



In-Network Computing Technology and Performance Advantages

September 2019



SUPERCONNECTING the #1 Supercomputers



OAK RIDGE
National Laboratory

Lawrence Livermore
National Laboratory

国家超级计算无锡中心
National Supercomputing Center in Wuxi

TACC
TEXAS ADVANCED COMPUTING CENTER

AIST
NATIONAL INSTITUTE OF
ADVANCED INDUSTRIAL SCIENCE
AND TECHNOLOGY (AIST)

Lawrence Livermore
National Laboratory



1 TOP 500
The List.

2 TOP 500
The List.

3 TOP 500
The List.

5 TOP 500
The List.

8 TOP 500
The List.

10 TOP 500
The List.

InfiniBand Accelerates 6 of Top 10 Supercomputers

SUPERCONNECTING the #1 Supercomputers



TACC
TEXAS ADVANCED COMPUTING CENTER



FRONTIER

5 **TOP 500**
The List.

MISSISSIPPI STATE UNIVERSITY



62 **TOP 500**
The List.



166 **TOP 500**
The List.

筑波大学
University of Tsukuba



264 **TOP 500**
The List.



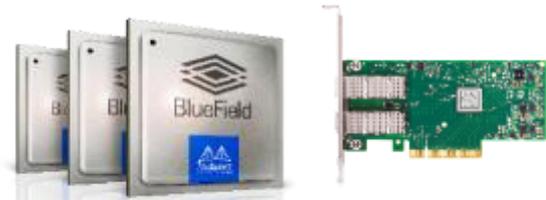
World's First
HDR InfiniBand
Supercomputer



India's National
Supercomputing
Program

HDR 200G InfiniBand Accelerated Supercomputers

HPC and AI Needs the Most Intelligent Interconnect



SmartNIC



System on a Chip

Higher

Data Speeds

Faster

Data Processing

Better

Data Security



Adapters



Switches



Cables & Transceivers



The Need for Intelligent and Faster Interconnect

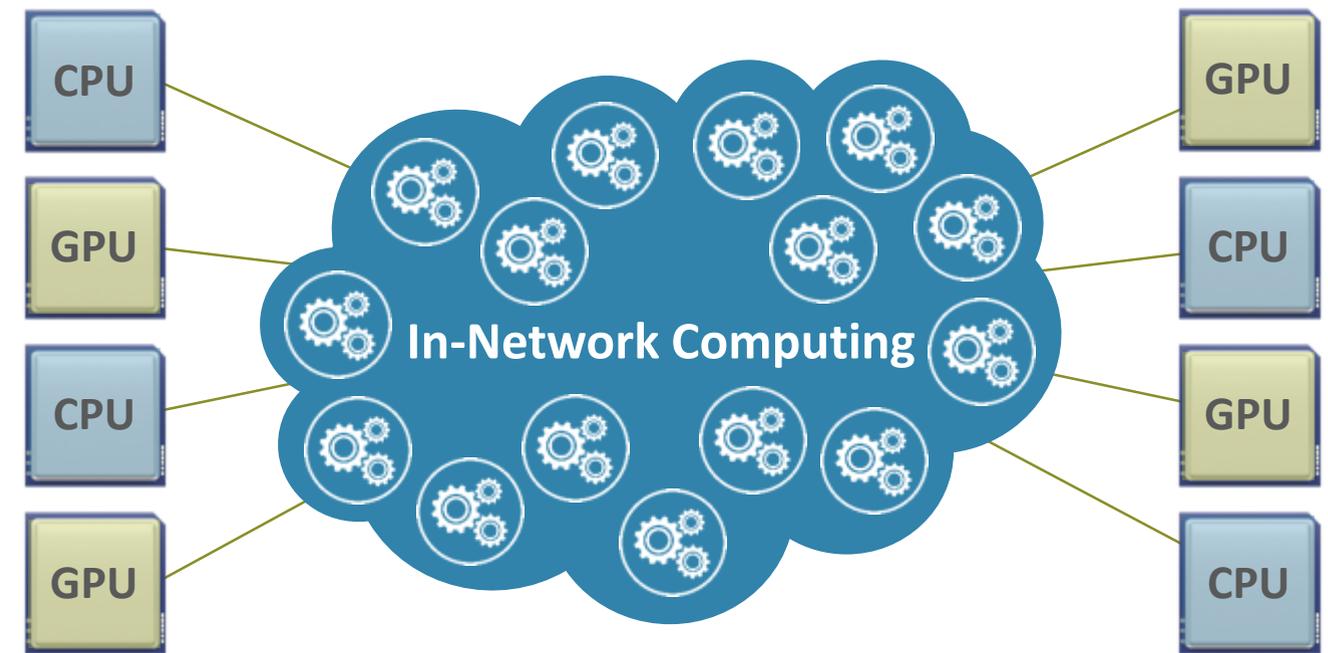
Faster Data Speeds and In-Network Computing
Enable Higher Performance and Scale

CPU-Centric (Onload)



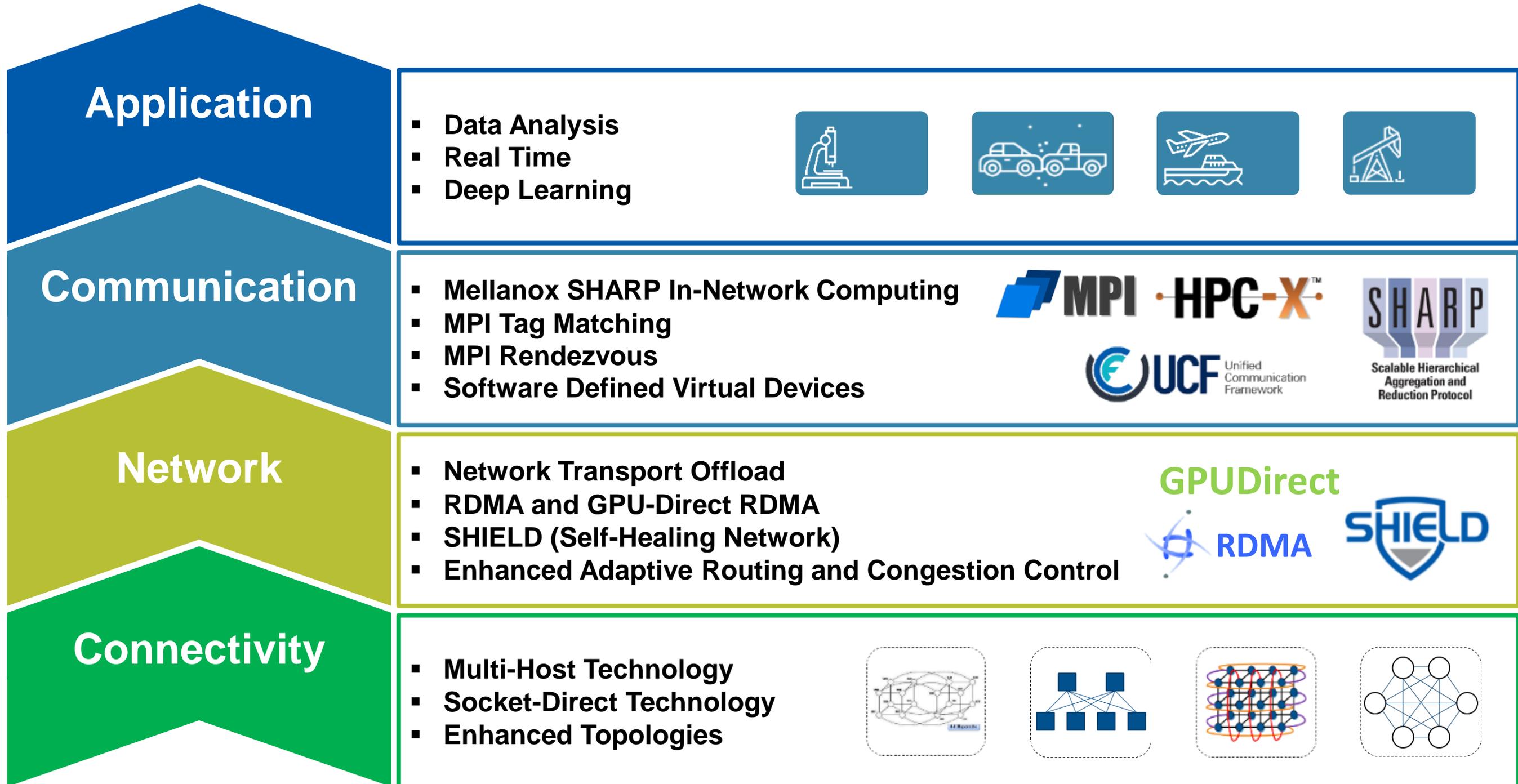
Must Wait for the Data
Creates Performance Bottlenecks

Data-Centric (Offload)



Analyze Data as it Moves!
Higher Performance and Scale

Accelerating All Levels of HPC / AI Frameworks



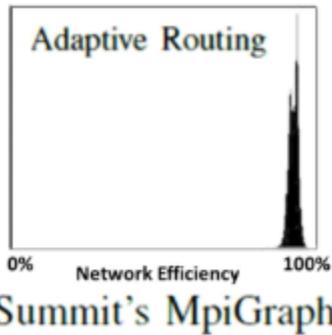
Highest Performance and Scalability for Exascale Platforms

