



Designing Towards Exascale Compute and Extreme Performance in the Cloud

BEAR Cloud Launch

October 21, 2016



Mellanox Connects the World's Fastest Supercomputer



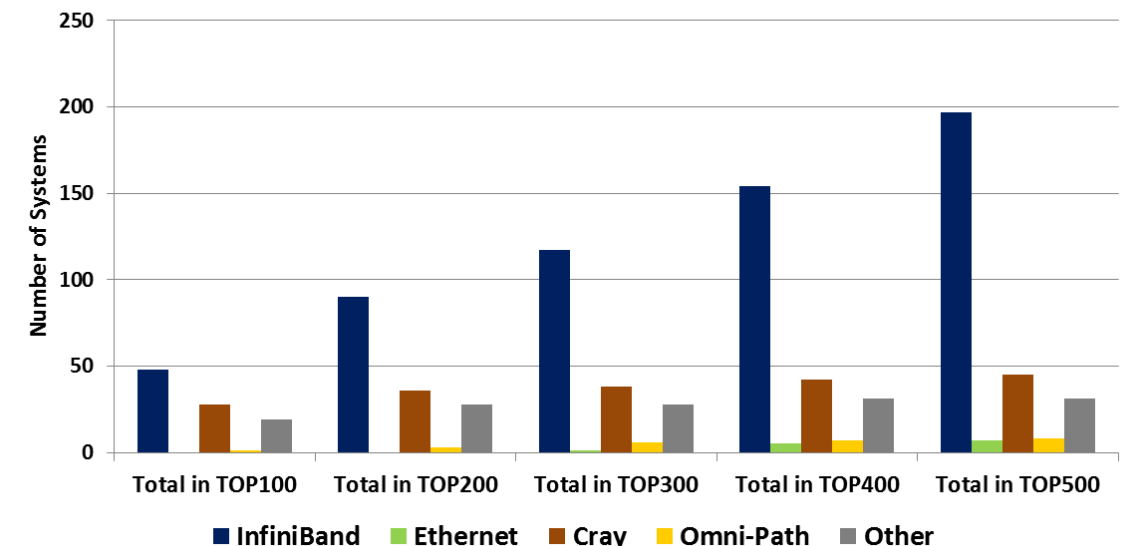
National Supercomputing Center in Wuxi, China #1 on the TOP500 Supercomputing List

- 93 Petaflop performance, 3X higher versus #2 on the TOP500
- Nearly 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

Technology Roadmap – One-Generation Lead over the Competition



Mellanox → 20G → 40G → 56G → 100G → 200G → 400G

Terascale

3rd



TOP500 2003
Virginia Tech (Apple)

1st



“Roadrunner”
Mellanox Connected

Petascale



Exascale

OAK RIDGE
National Laboratory
“Summit” System

Lawrence Livermore
National Laboratory
“Sierra” System

2000

2005

2010

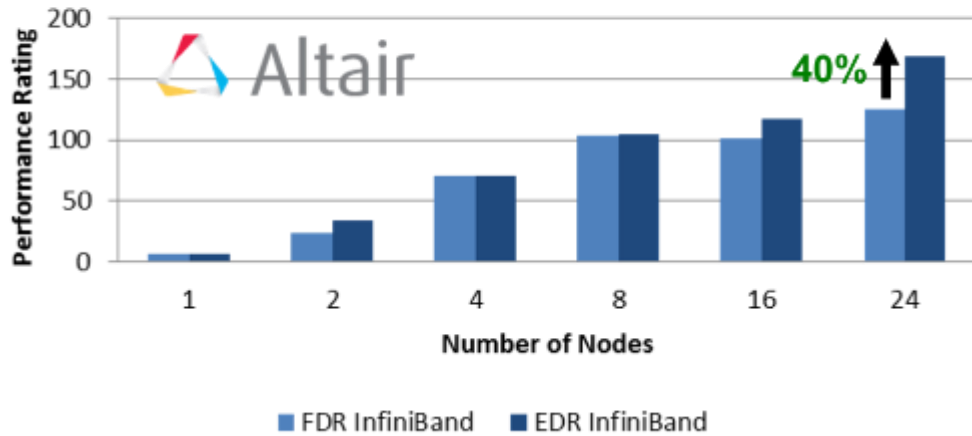
2015

2020

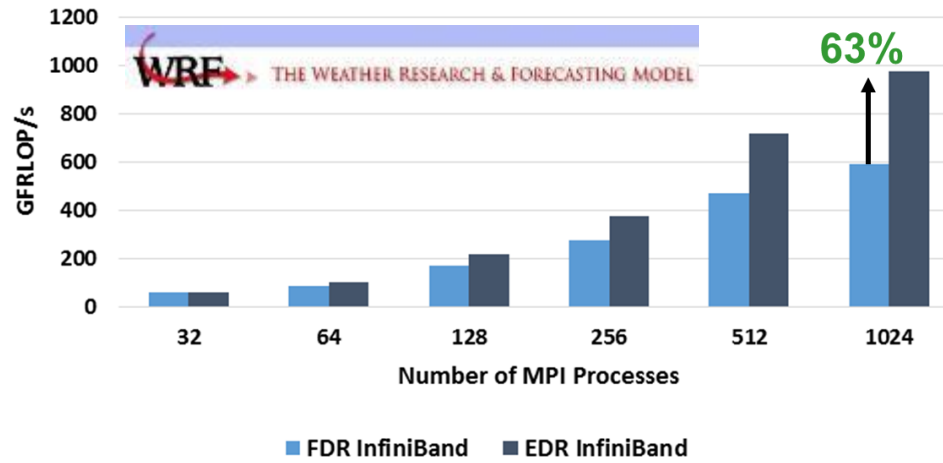
The Performance Advantage of EDR 100G InfiniBand (28-80%)



