



## Interconnect Your Future

Smart Interconnect for Next Generation HPC Platforms

Gilad Shainer, August 2016, 4th Annual MVAPICH User Group (MUG) Meeting





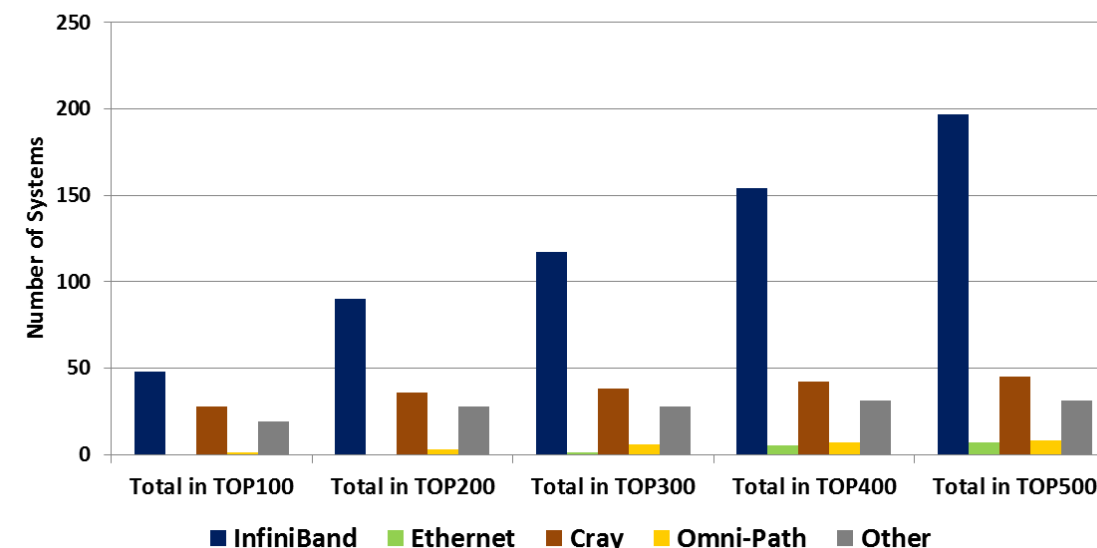
## National Supercomputing Center in Wuxi, China #1 on the TOP500 Supercomputing List

- 93 Petaflop performance, 3X higher versus #2 on the TOP500
- 41K nodes, 10 million cores, 256 cores per CPU
- Mellanox adapter and switch solutions

- The TOP500 list has evolved, includes HPC & Cloud / Web2.0 Hyperscale systems
- Mellanox connects 41.2% of overall TOP500 systems
- Mellanox connects 70.4% of the TOP500 HPC platforms
- Mellanox connects 46 Petascale systems, Nearly 50% of the total Petascale systems

**InfiniBand is the Interconnect of Choice for  
HPC Compute and Storage Infrastructures**

TOP500 - TOP 100, 200, 300, 400, 500 Systems Distribution  
HPC Systems Only



 **OAK RIDGE**  
National Laboratory

**“Summit” System**



 **Lawrence Livermore**  
**National Laboratory**

**“Sierra” System**



**Proud to Pave the Path to Exascale**

## Performance Development

Terascale



Petascale

1<sup>st</sup>



“Roadrunner”



Exascale



2000

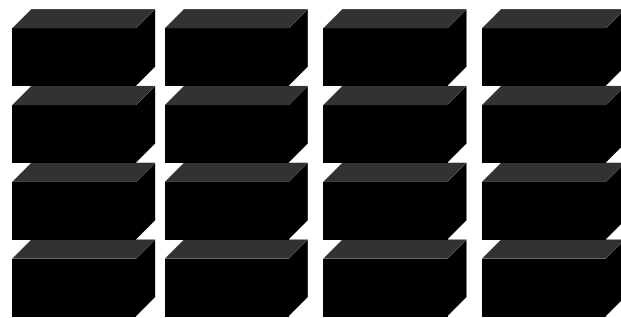
2005

2010

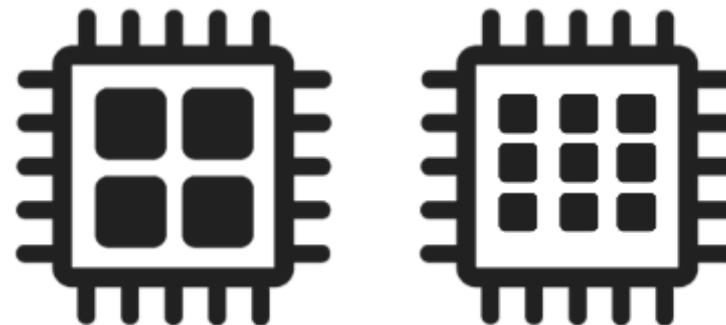
2015

2020

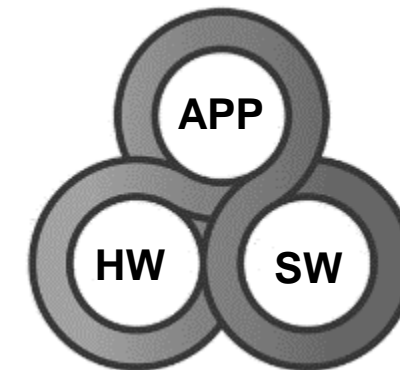
## The Interconnect is the Enabling Technology



