



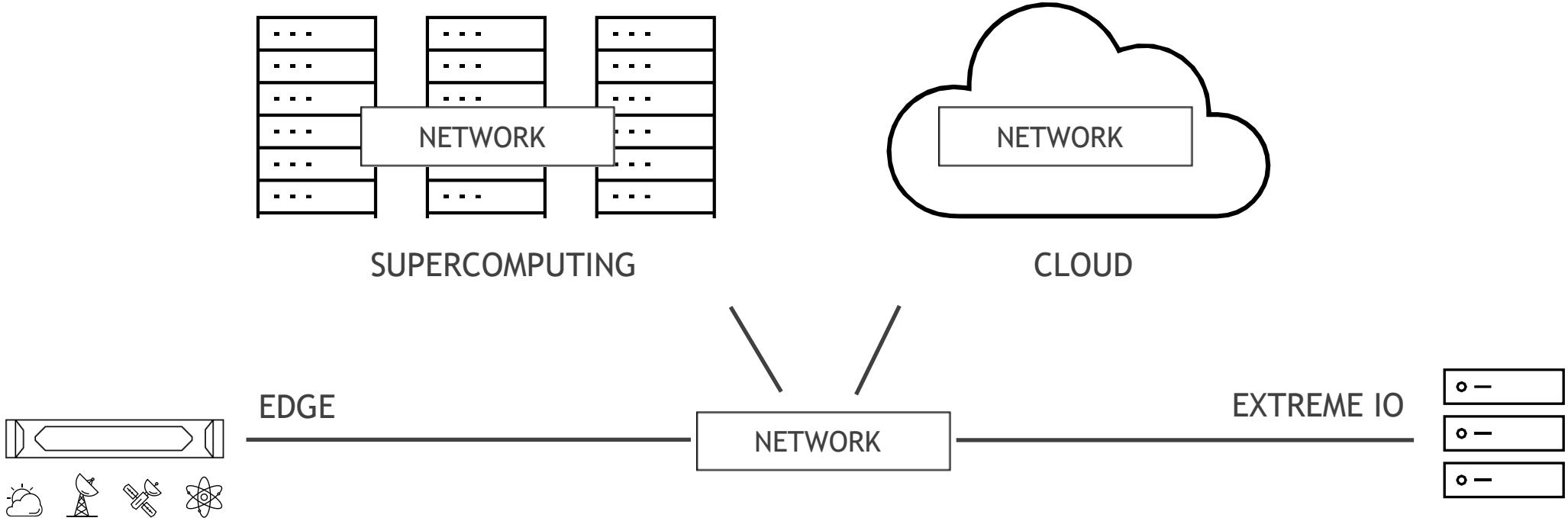
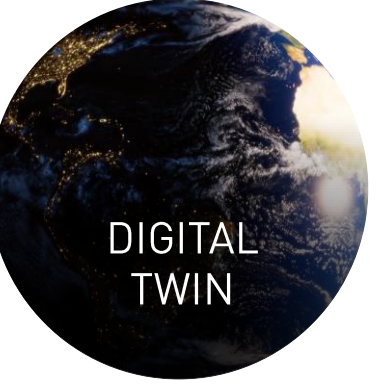
# ACCELERATING SCIENTIFIC COMPUTING

NVIDIA InfiniBand In-Network Computing Technology

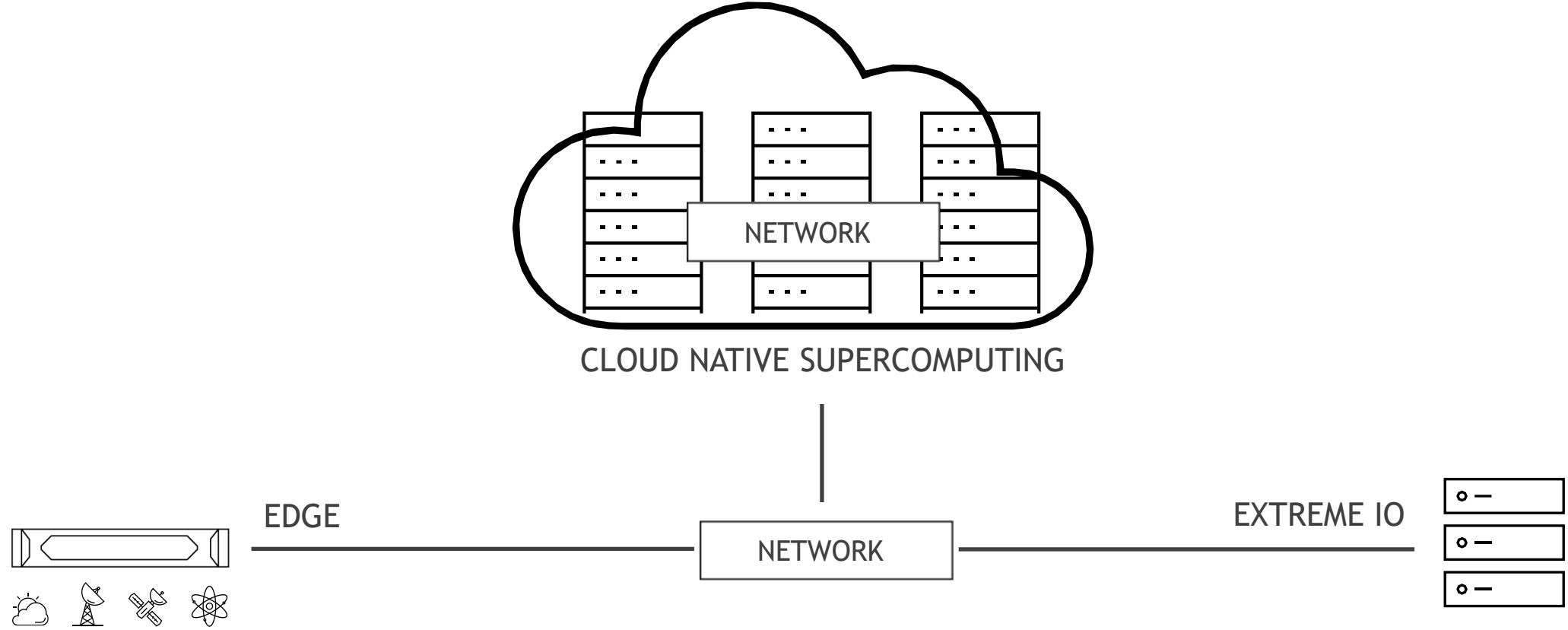
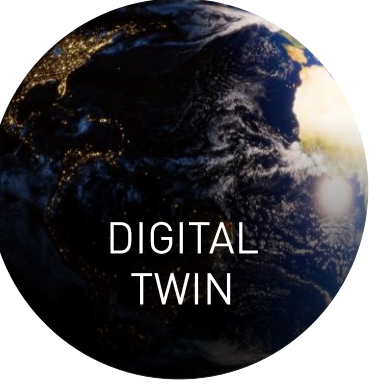
19th Workshop on High Performance Computing in Meteorology, 2021



