# HPC Interconnect Technology update

Paving the Road to Exascale – HPC User Forum

2017

Mellanox TECHNOLOGIES

Connect. Accelerate. Outperform.™

# The Ever-Growing Demand for Higher Performance



## Performance Development

| Terascale | Petascale | | Exascale |
|---|---|---|---|

**1st** LANL

**1st** Wuxi

OAK RIDGE National Laboratory "Summit" System

Lawrence Livermore National Laboratory "Sierra" System

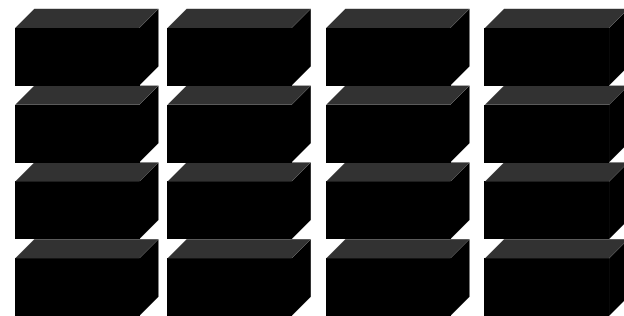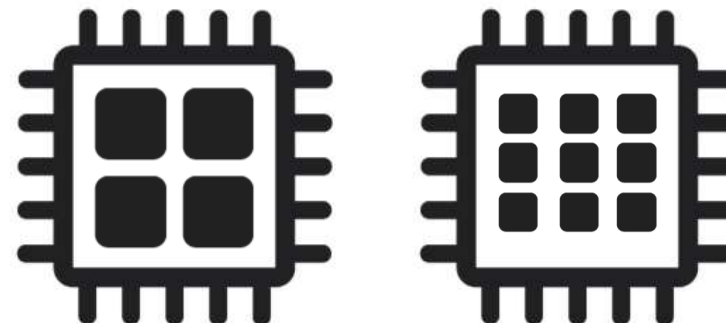2000        2005        2010        2015        2020
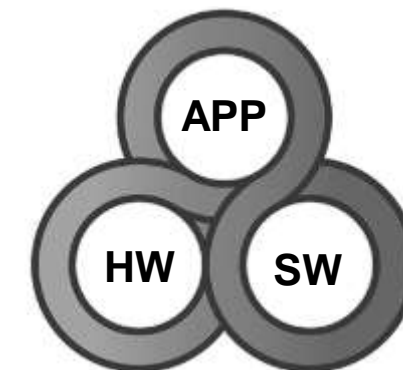
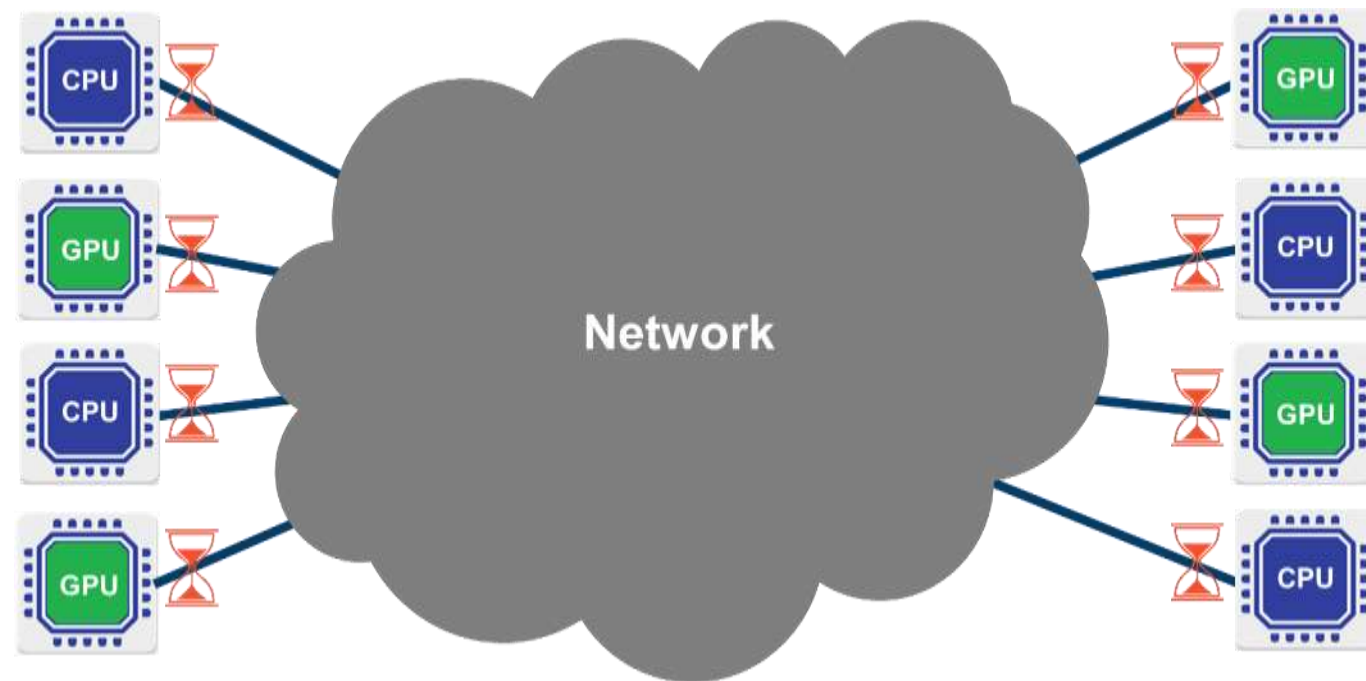## The Interconnect is the Enabling Technology

