



## Paving the Road to Exascale

Interconnecting Your Future

February 19, 2016 | PC Cluster Workshop

 **Mellanox**  
TECHNOLOGIES  
Connect. Accelerate. Outperform.™

本資料に関するお問い合わせ先



メラノックステクノロジーズジャパン株式会社

[japan\\_sales@mellanox.com](mailto:japan_sales@mellanox.com)

# Leading Supplier of End-to-End Interconnect Solutions



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

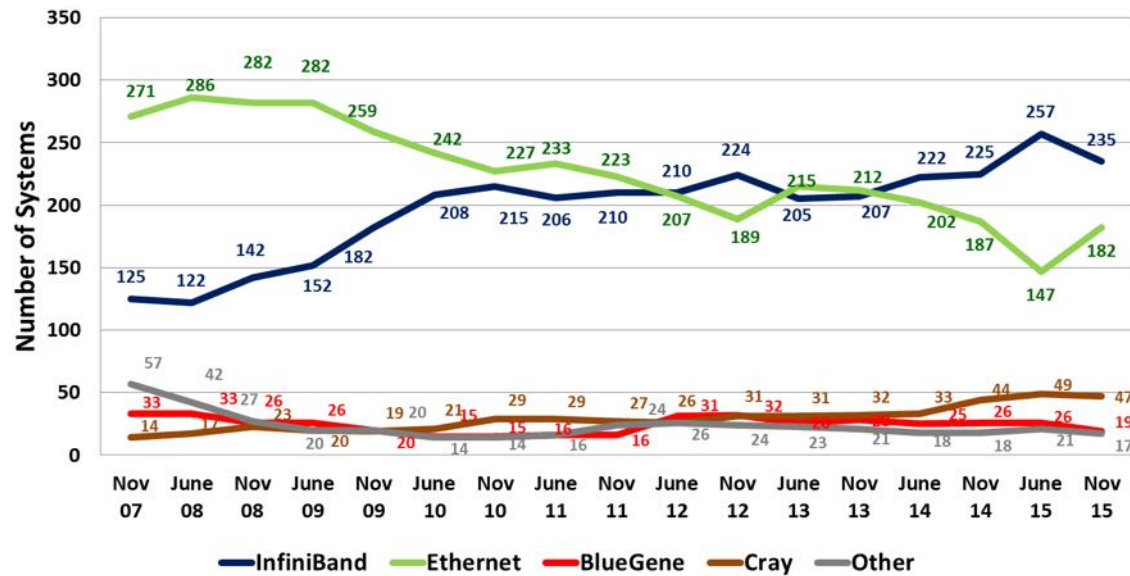
ICs	Adapter Cards	Switches/Gateways	Software and Services	Metro / WAN	Cables/Modules