# EXPANDING UNIVERSE OF HIGH PERFORMANCE COMPUTING



# EXPANDING UNIVERSE OF HIGH PERFORMANCE COMPUTING



# NVIDIA INFINIBAND INFRASTRUCTURE

In-Network Computing Accelerated Network for Supercomputing



Metrox Long-haul



Skyway Gateway



UFM Cyber-AI



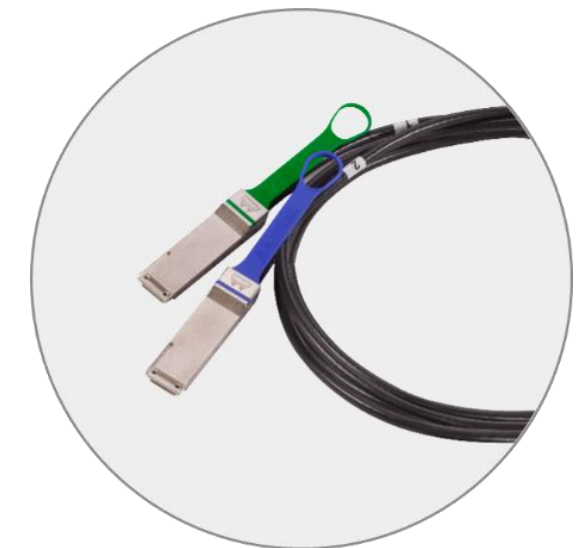
ConnectX Adapter



BlueField DPU



Quantum Switch



Linkx

# IN-NETWORK COMPUTING ACCELERATED SUPERCOMPUTING

Software-Defined, Hardware-Accelerated, InfiniBand Network

Most 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

# SCALABLE HIERARCHICAL AGGREGATION AND REDUCTION PROTOCOL (SHARP)

In-network Tree Based Aggregation Mechanism

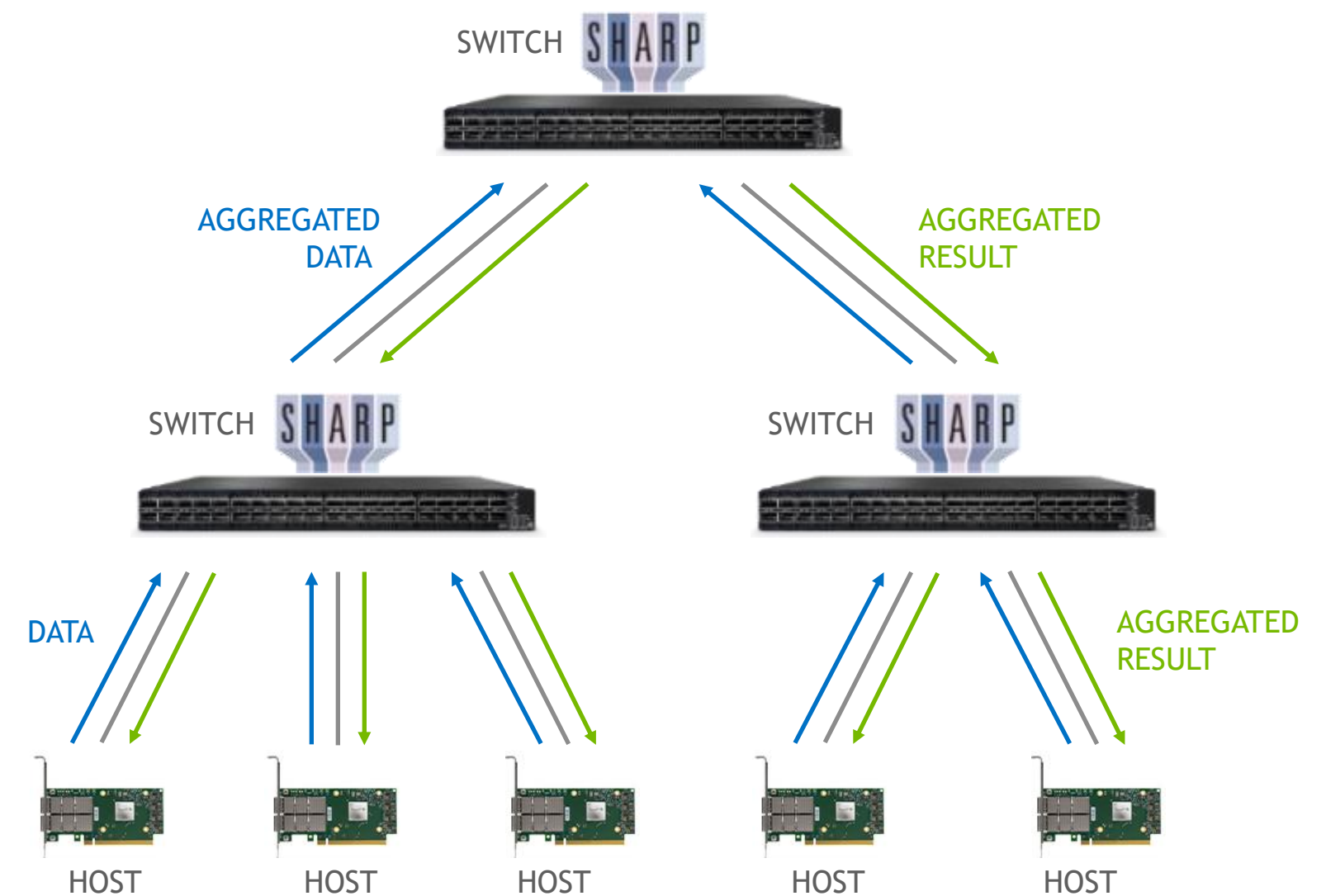
Multiple Simultaneous Outstanding Operations

Small Message and Large Message Reduction

Barrier, Reduce, All-Reduce, Broadcast and More

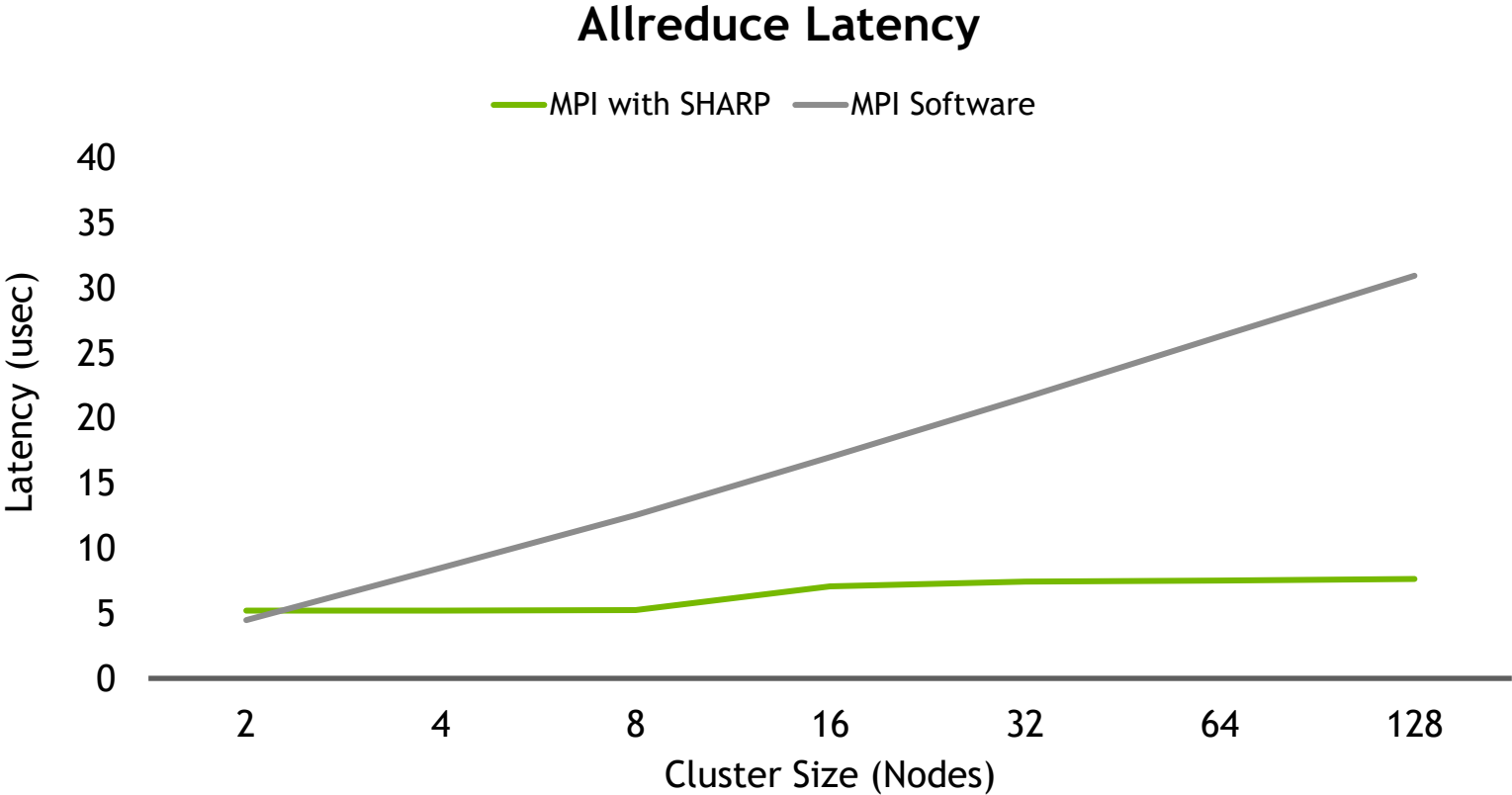
Sum, Min, Max, Min-loc, max-loc, OR, XOR, AND