**SMP to Clusters**

**Single-Core to Many-Core**
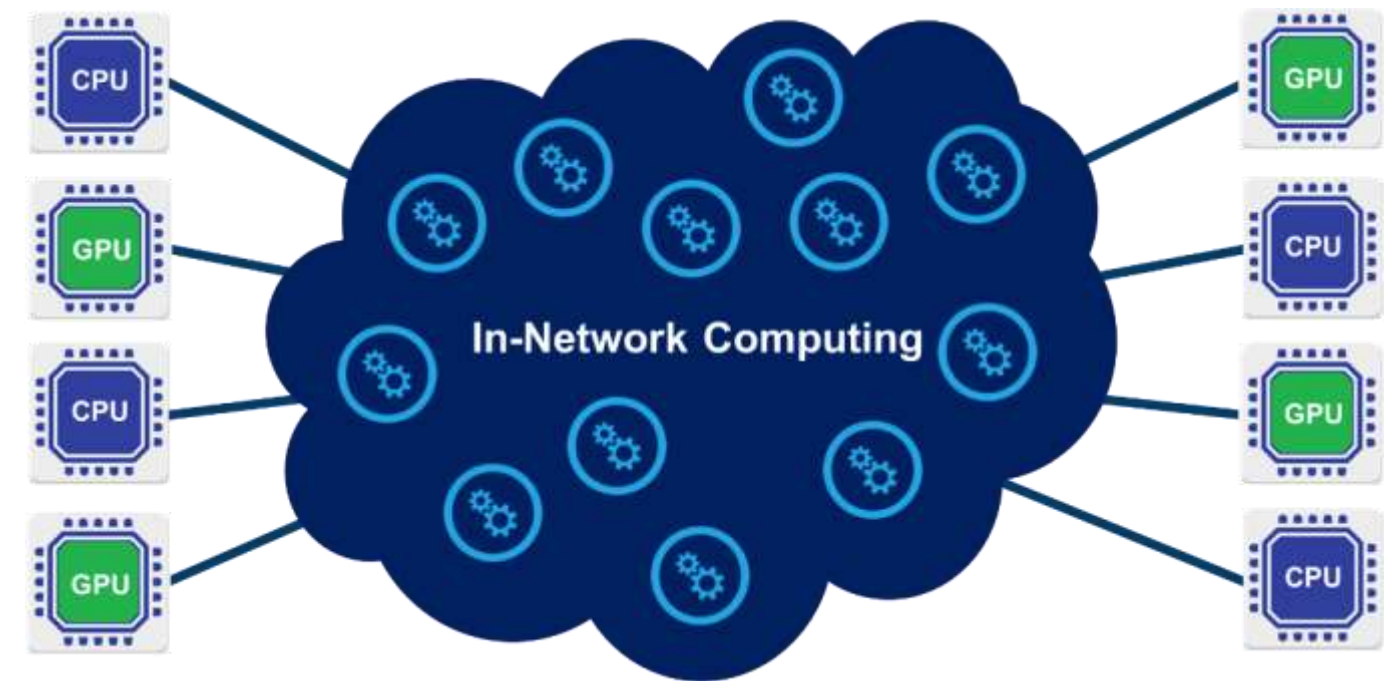
APP

HW    SW

Application
Software
Hardware

**Co-Design**

# CPU-Centric (Onload)

# Data-Centric (Offload)



**Must Wait for the Data**
**Creates Performance Bottlenecks**
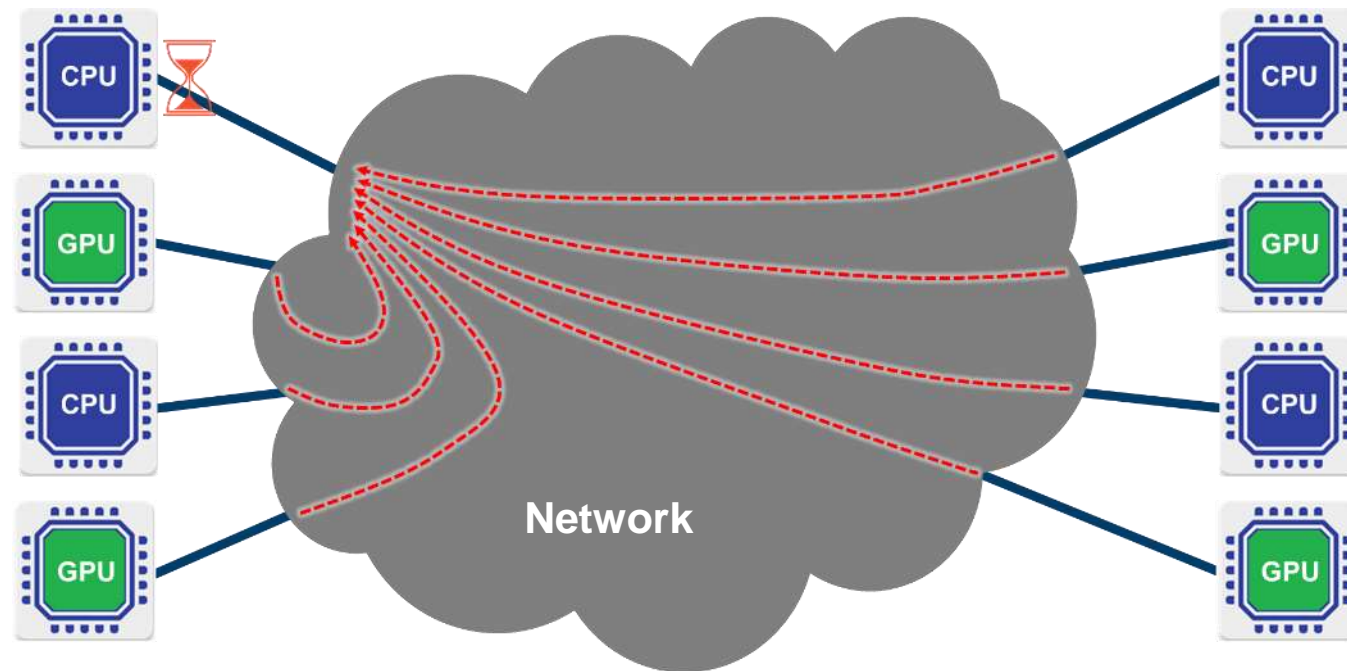
**Analyze Data as it Moves!**

