

Advancing Applications Performance With InfiniBand

Pak Lui, Application Performance Manager

September 12, 2013



- **Leading provider of high-throughput, low-latency server and storage interconnect**
 - FDR 56Gb/s InfiniBand and 10/40/56GbE
 - Reduces application wait-time for data
 - Dramatically increases ROI on data center infrastructure
- **Company headquarters:**
 - Yokneam, Israel; Sunnyvale, California
 - ~1,200 employees* worldwide
- **Solid financial position**
 - Record revenue in FY12; \$500.8M, up 93% year-over-year
 - Q2'13 revenue of \$98.2M
 - Q3'13 guidance ~\$104M to \$109M
 - Cash + investments @ 6/30/13 = \$411.3M



* As of June 2013

Comprehensive End-to-End Software Accelerators and Management

<div data-bbox="126 453 403 656"> <p>MPI</p> </div> <div data-bbox="418 453 696 656"> <p>SHMEM</p> </div> <div data-bbox="711 453 988 656"> <p>PGAS</p> </div> <div data-bbox="126 673 548 887"> <p>MXM Mellanox Messaging Acceleration</p> </div> <div data-bbox="563 673 988 887"> <p>FCA Fabric Collectives Acceleration</p> </div>	<div data-bbox="1111 453 1965 656"> <p>Management</p> </div> <div data-bbox="1111 673 1965 887"> <p>UFM Unified Fabric Management</p> </div>	<div data-bbox="2097 453 2950 656"> <p>Storage and Data</p> </div> <div data-bbox="2097 673 2519 887"> <p>VSA Storage Accelerator (iSCSI)</p> </div> <div data-bbox="2534 673 2981 887"> <p>UDA Unstructured Data Accelerator</p> </div>
---	--	---

Comprehensive End-to-End InfiniBand and Ethernet Solutions Portfolio

<p>ICs</p>	<p>Adapter Cards</p>	<p>Switches/Gateways</p>	<p>Long-Haul Systems</p> <p>metroX™</p>	<p>Cables/Modules</p>
-------------------	-----------------------------	---------------------------------	---	------------------------------

ConnectX[®] 3 VPI Adapter



Applications

Networking

Storage

Clustering

Management

Acceleration Engines

PCI EXPRESS[™] 3.0



Ethernet: 10/40/56 Gb/s

InfiniBand: 10/20/40/56 Gb/s



LOM



Adapter Card



Mezzanine Card



SwitchX[®] 2 VPI Switch

Unified Fabric Manager

Switch OS Layer

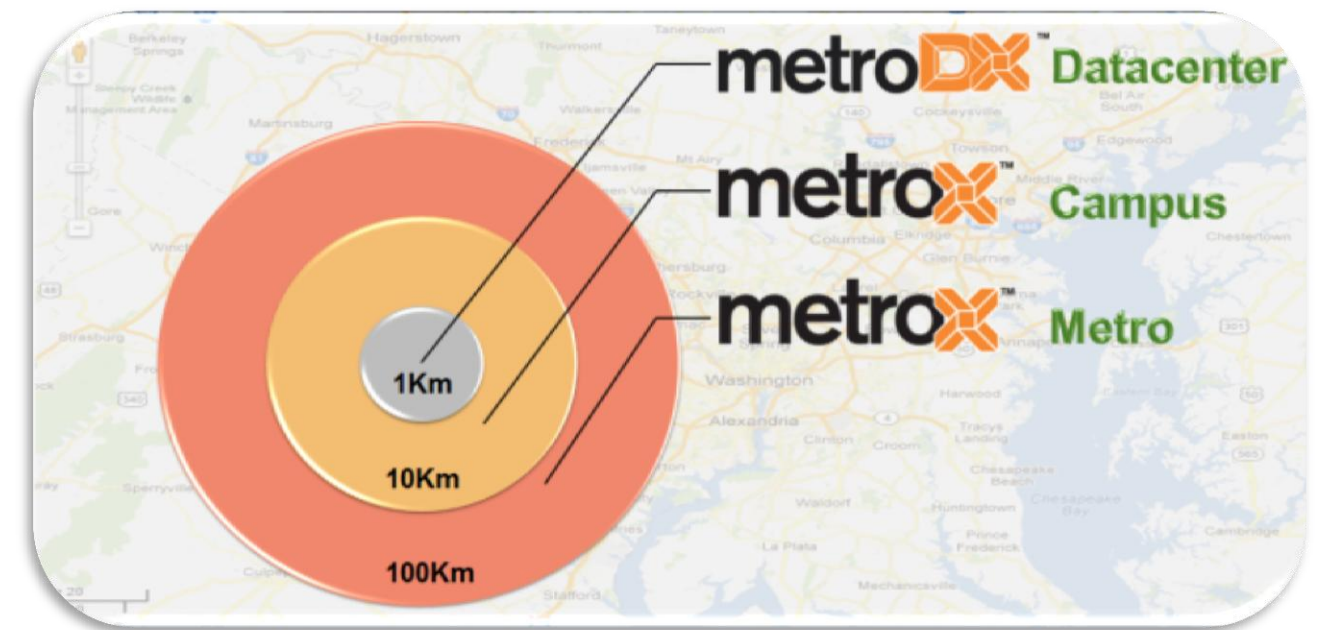


64 ports 10GbE
 36 ports 40/56GbE
 48 10GbE + 12 40/56GbE
 36 ports IB up to 56Gb/s
 8 VPI subnets



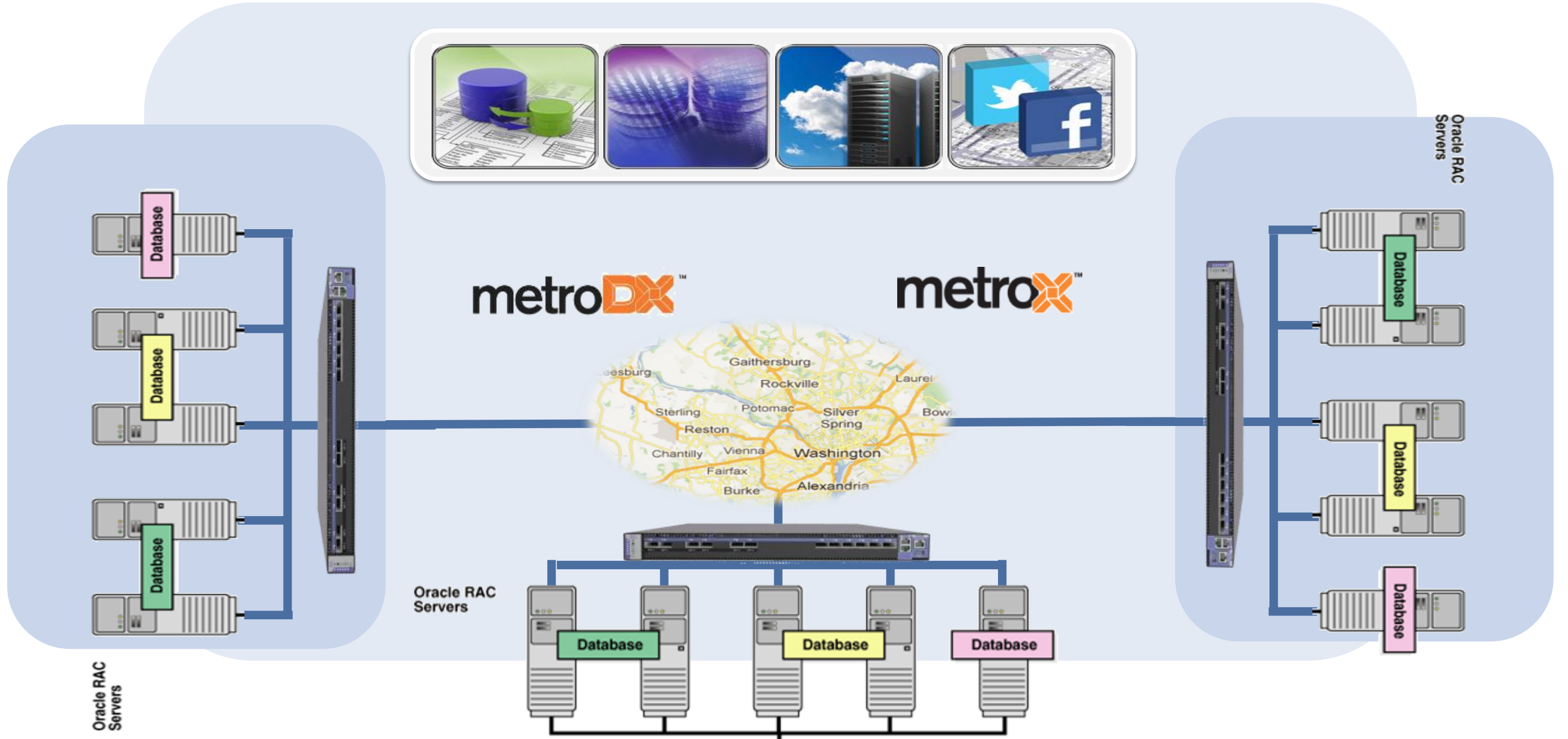
From data center to
 campus and metro
 connectivity

- MetroX™ and MetroDX™ extends InfiniBand and Ethernet RDMA reach
- Fastest interconnect over 40Gb/s InfiniBand or Ethernet links
- Supporting multiple distances
- Simple management to control distant sites
- Low-cost, low-power , long-haul solution

