



## The Future Of Interconnect



# Leading Supplier of End-to-End Interconnect Solutions



## Comprehensive End-to-End InfiniBand and Ethernet Portfolio

ICs	Adapter Cards	Switches/Gateways	Host/Fabric Software	Metro / WAN	Cables/Modules



# Mellanox InfiniBand Connected Petascale Systems

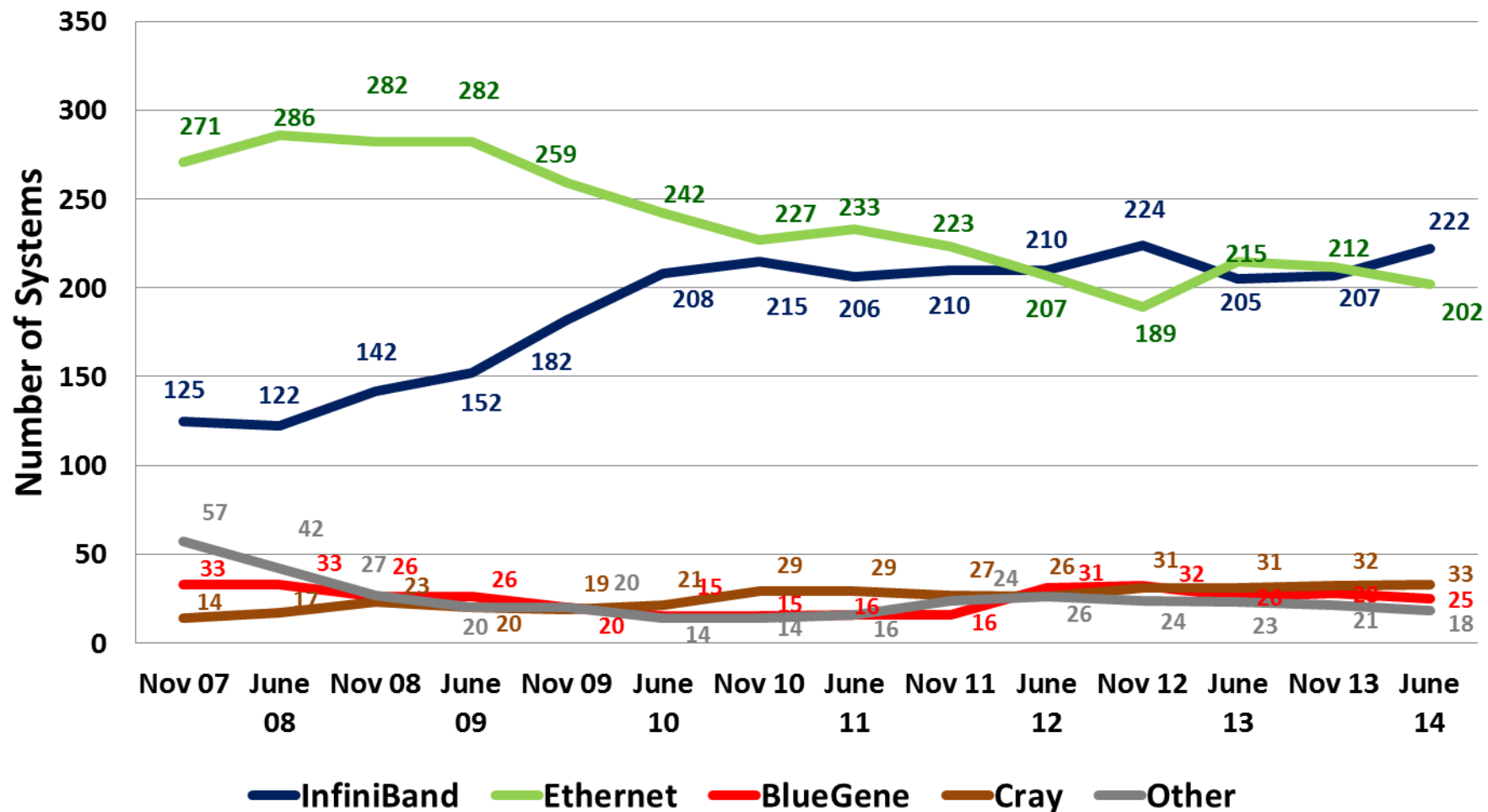


**Connecting Half of the World's Petascale Systems**  
Mellanox Connected Petascale System Examples

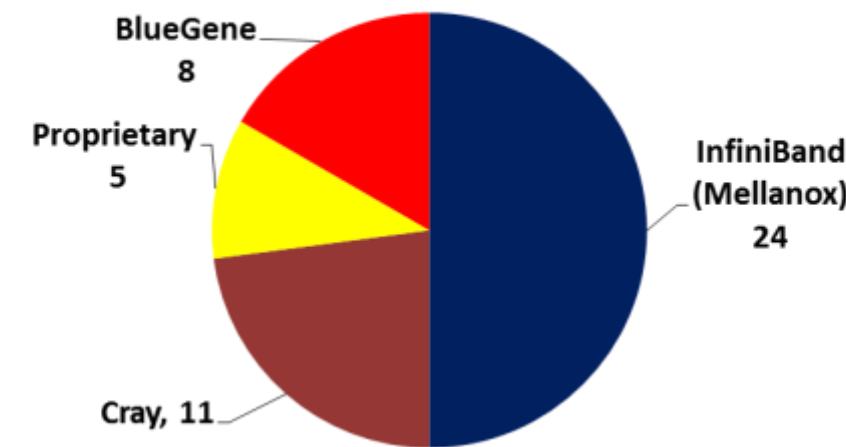
# Mellanox Enables Most Efficient HPC Platforms