**Faster Data Speeds and In-Network Computing Enable Higher Performance and Scale**

## CPU-Centric (Onload)

## Data-Centric (Offload)

Network

In-Network Computing

**HPC / Machine Learning
Communications Latencies of 30-40us**

**HPC / Machine Learning
Communications Latencies of 3-4us**

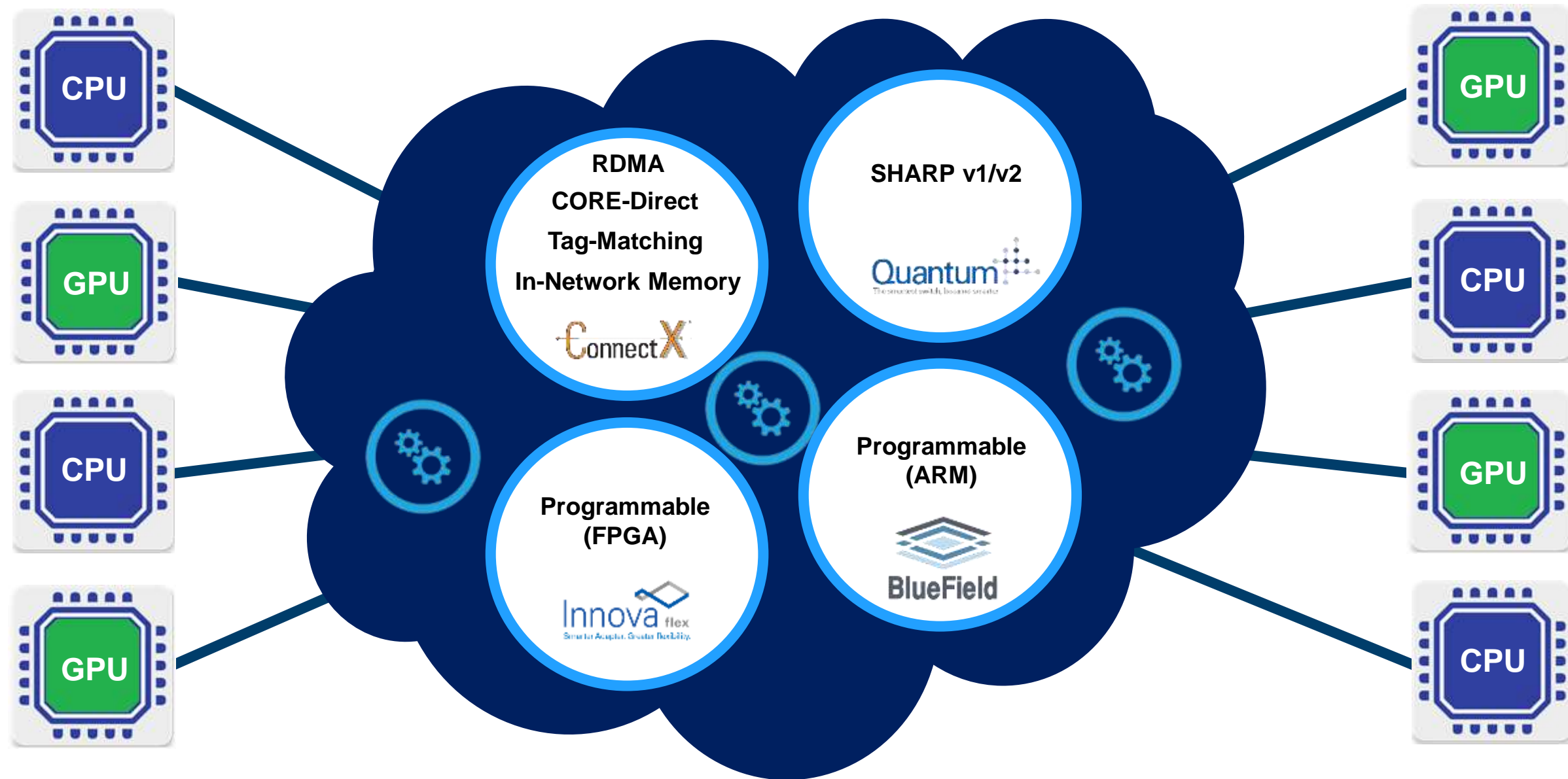## Intelligent Interconnect Paves the Road to Exascale Performance

# In-Network Computing to Enable Data-Centric Data Center

**In-Network Computing Key for Highest Return on Investment**

# In-Network Computing to Enable Data-Centric Data Center



**RDMA**
**CORE-Direct**
**Tag-Matching**
**In-Network Memory**

ConnectX

**SHARP v1/v2**

Quantum

**Programmable (FPGA)**

Innova flex

**Programmable (ARM)**

BlueField

## In-Network Computing Key for Highest Return on Investment

# In-Network Computing and Acceleration Engines

## RDMA GPUDirect

**Most Efficient Data Access and Data Movement for Compute and Storage platforms, SRIOV for HPC Clouds**

**200G with <1%CPU Utilization
10X Performance Improvement with GPUDirect**

## Collectives

**CORE-Direct and SHARP Technologies Executes and Manages Data Aggregation and Reduction Algorithms**

**Accelerates MPI, PGAS/SHMEM and UPC Communication Performance, Accelerates Machine Learning Training Algorithms**

## Storage

**NVMe over Fabrics Offloads, T10-DIF and Erasure Coding offloads**

**Efficient End-to-End Data Protection, Background Check-Pointing (burst-buffer) and More. Increase System Performance and CPU Availability**

## Network Transport

**All Communications Managed and Operated by the Network Hardware; Adaptive Routing and Congestion Management, Dynamic Connected Transport (DCT)**

