



2023 OFA Virtual Workshop

# ACCELERATING SCIENTIFIC COMPUTING WORKLOADS WITH INFINIBAND DPUS

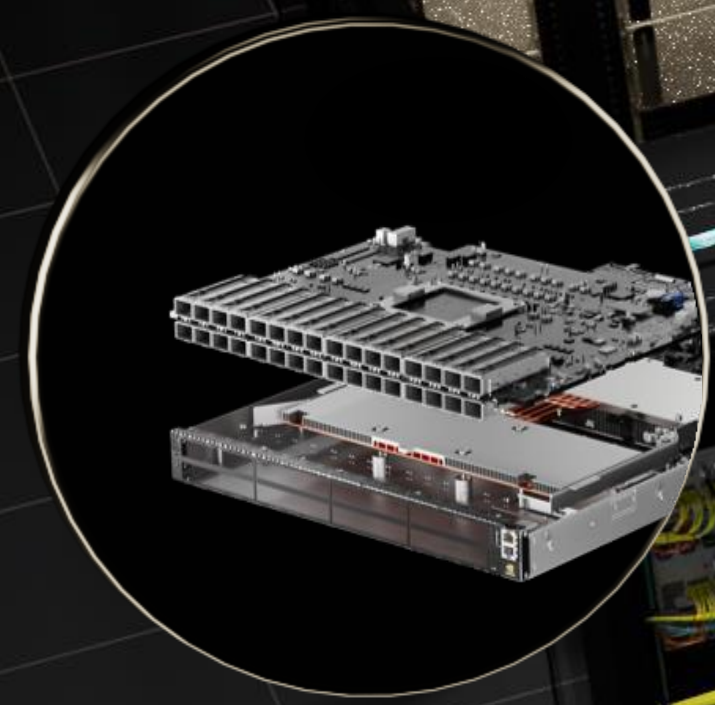
Richard Graham

NVIDIA

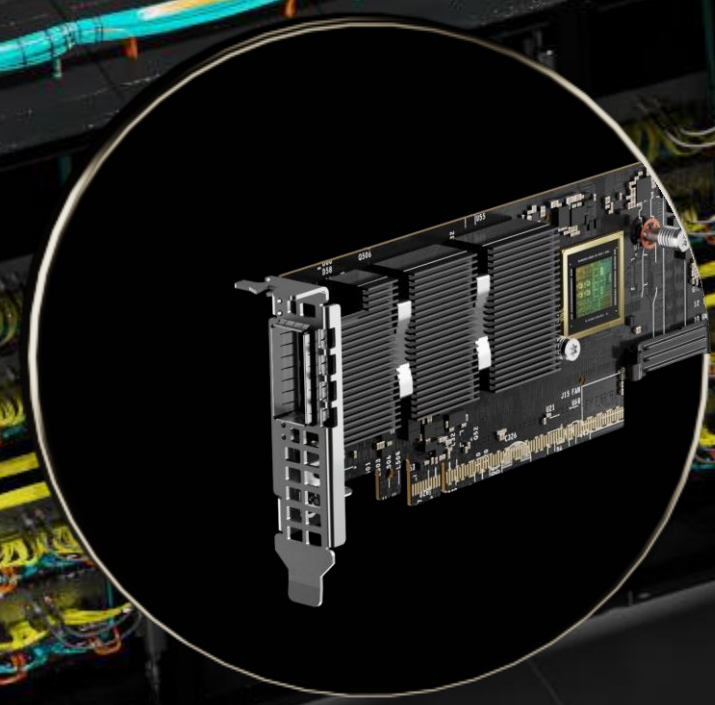


## NVIDIA CLOUD-NATIVE HPC PLATFORM

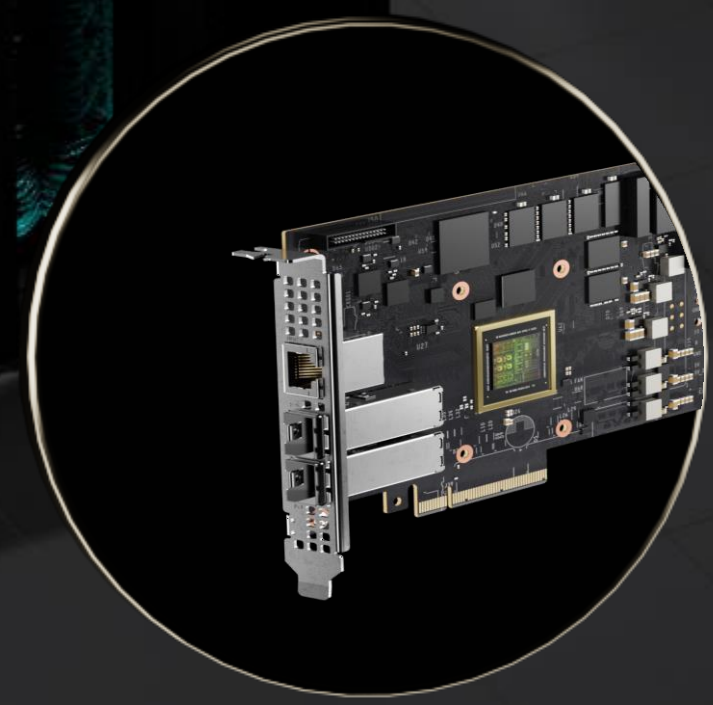
In-Network Computing  
Zero Trust Security  
Performance Isolation  
Computational Storage  
Enhanced Telemetry