At the Speeds of 10, 25, 40, 50, 56 and 100 Gigabit per Second

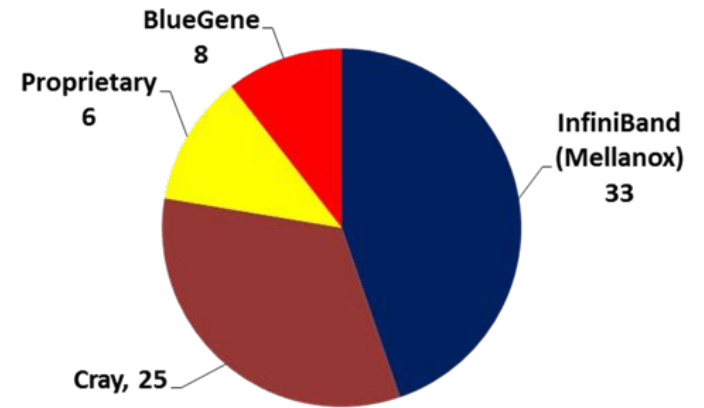
# TOP500 Interconnect Trends



TOP500 Interconnect Trends



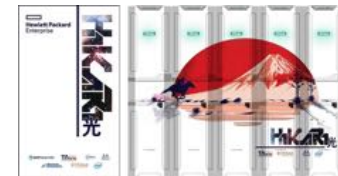
PetaFlops Systems on the TOP500 list



“LENOX“ EDR InfiniBand system  
Lenovo HPC innovation center



Shanghai Supercomputer Center  
Magic Cube II supercomputer



# InfiniBand Switch Portfolio



## Modular Switches



648 port

324 port

216 port

108 port

## Edge Switches



36 ports externally managed



18 port externally managed



8-12 ports externally managed



36 ports managed



18 port managed



12 port managed

## Long Distance



metroX™

## Bridge - VPI



SwitchX.2

## Management



# High-Performance Designed 100Gb/s Interconnect Solutions



<b>Adapters</b>		<b>100Gb/s Adapter, 0.7us latency</b> <b>150 million messages per second</b> <b>(10 / 25 / 40 / 50 / 56 / 100Gb/s)</b>		
<b>Switch</b>		<b>36 EDR (100Gb/s) Ports, &lt;90ns Latency</b> <b>Throughput of 7.2Tb/s</b> <b>7.02 Billion msg/sec (195M msg/sec/port)</b>		
<b>Switch</b>		<b>32 100GbE Ports, 64 25/50GbE Ports</b> <b>(10 / 25 / 40 / 50 / 100GbE)</b> <b>Throughput of 6.4Tb/s</b>		
<b>Interconnect</b>		<b>Transceivers</b> <b>Active Optical and Copper Cables</b> <b>(10 / 25 / 40 / 50 / 56 / 100Gb/s)</b>		<b>VCSELs, Silicon Photonics and Copper</b>

# Enter the Word of Scalable Performance – 100Gb/s Switch



## Best Performance

- 90ns switch latency
- Throughput of 7.2 Tb/s in 1U
- 195M messages per second
- 136W ATIS reported power

## Enhanced Capabilities

- InfiniBand Router
- Adaptive Routing (AR)
- Fault Routing (FR)



## High Resiliency

- High efficiency power supplies
- AC / DC / BBU power supplies option
- Class 4 (3.5W) supported on all ports

## x86 Powerful CPU

- Improved software upgrade time
- Up to 2048 nodes in cluster
- Run Virtual Machine (VM)

# Benefits of InfiniBand Router



- **Enable scaling an HPC cluster over 48K LIDs**
- **Enable sharing a common storage network by multiple disconnected subnets**
  - Limit congestion spread to source subnet
- **Allow running HPC/MPI jobs efficiently on the joint network**
  - Maintaining large bisectional bandwidth between the subnets
  - Low latency penalty for crossing subnets
- **Isolation of SM responsibilities**
- **Simple administration and out-of-the-box experience**



## ConnectX-4: Highest Performance Adapter in the Market

InfiniBand: SDR / DDR / QDR / FDR / EDR

Ethernet: 10 / 25 / 40 / 50 / 56 / 100GbE

100Gb/s, <0.7us latency

150 million messages per second

OpenPOWER CAPI Technology

CORE-Direct Technology

GPUDirect RDMA

Dynamically Connected Transport (DCT)

Ethernet Offloads (HDS, RSS, TSS, LRO, LSOv2)



# Shattering The World of Interconnect Performance!



## ConnectX-4 EDR 100G InfiniBand

<b>Uni-Directional Throughput</b>	<b>100 Gb/s</b>
<b>Bi-Directional Throughput</b>	<b>195 Gb/s</b>
<b>Latency</b>	<b>0.61 us</b>
<b>Message Rate</b>	<b>149.5 Million/sec</b>

# InfiniBand Adapters Performance Comparison



<b>Mellanox Adapters</b> <b>Single Port Performance</b>	<b>ConnectX-4</b> PCI Express 3.0 x16 <b>EDR 100Gb/s</b>	<b>Connect-IB</b> PCI Express 3.0 x16 <b>FDR 56Gb/s</b>	<b>ConnectX-3 Pro</b> PCI Express 3.0 x8 <b>FDR 56Gb/s</b>
<b>Uni-Directional Throughput</b>	<b>100 Gb/s</b>	<b>54.24 Gb/s</b>	<b>51.1 Gb/s</b>
<b>Bi-Directional Throughput</b>	<b>195 Gb/s</b>	<b>107.64 Gb/s</b>	<b>98.4 Gb/s</b>
<b>Latency</b>	<b>0.61 us</b>	<b>0.63 us</b>	<b>0.64 us</b>
<b>Message Rate</b>	<b>149.5 Million/sec</b>	<b>105 Million/sec</b>	<b>35.9 Million/sec</b>