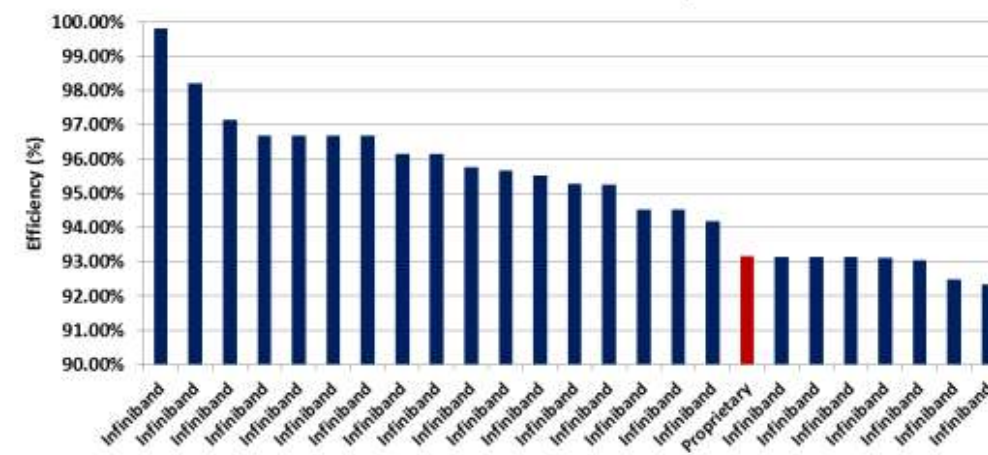
## TOP500 Interconnect Trends



## PetaFlop Capable Systems on the TOP500 list



## TOP500 - TOP25 Most Efficient System

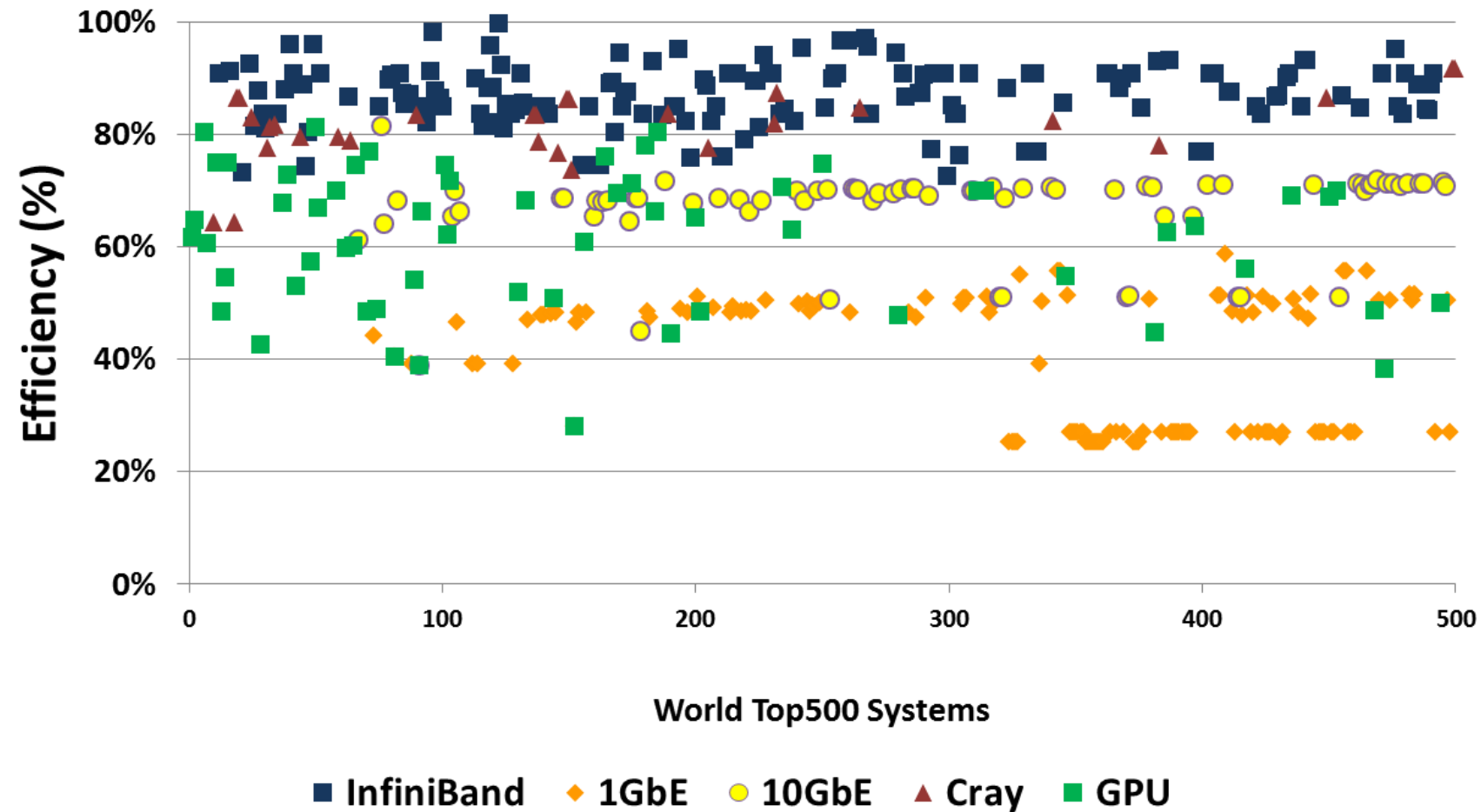




# InfiniBand's Unsurpassed System Efficiency



## World Leading Compute Systems Efficiency Comparison



### Average Efficiency

- **InfiniBand: 87%**
- **Cray: 79%**
- **10GbE: 67%**
- **GigE: 40%**

- TOP500 systems listed according to their efficiency
- InfiniBand is the key element responsible for the highest system efficiency; in average 30% higher than 10GbE
- Mellanox delivers efficiencies of more than 99% with InfiniBand

Connect **IB**

# Architectural Foundation for Exascale Computing

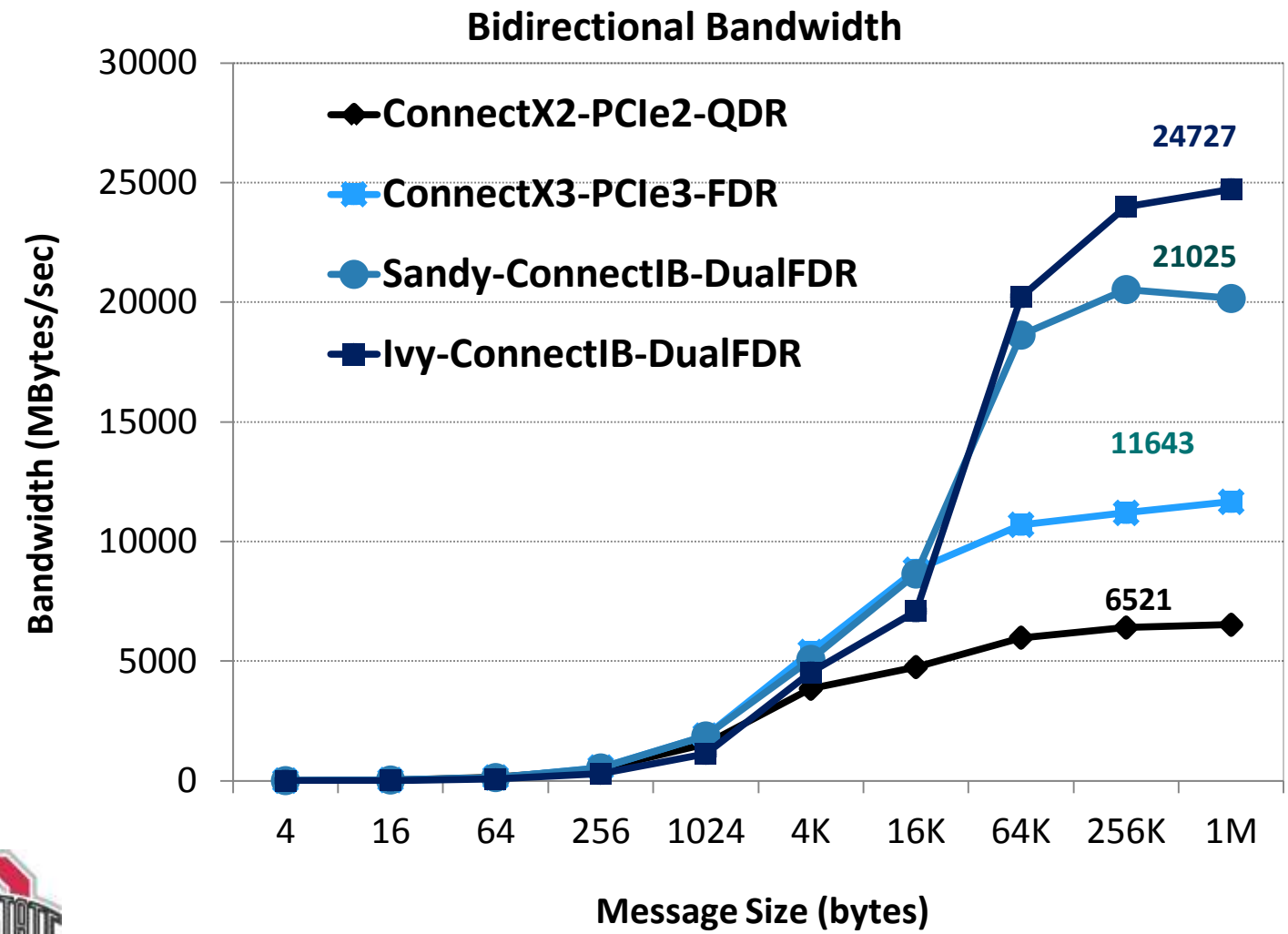
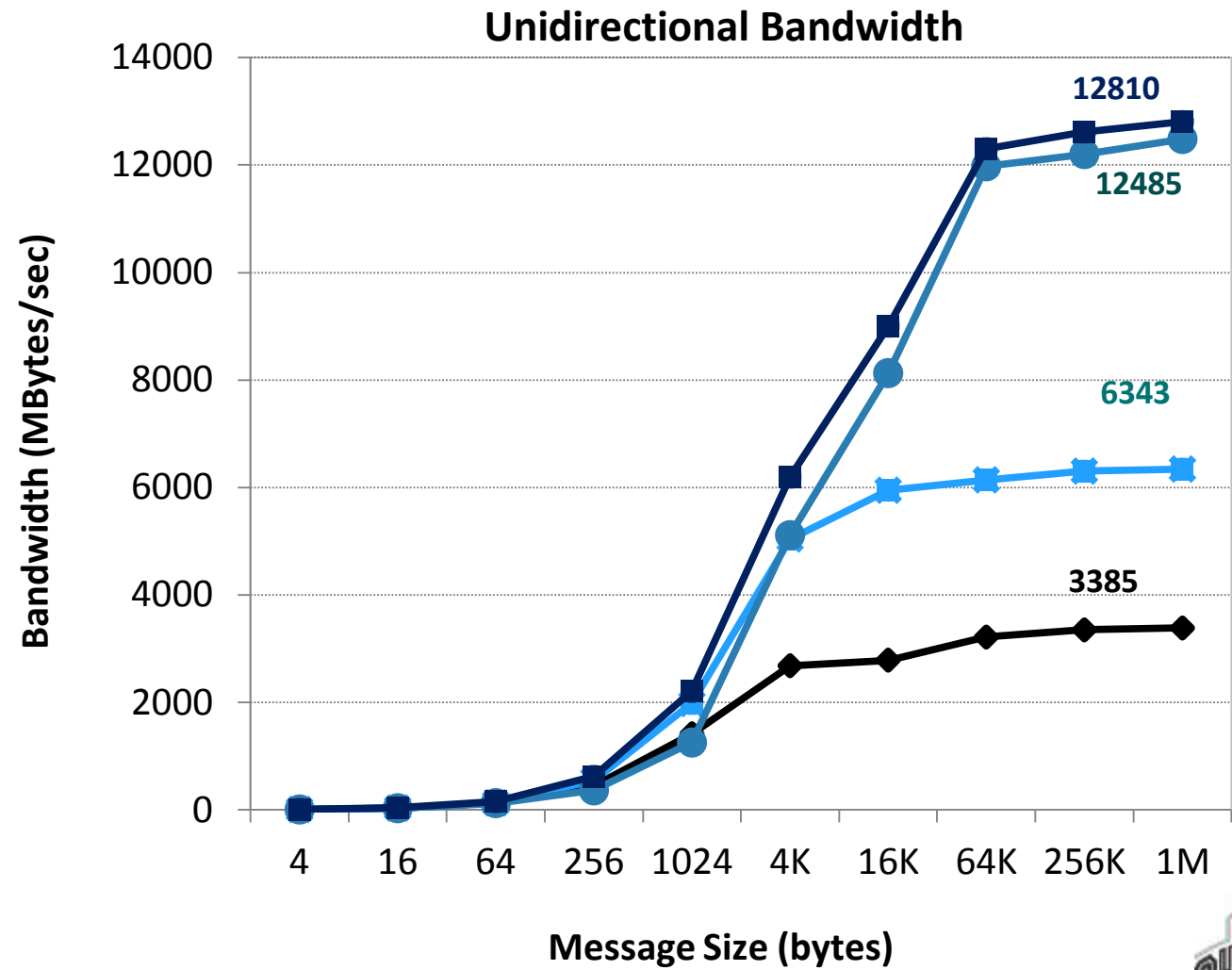
- The 7<sup>th</sup> generation of Mellanox interconnect adapters
- World's first 100Gb/s interconnect adapter (dual-port FDR 56Gb/s InfiniBand)
- Delivers 137 million messages per second – 4X higher than competition
- Support the new innovative InfiniBand scalable transport – Dynamically Connected