SMP to Clusters



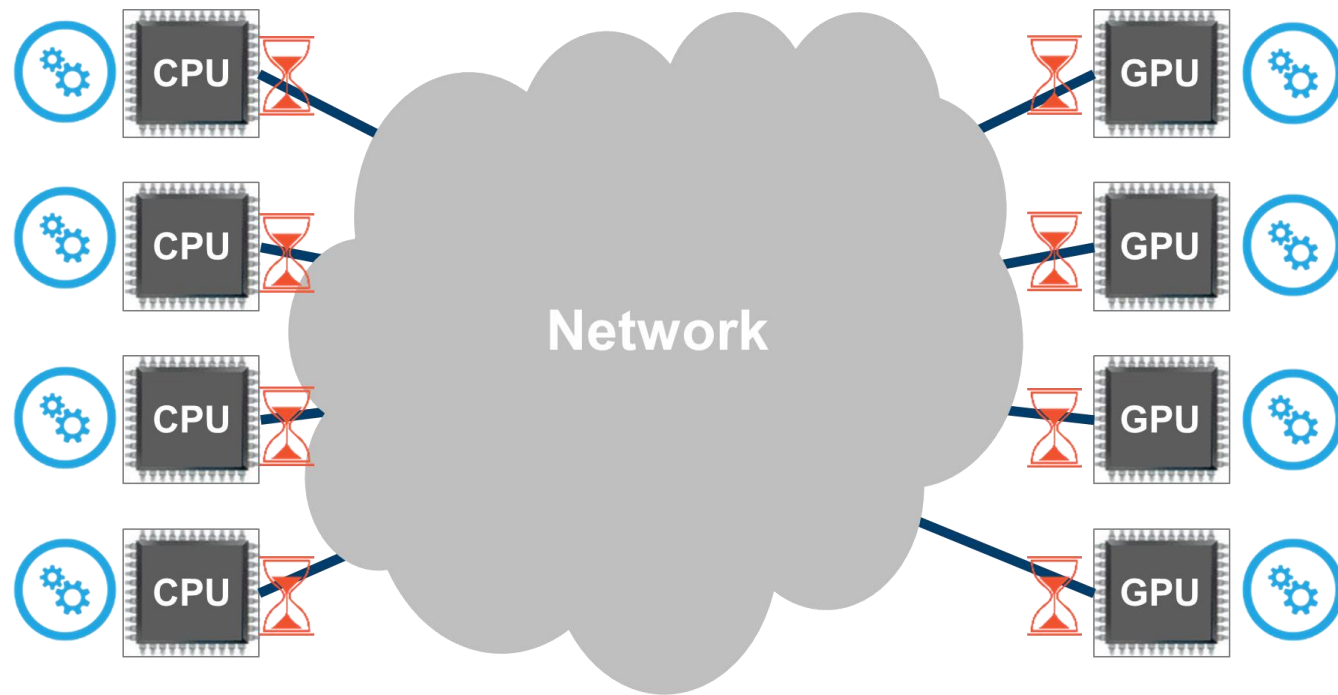
Single-Core to Many-Core



Application  
Software  
Hardware

Co-Design

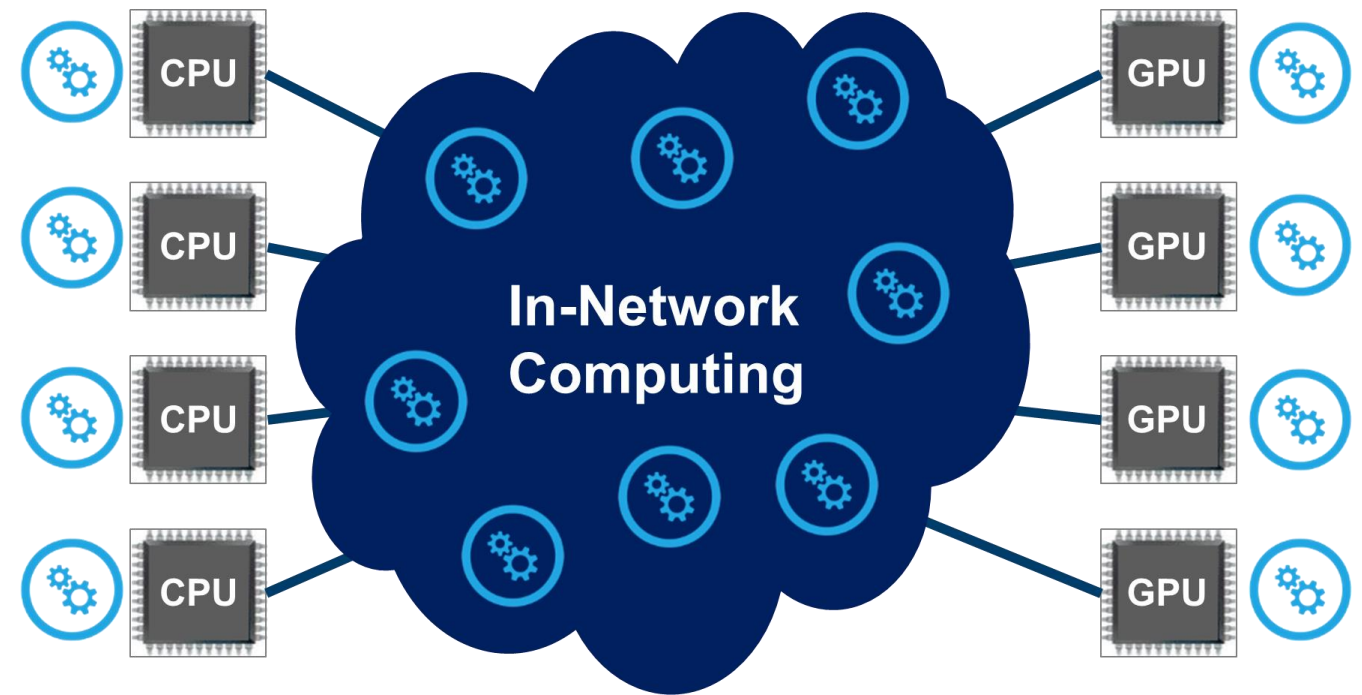
## CPU-Centric



Limited to Main CPU Usage  
Results in Performance Limitation

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

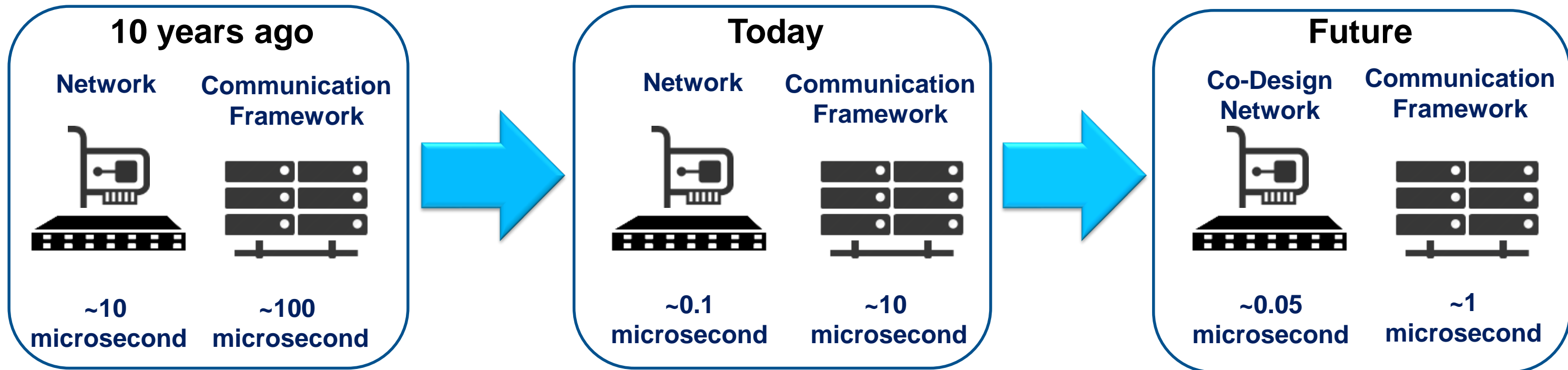
## Co-Design



Creating Synergies  
Enables Higher Performance and Scale

**Work on The Data as it Moves  
Enables Performance and Scale**

# Breaking the Application Latency Wall



- Today: Network device latencies are on the order of 100 nanoseconds
- Challenge: Enabling the next order of magnitude improvement in application performance
- Solution: Creating synergies between software and hardware – intelligent interconnect

**Intelligent Interconnect Paves the Road to Exascale Performance**

# State of the Smart

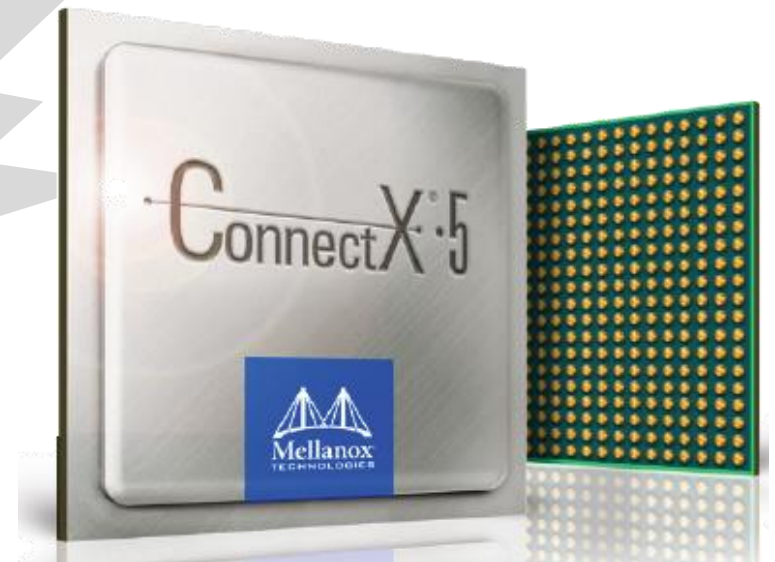
a new generation of co-processors emerges

# Mellanox Smart Interconnect

Switch-IB™ 2



ConnectX® 5



# Switch-IB 2 and ConnectX-5 Smart Interconnect Solutions



SHArP Enables Switch-IB 2 to Execute Data Aggregation / Reduction Operations in the Network

Barrier, Reduce, All-Reduce, Broadcast  
Sum, Min, Max, Min-loc, max-loc, OR, XOR, AND  
Integer and Floating-Point, 32 / 64 bit

Delivering **10X** Performance Improvement for MPI  
and SHMEM/PAGS Communications

100Gb/s Throughput  
0.6usec Latency (end-to-end)  
200M Messages per Second

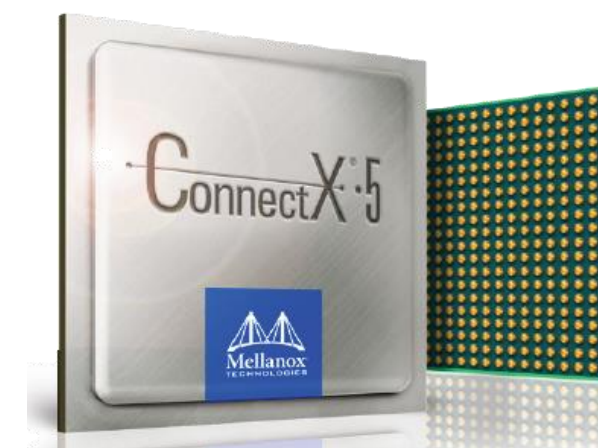
MPI Collectives in Hardware  
MPI Tag Matching in Hardware  
In-Network Memory

