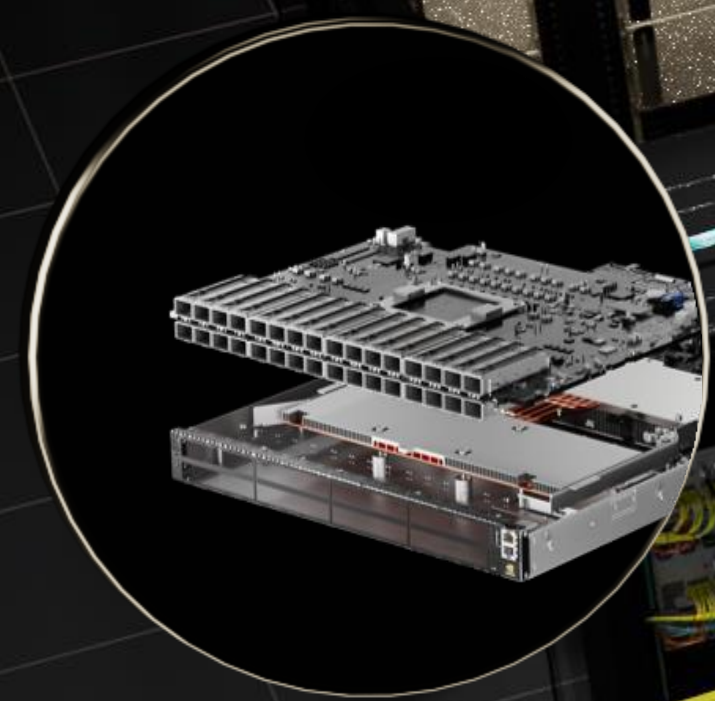


HARDWARE

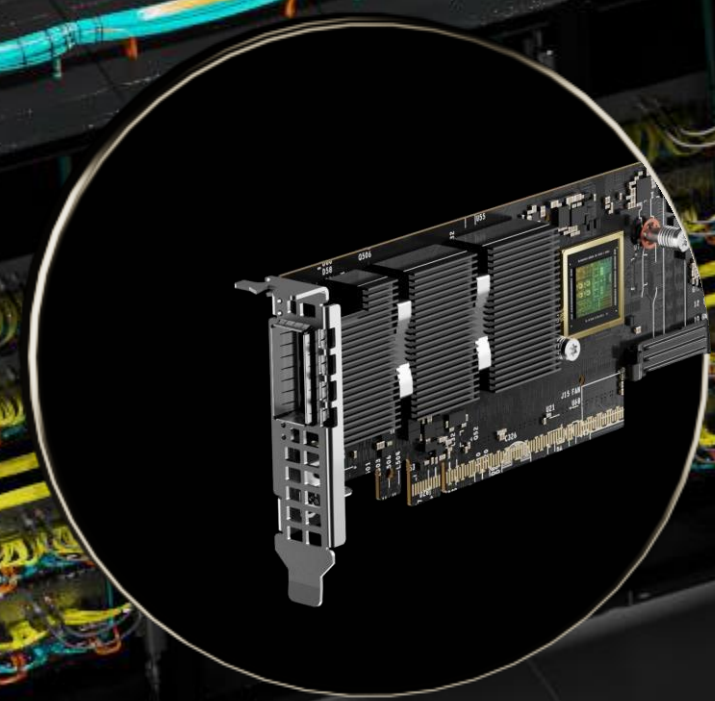




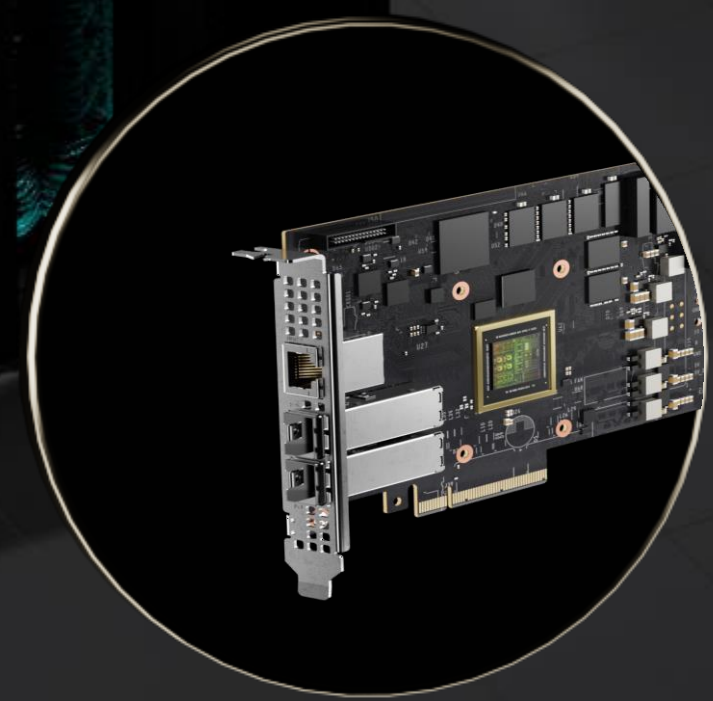
# NVIDIA Networking Platforms



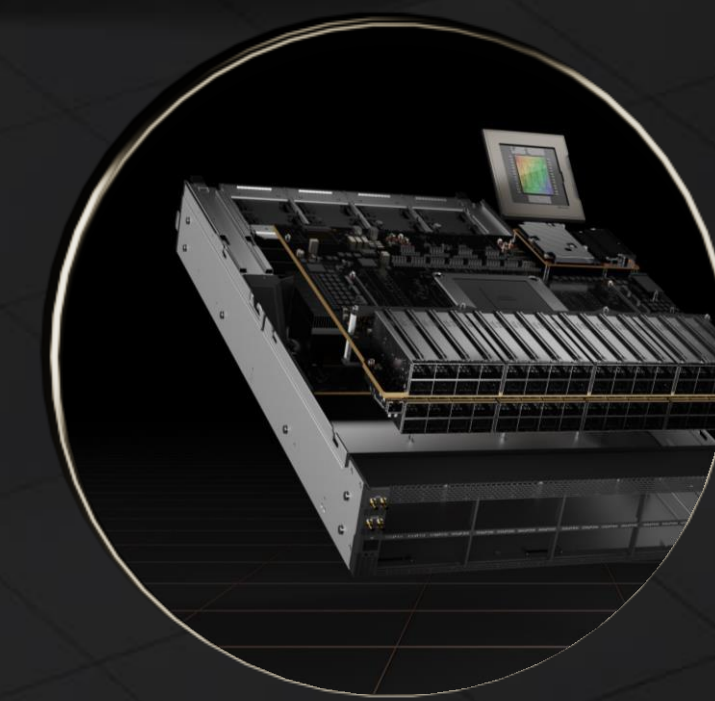
QUANTUM-2 INFINIBAND SWITCH



CONNECTX-7 SMARTNIC



BLUEFIELD-3 DPU



SPECTRUM-4 ETHERNET SWITCH



MANAGEMENT



# In-Network Computing Accelerated Supercomputing

Software-Defined, Hardware-Accelerated, InfiniBand Network

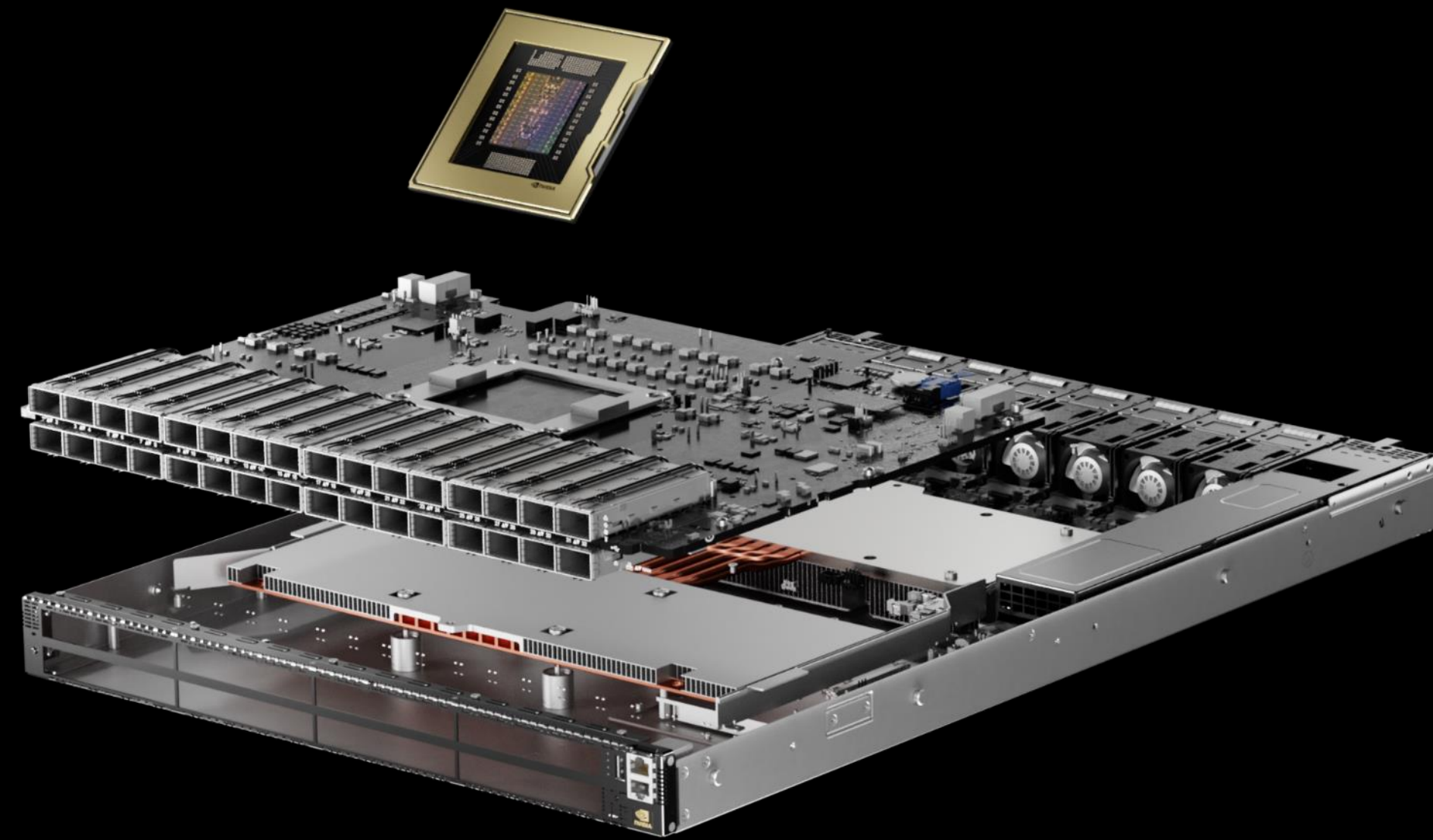
Advanced Networking

End-to-End	High Throughput	Extremely Low Latency	High Message Rate
	RDMA	GPUDirect RDMA	GPUDirect Storage
	Adaptive Routing	Congestion Control	Smart Topologies