# Connect-IB Provides Highest Interconnect Throughput



Higher is Better



Source: Prof. DK Panda

**Gain Your Performance Leadership With Connect-IB Adapters**

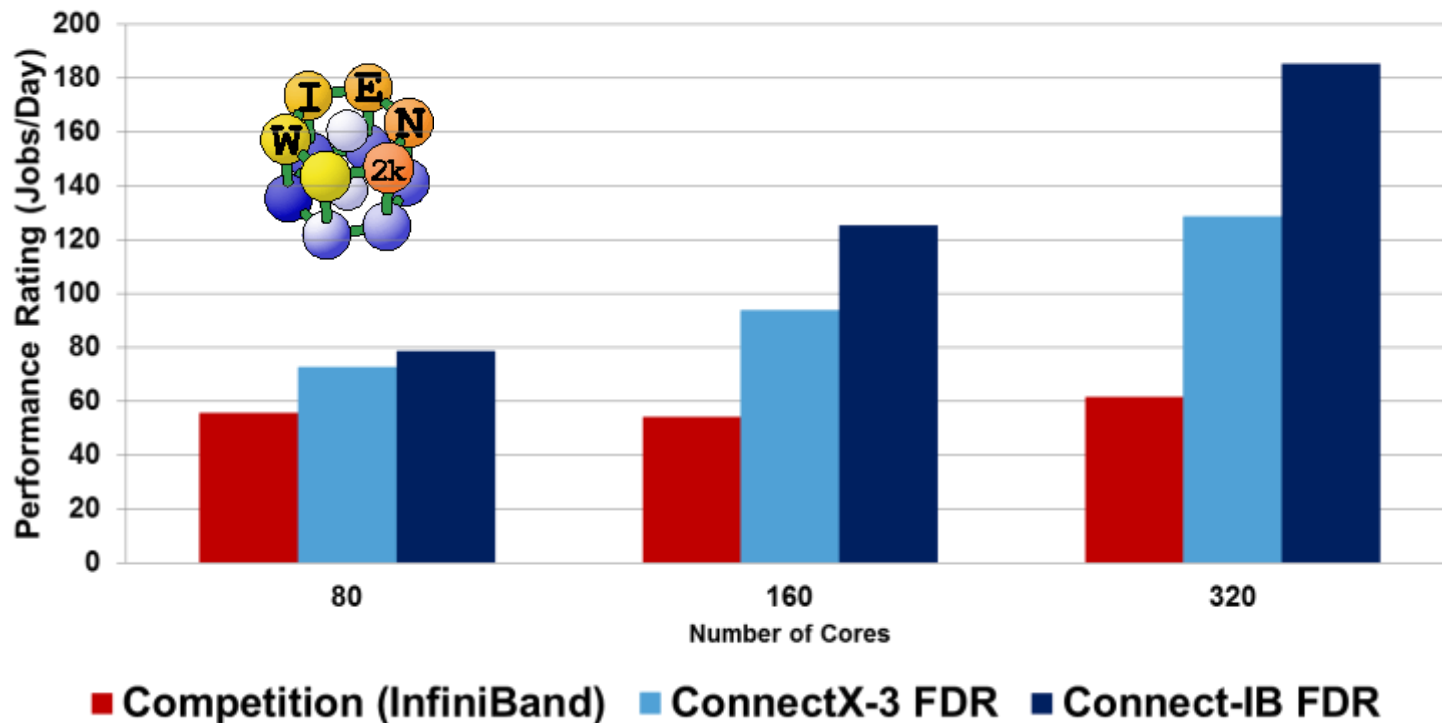


# Connect-IB Delivers Highest Application Performance

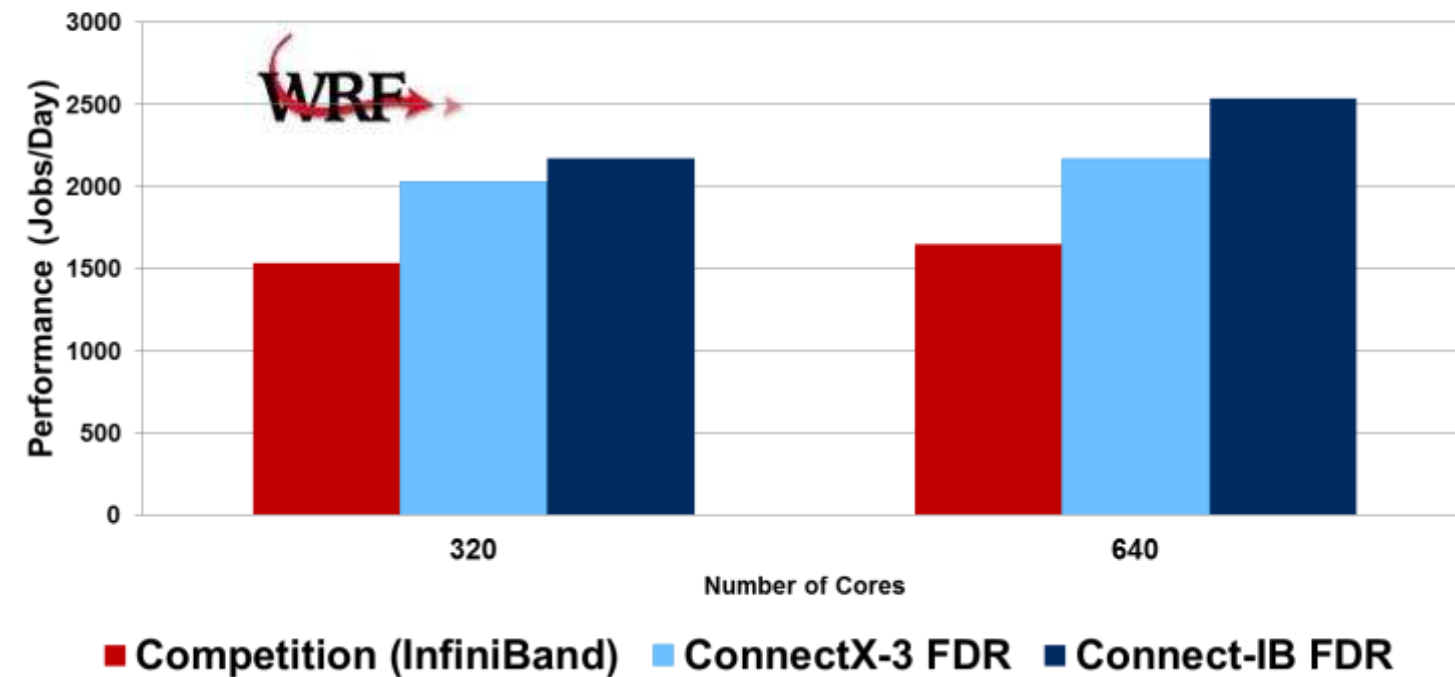


**200% Higher Performance Versus Competition, with Only 32-nodes**  
**Performance Gap Increases with Cluster Size**

### WIEN2k Performance



### WRF Performance (conus12km)



NEW

Switch IB™



## 100Gb/s Cables Demonstrated March '14



Copper (Passive, Active)



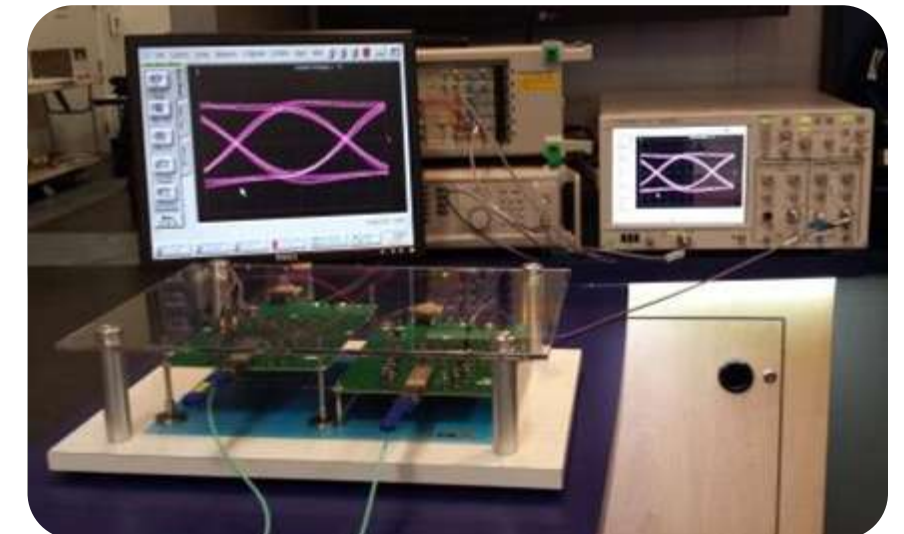
Optical Cables (VCSEL)