**Maximizes CPU Availability for Applications, increases Network Efficiency and Scalability**

## Tag Matching

**MPI Tag-Matching Offload
MPI Rendezvous Protocol Offload**

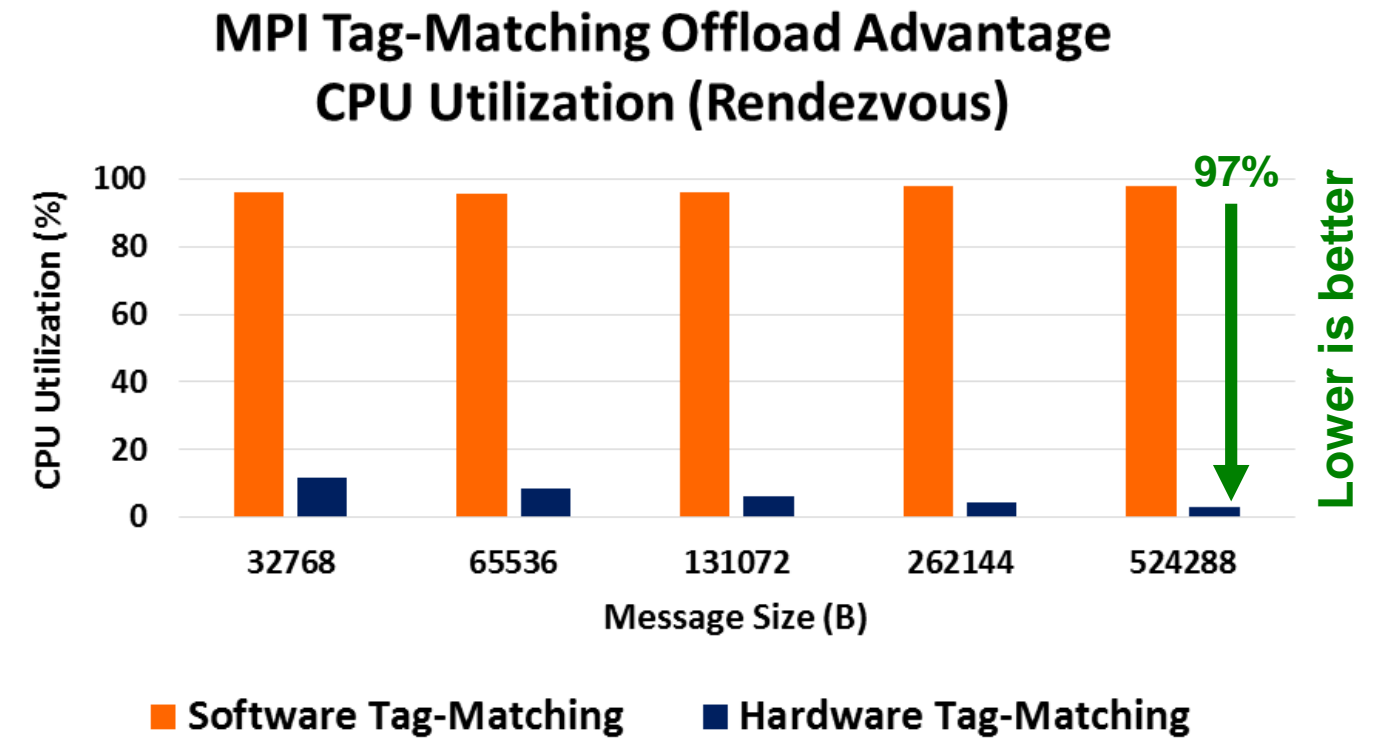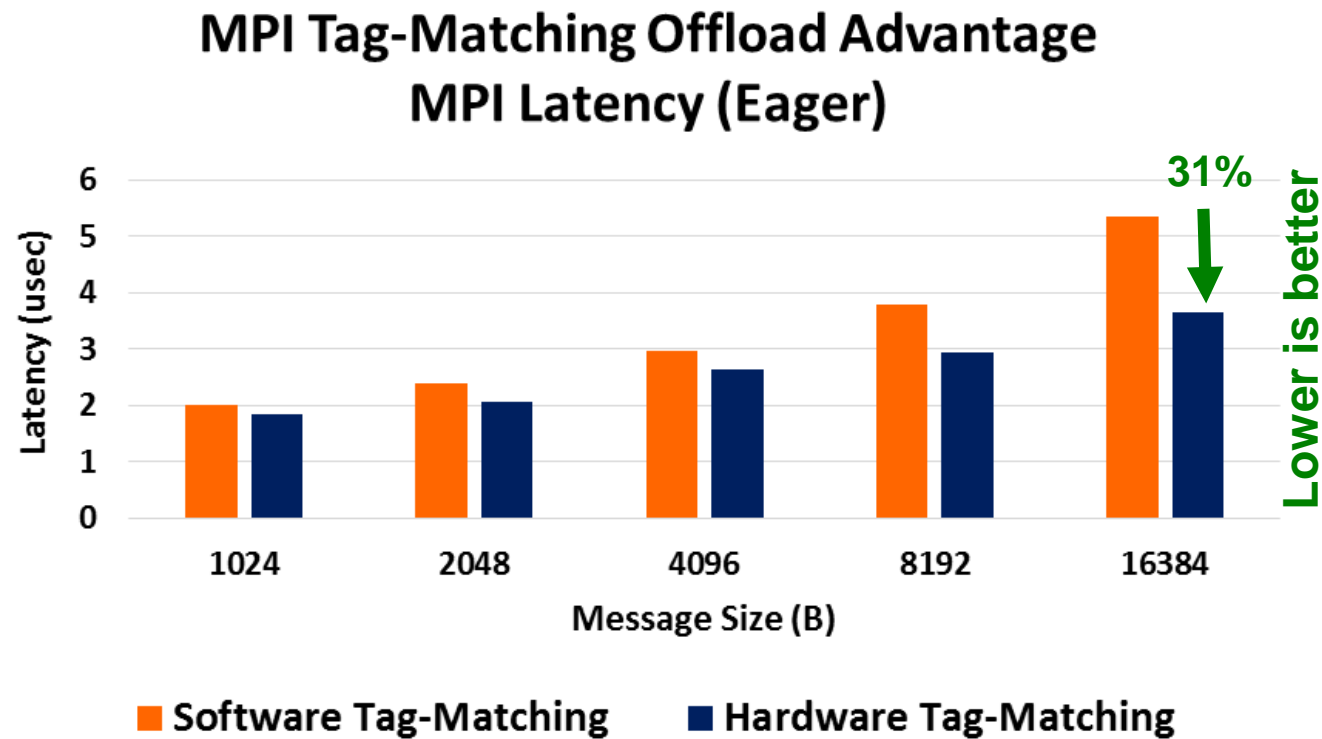**Accelerates MPI Application Performance**

