# In-Network Computing

## Paving the Road to Exascale

August 2018

# Mellanox Accelerates Leading HPC and AI Systems

## World's Top 3 Supercomputers
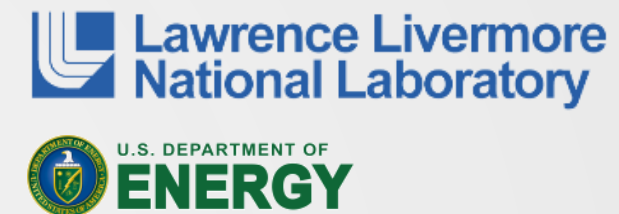


**1**

Summit CORAL System
World's Fastest HPC / AI System
9.2K InfiniBand Nodes



**2**

Wuxi Supercomputing Center
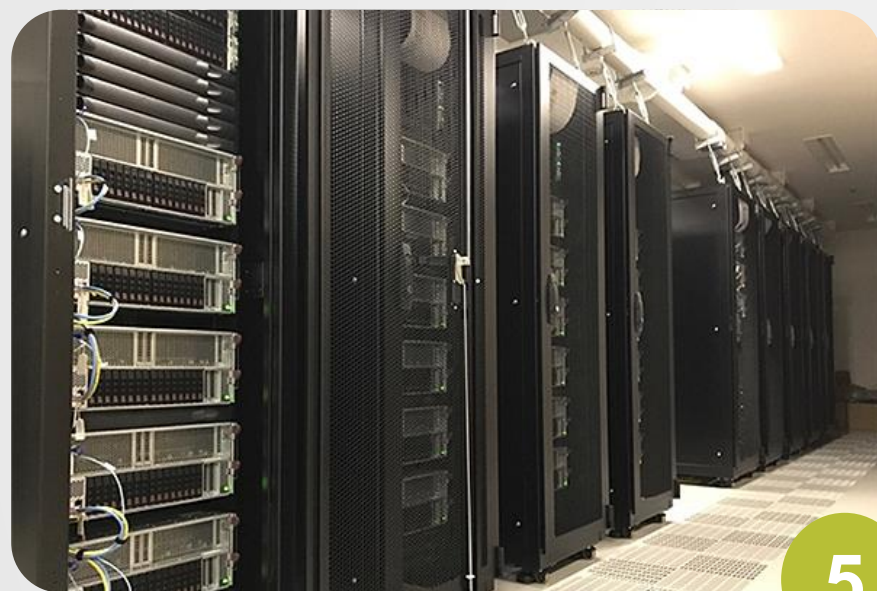Fastest Supercomputer in China
41K InfiniBand Nodes



**3**

Sierra CORAL System
#2 USA Supercomputer
8.6K InfiniBand Nodes
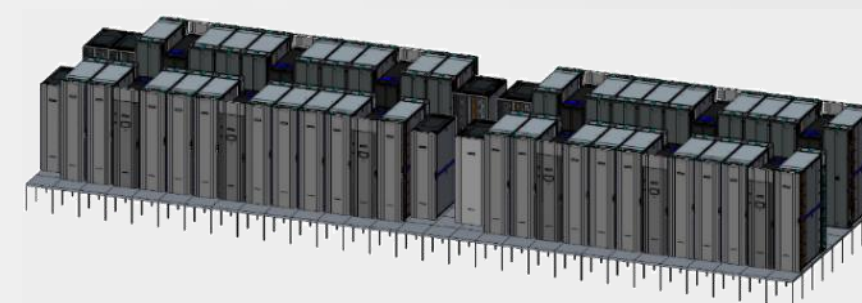
# Mellanox Accelerates Leading HPC and AI Systems

(Examples)