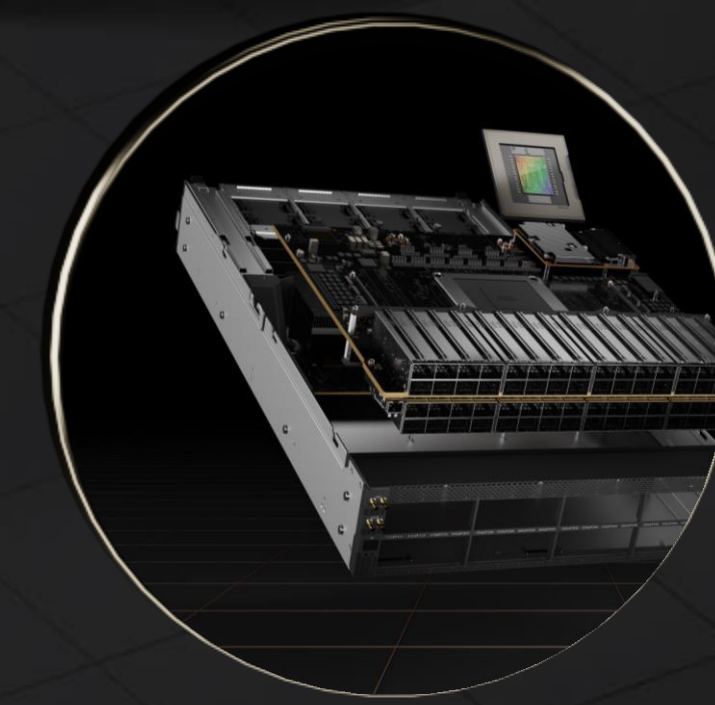
QUANTUM INFINIBAND SWITCH



CONNECTX SMARTNIC



BLUEFIELD DPU



SPECTRUM ETHERNET SWITCH

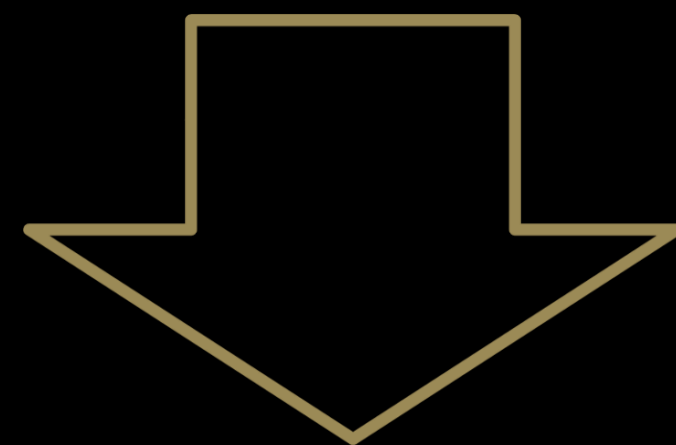


UFM/NetQ



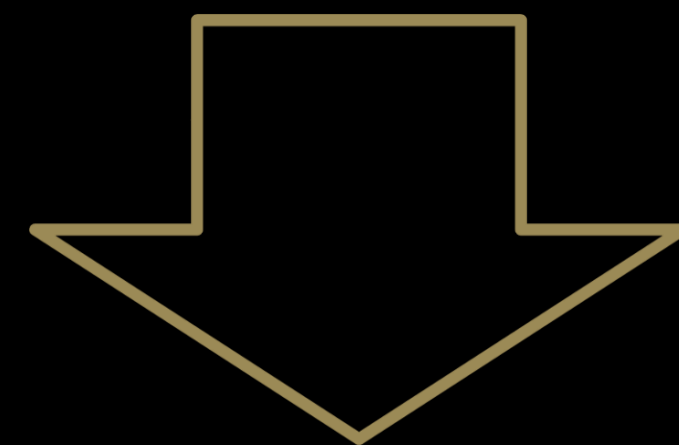
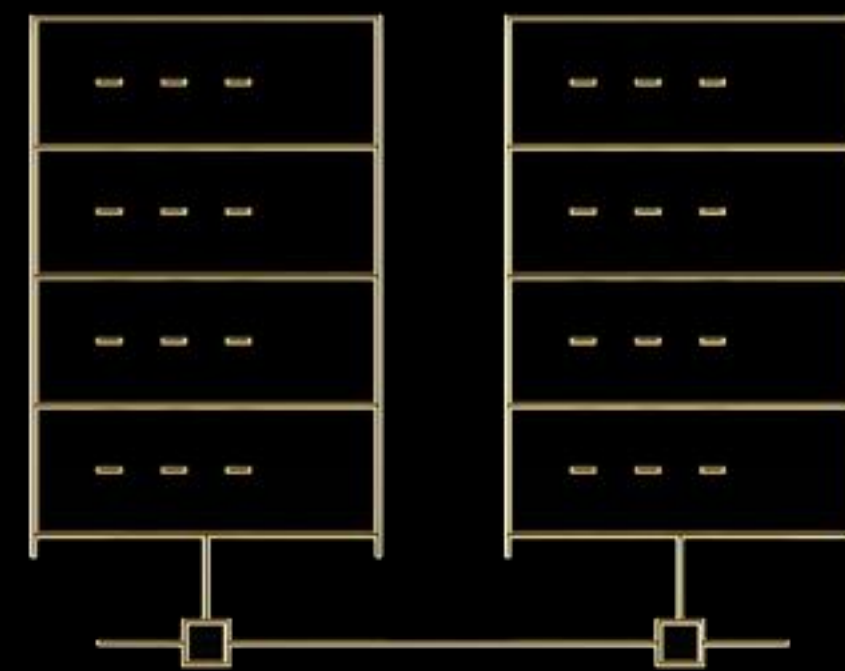
# HPC Performance Bottlenecks

Overlapping



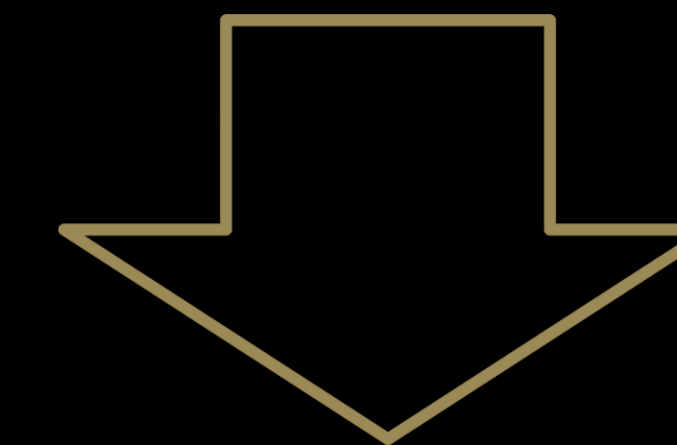
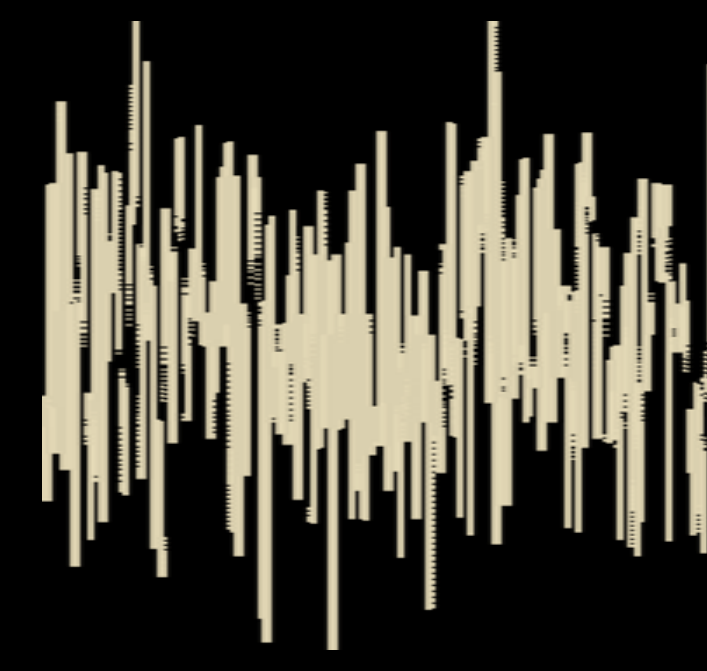
In-Network Computing  
Asynchronous Progress  
(Compute – Communication Overlap)

Load Imbalanced



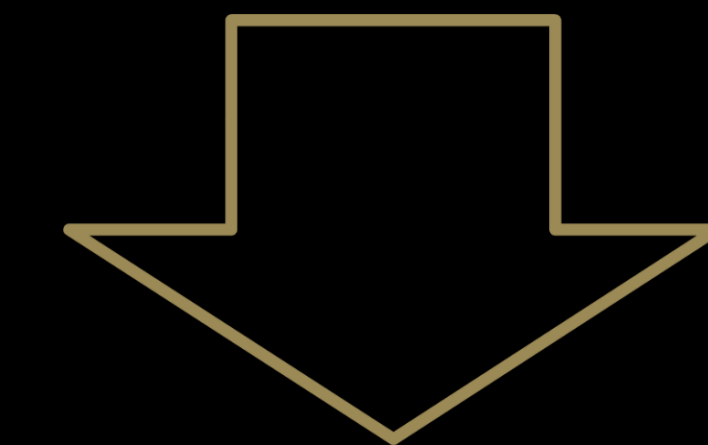
In-Network Computing  
and DPU Synchronization

