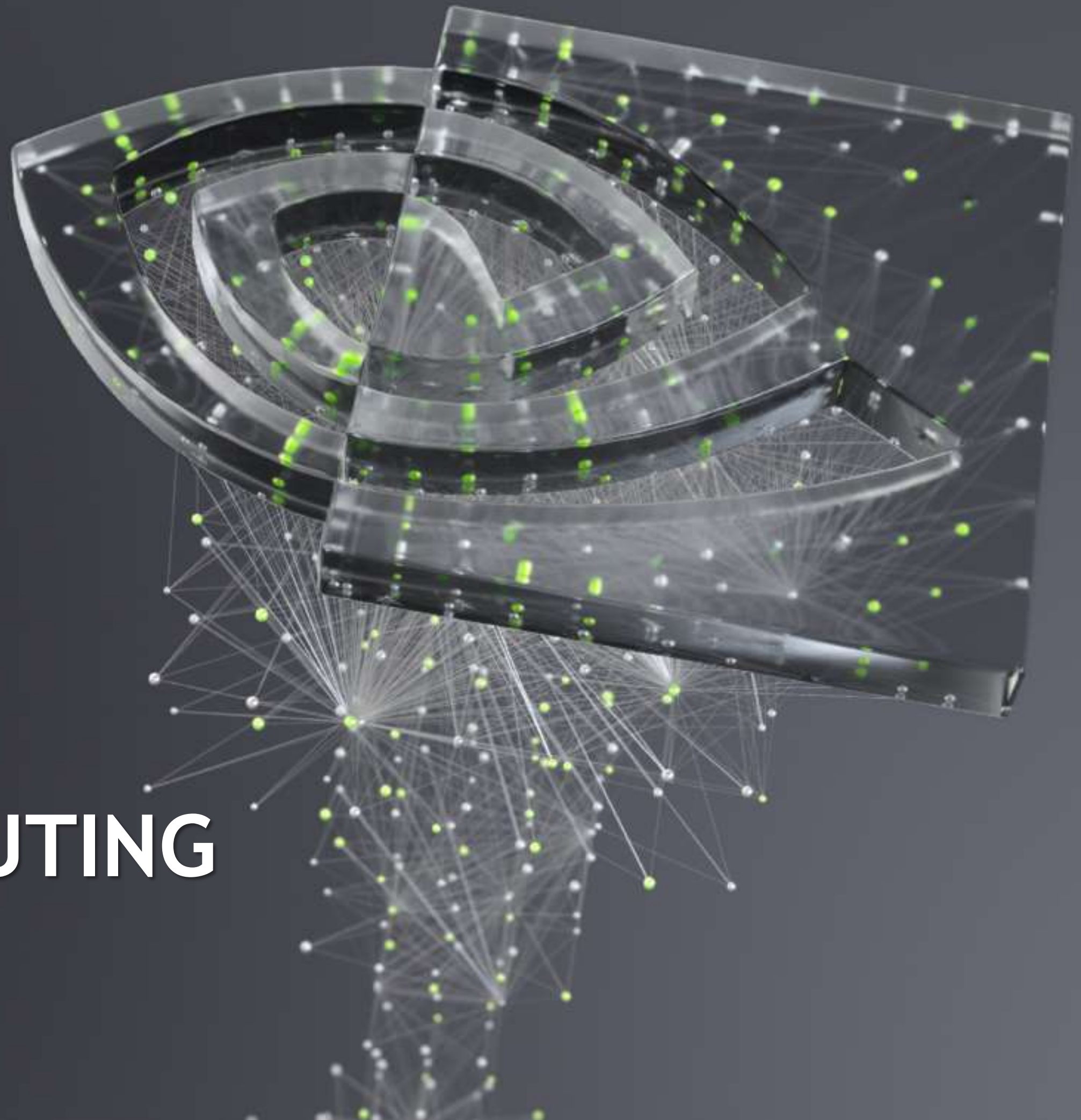


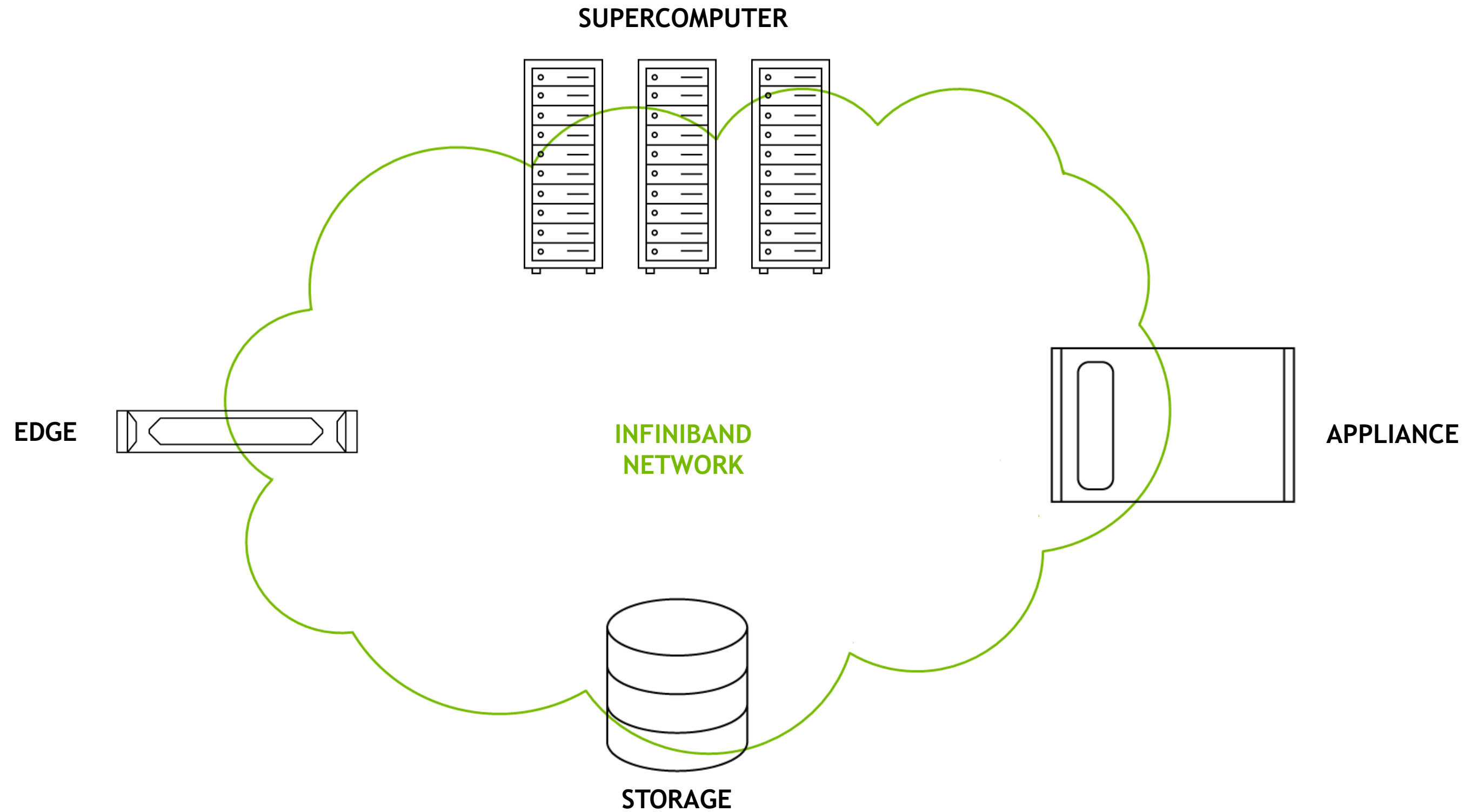


INFINIBAND IN-NETWORK COMPUTING

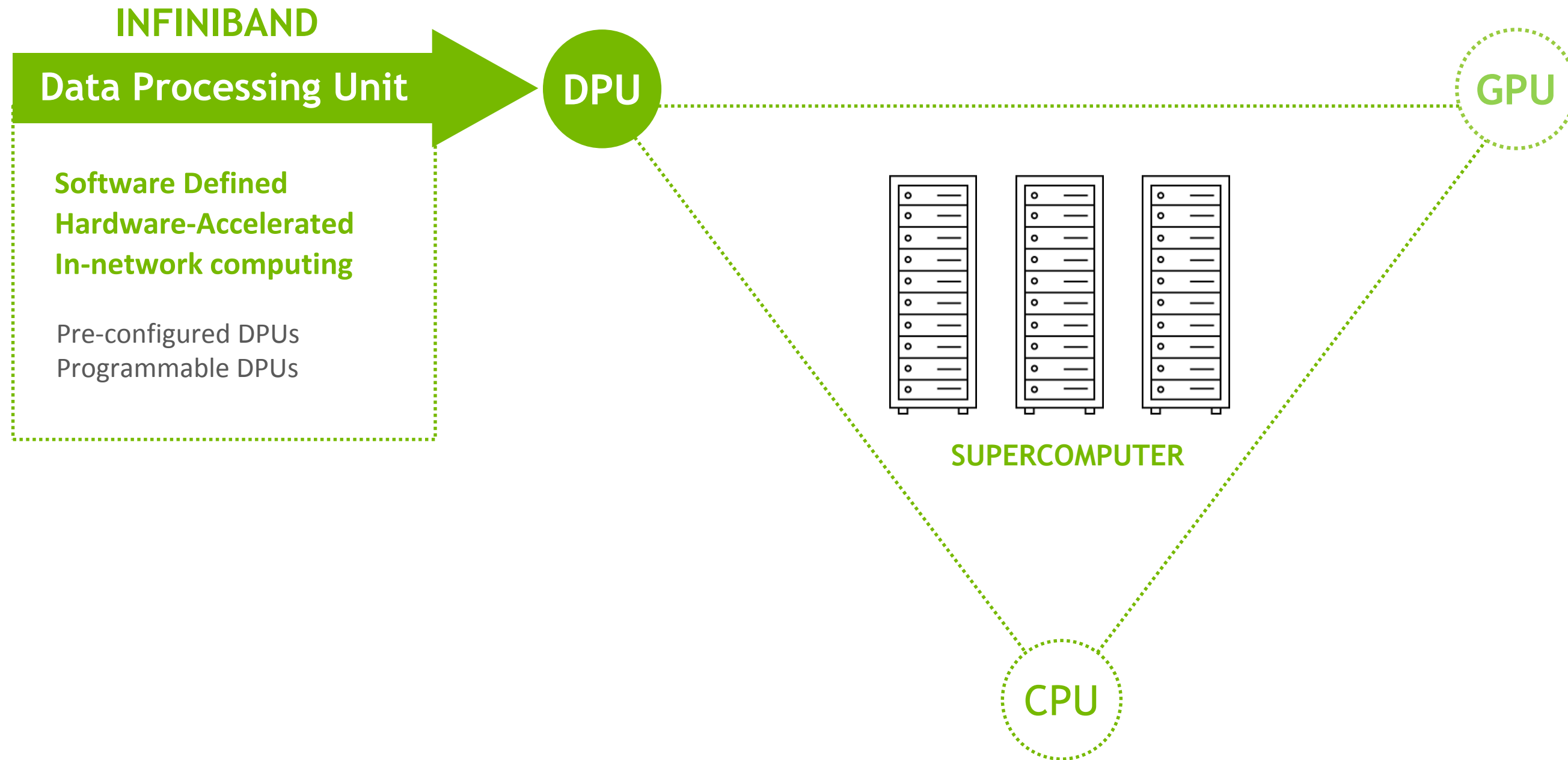
MUG, August 2020



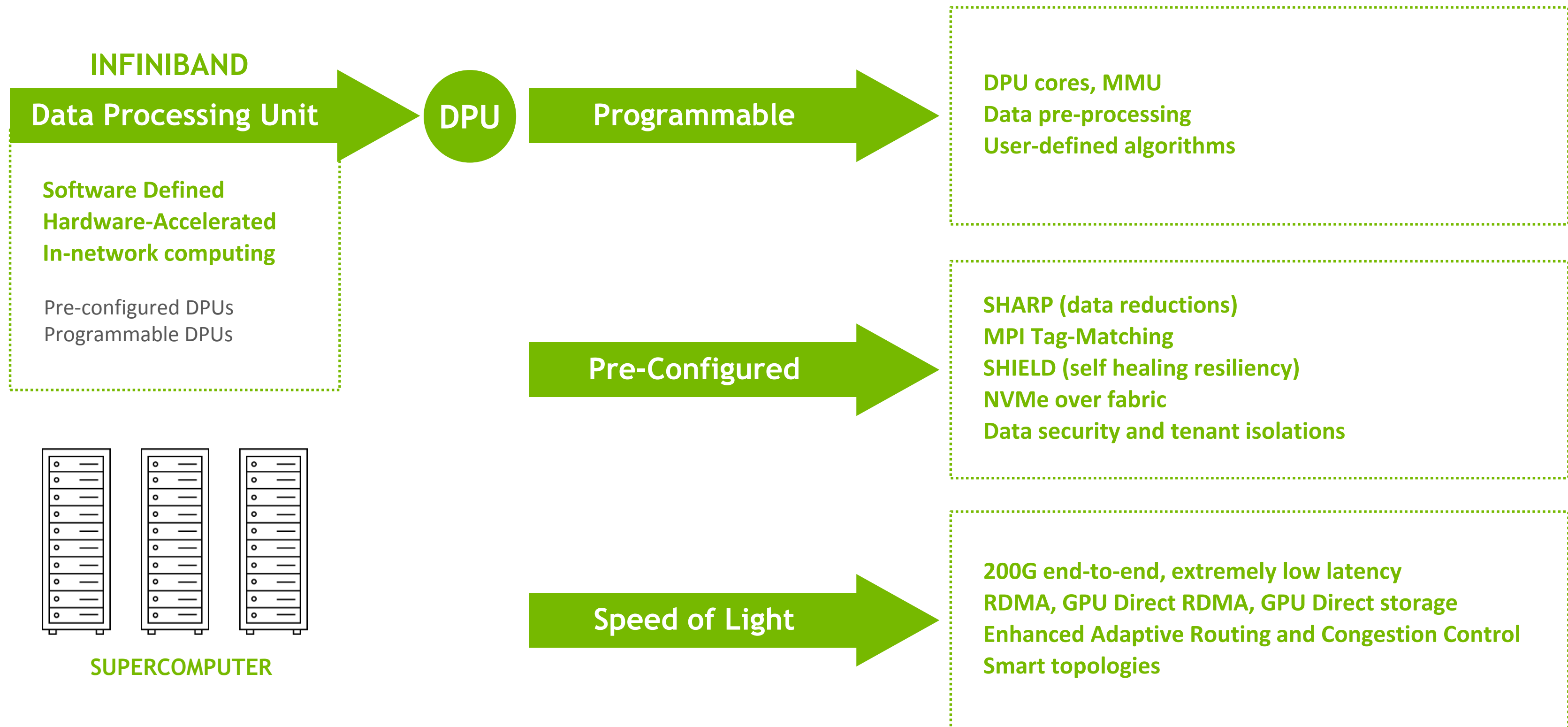
THE NEW SCIENTIFIC COMPUTING WORLD



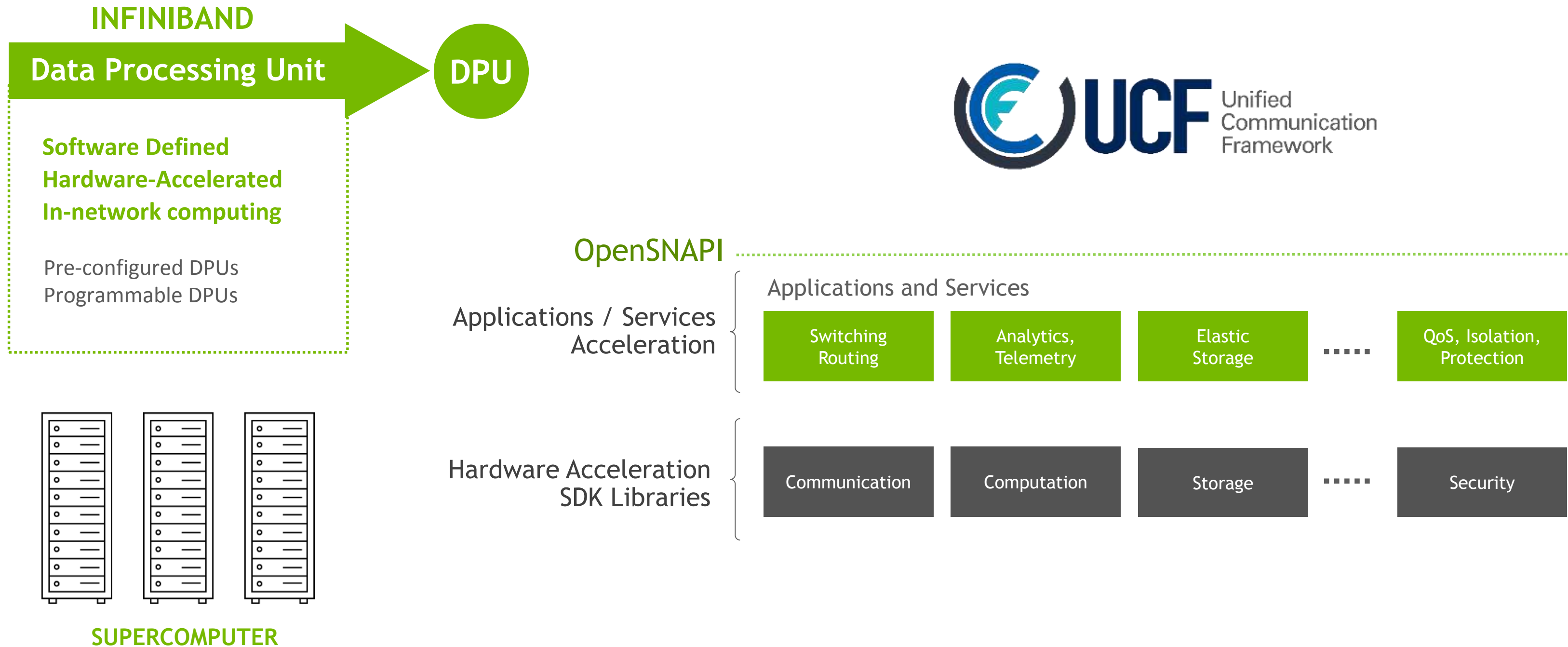
IN-NETWORK COMPUTING ACCELERATED SUPERCOMPUTING



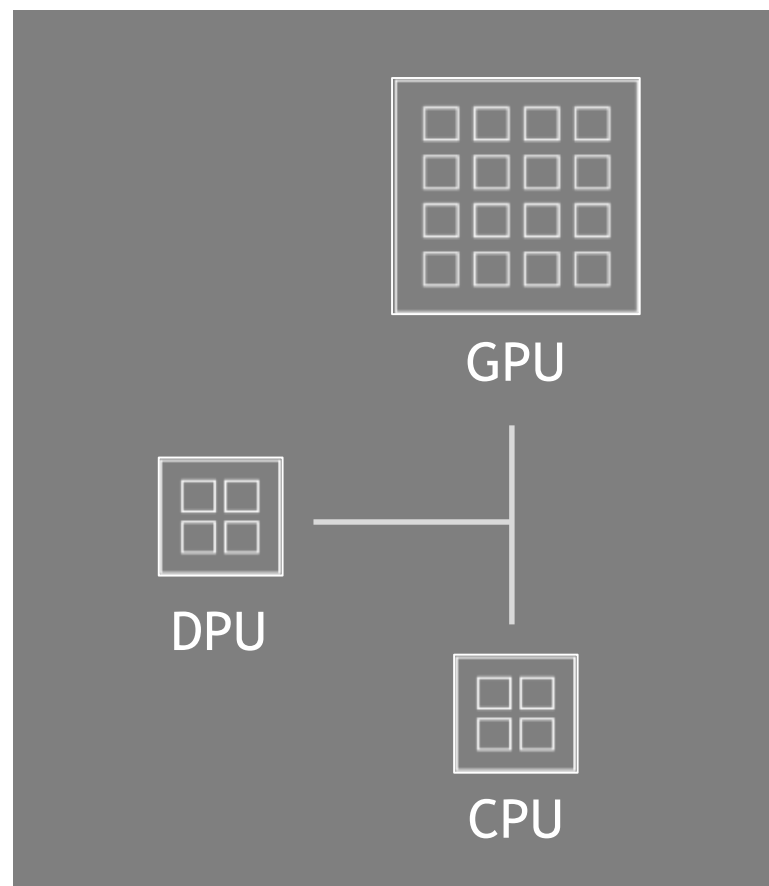