PCIe Gen3 and Gen4  
Integrated PCIe Switch  
Advanced Dynamic Routing

Switch-IB™ 2 SHArP

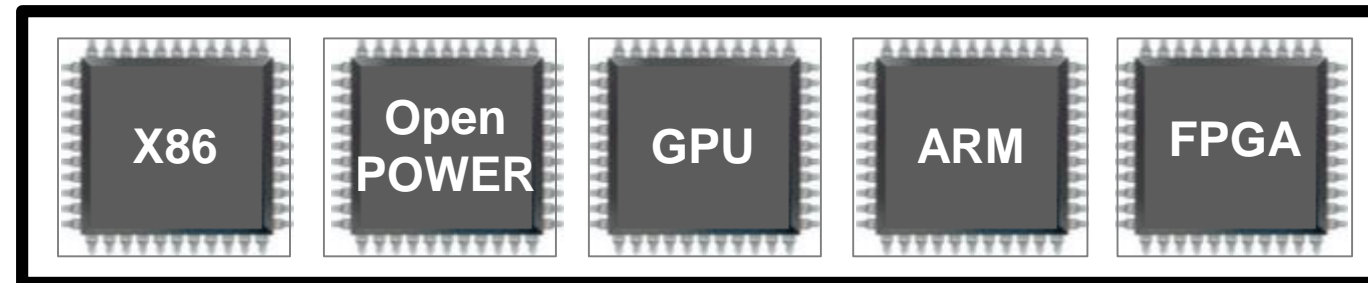


ConnectX® 5





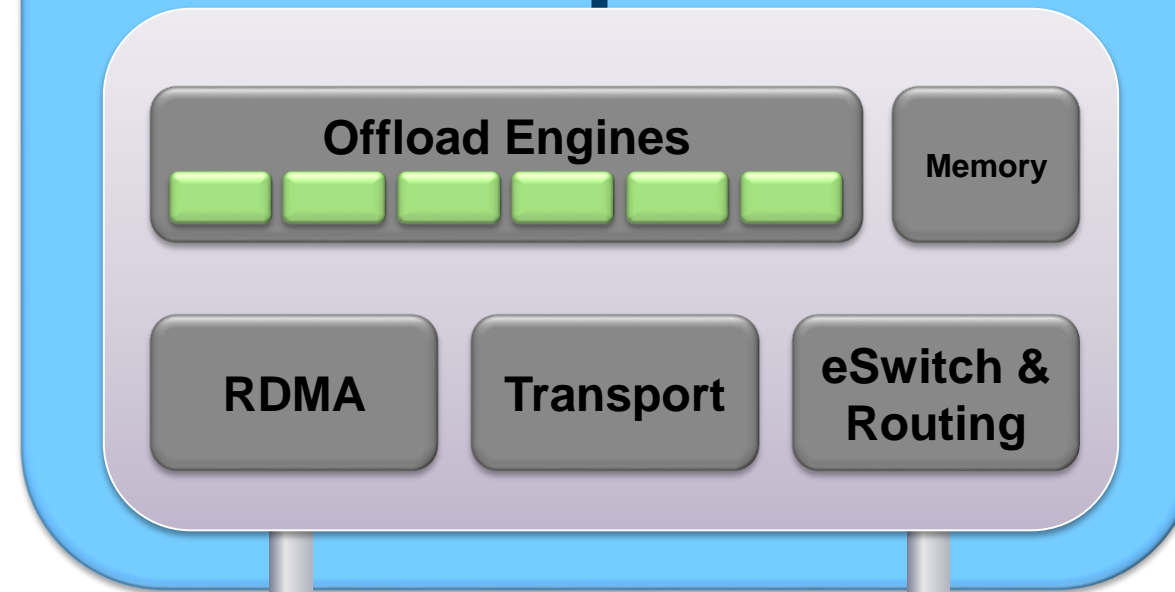
# ConnectX-5 EDR 100G Architecture



PCIe Gen4

Multi-Host Technology  
PCIe Switch

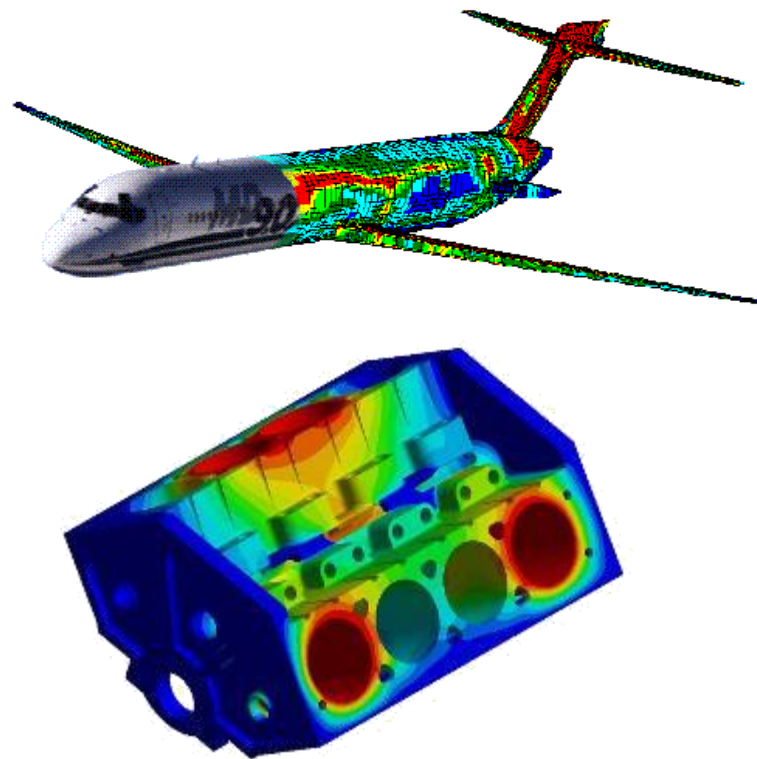
PCIe Gen4



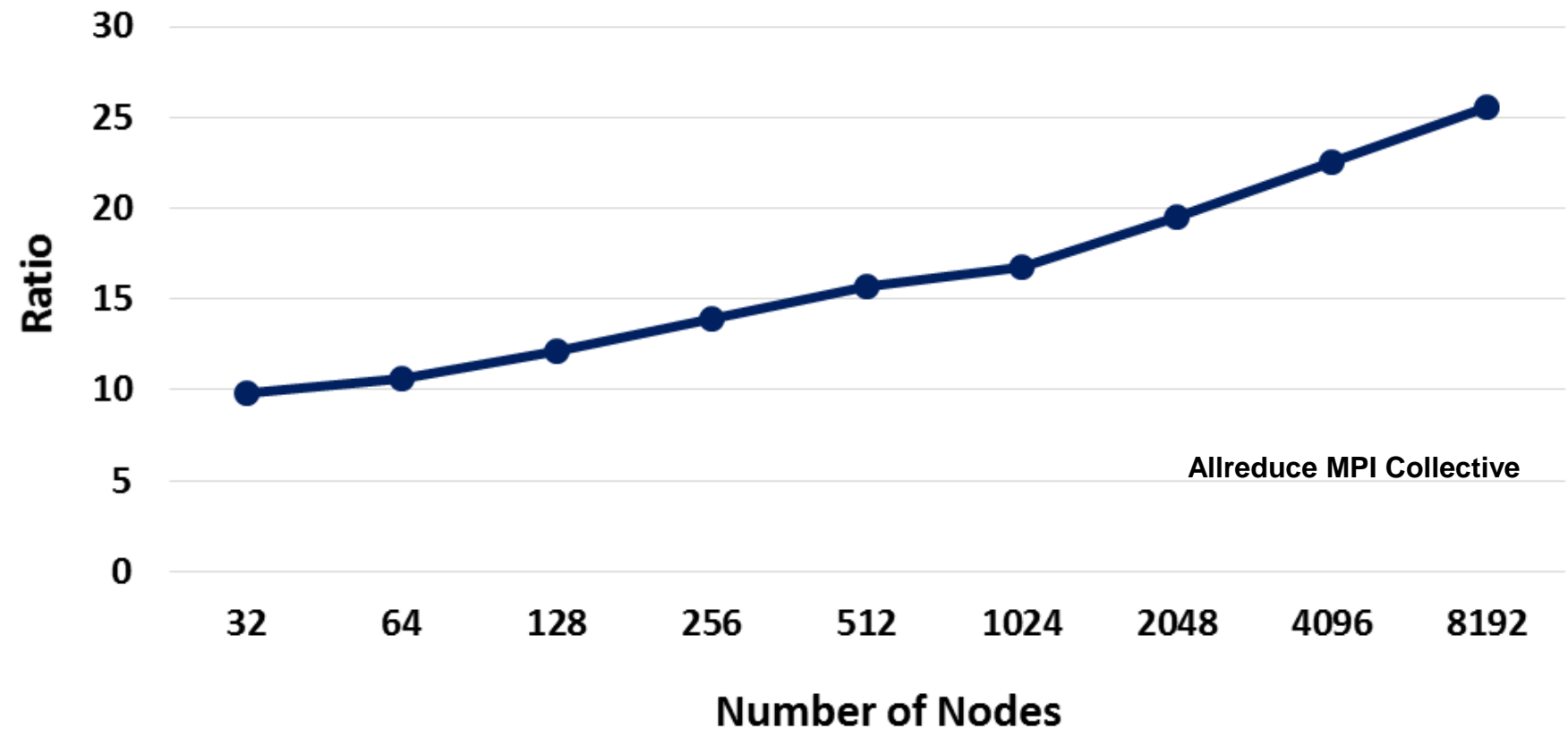
InfiniBand / Ethernet  
10,20,40,50,56,100G