OAK RIDGE
National Laboratory

SHARP SHIELD
SELF-HEALING



96%
Network
Utilization

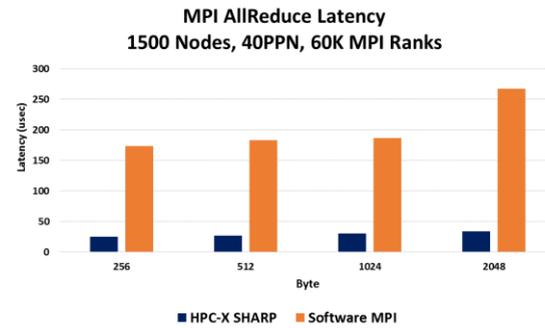


UNIVERSITY OF
TORONTO
SciNet

SHARP

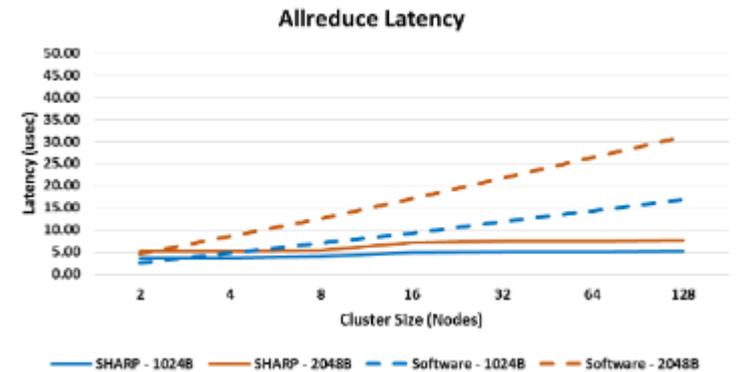


7X
Higher
Performance

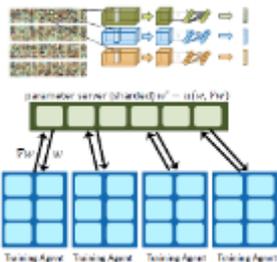


SHARP
Scalable Hierarchical
Aggregation and
Reduction Protocol

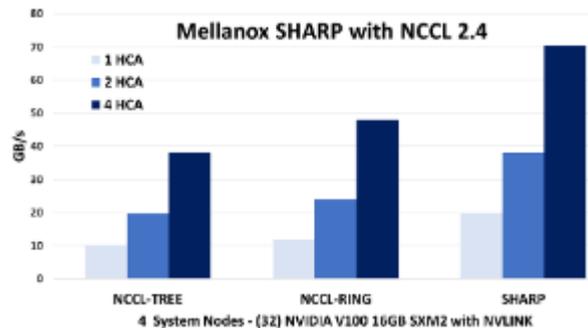
**Flat
Latency**



**Deep
Learning**

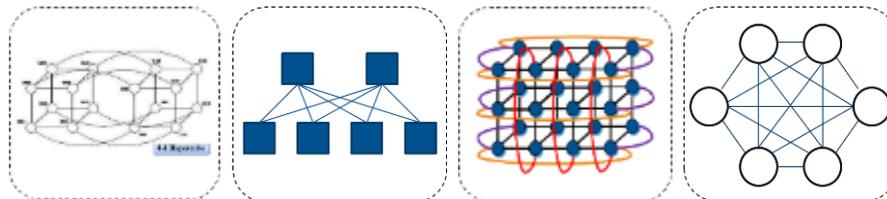


2X
Higher
Performance



SHIELD
SELF-HEALING

5000X
Higher
Resiliency



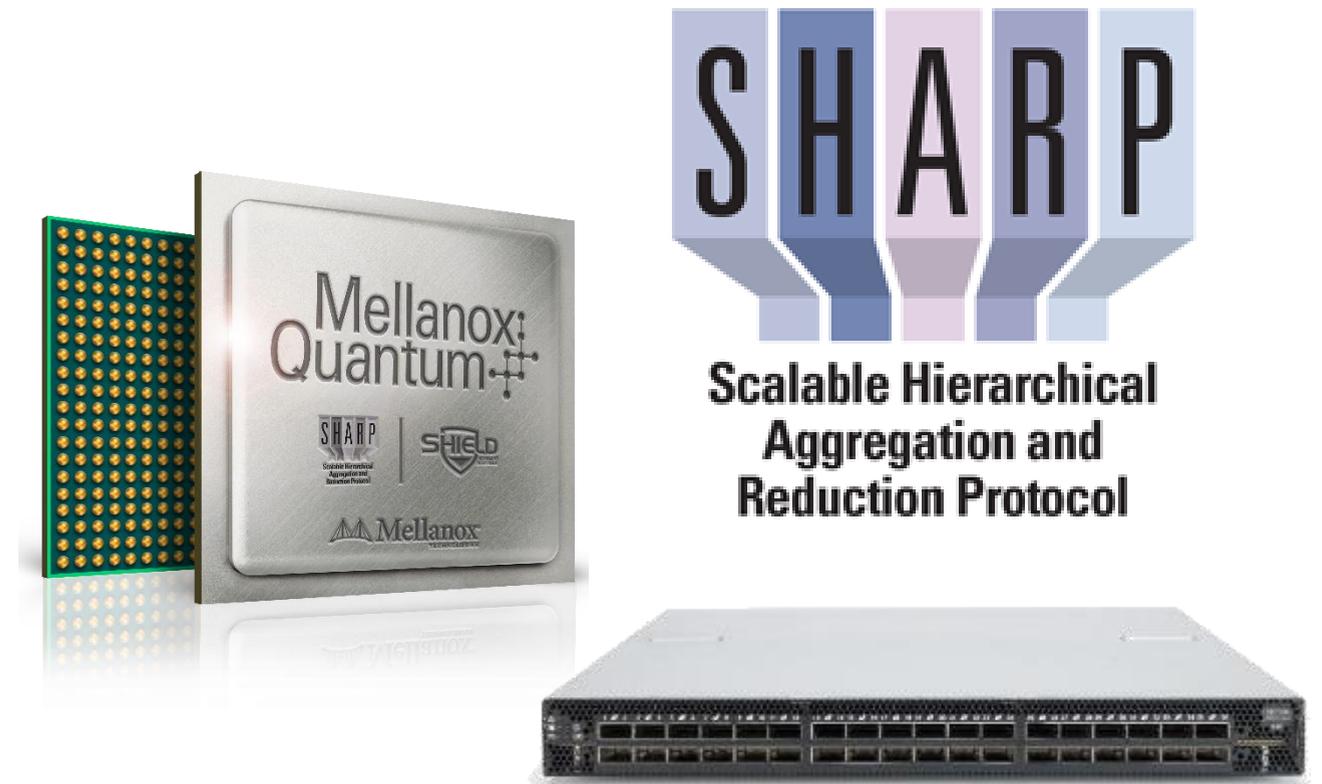
XDR 1000G

NDR 400G

HDR 200G



Scalable Hierarchical Aggregation and Reduction Protocol (SHARP)

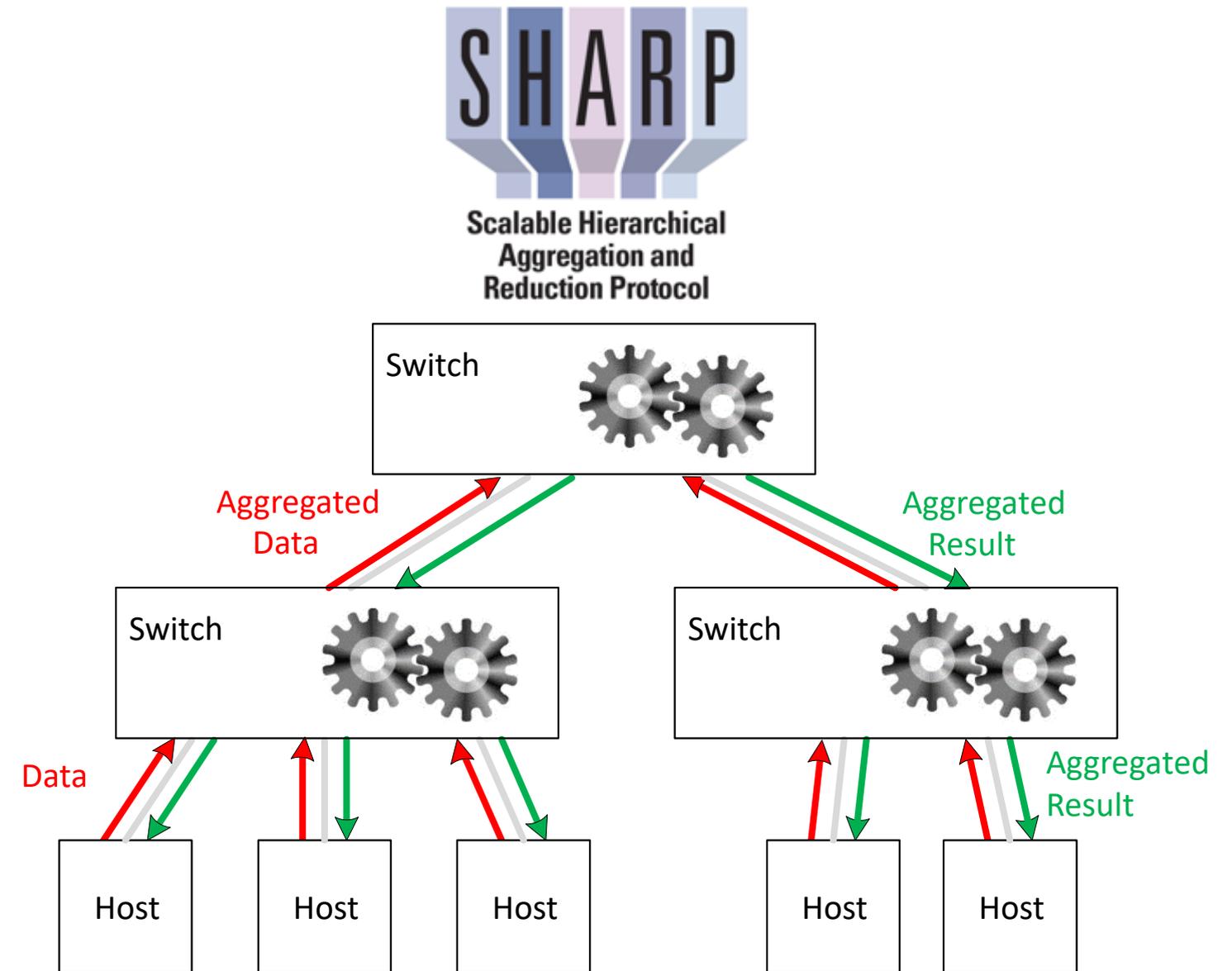


Scalable Hierarchical Aggregation and Reduction Protocol (SHARP)