**5**

Fastest HPC / AI System in Japan
1.1K InfiniBand Nodes

**13**

The world's Fastest Industry Supercomputer
1.6K InfiniBand Nodes

'Astra' Arm-Based Supercomputer
NNSA Vanguard Program
2.6K InfiniBand Nodes

To be Listed Nov'18 (TOP100)
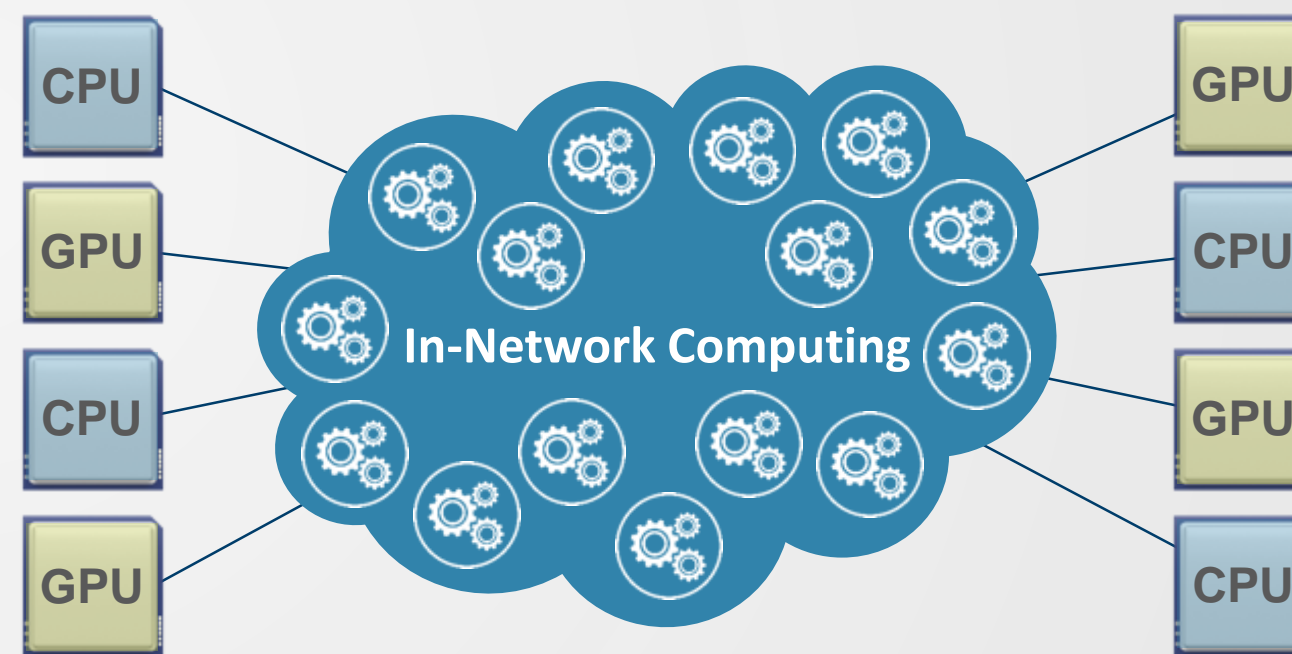
# The Need for Intelligent and Faster Interconnect

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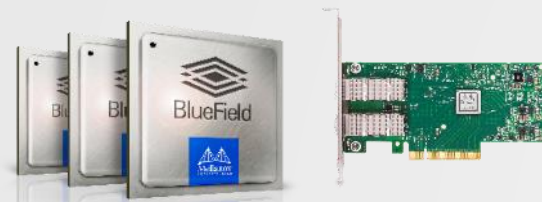Enable Higher Performance and Scale