## Security

**Data Encryption / Decryption (IEEE XTS standard) and Key Management; Federal Information Processing Standards (FIPS) Compliant**

**Enhances Data Security Options, Enables Protection Between Users Sharing the Same Resources (Different Keys)**
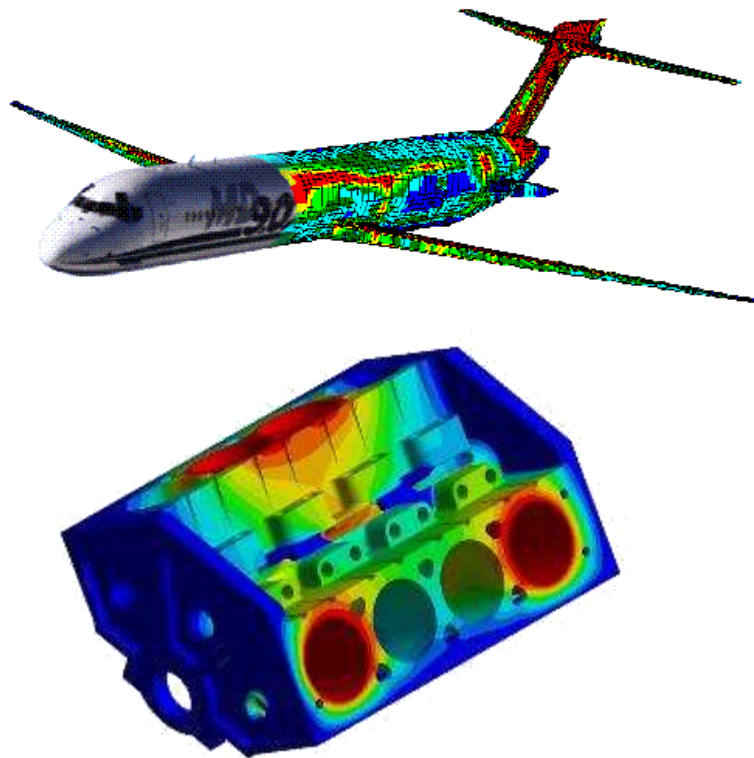
# MPI Tag-Matching Offload Advantages