IN-NETWORK COMPUTING ACCELERATED SUPERCOMPUTING



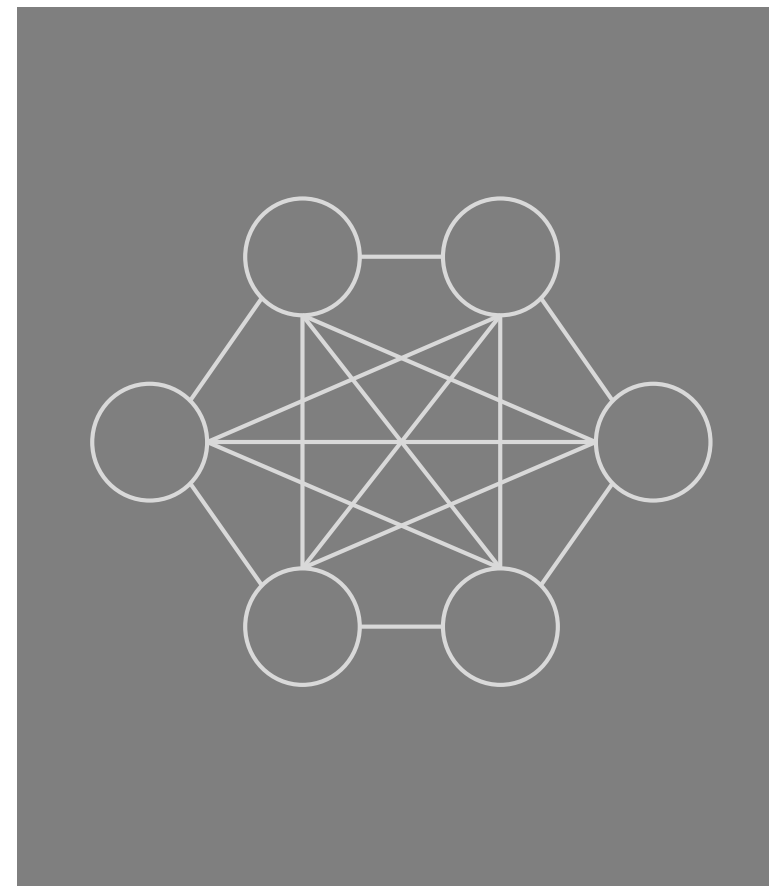
IN-NETWORK COMPUTING ACCELERATED SUPERCOMPUTING



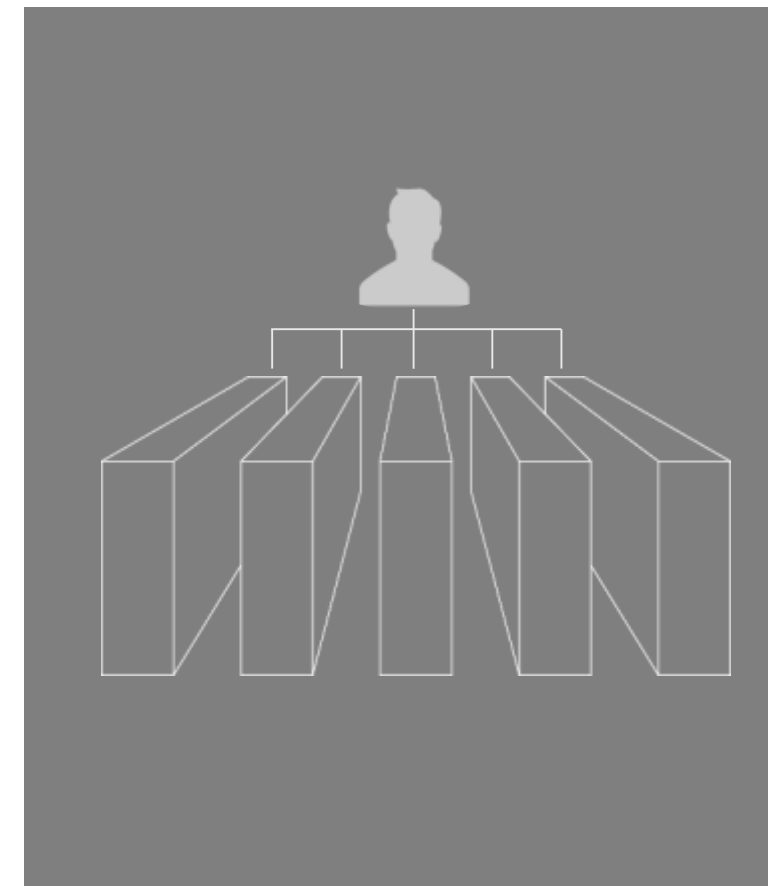
INFINIBAND TECHNOLOGY FUNDAMENTALS



Smart End-Point



Architected to Scale



Centralized Management



Standard

HDR 200G INFINIBAND ACCELERATES NEXT GENERATION HPC AND AI SUPERCOMPUTERS (EXAMPLES)