# Mellanox QSFP 100Gb/s Cables



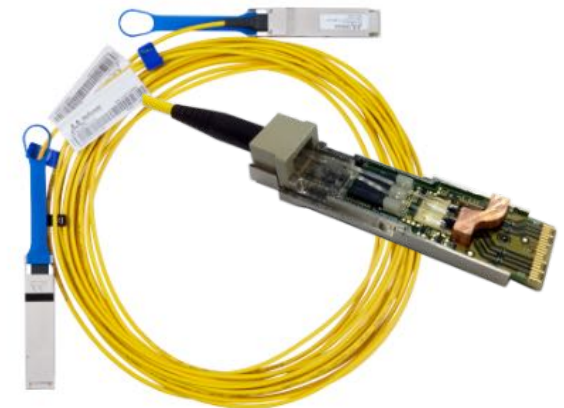
## Complete Solution of 100Gb/s Copper and Fiber Cables



**Copper Cables**



**VCSEL AOCs**



**Silicon Photonics AOCs**

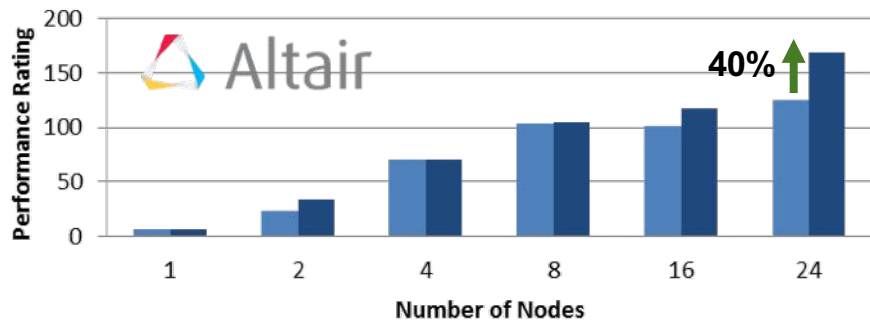


**Making 100Gb/s Deployments as Easy as 10Gb/s**

# EDR InfiniBand Performance Leadership

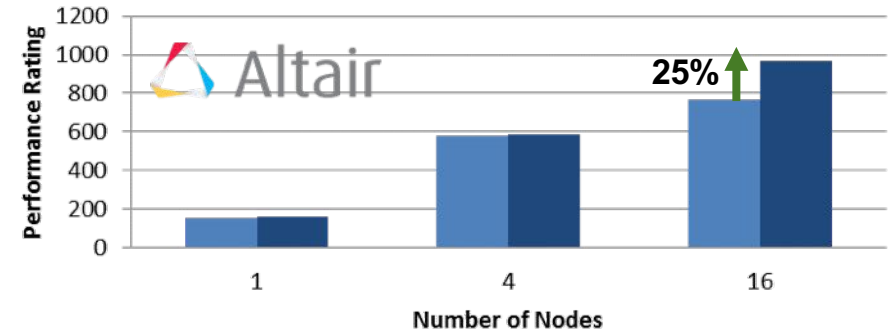


## OptiStruct Performance (Engine\_Assy.fem)



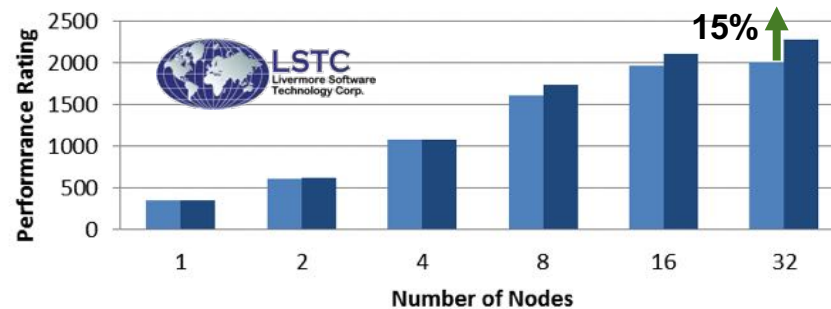
■ FDR InfiniBand ■ EDR InfiniBand

## RADIOSS 13.0 Performance (NEON1M11, MPP)



■ FDR InfiniBand ■ EDR InfiniBand

## LS-DYNA Performance (neon\_refined\_revised)



■ FDR InfiniBand ■ EDR InfiniBand



For all graphs: higher is better

# Introducing Switch-IB 2 World's First Smart Switch



Switch-IB™ 2

# Introducing Switch-IB 2 World's First Smart Switch



## Switch-IB™ 2

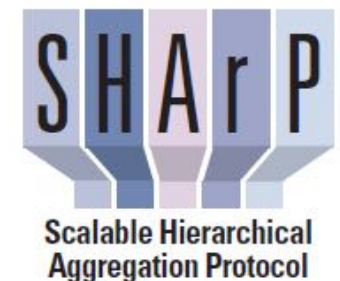