Silicon Photonics



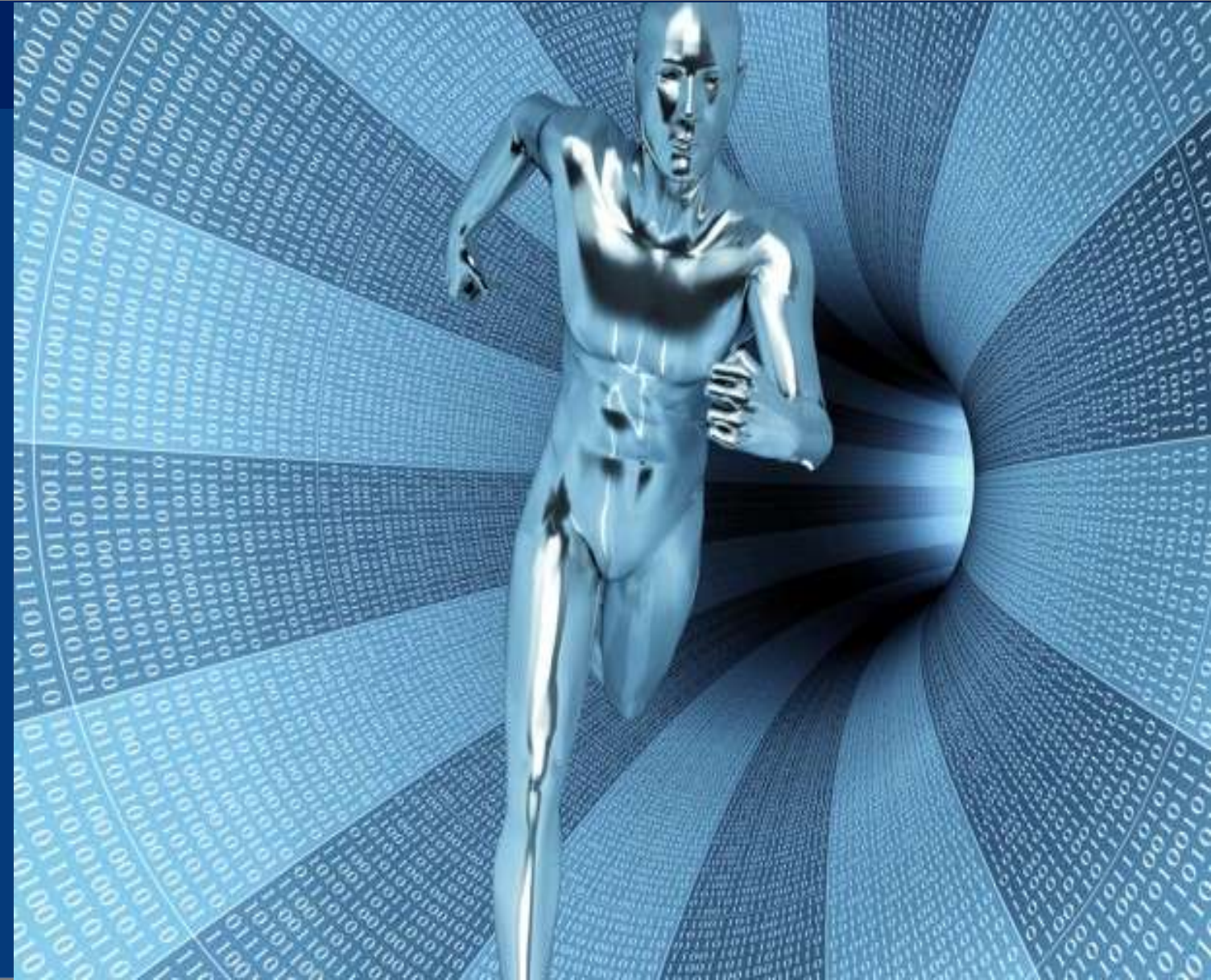
100Gb/s cables demo  
at OFC conference  
March '14





## Switch-IB EDR 100G InfiniBand

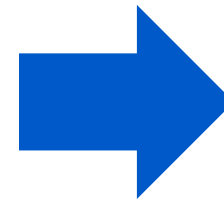
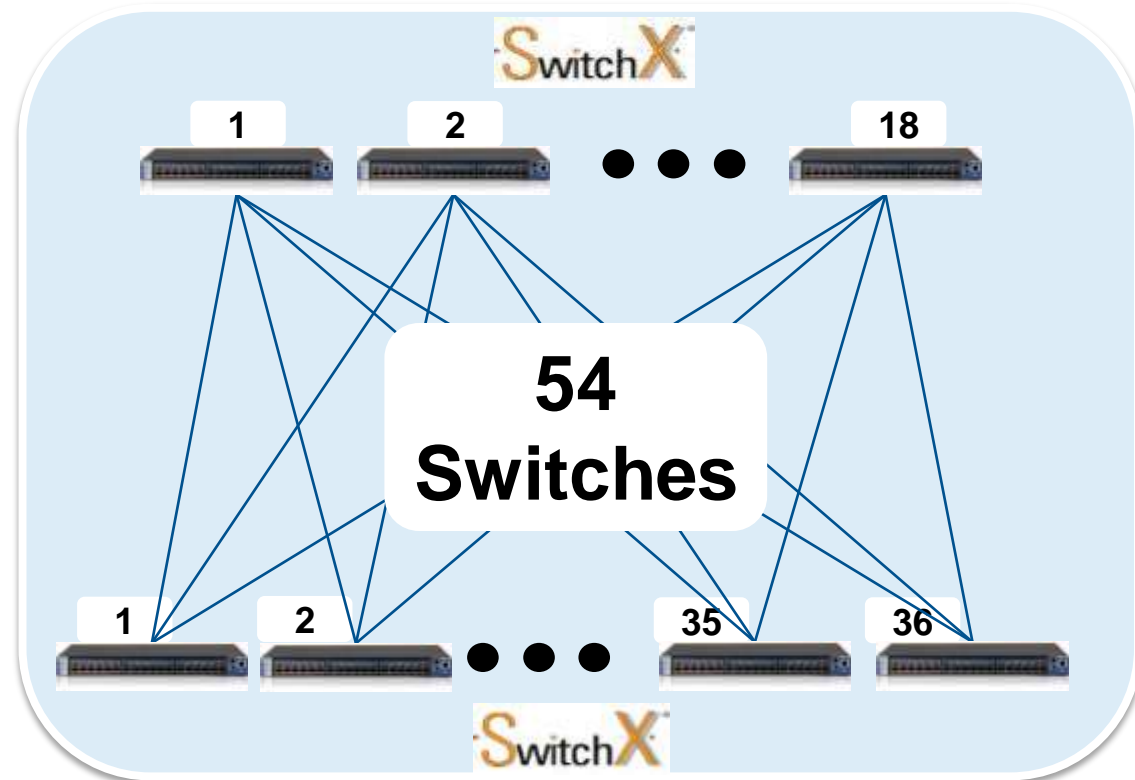
- 7<sup>th</sup> generation of Mellanox interconnect switch
- 36 EDR (100Gb/s) Ports
- <130ns latency
- Throughput - 7.2 Tb/s
- Software compatible with FDR InfiniBand
- x86 CPU management
- InfiniBand router
- Multiple topologies (Fat-Tree, Torus, Dragonfly +)