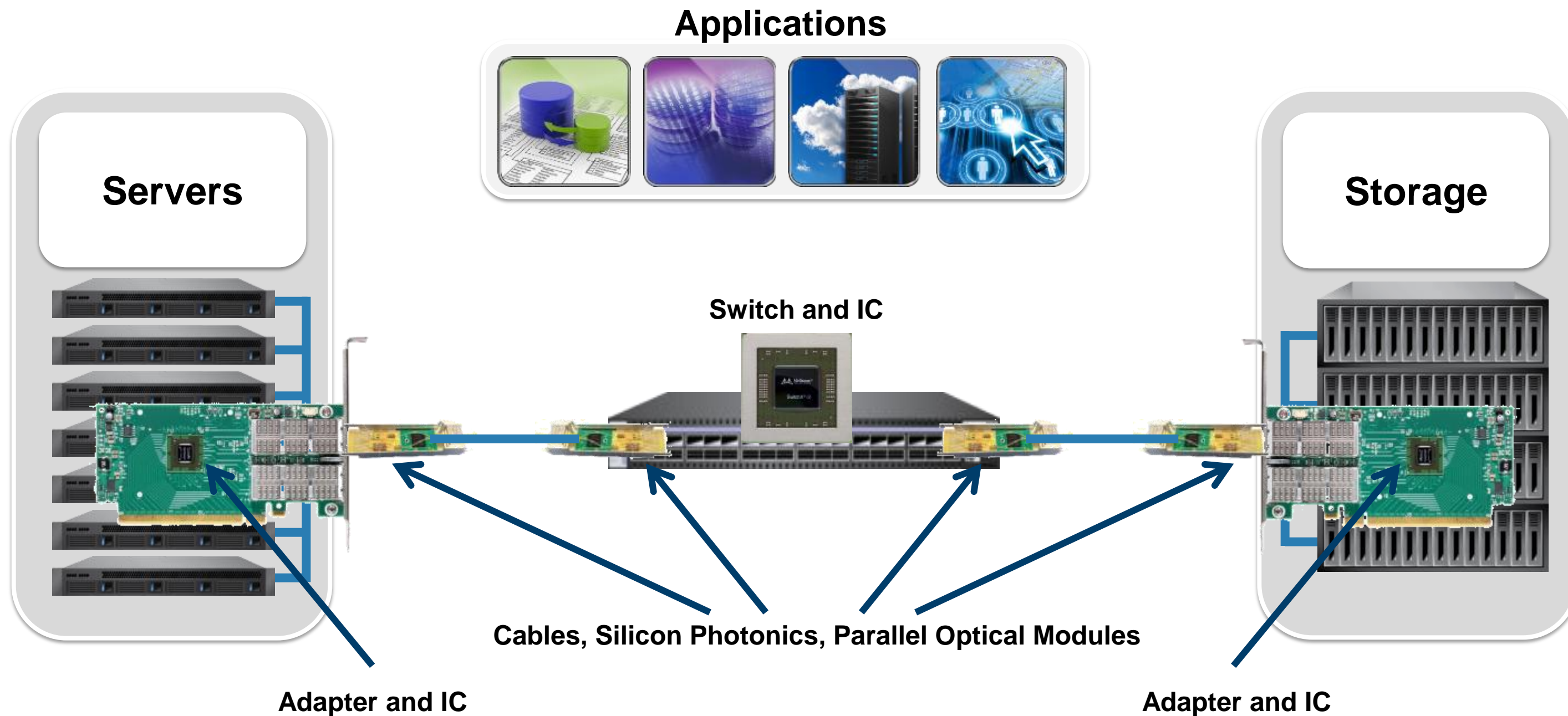


40Gb/s over Campus and Metro

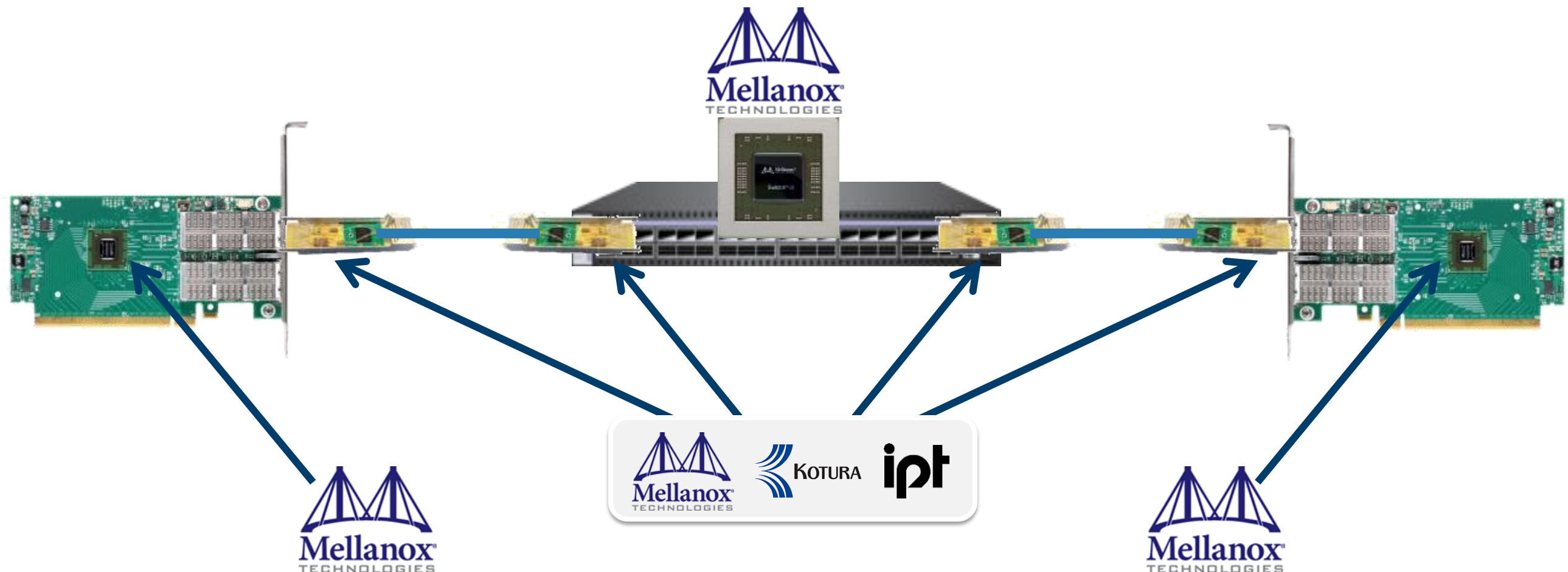
Data Center Expansion Example – Disaster Recovery



Key Elements in a Data Center Interconnect



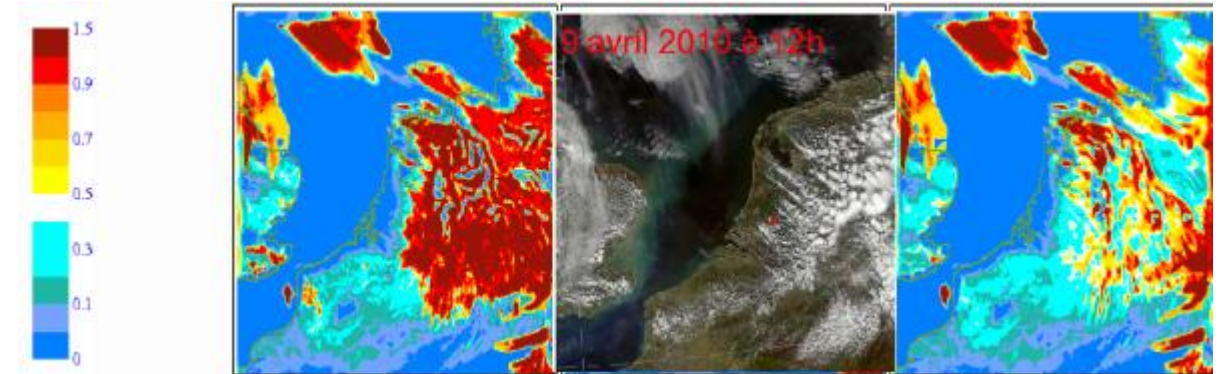
Recent Acquisitions of Kotura and IPtronics Enable Mellanox to Deliver Complete High-Speed Optical Interconnect Solutions for 100Gb/s and Beyond



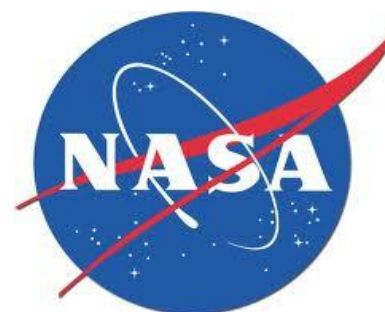
Mellanox InfiniBand Paves the Road to Exascale



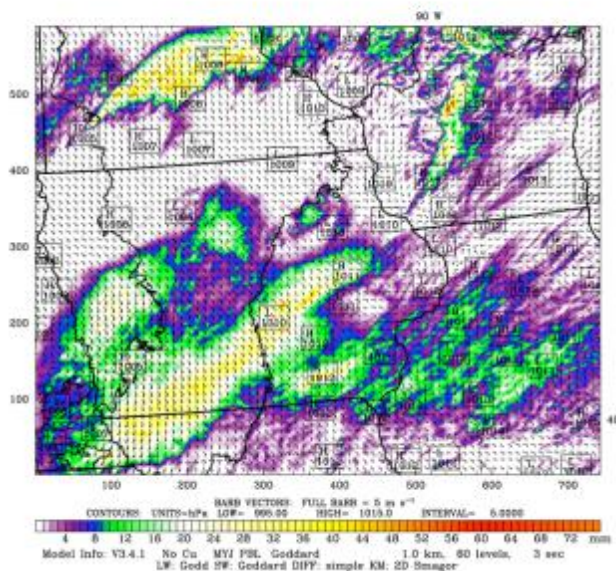
	2013	2014	2015	2016
Research	<p>522 TFlops peak performance 56 racks bullx DLC 1008 nodes Fat Tree InfiniBand FDR Lustre 2 Po, 69 GB/s Disks storage 209 TB</p> <p><i>Computer C1 (03/2013)</i></p>	<p>Centre National de Calcul Météopole, Toulouse</p>	<p>2,85 PFlops peak performance 56+45 racks bullx DLC 1800 nodes Fat Tree InfiniBand FDR Lustre 3,57 Po, 138 GB/s Disks storage 400 TB</p> <p><i>Computer C3 (08/2015)</i></p>	
Operational	<p>522 TFlops peak performance 56 racks bullx DLC 1008 nodes Fat Tree InfiniBand FDR Lustre 1,53 Po, 46 GB/s Disks storage 135 TB</p> <p><i>Computer C2 (11/2013)</i></p>	<p>Espace Clément Ader Montaudran</p>	<p>2,85 PFlops peak performance 56+45 racks bullx DLC 1800 nodes Fat Tree InfiniBand FDR Lustre 2,55 Po, 92 GB/s Disks storage 135 TB</p> <p><i>Computer C4 (04/2016)</i></p>	



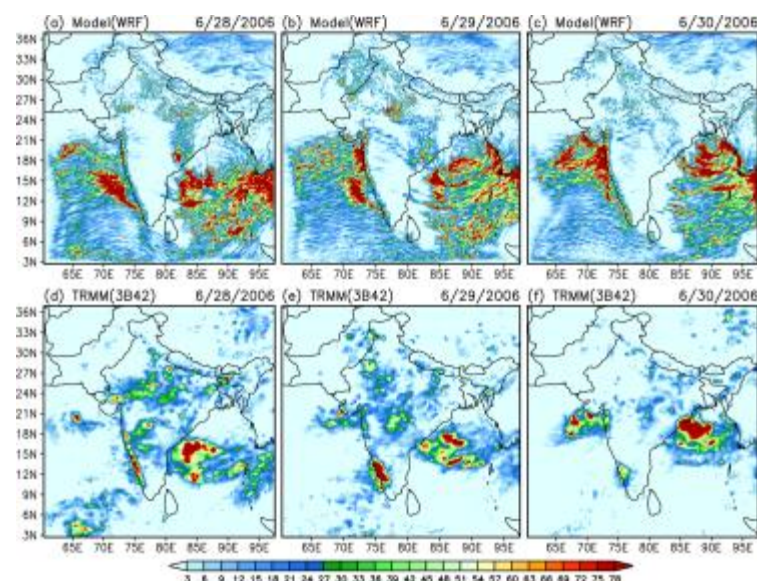
- 20K InfiniBand nodes
- Mellanox end-to-end FDR and QDR InfiniBand
- Supports variety of scientific and engineering projects
 - Coupled atmosphere-ocean models
 - Future space vehicle design
 - Large-scale dark matter halos and galaxy evolution