### World's First Smart Switch

### Build for Scalable Compute and Storage Infrastructures

## 10X Higher Performance with The New Switch SHArP Technology

- The world fastest switch with <90 nanosecond latency
- 36-ports, 100Gb/s per port, 7.2Tb/s throughput, 7.02 Billion messages/sec
- Adaptive Routing, Congestion Control, support for multiple topologies



# SHArP (Scalable Hierarchical Aggregation Protocol) Technology



SHArP Enables Switch-IB 2 to Manage and Execute MPI Operations in the Network

Switch-IB 2 Enables the Switch Network to Operate as a Co-Processor

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

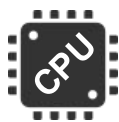




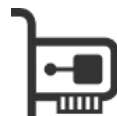
# The Intelligence is Moving to the Interconnect



## Communication Frameworks (MPI, SHMEM/PGAS)



Applications



Transport

RDMA

SR-IOV

Collectives

Peer-Direct

GPUDirect

More...



Scalable Hierarchical  
Aggregation Protocol

MPI / SHMEM Offloads



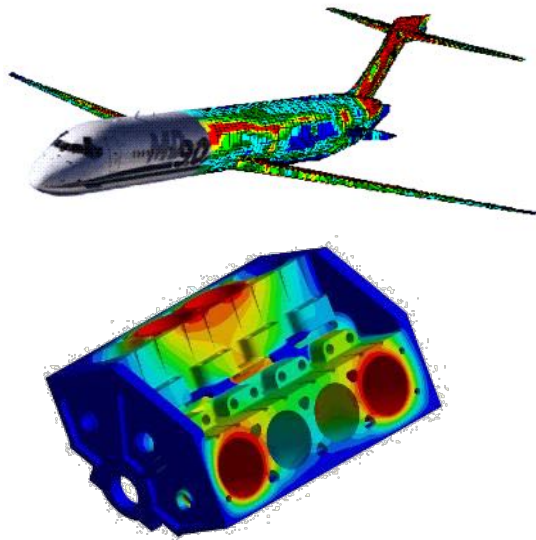
The Only Approach to Deliver **10X** Performance Improvements