**CPU-Centric (Onload)**　　　　　　　　　**Data-Centric (Offload)**



**Must Wait for the Data**
**Creates Performance Bottlenecks**

Analyze Data as it Moves!
Higher Performance and Scale

# HPC and AI Needs the Most Intelligent Interconnect

**SmartNIC**

**System on a Chip**

BlueField

**Adapters**

Innova ConnectX

**Switches**

Spectrum Quantum

**Cables & Transceivers**

LinkX

**Higher** Data Speeds

**Faster** Data Processing

**Better** Data Security
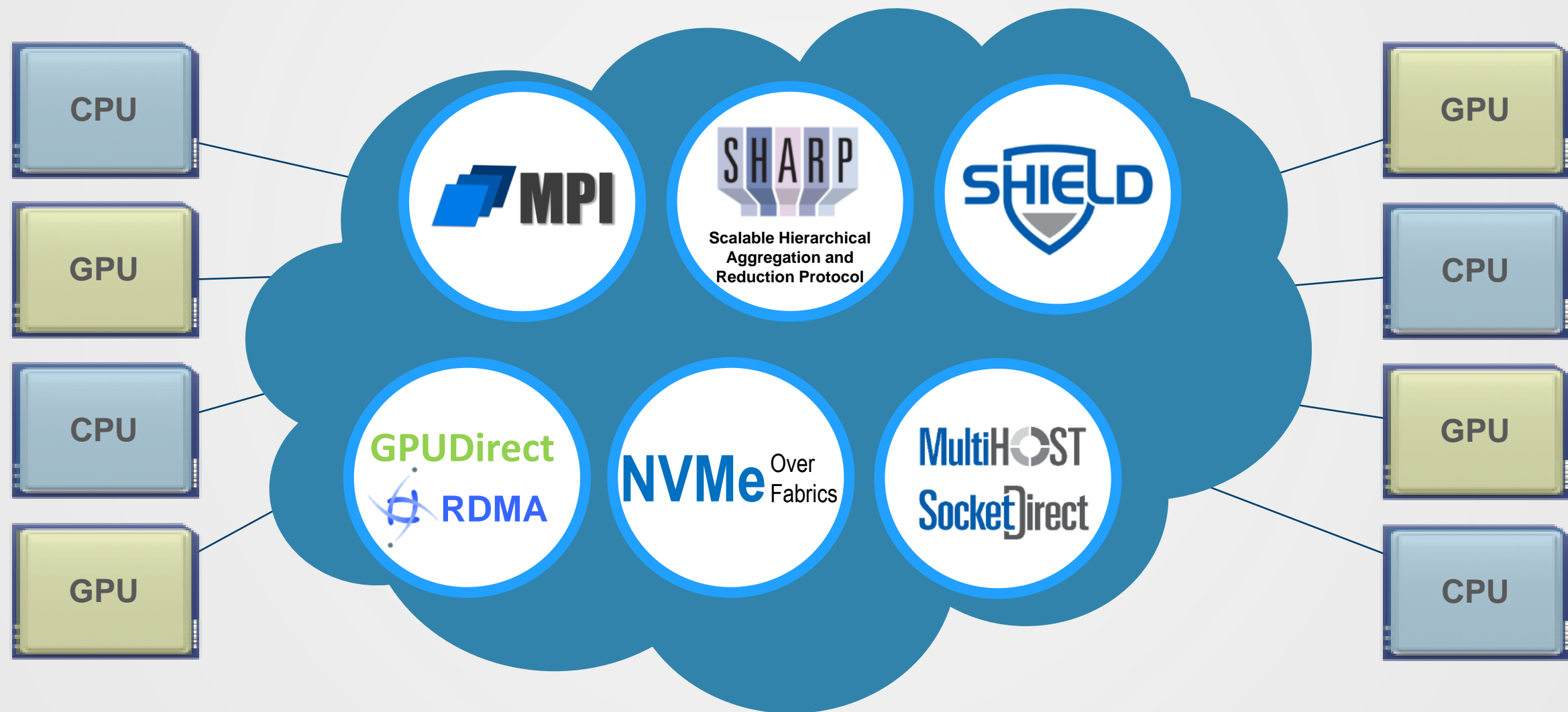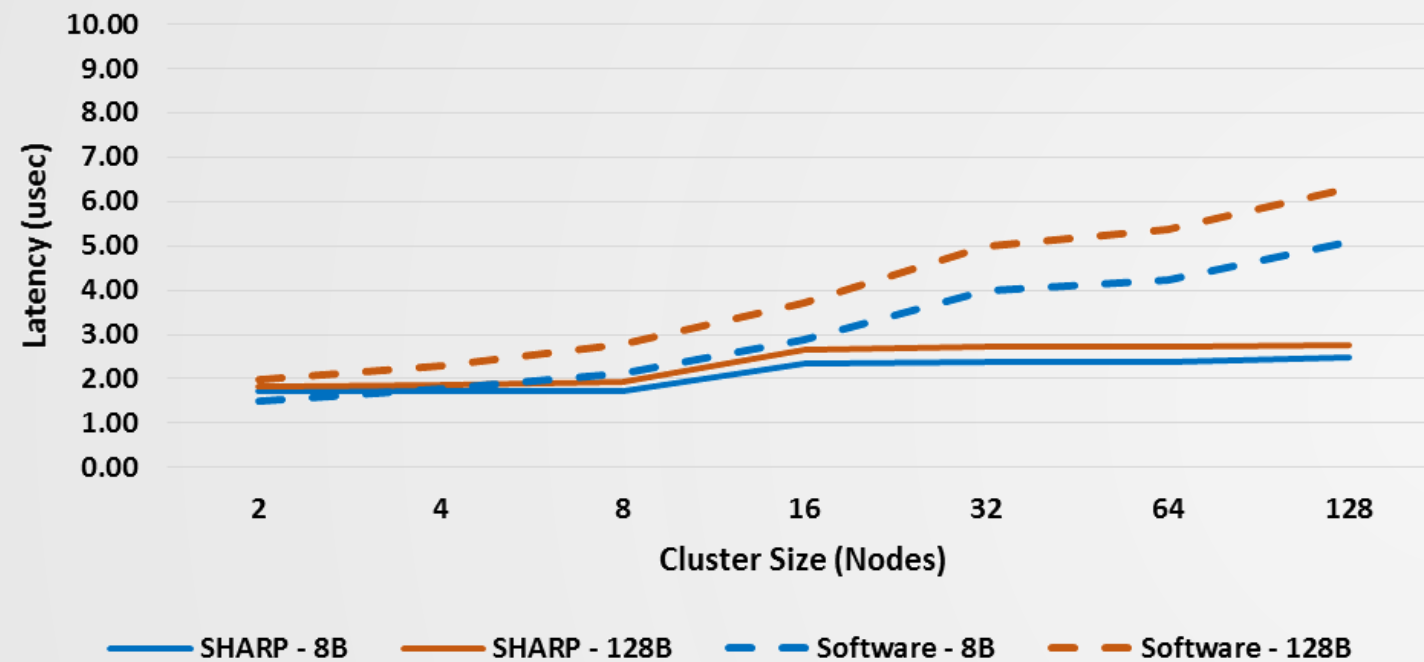
# In-Network Computing to Enable Data-Centric Data Centers
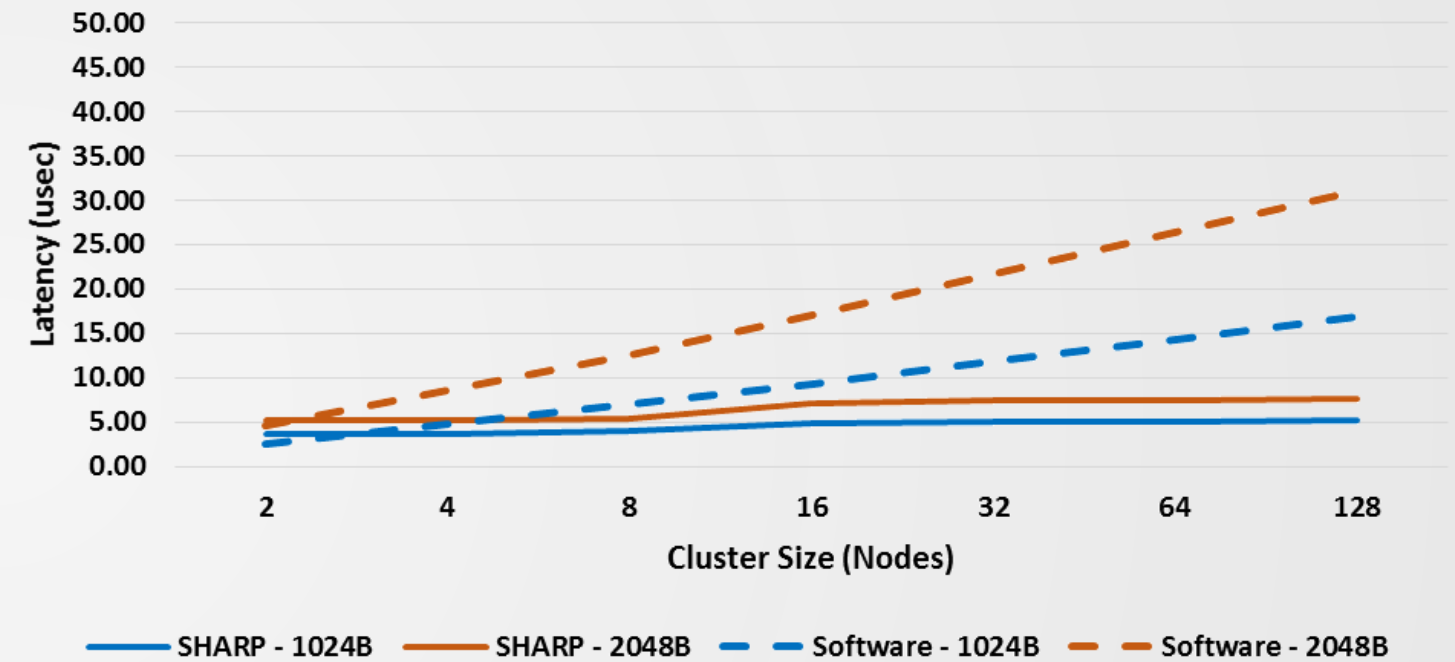
# SHARP AllReduce Performance Advantages (128 Nodes)
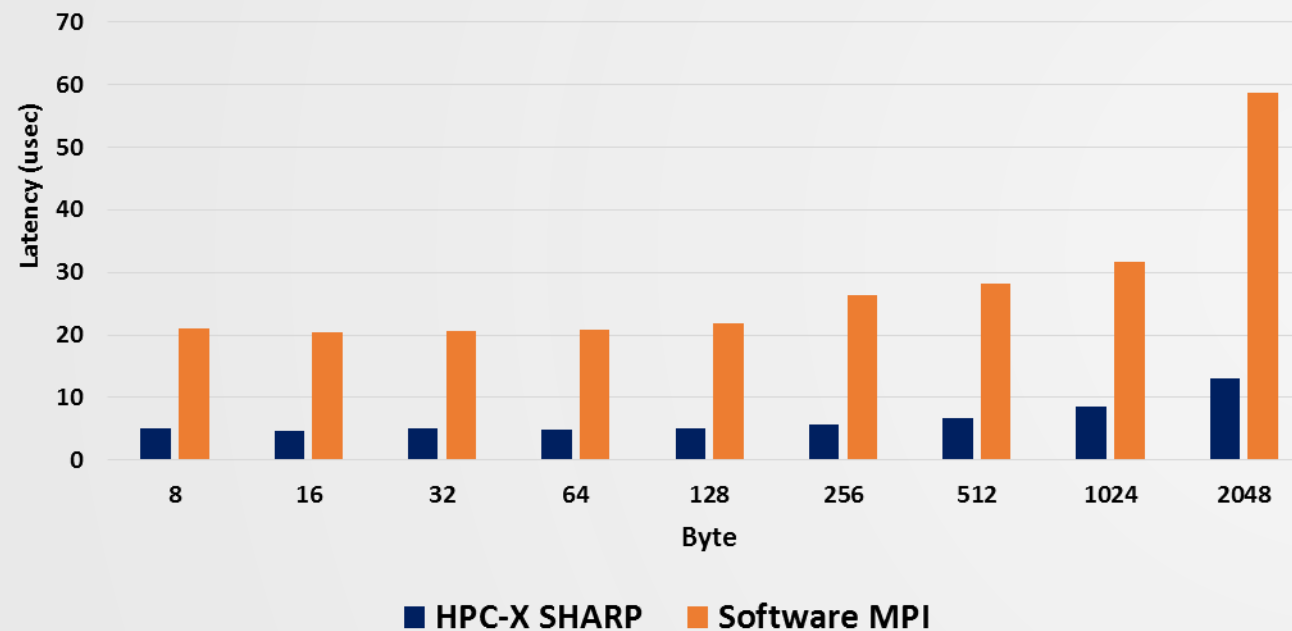