Jitter



In-Network Computing  
Infrastructure Processing

Multi-Job Performance

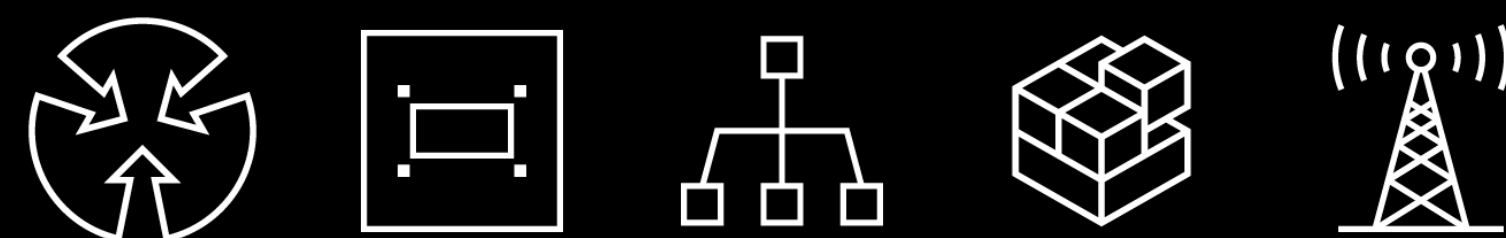


Adaptive Performance  
Isolation



# BlueField Data Processing Unit

SOFTWARE DEFINED NETWORKING



SOFTWARE DEFINED SECURITY



SOFTWARE DEFINED STORAGE



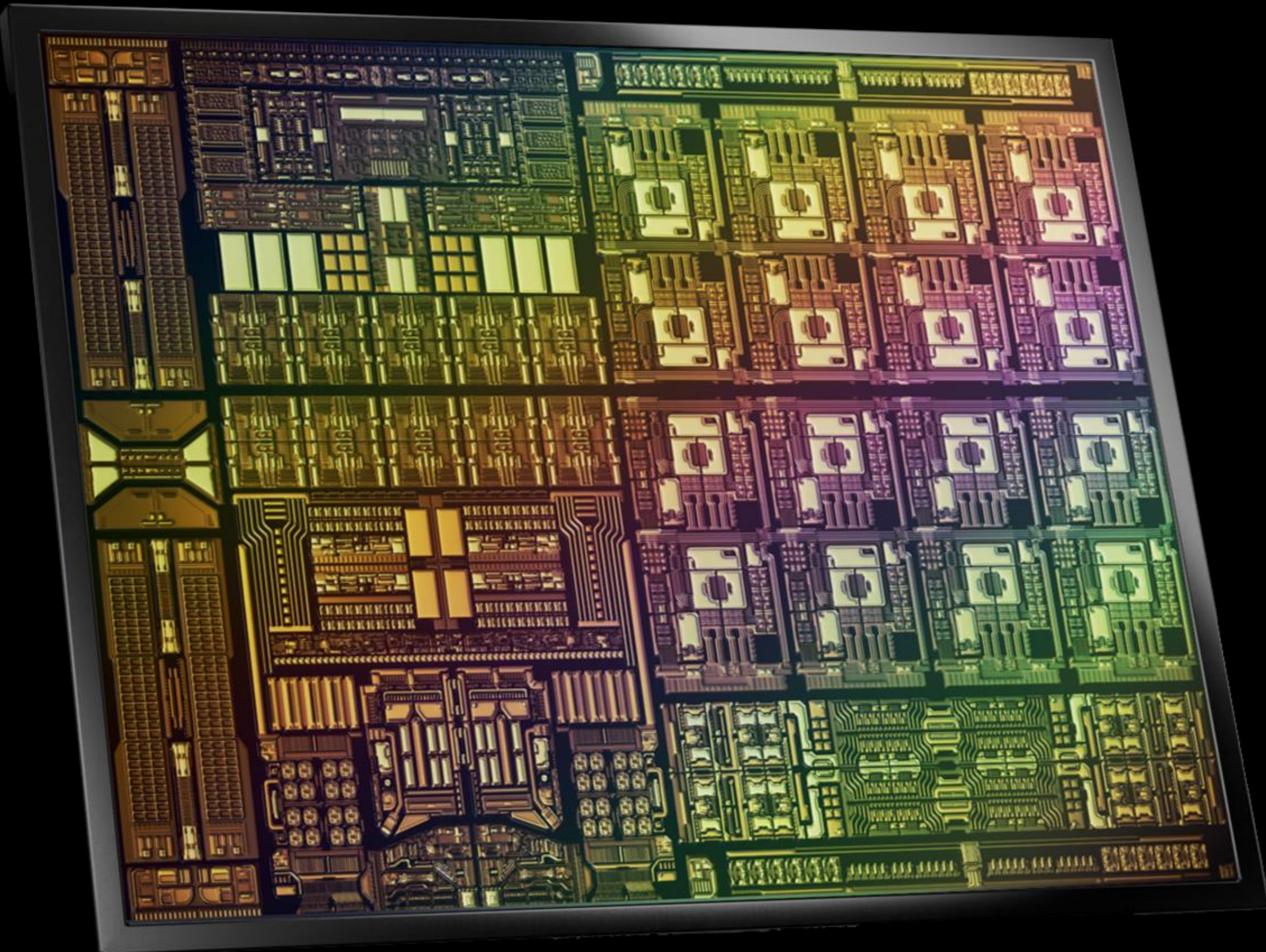
Infrastructure Services

## Data Center on a Chip

- 16 Arm 64-Bit Cores
- 16 Core / 256 Threads Datapath Accelerator
- ConnectX InfiniBand / Ethernet
- DDR memory interface
- PCIe switch



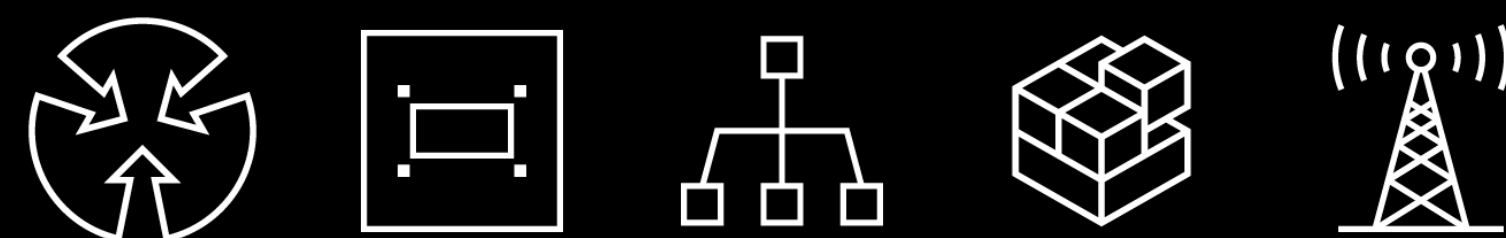
BlueField Infrastructure  
Compute Platform





