

# Mellanox Technologies

## Maximize Cluster Performance and Productivity

Gilad Shainer, [shainer@mellanox.com](mailto:shainer@mellanox.com)

October, 2007

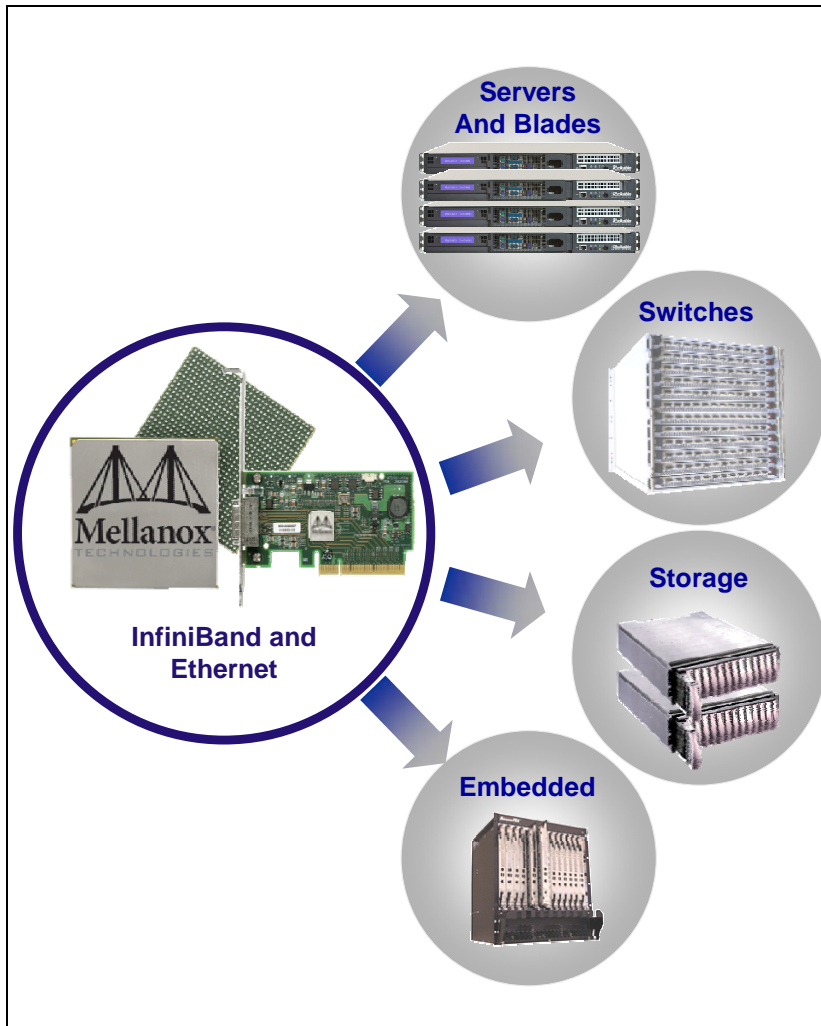
**WOLFRAM**  
**TECHNOLOGY**  
**CONFERENCE 2007**  
OCTOBER 11-13 | CHAMPAIGN, ILLINOIS



# Mellanox Technologies



## Hardware OEMs



## Applications

Logos of application providers:

- Schlumberger
- LSTC (Livermore Software Technology Corp.)
- ESI GROUP
- ANSYS
- FLUENT
- DREAMWORKS ANIMATION SKG
- WOLFRAMRESEARCH (MAKERS OF MATHEMATICA)
- Autodesk
- SYNOPSYS
- vmware
- TIBCO (The Power of Now)

## End-Users

### Enterprise Data Centers

Logos of Enterprise Data Center users:

- EDS
- JPMorgan
- Burlington
- KELLY SERVICES
- Prudential Financial
- BT
- myspace.com (a place for friends)
- Canon
- Sports Illustrated

### High-Performance Computing

Logos of High-Performance Computing users:

- cea
- HONDA
- TOYOTA
- Chevron
- Los Alamos NATIONAL LABORATORY
- Sandia National Laboratories
- THE UNIVERSITY OF CHICAGO
- NCSA
- UNIVERSITY OF SOUTH ALABAMA
- QUSITE
- السعودية العربية Saudi Aramco
- VW

### Embedded

Logos of Embedded users:

- CABLEVISION
- Comcast
- COX COMMUNICATIONS
- TimeWarner

**Interconnect: A Competitive Advantage**

Mellanox Technologies

# Connecting The Most Powerful Clusters



1280 server nodes

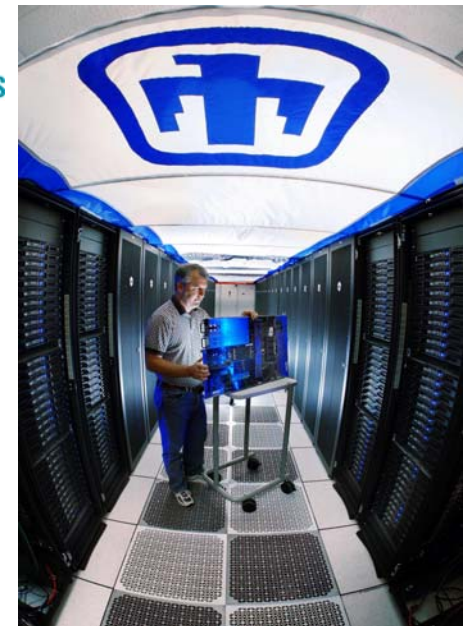


1300 server nodes



4500 server nodes

96 server nodes



2300 server nodes



1400 server nodes



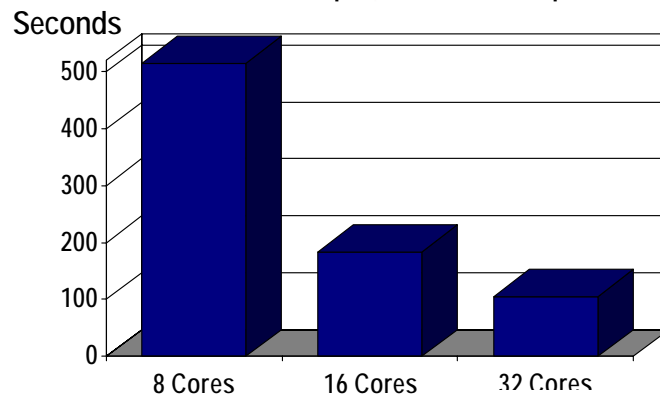
# But Not Only The Most Powerful Clusters



- Personal supercomputing (4-5 nodes)
  - Wolfram Air Pollution Simulation
  - Maximum utilization and efficiency



Time to Compute 10 Time Steps



TYAN



HP



- Sikorsky CH-53K program
  - 24 nodes, dual core CPUs
  - Reducing simulations duration from 4 days to several hours