**Breakthrough  
in Performance  
& Total Cost of  
Ownership!**

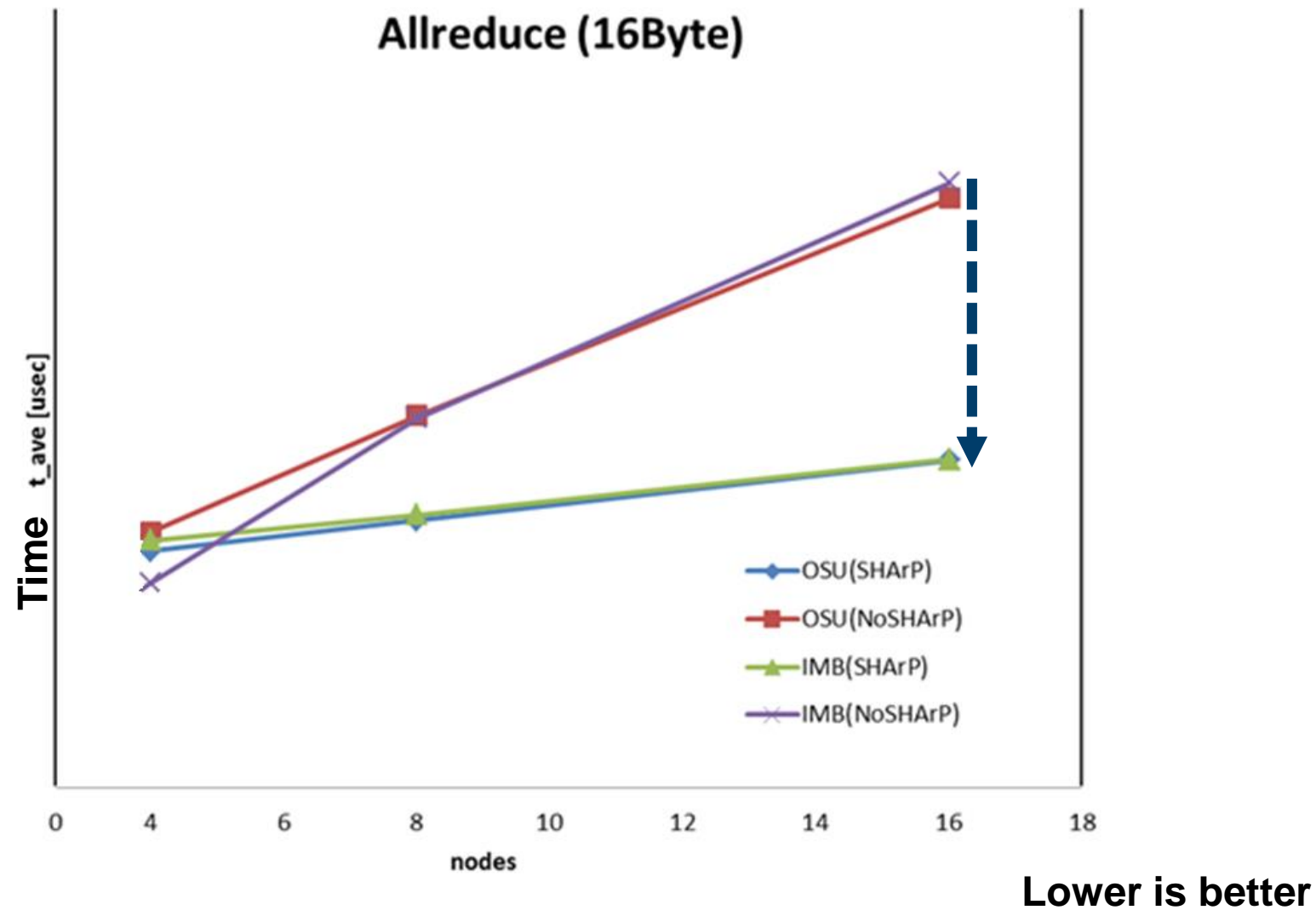
- MiniFE is a Finite Element mini-application
  - Implements kernels that represent implicit finite-element applications