## MPI Tag-Matching Offload Advantage
### MPI Latency (Eager)

**31%**

Lower is better

Latency (usec)

Message Size (B): 1024, 2048, 4096, 8192, 16384

■ Software Tag-Matching    ■ Hardware Tag-Matching

## MPI Tag-Matching Offload Advantage
### CPU Utilization (Rendezvous)

**97%**

Lower is better

CPU Utilization (%)

Message Size (B): 32768, 65536, 131072, 262144, 524288

■ Software Tag-Matching    ■ Hardware Tag-Matching

- 31% lower latency and 97% lower CPU utilization for MPI operations
- Performance comparisons based on ConnectX-5

**Mellanox In-Network Computing Technology Deliver Highest Performance**

# SHARP Performance Advantage

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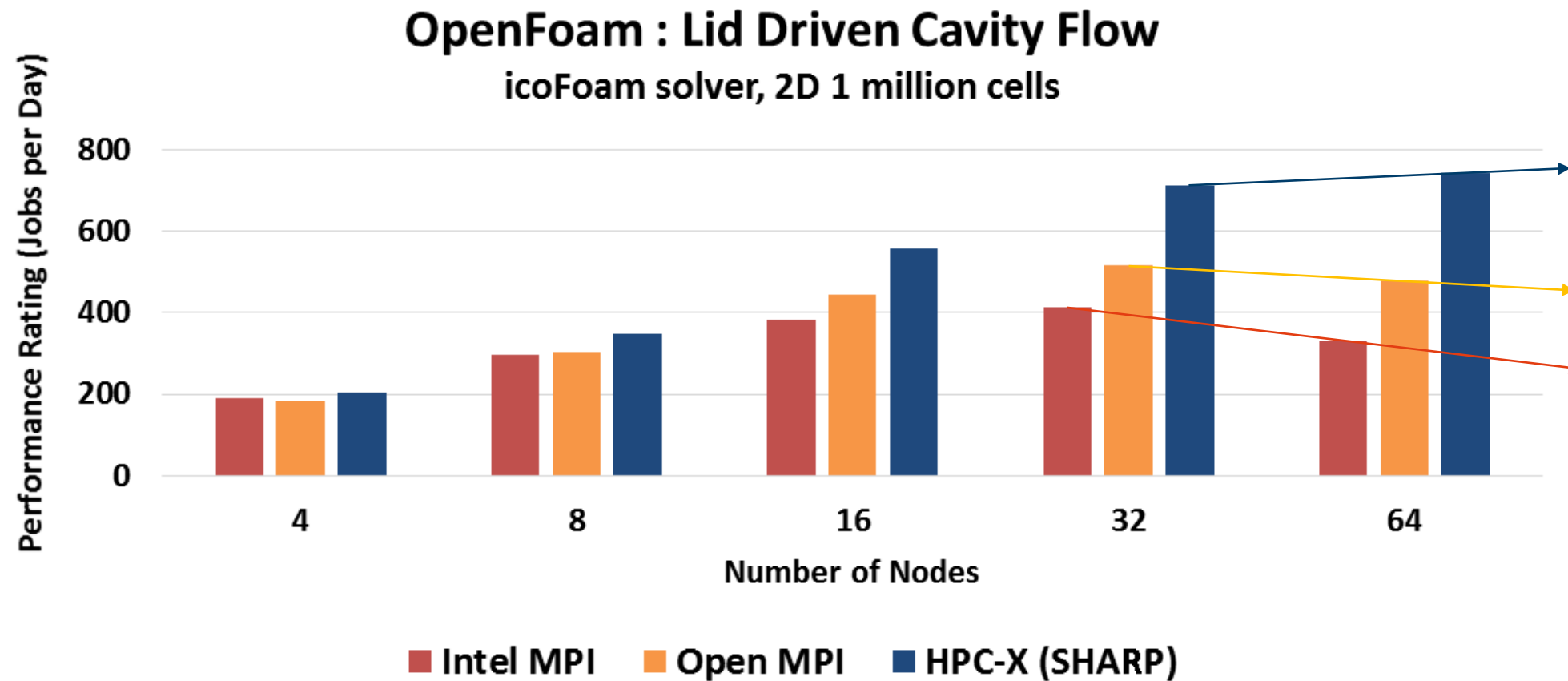