Integer and Floating-Point, 16/32/64 bits

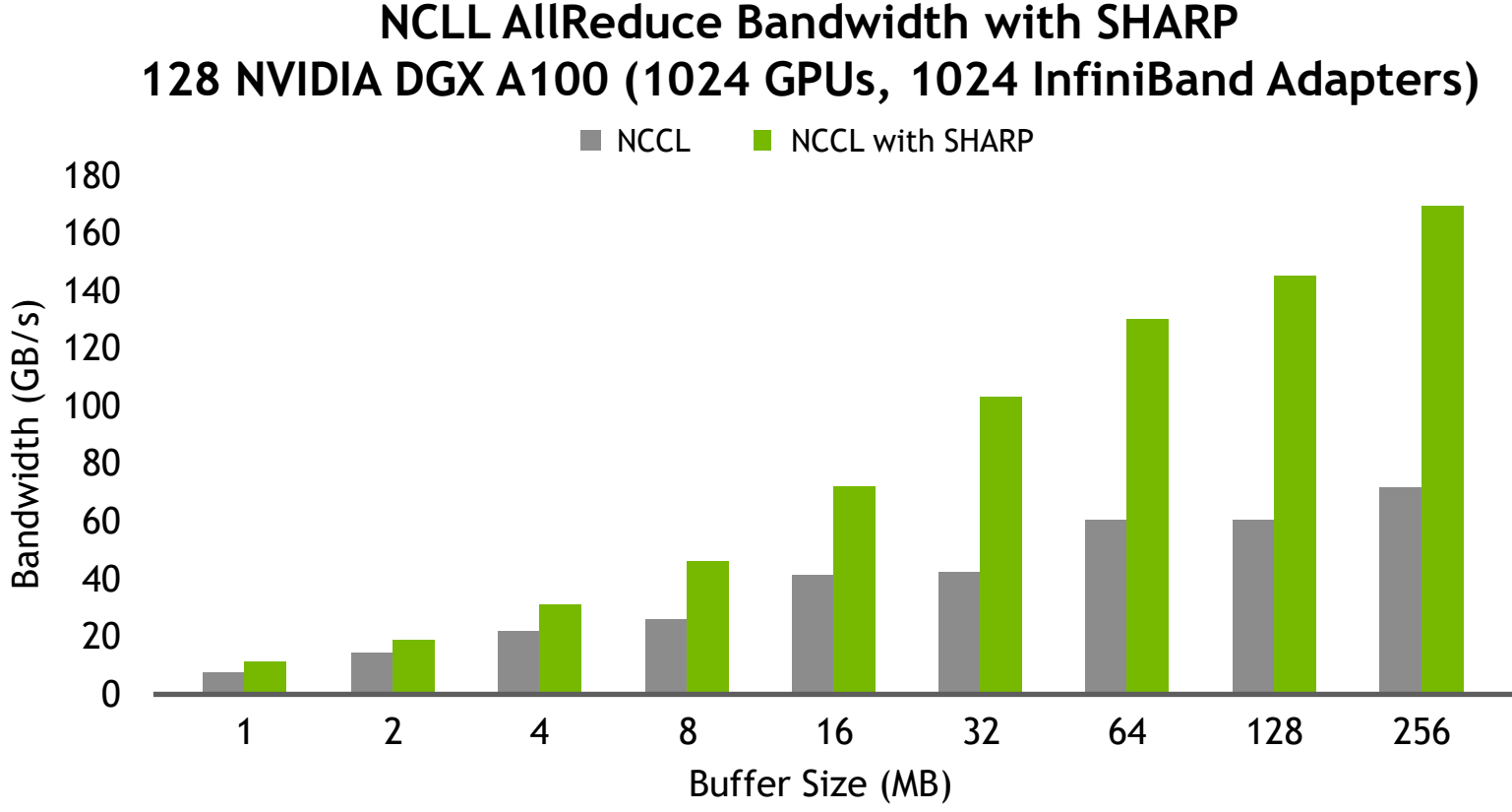


# SHARP ALLREDUCE PERFORMANCE ADVANTAGES

7x Higher MPI Performance, 2.5x Higher AI Performance

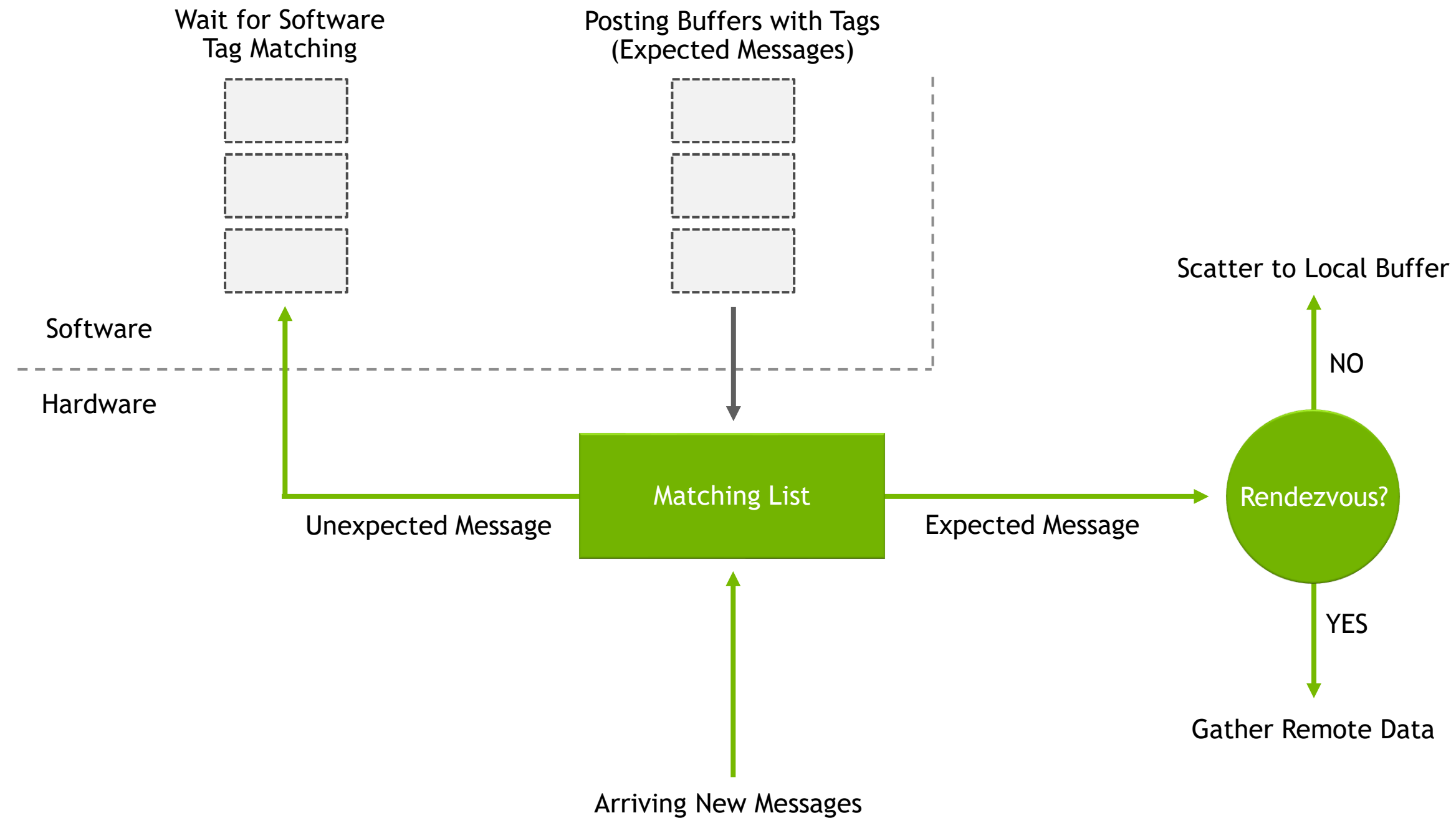


Lower is Better



Higher is Better

# INFINIBAND MPI TAG MATCHING HARDWARE ENGINE





# IN-NETWORK COMPUTING ACCELERATED SUPERCOMPUTING

Software-Defined, Hardware-Accelerated, InfiniBand Network

## In-Network Computing Acceleration Engines



New generations Introduce and Enhance Acceleration Technologies

NDR InfiniBand Includes SHARP v3 and All-to-all Engines

	Faster Data Communications	Higher Application Performance
Small Data Reduction SHARP v2	7x Faster All-Reduce	~15% higher Performance (Weather, CFD)
Large Data Reduction Sharp v2	2.5x Faster All-Reduce	15% Faster Deep Learning Recommendation 17% Faster Natural Language Processing
MPI Tag Matching	1.8x Faster MPI Iscatterv 100% Overlapping	Up to 40% Higher Performance
All-to-All (Introduced with NDR 400G)	4x Higher Throughput	Coming Soon with NVIDIA NDR InfiniBand!