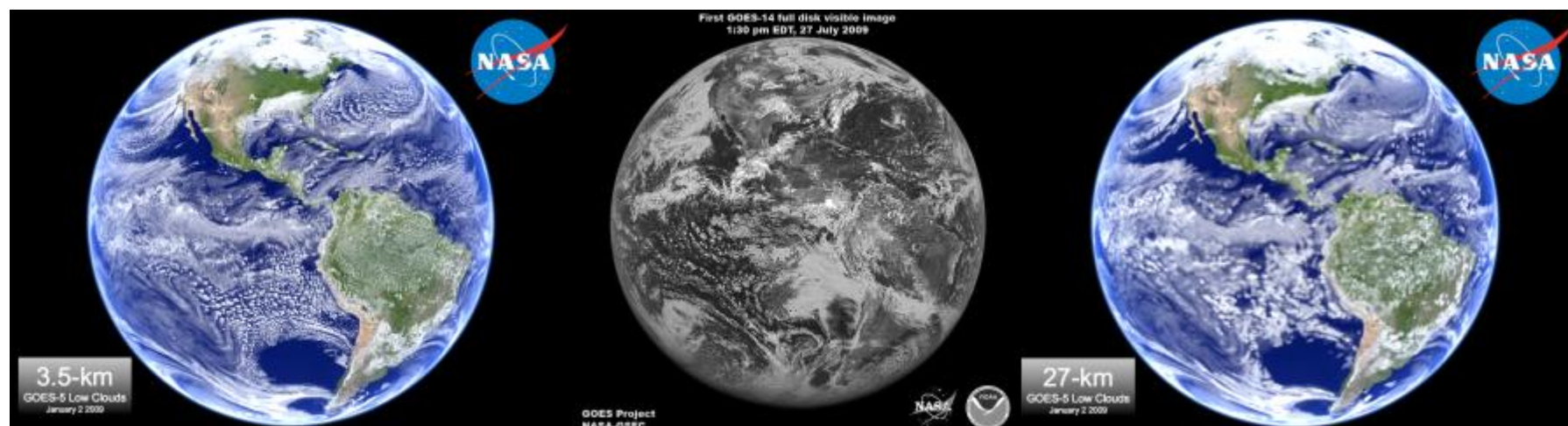
Dataset: ifloods d03 RIP: rip accp12h d03 Init: 0000 UTC Mon 24 Jun 13
 Fcst: 27.00 h Valid: 0300 UTC Tue 25 Jun 13 (2200 CDT Mon 24 Jun 13)
 Total precip. in past 12 h
 Sea-level pressure
 Horizontal wind vectors at k-index = 60



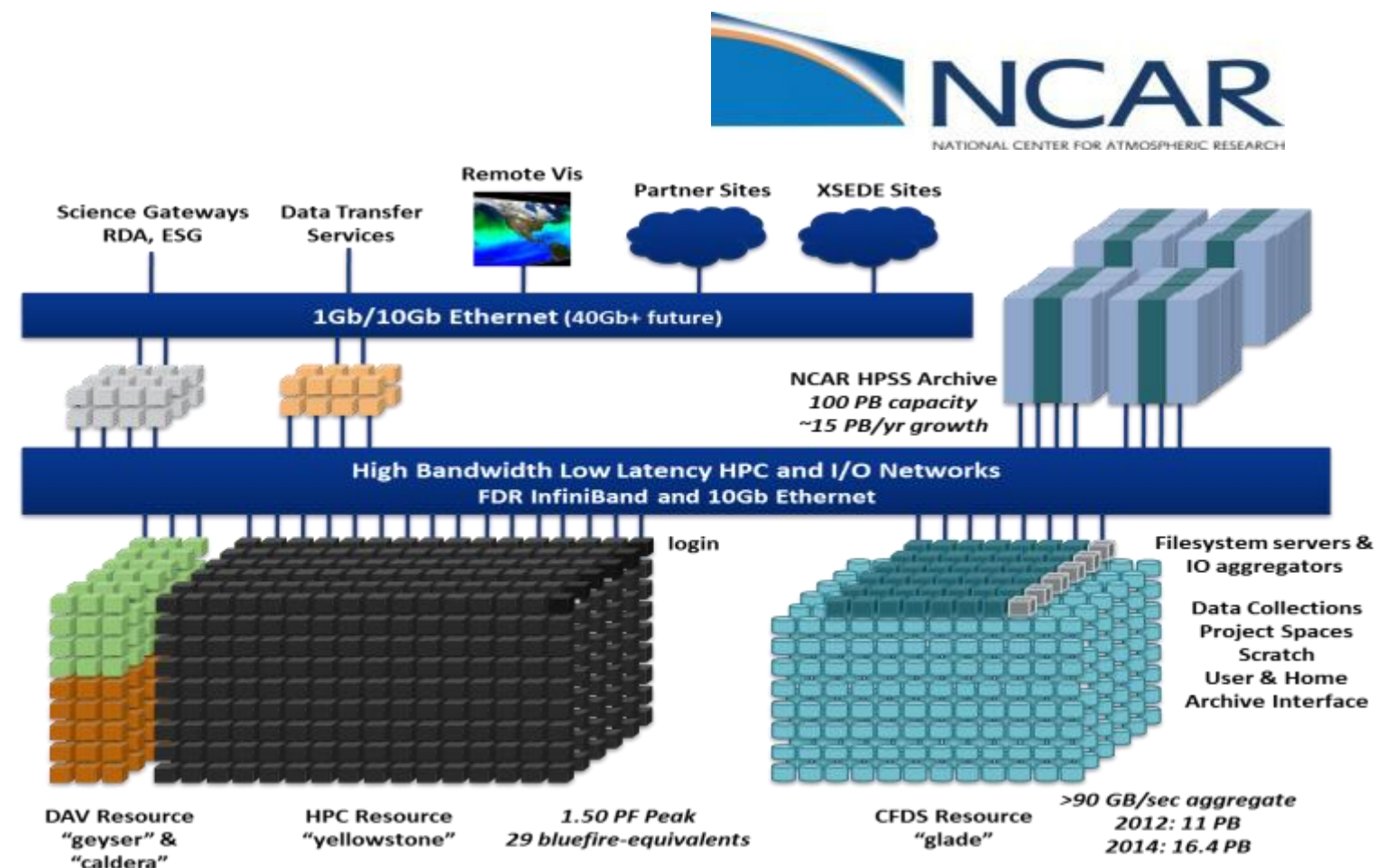
Asian Monsoon Water Cycle

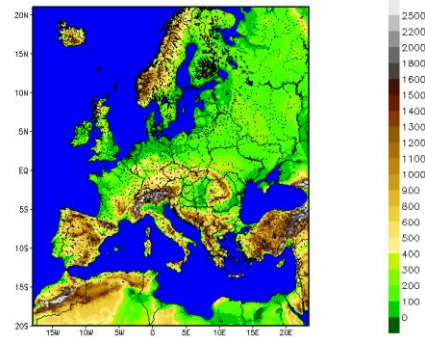


High-Resolution Climate Simulations

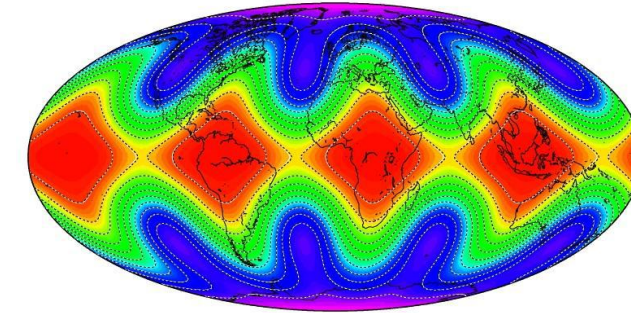
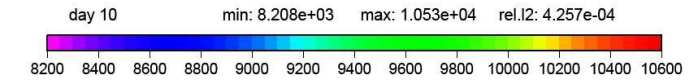
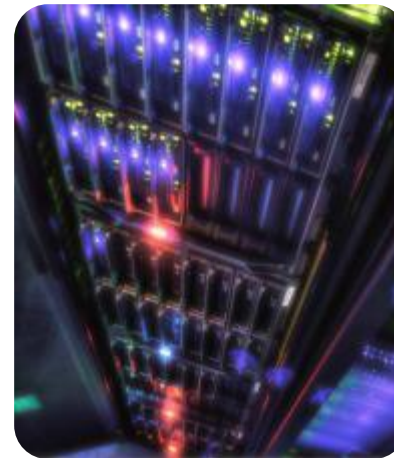
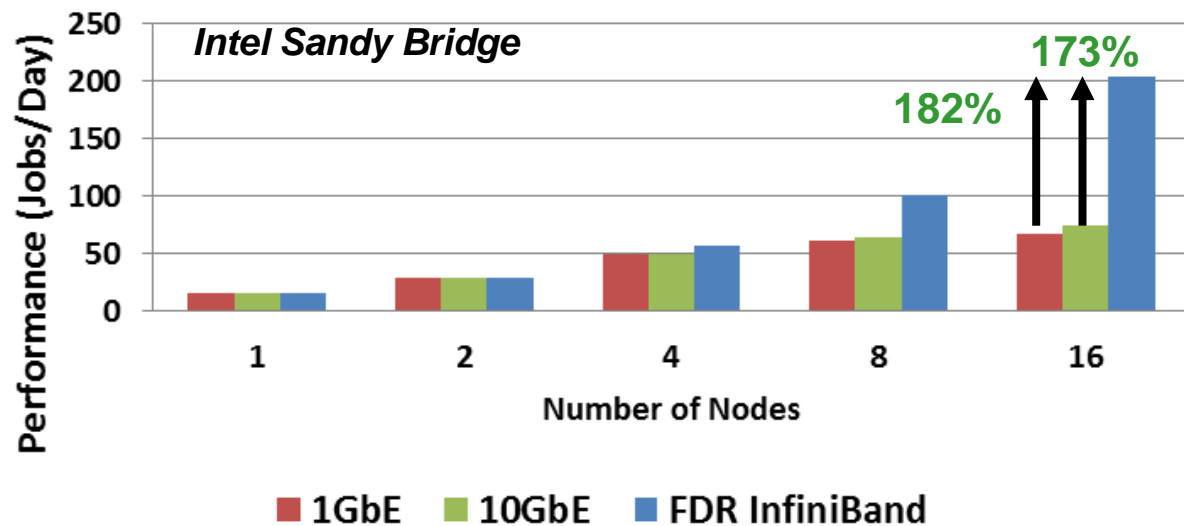


- “Yellowstone” system
- 72,288 processor cores, 4,518 nodes
- Mellanox end-to-end FDR InfiniBand, CLOS (full fat tree) network, single plane