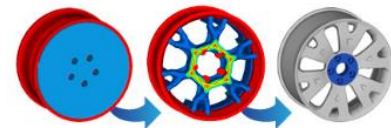
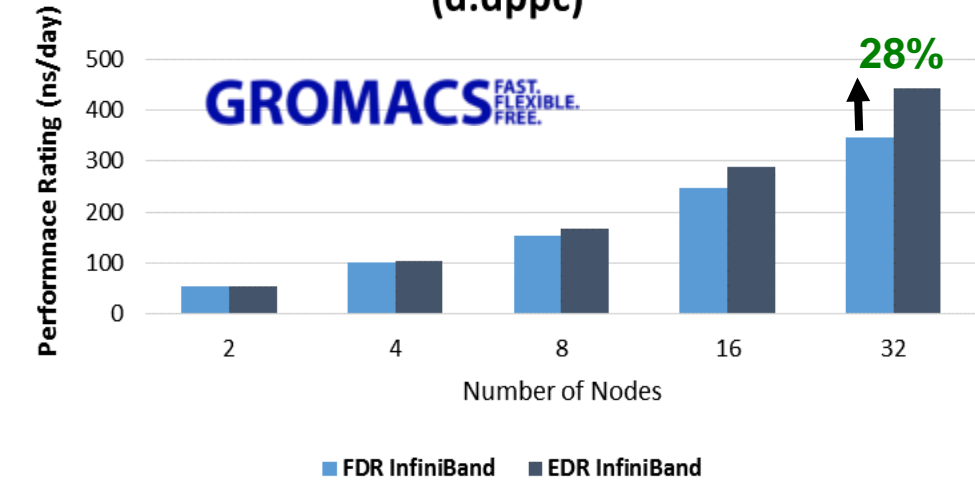
OptiStruct Performance (Engine_Assy.fem)



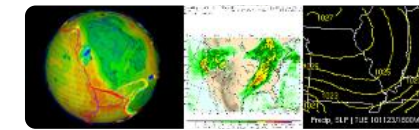
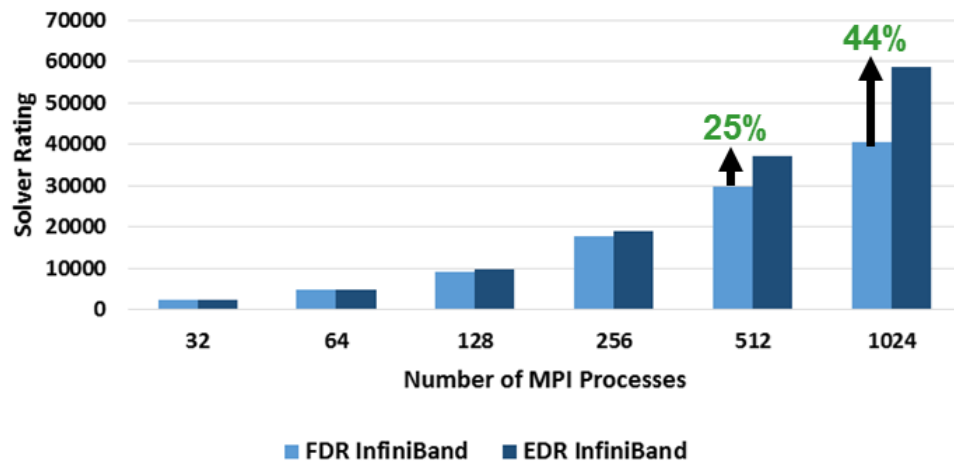
WRF Performance (conus12km)



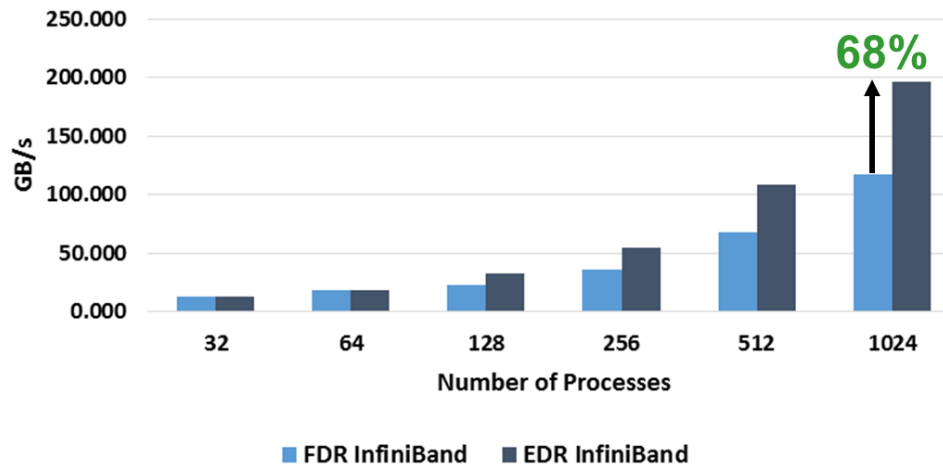
GROMACS Performance (d.dppc)



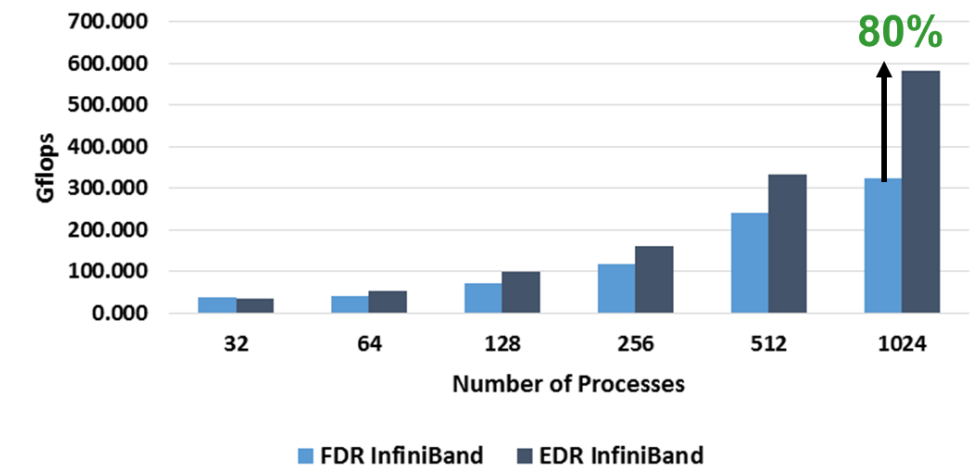
ANSYS Fluent 16.0 Performance (sedan_4m)