In-Network Computing

Adapter/DPU	All-to-All	MPI Tag Matching	Data Reductions (SHARP)	Switch
	Programmable Datapath Accelerator	Data processing units (Arm cores)	Self-Healing Network	
End-to-End	Data security / tenant isolation			End-to-End

# NVIDIA Quantum-2 400G In-Network Computing



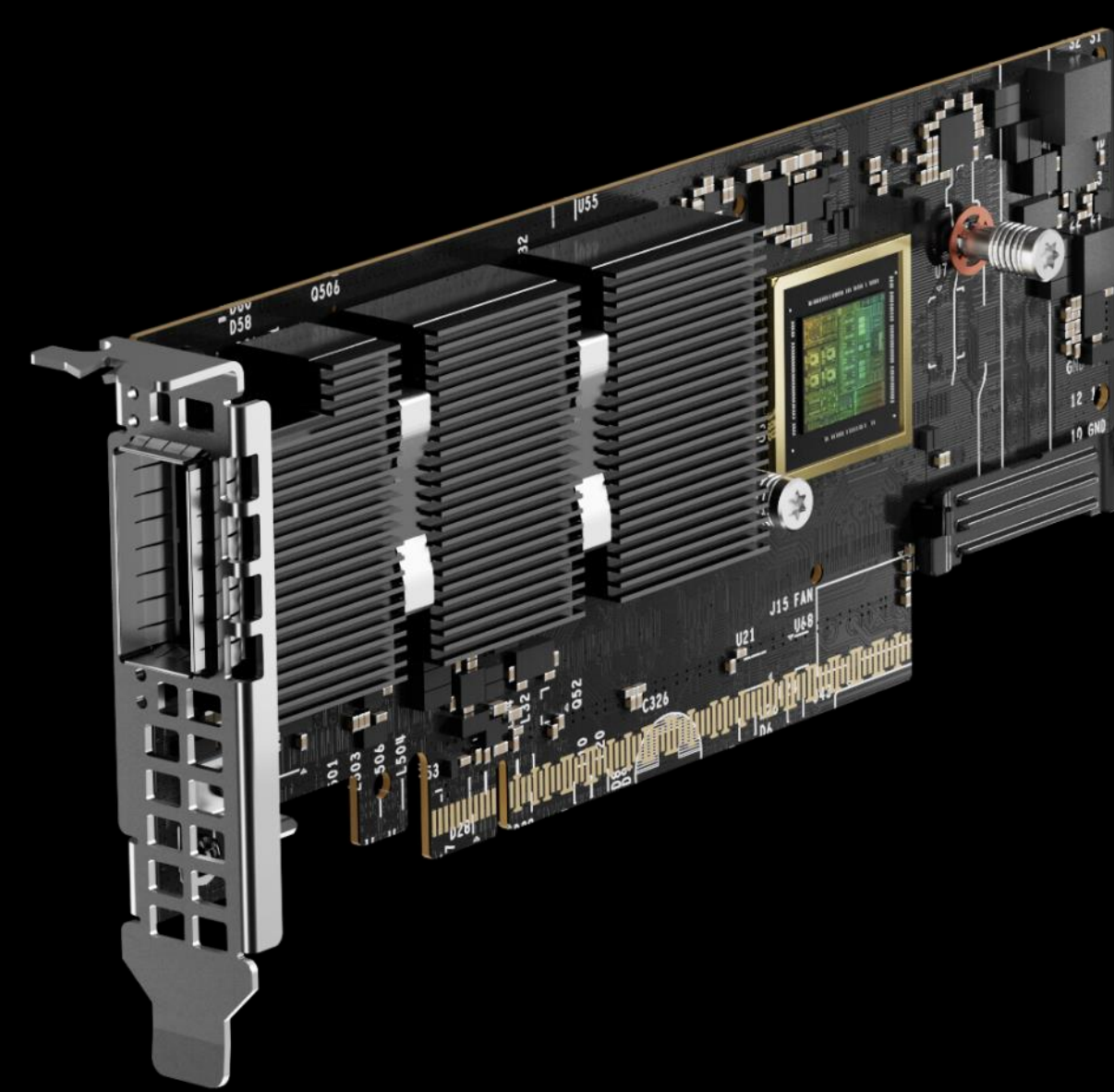
## QUANTUM-2 SWITCH

64-Ports of 400 Gbps or 128-Ports of 200 Gbps

SHARPV3 Small Message Data Reductions

SHARPV3 Large Message Data Reductions

32X More AI Acceleration Engines



## CONNECTX-7 INFINIBAND

16 Core / 256 Threads Datapath Accelerator

Full Transport Offload and Telemetry

Hardware-Based RDMA / GPUDirect

MPI Tag Matching and All-to-All



## BLUEFIELD-3 INFINIBAND

16 Arm 64-Bit Cores

16 Core / 256 Threads Datapath Accelerator

Full Transport Offload and Telemetry

Hardware-Based RDMA / GPUDirect

MPI and NCCL Accelerations

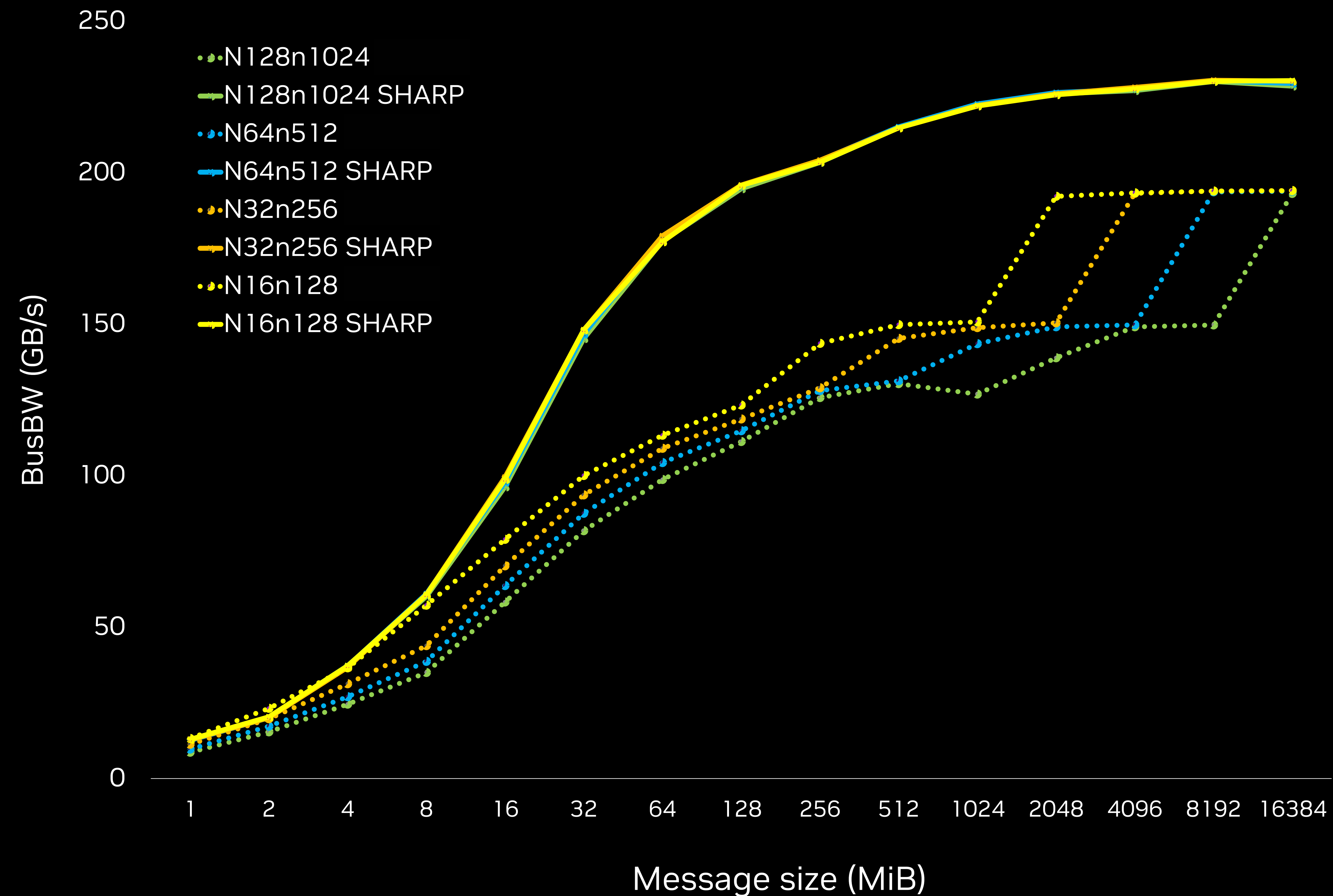
Computational Storage

Security Engines



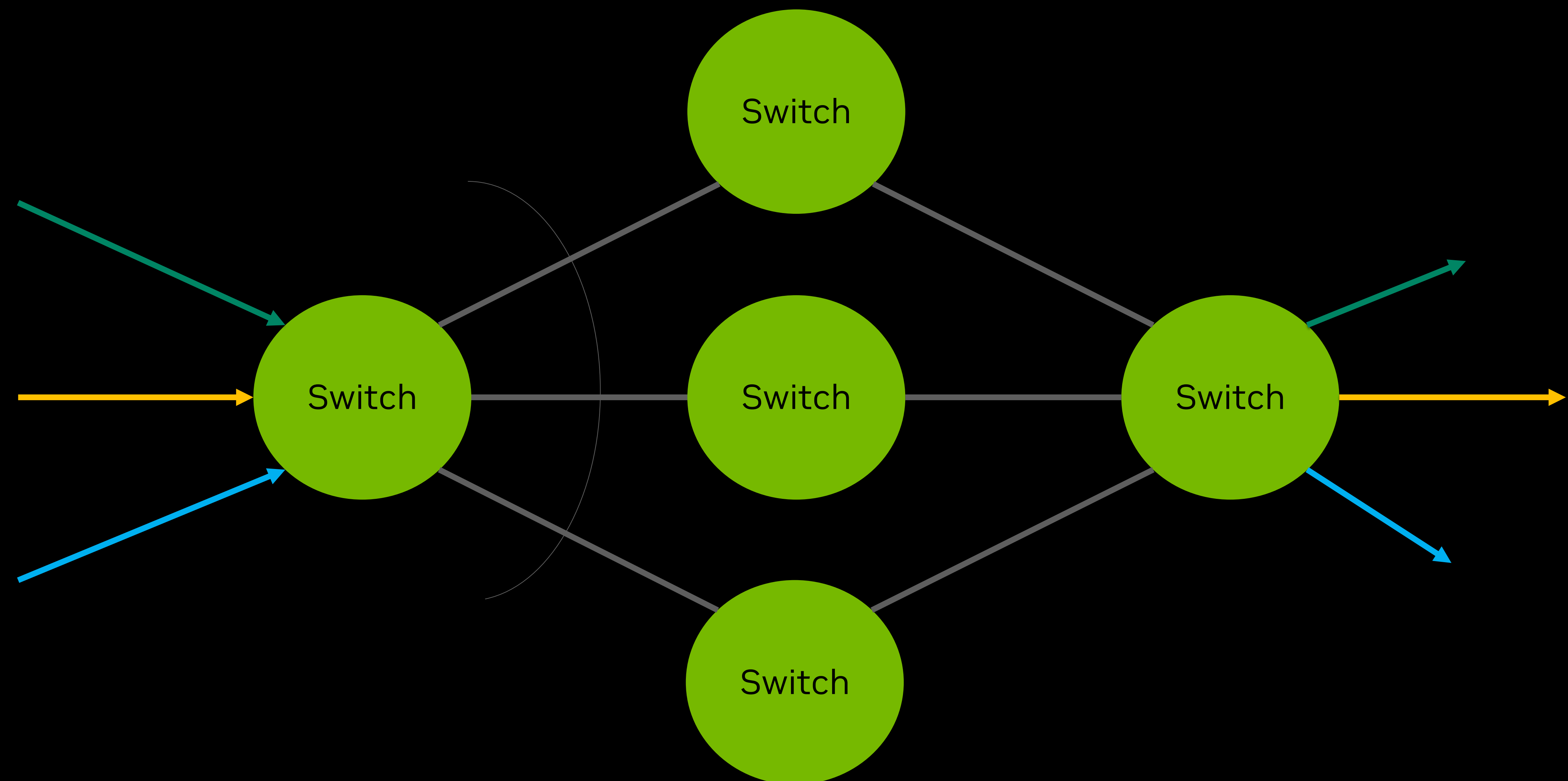
# NCCL All Reduce Performance with InfiniBand SHARP