## CPU-based versus Switch Collectives Offloads MiniFE Application - Latency Ratio (8 Bytes)



**10X to 25X** Performance Improvement

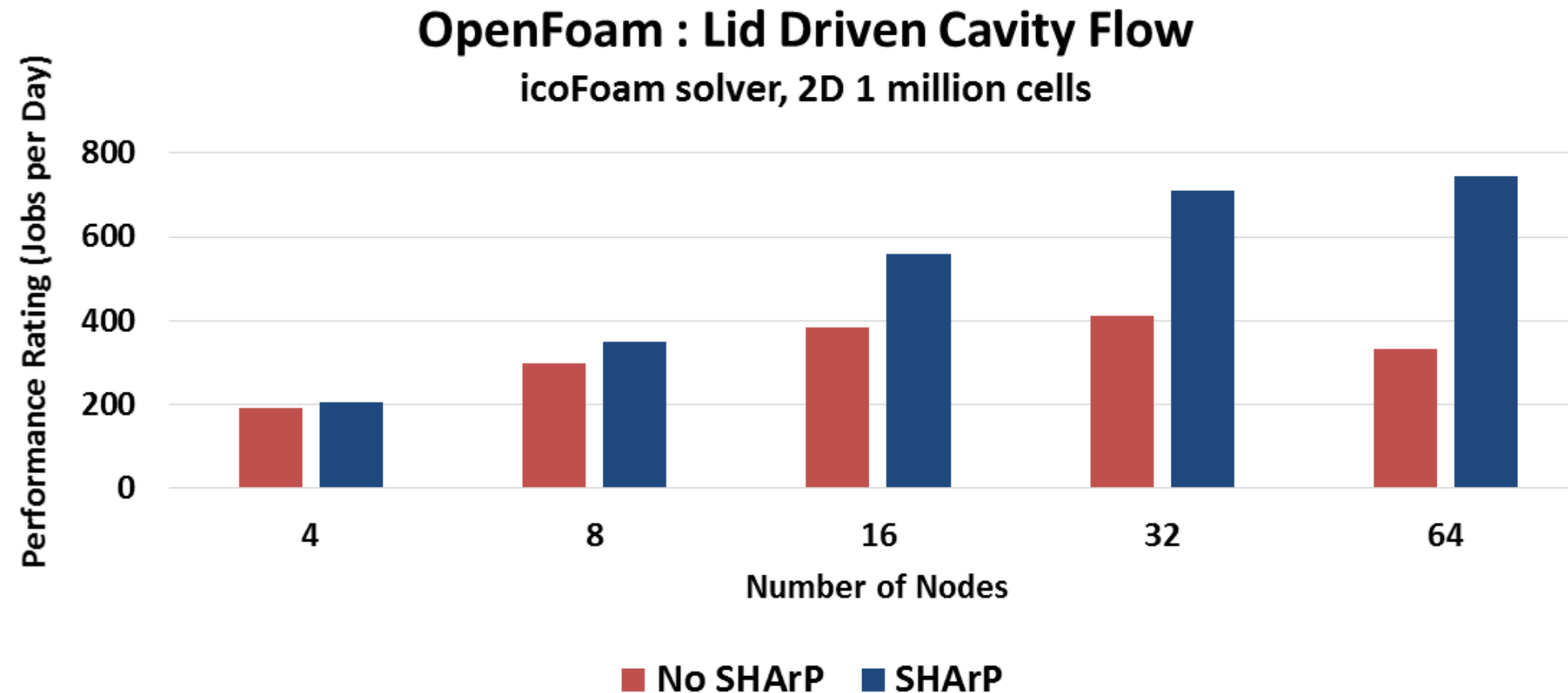


OSU - OSU MPI benchmark; IMB - Intel MPI Benchmark

**Maximizing KNL Performance – 50% Reduction in Run Time  
(Customer Results)**



OpenFOAM is a popular computational fluid dynamics application



SHArP Delivers **2.2X** Higher Performance

# BlueField System-on-a-Chip (SoC) Solution



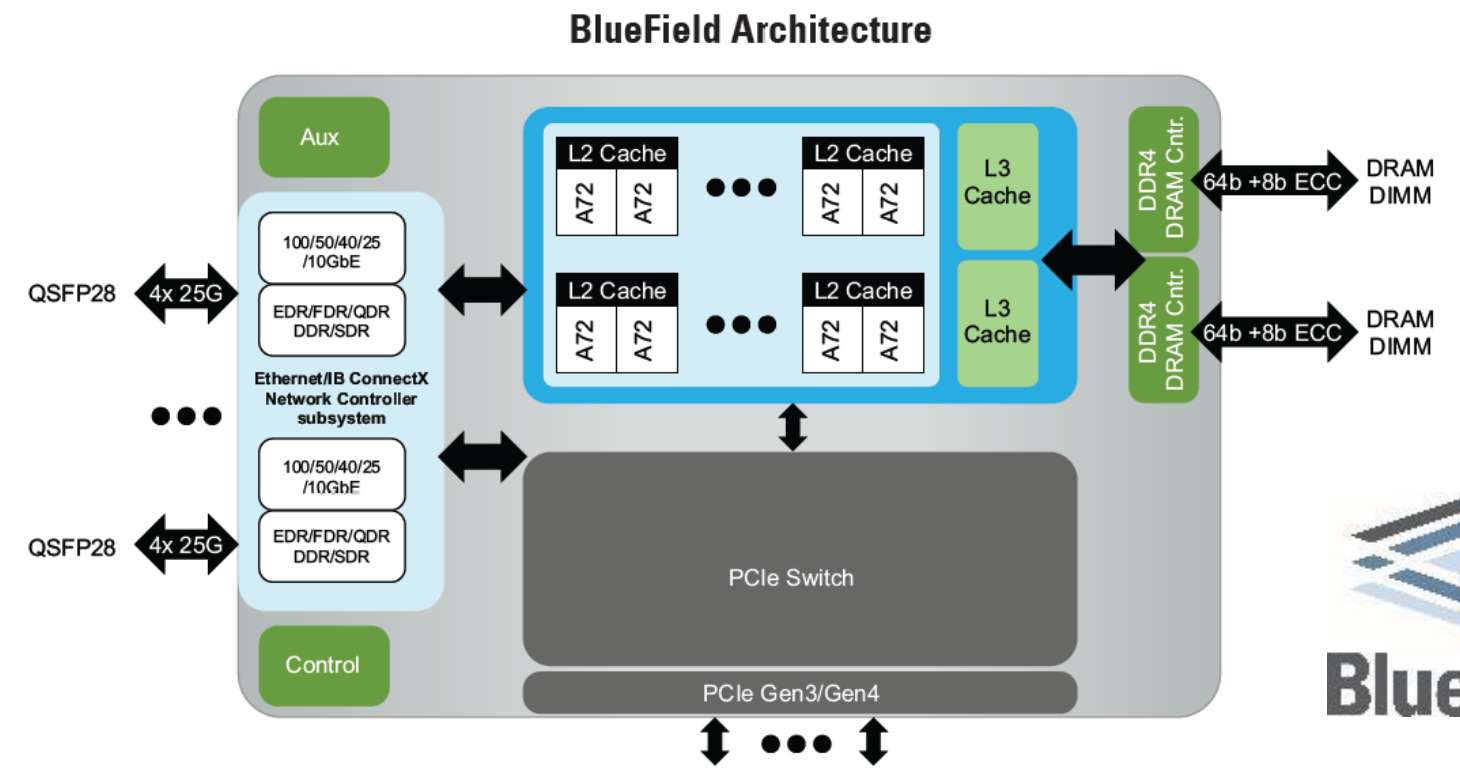
**Storage**

**NVMe Flash Storage Arrays**  
**Scale-Out Storage (NVMe over Fabric)**

**NFV**

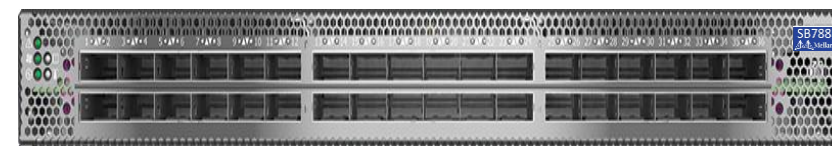
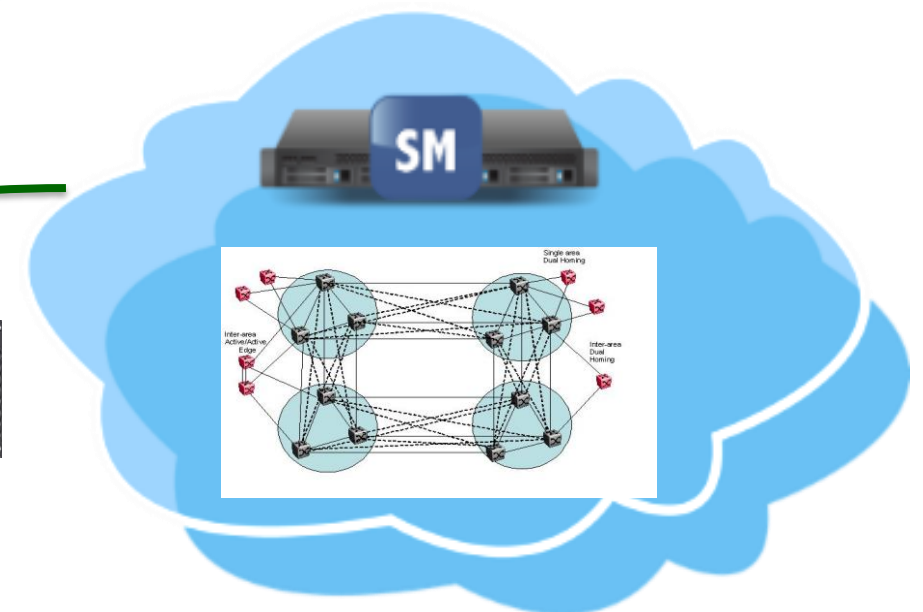
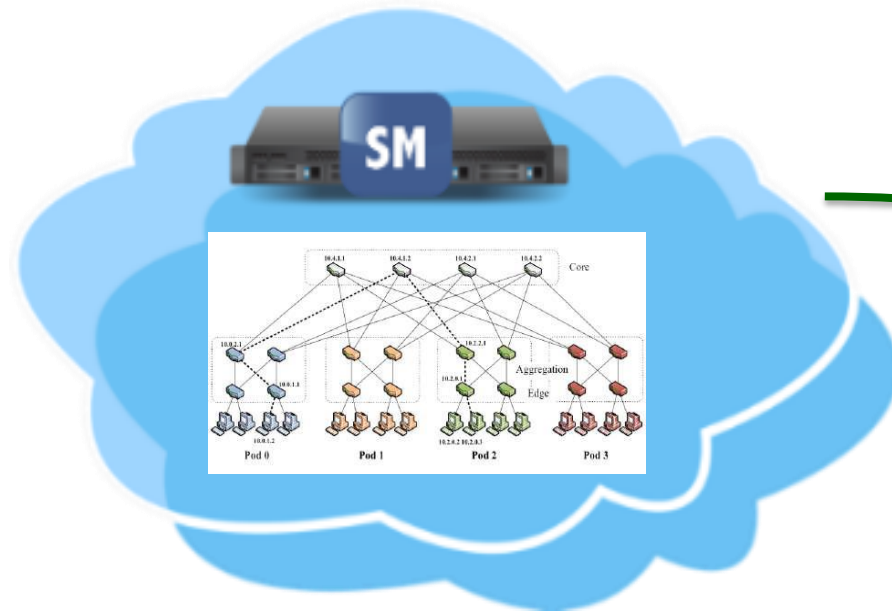
**Accelerating & Virtualizing VNFs**  
**Open vSwitch (OVS), SDN**  
**Overlay networking offloads**

- Integration of ConnectX5 + Multicore ARM
- State of the art capabilities
  - 10 / 25 / 40 / 50 / 100G Ethernet & InfiniBand
  - PCIe Gen3/Gen4
  - Hardware acceleration offload
    - RDMA, RoCE, NVMeF, RAID
- Family of products
  - Range of ARM core counts and I/O ports/speeds
  - Price/Performance points