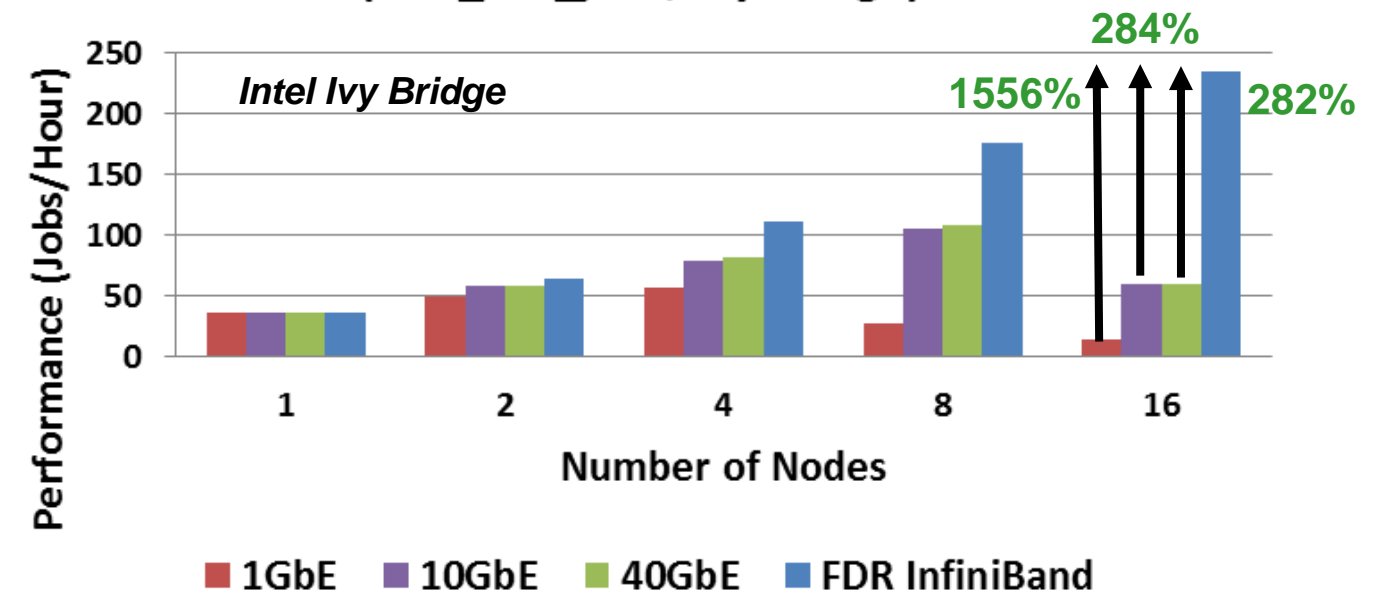




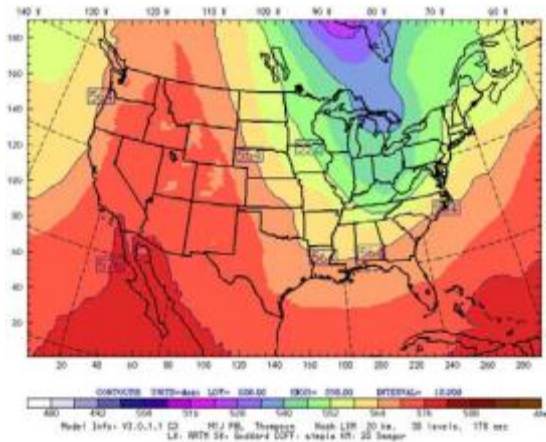
COSMO RAPS 5.1 Performance (COSMO_EU)



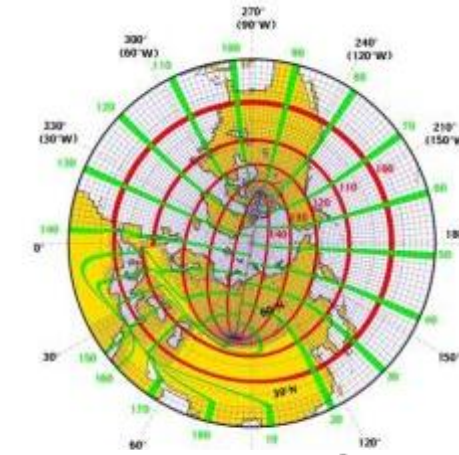
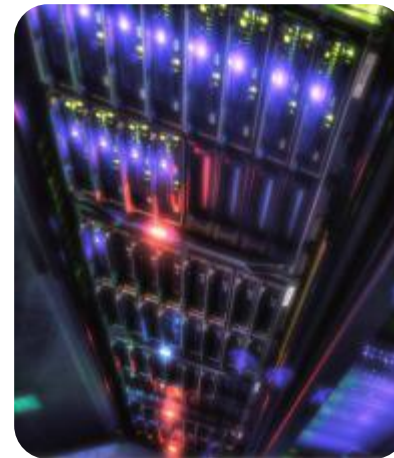
ICON RAPS 2.0 Performance (test_hat_jww, Ivy Bridge)



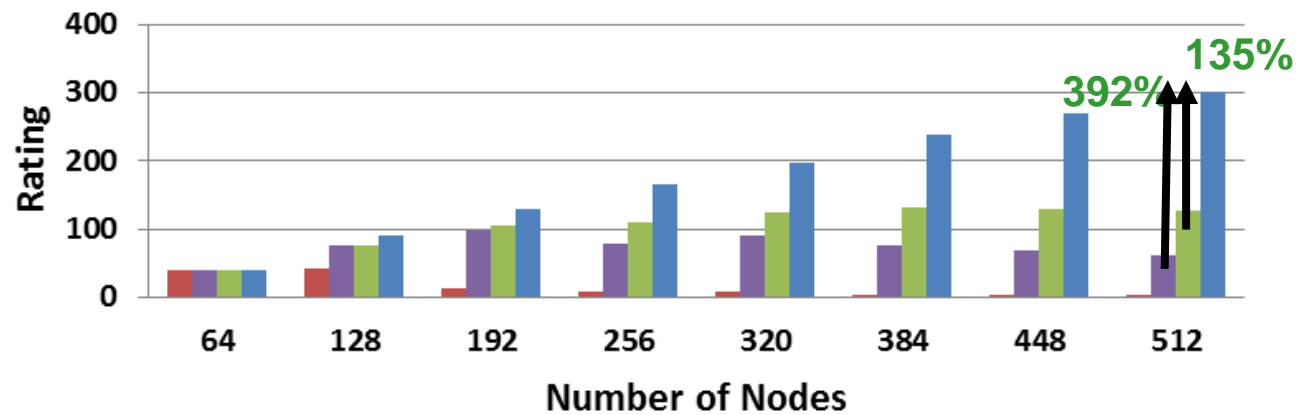
Applications Performance (Courtesy of the HPC Advisory Council)



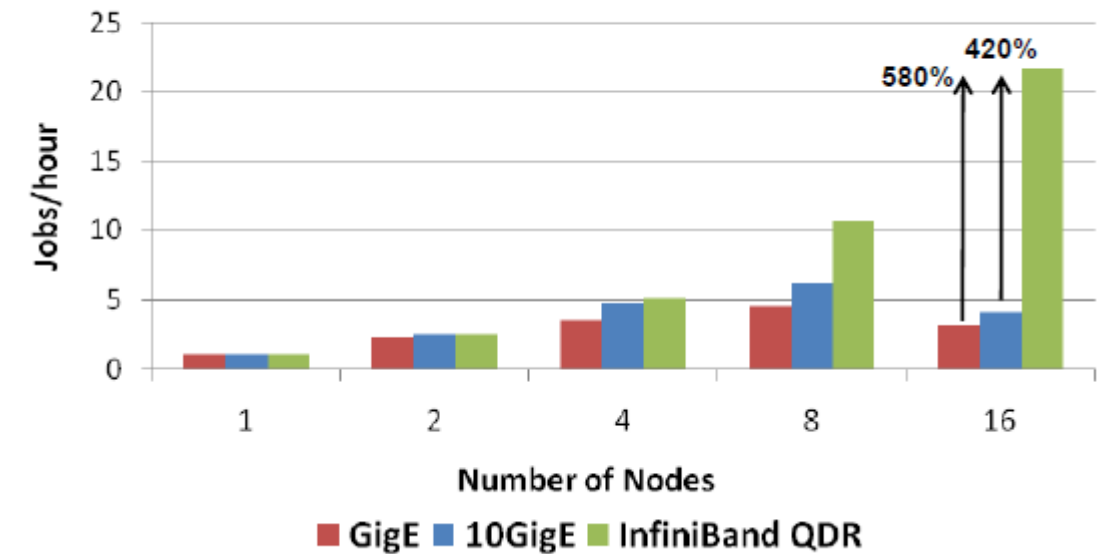
WRF Benchmark
(conus12km)