# CLOUD-NATIVE SUPERCOMPUTING

Bare-metal Secured Infrastructure

Higher Application Performance

From the Edge to the Main Data Center



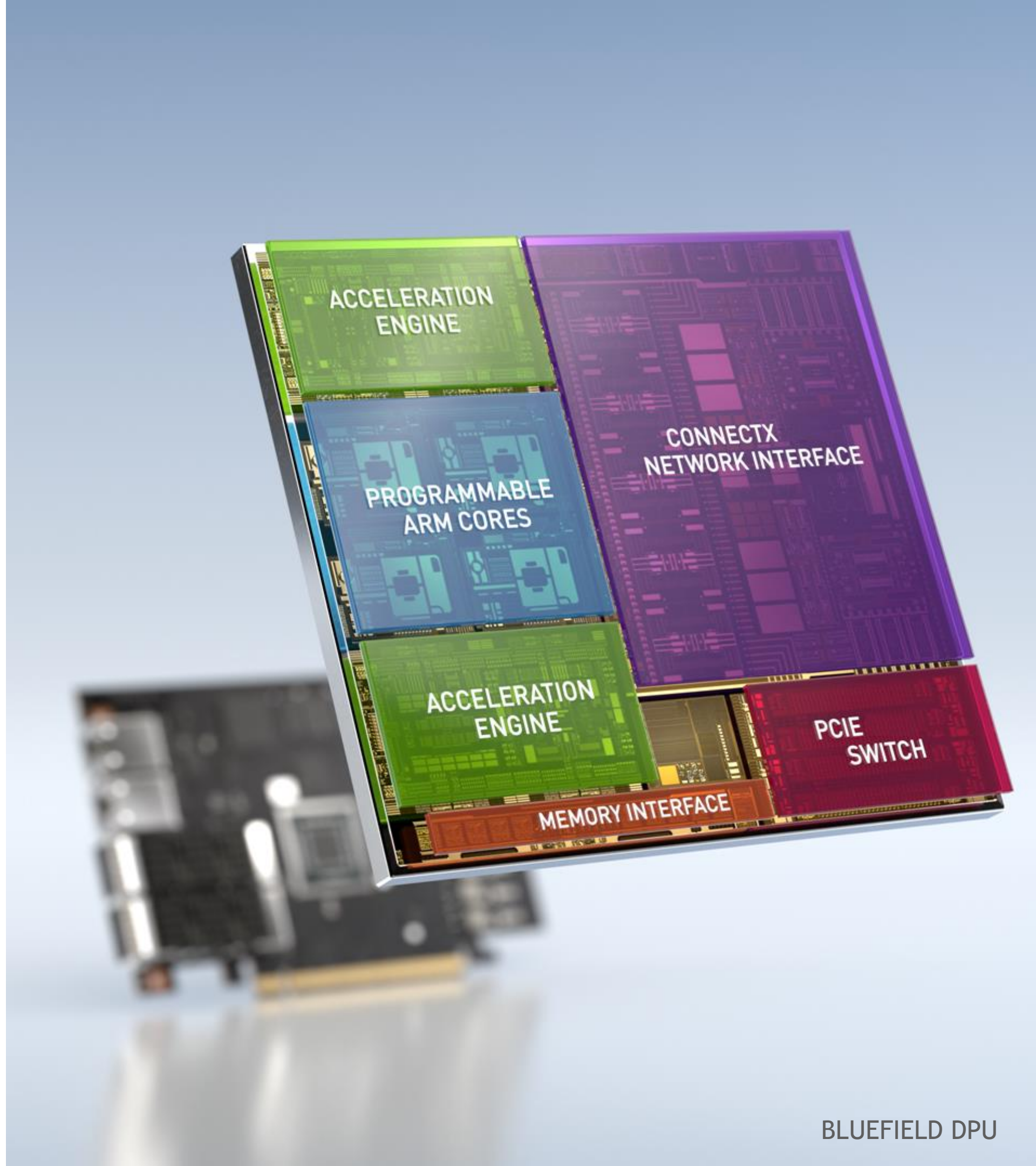
BARE-METAL  
PERFORMANCE



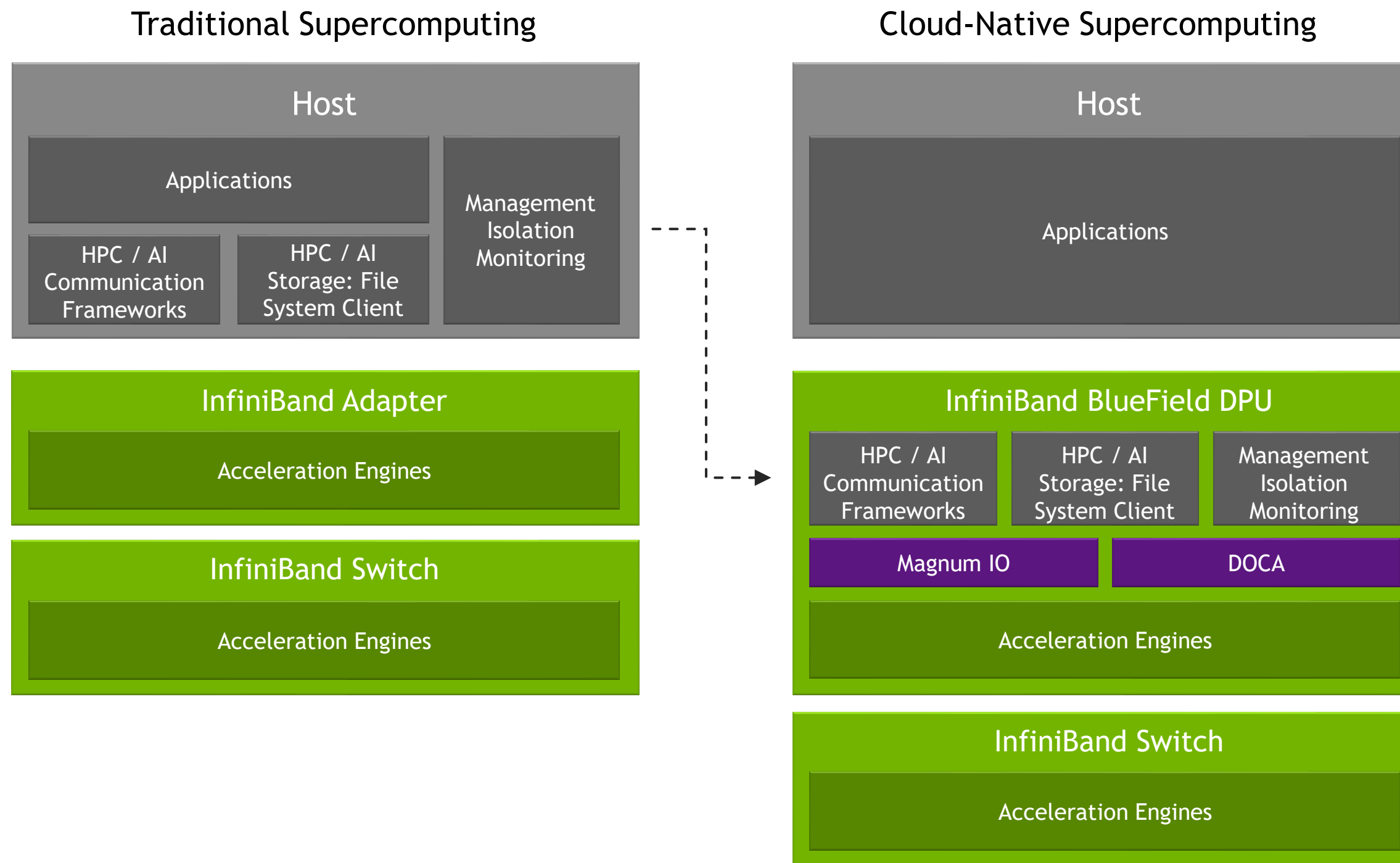
MULTI  
TENANCY



EDGE  
COMPUTING



# CLOUD-NATIVE SUPERCOMPUTING INFRASTRUCTURE



# MULTI-TENANT ISOLATION

## Zero-Trust Architecture

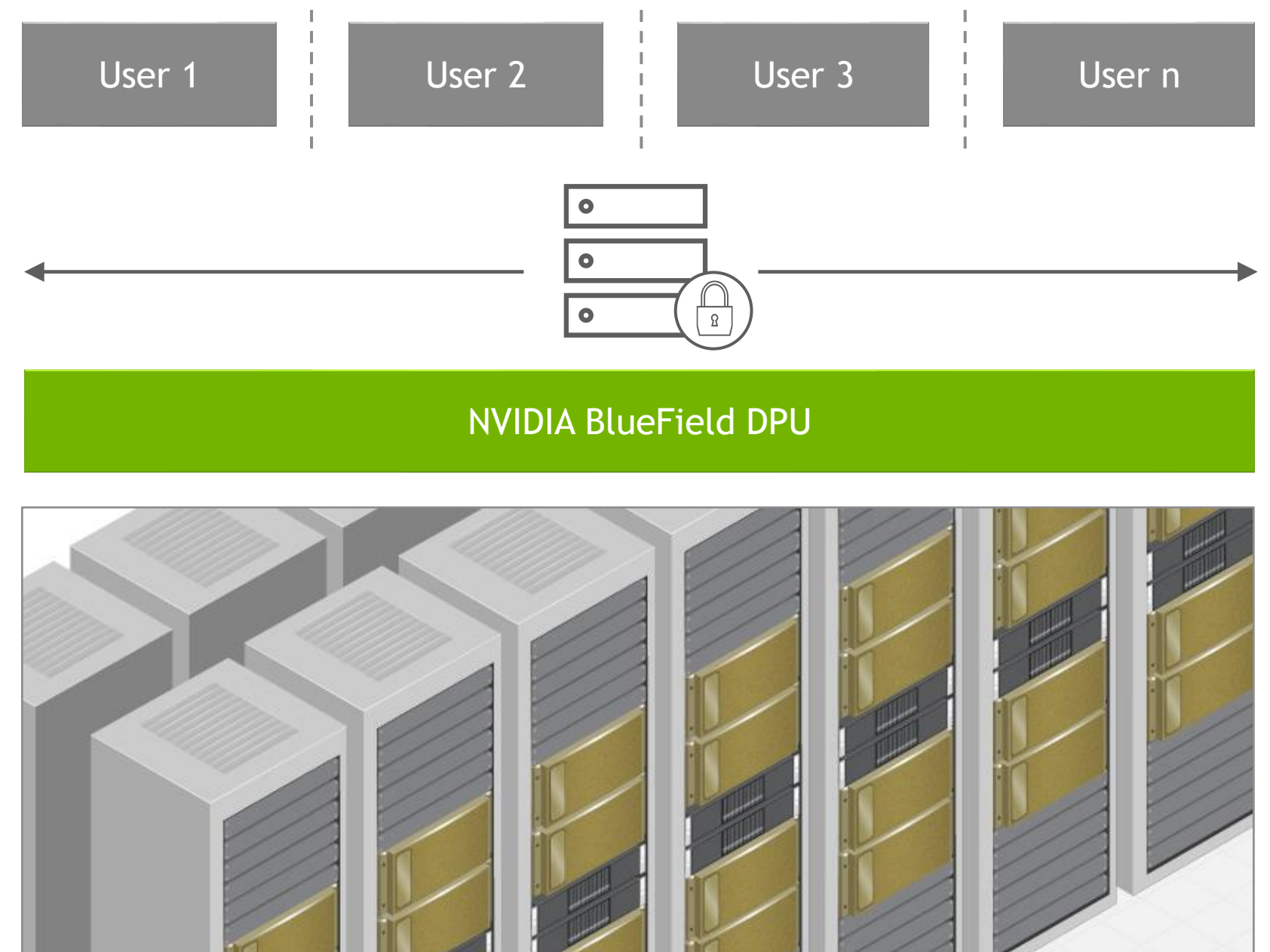
Secured Network Infrastructure and Configuration