9 PetaFLOPS
3K HDR Nodes
Dragonfly+ Topology



7.5 PetaFLOPS
2K HDR Nodes
Dragonfly+ Topology



3K HDR Nodes
16 PetaFLOPS
Dragonfly+ Topology



ECMWF
EUROPEAN CENTRE FOR MEDIUM RANGE WEATHER FORECASTS

8K HDR Nodes
Dragonfly+ Topology



35.5 PetaFLOPS
2K HDR Nodes
Fat-Tree Topology



23 PetaFLOPS
5.6K HDR Nodes
Dragonfly+ Topology



27.6 PetaFLOPS
3K HDR Nodes
Fat-Tree Topology



HPC/AI Cloud
HDR InfiniBand



SDSC
SAN DIEGO SUPERCOMPUTER CENTER



PURDUE
UNIVERSITY

HDR
Supercomputers



23.5 PetaFLOPS
8K HDR Nodes
Fat-Tree Topology



SCALABLE HIERARCHICAL AGGREGATION AND REDUCTION PROTOCOL (SHARP)

SCALABLE HIERARCHICAL AGGREGATION AND REDUCTION PROTOCOL (SHARP)

In-network Tree based aggregation mechanism

Multiple simultaneous outstanding operations

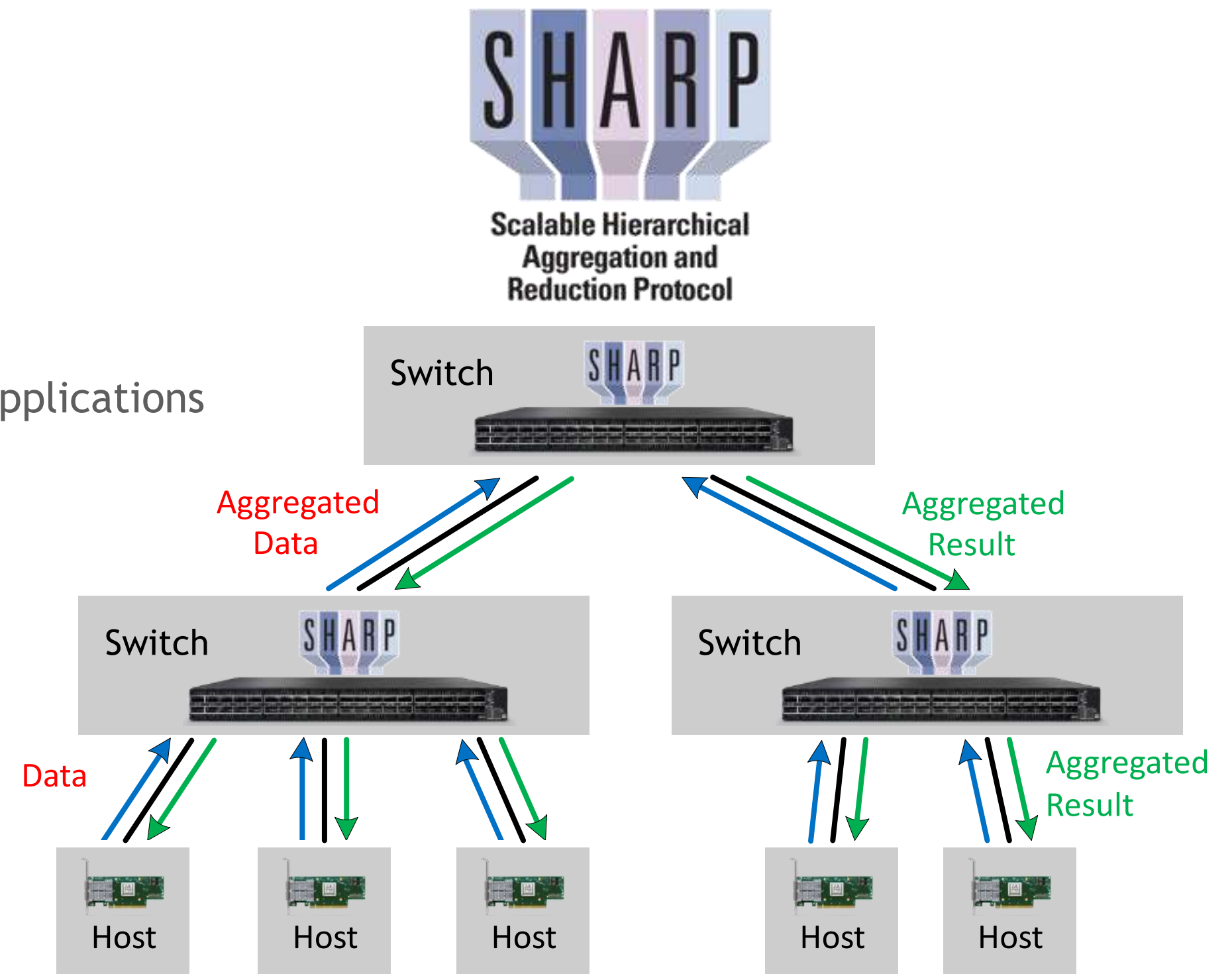
For HPC (MPI / SHMEM) and Distributed Machine Learning applications

Scalable High Performance Collective Offload

Barrier, Reduce, All-Reduce, Broadcast and more

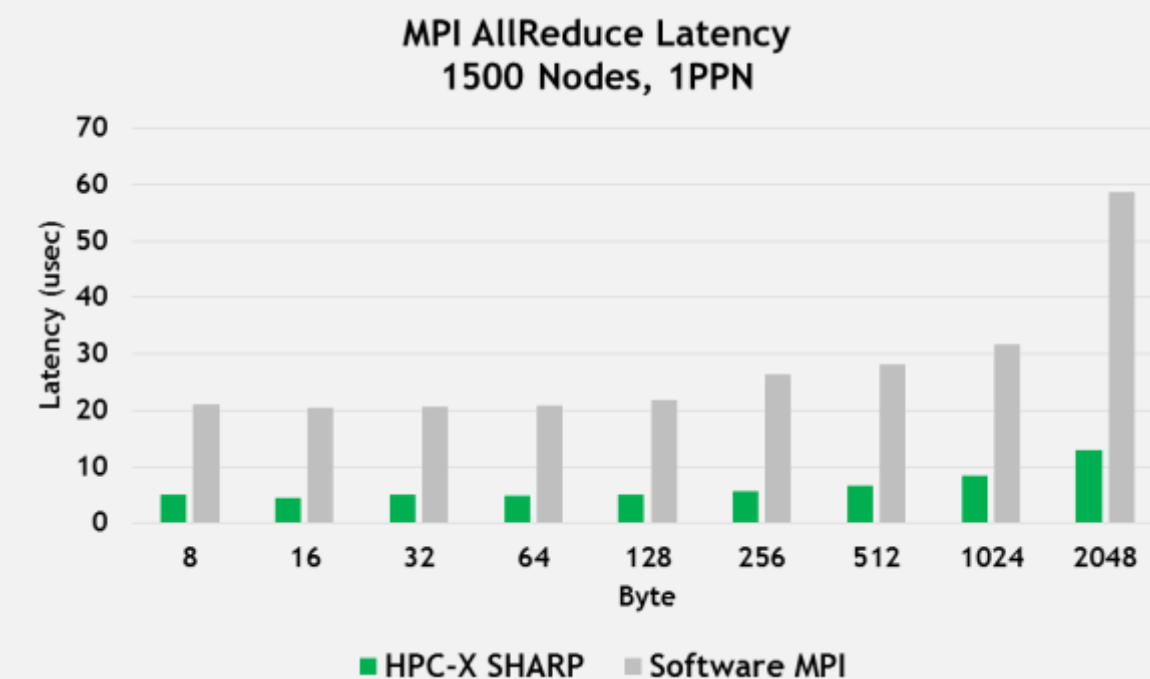
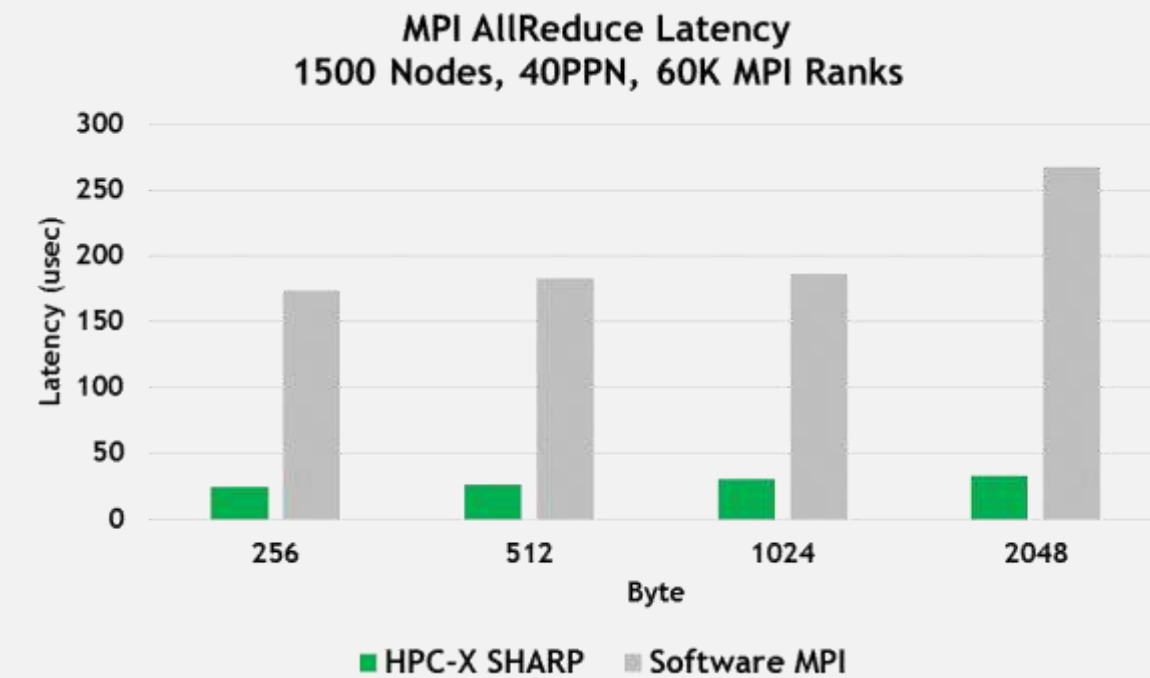
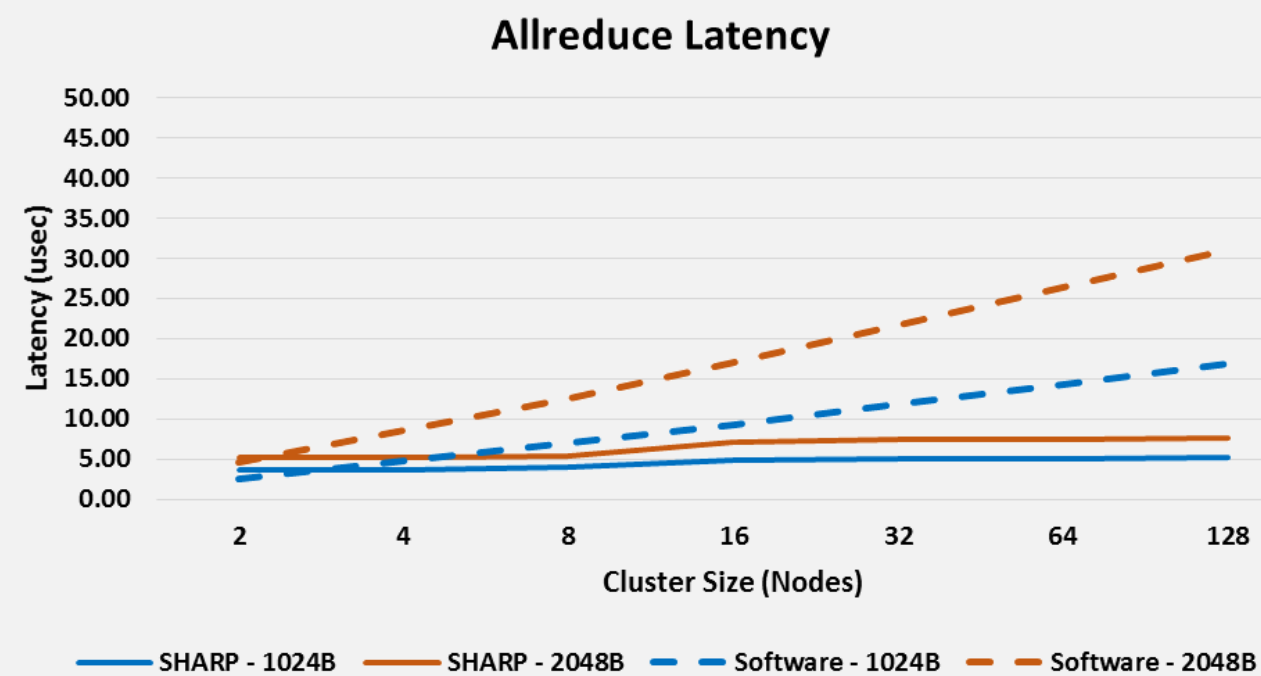
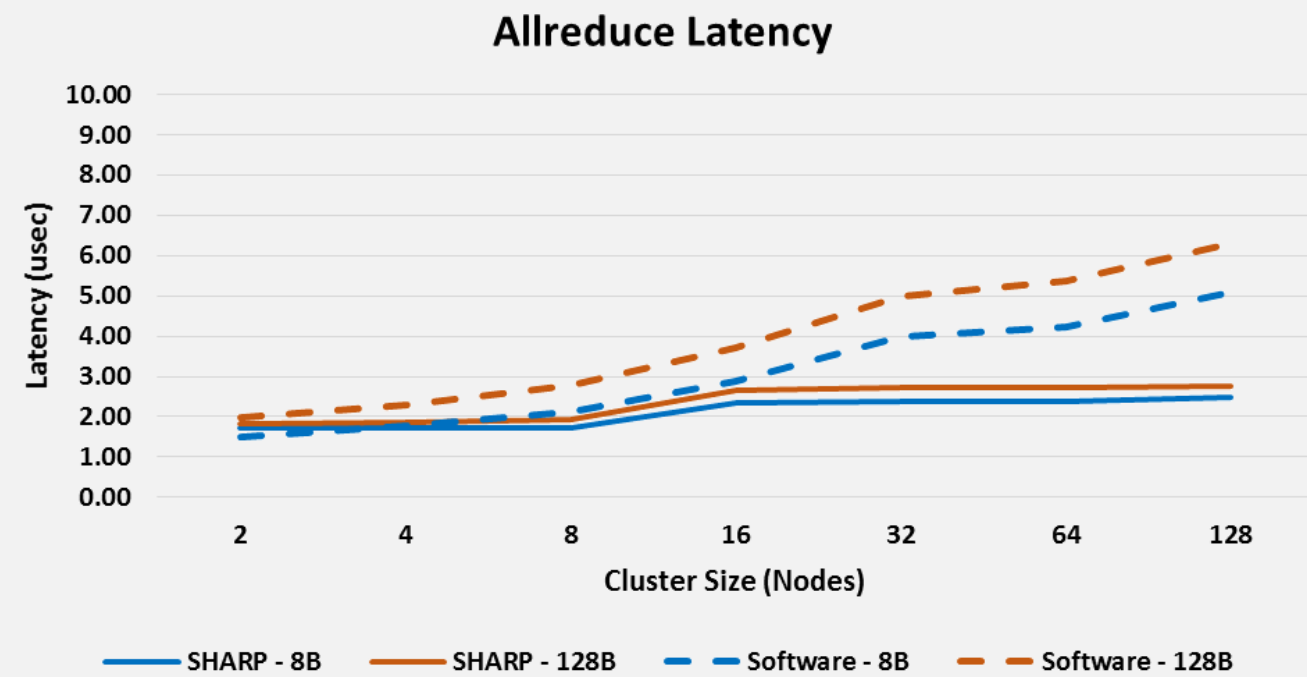
Sum, Min, Max, Min-loc, max-loc, OR, XOR, AND