HPCC Performance (PTRANS_GB)



HPCC Performance (MPIFFT)



Performance Development

Terascale



Petascale

1st



“Roadrunner”



Exascale



2000

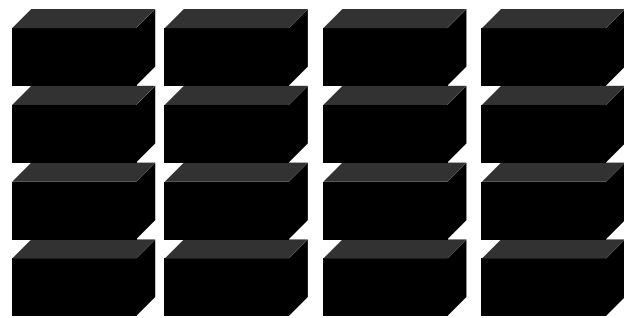
2005

2010

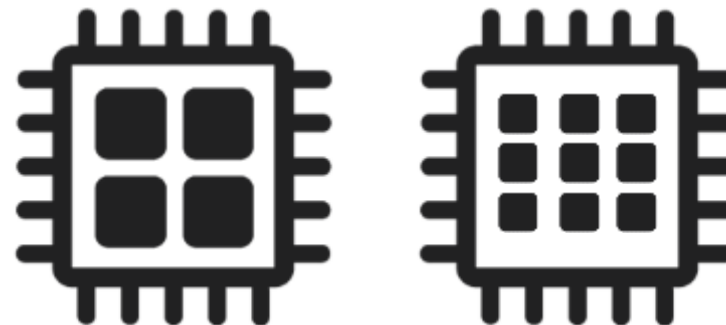
2015

2020

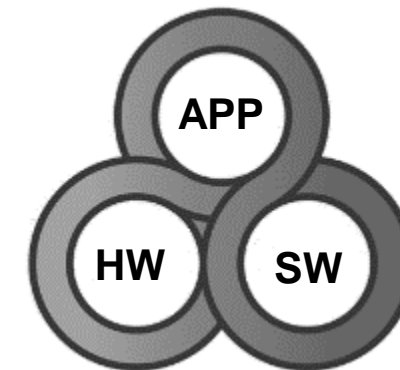
The Interconnect is the Enabling Technology