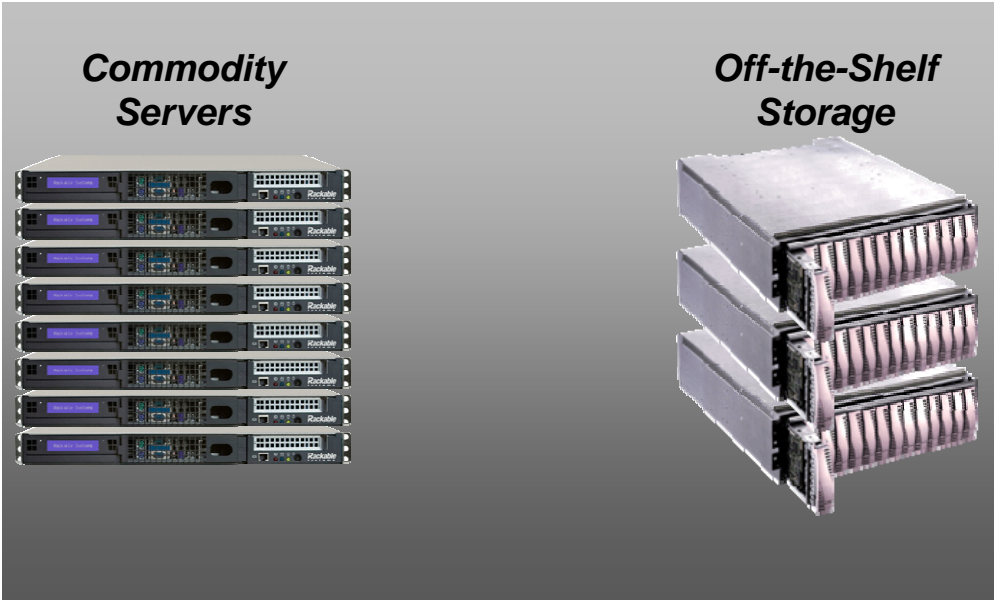


CH-53 HEAVY LIFT

# InfiniBand For Clustering



*Industry  
Megatrend*



## Proprietary Systems

- Expensive
- Not flexible

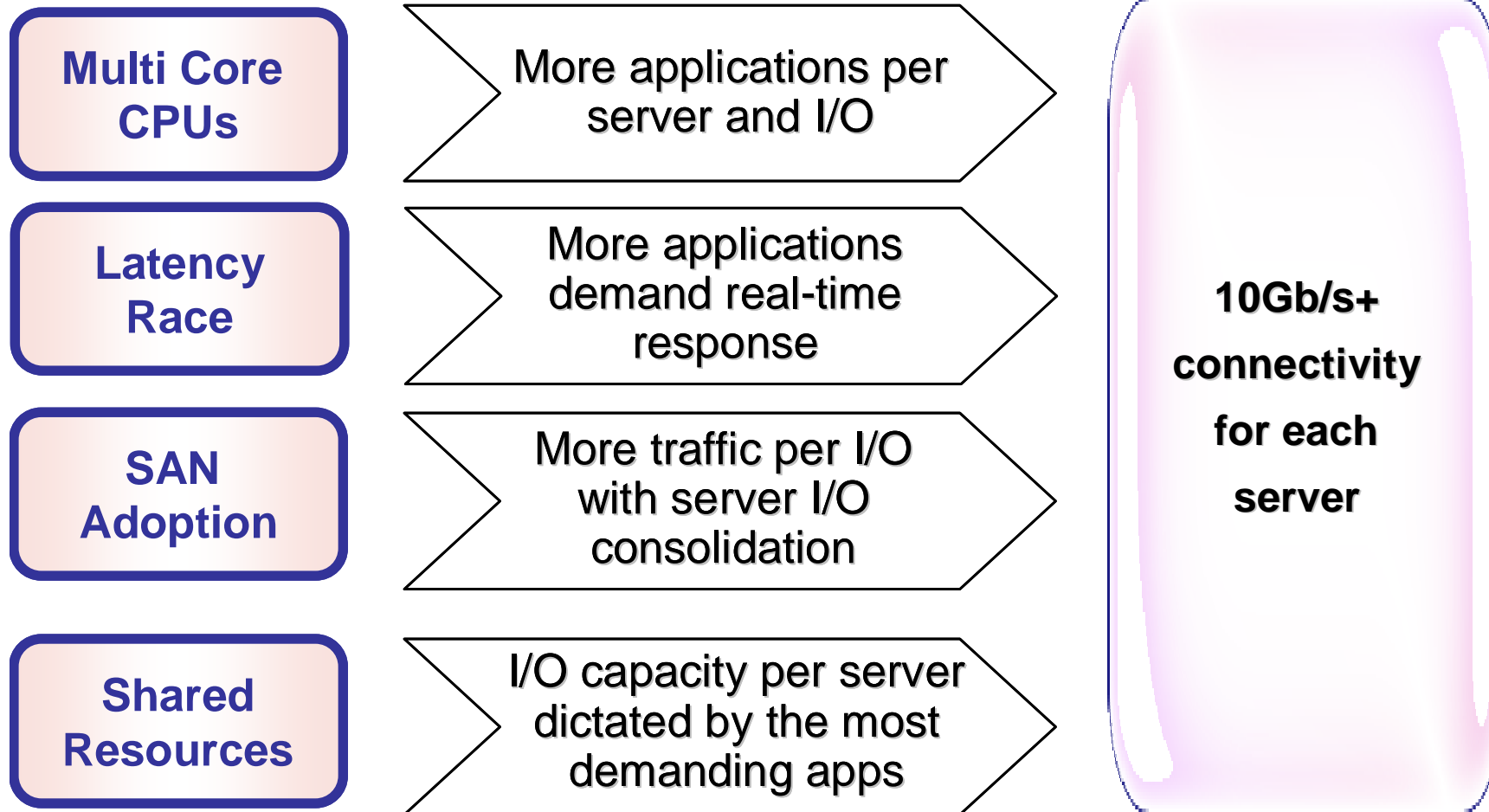
## Clusters

- Commodity
- Very flexible

## InfiniBand Clusters

- Maximum performance
- Scalability, large-scale clusters
- Multiple I/O traffics
- Flexible and easy to manage