Storage Virtualization

Tenant Service Level Agreement (SLA)

32K Concurrent Isolated Users on Single Subnet

### Secure Partitioning with Bare-Metal Performance



# HIGHER APPLICATION PERFORMANCE

## DPU-Accelerated HPC Communications

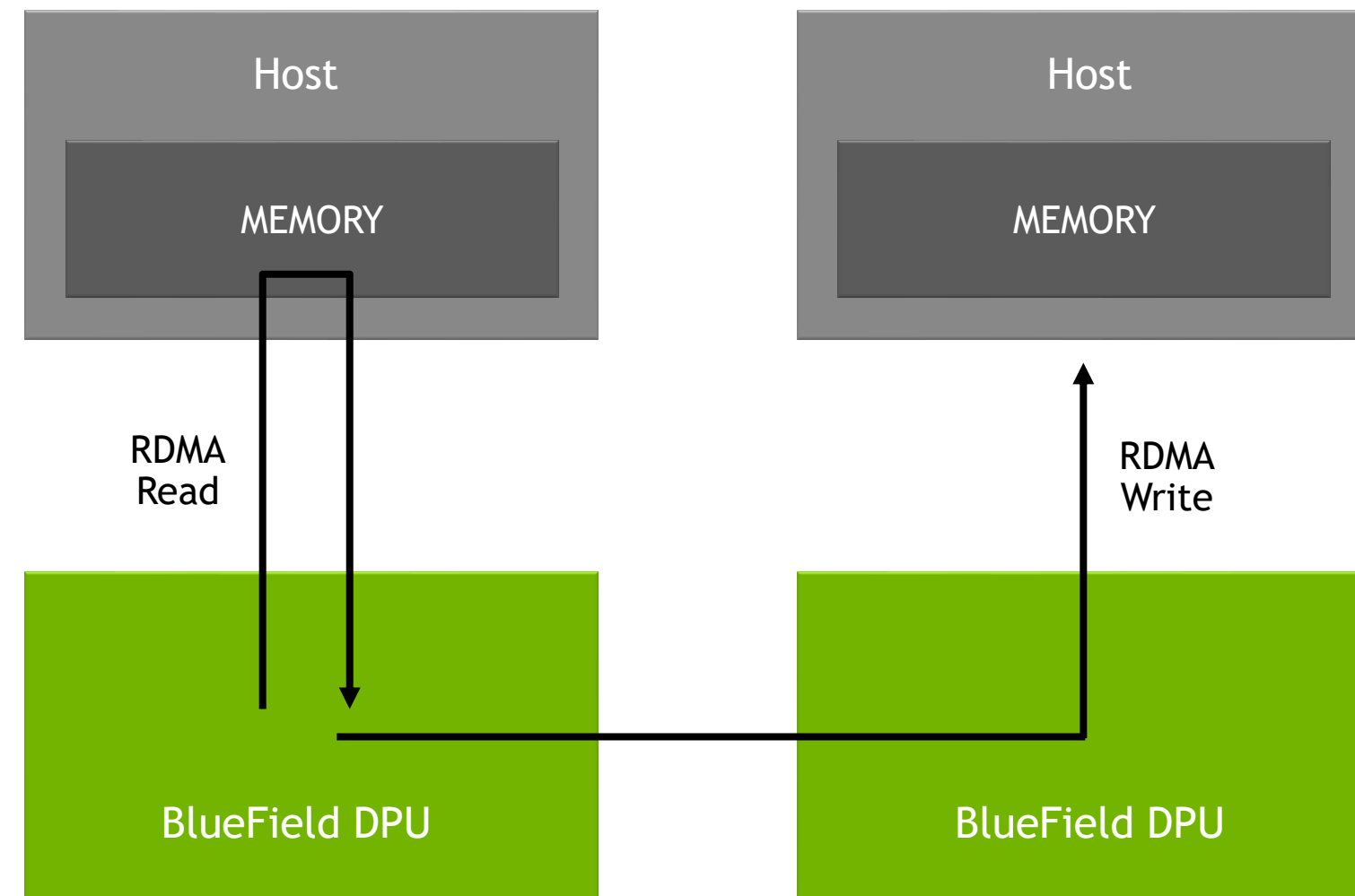
Collective Offloads

Active Messages

Smart MPI Progression

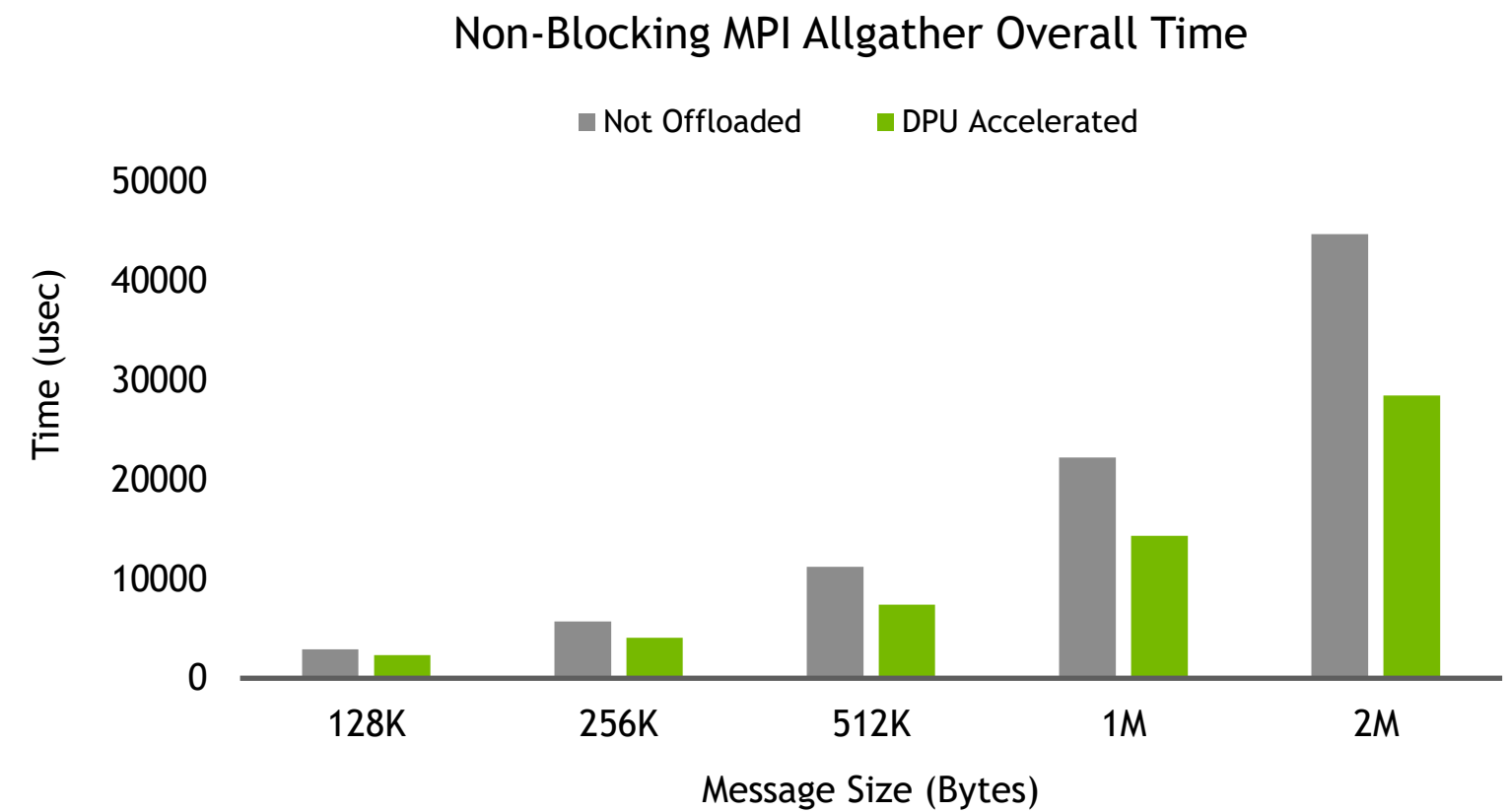
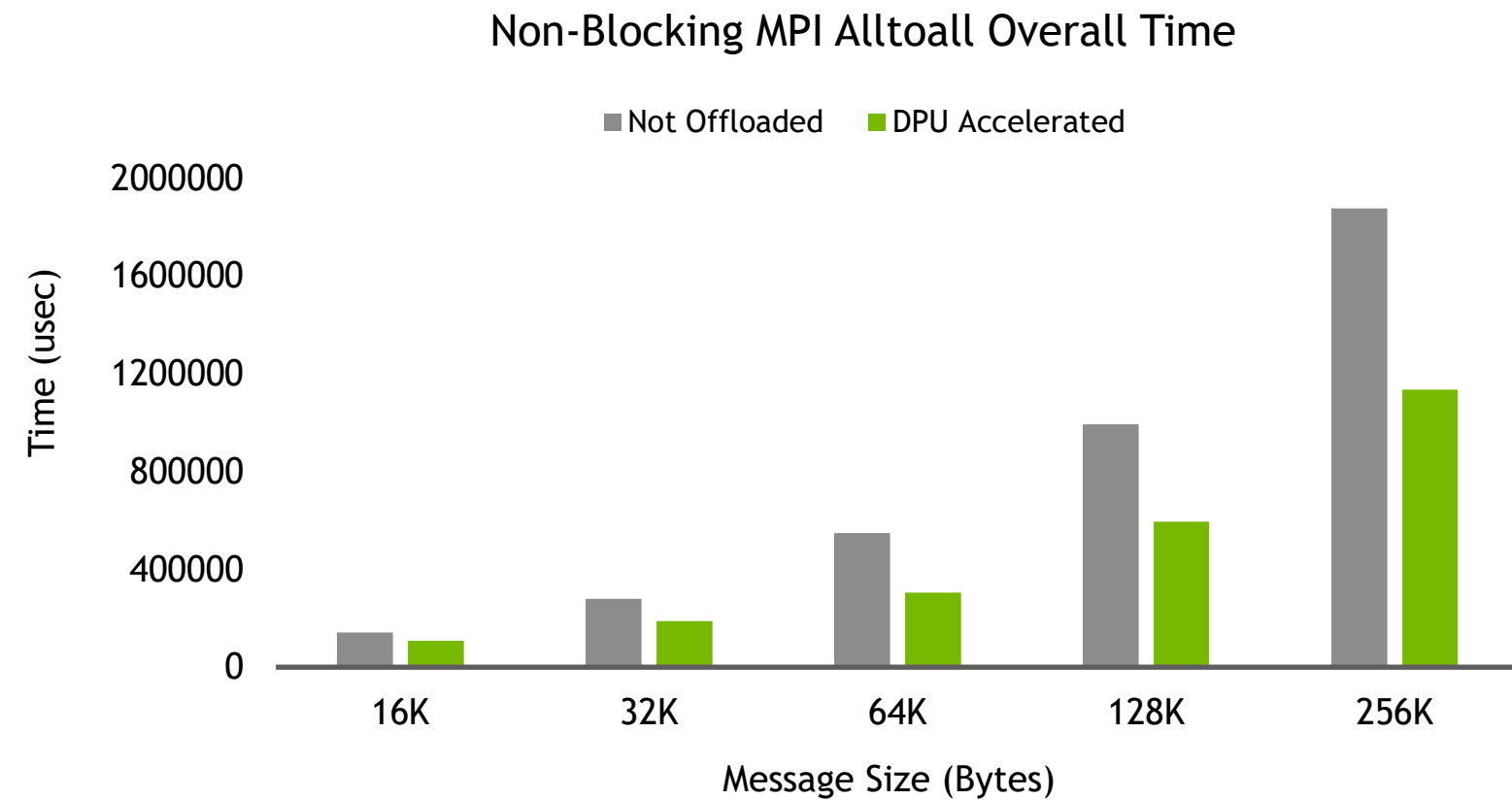
Data Compression