SHARP enables 75% Reduction in Latency
Providing Scalable Flat Latency
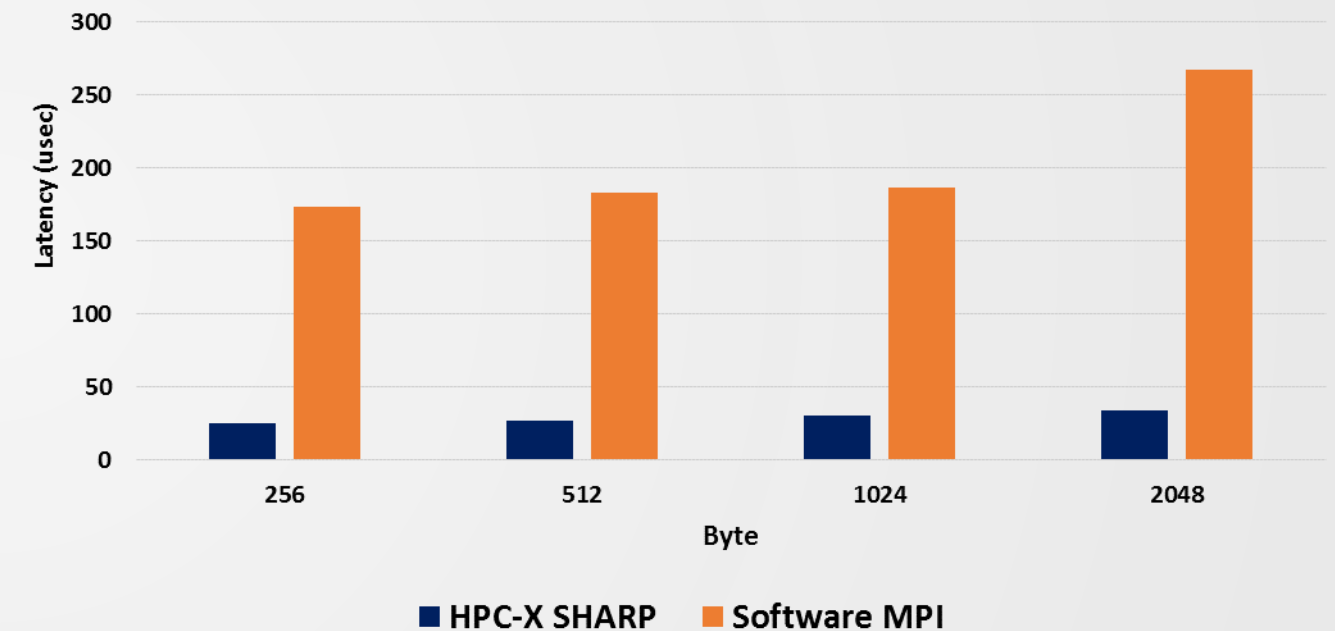
# SHARP AllReduce Performance Advantages
## 1500 Nodes, 60K MPI Ranks, Dragonfly+ Topology