# SHArP Performance Advantage – MiniFE

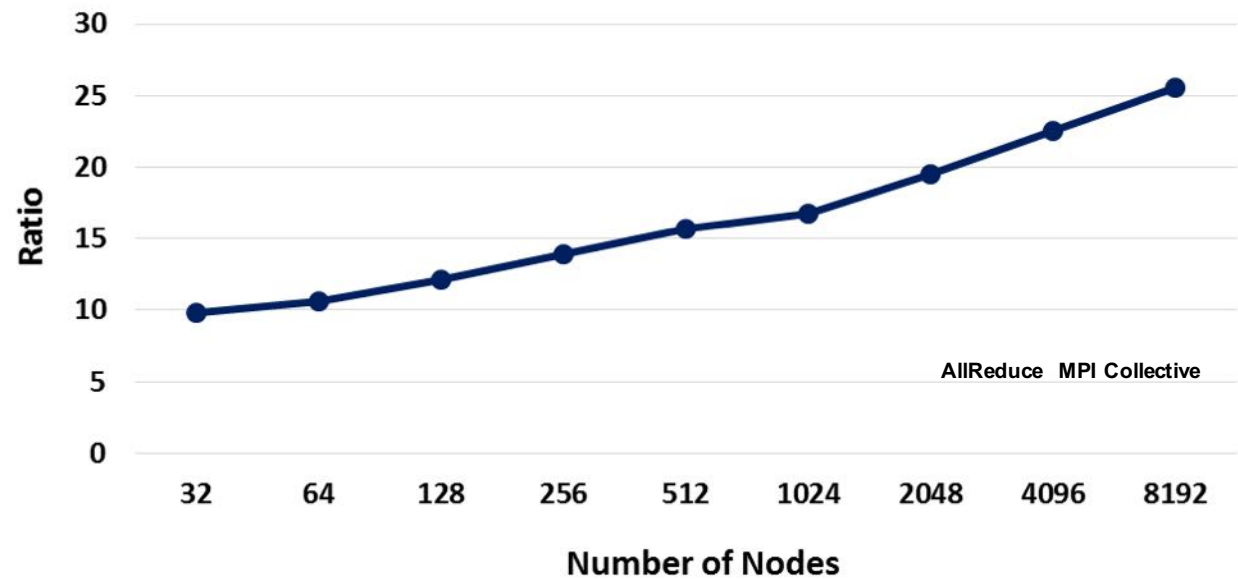


## MiniFE is a Finite Element mini-application

Implements kernels that represents implicit finite-element applications



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



**10X to 25X** Performance Improvement

# SHArP Performance Advantage – MiniFE Details

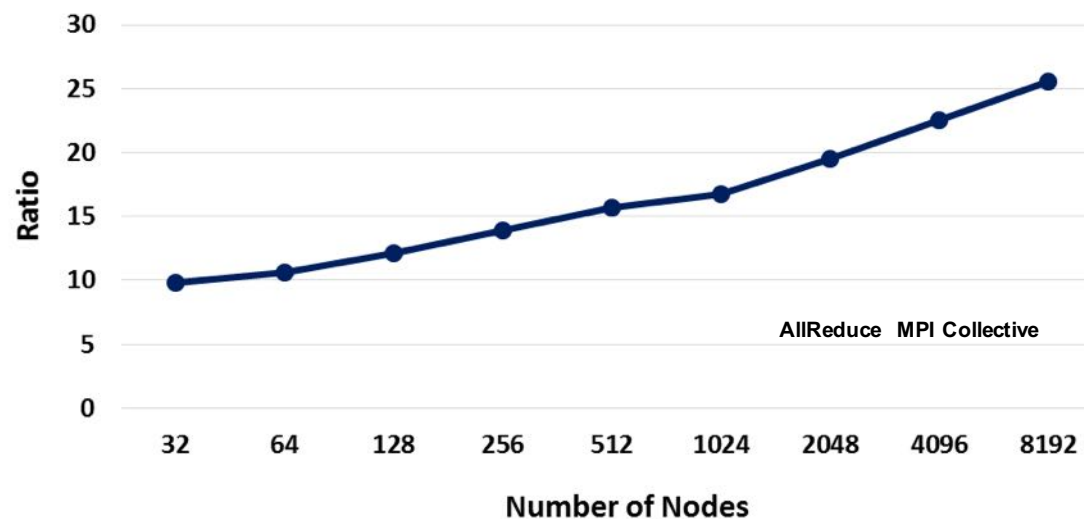


## MiniFE is a Finite Element mini-application

Implements kernels that represents implicit finite-element applications

Number of Nodes	CPU-based Latency (usec)	SHArP-based Latency (usec)	Ratio
32	41.7	4.24	9.9
64	49.08	4.63	10.6
128	57.67	4.76	12.1
256	67.76	4.87	13.9
512	79.62	5.09	15.6
1024	93.55	5.58	16.8
2048	109.92	5.63	19.5
4096	129.16	5.73	22.5
8192	151.76	5.94	25.5