**Isolation Between Different InfiniBand Networks (Each Network can be Managed Separately)**

**Native InfiniBand Connectivity Between Different Network Topologies (Fat-Tree, Torus, Dragonfly, etc.)**

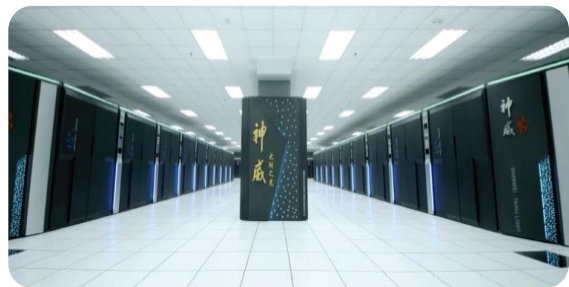


**SB7780 Router 1U  
Supports up to 6 Different Subnets**



## Petascale


#1 TOP500, 100Petaflop



2016



OAK RIDGE  
National Laboratory  
"Summit" System



Lawrence Livermore  
National Laboratory  
"Sierra" System

2019

## Exascale



World-wide  
Programs

2021

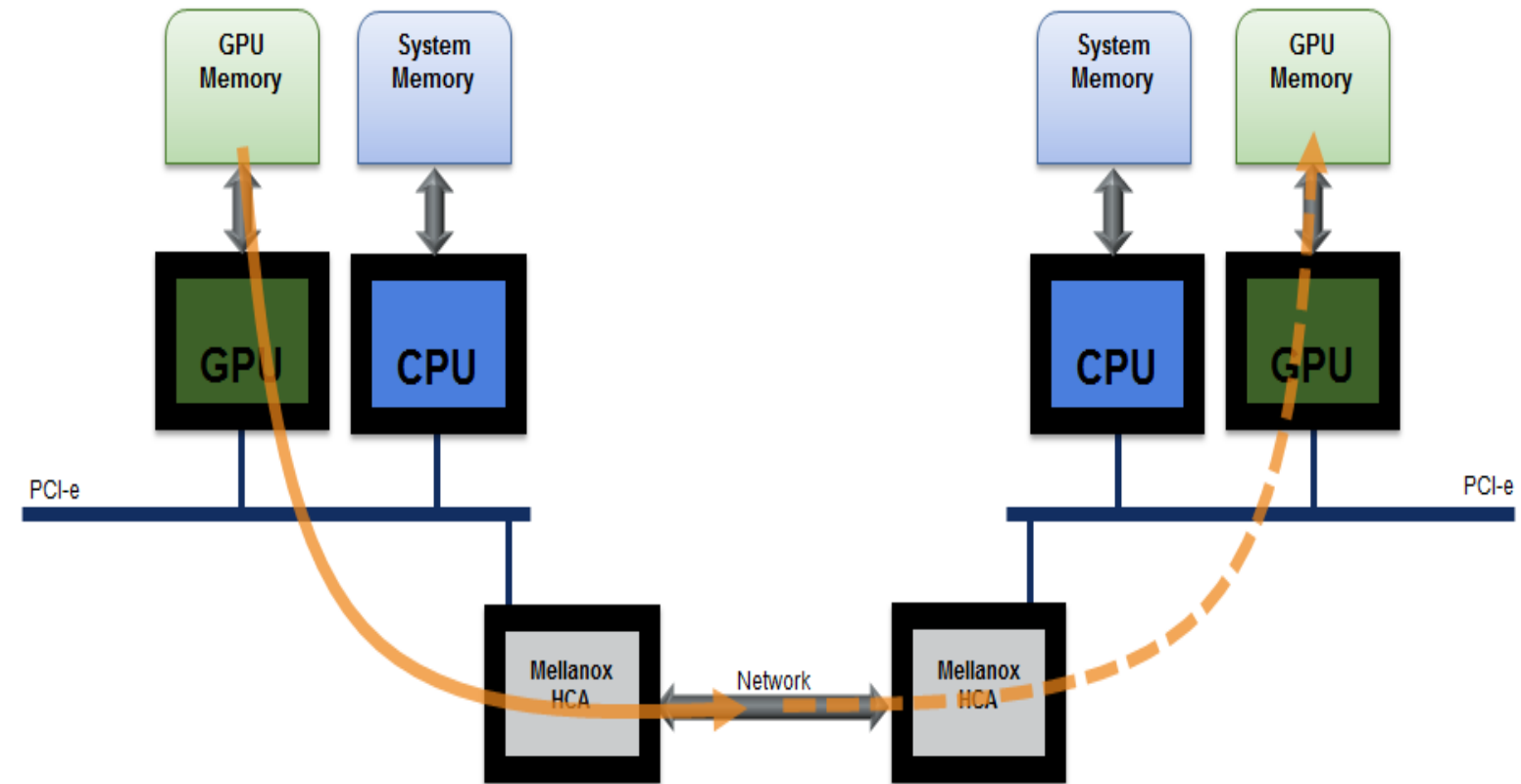
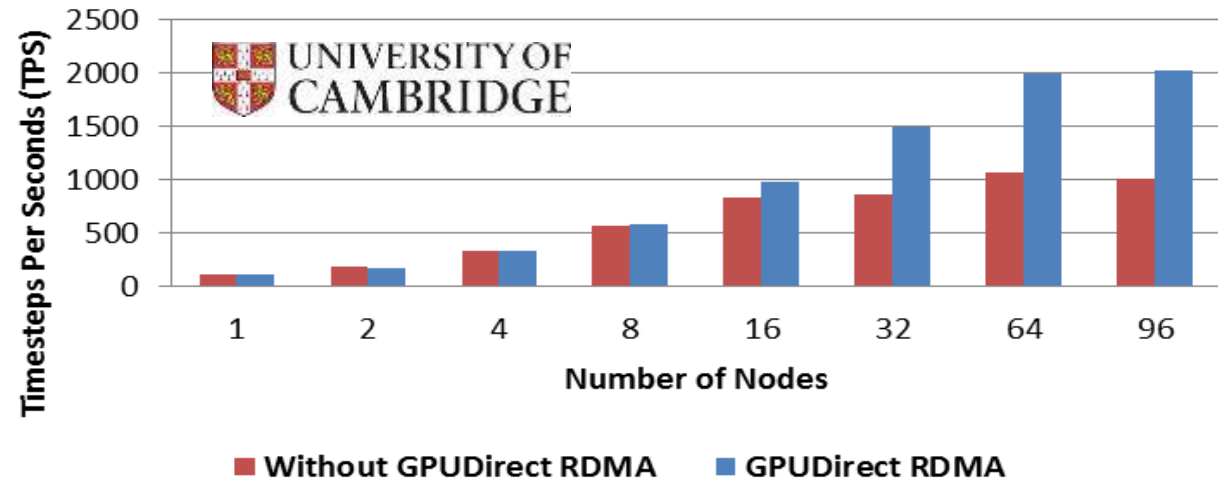
# GPUDirect RDMA Technology