SMP to Clusters



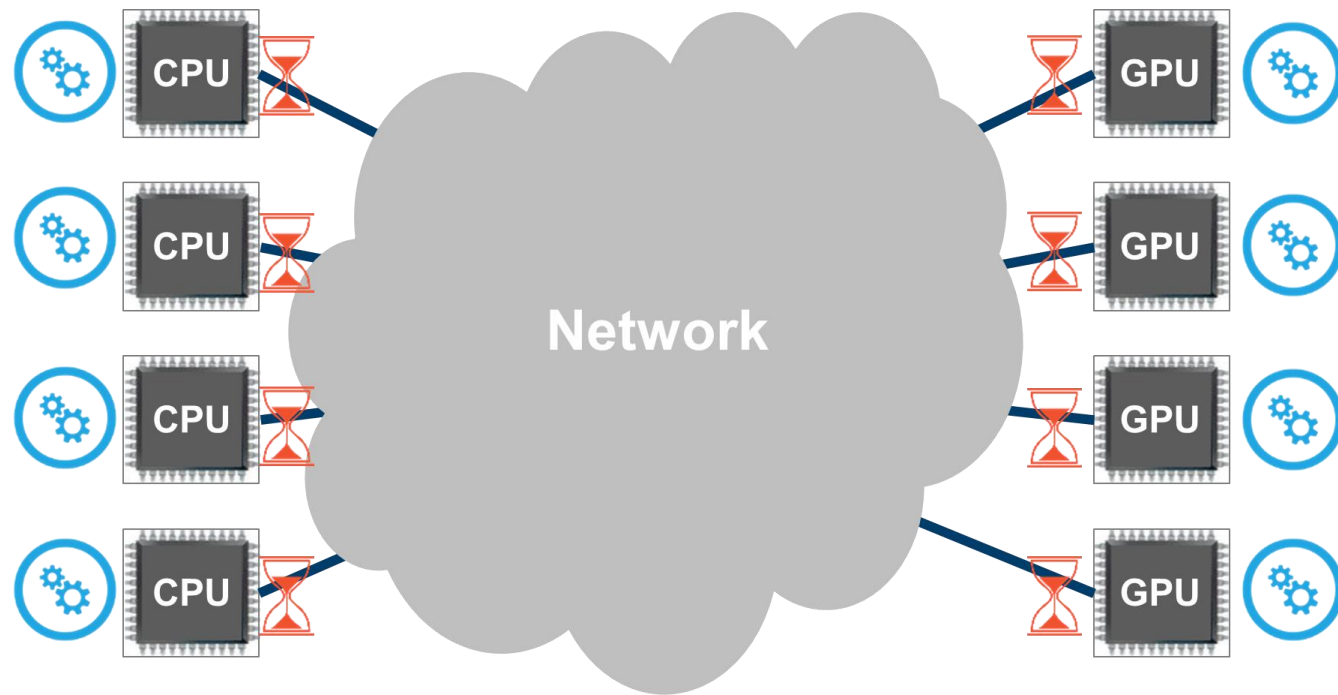
Single-Core to Many-Core



Application
Software
Hardware

Co-Design

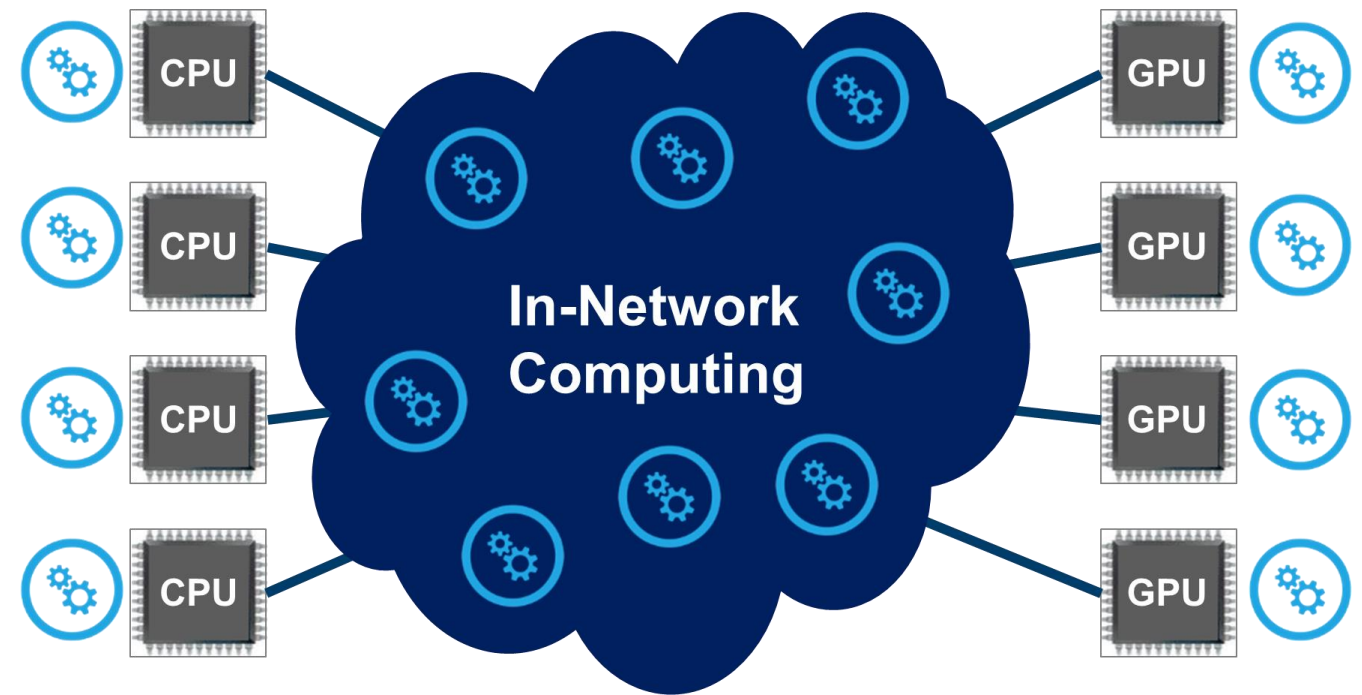
CPU-Centric



Limited to Main CPU Usage
Results in Performance Limitation

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

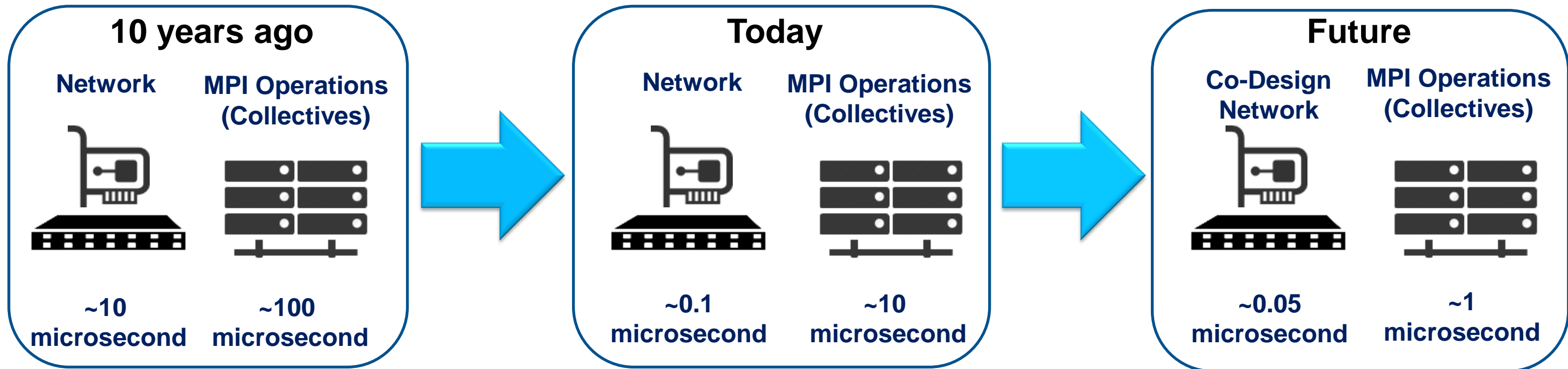
Co-Design



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

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/PGAS Communications

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

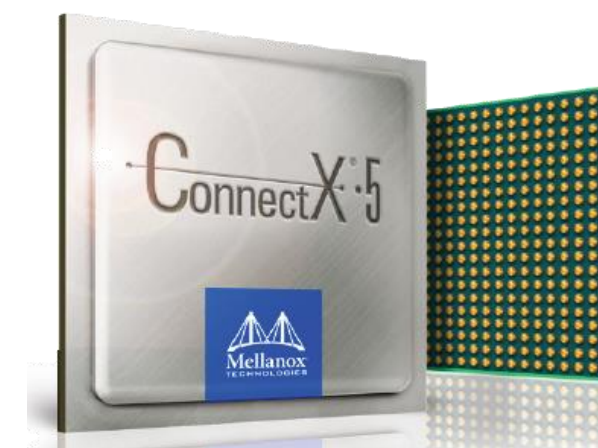
MPI Collectives in Hardware
MPI Tag Matching in Hardware
In-Network Memory