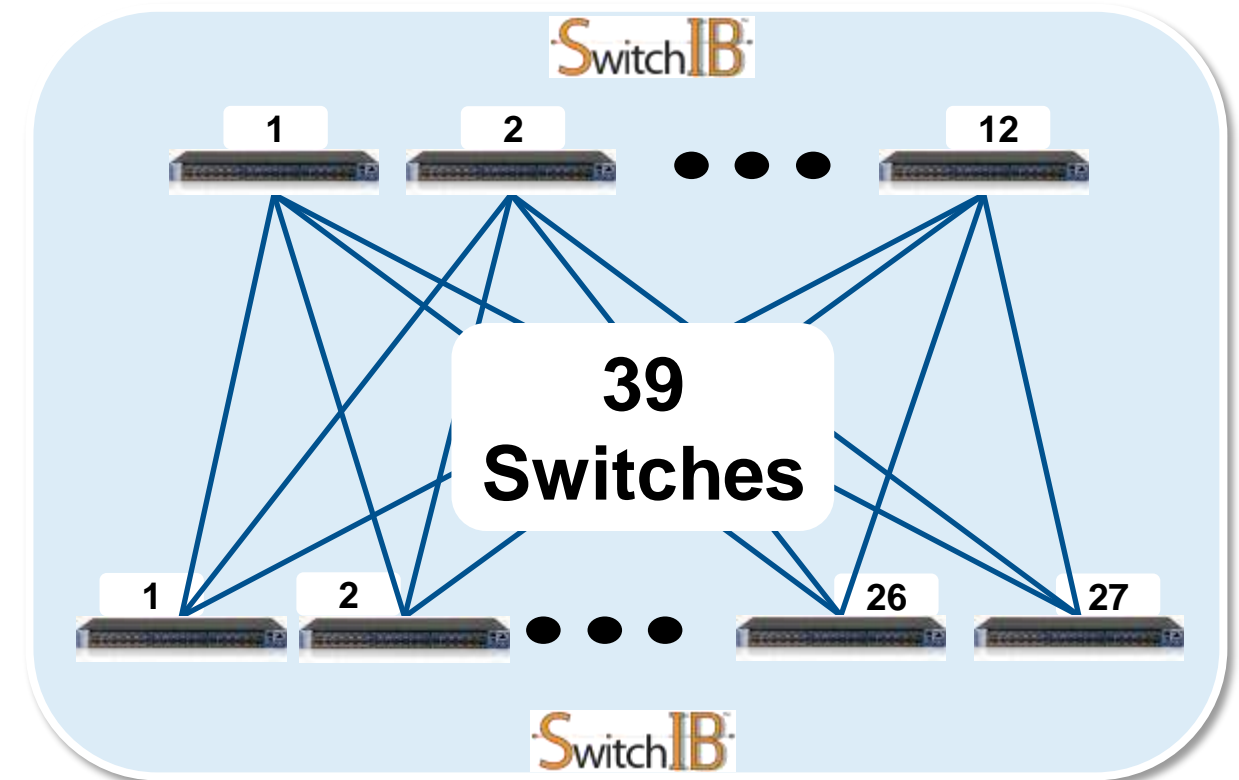


# Take Advantage of EDR Aggregation for FDR Clusters

648-node FDR cluster –1:1



648-node EDR cluster – 2:1 (FDR2EDR)



## EDR Network Aggregation Improves Cost-Performance

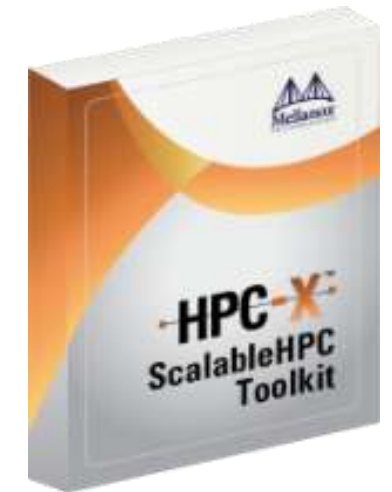
- Future proof
- Lower latency
- Less real estate (less switches, less cables)
- Wider pipes reduces congestion

NEW

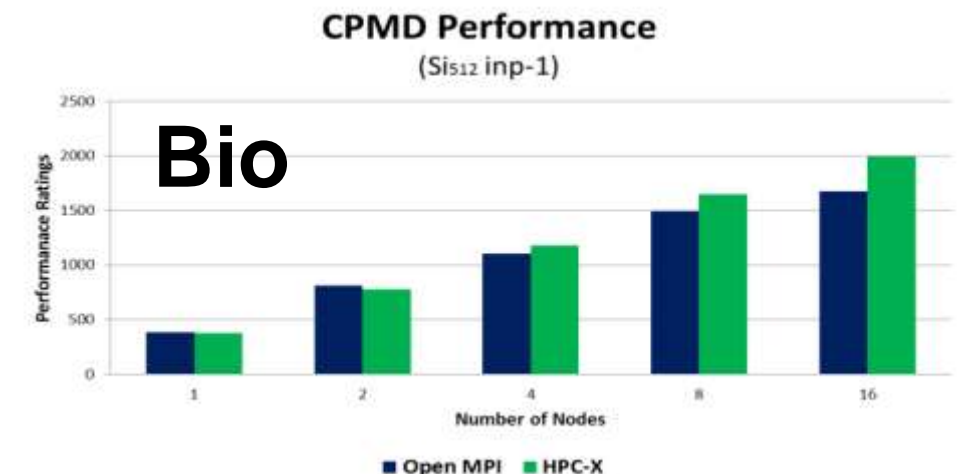
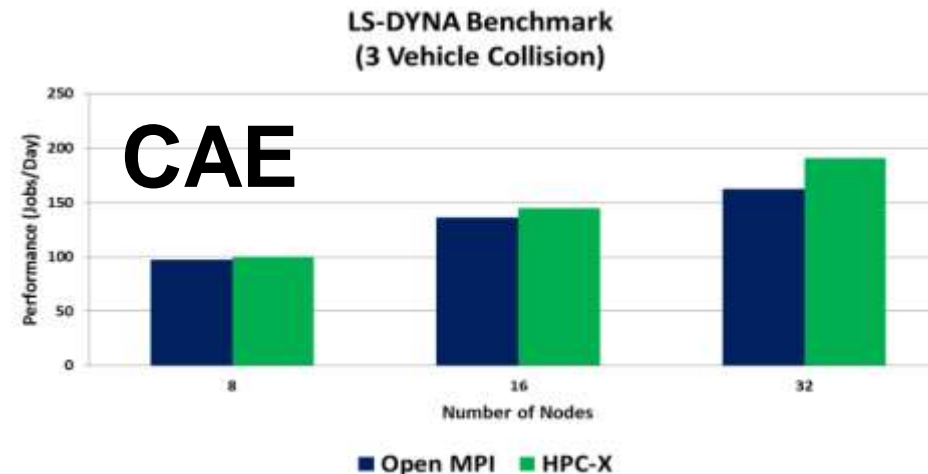
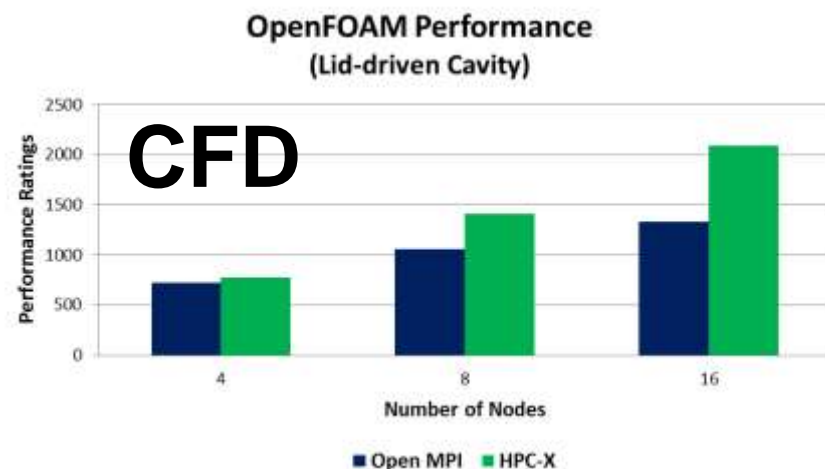
• HPC-X™ •



- Complete MPI, PGAS/OpenSHMEM/UPC package for HPC environments
  - Fully optimized for Mellanox InfiniBand and VPI interconnect solutions
  - Supports 3<sup>rd</sup> party solutions
  
- Components
  - Communication libraries: ScalableMPI, ScalableSHMEM, ScalableUPC
  - Acceleration libraries: MXM – Messaging Accelerator, FCA – Fabric Collectives Accelerator
  - Tools: Integrated Performance Monitoring Tool (IPM), benchmarks etc.