Maximize Performance via Accelerator and GPU Offloads

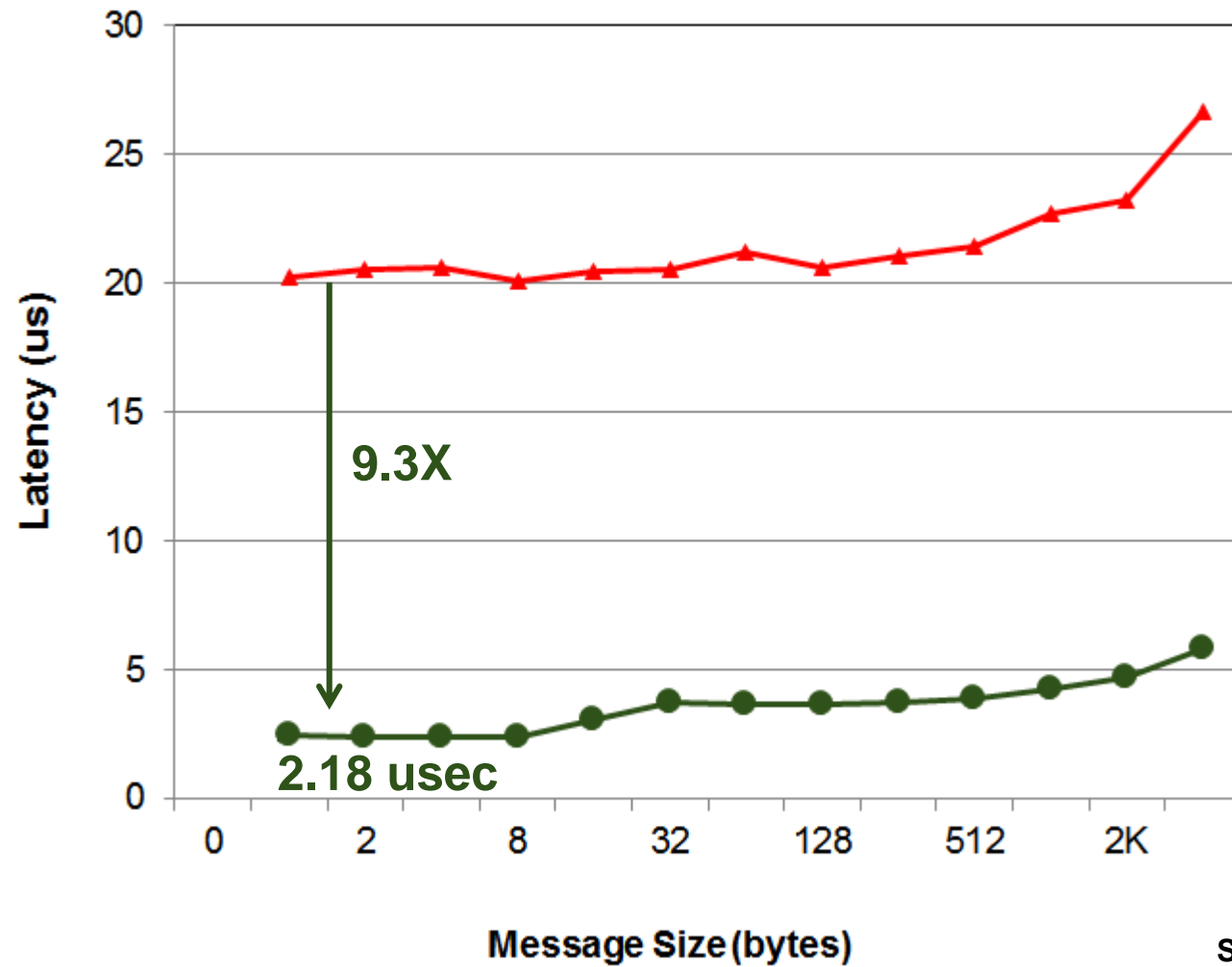


# GPUs are Everywhere!

## HOOMD-blue Performance (LJ Liquid Benchmark, 512K Particles)



## GPU-GPU Internode MPI Latency

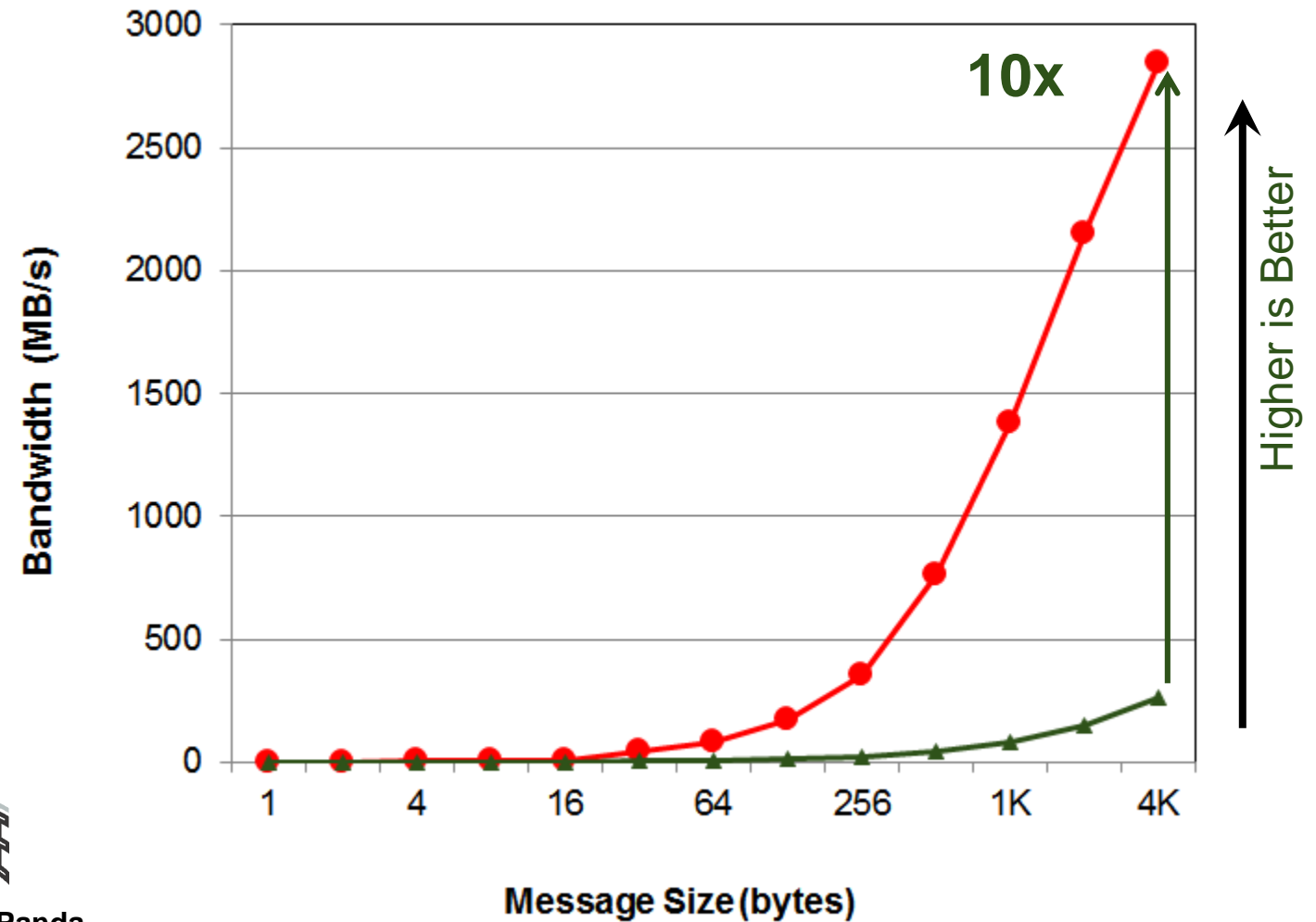


Lower is Better



Source: Prof. DK Panda

## GPU-GPU Internode MPI Bandwidth



Higher is Better

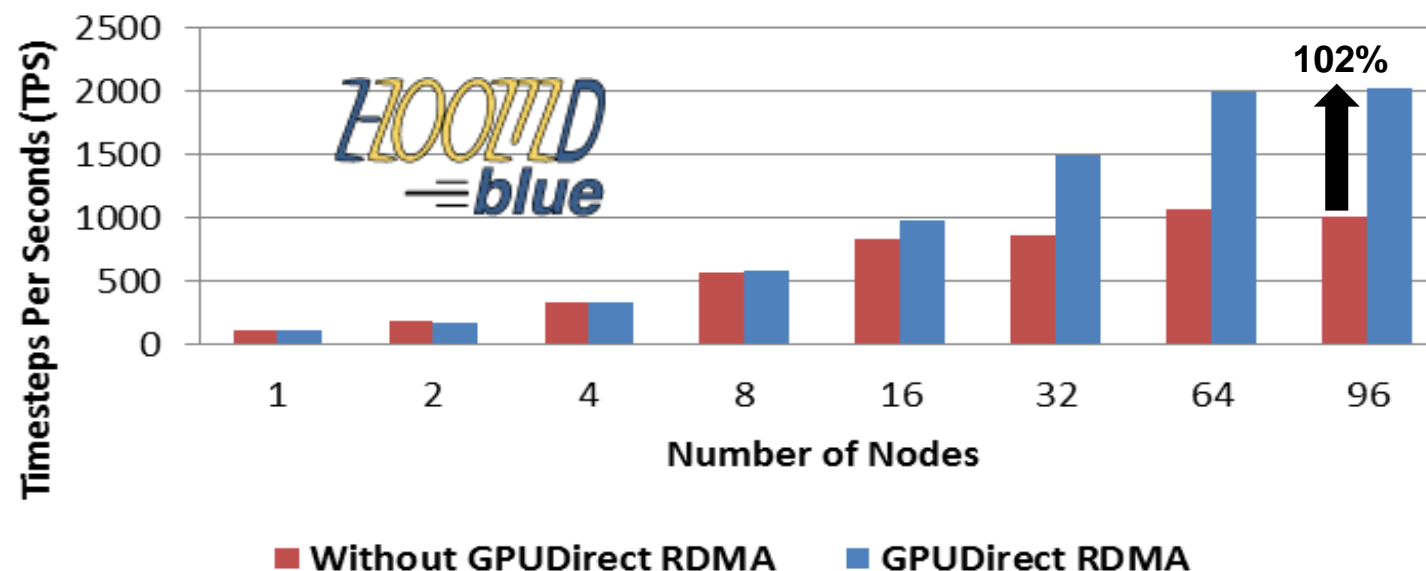
88% Lower Latency

10X Increase in Throughput

- HOOMD-blue is a general-purpose Molecular Dynamics simulation code accelerated on GPUs
- GPUDirect RDMA allows direct peer to peer GPU communications over InfiniBand
  - Unlocks performance between GPU and InfiniBand
  - This provides a significant decrease in GPU-GPU communication latency
  - Provides complete CPU offload from all GPU communications across the network