# The New Computing Platform for the Data Center Infrastructure

## SOFTWARE DEFINED NETWORKING



## SOFTWARE DEFINED SECURITY



## SOFTWARE DEFINED STORAGE

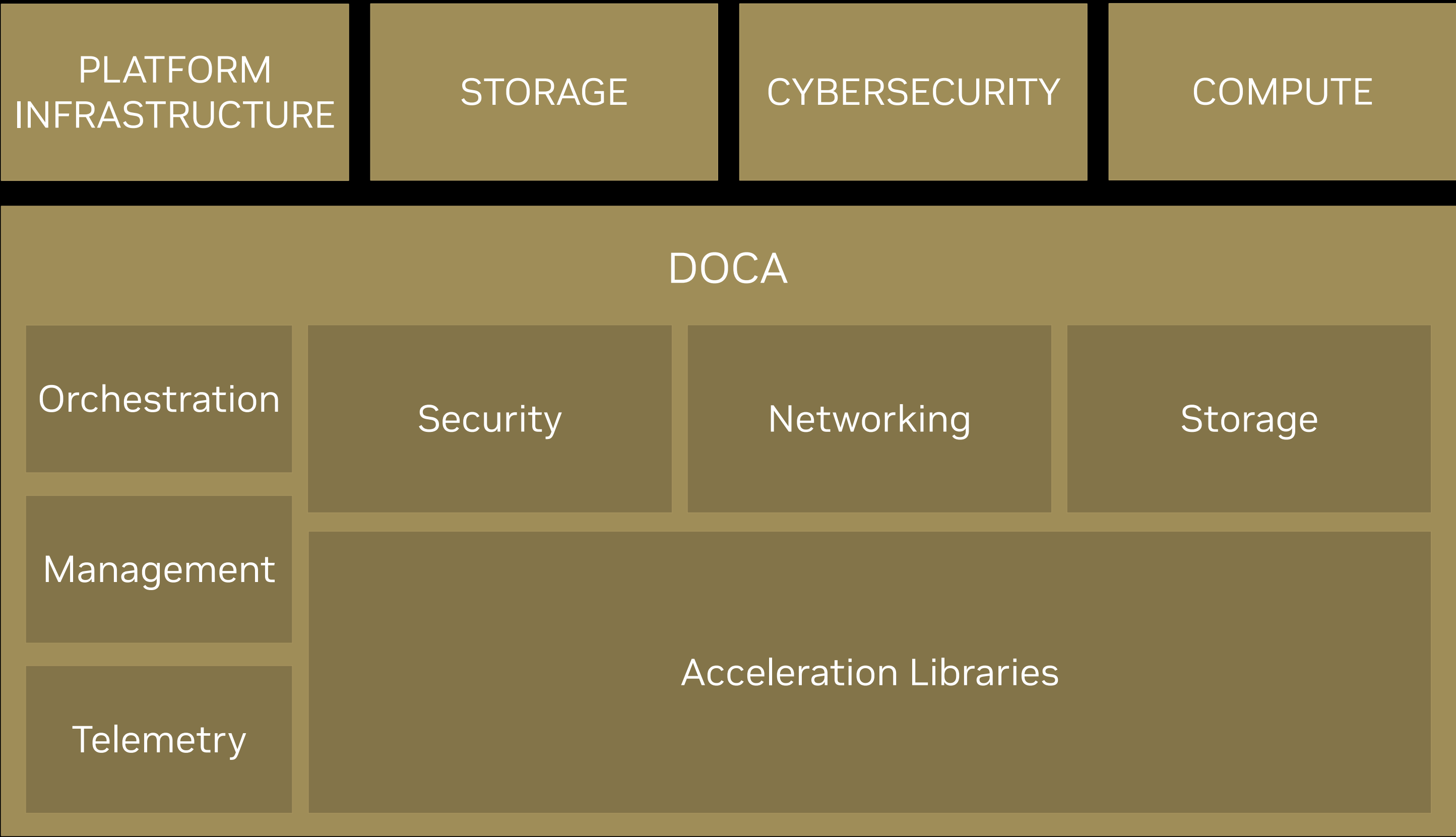


Infrastructure Services

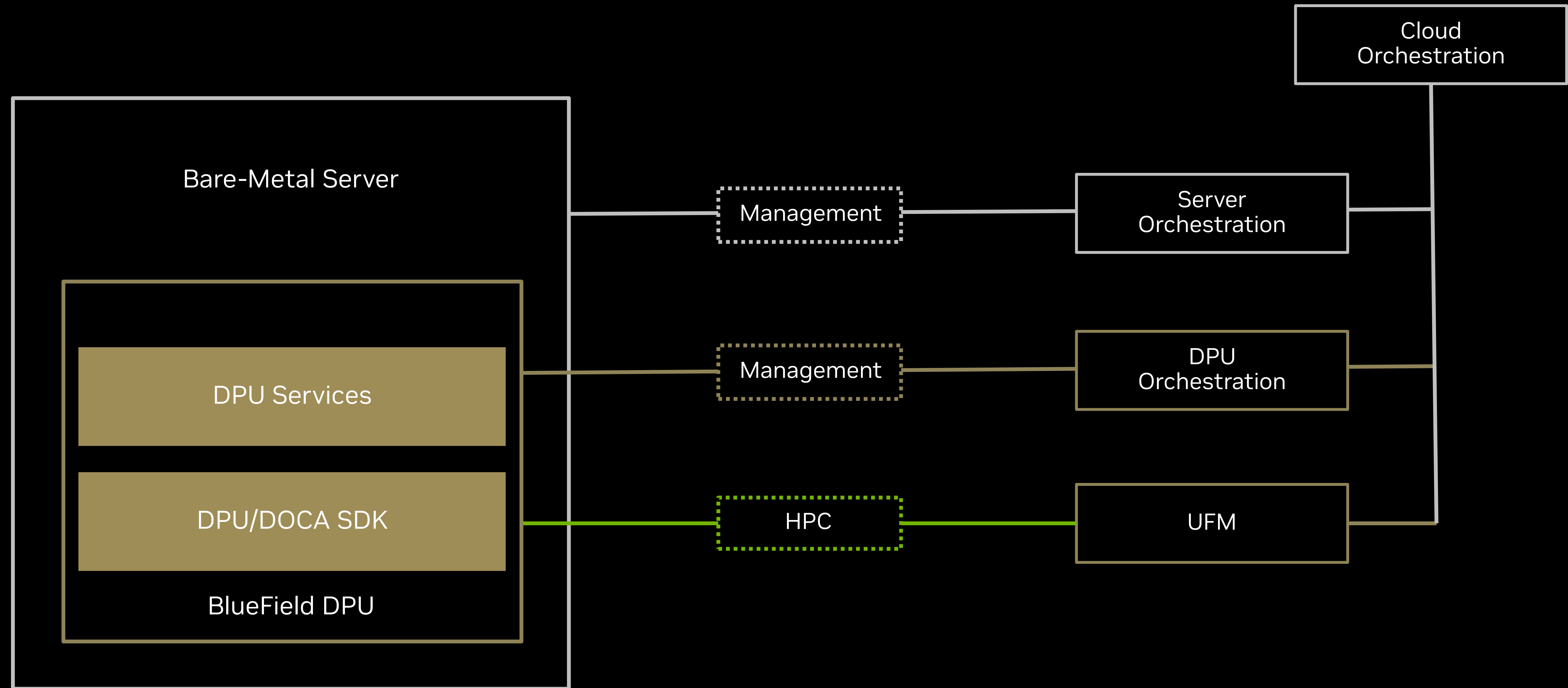
	BlueField-2	BlueField-3
Network Bandwidth	200Gb/s	400Gb/s
RDMA msg rate	215Mpps	370Mpps
Compute	SPECINT2K17: 9.8	SPECINT2K17: 42
Memory Bandwidth	17GB/s	80GB/s



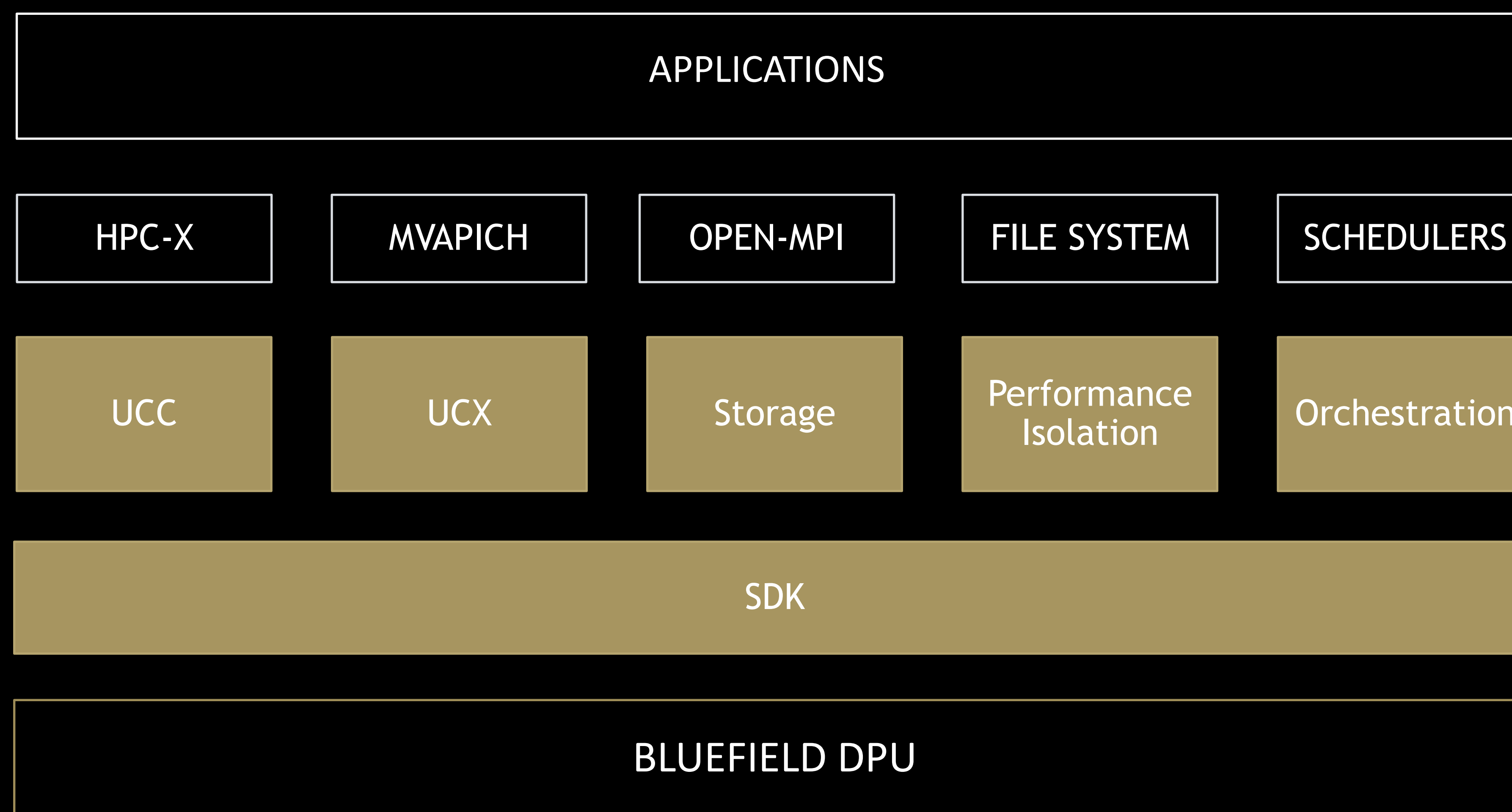
BlueField Infrastructure Compute Platform