InfiniBand SHARP maintains high bandwidth on large scale reaching up to 2x advantage



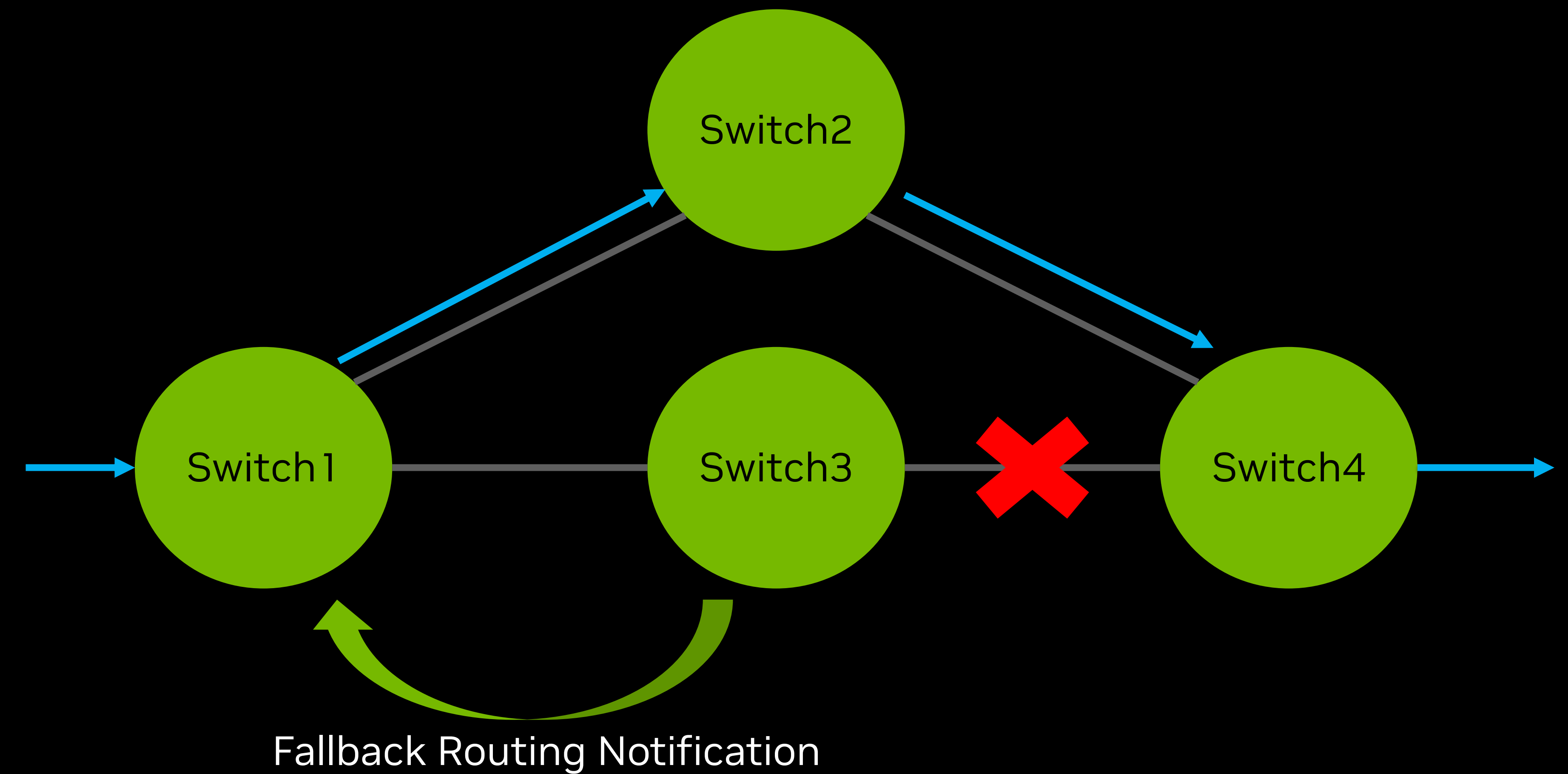
# Adaptive Routing

- Adaptive Routing is the switch's capability to:
  - Dynamically select between available paths to the destination
  - Using pre-configured options (by the SM) and real time data from the output queue state
- AR allows better utilization of the network resources available
- In-order data placement is done by the NIC/DPU (no buffer copy)



# SHIELD

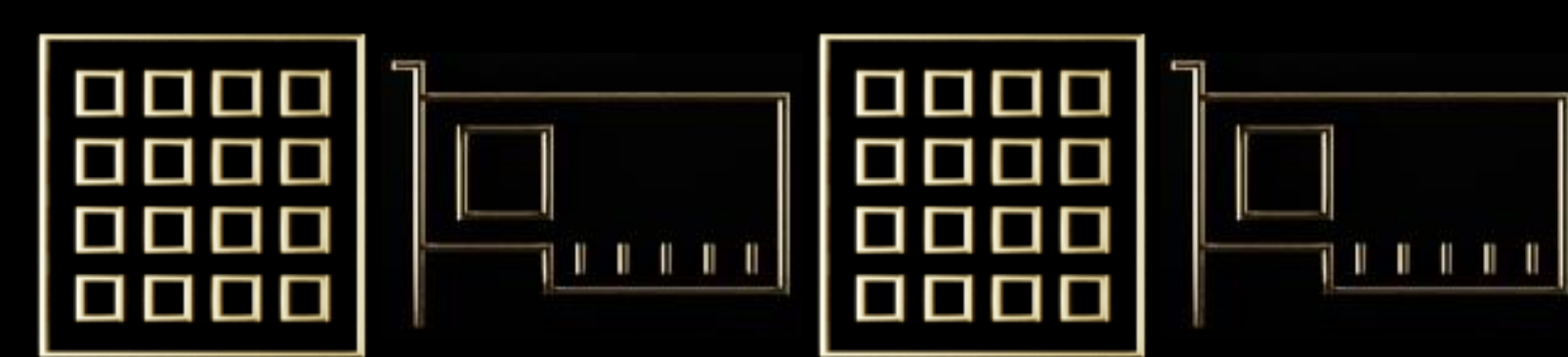
- Self Healing Interconnect Enhancement for Intelligent Datacenters (SHIELD) allows the switch to:
  - Dynamically bypass failed link
  - Using pre-configured egress port (by the SM) and local link status information
  - Link failures propagate from switch to switch using Proactive Failed Route Notification (PFRN) packets
- SHIELD allows faster recovery of connectivity, losing less network utilization
  - Considering the complexity of a modern Data Center, replacing a failed link is no easy feat
  - The time for the SM to bypass a failed link (reconfiguring all switches involved) is considerably more than doing it via HW propagated packets





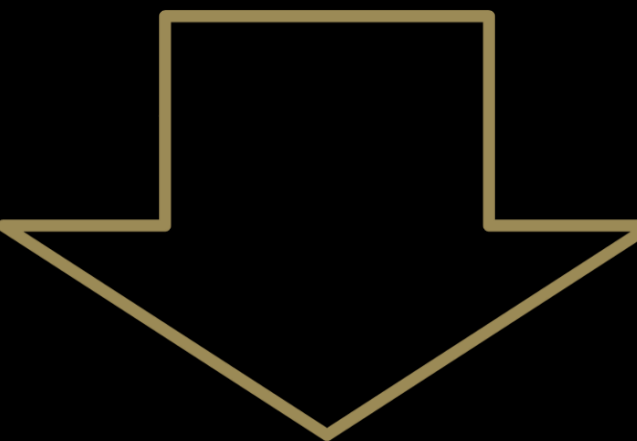
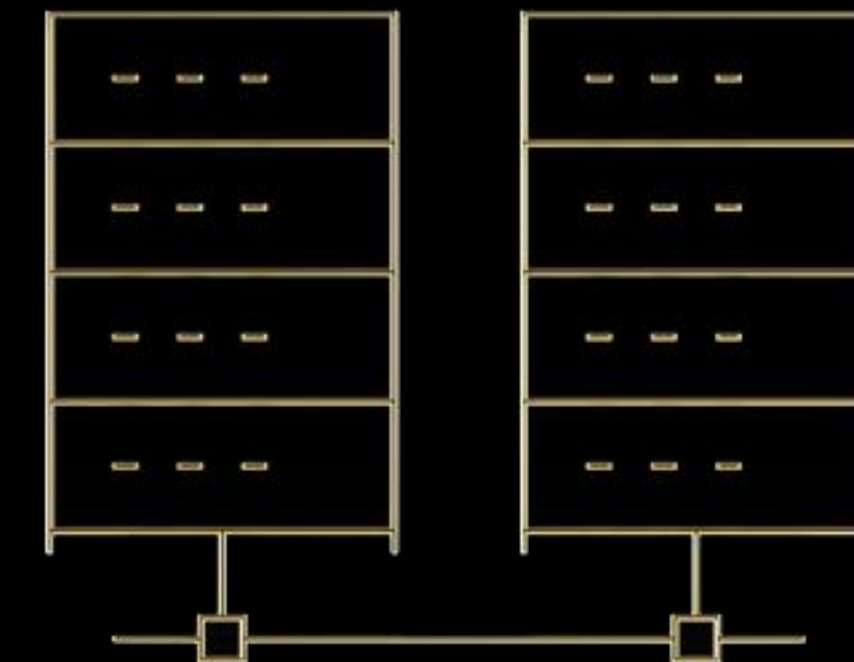
# In-Network Computing to Solve Performance Bottlenecks