**MPI AllReduce Latency**
**1500 Nodes, 1PPN**

■ HPC-X SHARP   ■ Software MPI

**MPI AllReduce Latency**
**1500 Nodes, 40PPN, 60K MPI Ranks**

■ HPC-X SHARP   ■ Software MPI

## SHARP Enables Highest Performance

# SHARP Accelerates AI Performance

The CPU in a parameter server
becomes the bottleneck



**SHARP**
Scalable Hierarchical
Aggregation and
Reduction Protocol

parameter server (sharded) $w' = u(w, \nabla w)$

$\nabla w$ | $w$

Training Agent    Training Agent    Training Agent    Training Agent

Performs the Gradient Averaging
Replaces all physical parameter servers
Accelerate AI Performance

# Network Topologies

# Supporting Variety of Topologies



**Fat Tree**

**Hypercube**

**Torus**

**Dragonfly**

# Dragonfly+ Topology

- Several "groups", connected using all to all links
- The topology inside each group can be any topology
- Reduce total cost of network (fewer long cables)
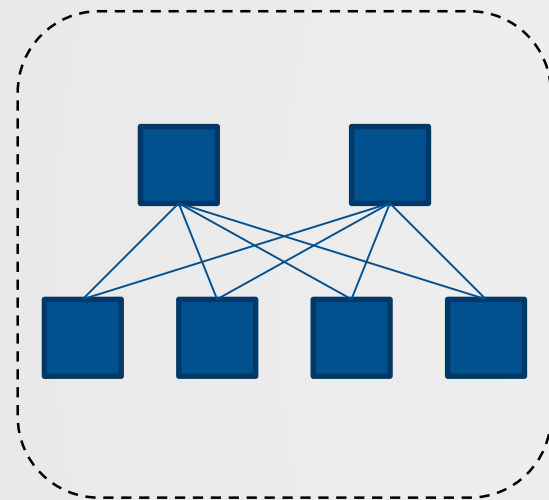- Utilizes Adaptive Routing to for efficient operations
- Simplifies future system expansion

**Full-Graph connecting every group to all other groups**

**1200-Nodes Dragonfly+ Systems Example**

ENABLER OF CO-DESIGN

UCF

Unified Communication X (UCX)

August 2018

# UCF Consortium

- Mission:
  - Collaboration between industry, laboratories, and academia to create production grade communication frameworks and open standards for data centric and high-performance applications

- Projects
  - UCX – Unified Communication X
  - Open RDMA

- Board members
  - **Jeff Kuehn**, UCF Chairman (Los Alamos National Laboratory)
  - **Gilad Shainer**, UCF President (Mellanox Technologies)
  - **Pavel Shamis**, UCF treasurer (ARM)
  - **Brad Benton**, Board Member (AMD)
  - **Duncan Poole**, Board Member (Nvidia)
  - **Pavan Balaji**, Board Member (Argonne National Laboratory)
  - **Sameh Sharkawi**, Board Member (IBM)
  - **Dhabaleswar K. (DK) Panda**, Board Member (Ohio State University)
  - **Steve Poole**, Board Member (Open Source Software Solutions)

# UCX Framework Mission

- Collaboration between industry, laboratories, government (DoD, DoE), and academia
- Create open-source production grade communication framework for HPC applications
- Enable the highest performance through co-design of software-hardware interfaces

| API | Performance oriented | Production quality |
|---|---|---|
| Exposes broad semantics that target data centric and HPC programming models and applications | Optimization for low-software overheads in communication path allows near native-level performance | Developed, maintained, tested, and used by industry and researcher community |

| Community driven | Research | Cross platform |
|---|---|---|
| Collaboration between industry, laboratories, and academia | The framework concepts and ideas are driven by research in academia, laboratories, and industry | Support for Infiniband, Cray, various shared memory (x86-64, Power, ARMv8), GPUs |

**Co-design of Exascale Network APIs**

# UCX Framework

- UCX is a framework for network APIs and stacks

- UCX aims to unify the different network APIs, protocols and implementations into a single framework that is portable, efficient and functional

- UCX doesn't focus on supporting a single programming model, instead it provides APIs and protocols that can be used to tailor the functionalities of a particular programming model efficiently

- When different programming paradigms and applications use UCX to implement their functionality, it increases their portability. As just implementing a small set of UCX APIs on top of a new hardware ensures that these applications can run seamlessly without having to implement it themselves

# UCX High-level Overview

**Applications**

MPICH, Open-MPI, etc.

OpenSHMEM, UPC, CAF, X10, Chapel, etc.

Parsec, OCR, Legions, etc.

Burst buffer, ADIOS, etc.

## UCX

### UC-P (Protocols) - High Level API
Transport selection, cross-transrport multi-rail, fragmentation, operations not supported by hardware

| Message Passing API Domain: tag matching, randevouze | PGAS API Domain: RMAs, Atomics | Task Based API Domain: Active Messages | I/O API Domain: Stream |

### UC-T (Hardware Transports) - Low Level API
RMA, Atomic, Tag-matching, Send/Recv, Active Message

| Transport for InfiniBand VERBs driver | Transport for Gemini/Aries drivers | Transport for intra-node host memory communication | Transport for Accelerator Memory communcation |
| RC  UD  XRC  DCT | GNI | SYSV  POSIX  KNEM  CMA  XPMEM | GPU |

### UC-S (Services)
Common utilities

Utilities  Data stractures

Memory Management

OFA Verbs Driver

Cray Driver

OS Kernel

Cuda

**Hardware**

# Thank You