# Delivering Cloud Native Supercomputing

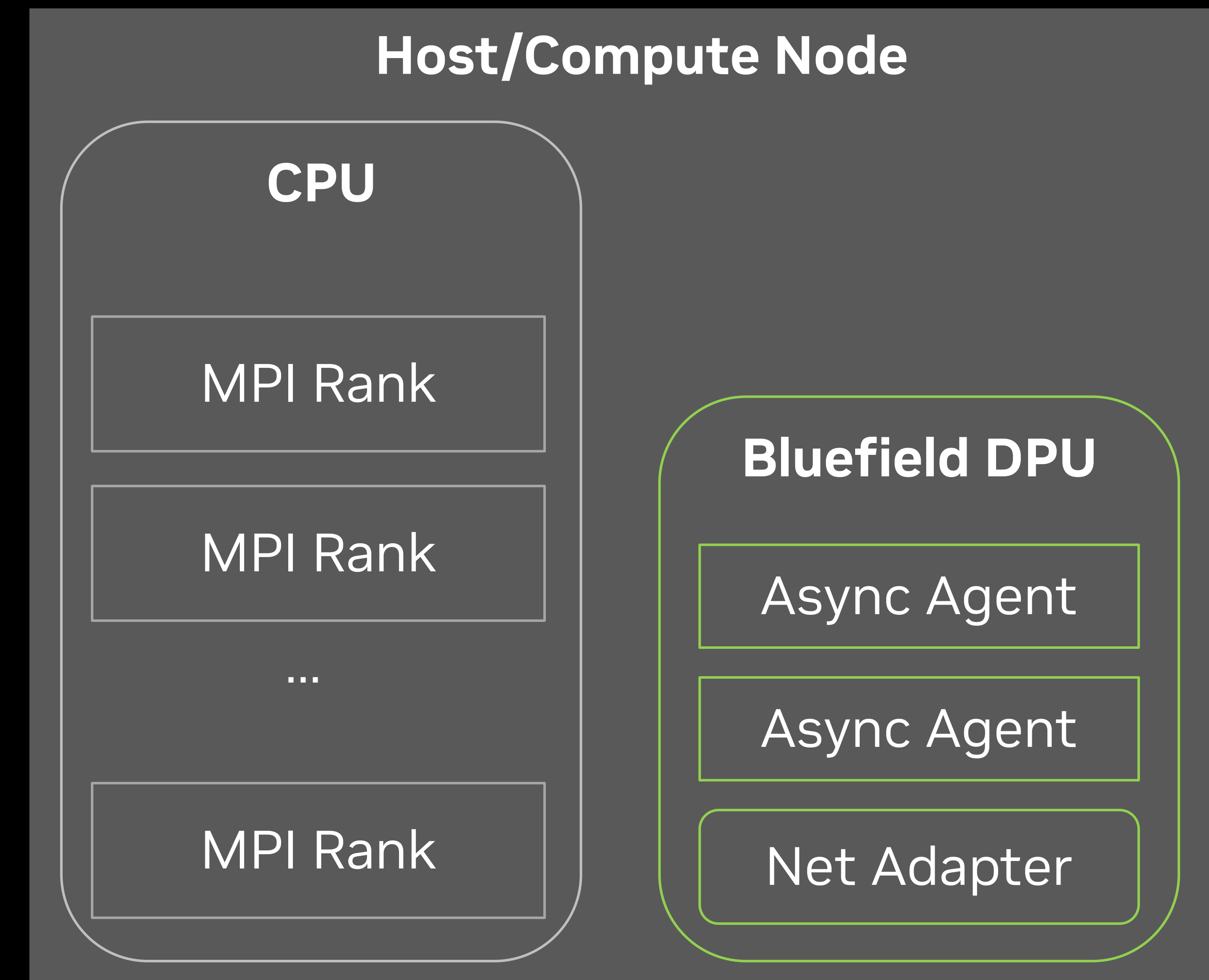


# Accelerating HPC Applications with DPU/DOCA Services



# High Level System Components from Software's Perspective

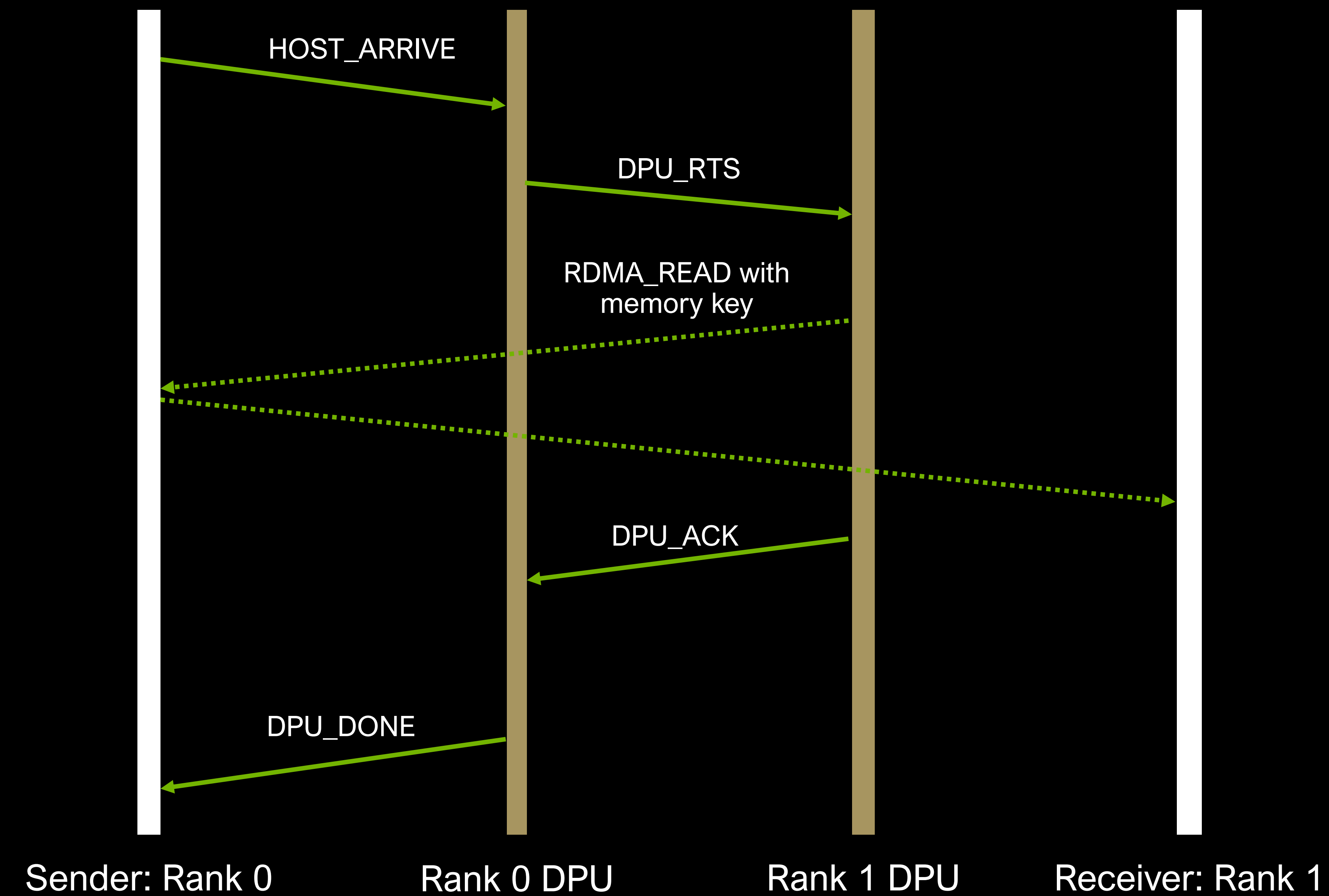
- Host paired with local DPU
- Local DPU runs service processes (SP)
  - Each local user process (such as MPI process) has a service process that it is pair with
  - Each service process serves multiple local processes
  - Algorithm is split between host and DPU – blocking and nonblocking may have different 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