- Reliable Scalable General Purpose Primitive
 - In-network Tree based aggregation mechanism
 - Large number of groups
 - Multiple simultaneous outstanding operations

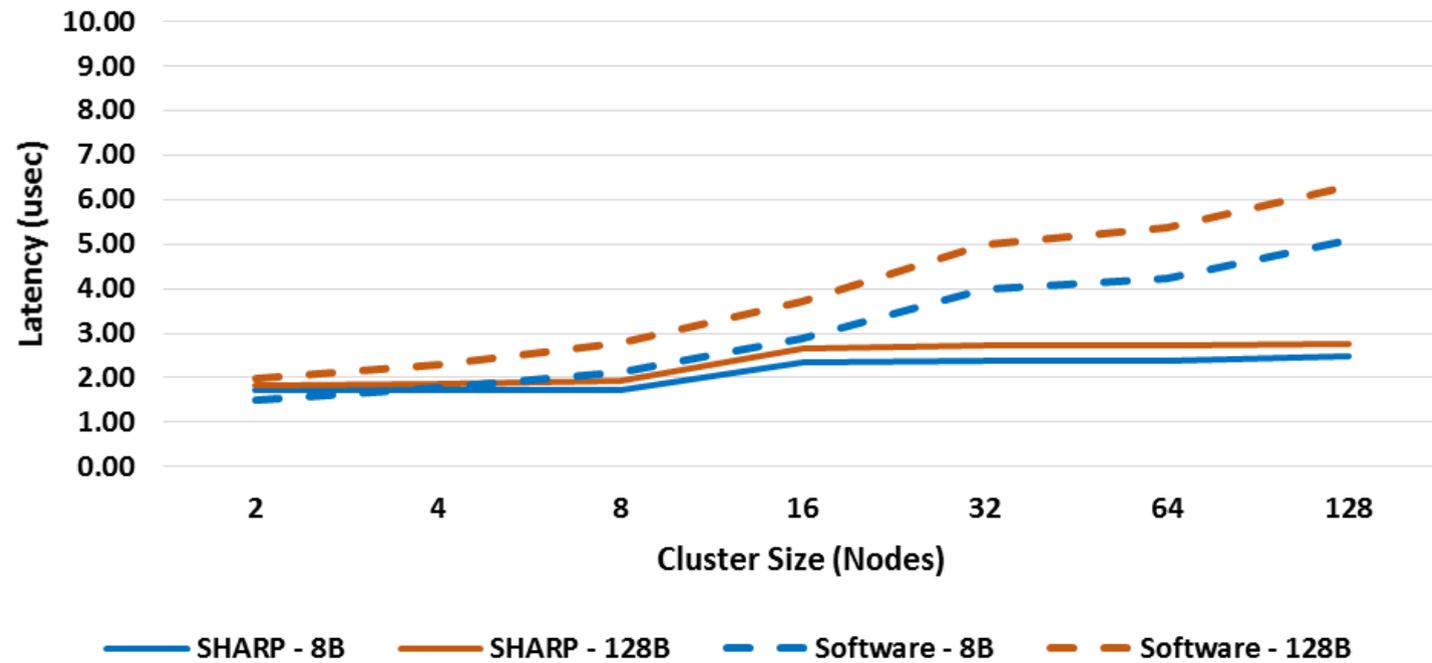
- Applicable to Multiple Use-cases
 - HPC Applications using MPI / SHMEM
 - Distributed Machine Learning applications

- 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

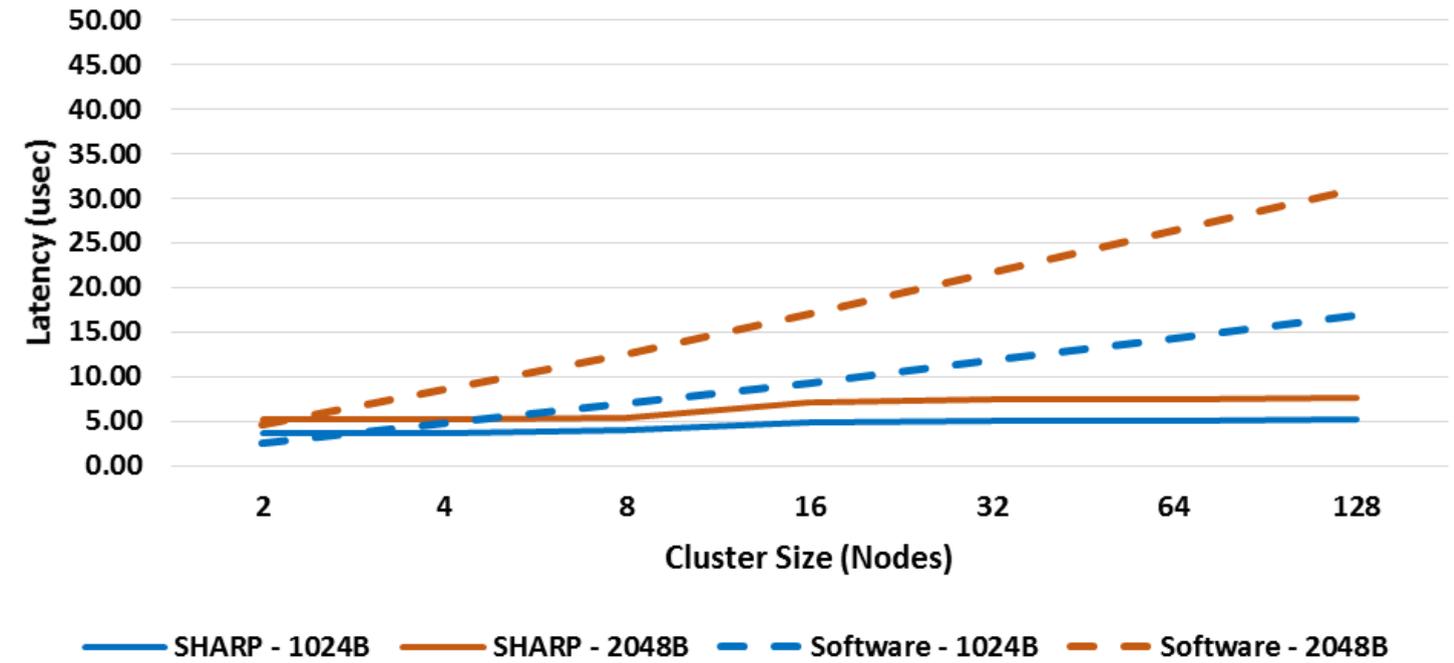


SHARP AllReduce Performance Advantages (128 Nodes)