Overlapping



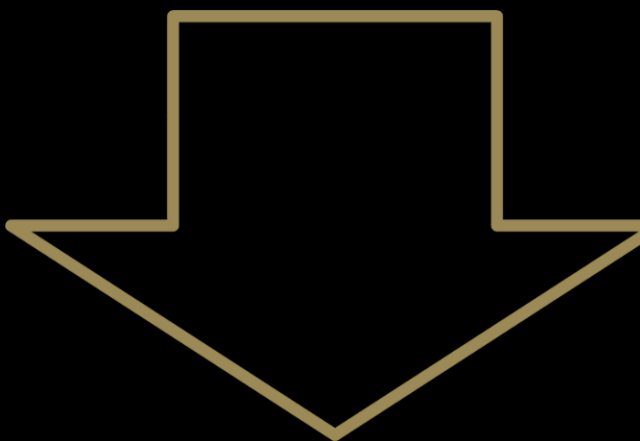
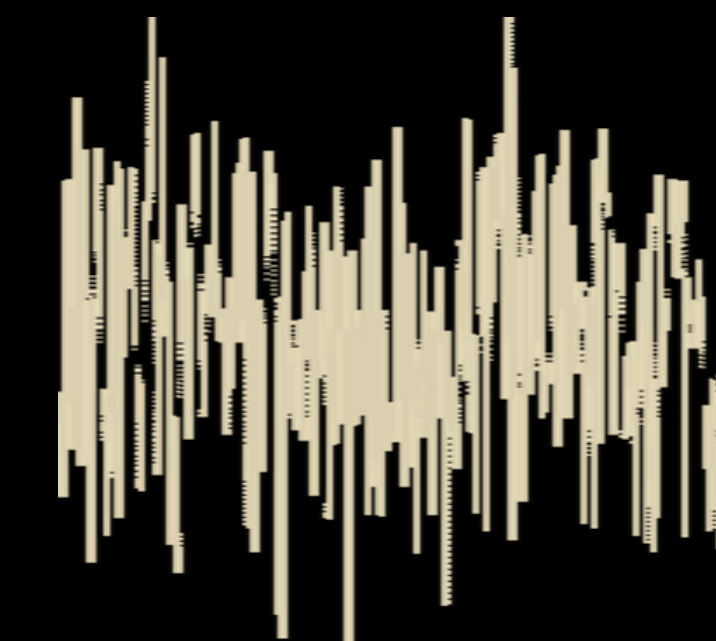
In-Network Computing  
Asynchronous Progress  
(Compute – Communication Overlap)

Load Imbalanced



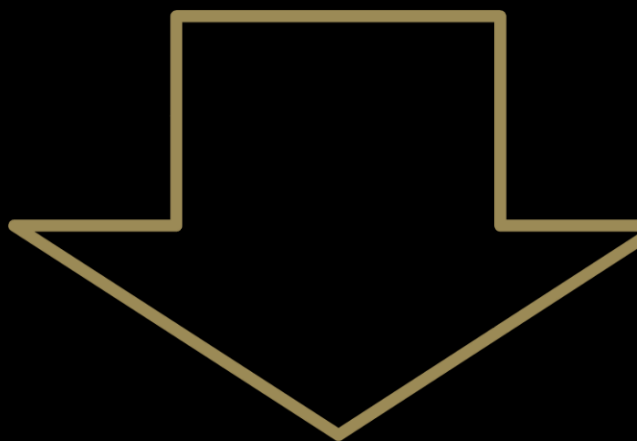
In-Network Computing  
and DPU Synchronization