PCIe Gen3 and Gen4
Integrated PCIe Switch
Advanced Dynamic Routing

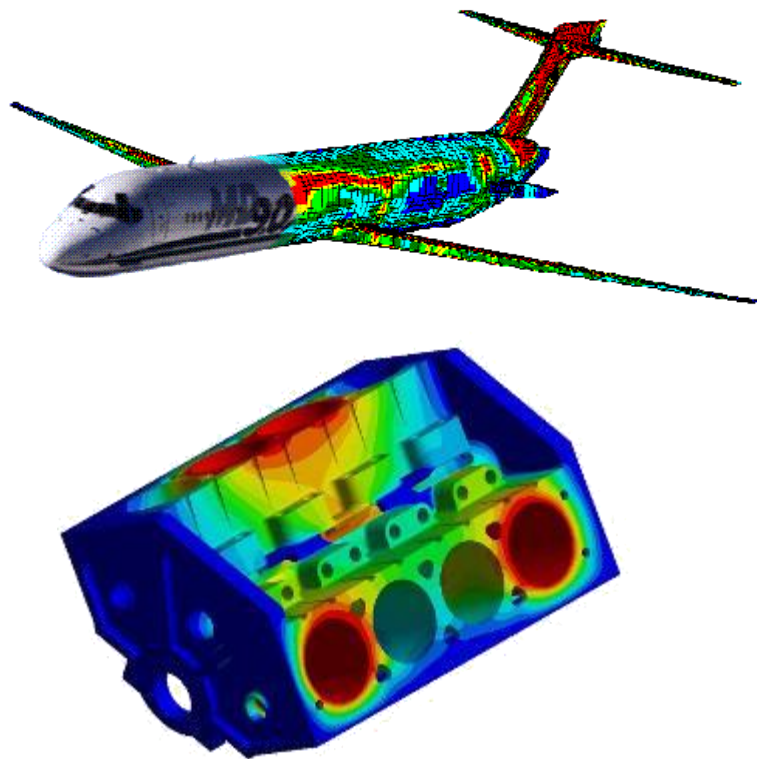
Switch-IB™ 2 SHArP



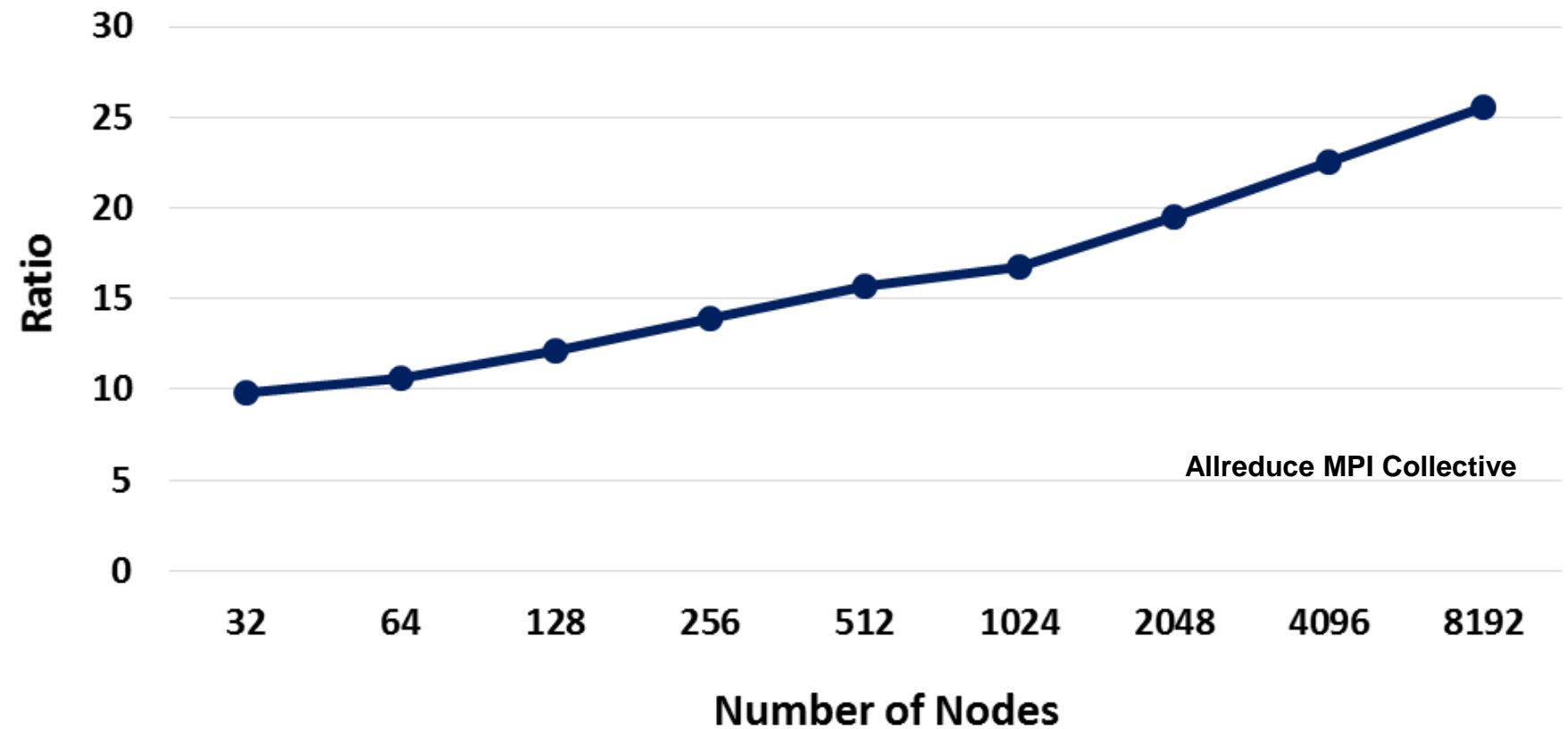
ConnectX® 5



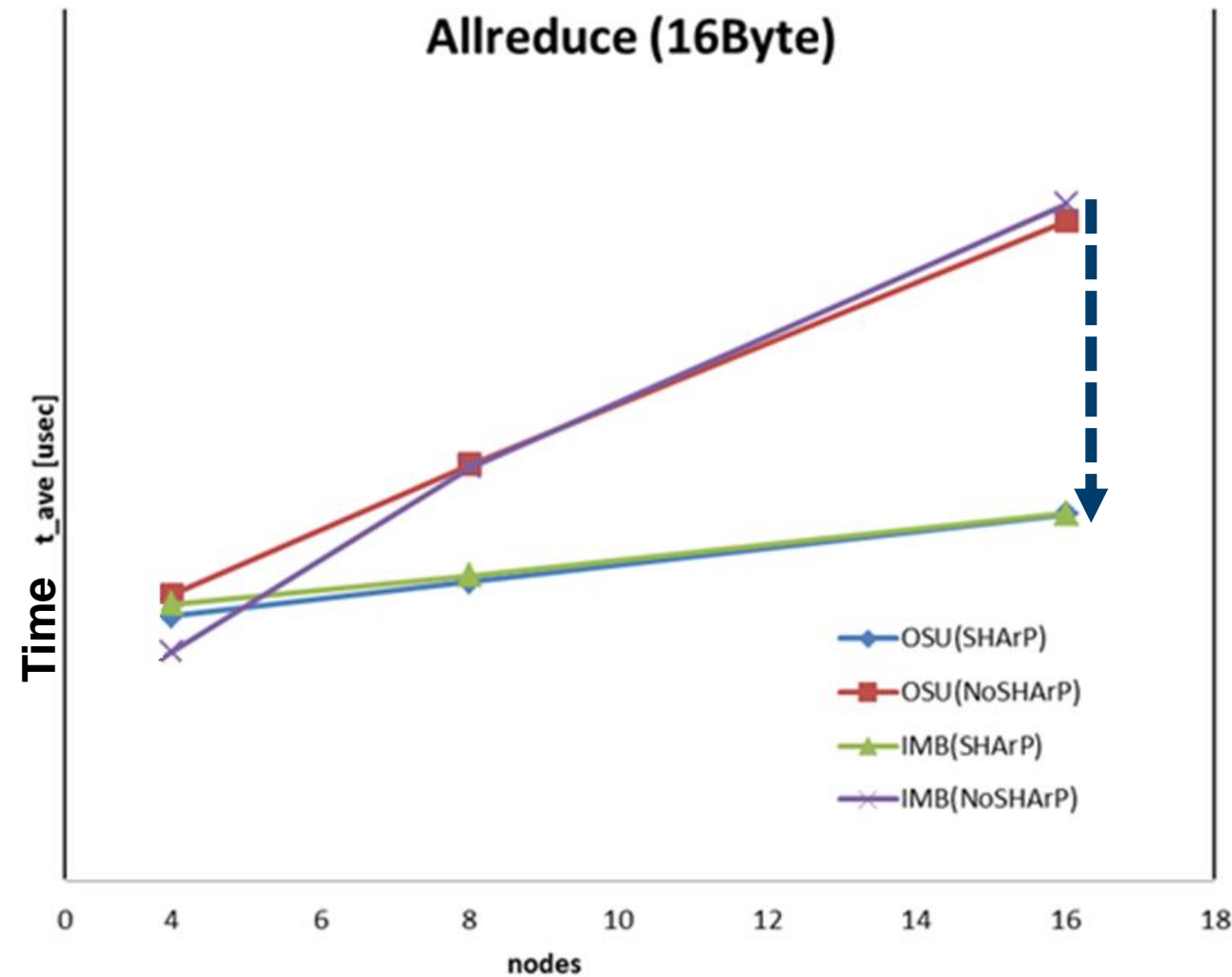
- MiniFE is a Finite Element mini-application