Allreduce Latency



Allreduce Latency

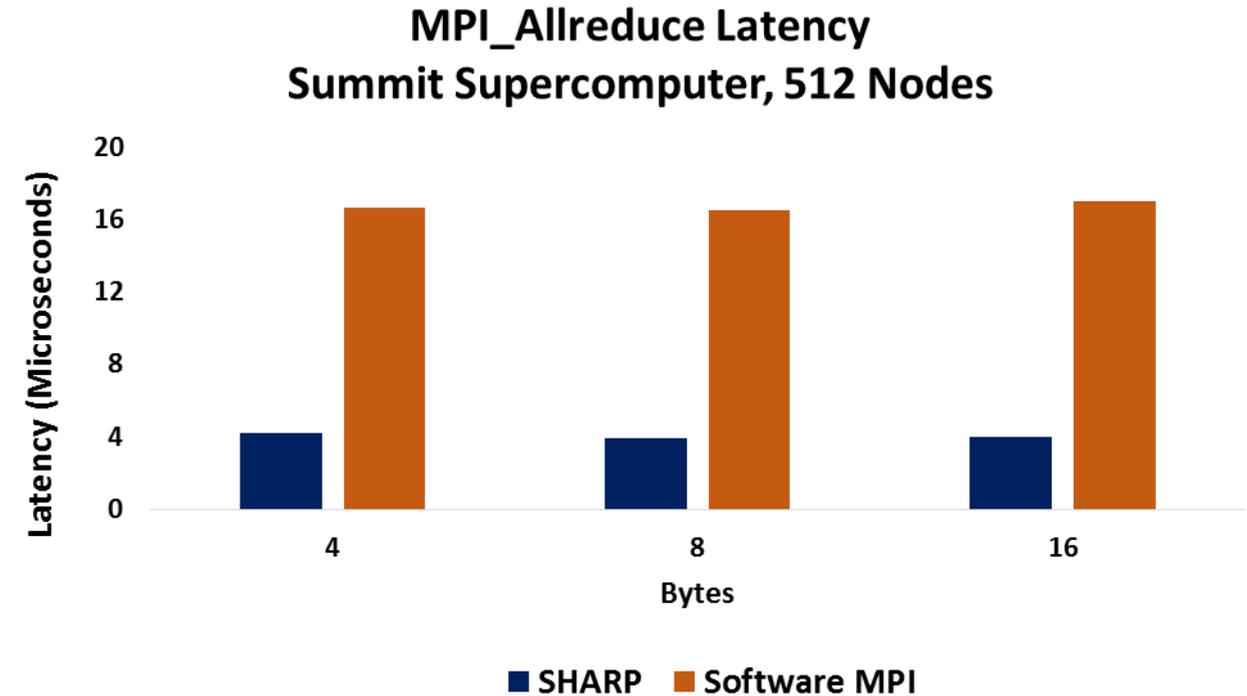
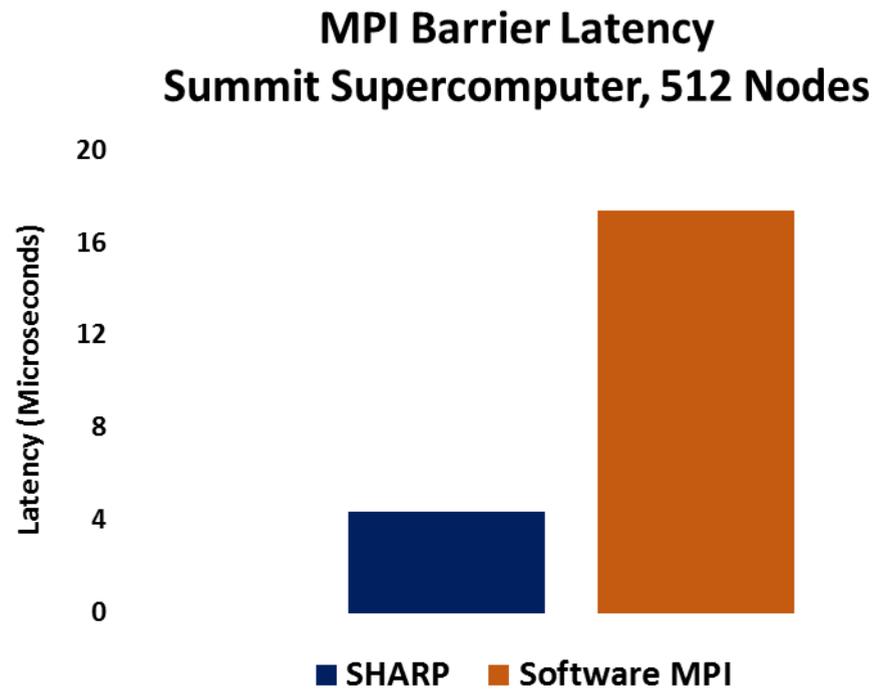


Scalable Hierarchical
Aggregation and
Reduction Protocol

SHARP enables 75% Reduction in Latency
Providing Scalable Flat Latency

SHARP AllReduce Performance Advantages

Oak Ridge National Laboratory – Coral Summit Supercomputer

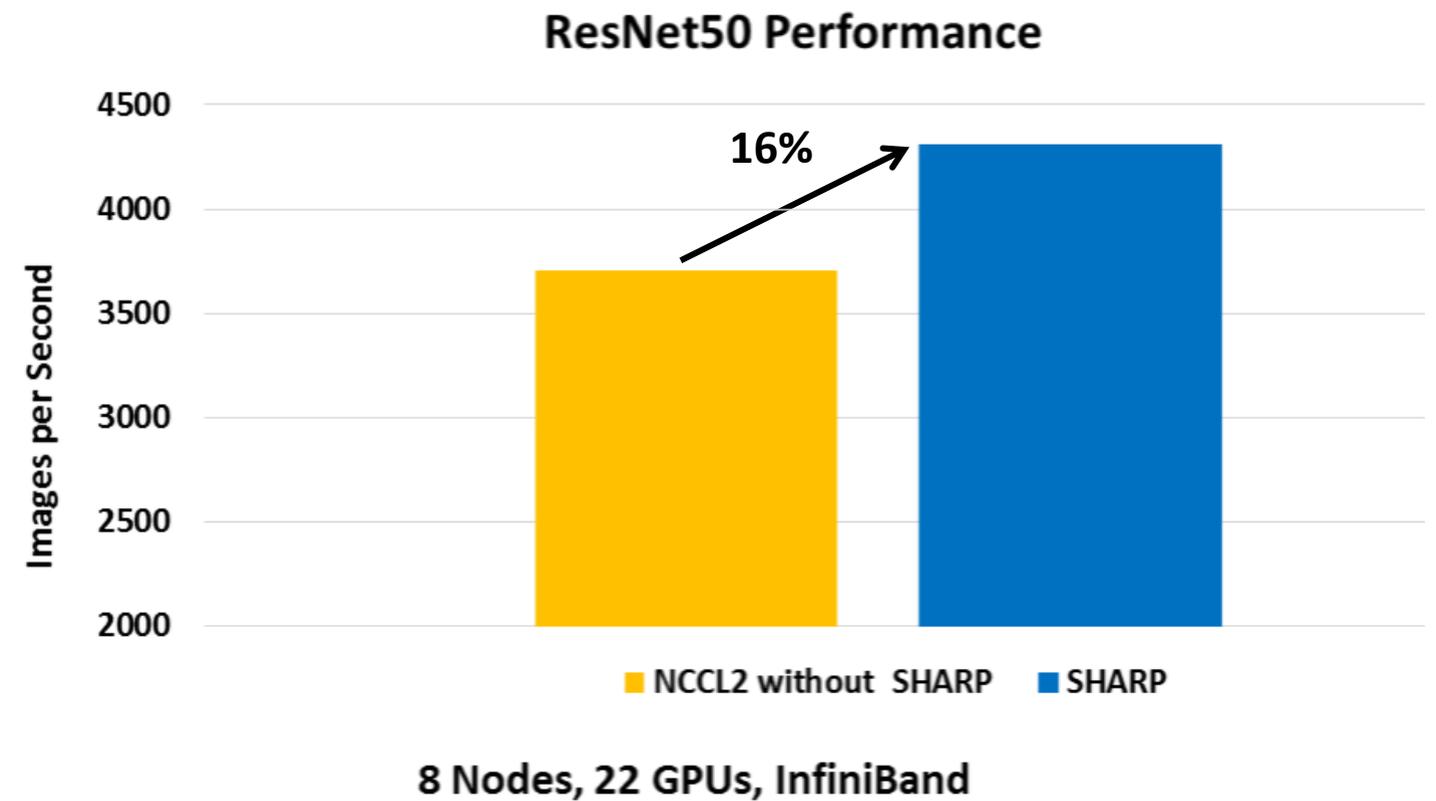
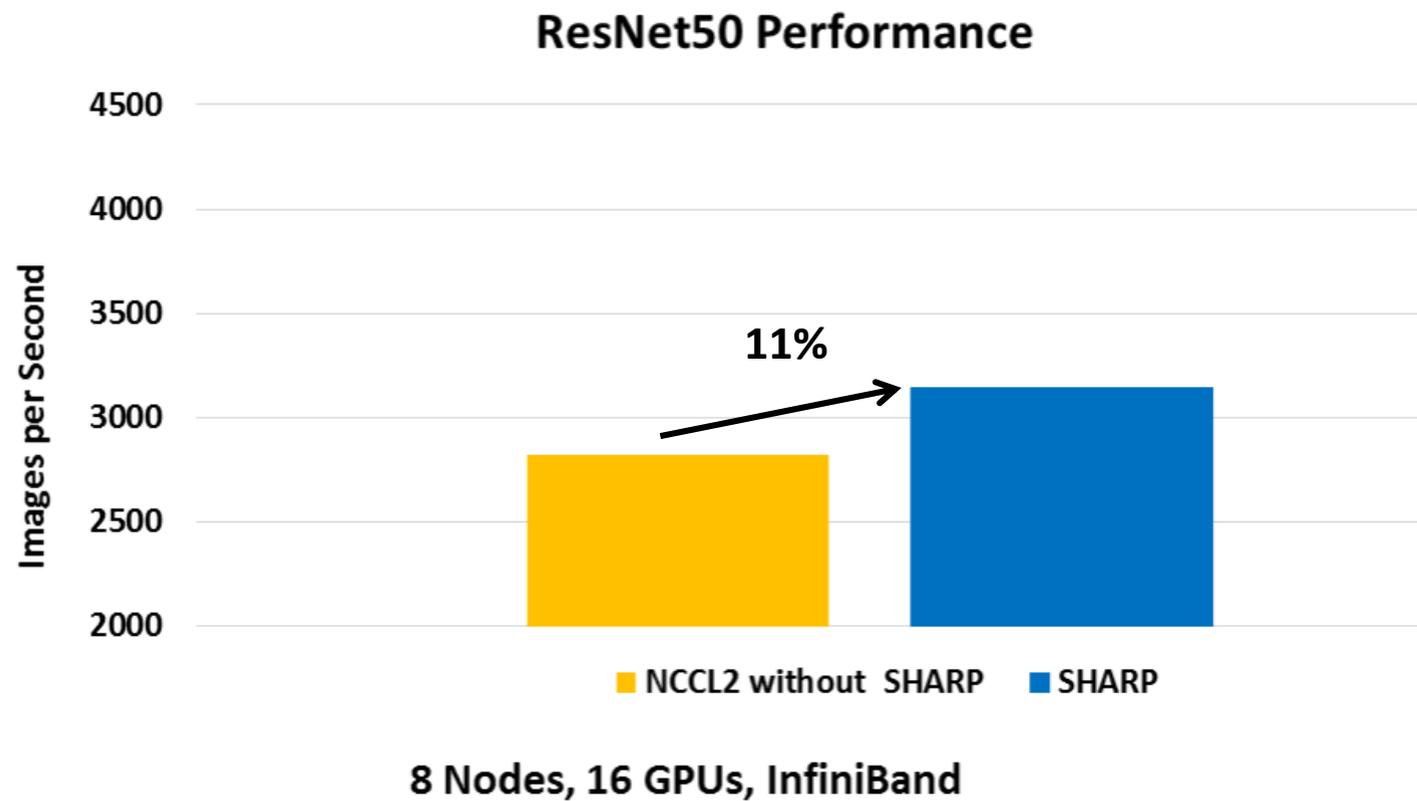