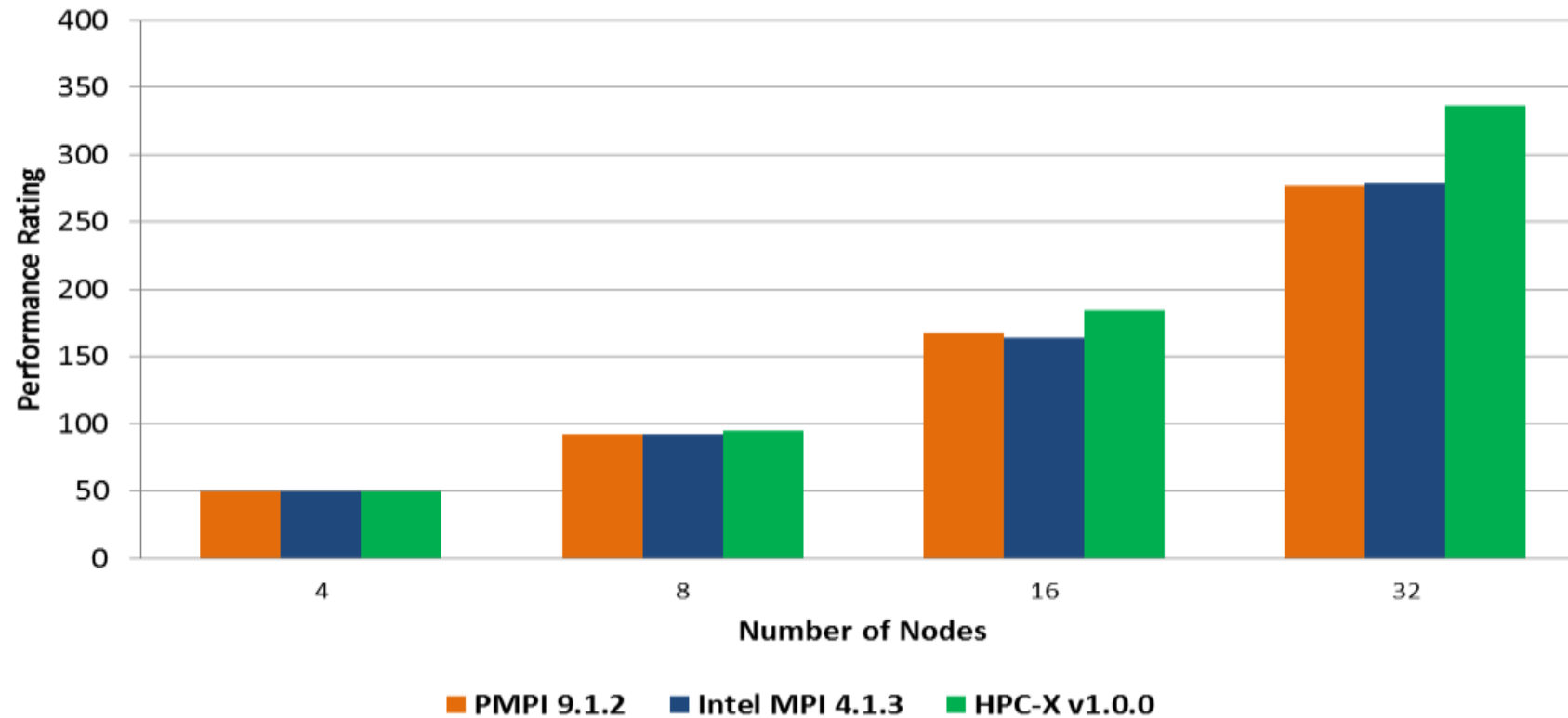
## 20%-70% Performance Improvement



**21% Performance Advantage!**



### STAR-CCM+ Performance

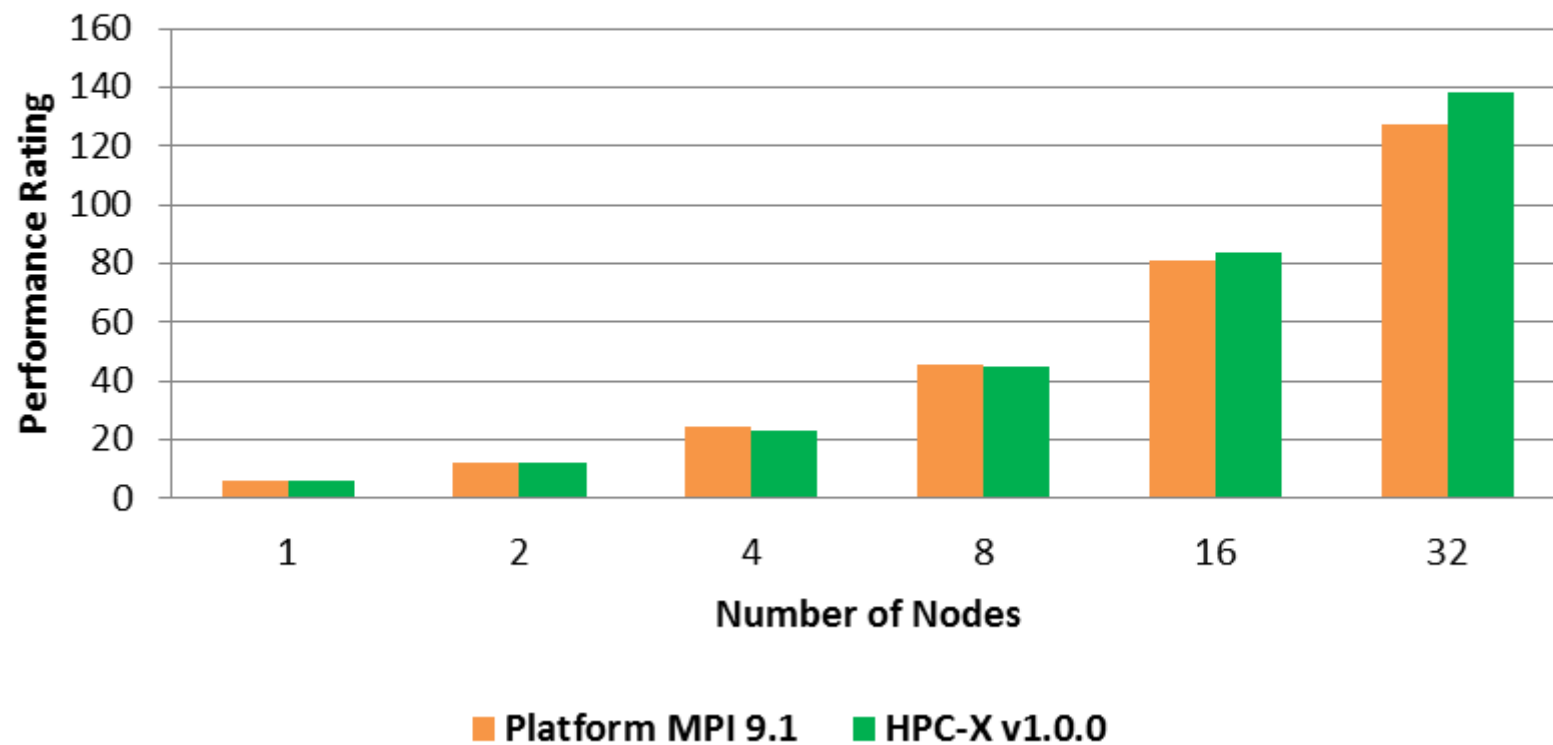


## HPC-X™ Performance for Commercial CFD



**8% Performance Advantage!**

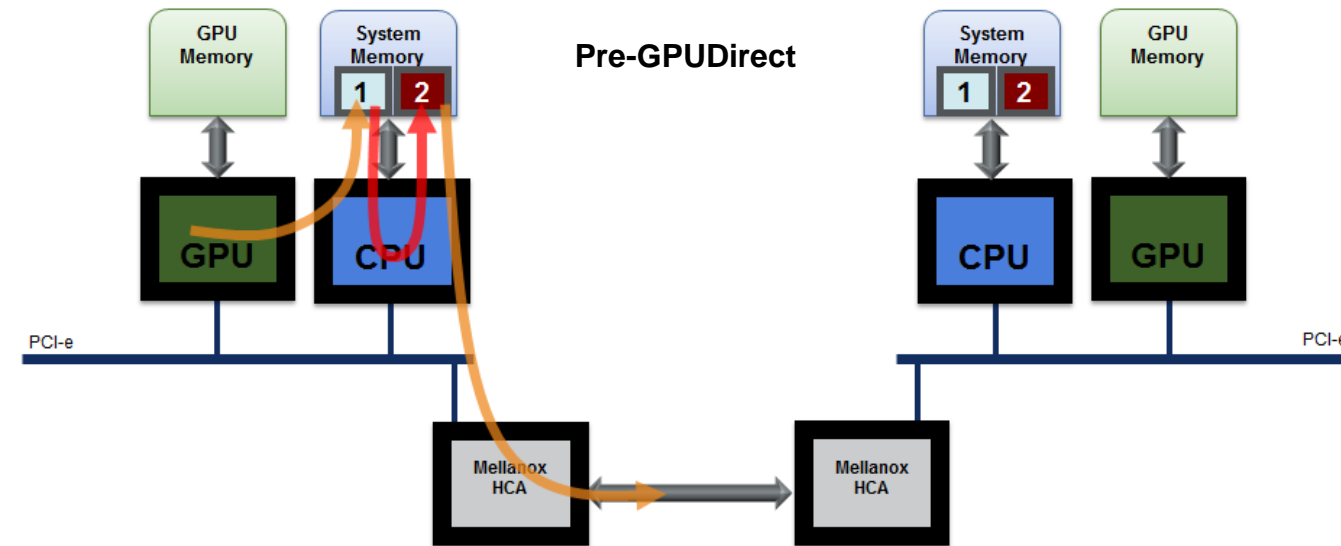
### VPS 2013.01 Performance (NEON\_FINE\_CAR2CAR)



## HPC-X™ Performance for Commercial FEA

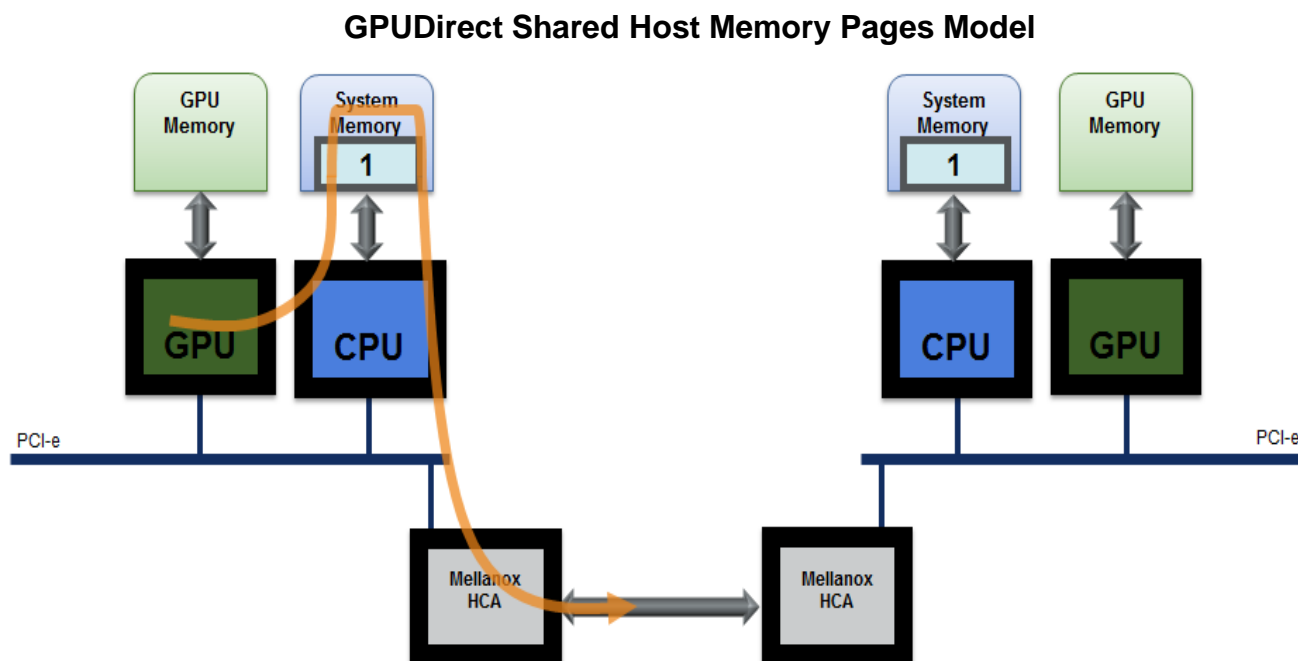
# GPUDirect RDMA

Native support for peer-to-peer communications  
between Mellanox HCA adapters and NVIDIA GPU devices