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

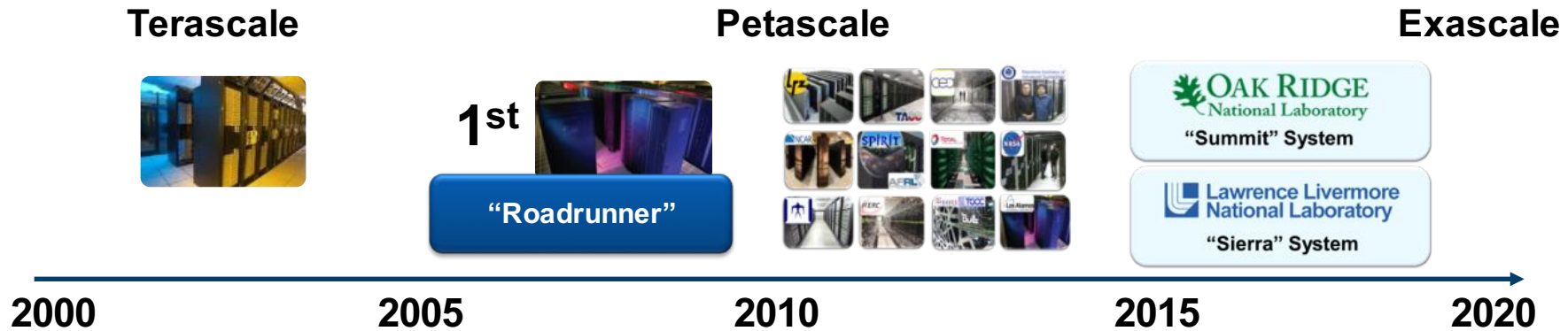


**10X to 25X** Performance Improvement

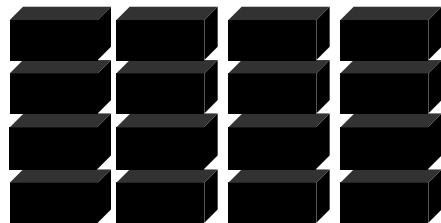
# The Ever Growing Demand for Higher Performance



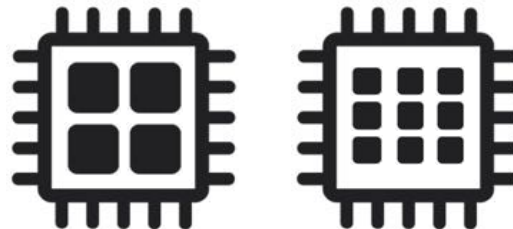
## Performance Development



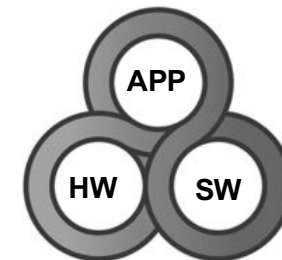
## The Interconnect is the Enabling Technology