User-defined Algorithms



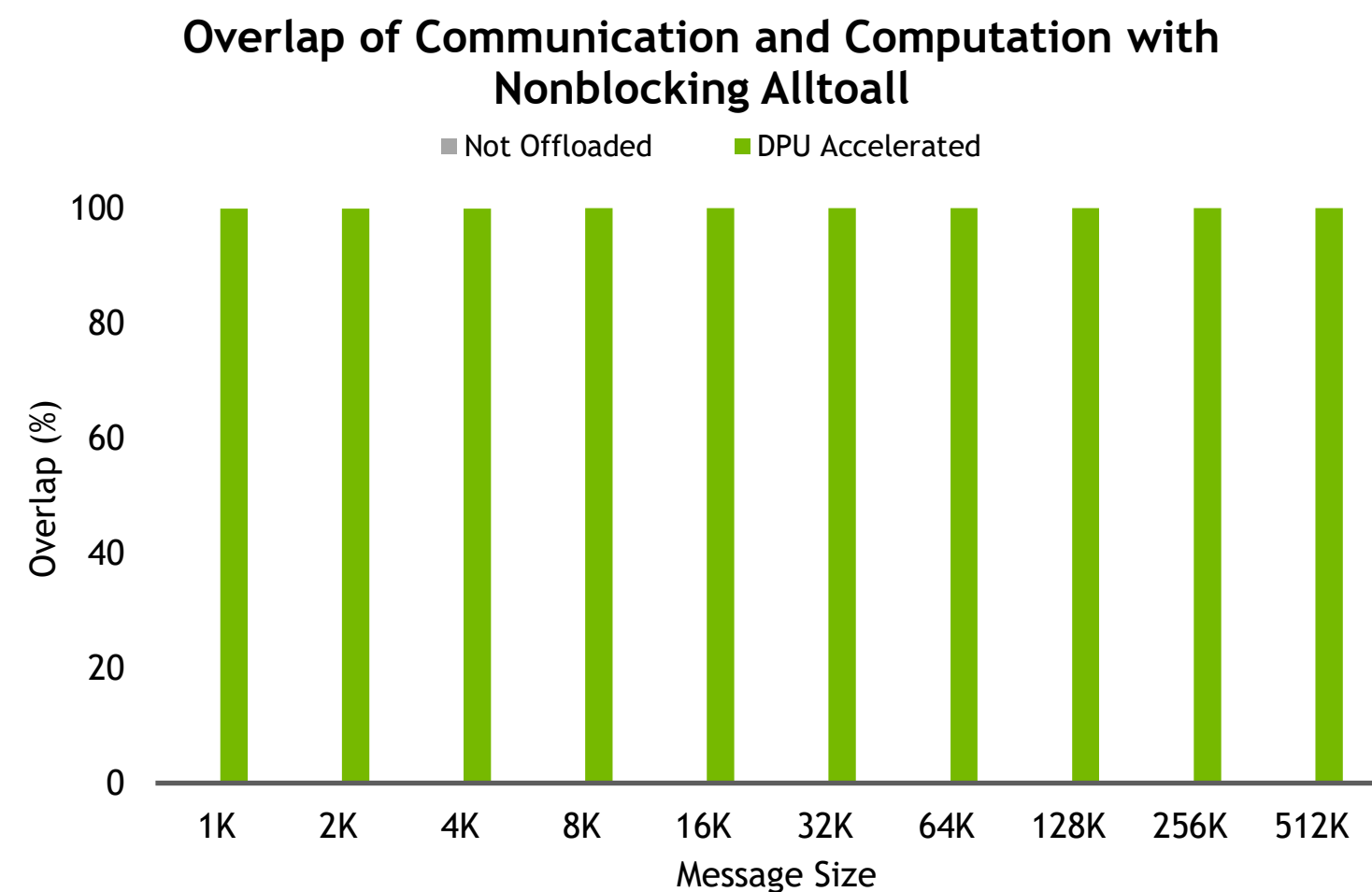
# NON-BLOCKING MPI PERFORMANCE

44% Performance Increase for MPI iAlltoall, 36% Performance Increase for MPI iAllgather