## Before GPUDirect

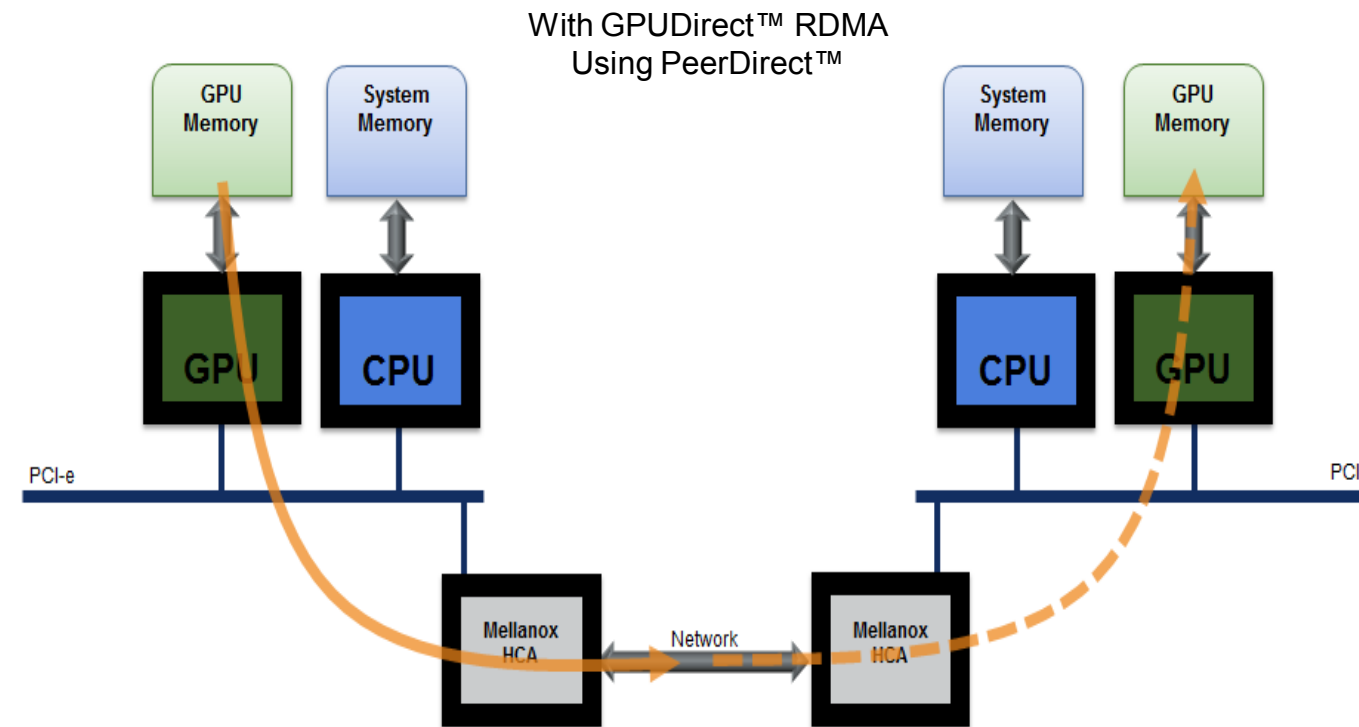
- Network and third-party device drivers, did not share buffers, and needed to make a redundant copy in host memory.



## With GPUDirect Shared Host Memory Pages

- The network and GPU can share “pinned” (page-locked) buffers, eliminating the need to make a redundant copy in host memory.



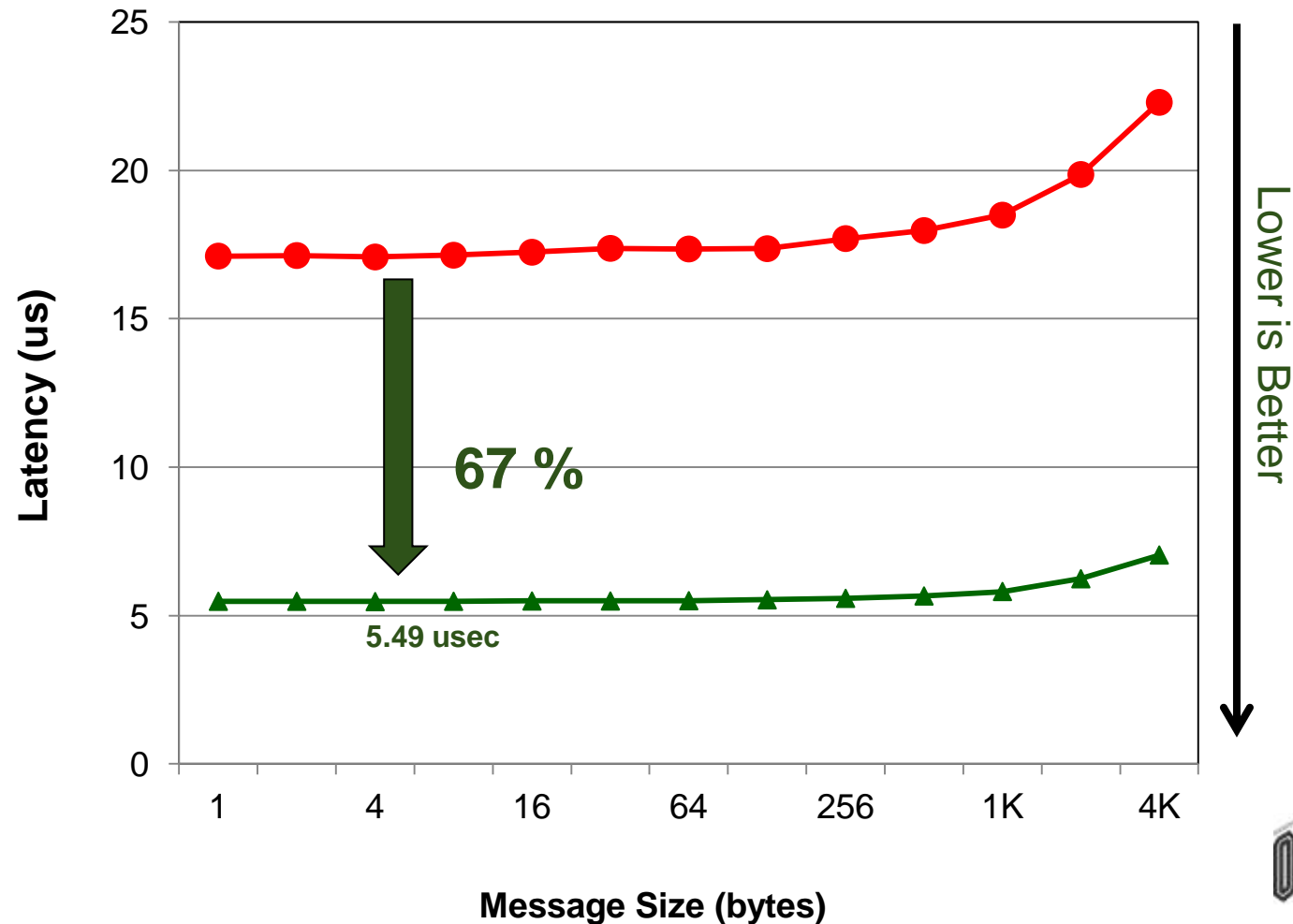


- Eliminates CPU bandwidth and latency bottlenecks
- Uses remote direct memory access (RDMA) transfers between GPUs
- Resulting in significantly improved MPI SendRecv efficiency between GPUs in remote nodes