# InfiniBand For I/O Growth Demands

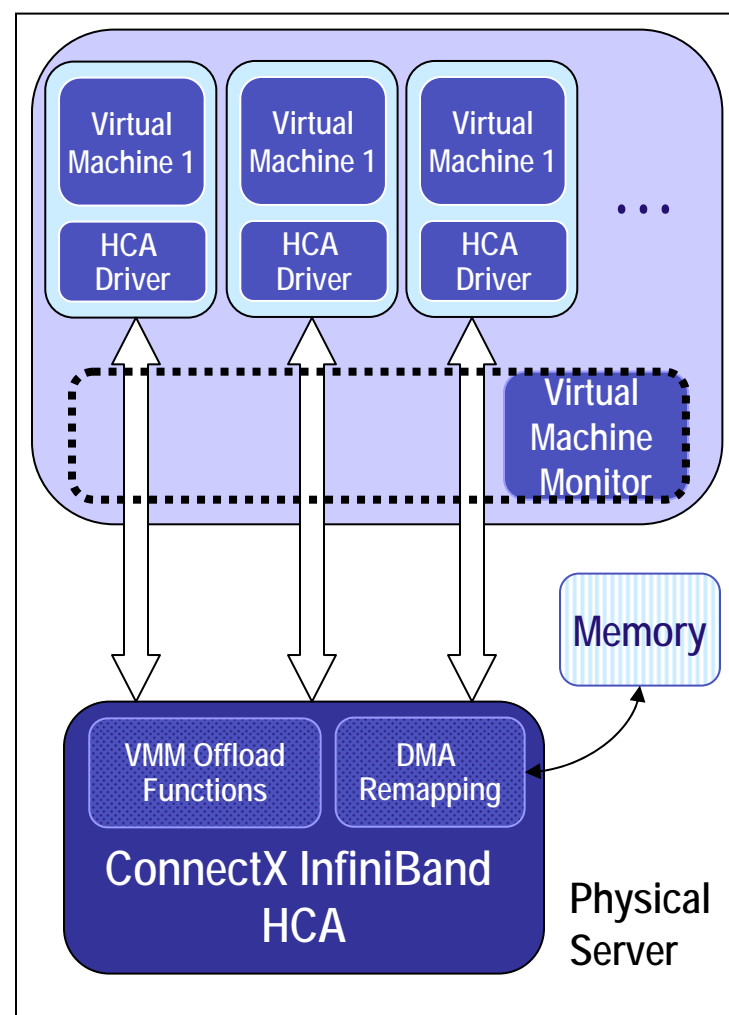


**Multi-core CPUs mandating 10Gb/s+ connectivity**

# InfiniBand For Virtualization



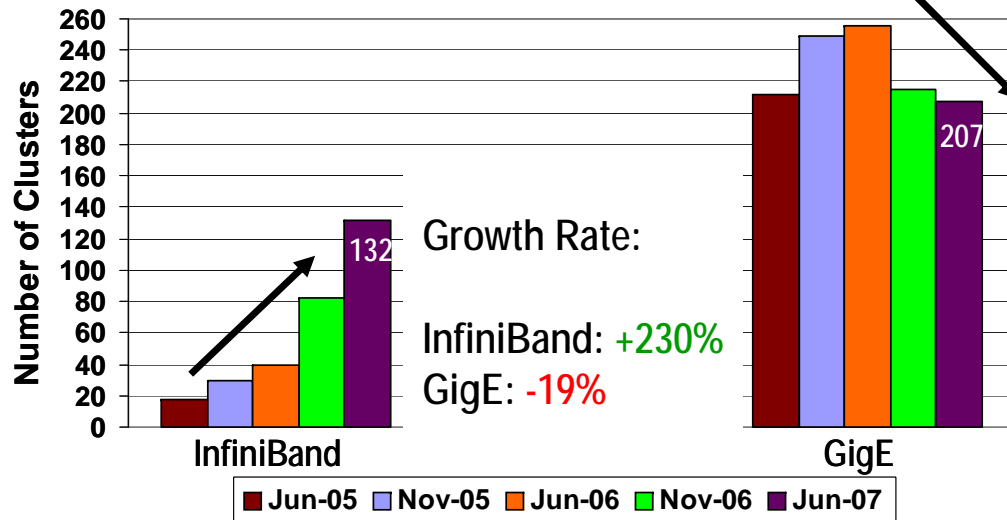
- **Hardware-based I/O virtualization**
  - Multiple VMs, multiple traffic types per VM
- **Supports current and future servers**
  - AMD and Intel IOV, PCI-SIG IOV
- **Better resource utilization**
  - Frees up CPU through hypervisor offload
  - Enable significantly more VMs per CPU
- **Native OS performance**
  - VMs enjoy native InfiniBand performance