NEMO Benchmark



■ 1GbE ■ 10GbE ■ 40GbE ■ InfiniBand



Dominant in Enterprise Back-End Storage Interconnects

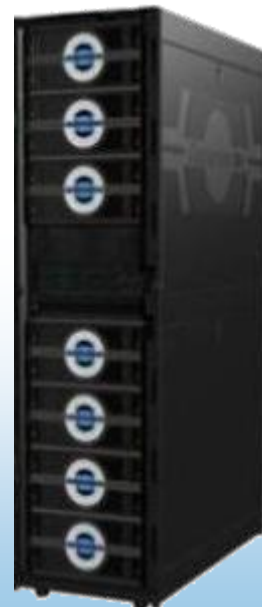
 Microsoft
SMB Direct




NetApp















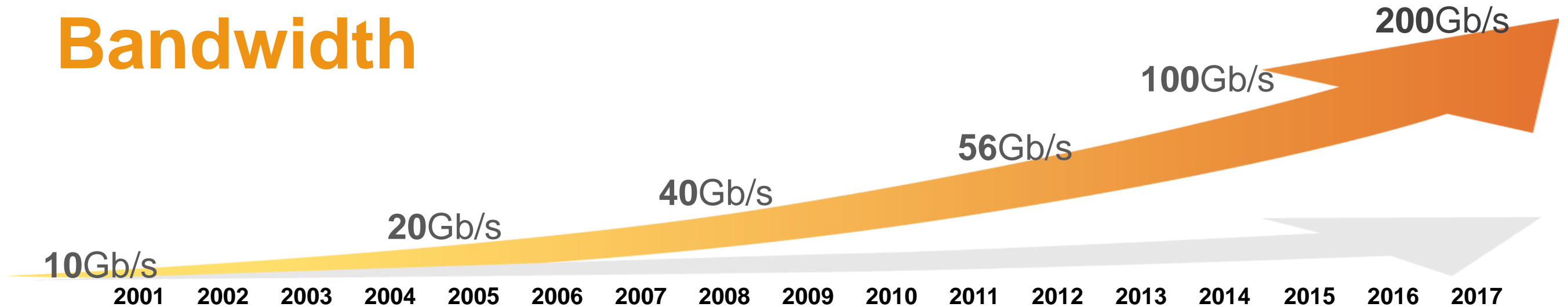




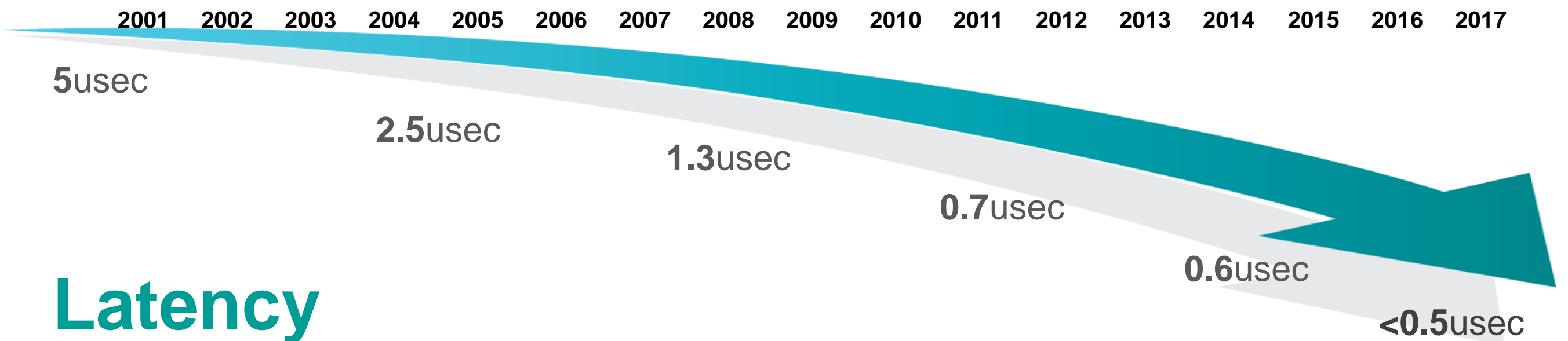





Bandwidth



Same Software Interface



Latency

Connect-IB

Architectural Foundation for Exascale Computing

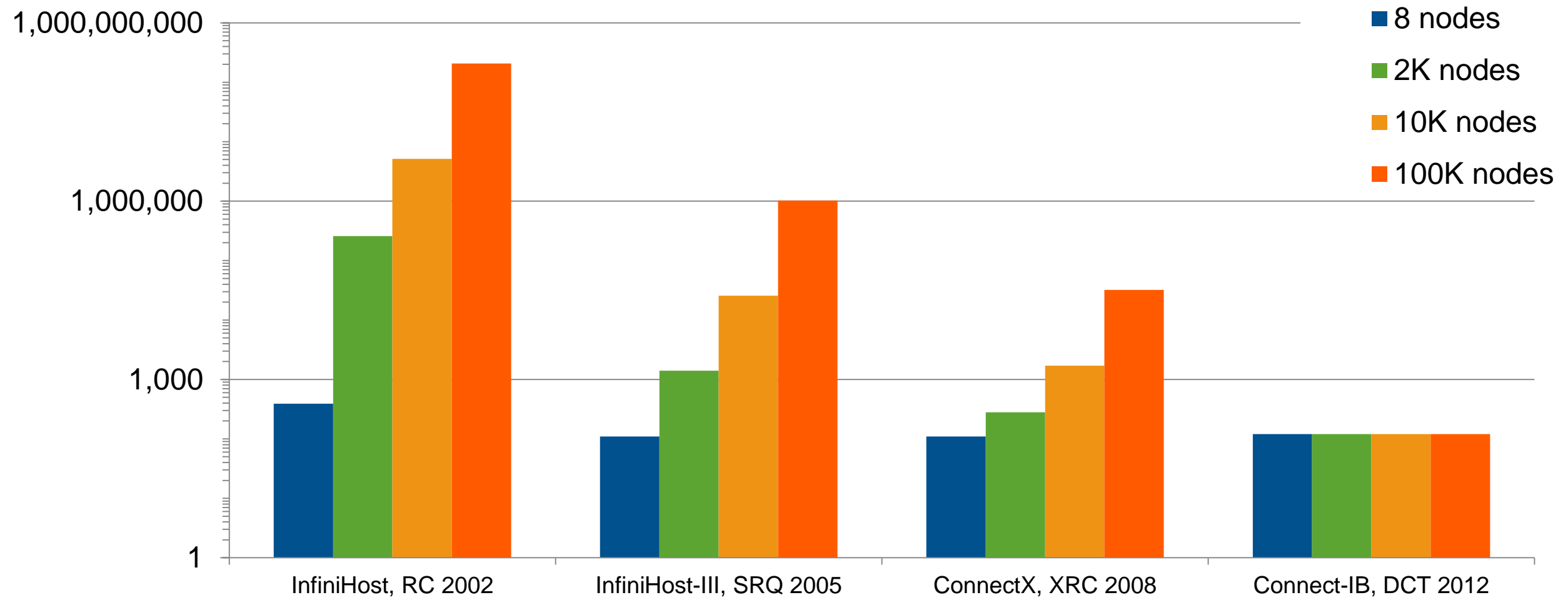
- World's first 100Gb/s interconnect adapter
 - PCIe 3.0 x16, dual FDR 56Gb/s InfiniBand ports to provide >100Gb/s
- Highest InfiniBand message rate: 137 million messages per second
 - 4X higher than other InfiniBand solutions
- <0.7 micro-second application latency
- Supports GPUDirect RDMA for direct GPU-to-GPU communication
- Unmatchable Storage Performance
 - 8,000,000 IOPs (1QP), 18,500,000 IOPs (32 QPs)
- New Innovative Transport – Dynamically Connected Transport Service
- Supports Scalable HPC with MPI, SHMEM and PGAS/UPC offloads



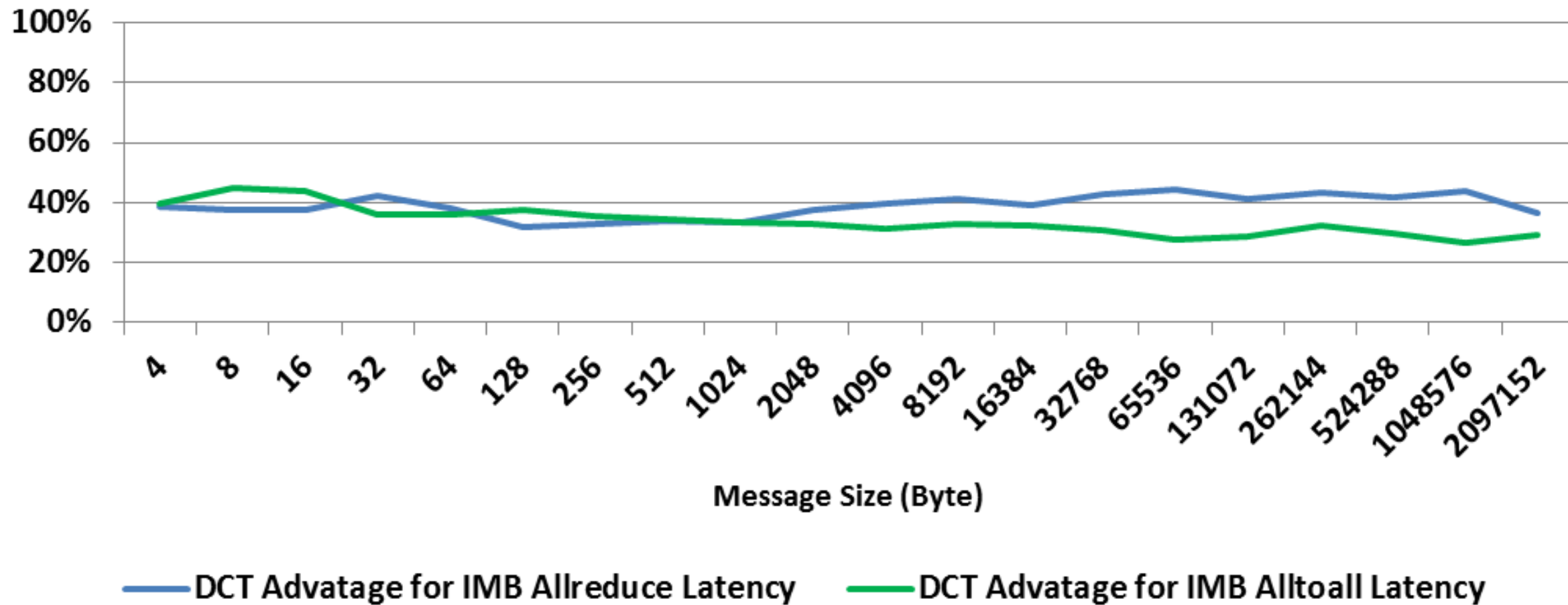
Connect **IB**[™]

Enter the World of Boundless Performance

Host Memory Consumption (MB)



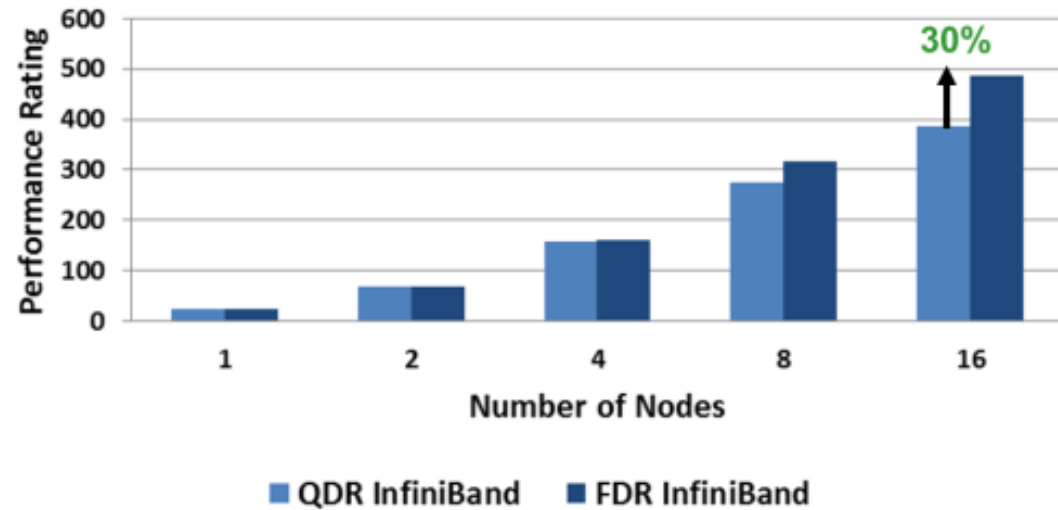
Dynamically Connected Transport Advantages (MPI Latency Reduction %)



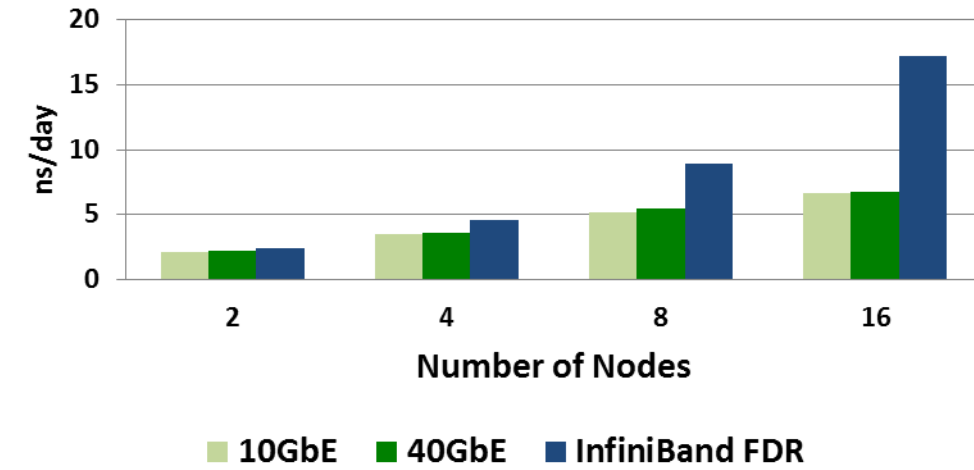
FDR InfiniBand Delivers Highest Application Performance



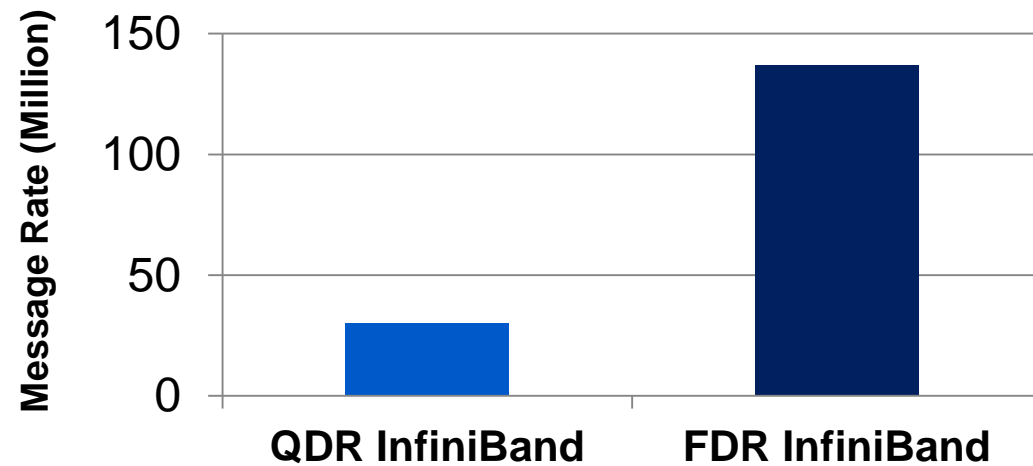
OpenFOAM Performance (Lid-driven Cavity)