Scalable Hierarchical
Aggregation and
Reduction Protocol

SHARP Enables Highest Performance

SHARP Performance Advantage for AI

- SHARP provides 16% Performance Increase for deep learning, initial results
- TensorFlow with Horovod running ResNet50 benchmark, HDR InfiniBand (ConnectX-6, Quantum)



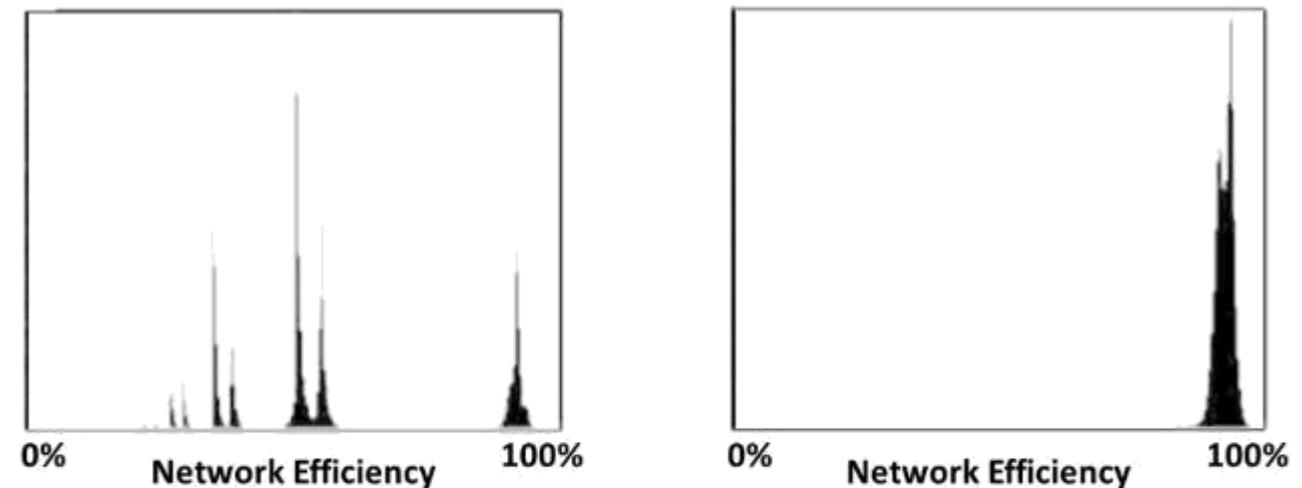
Adaptive Routing



InfiniBand Proven Adaptive Routing Performance

- Oak Ridge National Laboratory – Coral Summit supercomputer
- Bisection bandwidth benchmark, based on mpiGraph
 - Explores the bandwidth between possible MPI process pairs
- AR results demonstrate an average performance of 96% of the maximum bandwidth measured

mpiGraph explores the bandwidth between possible MPI process pairs. In the histograms, the single cluster with AR indicates that all pairs achieve nearly maximum bandwidth while single-path static routing has nine clusters as congestion limits bandwidth, negatively impacting overall application performance.



Without Adaptive Routing

With Adaptive Routing

Summit's MpiGraph Output

*“The Design, Deployment, and Evaluation of the CORAL Pre-Exascale Systems”,
Sudharshan S. Vazhkudai, Arthur S. Bland, Al Geist, Christopher J. Zimmer, Scott Atchley, Sarp Oral, Don E. Maxwell, Veronica G. Vergara Larrea, Wayne Joubert, Matthew A. Ezell, Dustin Leverman, James H. Rogers, Drew Schmidt, Mallikarjun Shankar, Feiyi Wang, Junqi Yin (Oak Ridge National Laboratory) and Bronis R. de Supinski, Adam Bertsch, Robin Goldstone, Chris Chambreau, Ben Casses, Elsa Gonsiorowski, Ian Karlin, Matthew L. Leininger, Adam Moody, Martin Ohmacht, Ramesh Pankajakshan, Fernando Pizzano, Py Watson, Lance D. Weems (Lawrence Livermore National Laboratory) and James Sexton, Jim Kahle, David Appelhans, Robert Blackmore, George Chochia, Gene Davison, Tom Gooding, Leopold Grinberg, Bill Hanson, Bill Hartner, Chris Marroquin, Bryan Rosenberg, Bob Walkup (IBM)*

HDR InfiniBand



Highest-Performance 200Gb/s InfiniBand Solutions

Adapters		<p>200Gb/s Adapter, 0.6us latency 215 million messages per second (10 / 25 / 40 / 50 / 56 / 100 / 200Gb/s)</p>	
Switch		<p>40 HDR (200Gb/s) InfiniBand Ports 80 HDR100 InfiniBand Ports Throughput of 16Tb/s, <90ns Latency</p>	
SoC		<p>System on Chip and SmartNIC Programmable adapter Smart Offloads</p>	
Interconnect		<p>Transceivers Active Optical and Copper Cables (10 / 25 / 40 / 50 / 56 / 100 / 200Gb/s)</p>	
Software		<p>MPI, SHMEM/PGAS, UPC For Commercial and Open Source Applications Leverages Hardware Accelerations</p>	

ConnectX-6 HDR InfiniBand Adapter

Leading Connectivity

- 200Gb/s InfiniBand and Ethernet
 - HDR, HDR100, EDR (100Gb/s) and lower speeds
 - 200GbE, 100GbE and lower speeds
- Single and dual ports
- 50Gb/s PAM4 SerDes

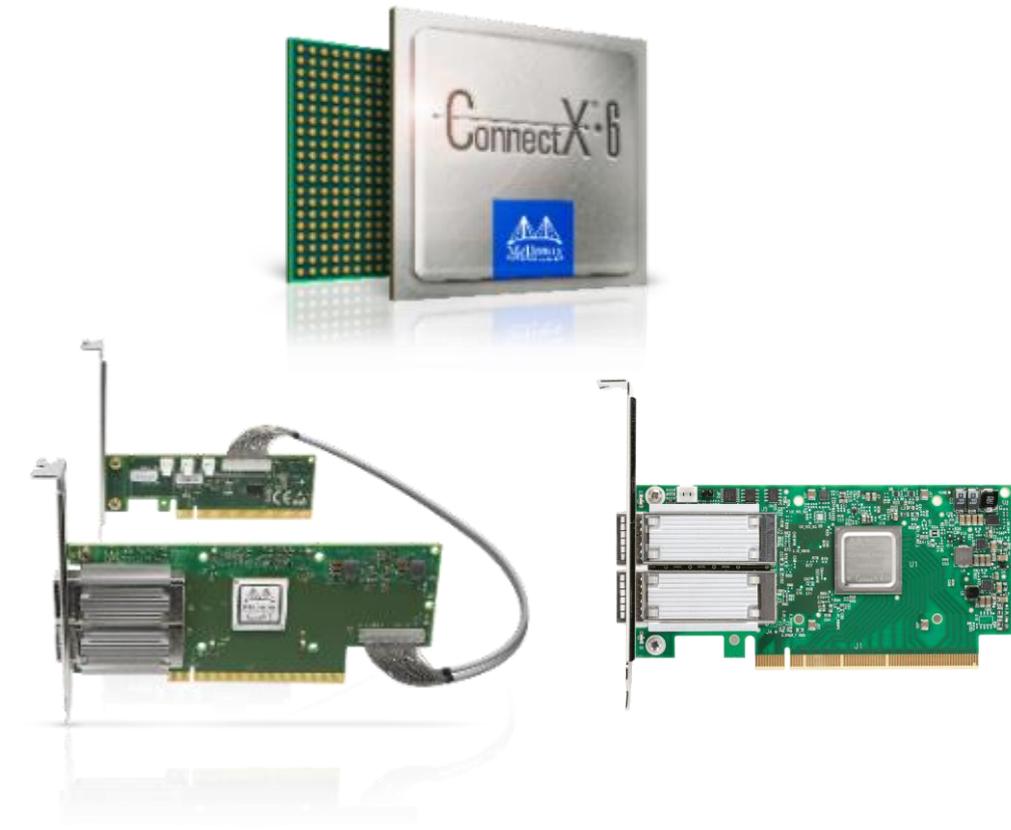
Leading Performance

- 200Gb/s throughput, 0.6usec latency, 215 million message per second
- PCIe Gen3 / Gen4, 32 lanes
- Integrated PCIe switch
- Multi-Host - up to 8 hosts, supporting 4 dual-socket servers

Leading Features

- In-network computing and memory for HPC collective offloads
- Security – Block-level encryption to storage, key management, FIPS
- Storage – NVMe Emulation, NVMe-oF target, Erasure coding, T10/DIF