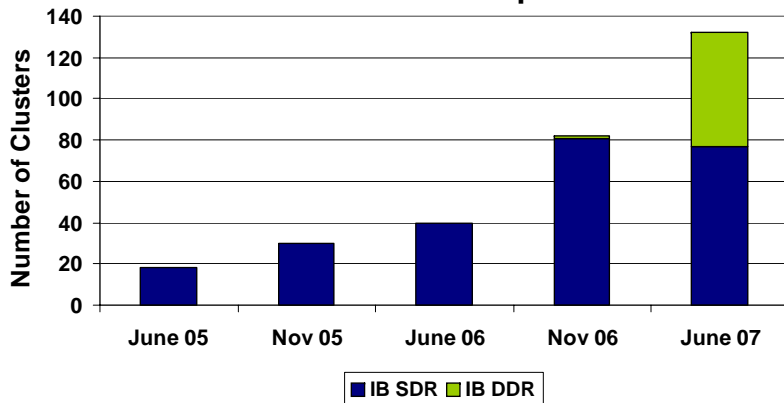
# TOP500 Interconnect Trends



## Top500 Interconnect Trends



### InfiniBand In The Top500

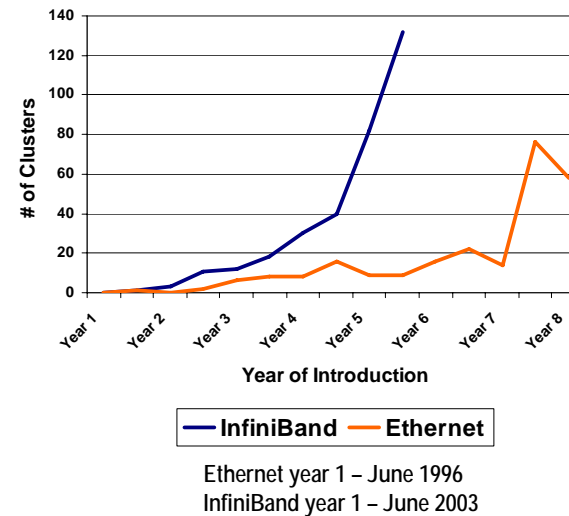


- Explosive growth of InfiniBand 20Gb/s
  - 55 clusters = 42% of the InfiniBand clusters



InfiniBand adoption is faster than Ethernet

### Top500 Interconnect Penetration



LANL



SANDIA



AERONAUTICAL SYSTEMS CENTER





# InfiniBand Value Proposition



## Automotive



- 3X simulation efficiency increase
- Car crash simulations eliminates the need for physical testing

## Oil and Gas



- Reduce reservoir modeling simulation runtime up to 55%
- Improves interpretation and modeling accuracy

## Fluid Dynamics



- 3X performance improvement and near linear scaling
- Intensive simulation for computer-aided engineering

## Digital Media



- 5X the bandwidth for real-time color grading
- High-resolution commercials and feature films

## Electronic Design Automation



- 4X improvement for photomask manufacturing - CATS

## Computational Science



- Monte Carlo simulation, astronomy, bioinformatics, chemistry and drug research
- Accelerate parallel execution of matrix operations

\* Partial List

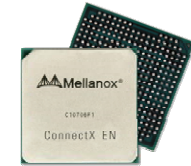
# Mellanox Product Roadmap



## 4th Generation

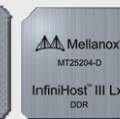
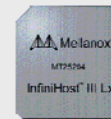
ConnectX™

Two 10/20/40 Gb/s InfiniBand or Two 1/10Gb/s Ethernet



PCI EXPRESS™  
2.0  
Adapter

## 3rd Generation



One 10/20Gb/s

PCI EXPRESS™  
Adapter



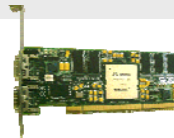
Two 10/20Gb/s IB

PCI EXPRESS™  
Adapter



480, 960Gb/s Total  
Switch

## 2nd Generation



Two 10Gb/s



Adapter



160Gb/s Total  
Switch

## 1st Gen.



40Gb/s Total



Adapter + Switch

2000

2001

2002

2003

2004

2005

2006

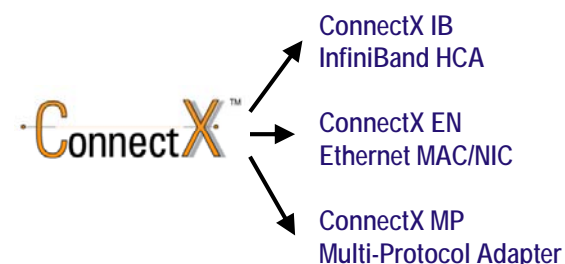
2007

\*Adapter Card Products Based on Adapter Silicon Not Shown

# Leading InfiniBand and 10GigE Adapters



- **Single-chip server and storage adapters**
  - Optimized cost, power, footprint, reliability
- **Highest performing InfiniBand adapters**
  - 20Gb/s (40Gb/s in 2008), 1us application latency
- **Highest performing 10GigE NICs**
  - 17.6Gb/s throughput, < 7us application latency
  - Supports OpenFabrics RDMA software stacks
- **Multi-core CPU optimized**
- **Virtualization acceleration**
- **Combination IB/Ethernet 4Q07**
- **First adapters to support PCI Express 2.0**

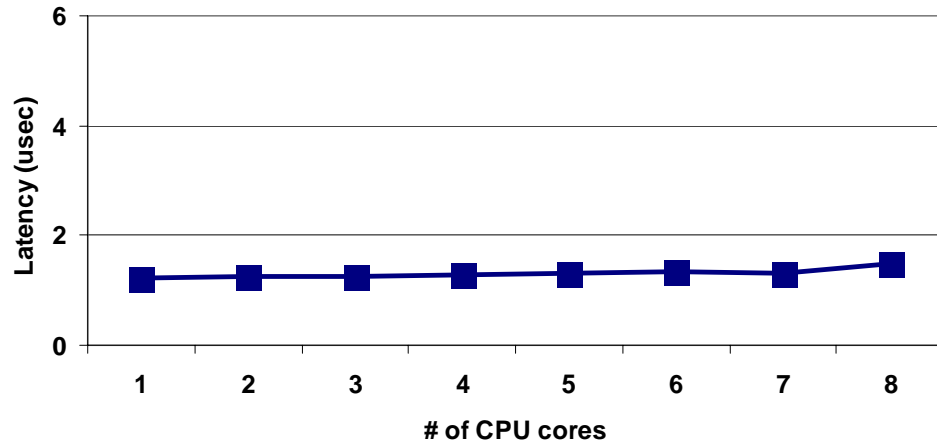


10 & 20Gb/s InfiniBand  
10 Gigabit Ethernet  
(Copper)