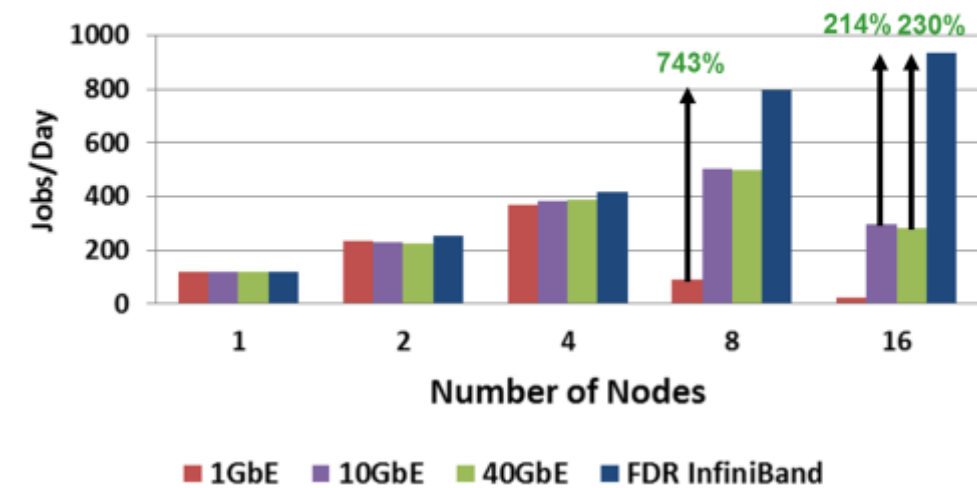
NAMD Benchmark (Platform MPI, ApoA1)



Message Rate



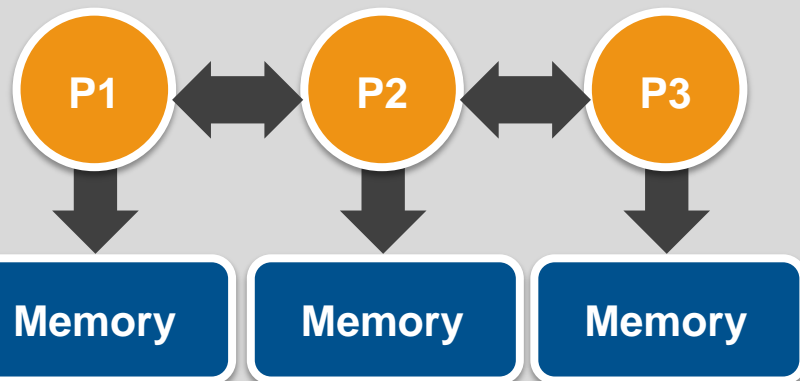
RADIOSS Benchmark (NEON1M11, MPI)



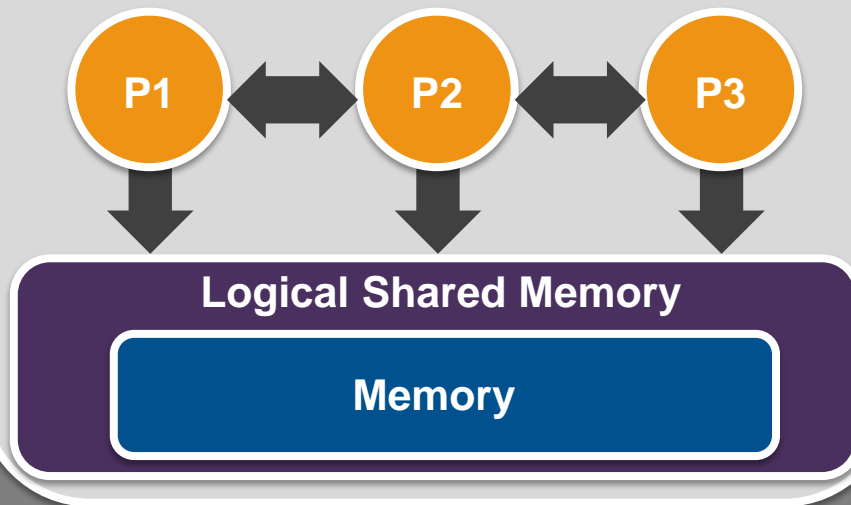
Scalable Communication

MXM, FCA

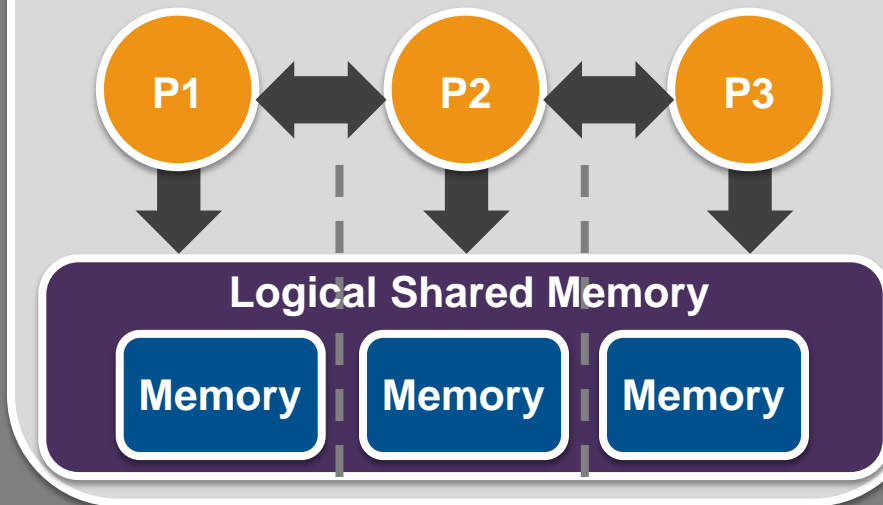
MPI



SHMEM



PGAS



MXM

- Reliable Messaging Optimized for Mellanox HCA
- Hybrid Transport Mechanism
- Efficient Memory Registration
- Receive Side Tag Matching

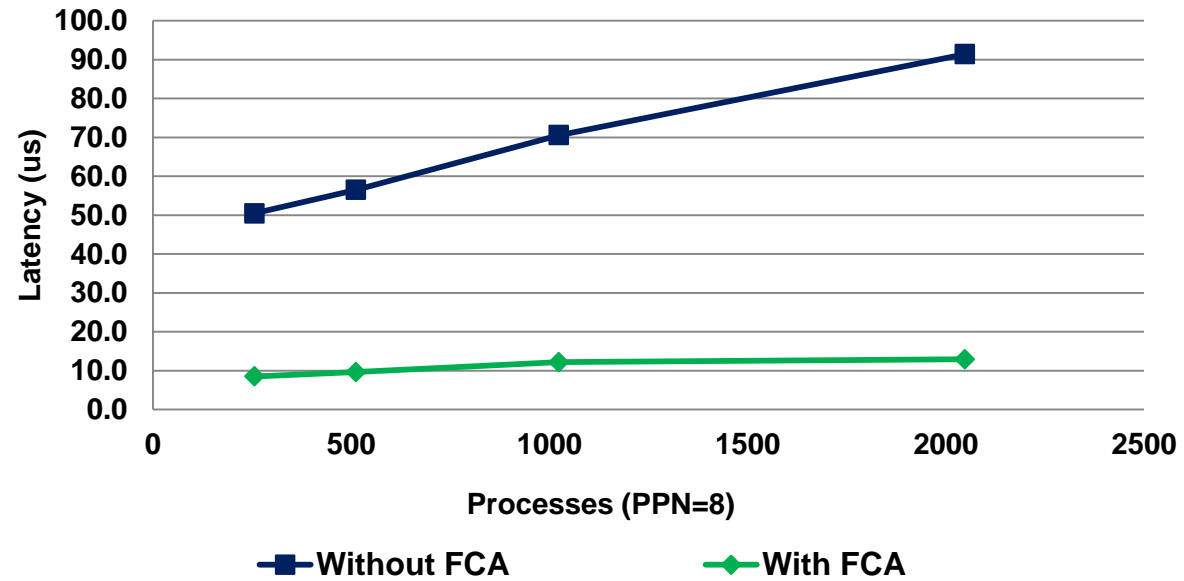
FCA

- Topology Aware Collective Optimization
- Hardware Multicast
- Separate Virtual Fabric for Collectives
- CORE-Direct Hardware Offload

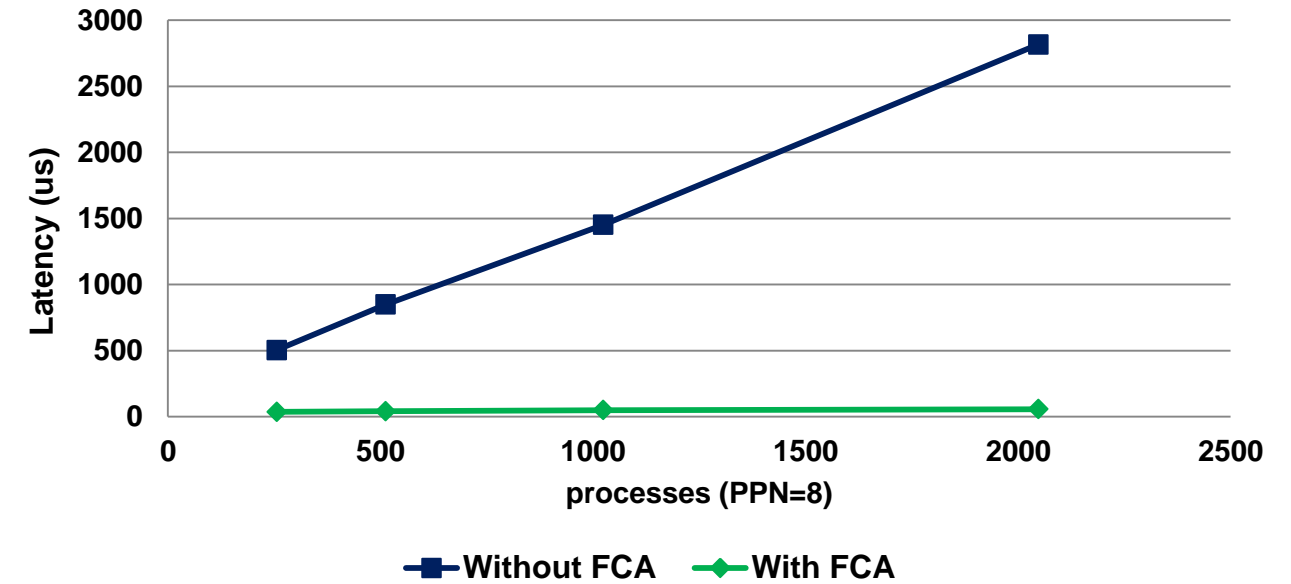
InfiniBand Verbs API

- Transport library integrated with OpenMPI, OpenSHMEM, BUPC, Mvapich2
 - More solutions will be added in the future
 - Utilizing Mellanox offload engines
- Supported APIs (both sync/async): AM, p2p, atomics, synchronization
- Supported transports: RC, UD, DC, RoCE, SHMEM
- Supported built-in mechanisms: tag matching, progress thread, memory registration cache, fast path send for small messages, zero copy, flow control
- Supported data transfer protocols: Eager Send/Recv, Eager RDMA, Rendezvous