ConnectX[®]·6



HDR InfiniBand Switches

40 QSFP56 ports

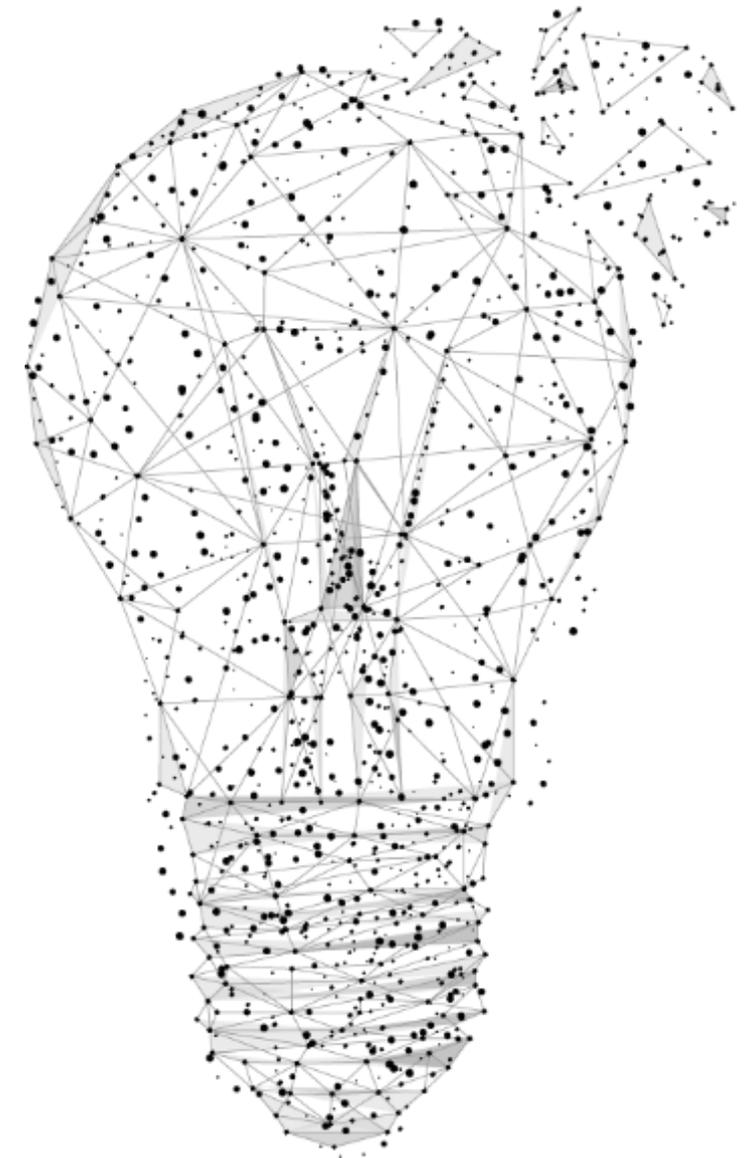
- 40 ports of HDR, 200G
- 80 ports of HDR100, 100G

800 QSFP56 ports

- 800 ports of HDR, 200G
- 1600 ports of HDR100, 100G

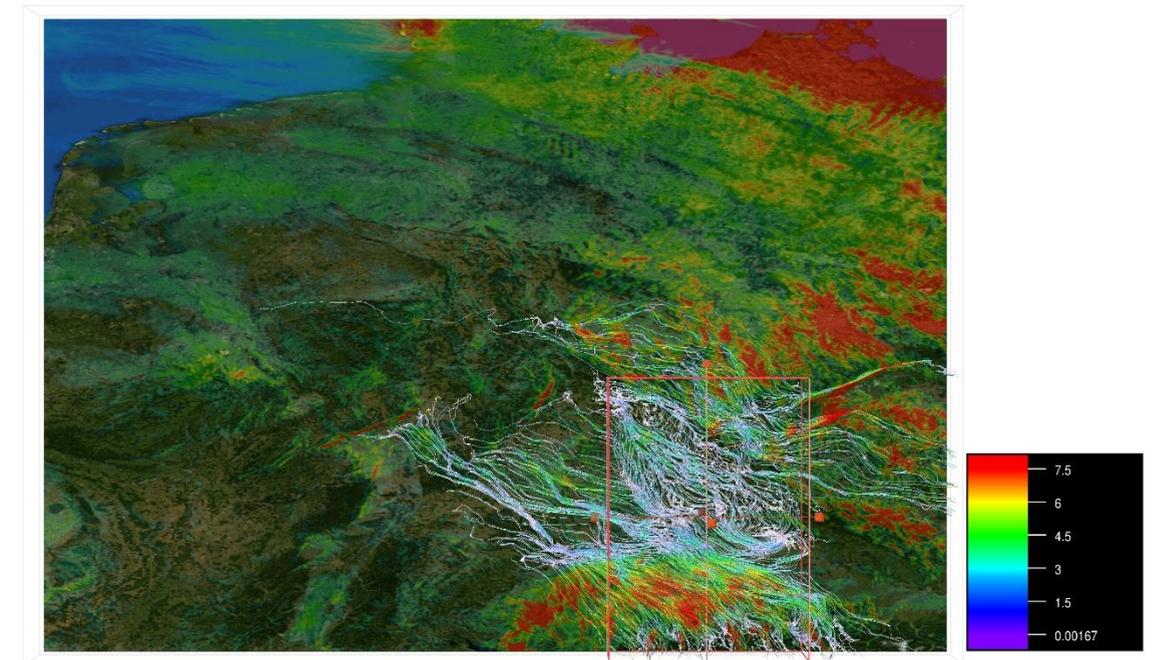
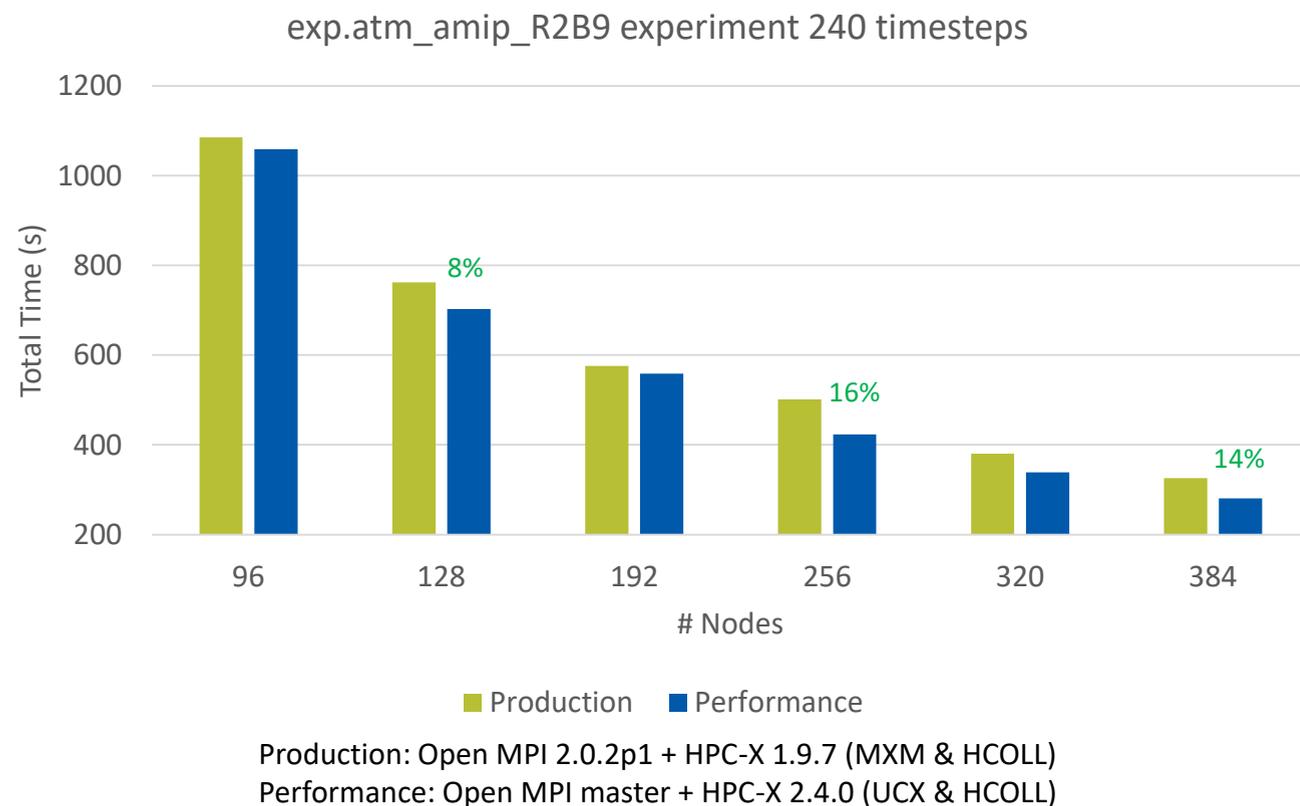


Weather and Climate Apps



ICON (ICOsahedral Non-hydrostatic model)