 - Implements kernels that represent implicit finite-element applications



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

Lower is better

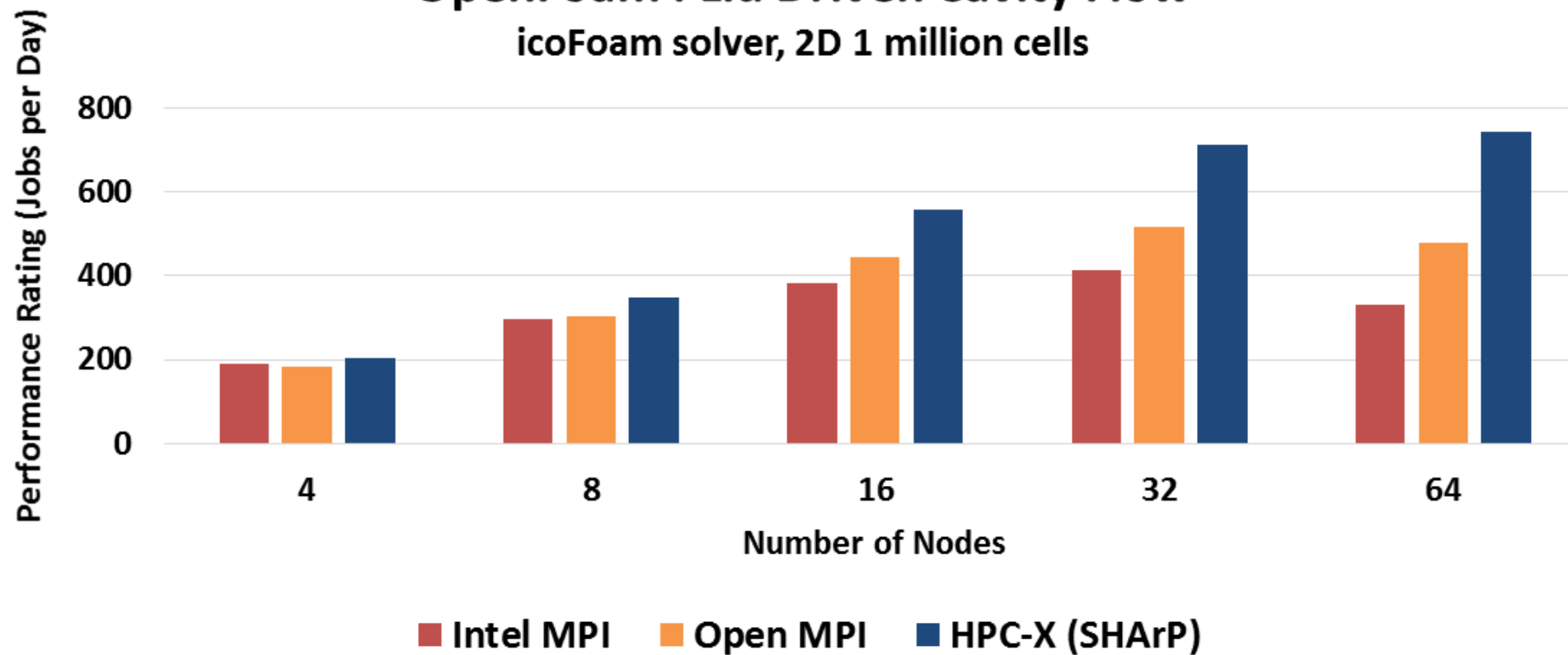
**Maximizing KNL Performance – 50% Reduction in Run Time
(Customer Results)**