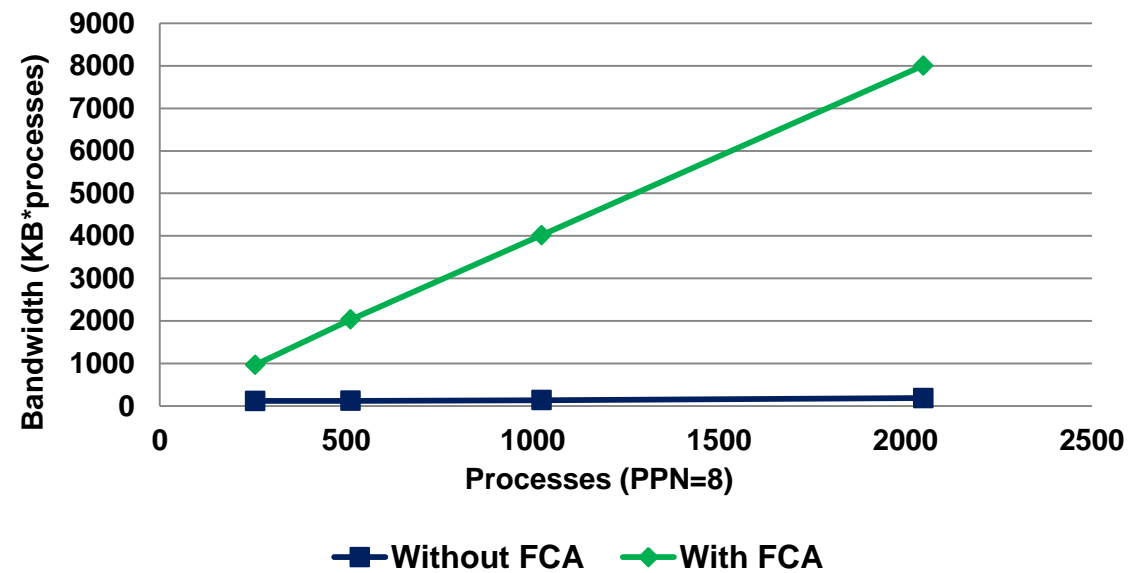
Barrier Collective



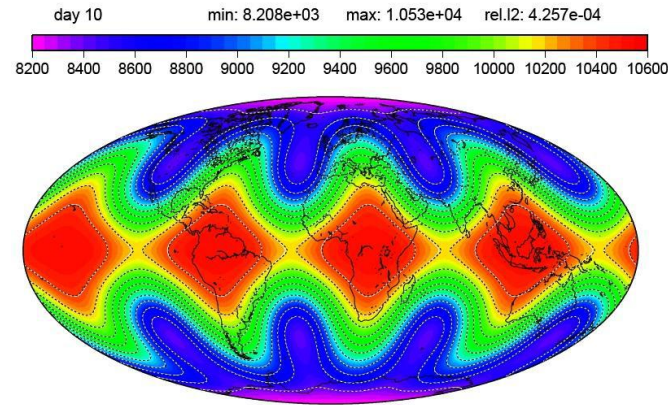
Reduce Collective



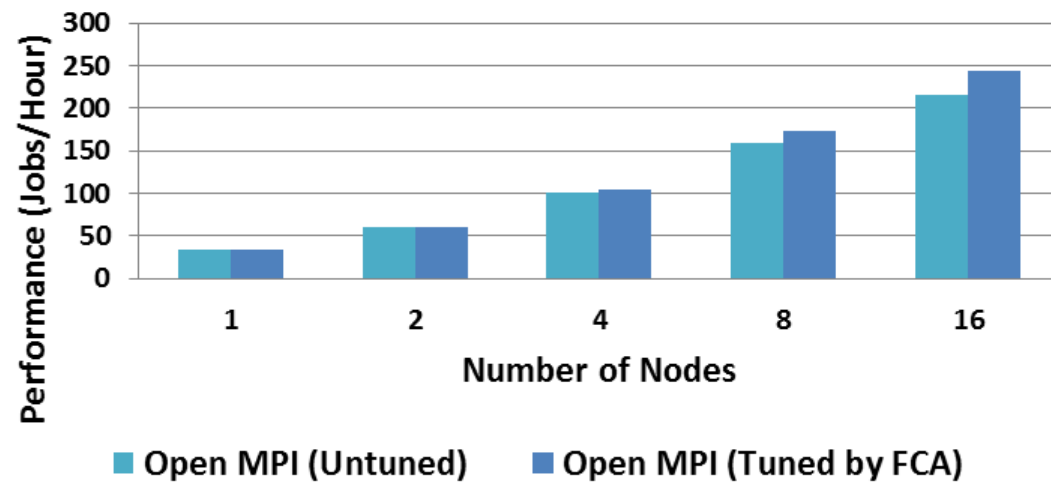
8-Byte Broadcast



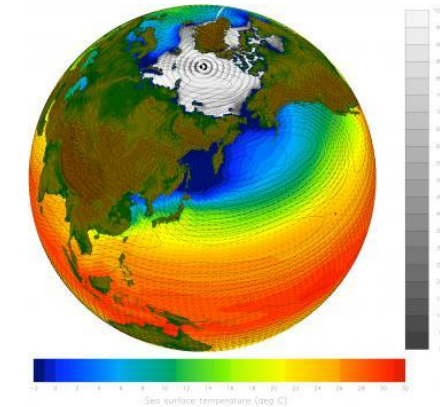
FDR InfiniBand Delivers Highest Application Performance



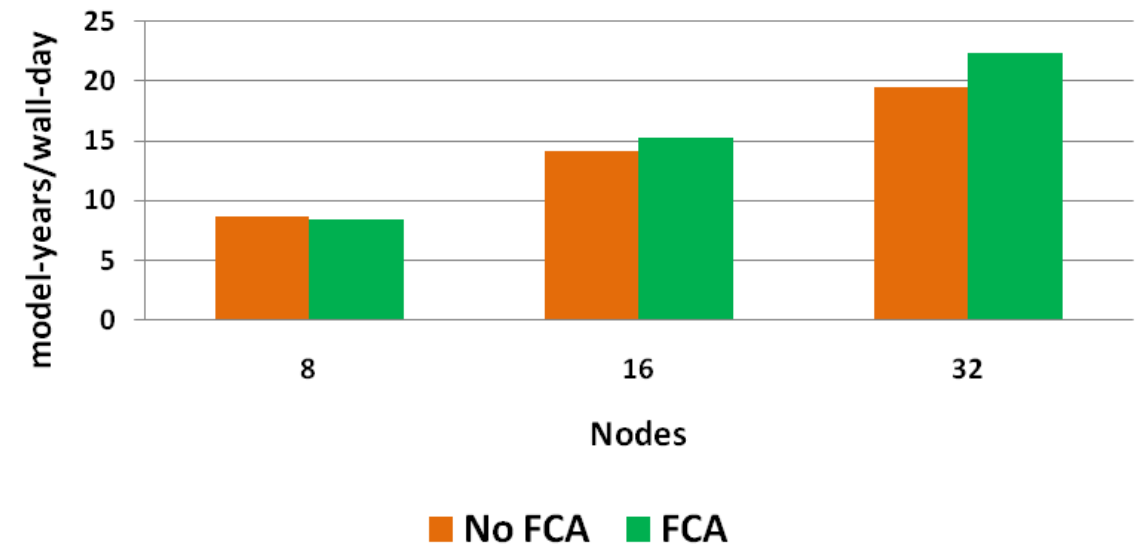
ICON 2.0 Performance
(test_hat_jww)



CESM
COMMUNITY EARTH SYSTEM MODEL



CESM Performance
(B1850CN, ATM)



GPU Direct

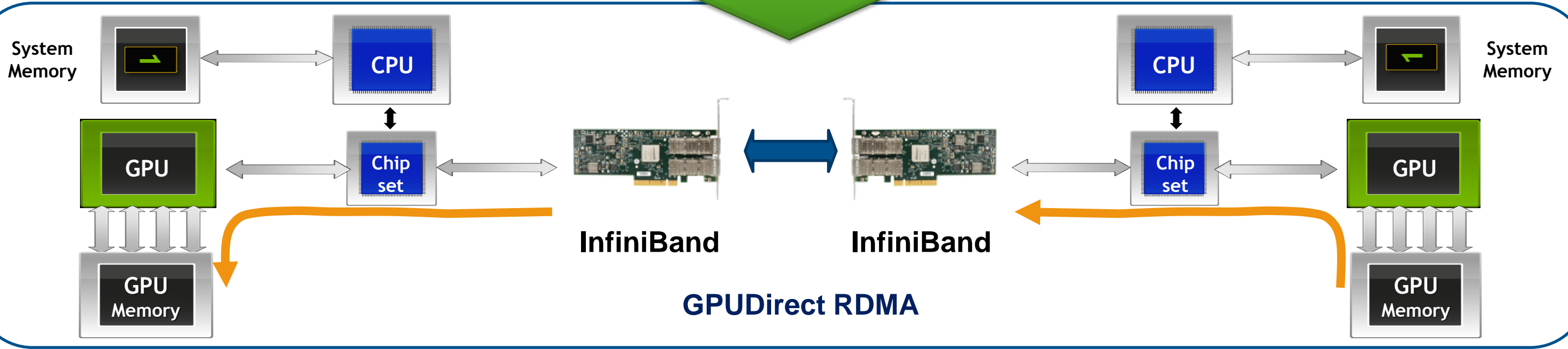
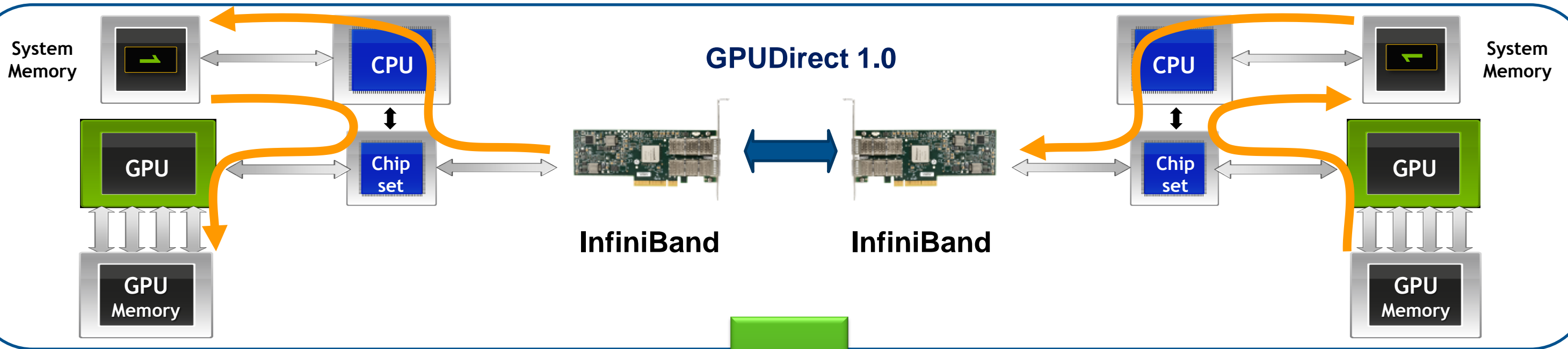
Receive