SMP to Clusters



Single-Core to Many-Core



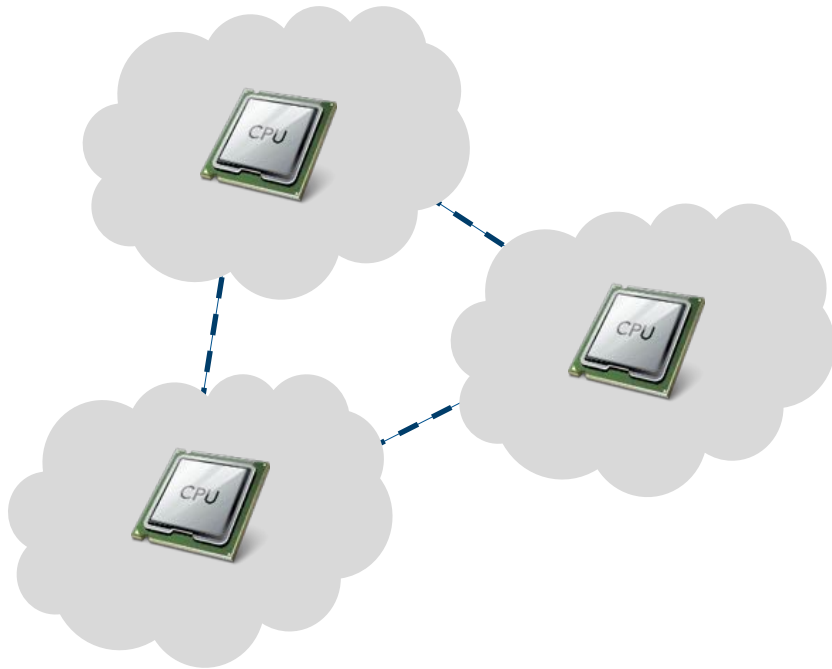
Application  
Software  
Hardware

Co-Design

# Co-Design Architecture to Enable Exascale Performance

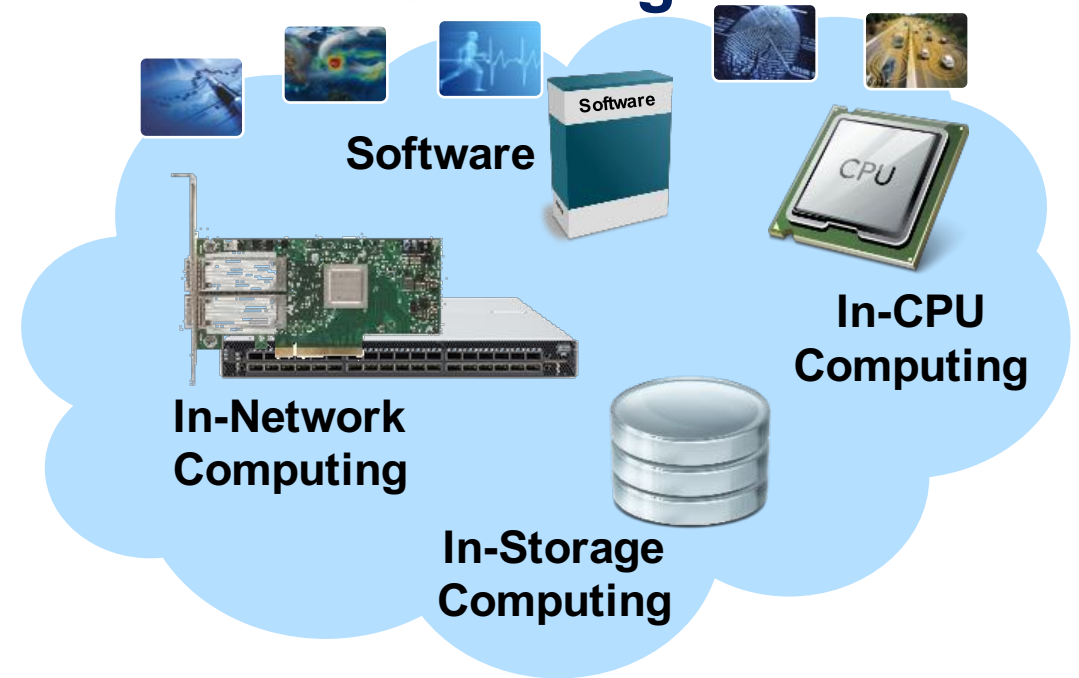


## CPU-Centric



**Limited to Main CPU Usage  
Results in Performance Limitation**

## Co-Design



**Creating Synergies  
Enables Higher Performance and Scale**

# Mellanox InfiniBand Solutions Deliver Highest ROI for Any Scale



**Smart Network For Smarter Systems**  
RDMA, Acceleration Engines, Programmability

Higher Performance



**100**  
Gb/s  
Link Speed



**200**  
Gb/s  
Link Speed

**Mellanox delivers the HIGHEST return on investment for ANY scale deployment!**

Protect Your Future

Power Consumption  
Per Switch Port

**25%**  
**Lower**

Message Rate

**44%**  
**Higher**

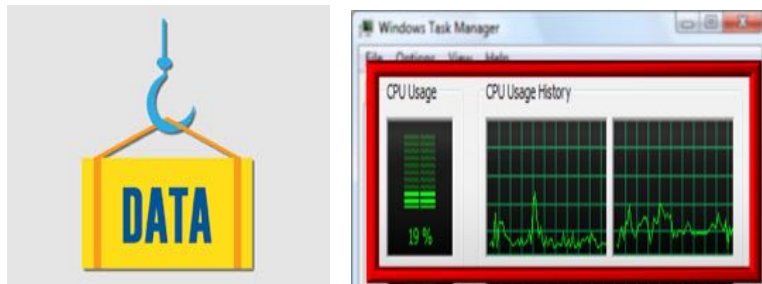
Switch Latency

**20%**  
**Lower**

Scalability  
CPU efficiency

**2X**  
**Higher**

## Offload Technology



## Proven Solutions



## Eco-System



## Standard



# Speed-Up Your Present, Protect Your Future

# Mellanox Delivers Best Interconnect



## Higher Performance

- 100Gb/s throughput at 0% CPU utilization