# HIGHER APPLICATION PERFORMANCE

100% Communication - Computation Overlap



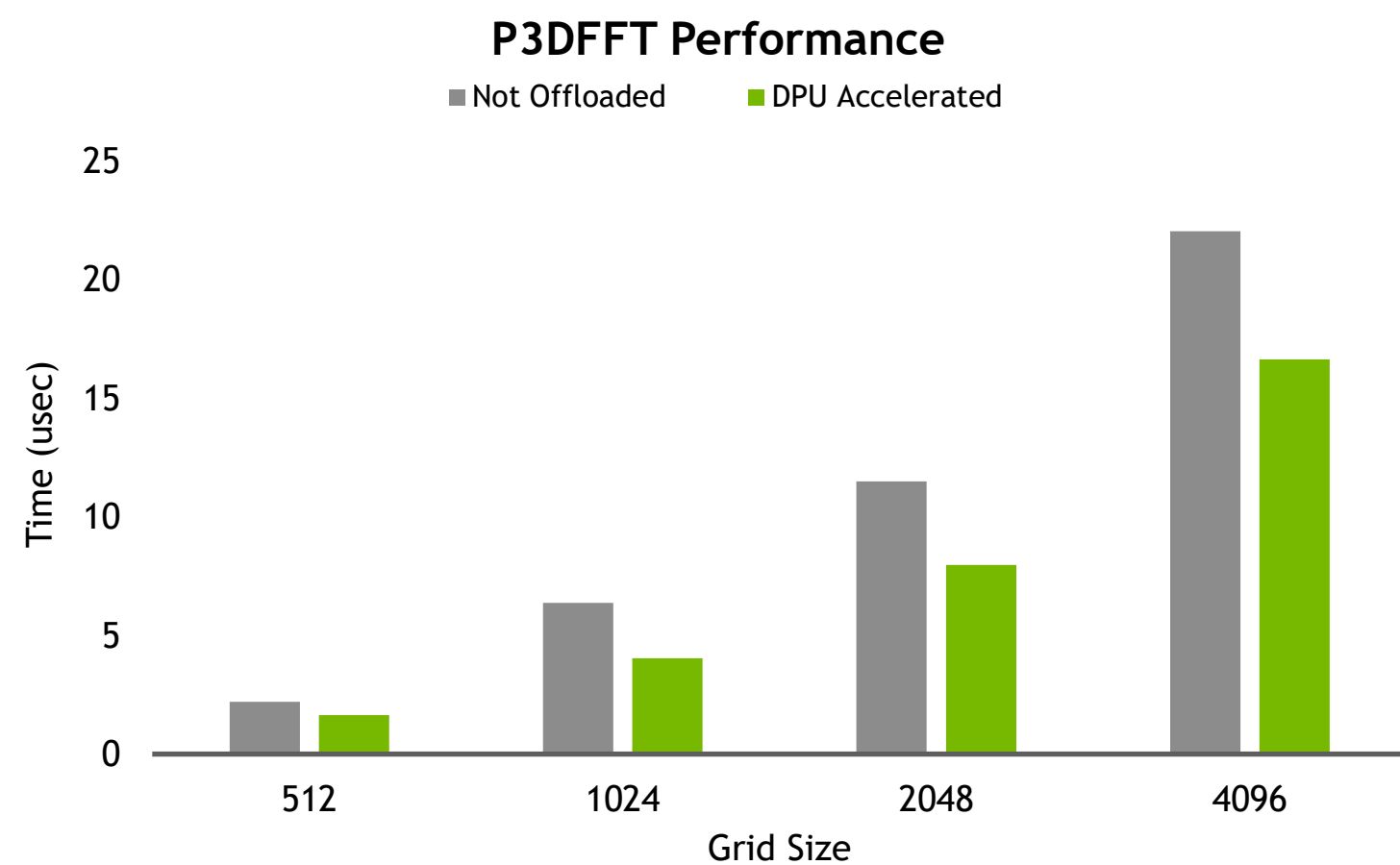
Courtesy of: Ohio State University MVAICH Team and X-ScaleSolutions



32 servers, Dual Socket Intel® Xeon® 16-core CPUs E5-2697A V4 @ 2.60 GHz (32 processes per node), NVIDIA BlueField-2 HDR100 DPUs and ConnectX-6 HDR100 adapters, NVIDIA HDR Quantum Switch QM7800 40-Port 200Gb/s HDR InfiniBand, 256GB DDR4 2400MHz RDIMMs memory and 1TB 7.2K RPM SATA 2.5" hard drive per node.

# HIGHER APPLICATION PERFORMANCE

Higher App Performance,  
MPI Collectives Offload



Courtesy of: Ohio State University MVAICH Team and X-ScaleSolutions



32 servers, Dual Socket Intel® Xeon® 16-core CPUs E5-2697A V4 @ 2.60 GHz (32 processes per node), NVIDIA BlueField-2 HDR100 DPUs and ConnectX-6 HDR100 adapters, NVIDIA HDR Quantum Switch QM7800 40-Port 200Gb/s HDR InfiniBand, 256GB DDR4 2400MHz RDIMMs memory and 1TB 7.2K RPM SATA 2.5" hard drive per node.



# NVIDIA QUANTUM NDR 400G INFINIBAND SYSTEMS

In-Network Computing Accelerated  
Network for Cloud-Native  
Supercomputing at Any Scale

**2x**

Data Throughput  
400 Gigabits per Second

**32x**

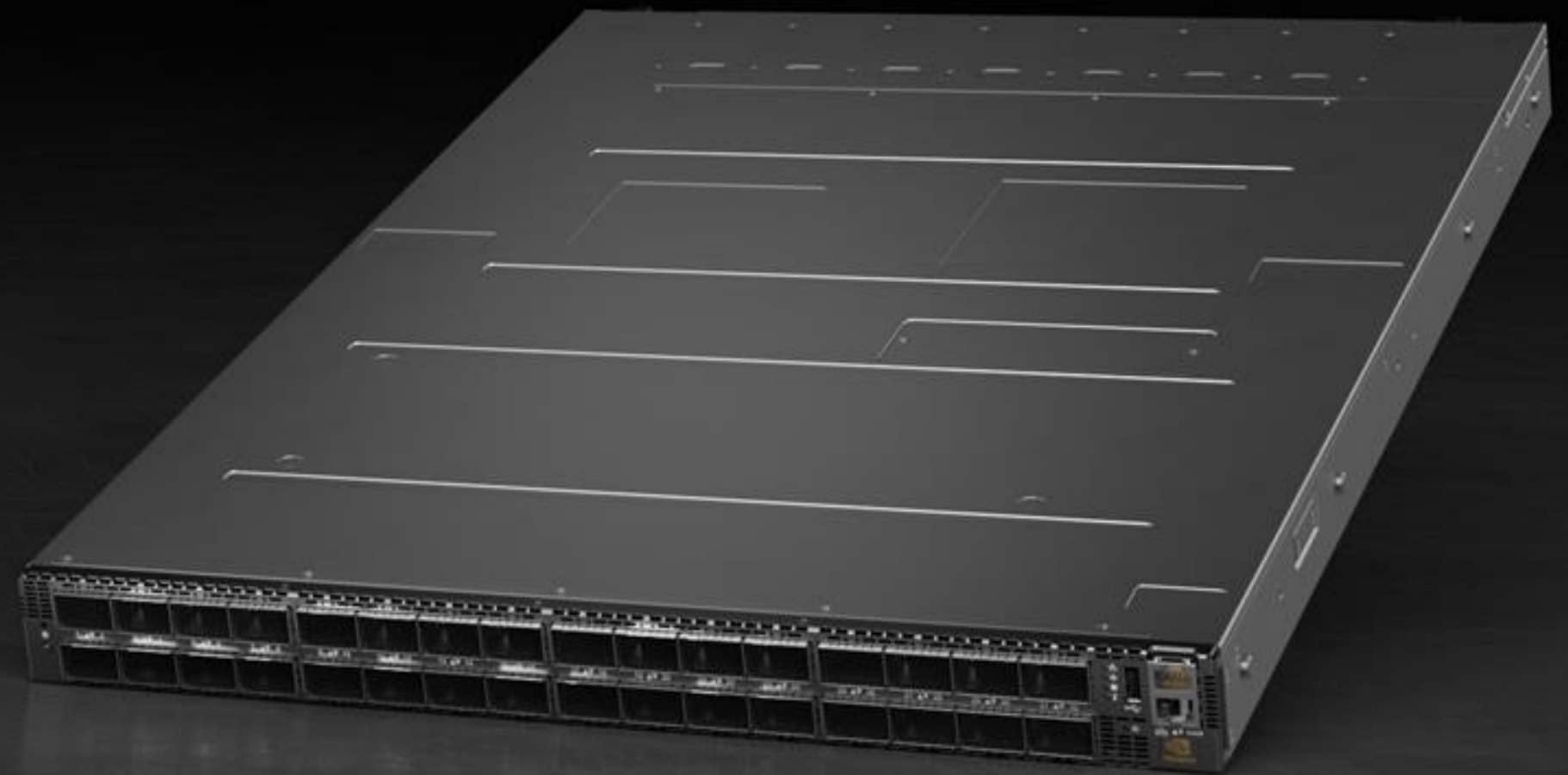
More AI Acceleration  
SHARP In-Network Computing