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



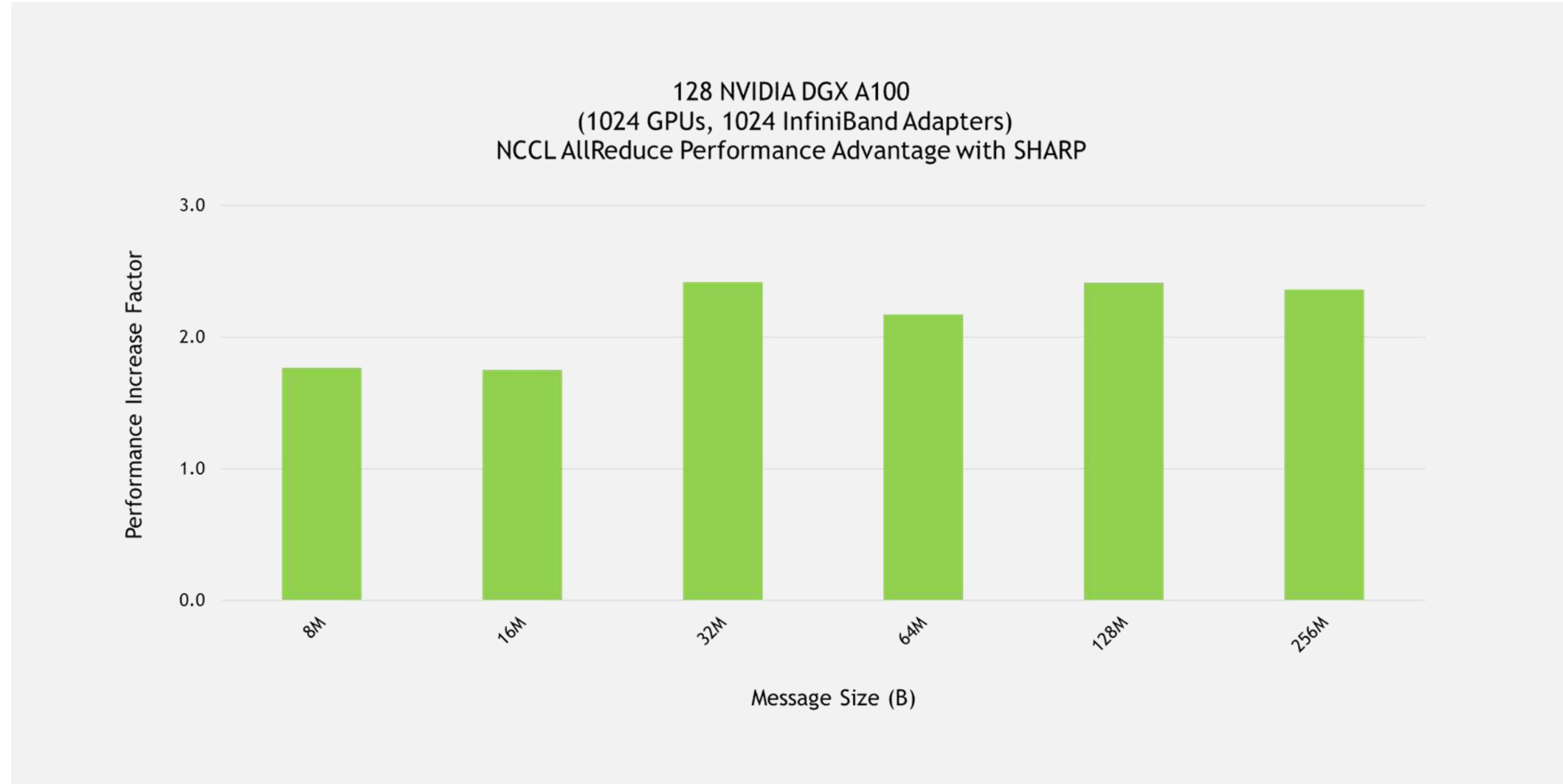
SHARP ALLREDUCE PERFORMANCE ADVANTAGES

Providing Flat Latency, 7X Higher Performance



INFINIBAND SHARP AI PERFORMANCE ADVANTAGE

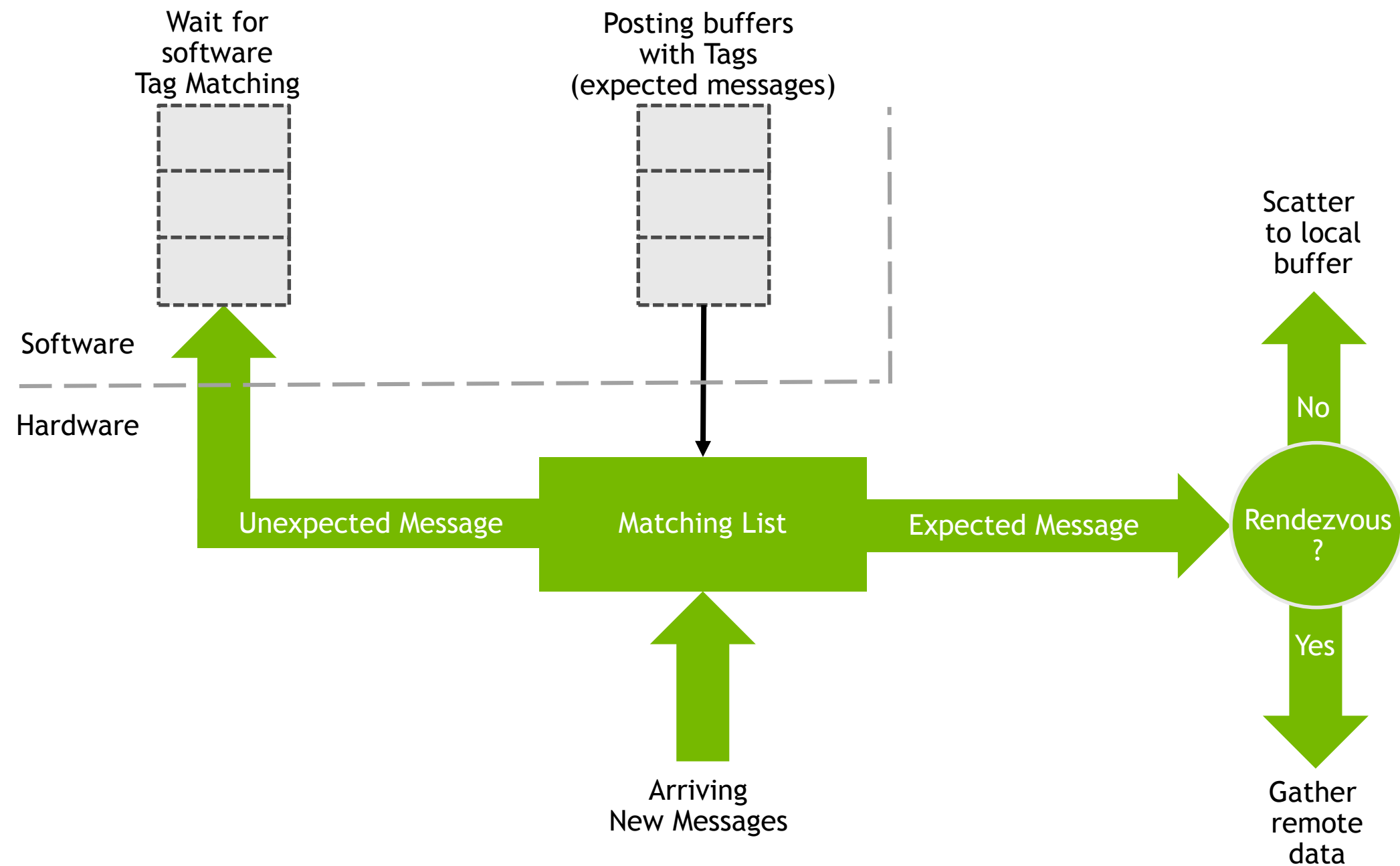
2.5X Higher Performance





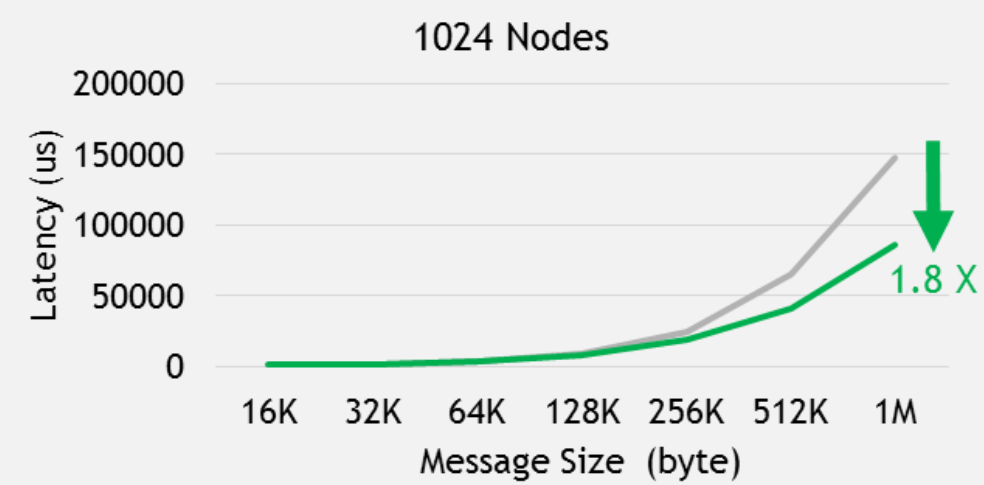