Jitter



In-Network Computing  
Infrastructure Processing

Multi-Job Performance

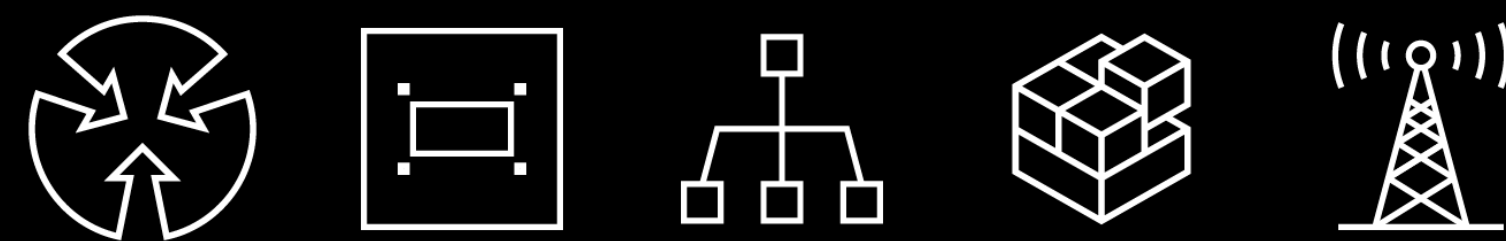


Adaptive Performance  
Isolation



# BlueField Data Processing Unit

## SOFTWARE DEFINED NETWORKING



## SOFTWARE DEFINED SECURITY



## SOFTWARE DEFINED STORAGE



## Infrastructure Services

### Data Center on a Chip

16 Arm 64-Bit Cores

16 Core / 256 Threads Datapath Accelerator

ConnectX InfiniBand / Ethernet

DDR memory interface

PCIe switch

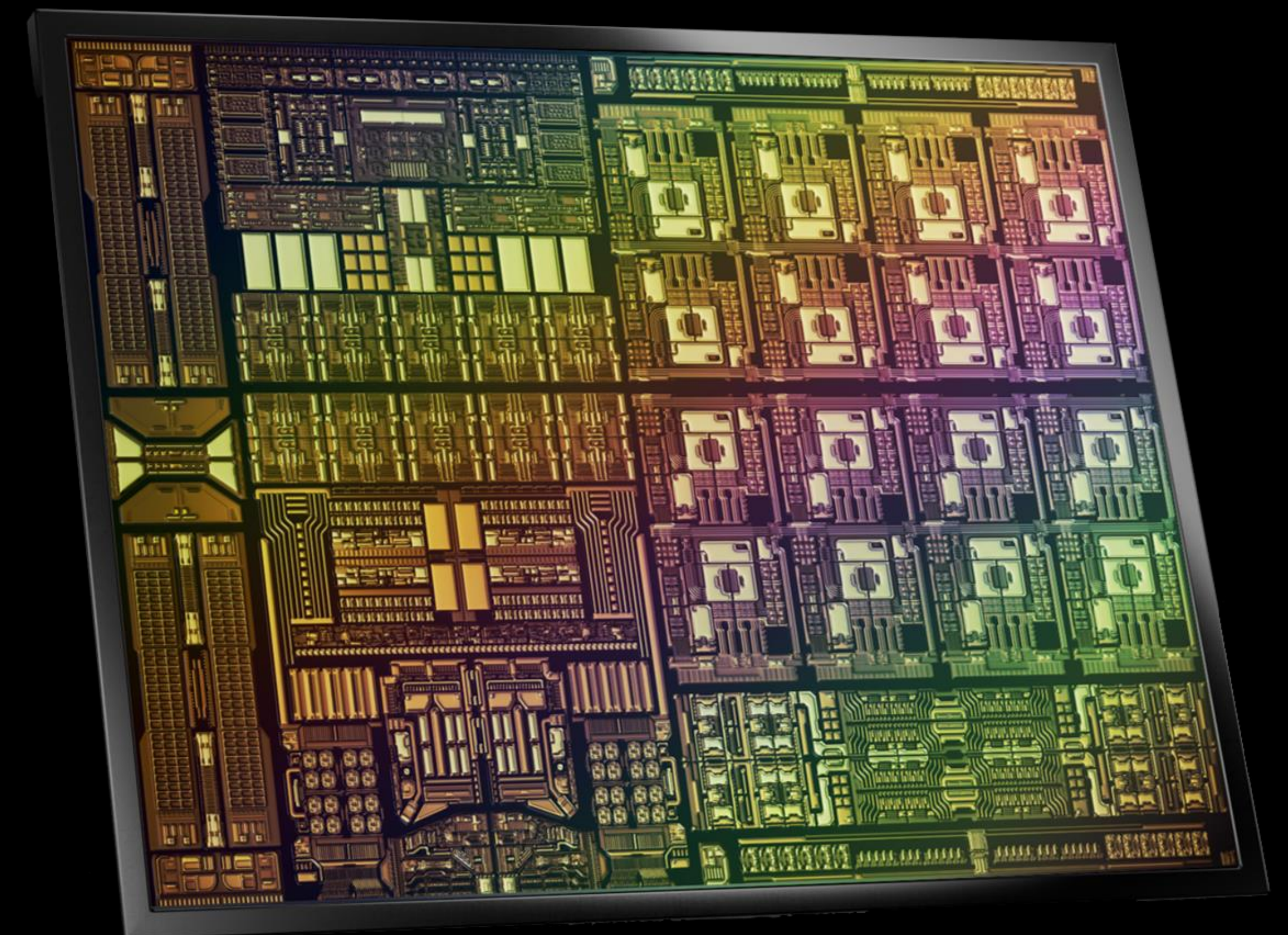
SPECINT2K17: 42

Memory Bandwidth: 80GB/s

NVMe SNAP: 10M IOPS @ 4KB

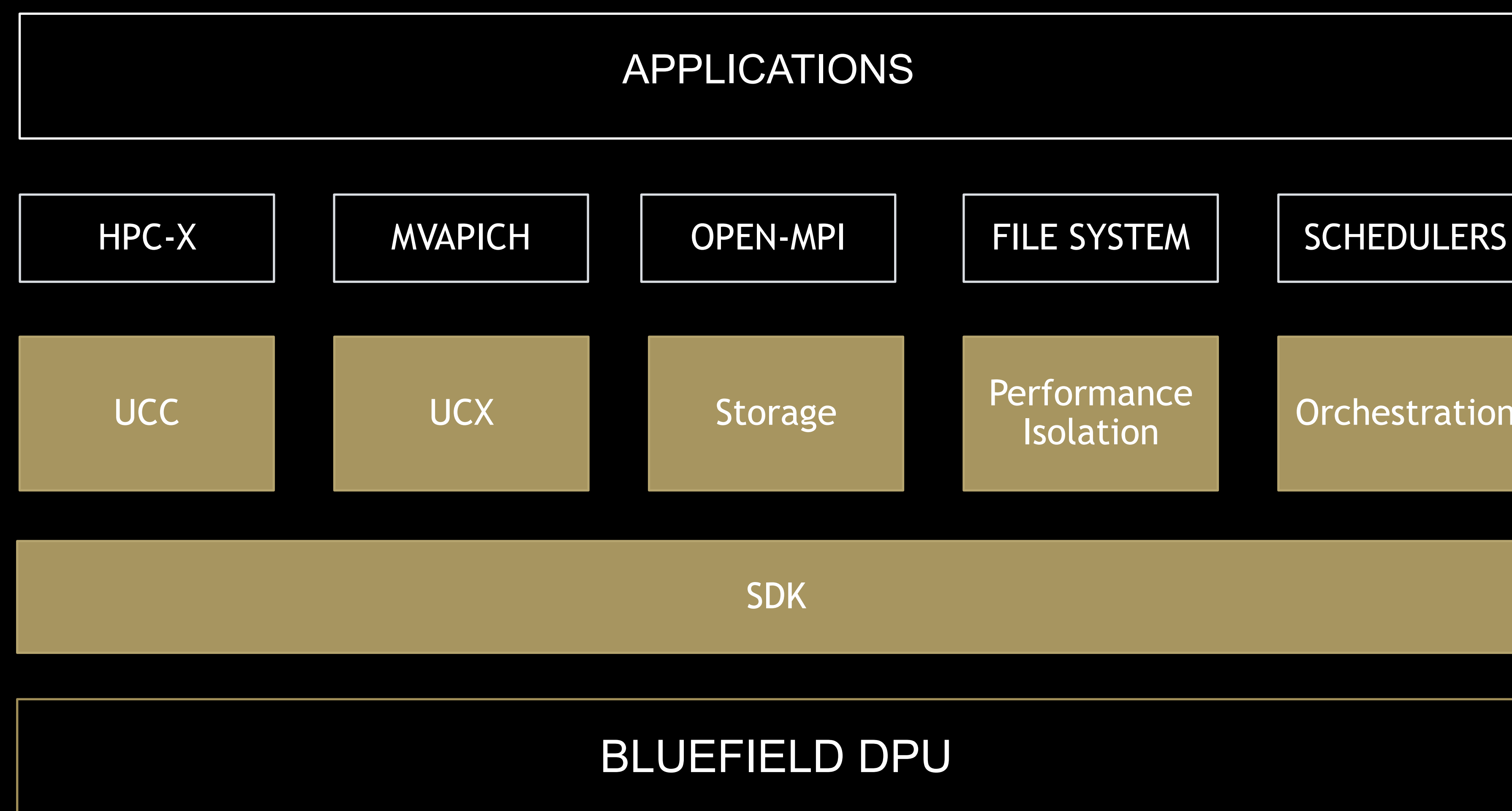


BlueField Infrastructure  
Compute Platform





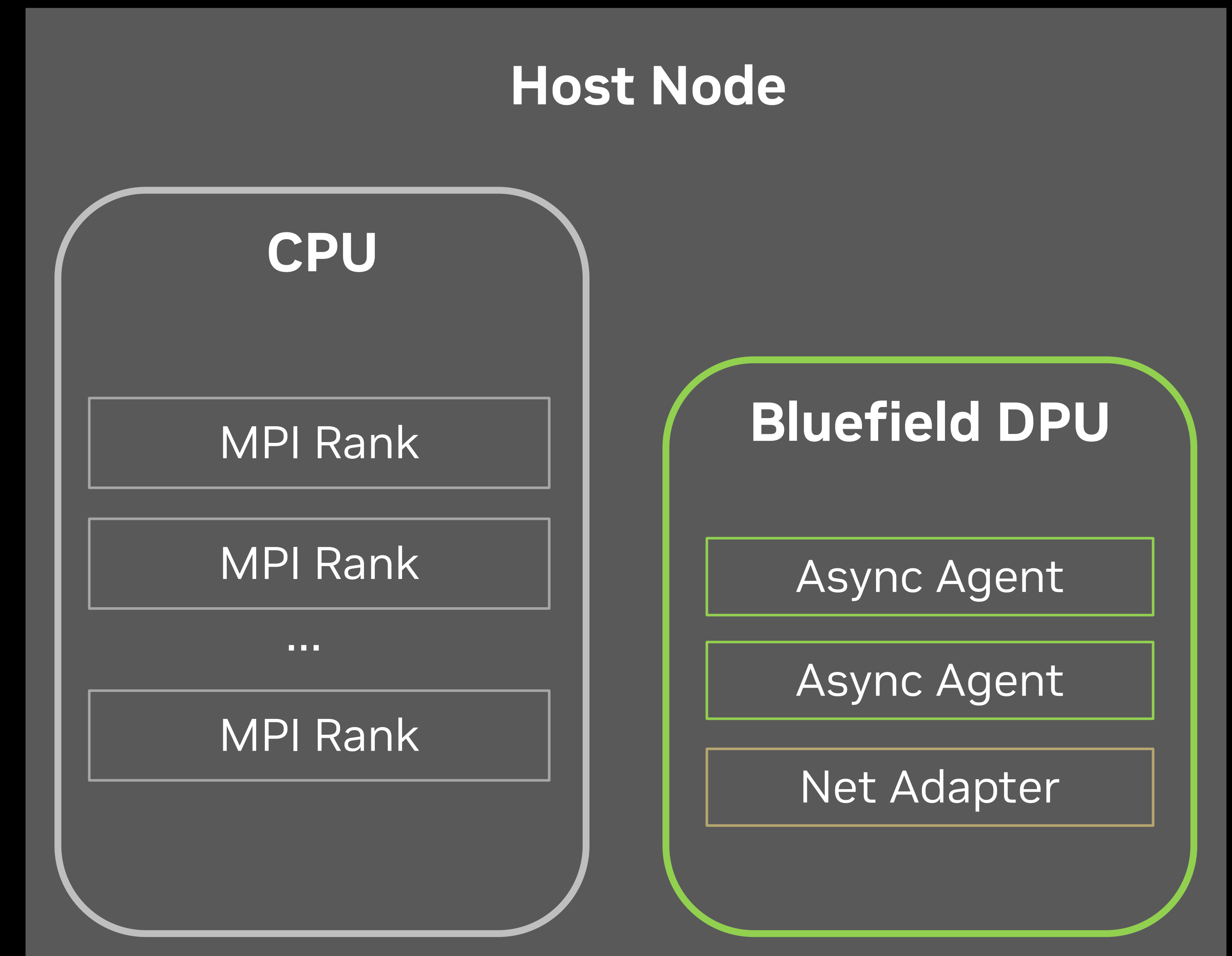
# Accelerating HPC Applications with DPU/DOCA Services





# High Level System Components from Software's Perspective

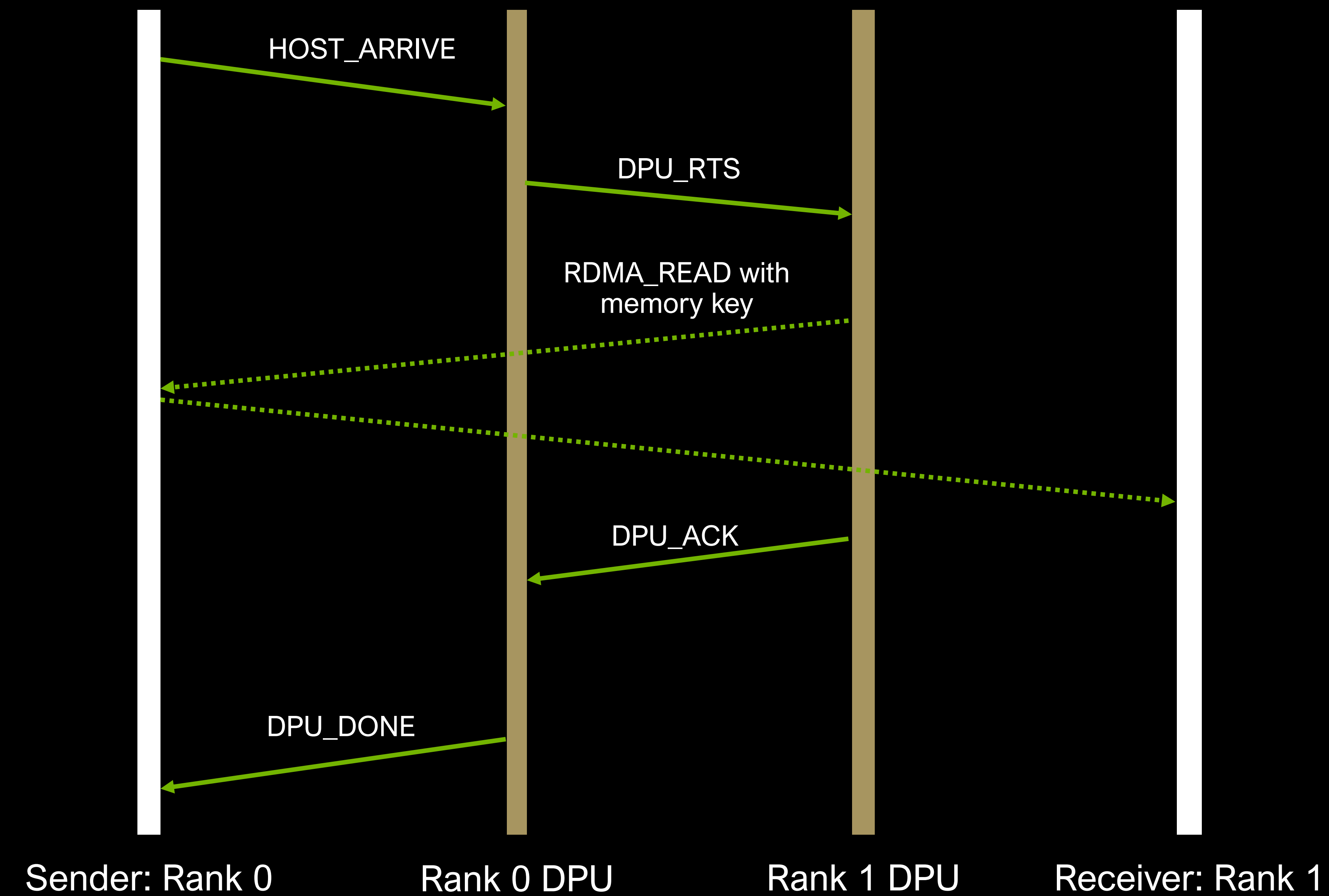
- Host paired with local DPU
- Local DPU runs service processes (SP)
  - Each local user process (such as MPI process) has a service process that it is pair with
  - Each service process serves multiple local processes
  - Algorithm is split between host and DPU
    - Blocking and nonblocking may have different split
- Hosts and SP's may communicate with other hosts and/or SP's
- Cross-GVMI (XGVMI) - The DPU can initiates RDMA operations on behalf of host resident memory
  - DPU memory is involved only if the data originates from or is targeted to DPU memory