**6.5x**

Higher Scalability  
>1M nodes with DF+ 3 hops

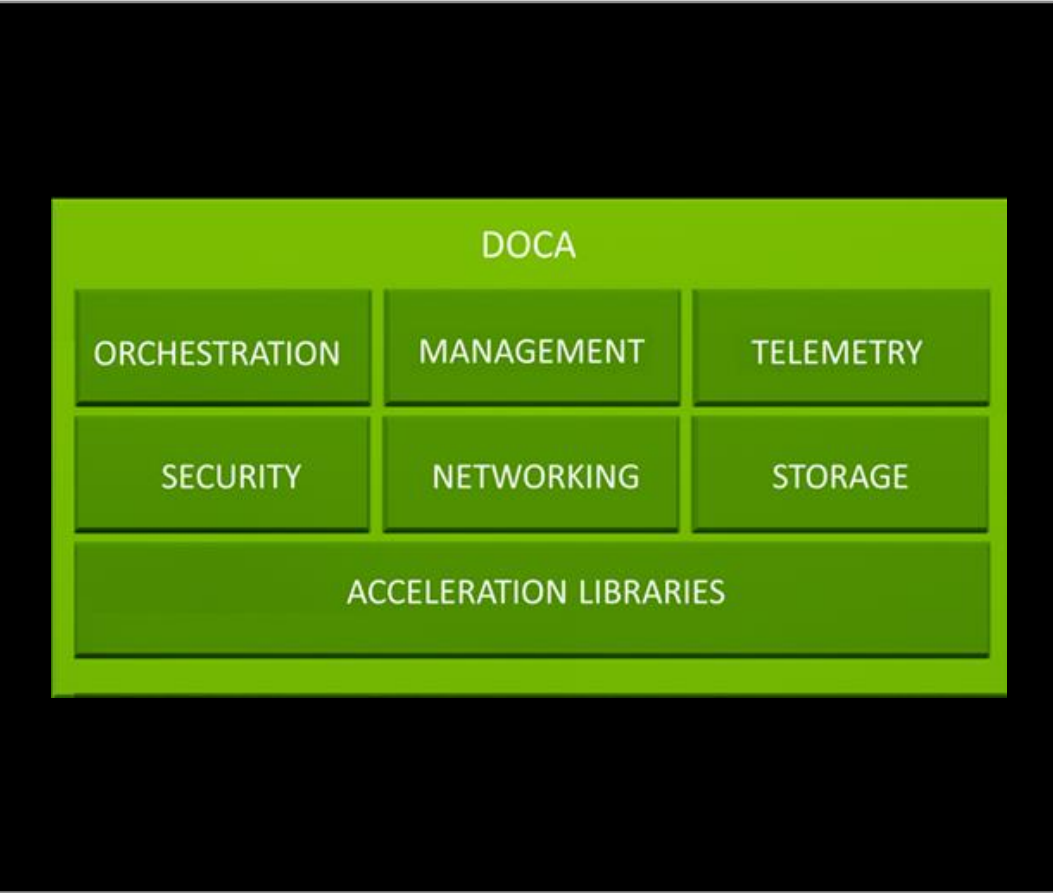
**5x**

Switch System Capacity  
>1.6 Petabit per Second

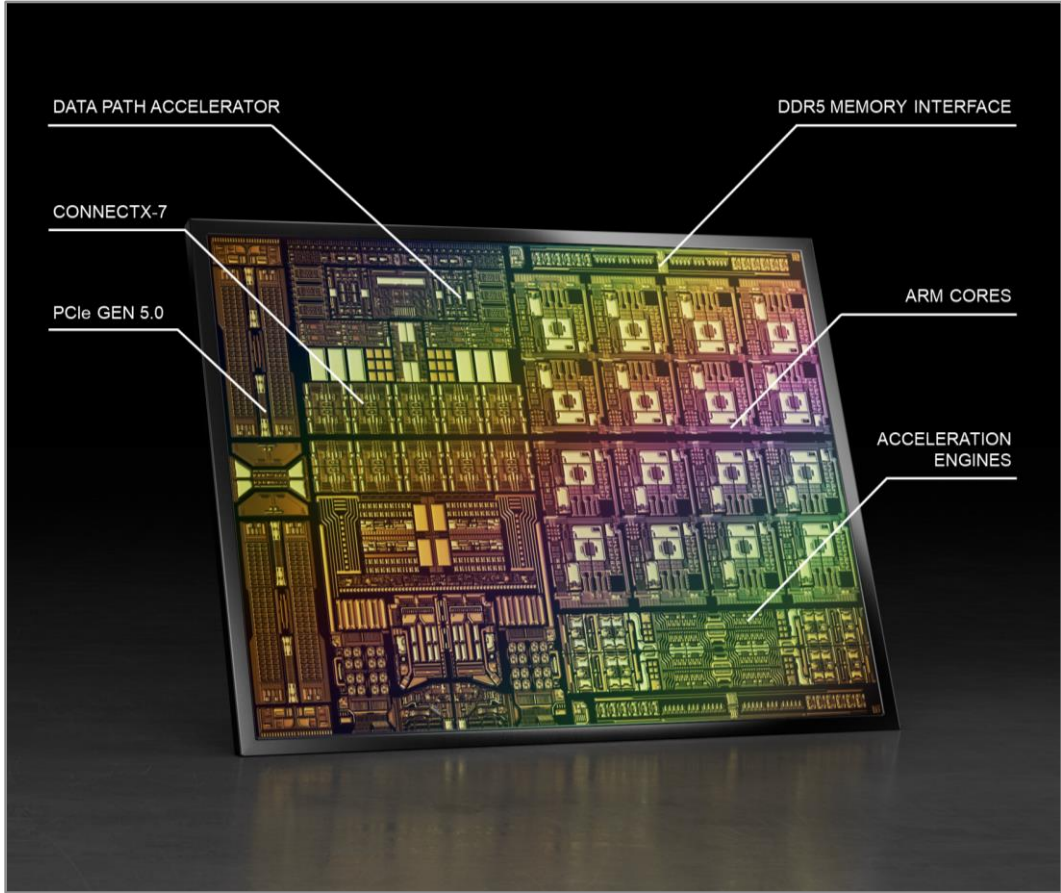


# FROM SUPERCOMPUTERS TO SUPERCLOUDS: CLOUD-NATIVE SUPERCOMPUTERS

Software-Defined, Hardware-Accelerated, InfiniBand Network



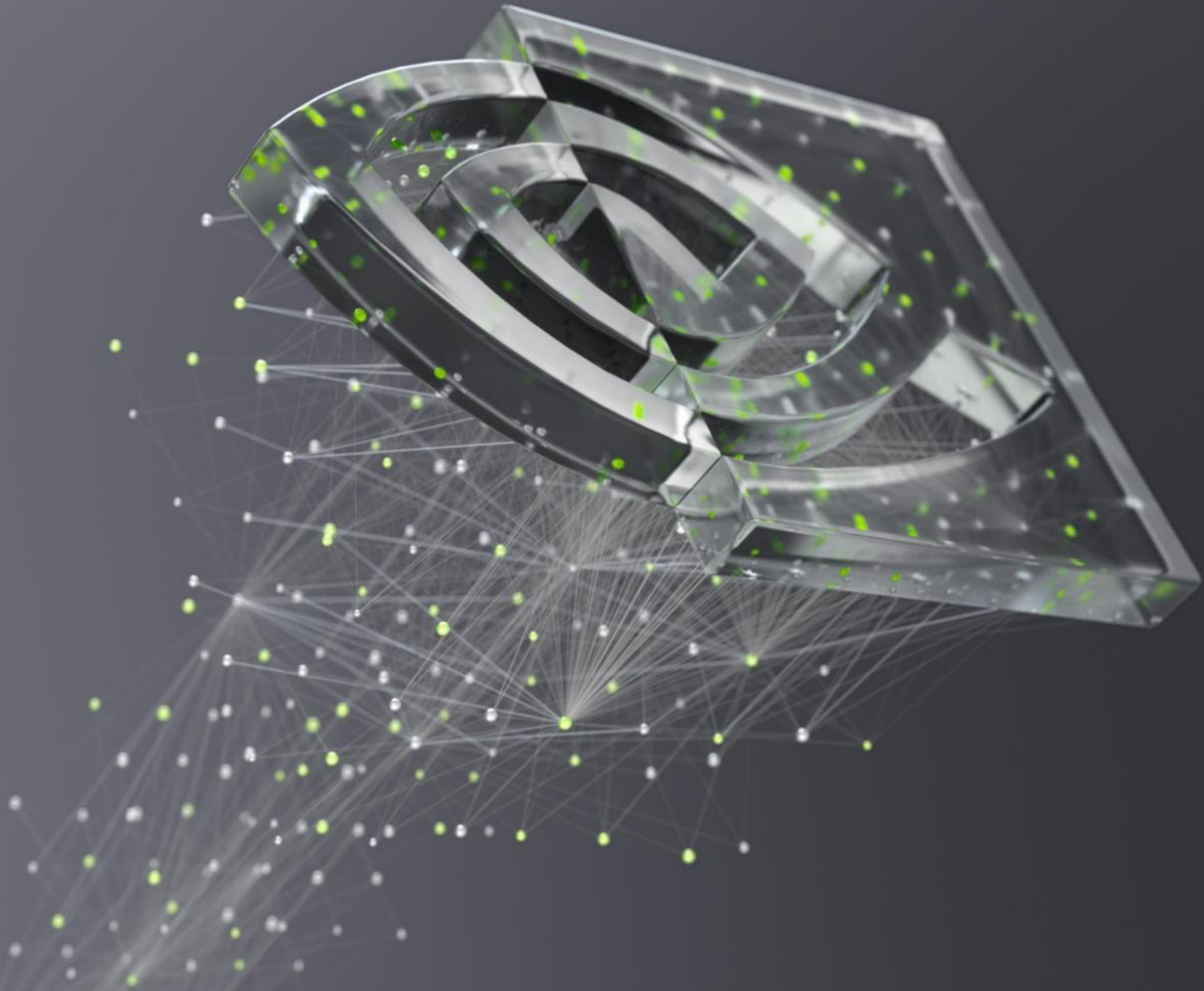
DOCA Enabling Growing  
Partner Ecosystem



Bluefield Data Center Infra Processor



NVIDIA Quantum NDR 400G InfiniBand  
In-network Computing Interconnect



**nVIDIA**