# Performance of MVAPICH2 with GPUDirect RDMA



## GPU-GPU Internode MPI Latency

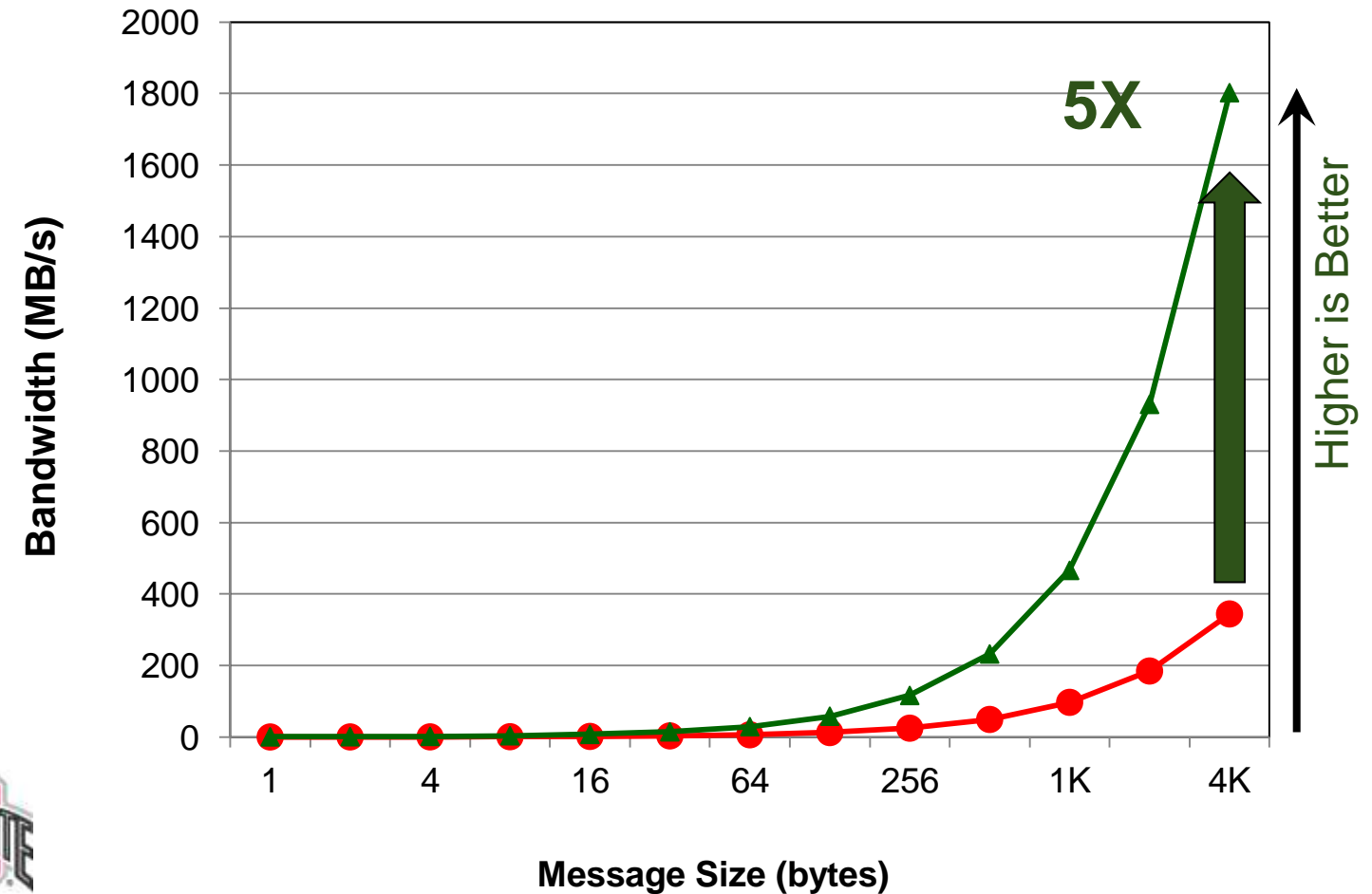


Lower is Better



Source: Prof. DK Panda

## GPU-GPU Internode MPI Bandwidth



Higher is Better

67% Lower Latency

5X Increase in Throughput

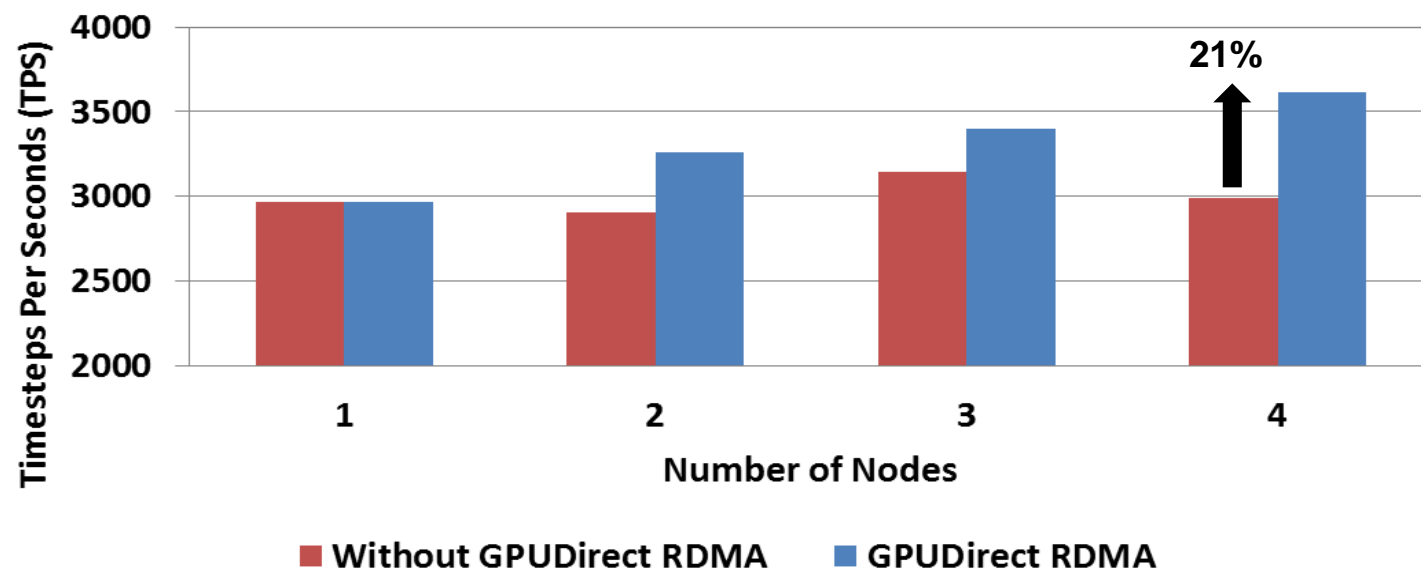
# Mellanox PeerDirect™ with NVIDIA GPUDirect RDMA



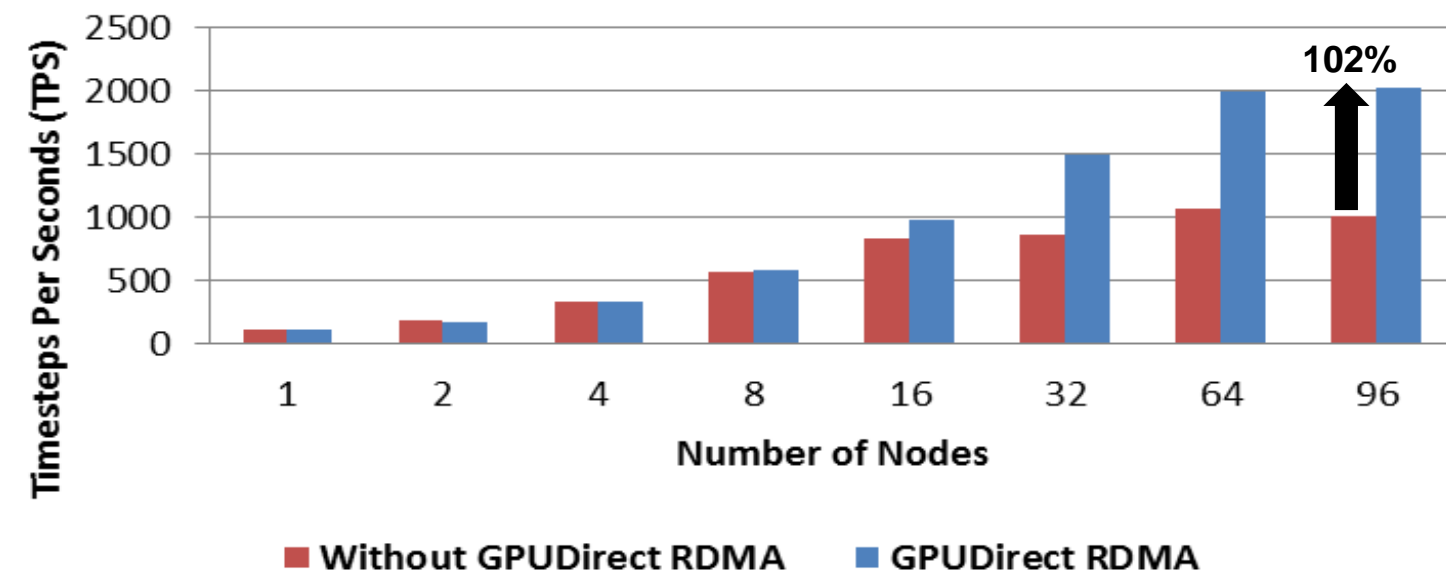
- HOOMD-blue is a general-purpose Molecular Dynamics simulation code accelerated on GPUs
- GPUDirect RDMA allows direct peer to peer GPU communications over InfiniBand
  - Unlocks performance between GPU and InfiniBand
  - This provides a significant decrease in GPU-GPU communication latency
  - Provides complete CPU offload from all GPU communications across the network
- Demonstrated up to 102% performance improvement with large number of particles



## HOOMD-blue Performance (LJ Liquid Benchmark, 16K Particles)



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





Questions?



Thank You