Transmit

GPUDirect 1.0

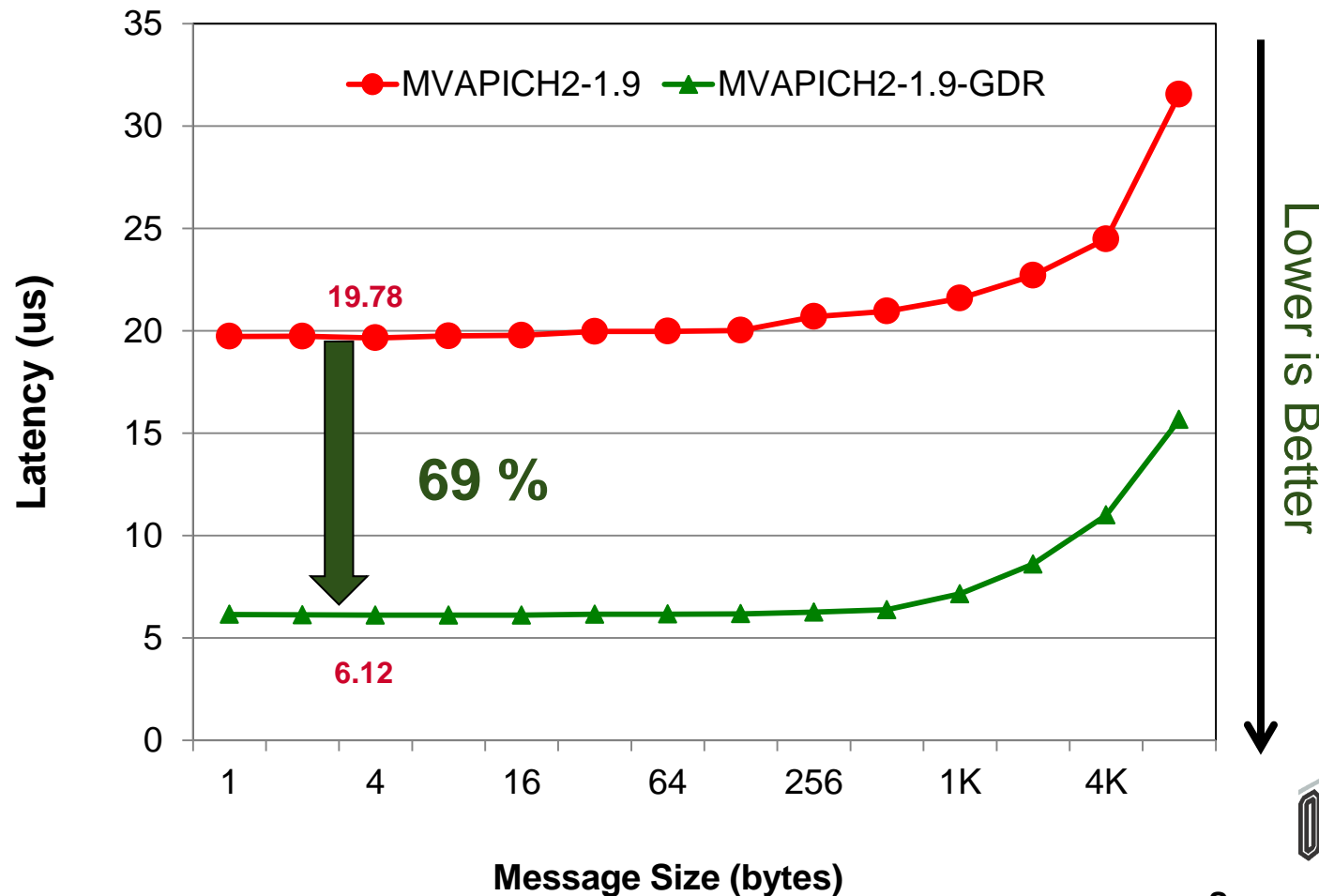


GPUDirect RDMA



GPU-GPU Internode MPI Latency

Small Message Latency



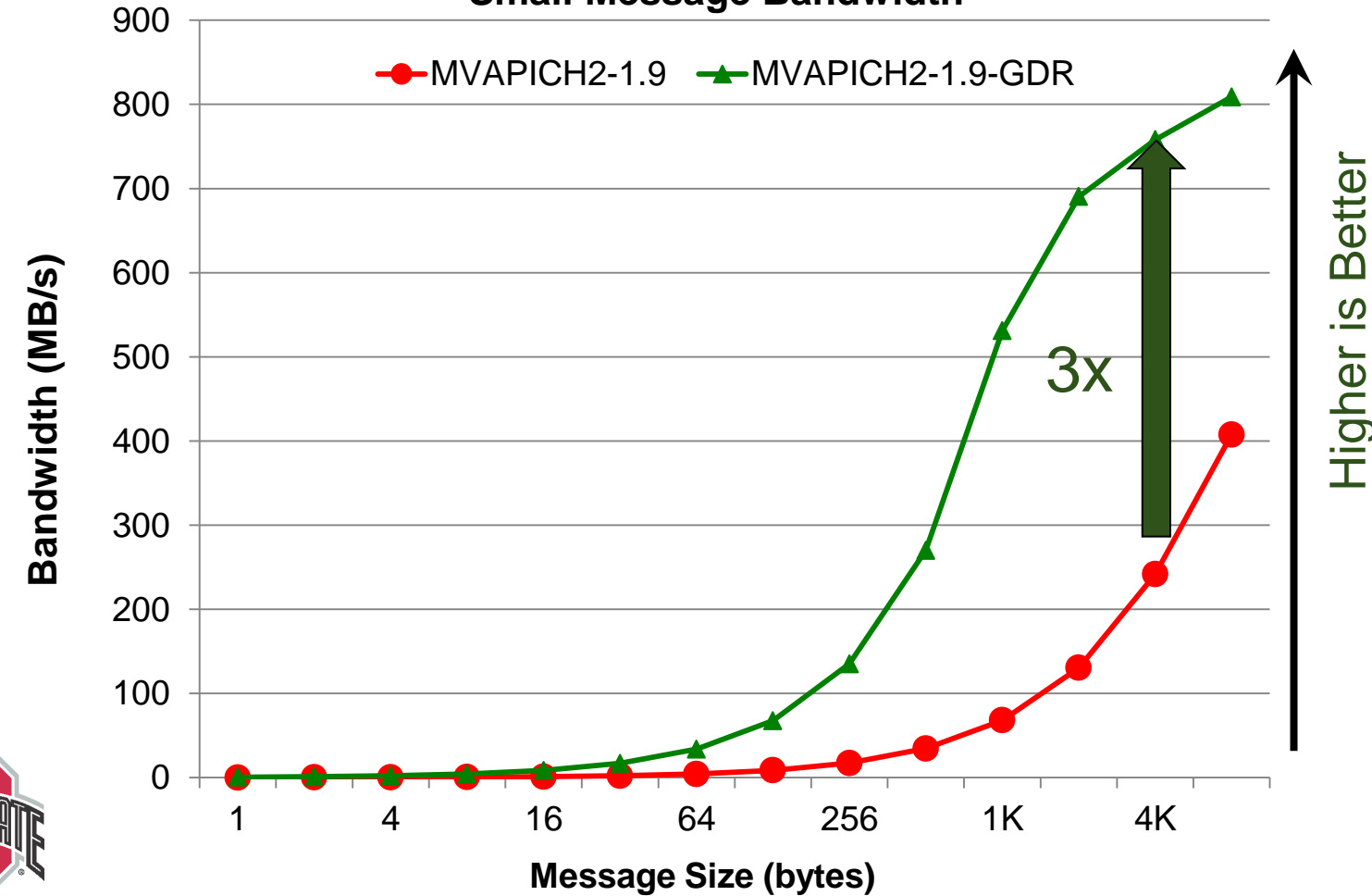
Lower is Better



Source: Prof. DK Panda

GPU-GPU Internode MPI Bandwidth

Small Message Bandwidth



Higher is Better

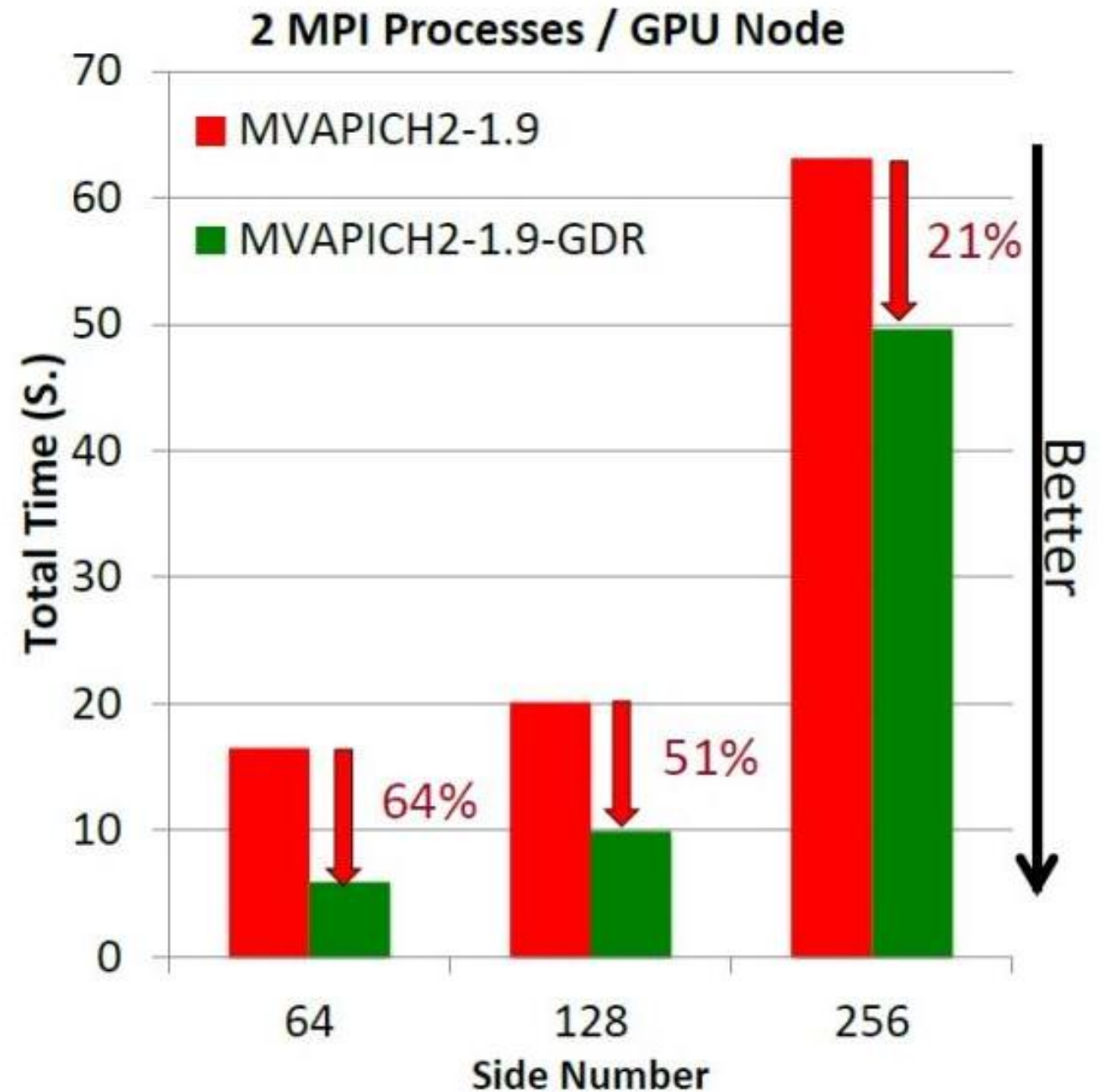
69% Lower Latency

3X Increase in Throughput

Execution Time of HSG (Heisenberg Spin Glass) Application with 2 GPU Nodes



Source: Prof. DK Panda



Thank You