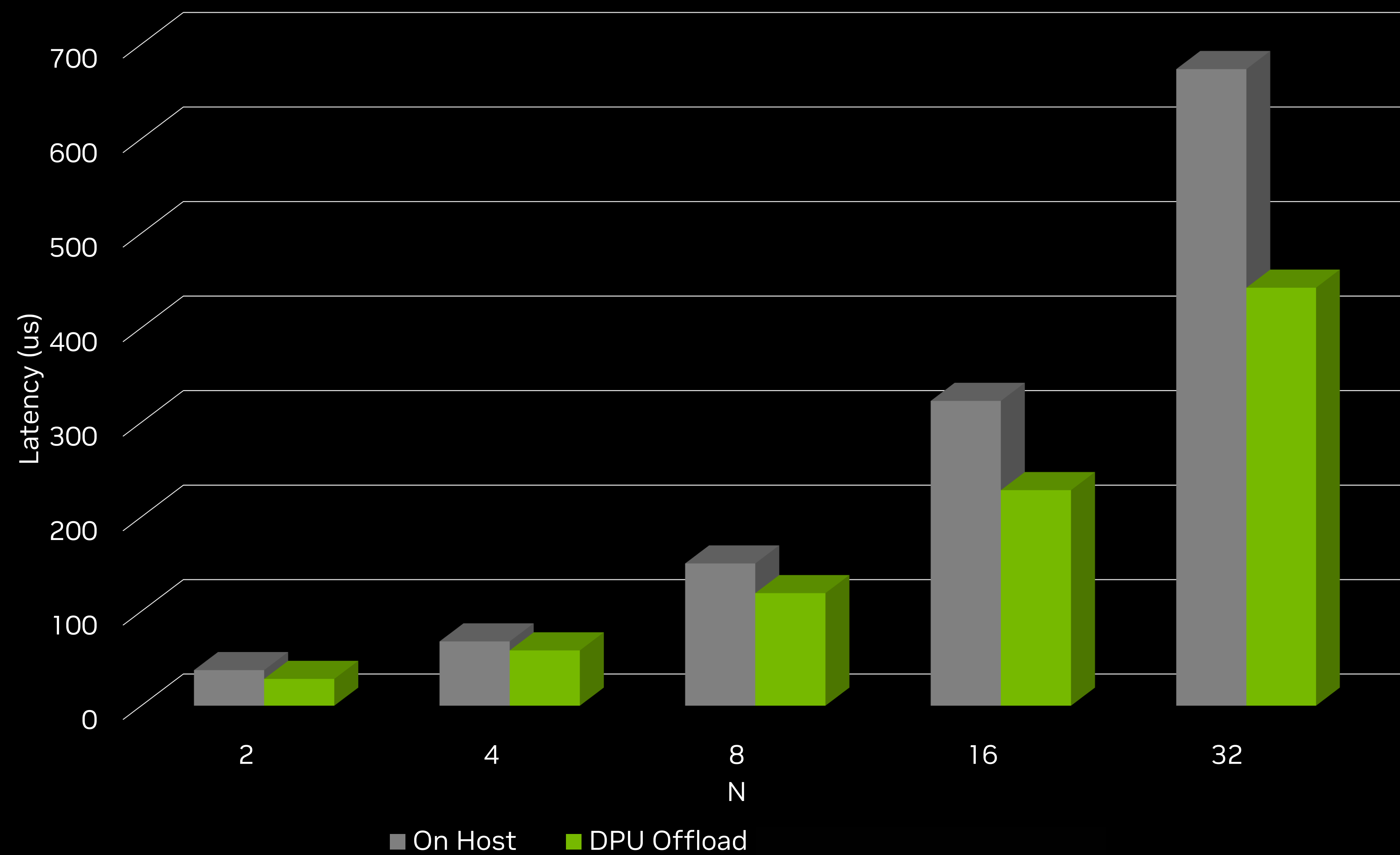
An Element of Collective Algorithm



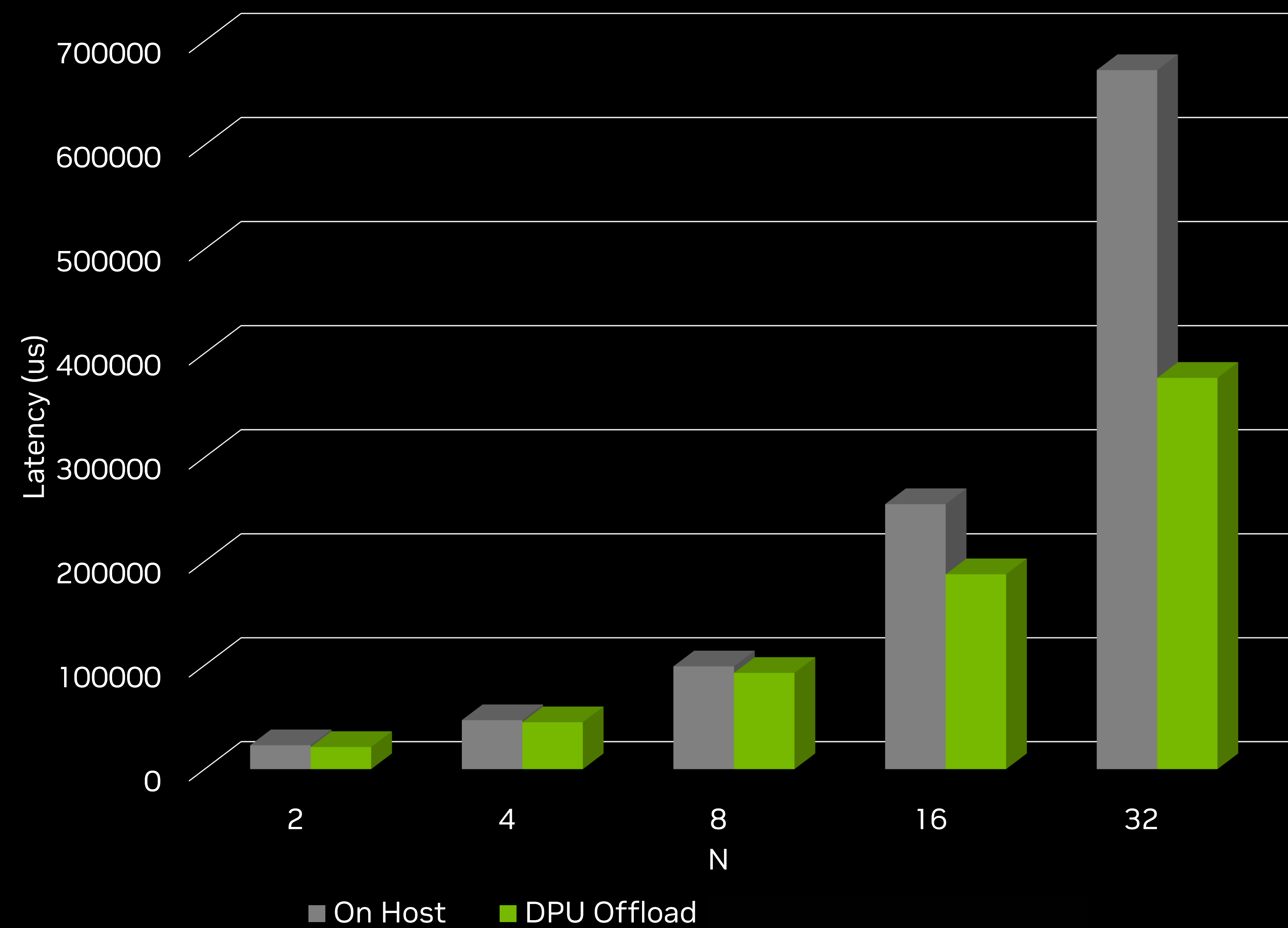


# Alltoallv Latency

OSU Alltoallv 1 PPN, Size = 128 KB



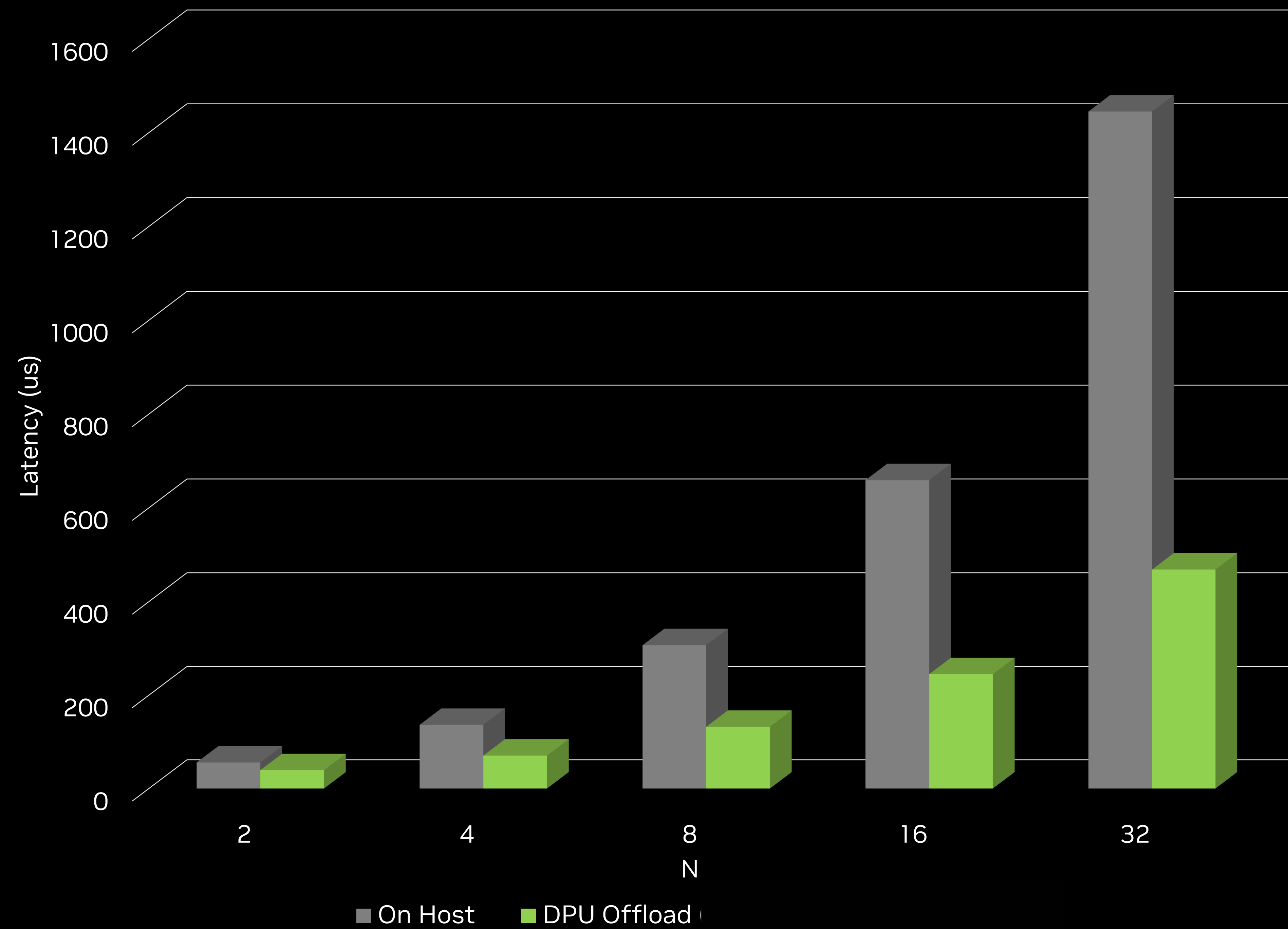
OSU Alltoallv 32 (full) PPN, Size = 128 KB