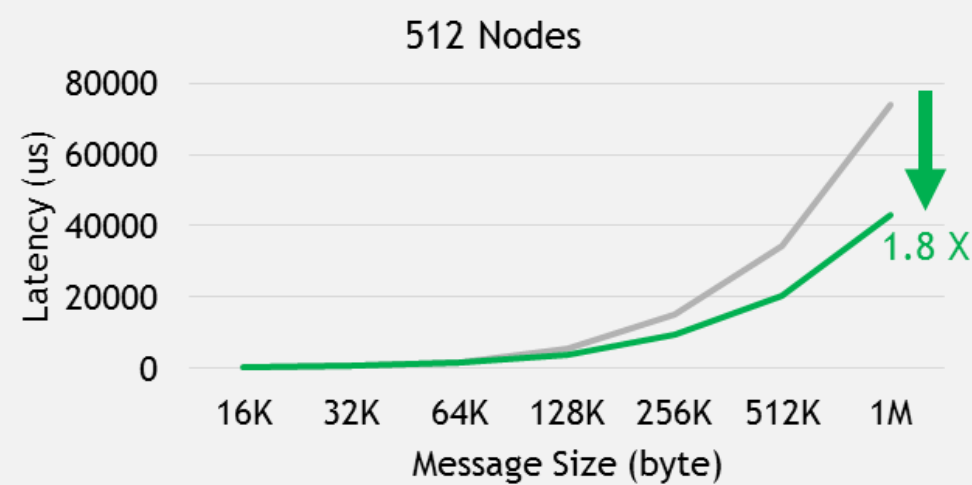
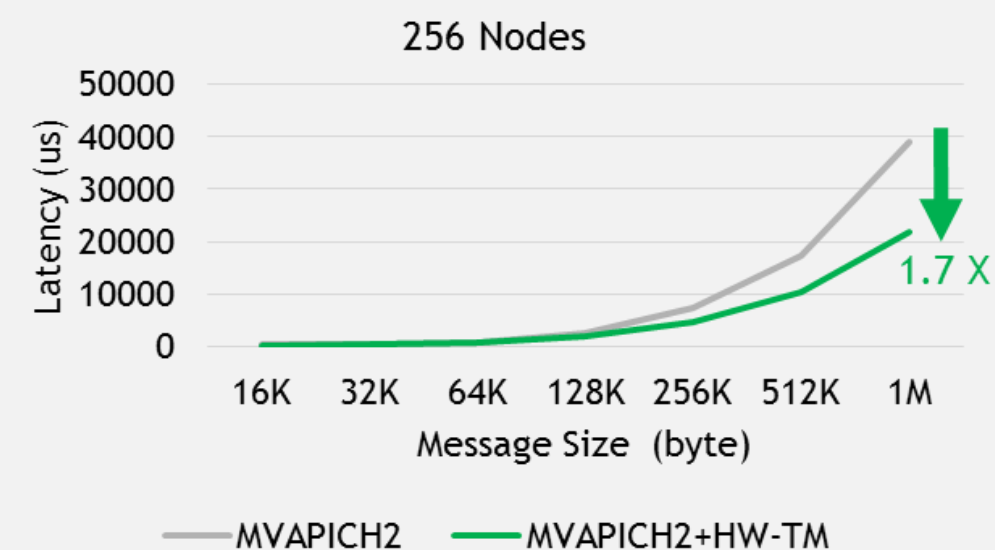
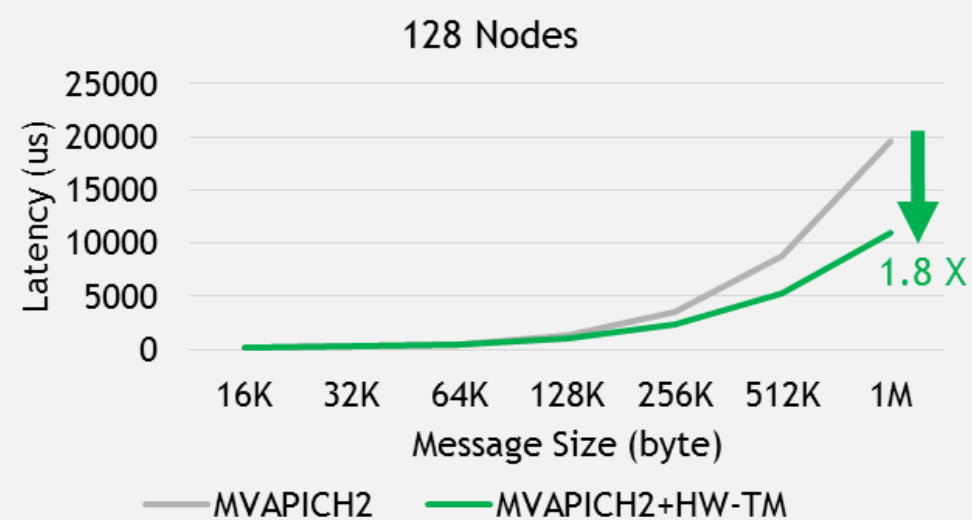
MPI TAG MATCHING HARDWARE ENGINE

INFINIBAND MPI TAG MATCHING HARDWARE ENGINE



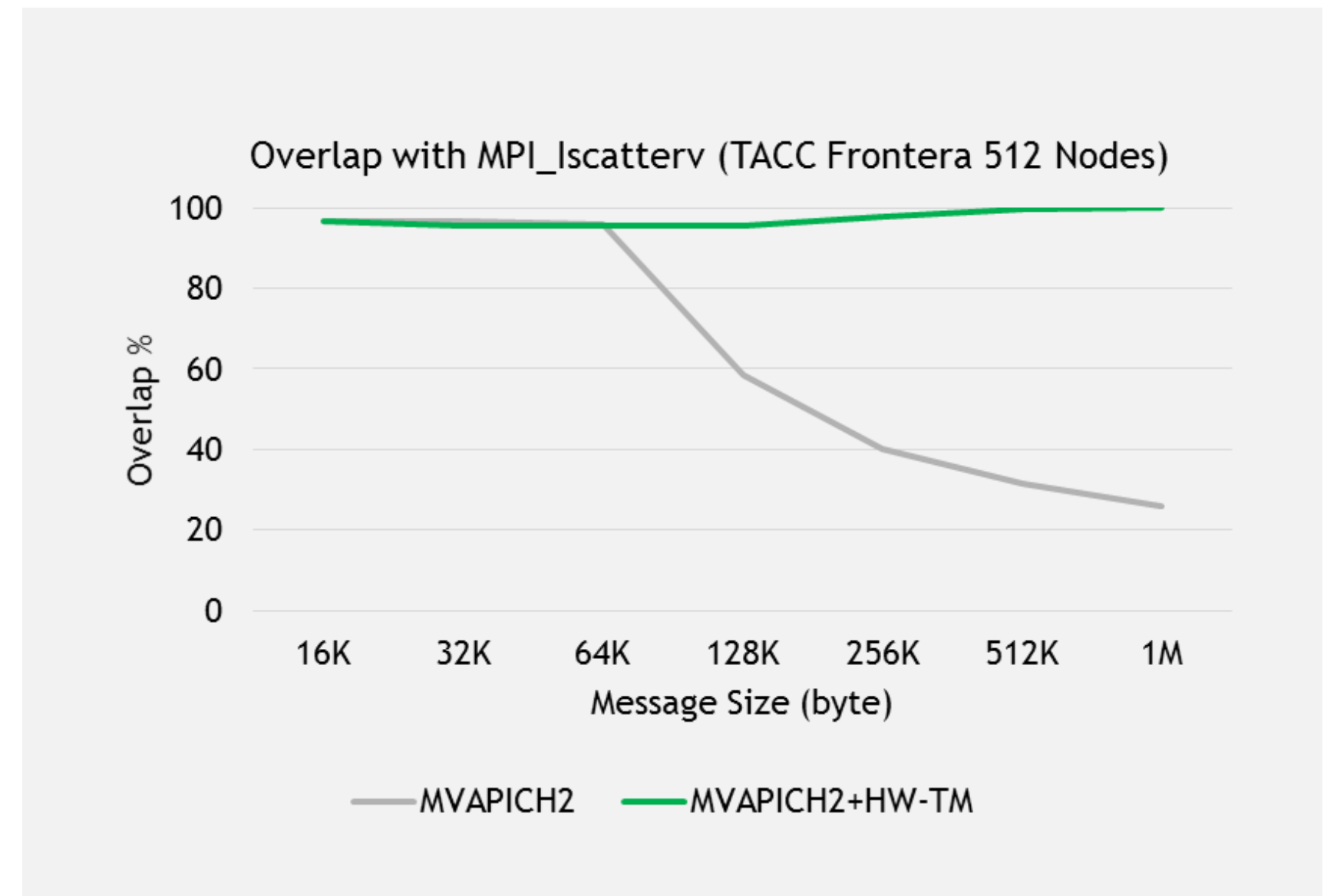
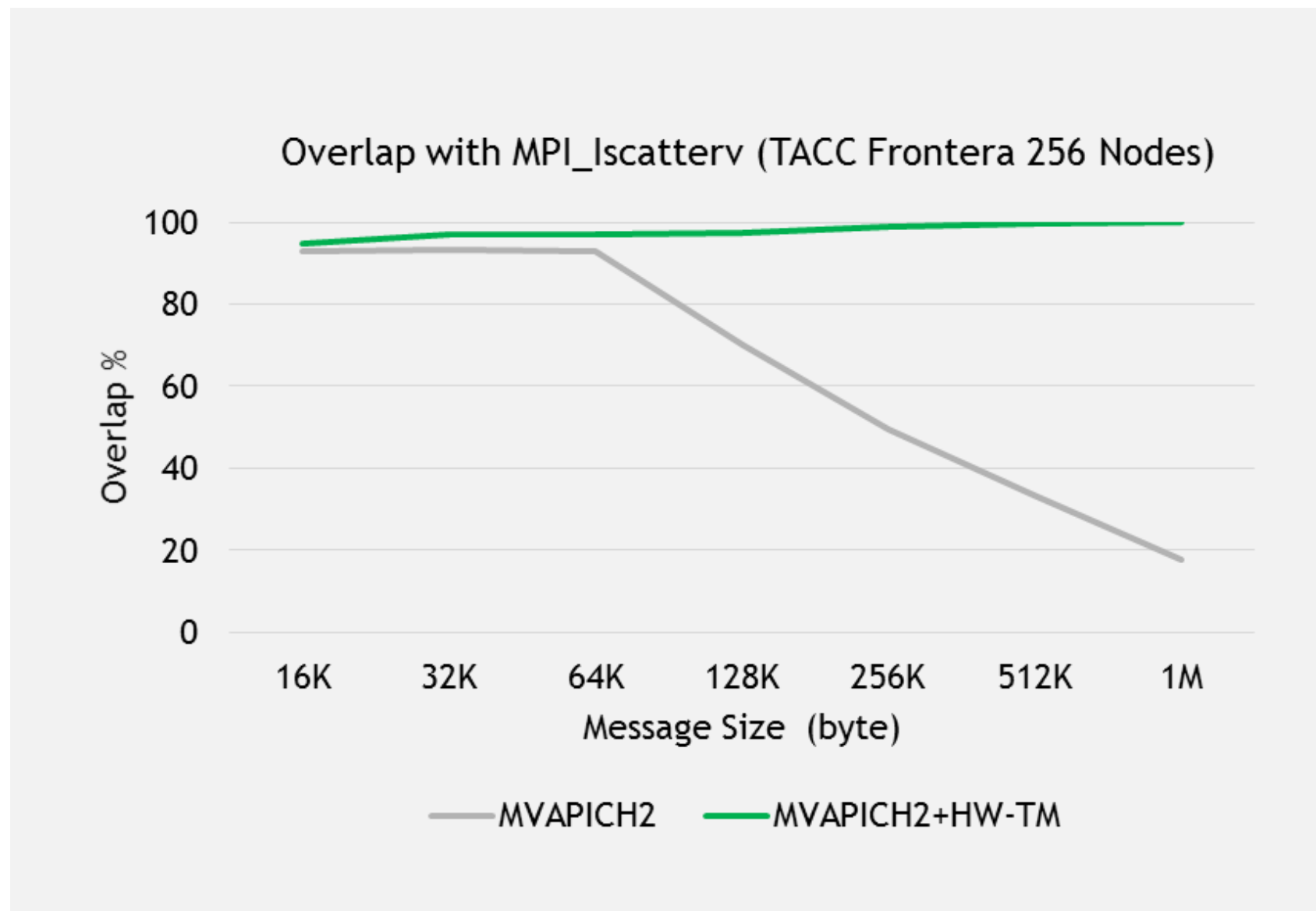
HARDWARE TAG MATCHING PERFORMANCE ADVANTAGES