## CPU-based versus Switch Collectives Offloads
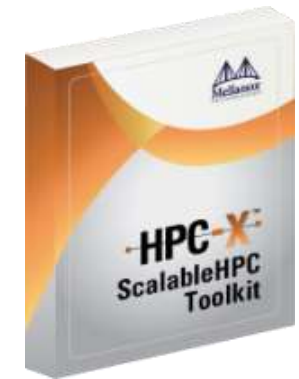## MiniFE Application - Latency Ratio (8 Bytes)

AllReduce MPI Collective

Ratio vs. Number of Nodes (32, 64, 128, 256, 512, 1024, 2048, 4096, 8192)

## 10X to 25X Performance Improvement

# Proven Advantages

- Scalable, flexible, high performance, high bandwidth, end-to-end connectivity

- Standards-based and supported by the largest eco-system

- Supports all compute architectures: x86, Power, ARM, GPU, FPGA etc.

- Native Offloading architecture

- RDMA, GPUDirect, rCUDA, SHARP and other accelerations

- Backward and future compatible

**Scalable HPC Depends on Mellanox**

# Media Resources

## OrionX Reports Position InfiniBand as the Leading HPI Technology and Mellanox the Leading Vendor

July 7, 2016 by staff — Leave a Comment

*In this special guest feature, Peter ffoulkes from OrionX outlines a series of new reports that show how InfiniBand continues to dominate the market for High Performance Interconnects.*

The OrionX Constellation reports published June 29th address the evolution, environment, evaluation and excellence ratings for the High Performance Interconnect (HPI) market. Defined as the very high end of the networking equipment market where high bandwidth and low latency are non-negotiable, HPI technologies support the most demanding workloads that are typical of extreme-scale systems in high performance computing (HPC), artificial intelligence, cloud computing, and web-scale deployments.

Peter ffoulkes, OrionX

**Link to the article**

## InfiniBand Enables Intelligent Networks

January 13, 2016 by staff

*In this special guest feature, Gilad Shainer from Mellanox writes that the network is the key to future scalable systems.*

### HPC Frequently Reinvents Itself to Keep Pace

In the world of high-performance computing, there is a constant and ever-growing demand for even higher performance. Technology providers have worked ceaselessly to keep up with that demand, with each new generation of faster, more reliable, and more efficient systems.

Ultimately, though, every technology reaches its limits, and progress can therefore stall unless there is a

Gilad Shainer, VP of Marketing, Mellanox

**Link to the article**

## Slidecast: Advantages of Offloading Architectures for HPC

April 19, 2016 by Rich Brueckner

Interconnect Your Future with InfiniBand

SB7780 Router 1U
Supports up to 6 Different Subnets

**Link to the article**

## Interview: Why Co-design is the Path Forward for Exascale Computing

March 4, 2016 by Rich Brueckner

Interview: Why Codesign is the Path Fo...

Gilad Shainer
VP of Marketing
Mellanox

**Link to the article**

## The Ultimate Debate – Interconnect Offloading Versus Onloading

April 12, 2016

Gilad Shainer, Mellanox

The high performance computing market is going through a technology transition – the Co-Design transition. As has already been discussed in many articles, this transition has emerged in order to solve the performance bottlenecks of today's infrastructures and applications, performance bottlenecks that were created by multi-core CPUs and the existing CPU-centric system architecture.

How are multi-core CPUs the source for today's performance bottlenecks? In order to understand that, we need to go back in time to the era of single-core CPUs. Back then, performance gains came from increases in CPU frequency and from the reduction of networking functions (network adapter and switches). Each new generation of product brought faster CPUs and lower-latency network adapters and

**Link to the article**

## Offloading vs. Onloading: The Case of CPU Utilization

June 18, 2016

Gilad Shainer, Mellanox

One of the primary conversations these days in the field of networking is whether it is better to onload network functions onto the CPU or better to offload these functions to the interconnect hardware.

Onloading interconnect technology is easier to build, but the issue becomes the CPU utilization; because the CPU must manage and execute network operations, it has less availability for applications, which is its primary purpose.

Offloading, on the other hand, seeks to overcome performance bottlenecks in the CPU by performing the network functions, as well as complex communications operations,

**Link to the article**

**Link to the article**

**Link to the webinar**

**Link to the article**

**Link to the article**

**Link to the article**

**Link to the Session Video**

Thank You

Mellanox TECHNOLOGIES

Connect. Accelerate. Outperform.™