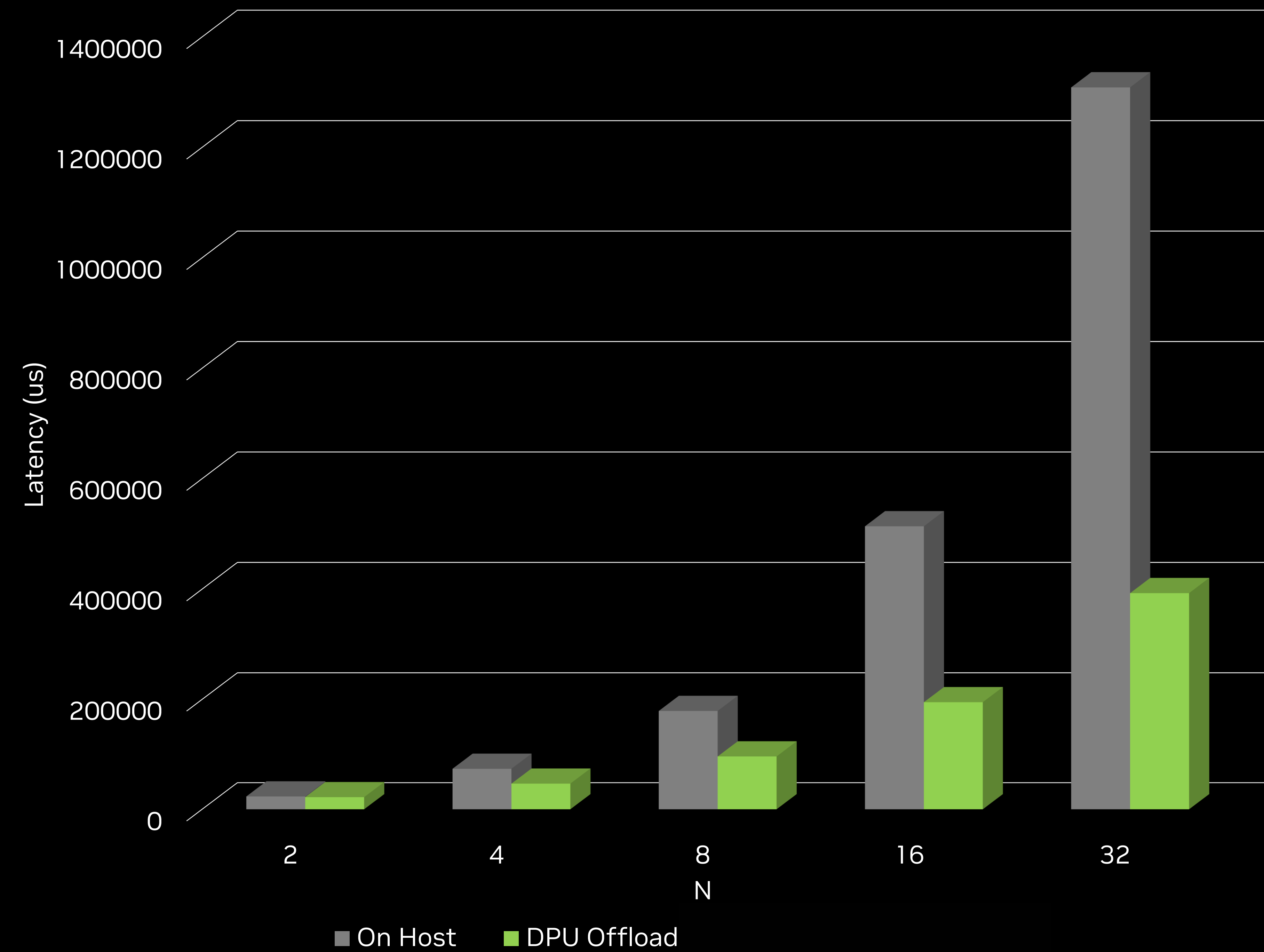


# iAlltoallv latency

OSU iAlltoallv 1 PPN, Size = 128 KB



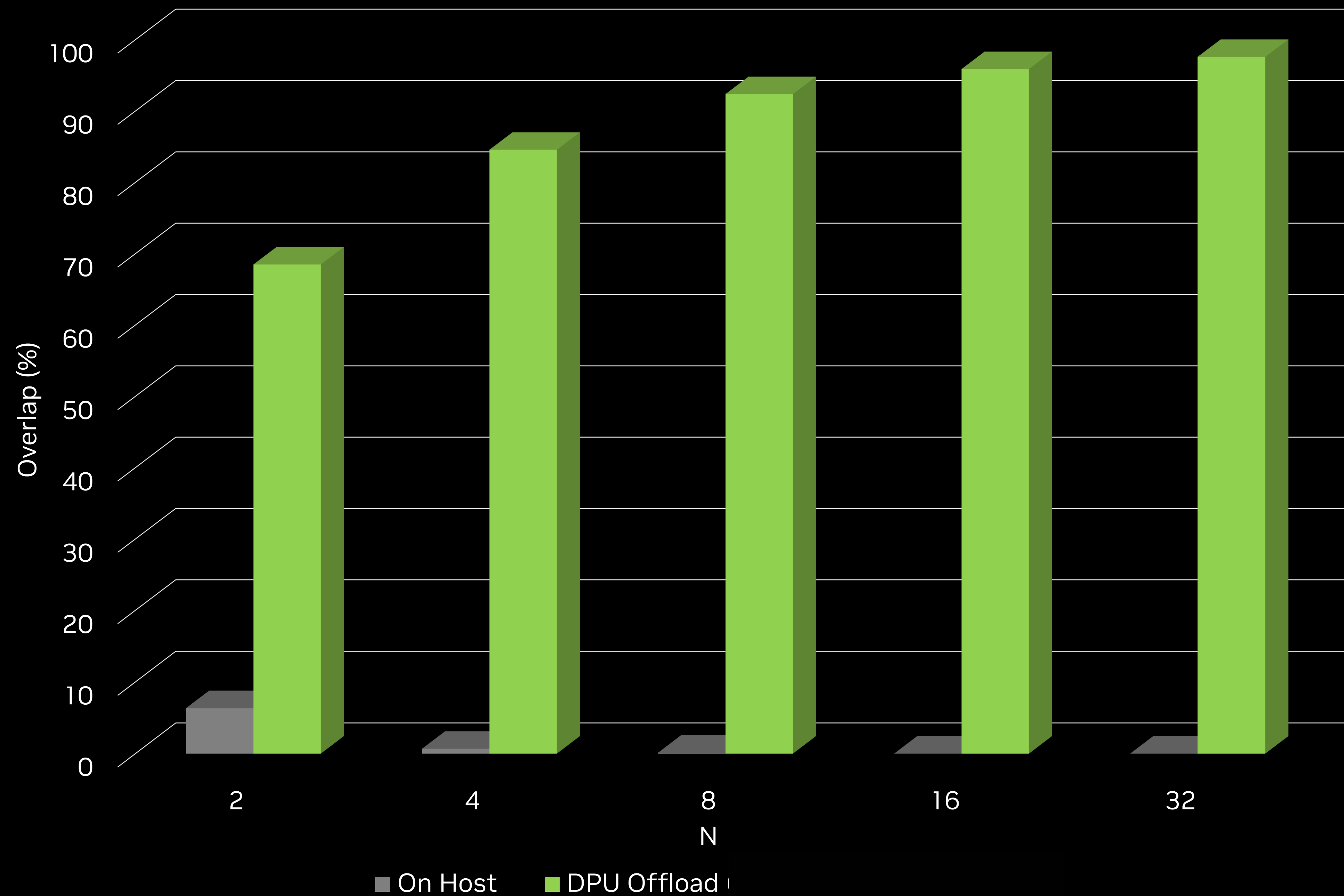
OSU iAlltoallv 32 (full) PPN, Size = 128 KB