1.8X Higher MPI_Isscatterv Performance on TACC Frontera



HARDWARE TAG MATCHING PERFORMANCE ADVANTAGES

Nearly 100% Compute - Communication Overlap





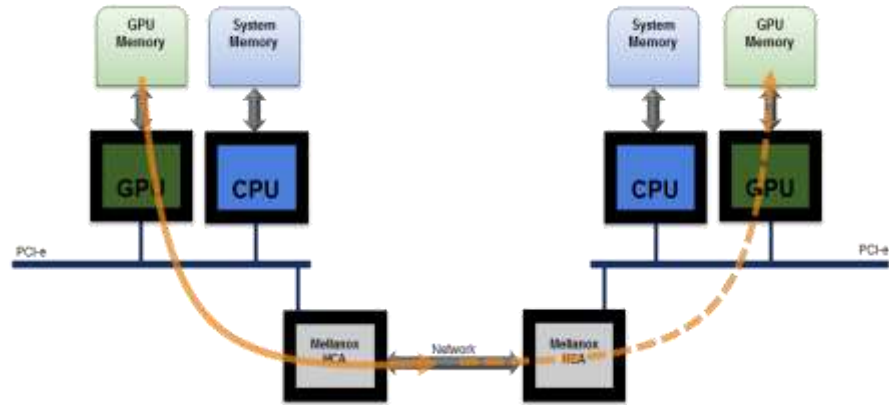
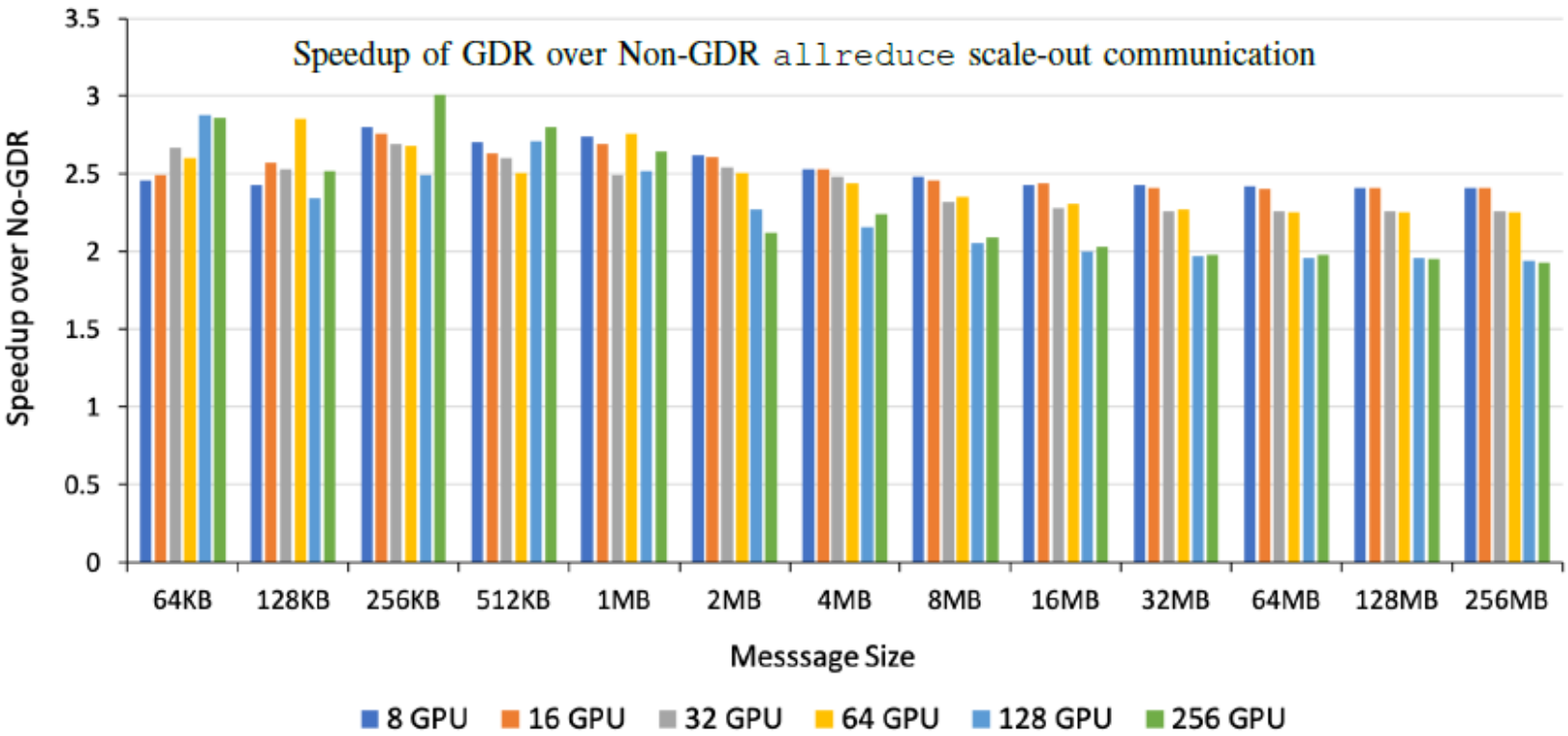
GPUDIRECT

EFFICIENT COMMUNICATION FOR ACCELERATED TRAINING

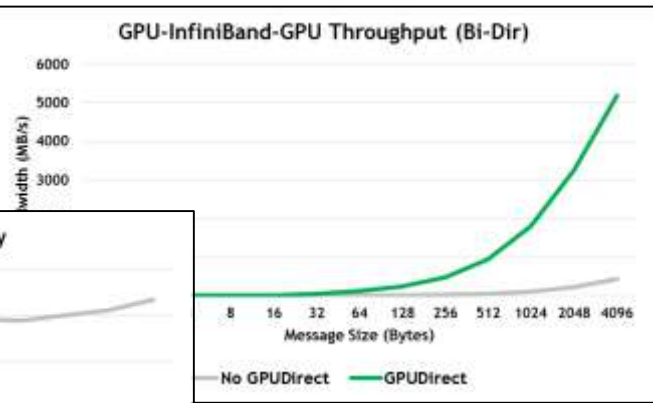
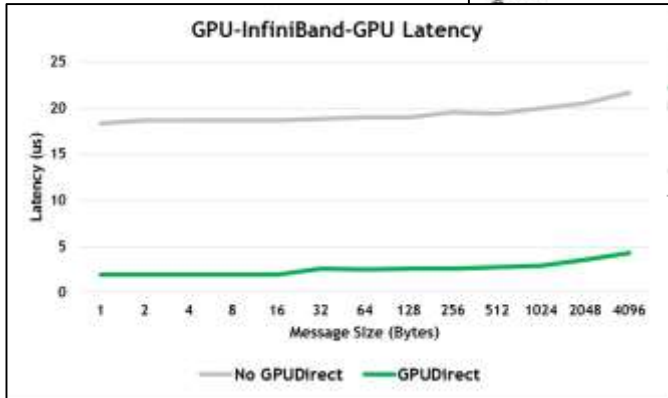
10X Better Latency & Bandwidth, 3X Faster Deep Learning

Deep Learning Training in Facebook Data Centers: Design of Scale-up and Scale-out Systems

Maxim Naumov*, John Kim†, Dheevatsa Mudigere‡, Srinivas Sridharan, Xiaodong Wang,
Whitney Zhao, Serhat Yilmaz, Changkyu Kim, Hector Yuen, Mustafa Ozdal, Krishnakumar Nair,
Isabel Gao, Bor-Yiing Su, Jiyan Yang and Mikhail Smelyanskiy
Facebook, 1 Hacker Way, Menlo Park, CA