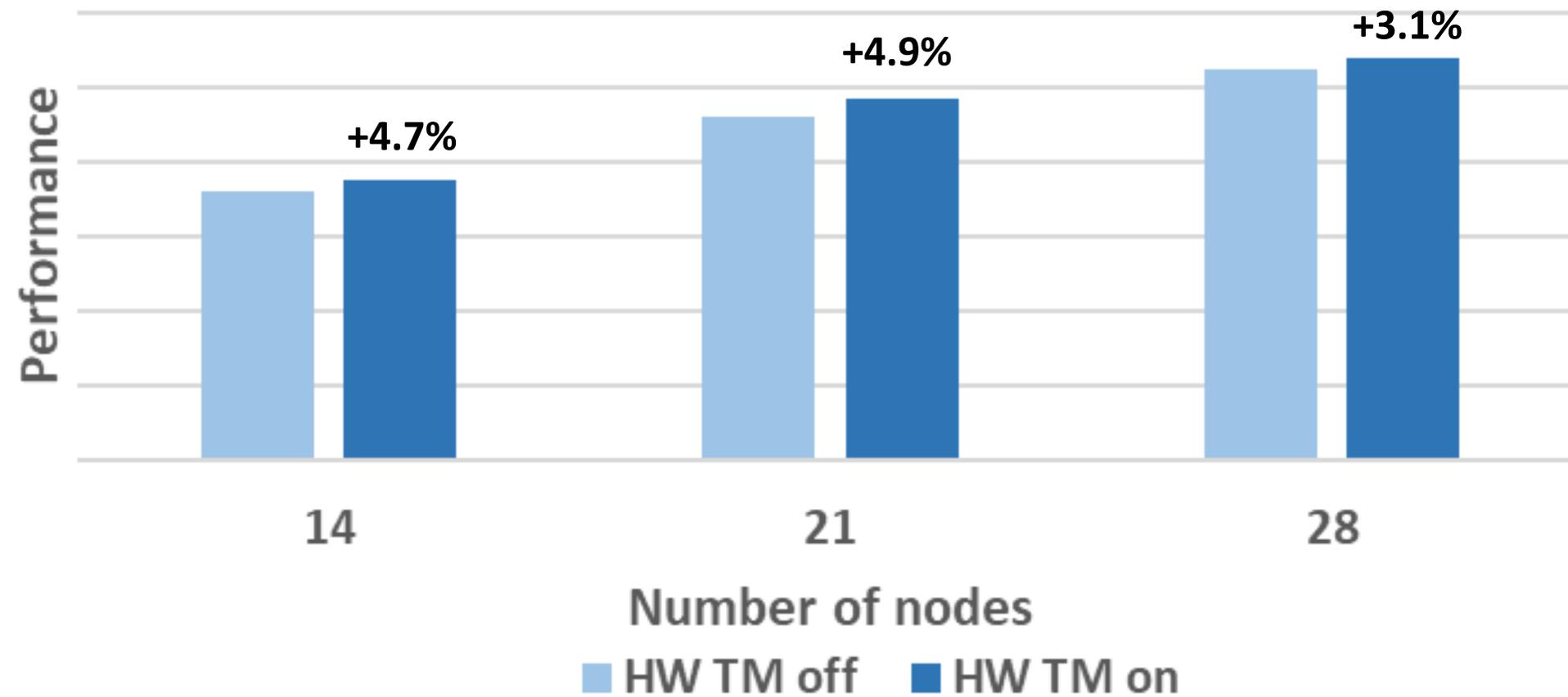
- New generation unified weather forecasting and climate model developed by MPI-M and DWD
- New data exchange module YAXT developed by DRKZ to replace traditional halo exchange mechanism
 - Main challenge lies on efficient handling of sparse data at scale
 - Improvement jointly developed by DKRZ, UTK, and Mellanox



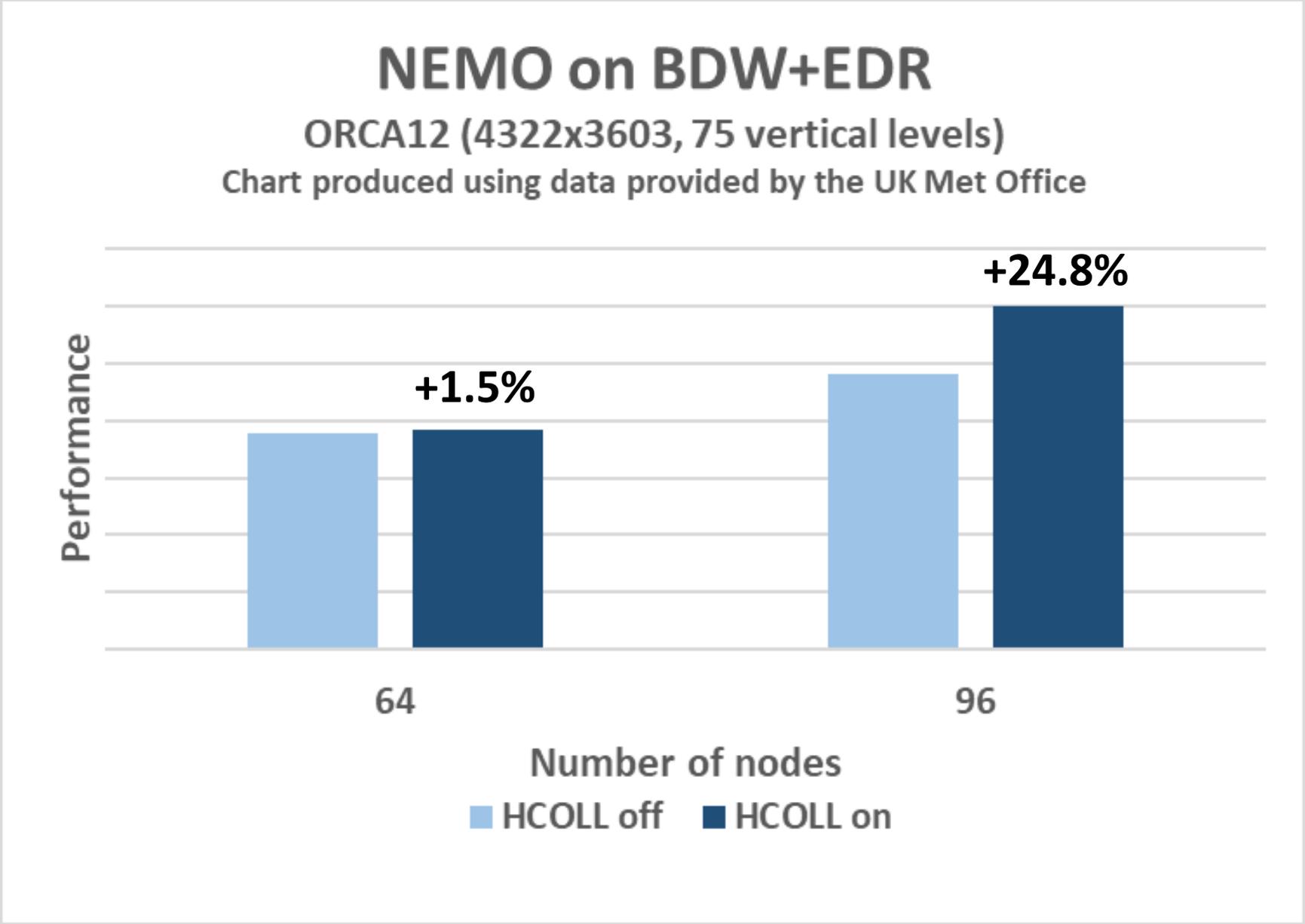
UM Global Climate on BDW+EDR

UKESM N96 (192x144, 85 lev)