OpenFOAM

OpenFOAM is a popular computational fluid dynamics application

HPC-X™

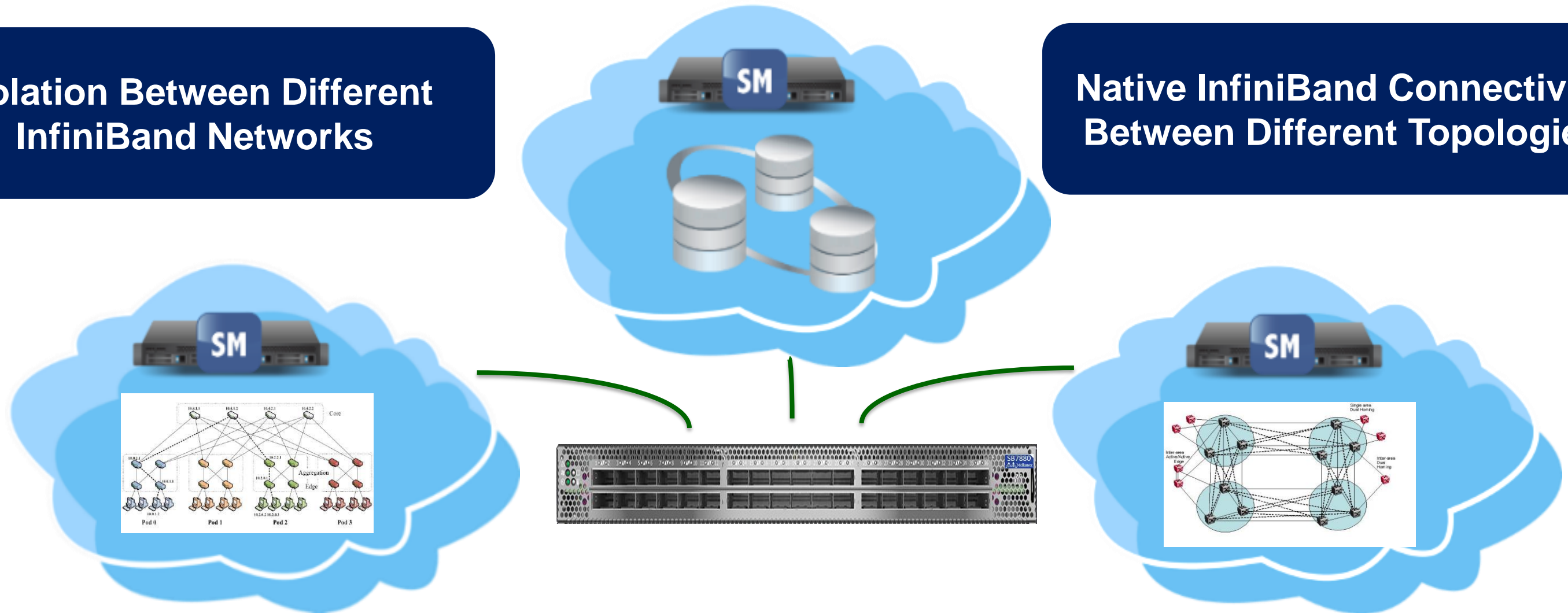
OpenFoam : Lid Driven Cavity Flow
icoFoam solver, 2D 1 million cells