OHIO STATE
MVAPICH
Courtesy of Dhabaleswar K. (DK) Panda
Ohio State University

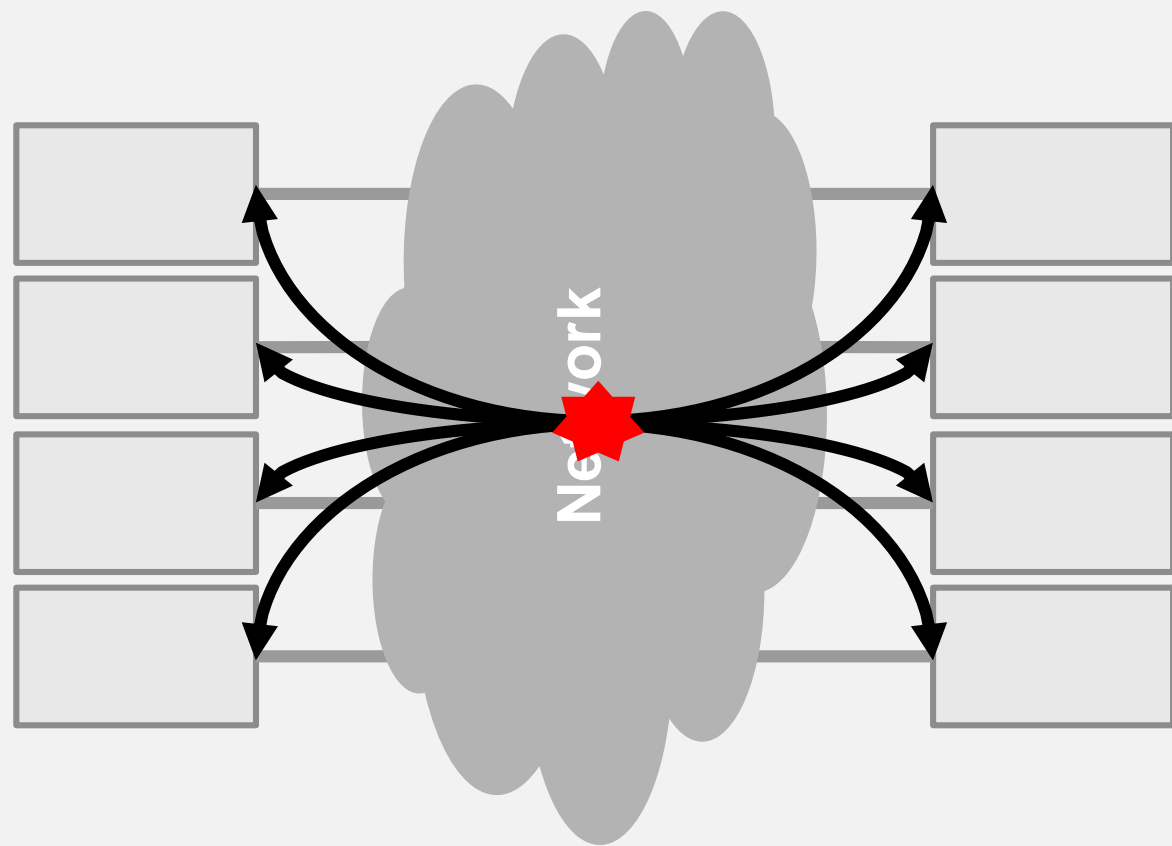




QUALITY OF SERVICE

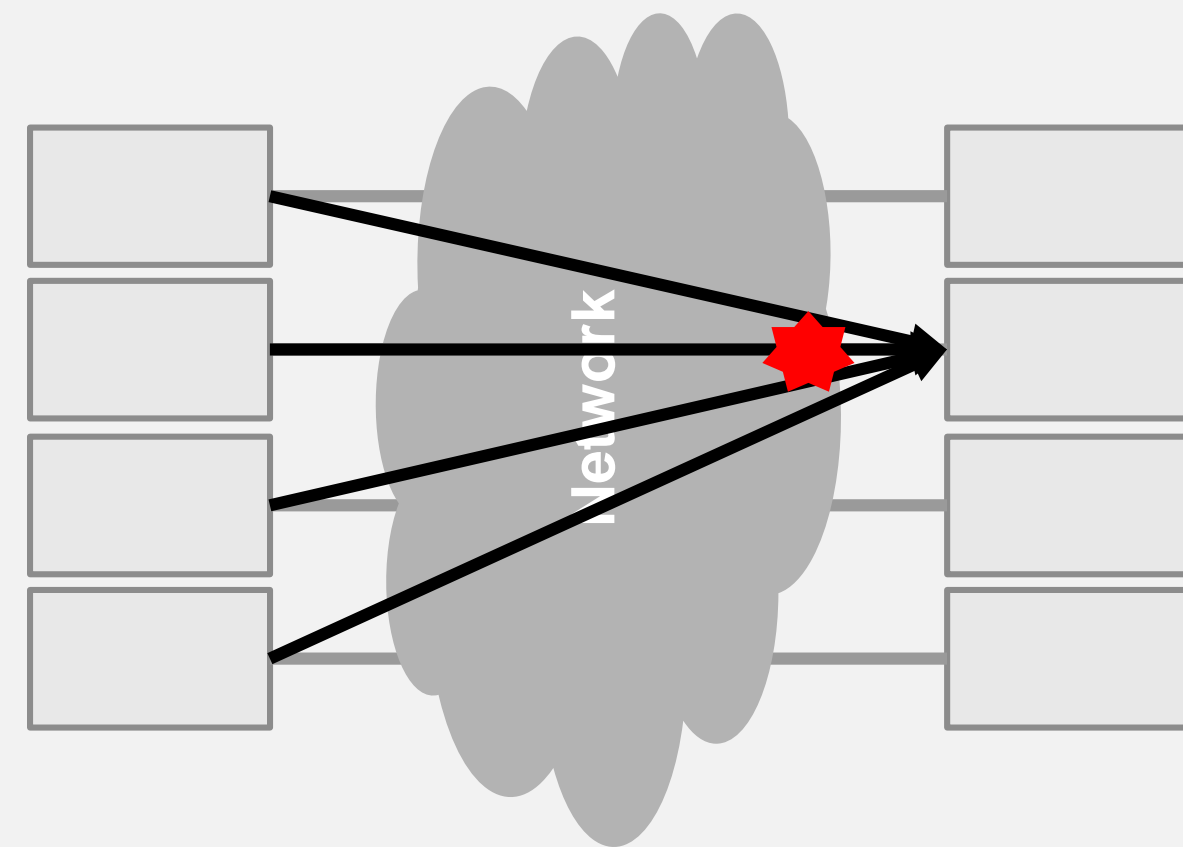
NETWORK CONGESTION TYPES

In-network Congestion



Solution: Adaptive Routing

In-cast Congestion



Solutions: Congestion Control

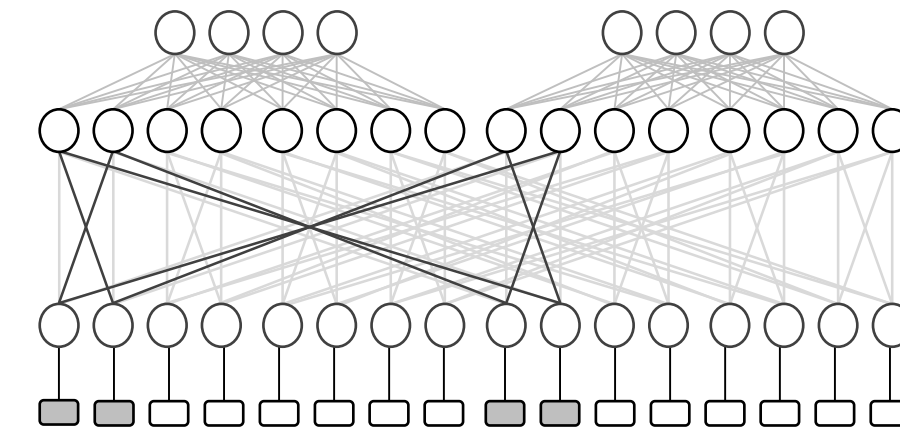
IN-NETWORK CONGESTION: ADAPTIVE ROUTING



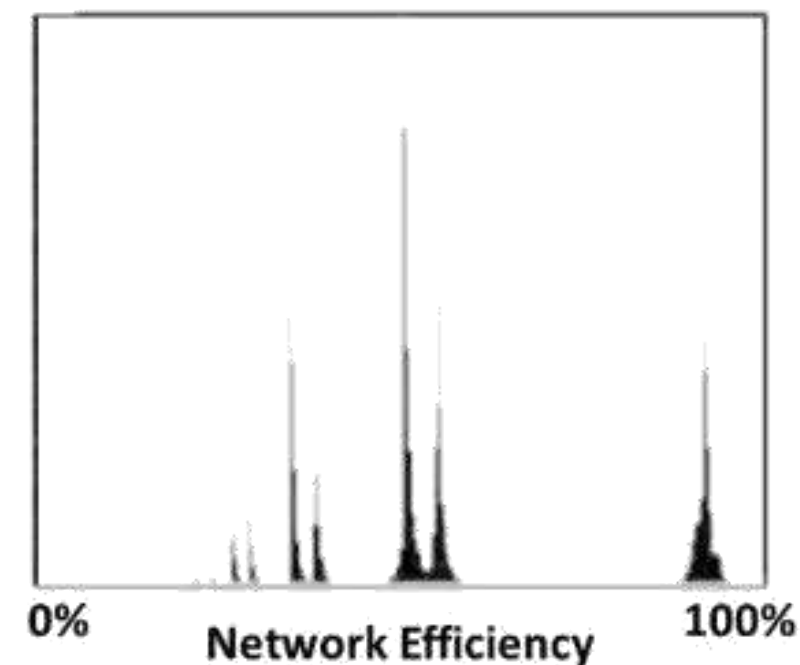
The Design, Deployment, and Evaluation of the CORAL Pre-Exascale Systems

Sudharshan S. Vazhkudai[†], Bronis R. de Supinski[‡], Arthur S. Bland[†], Al Geist[†], James Sexton*, Jim Kahle*, Christopher J. Zimmer[†], Scott Atchley[†], Sarp Oral[†], Don E. Maxwell[†], Veronica G. Vergara Larrea[†], Adam Bertsch[‡], Robin Goldstone[‡], Wayne Joubert[†], Chris Chamberau[†], David Appelhans*, Robert Blackmore*, Ben Casses[‡], George Chochia*, Gene Davison*, Matthew A. Ezell[†], Tom Gooding*, Elsa Gonsiorowski[‡], Leopold Grinberg*, Bill Hanson*, Bill Hartner*, Ian Karlin[†], Matthew L. Leininger[‡], Dustin Leverman[†], Chris Marroquin*, Adam Moody[‡], Martin Ohmacht*, Ramesh Pankajakshan[‡], Fernando Pizzano*, James H. Rogers[†], Bryan Rosenberg*, Drew Schmidt[†], Mallikarjun Shankar[†], Feiyi Wang[†], Py Watson[‡], Bob Walkup*, Lance D. Weems[‡], Junqi Yin[†]

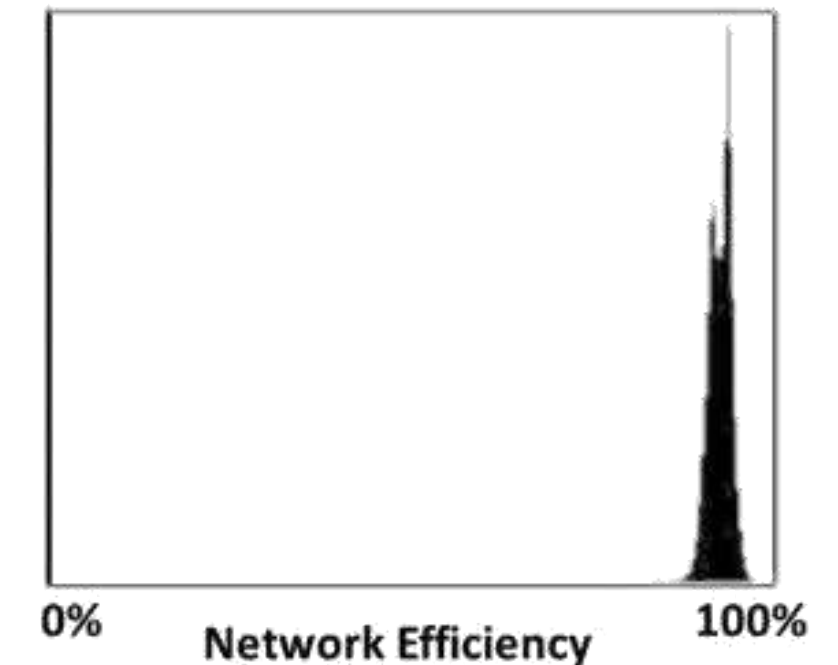
[†] Oak Ridge National Laboratory, [‡] Lawrence Livermore National Laboratory, * IBM
{vazhkudaiss@ornl.gov, bronis@llnl.gov}