# ConnectX Multi Core Performance



### ConnectX MPI Latency - Multi-core Scaling

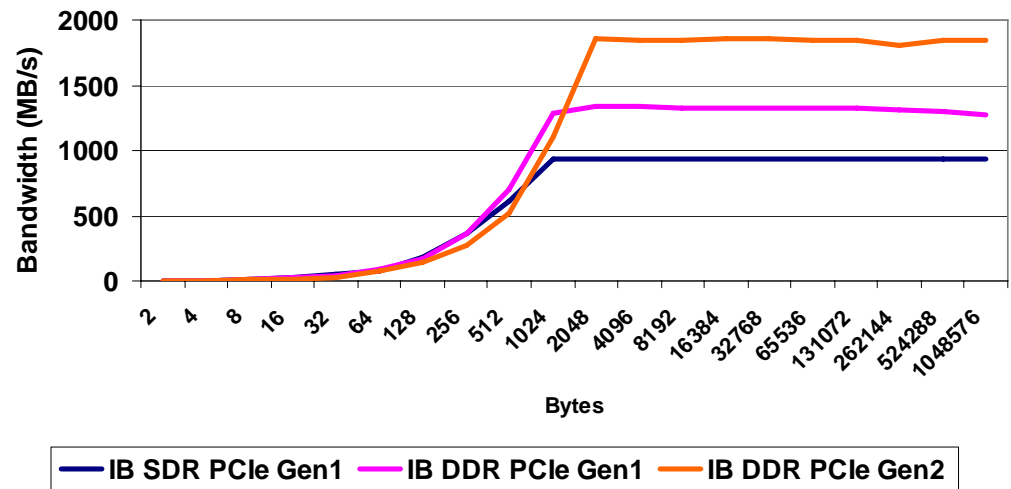


MPI (message passing)



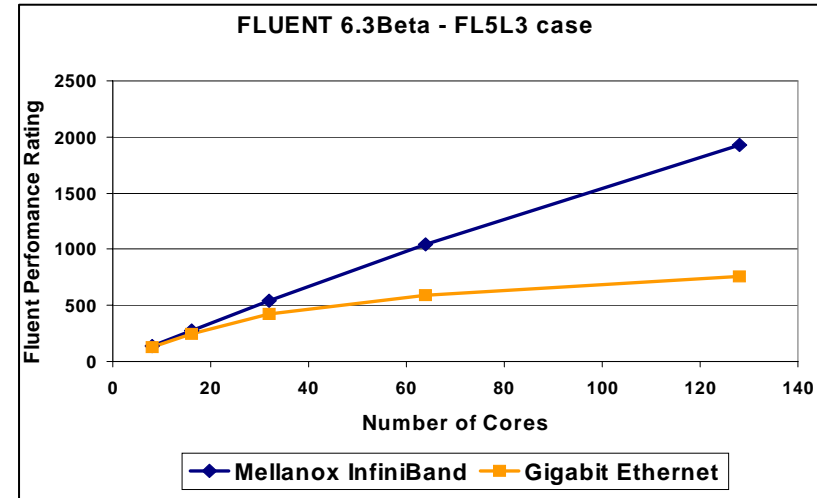
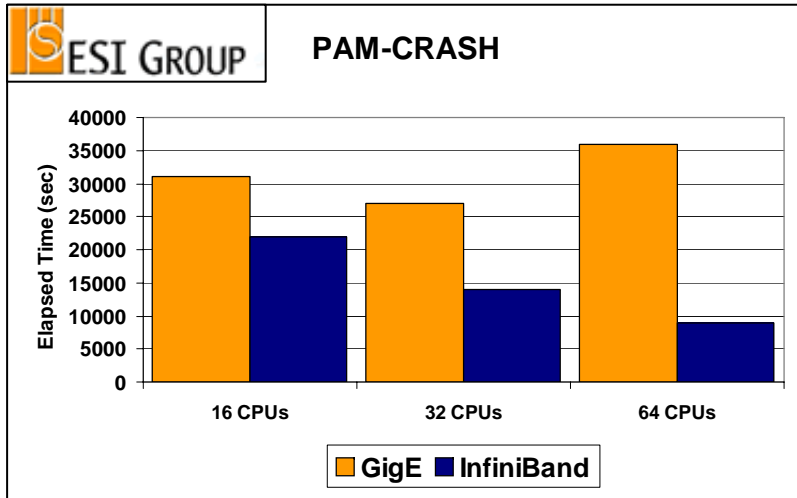
Sockets (TCP)