HPC-X with SHArP Delivers 2.2X Higher Performance over Intel MPI

**Isolation Between Different
InfiniBand Networks**

**Native InfiniBand Connectivity
Between Different Topologies**

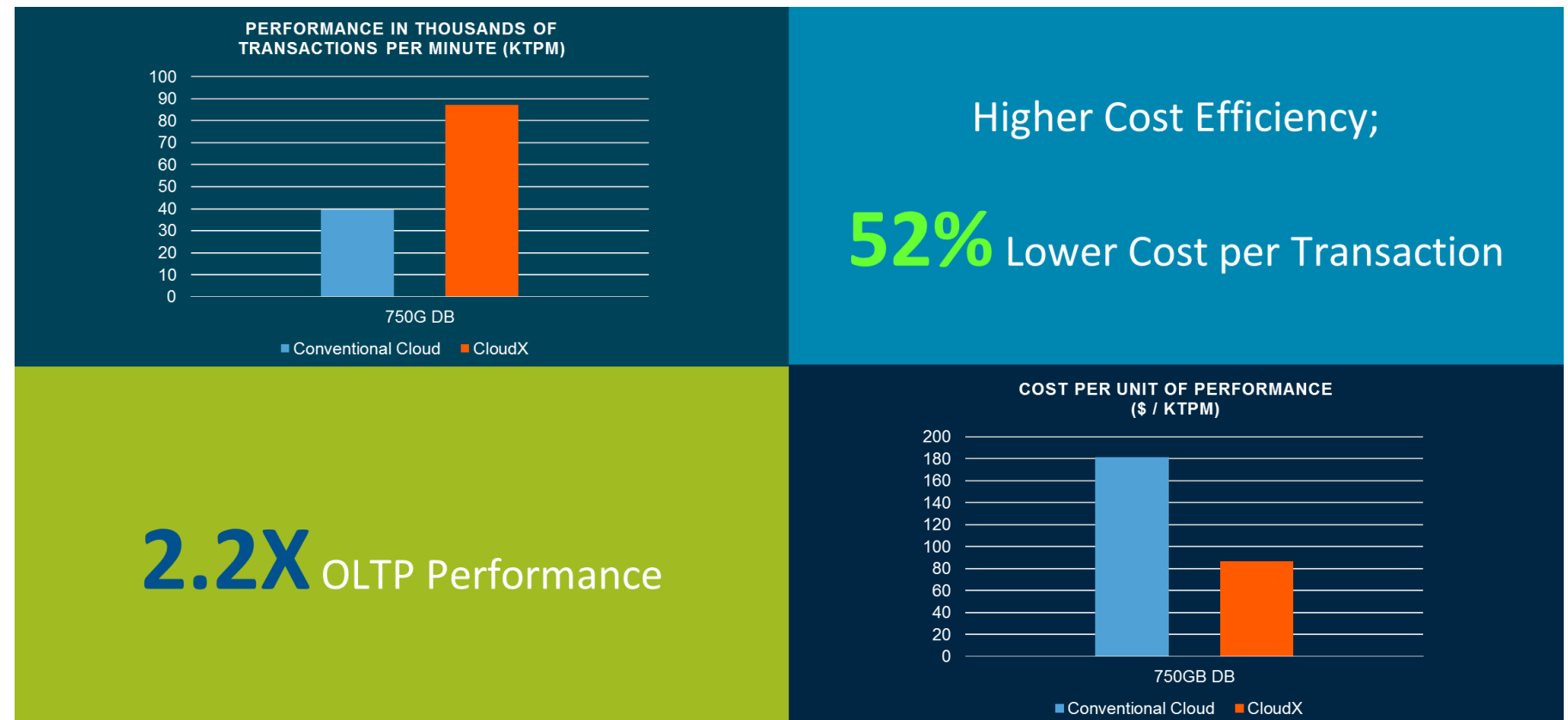
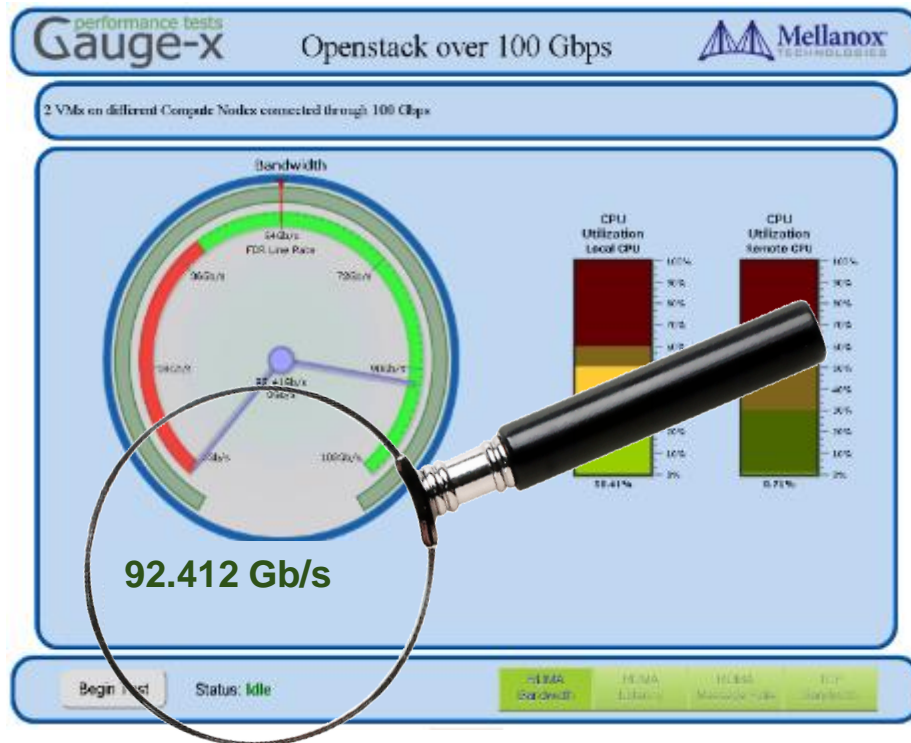


**SB7780 Router 1U
Supports up to 6 Different Subnets**

OpenStack Over InfiniBand – The 100G Cloud is Here!



- Transparent InfiniBand integration into OpenStack
- InfiniBand SDN network
- RDMA directly from VM



The Ideal Fit for HPC Clouds !

InfiniBand The Smart Choice for HPC Platforms and Applications



- *“We chose a co-design approach. This system was of course targeted at supporting in the best possible manner our key applications. The only interconnect that really could deliver that was Mellanox InfiniBand.”*



[Watch Video](#)

- *“One of the big reasons we use InfiniBand and not an alternative is that we’ve got backwards compatibility with our existing solutions.”*



UNIVERSITY OF BIRMINGHAM

[Watch Video](#)

- *“InfiniBand is the most advanced high performance interconnect technology in the world, with dramatic communication overhead reduction that fully unleashes cluster performance.”*



[Watch Video](#)

- *“InfiniBand is the best that is required for our applications. It enhancing and unlocking the potential of the system.”*



[Watch Video](#)

Scalable Performance

At the Speed of 100Gb/s!

Leading Supplier of End-to-End Interconnect Solutions



Comprehensive End-to-End InfiniBand and Ethernet Portfolio (VPI)

ICs	Adapter Cards	NPU & Multicore	Switches/Gateways	Software	Metro / WAN	Cables/Modules