Chart produced using Met Office Software

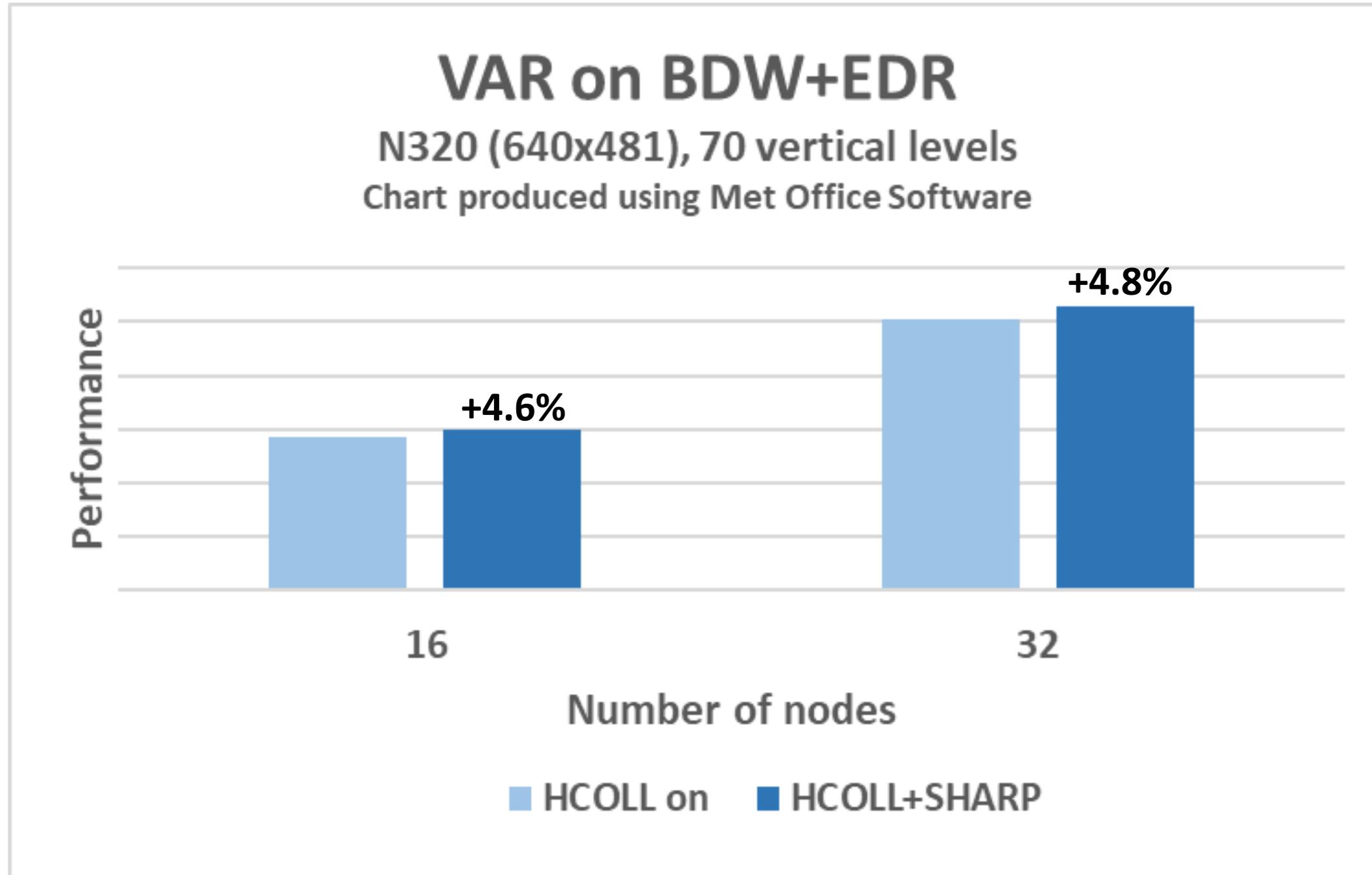


NEMO



NEMO ocean engine, Madec Gurvan and NEMO System Team, *Issue 27, Scientific Notes of Climate Modelling Center*, Institut Pierre-Simon Laplace (IPSL), ISSN 1288-1619

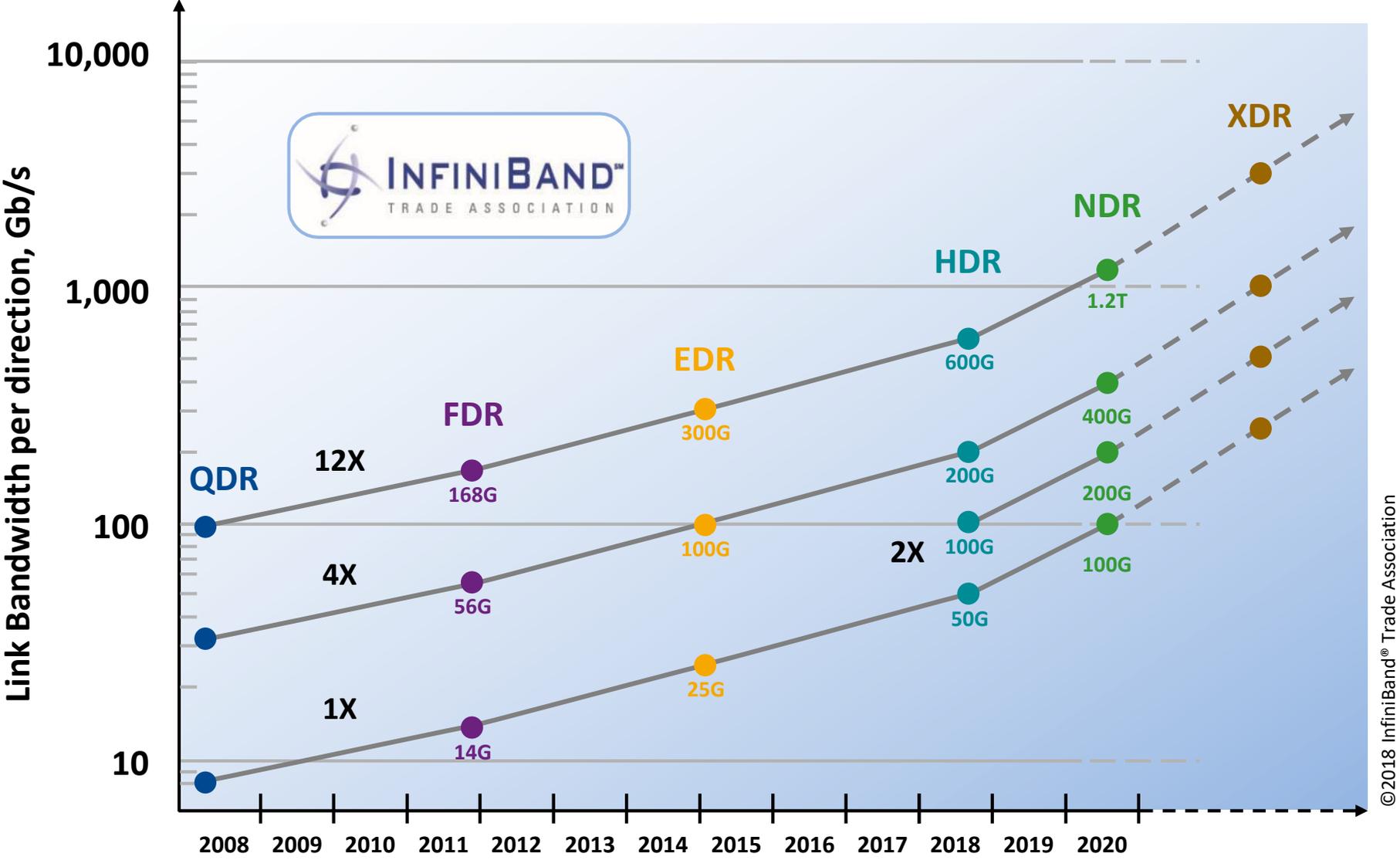
VAR

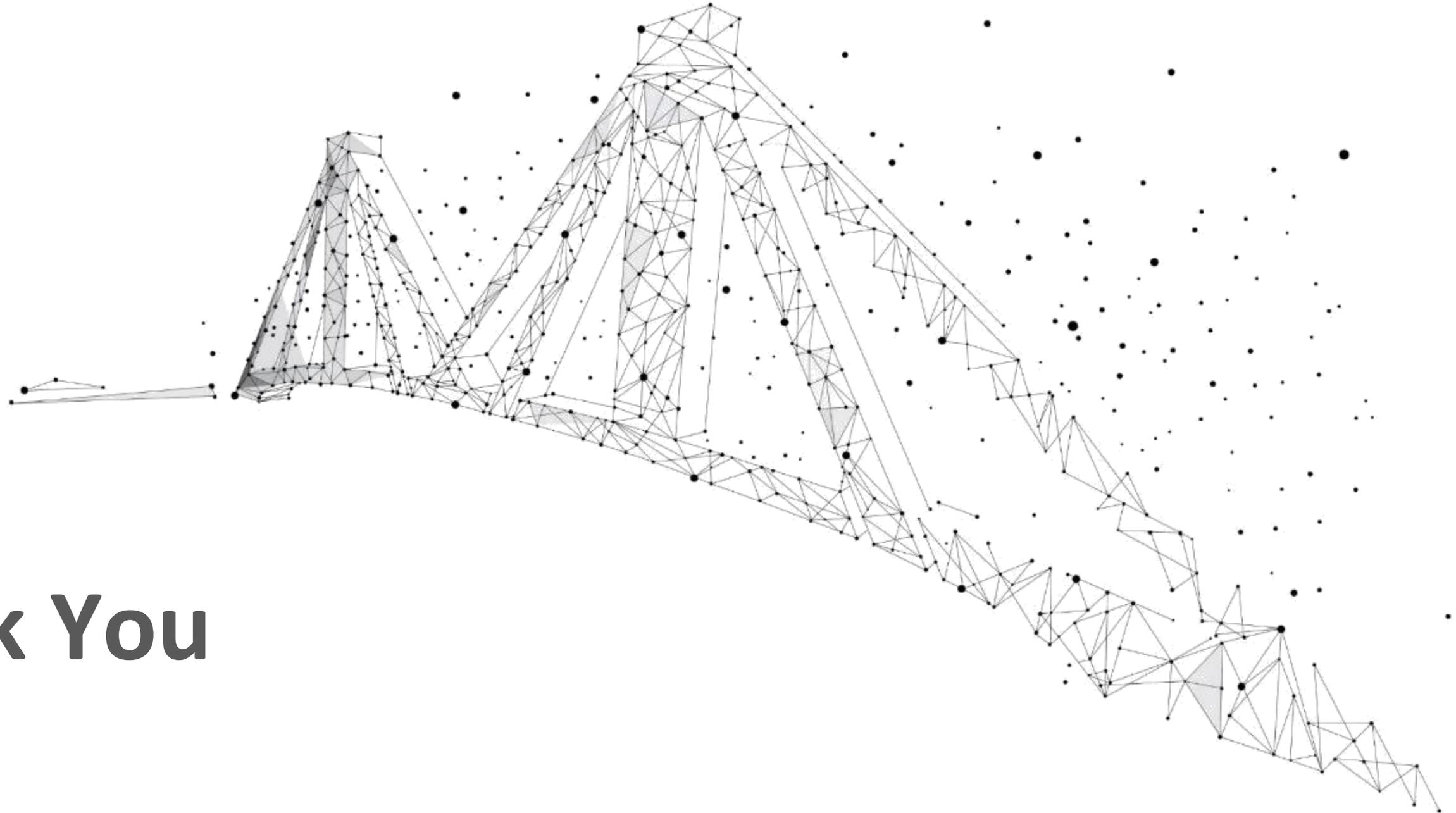


InfiniBand Roadmap



InfiniBand Roadmap (IBTA)





Thank You