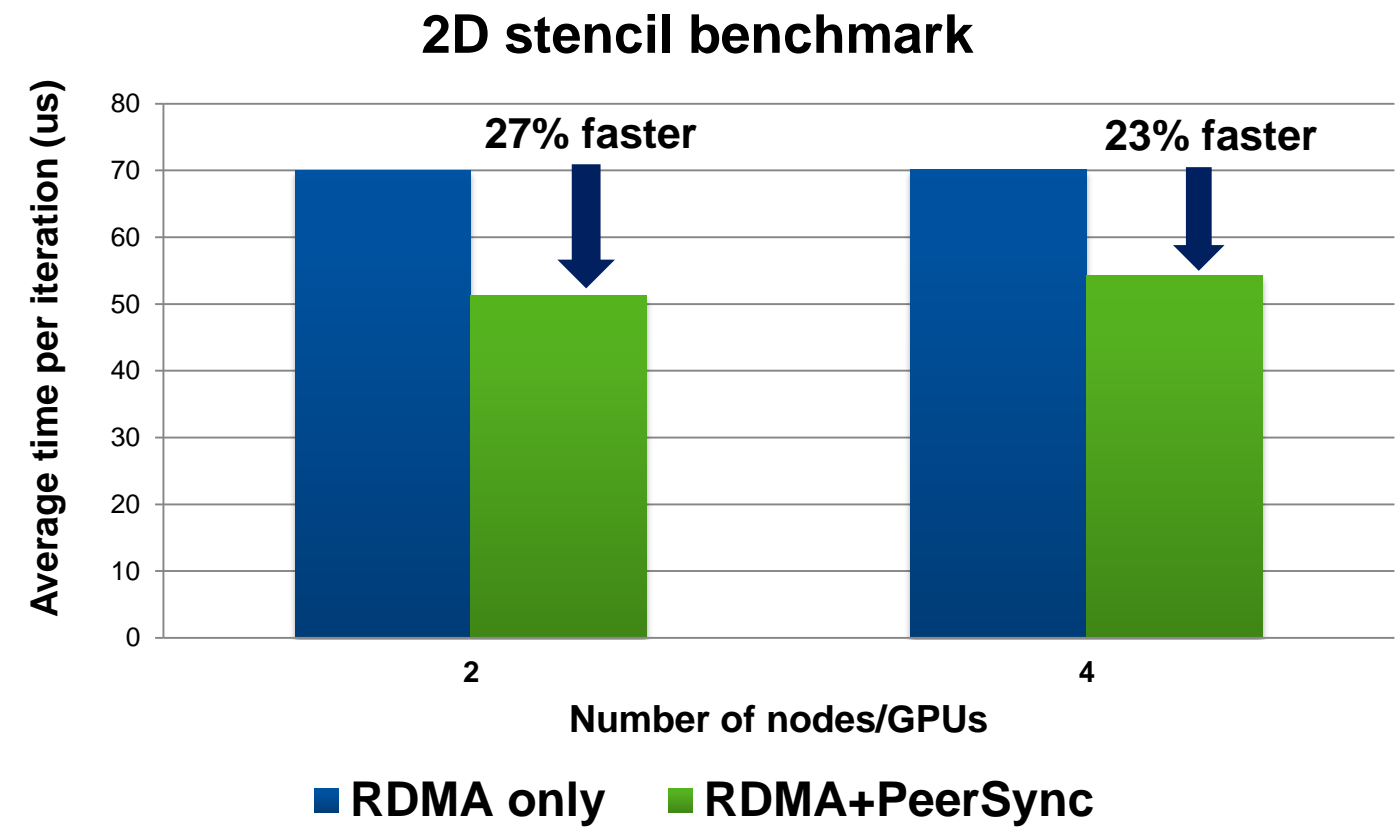
## HOOMD-blue Performance (LJ Liquid Benchmark, 512K Particles)



**2X Application  
Performance!**

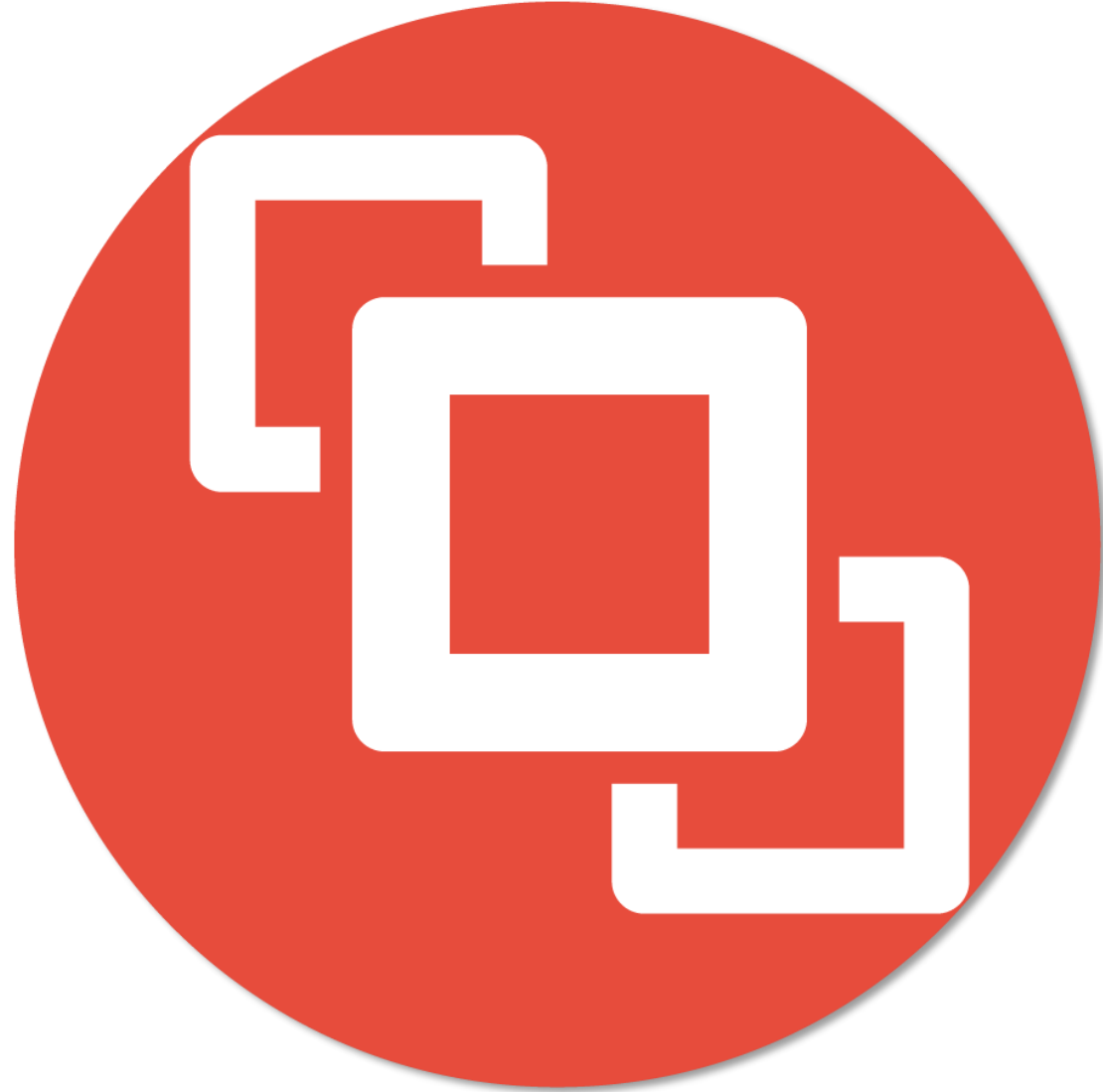
- GPUDirect RDMA (3.0) – direct data path between the GPU and Mellanox interconnect
  - Control path still uses the CPU
    - CPU prepares and queues communication tasks on GPU
    - GPU triggers communication on HCA
    - Mellanox HCA directly accesses GPU memory
- GPUDirect Sync (GPUDirect 4.0)
  - Both data path and control path go directly between the GPU and the Mellanox interconnect

**Maximum Performance  
For GPU Clusters**



# UCX

Unified Communication X



# UCX

Unified Communication - X  
Framework

[www.openucx.org](http://www.openucx.org)  
[ucx-group@email.ornl.gov](mailto:ucx-group@email.ornl.gov)





# Background

## MXM

- Developed by Mellanox Technologies
- HPC communication library for InfiniBand devices and shared memory
- Primary focus: MPI, PGAS

## UCCS

- Developed by ORNL, UH, UTK
- Originally based on Open MPI BTL and OPAL layers
- HPC communication library for InfiniBand, Cray Gemini/Aries, and shared memory
- Primary focus: OpenSHMEM, PGAS
- Also supports: MPI

## PAMI

- Developed by IBM on BG/Q, PERCS, IB VERBS
- Network devices and shared memory
- MPI, OpenSHMEM, PGAS, CHARM++, X10
- C++ components
- Aggressive multi-threading with contexts
- Active Messages
- Non-blocking collectives with hw acceleration support



# High-level Overview

## Applications

MPICH, Open-MPI, etc.

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

Parsec, OCR, Legions, etc.

Burst buffer, ADIOS, etc.

### UCP (Protocols)

Message Passing  
API Domain

PGAS  
API Domain

Task Based  
API Domain

I/O  
API Domain

### UCT (Transports)

InfiniBand VERBs

RC

UD

XRC

DCT

Gemini/Aries

GNI

Host Memory

SYSV

POSIX

KNEM

CMA

XPMEM

Accelerator Memory

CUDA

### UCS (Services)

Utilities

Platforms

UCX

Hardware/Driver





## Petascale

#1 TOP500, 100Petaflop



## Exascale



**OAK RIDGE**  
National Laboratory  
"Summit" System

**Lawrence Livermore**  
National Laboratory  
"Sierra" System



World-wide Programs

2016

2019

2021



Thank You