# Offloading and Accelerating Data Exchange Example

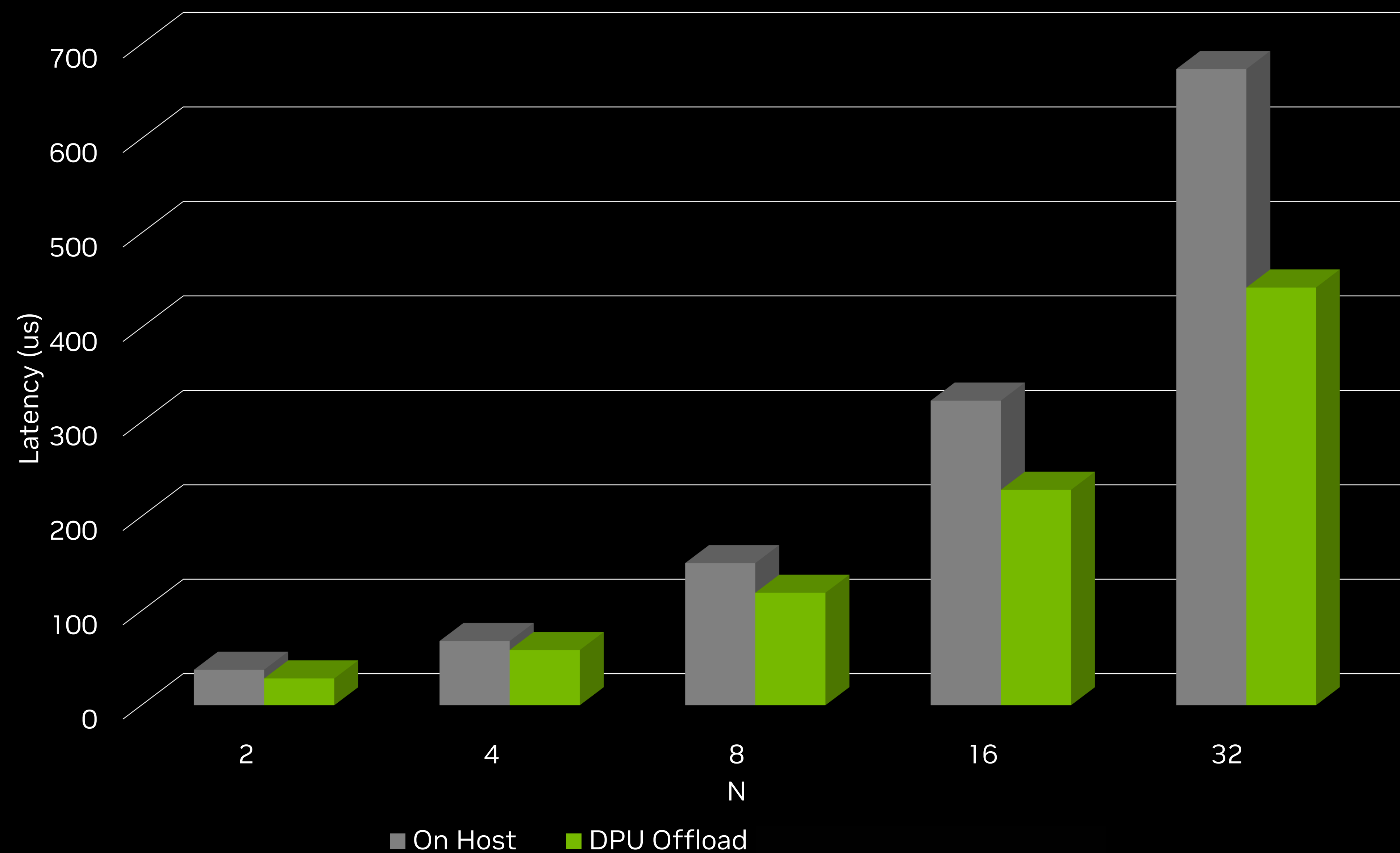
An Element of Collective Algorithm



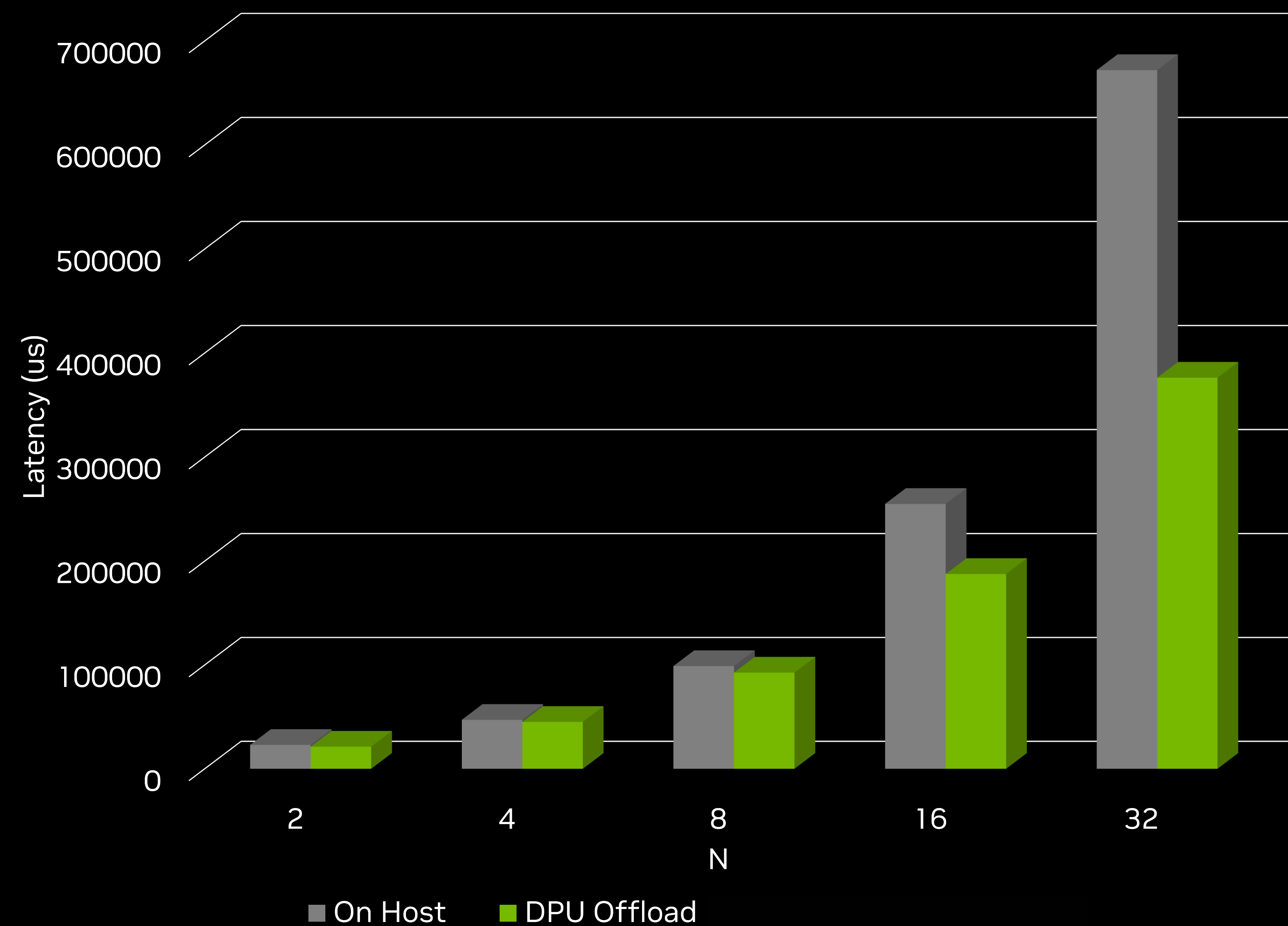


# Alltoallv Latency

OSU Alltoallv 1 PPN, Size = 128 KB



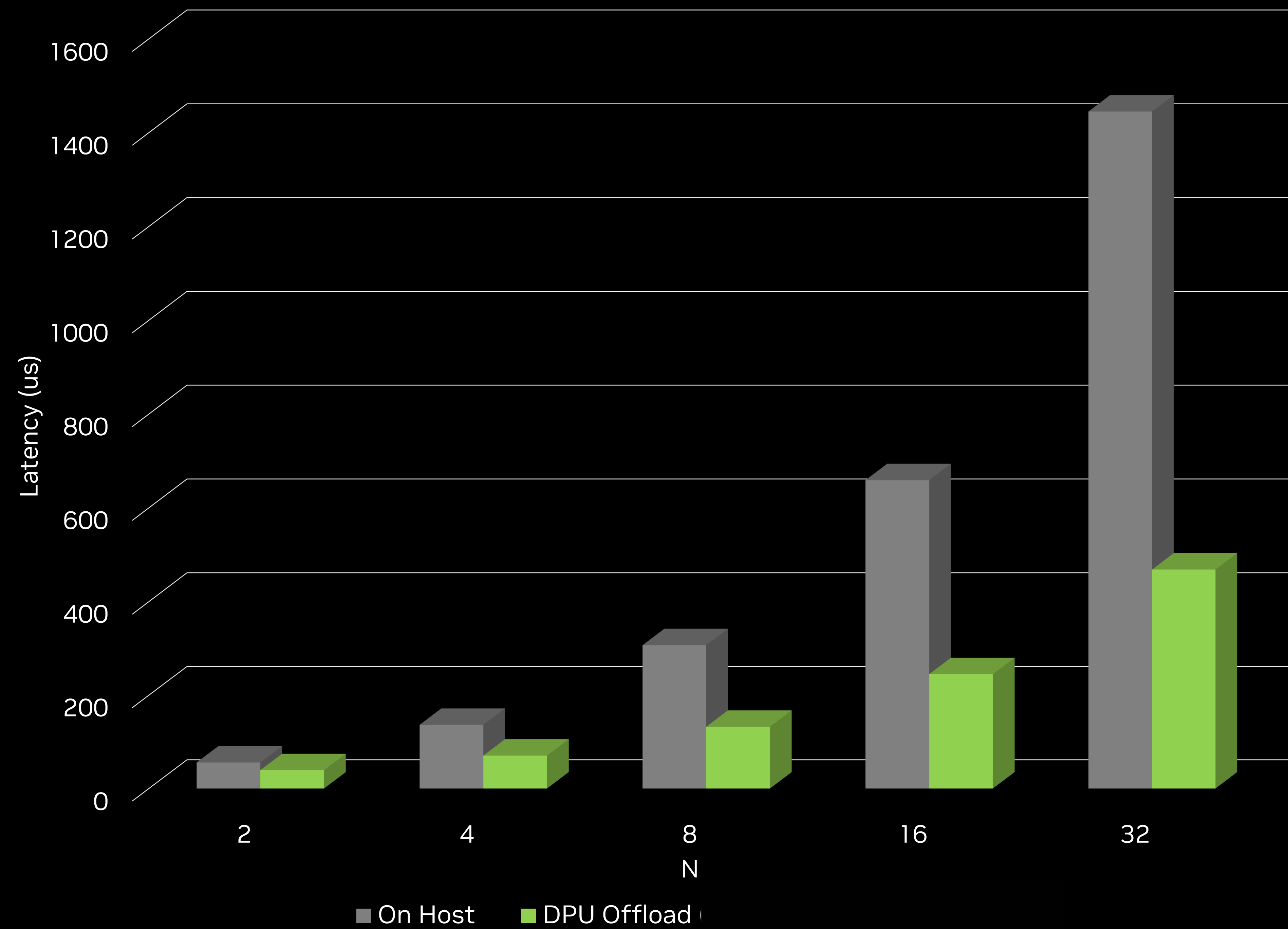
OSU Alltoallv 32 (full) PPN, Size = 128 KB