### IPoIB-CM ConnectX IB - SDR, DDR PCIe Gen1, DDR PCIe Gen2

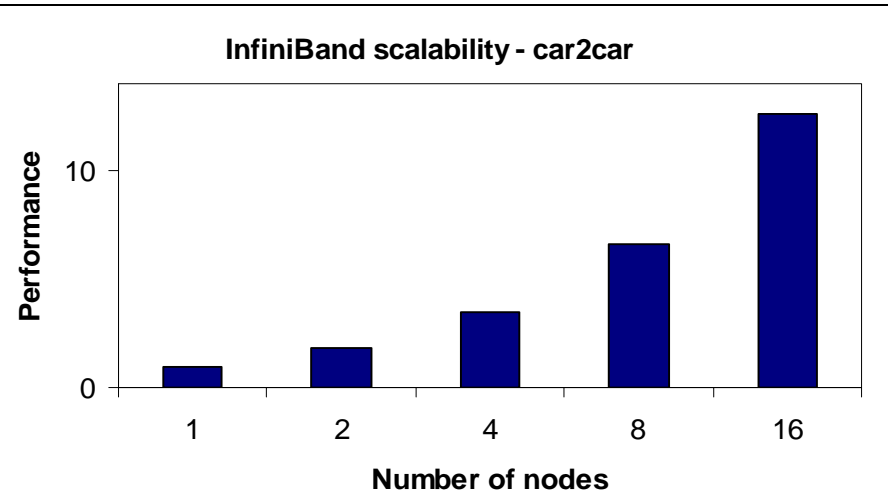


\*Optimizations on going

# Superior Application Productivity



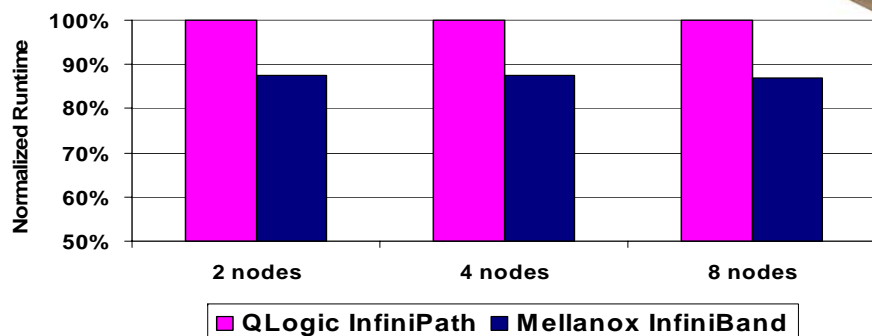
HP C-Class Blade System with Mellanox 20Gb/s InfiniBand I/O



# Superior Application Performance



**Mellanox InfiniBand versus QLogic InfiniPath - neon\_refined\_revised**



Lower is better

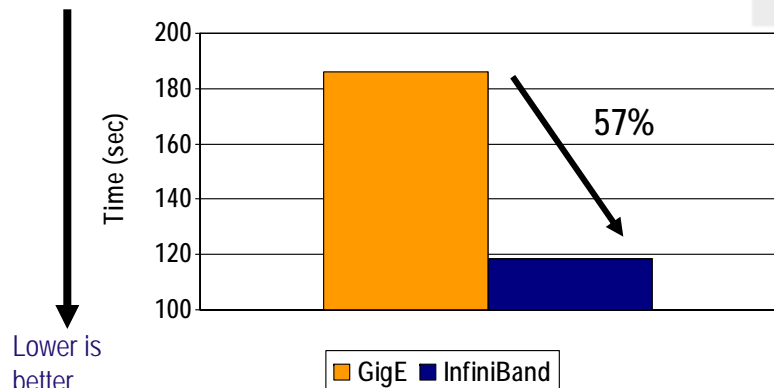
## Mellanox InfiniBand – full CPU offload

- Transport offload
- RDMA
- Robustness
- Scalability
- Efficiency

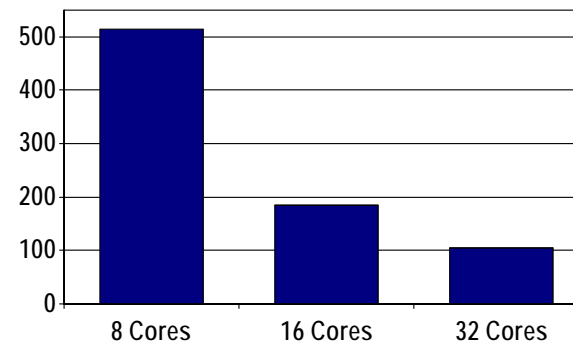
## CPU runs application, not network

Systems: Intel Woodcrest 3.0GHz  
Interconnects: Mellanox InfiniBand SDR, QLogic InfiniPath