mpiGraph: Static vs. Adaptive Routing

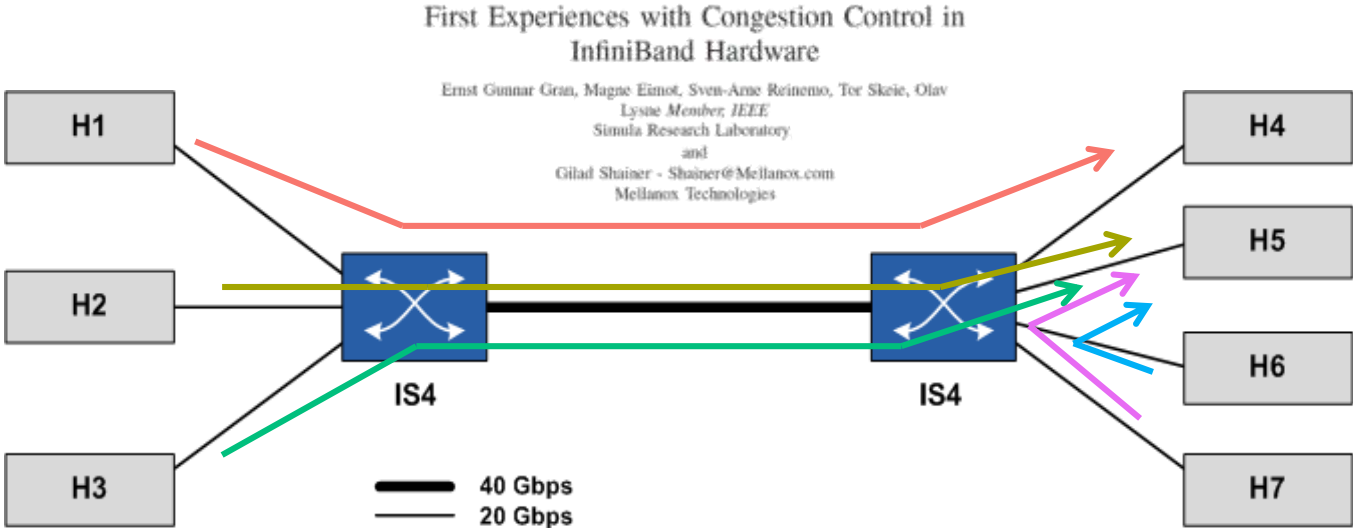


Static Routing

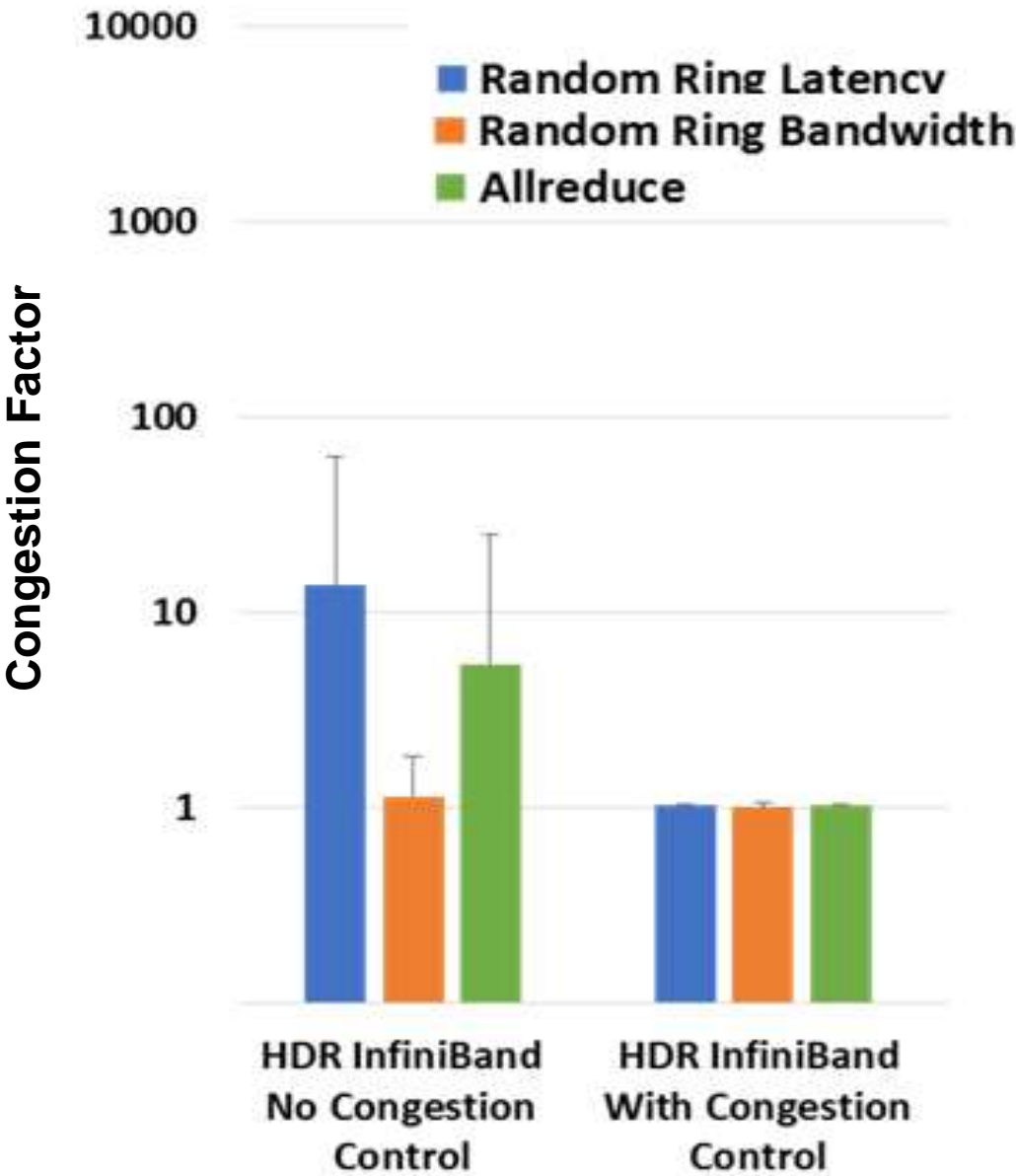


Adaptive Routing

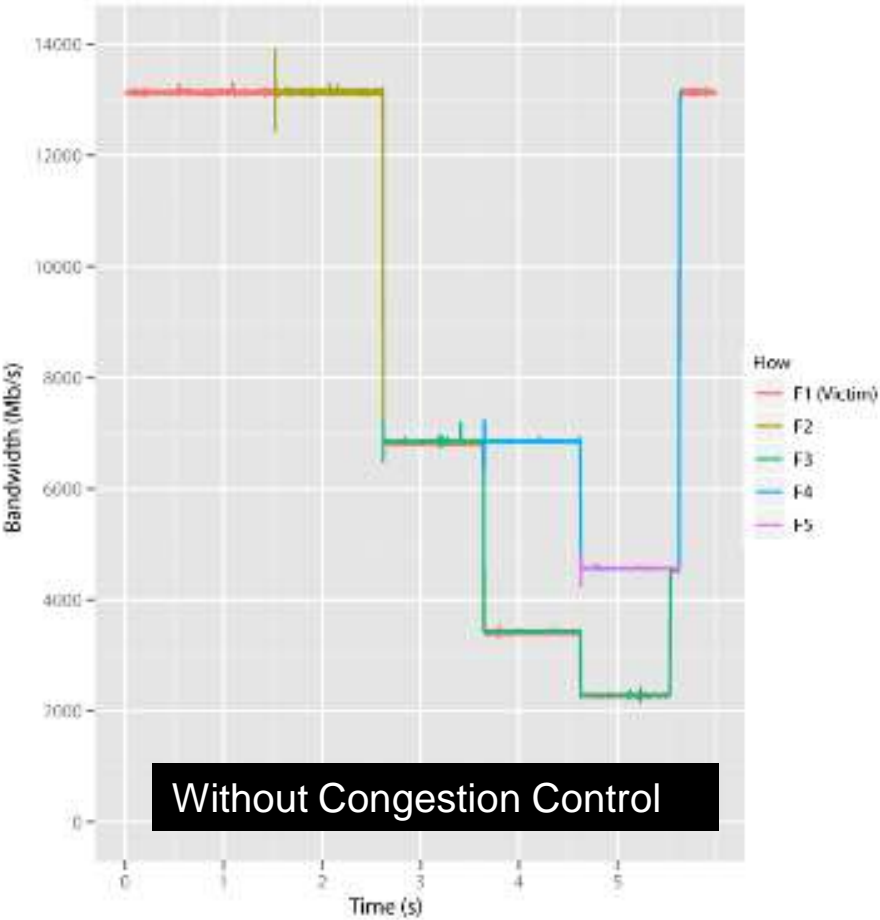
INFINIBAND CONGESTION CONTROL



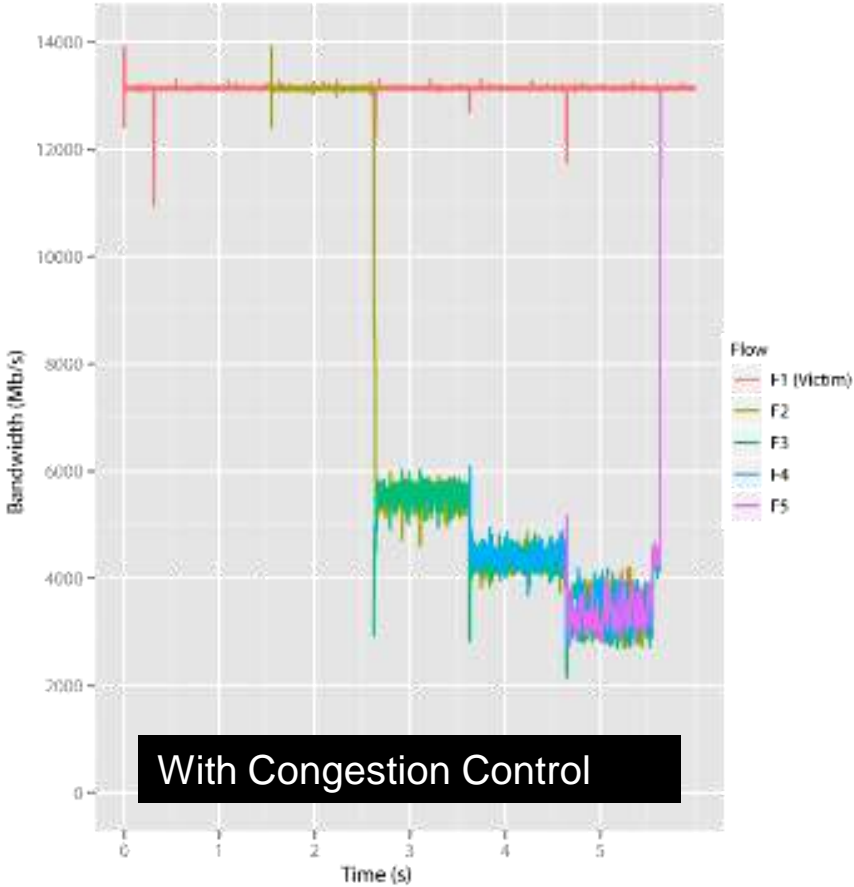
GPCNet Benchmark



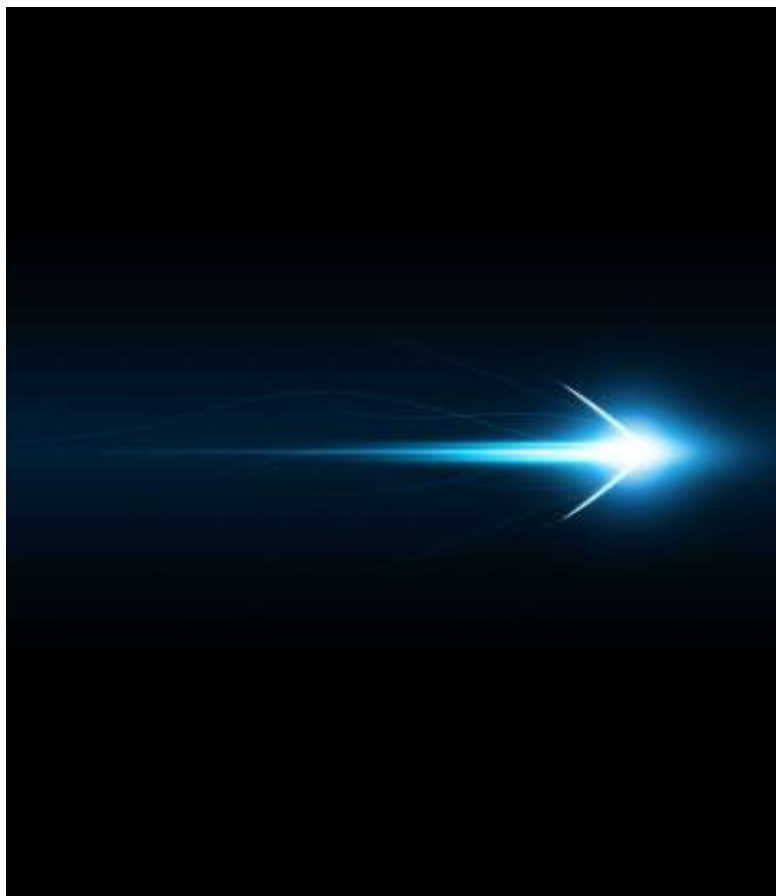
Congestion – Throughput loss



No congestion – highest throughput!

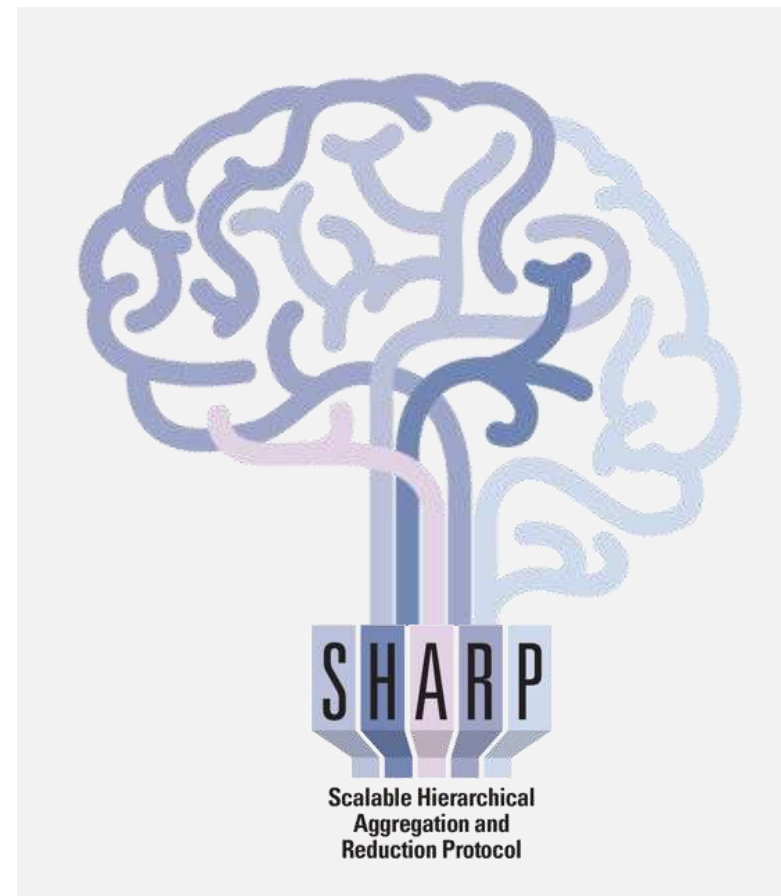


INFINIBAND ACCELERATED SUPERCOMPUTING



Speed of Light

200Gb/s Data Throughput
RDMA and GPUDirect RDMA
3X Better (Lower) Latency



SHARP AI Technology

AI Acceleration Engines
2.5X Higher AI Performance



SHIELD AI Technology

Self Healing Network
1000X Faster Recovery Time



UFM Cyber AI

Data Center Cyber Intelligence
and Analytics