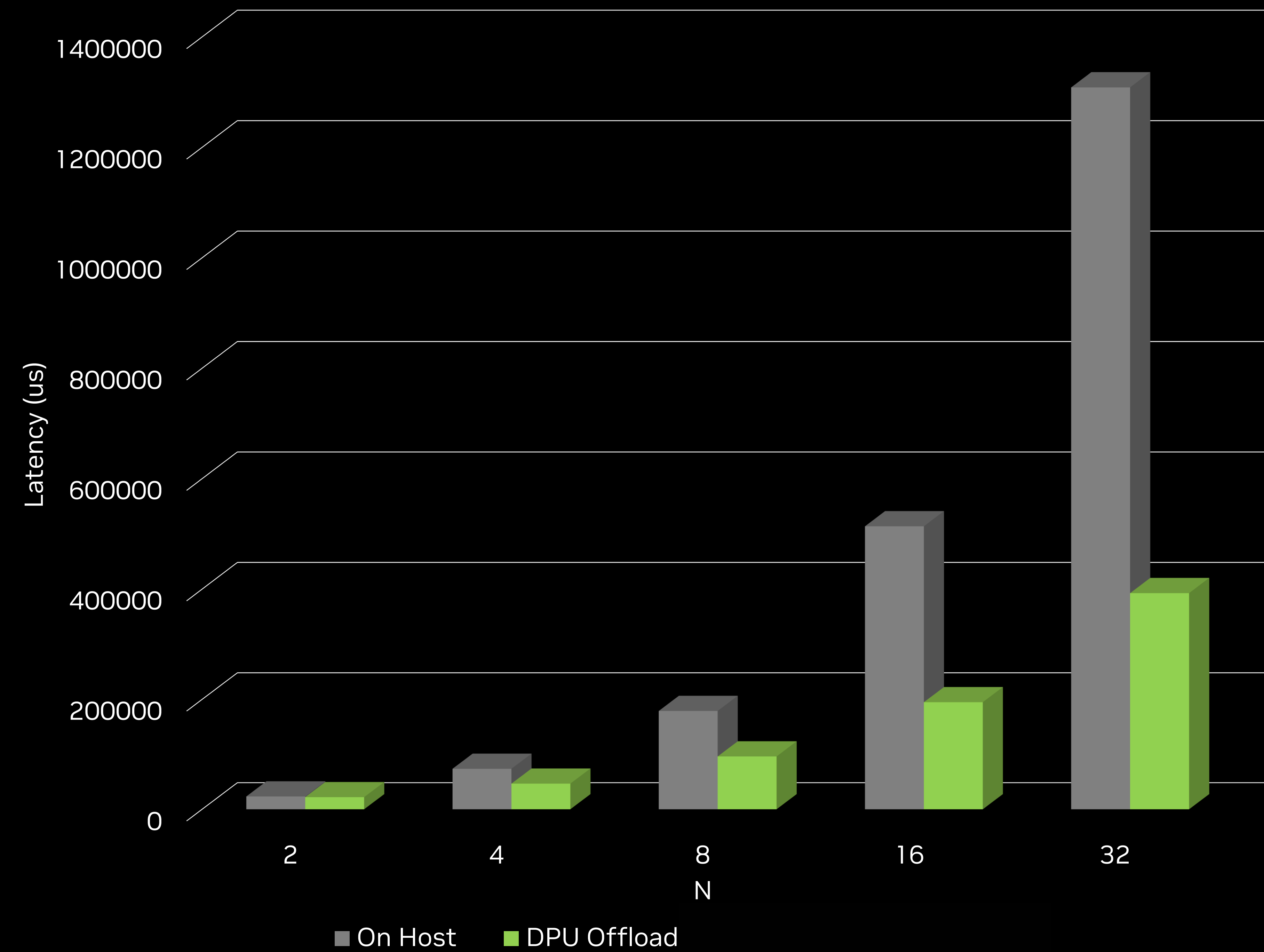


# iAlltoallv latency

OSU iAlltoallv 1 PPN, Size = 128 KB



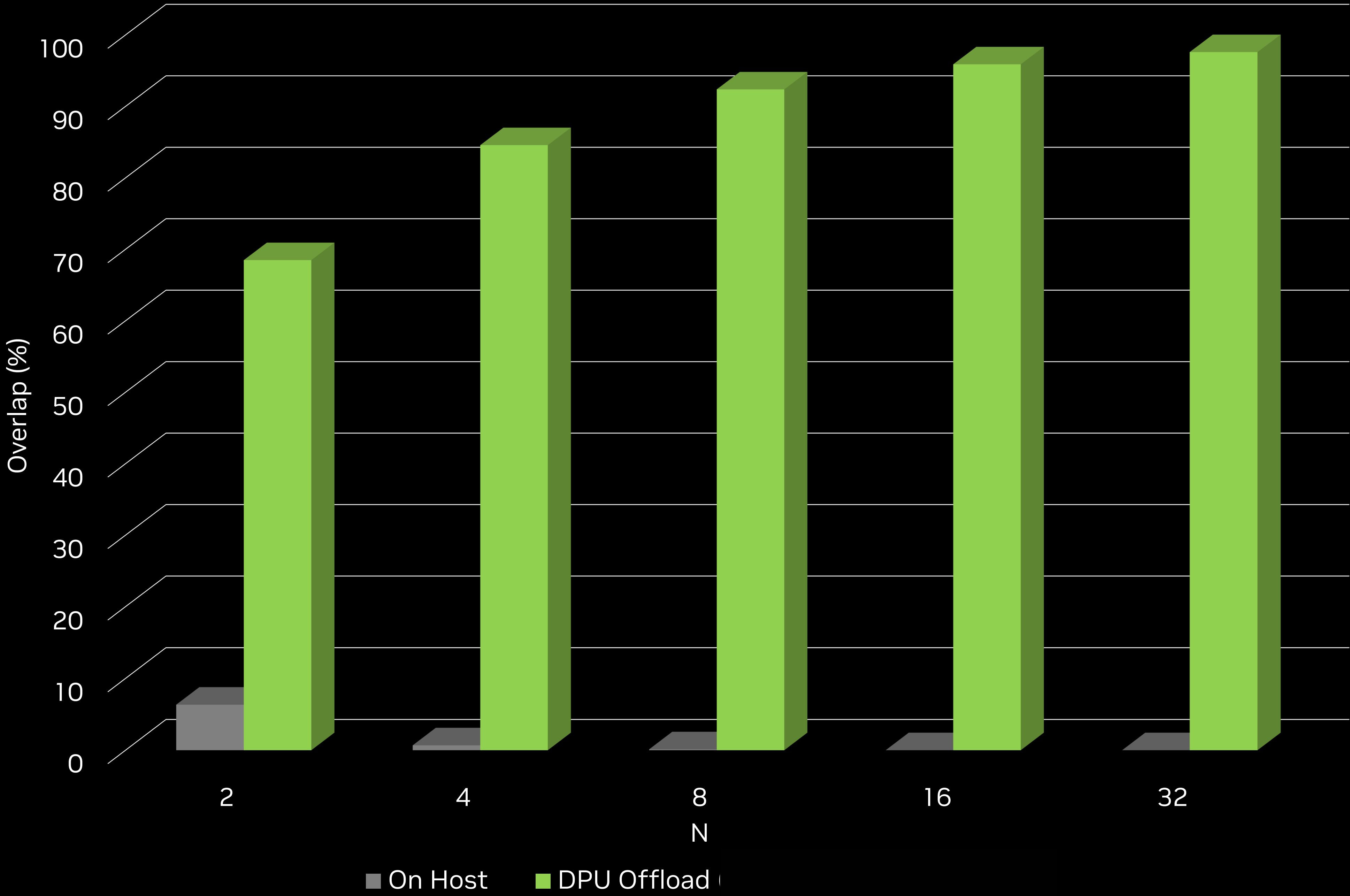
OSU iAlltoallv 32 (full) PPN, Size = 128 KB