**Vector Distribution Benchmark**



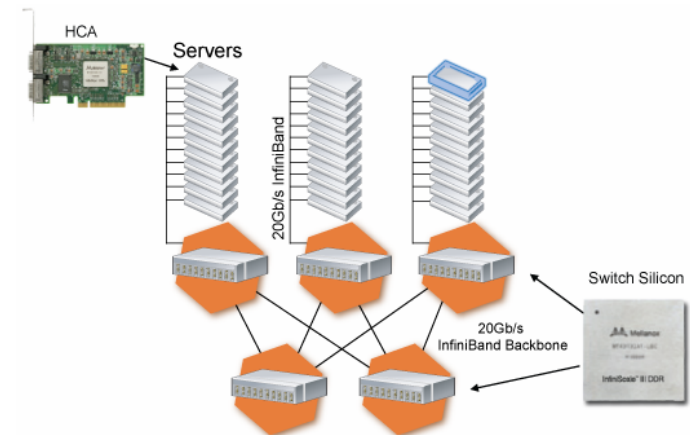
**Time to Compute 10 Time Steps**



# Mellanox Cluster Center



- <http://www.mellanox.com/applications/clustercenter.php>
- **Neptune cluster**
  - 32 nodes
  - Dual core AMD Opteron CPUs
- **Helios cluster**
  - 32 nodes
  - Quad core Intel Clovertown CPUs
- **Vulcan cluster – coming soon**
  - 32 nodes
  - Quad core AMD Barcelona CPUs
- **Utilizing “Fat Tree” network architecture (CBB)**
  - Non-blocking switch topology
  - Non-blocking bandwidth
- **ConnectX InfiniBand 20Gb/s**
- **InfiniBand based storage**
  - InfiniBand storage
  - NFS over RDMA, SRP

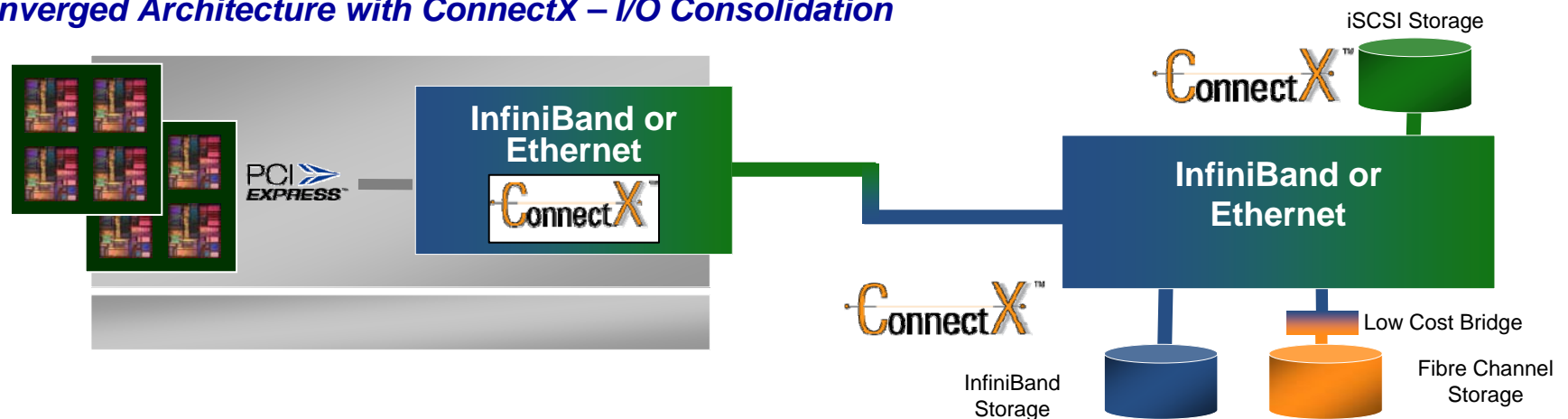


# Summary



- **Market-wide adoption of InfiniBand**
  - Servers/blades, storage and switch systems
  - Data Centers, High-Performance Computing, Embedded
  - Performance, Price, Power, Reliable, Efficient, Scalable
- **4<sup>th</sup> Generation adapter - connectivity to InfiniBand and Ethernet**
  - Market leading performance, capabilities and flexibility
- **Driving key trends in the market**
  - Clustering/blades, low-latency, I/O consolidation, multi-core, virtualization

## *Converged Architecture with ConnectX – I/O Consolidation*





Thank You