- Adapter: 150 Million messages/sec on today's systems, 44% higher
- Switch: 7.02 Billion messages/sec (195 Million per port)
- 20% lower switch latency, with deterministic latency!

## Lower TCO

- 25% lower power consumption per switch port
- Standard-based solutions, large eco-system support
- Backward and future compatibility – protect investments
- Offloading Architecture (RDMA, GPUDirect etc.) delivers highest system efficiency

## Higher Reliability

- 1,000X higher reliability - Mellanox delivers Bit Error Rate of  $10^{-15}$  versus  $10^{-12}$
- Superior signal integrity
- Support for Multiple data integrity mechanisms (FEC<sup>1</sup>, LLR<sup>2</sup>, COD<sup>3</sup> and more)

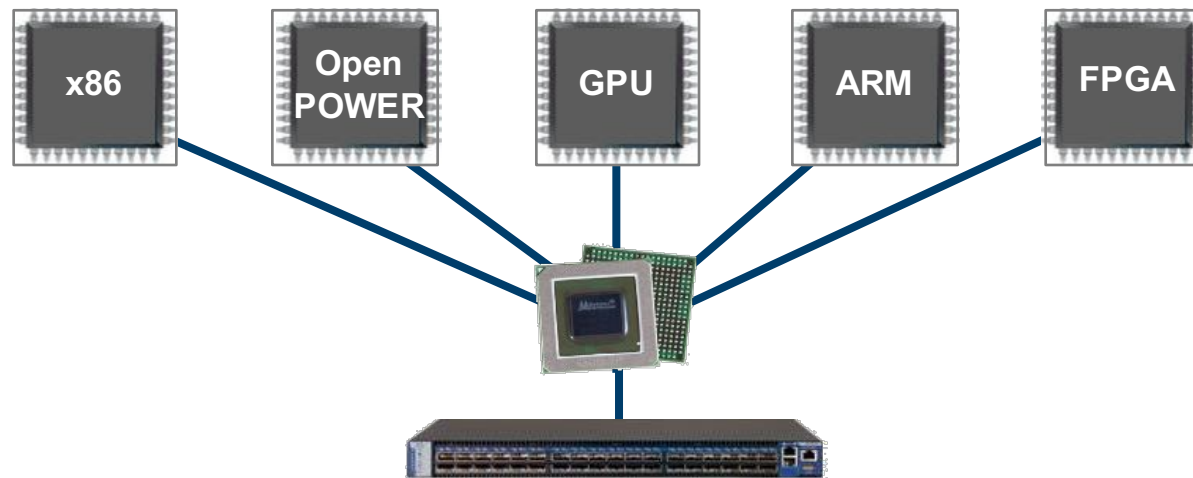
1 - Forward Error Correction; 2 - Link Level Retransmission; 3 - Correction on Demand



# End-to-End Interconnect Solutions for All Platforms



**Highest Performance and Scalability for  
x86, Power, GPU, ARM and FPGA-based Compute and Storage Platforms  
10, 20, 25, 40, 50, 56 and 100Gb/s Speeds**



**Smart Interconnect to Unleash The Power of All Compute Architectures**



Thank You

 **Mellanox**  
TECHNOLOGIES  
Connect. Accelerate. Outperform.™