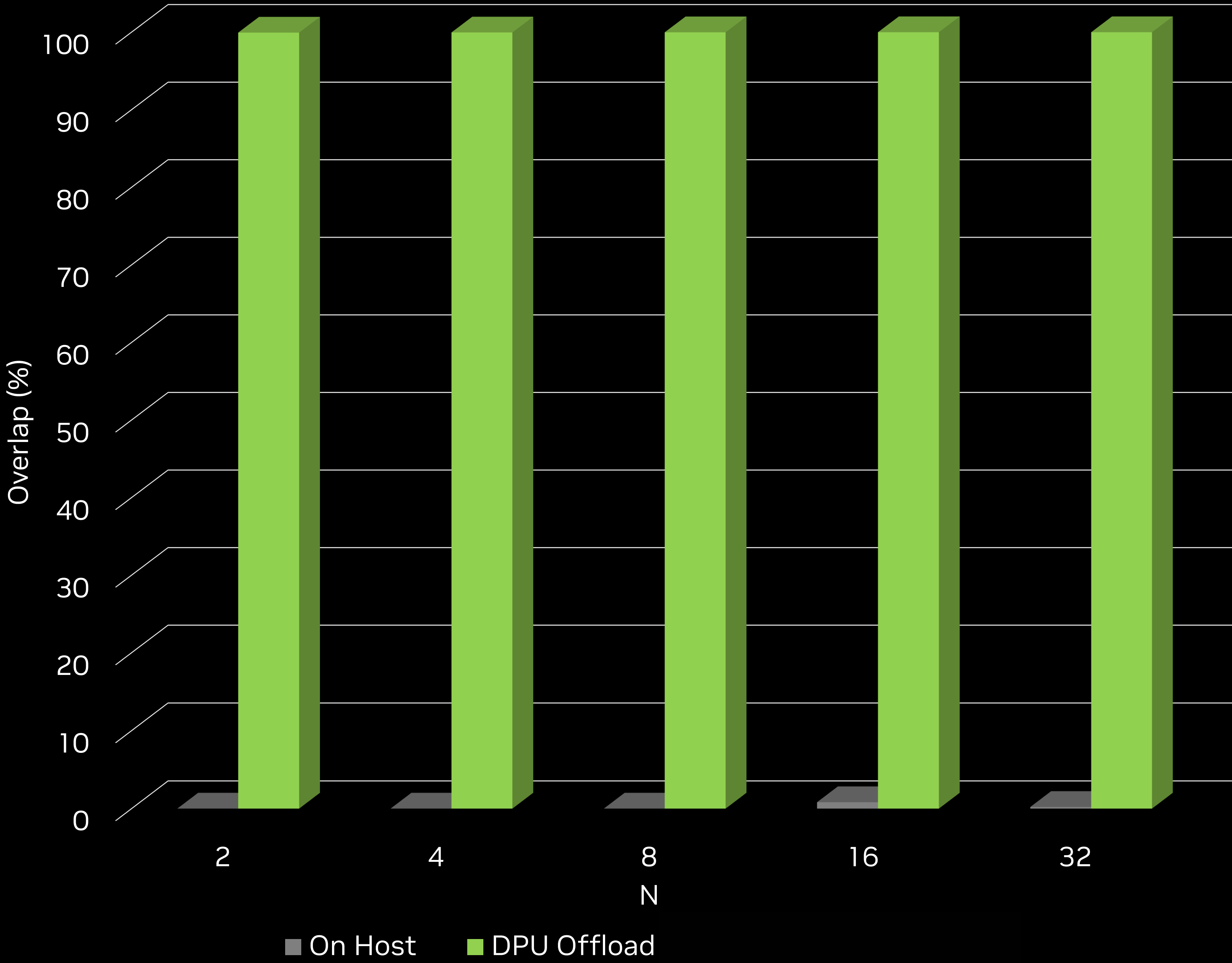


# iAlltoallv Compute/Communication Overlap

OSU iAlltoallv 1 PPN, Size = 128 KB

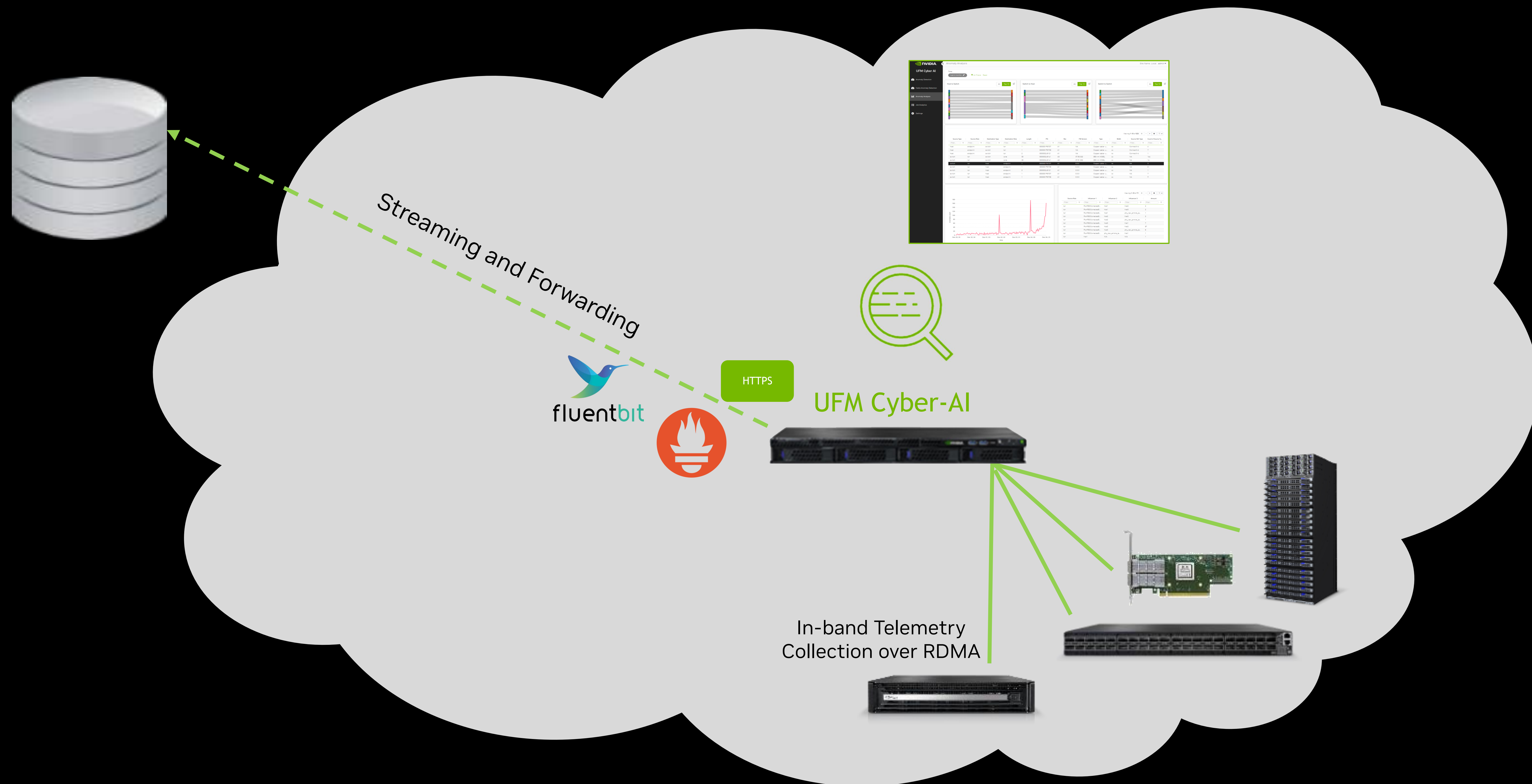


OSU iAlltoallv 32 (full) PPN, Size = 128 KB





# NVIDIA UFM Predictive Maintenance





# UFM SDK

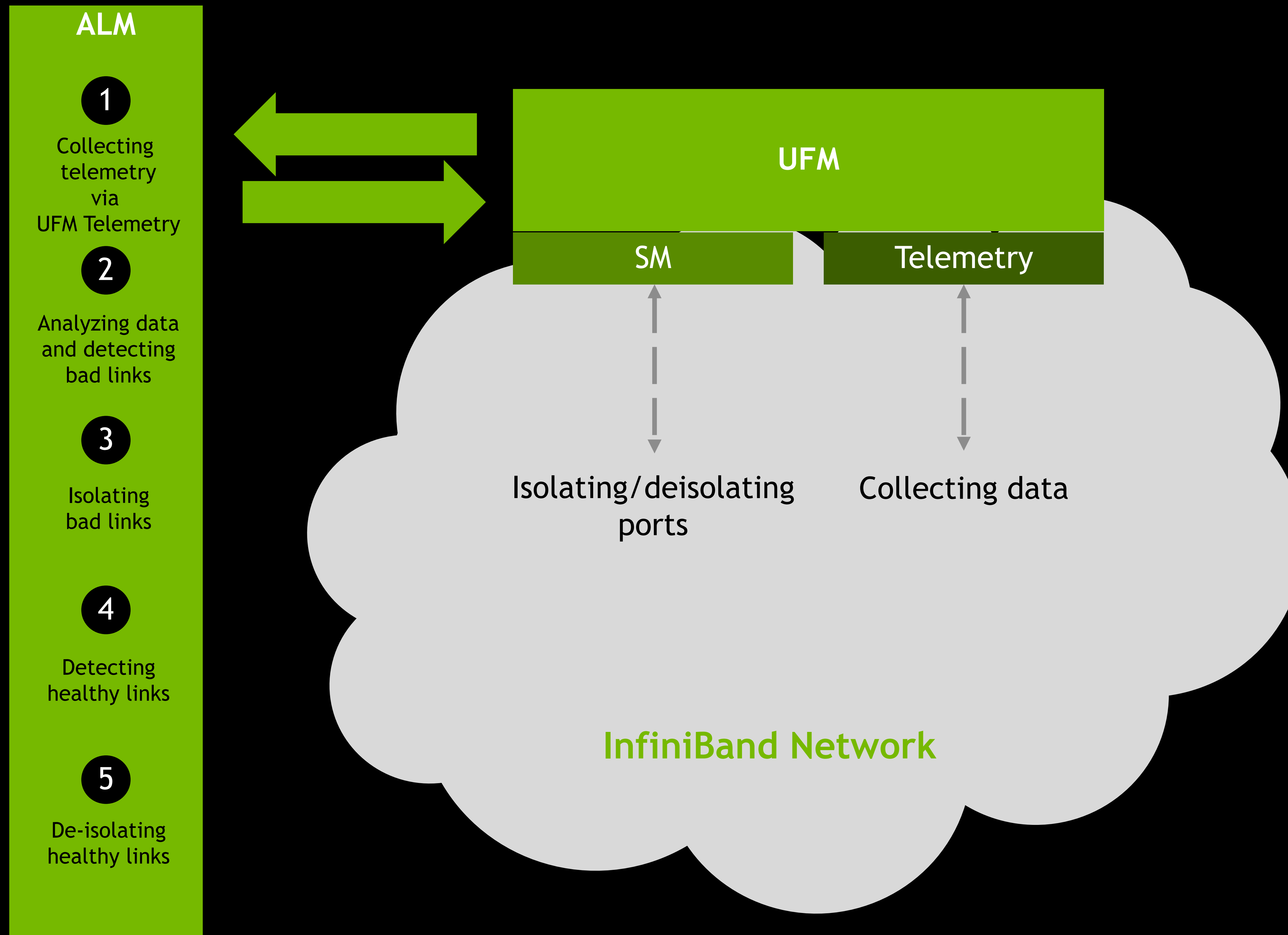
[https://github.com/Mellanox/ufm\\_sdk\\_3.0](https://github.com/Mellanox/ufm_sdk_3.0)

- A Comprehensive Suite of Plugins, Scripts and Tools for InfiniBand Network Management
- Combining enhanced, real-time network telemetry with AI-powered cyber intelligence and analytics to support scale-out data centers
- NVIDIA UFM SDK extends the capabilities of the UFM Platform with additional tools for easy third-party plugin integration



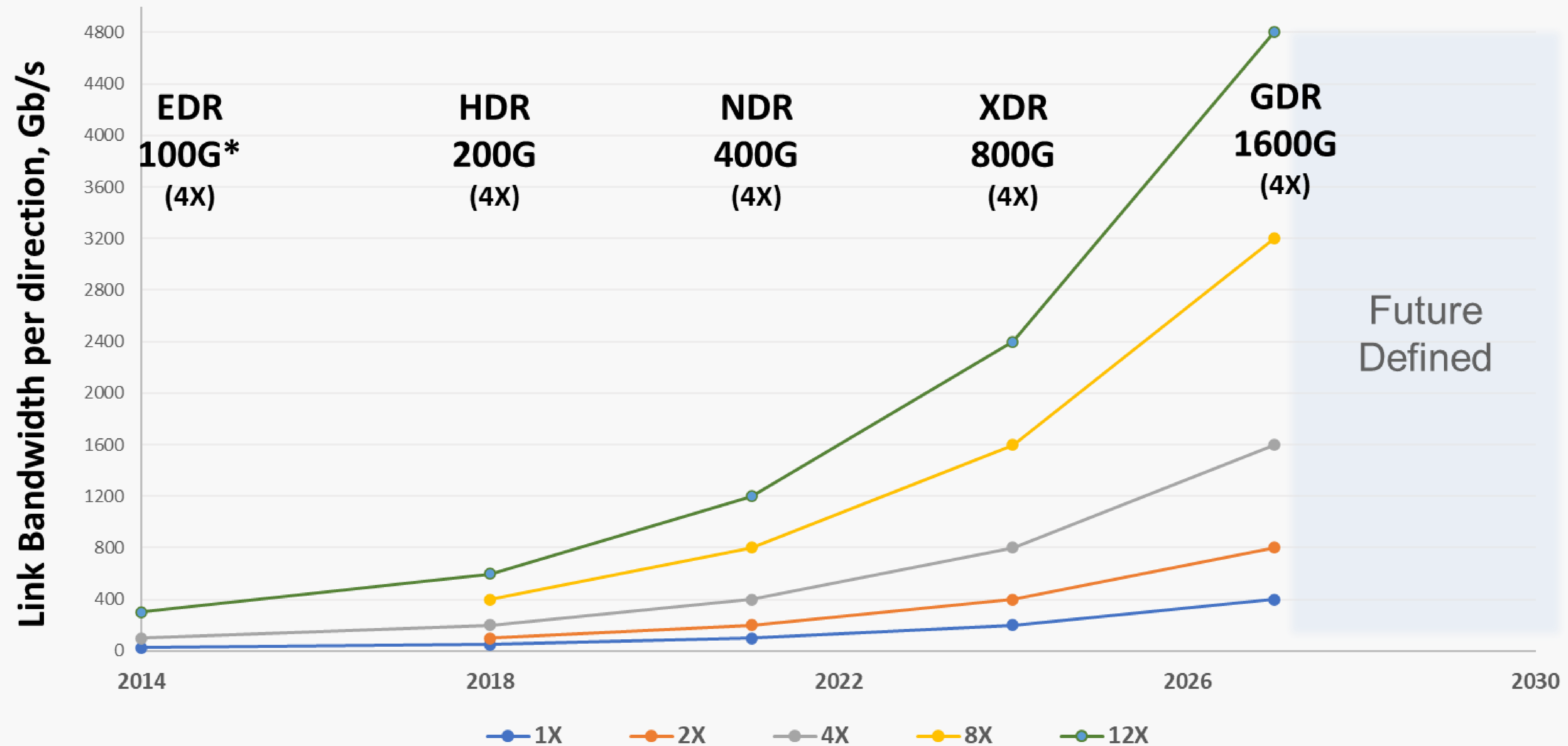


# UFM Plugin - ALM





# IBTA InfiniBand Roadmap



\*Link speeds specified in Gb/s at 4X (4 lanes)