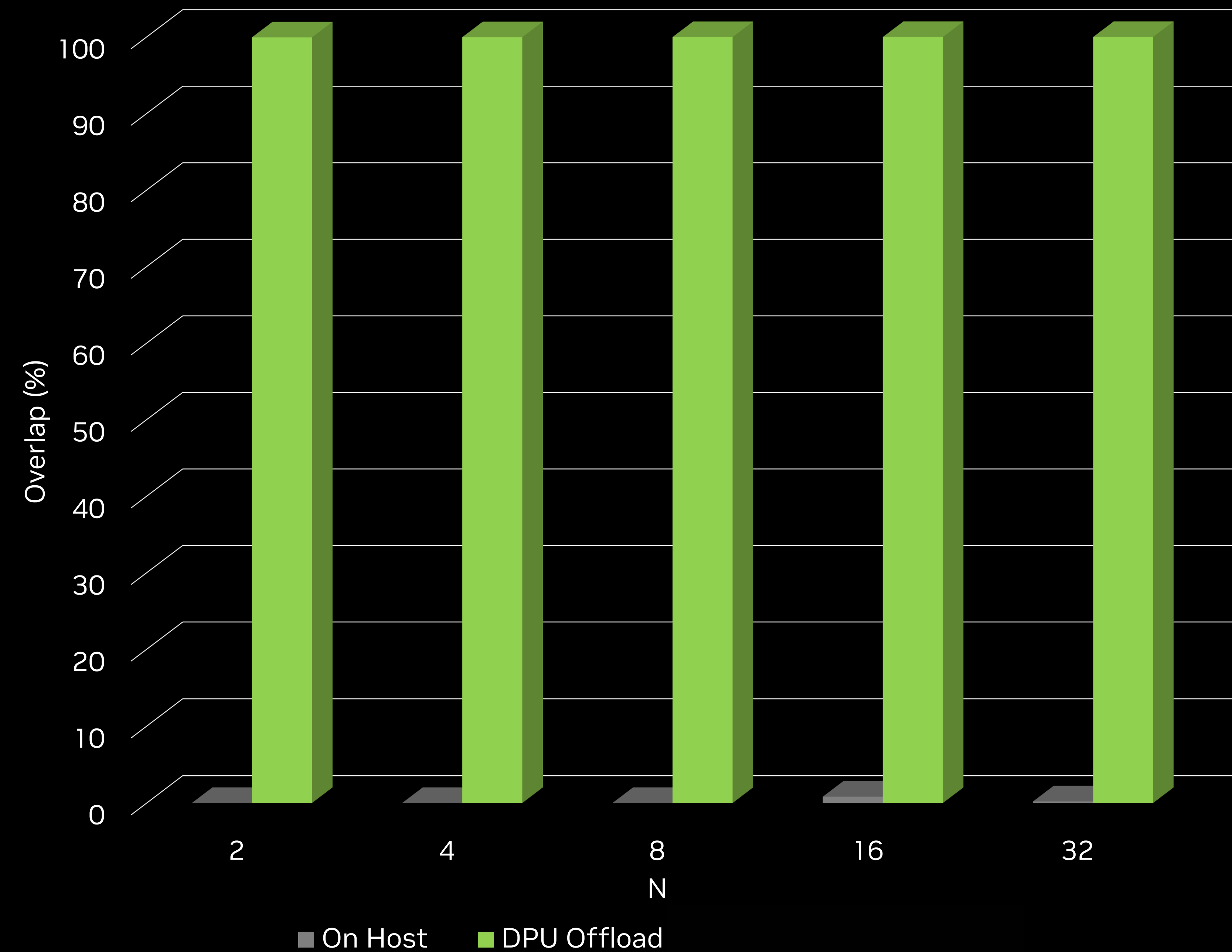


# iAlltoallv compute/communication overlap

OSU iAlltoallv 1 PPN, Size = 128 KB



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

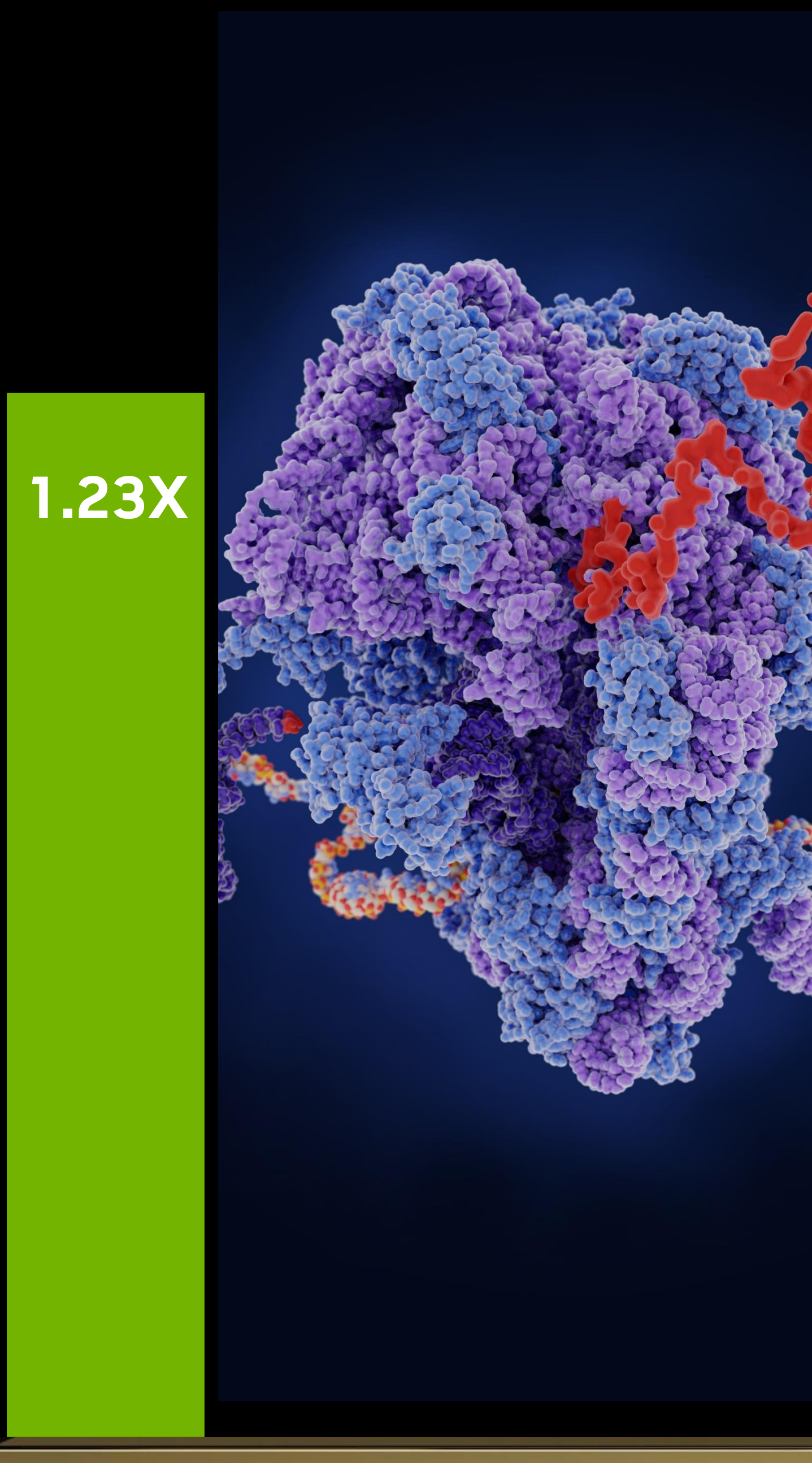




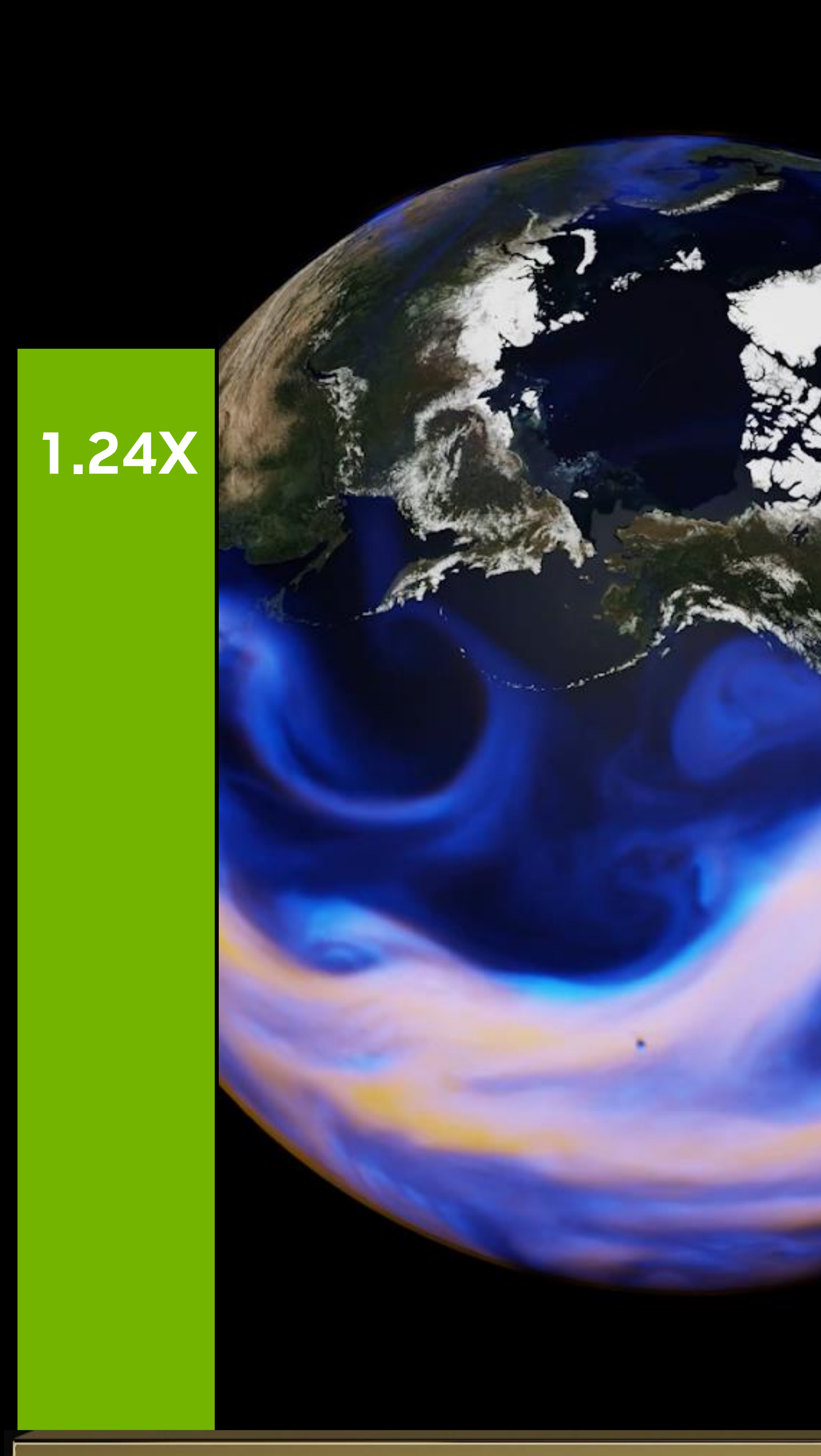
# HIGHER APPLICATION PERFORMANCE

With BlueField DPU and Quantum InfiniBand In-Network Computing

Octopus (Physics / Chemistry)



GSAM (Weather)



FFT





