



Interconnect Your Future

September 2016

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

## Performance Development

Terascale



Petascale

1<sup>st</sup>



“Roadrunner”



Exascale

OAK RIDGE  
National Laboratory

“Summit” System

Lawrence Livermore  
National Laboratory

“Sierra” System

2000

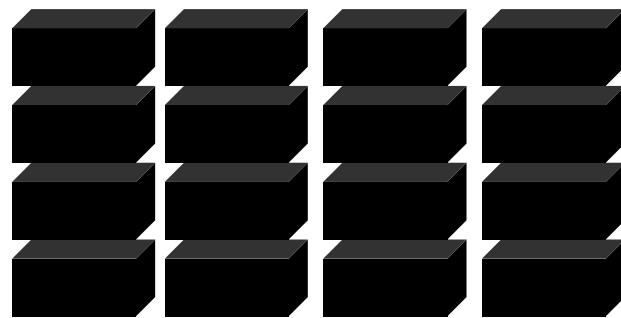
2005

2010

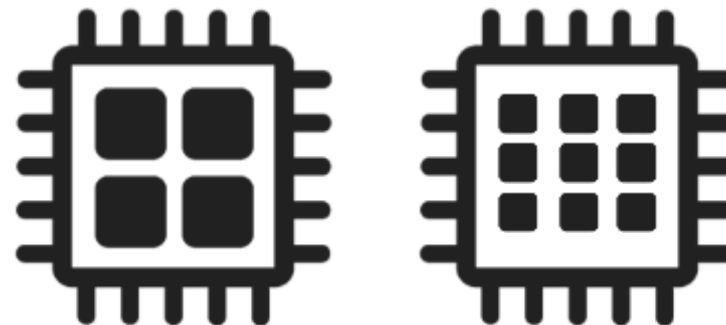
2015

2020

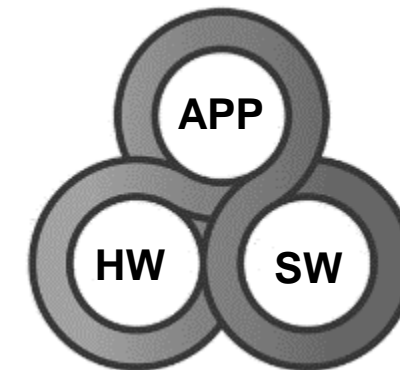
## The Interconnect is the Enabling Technology



SMP to Clusters



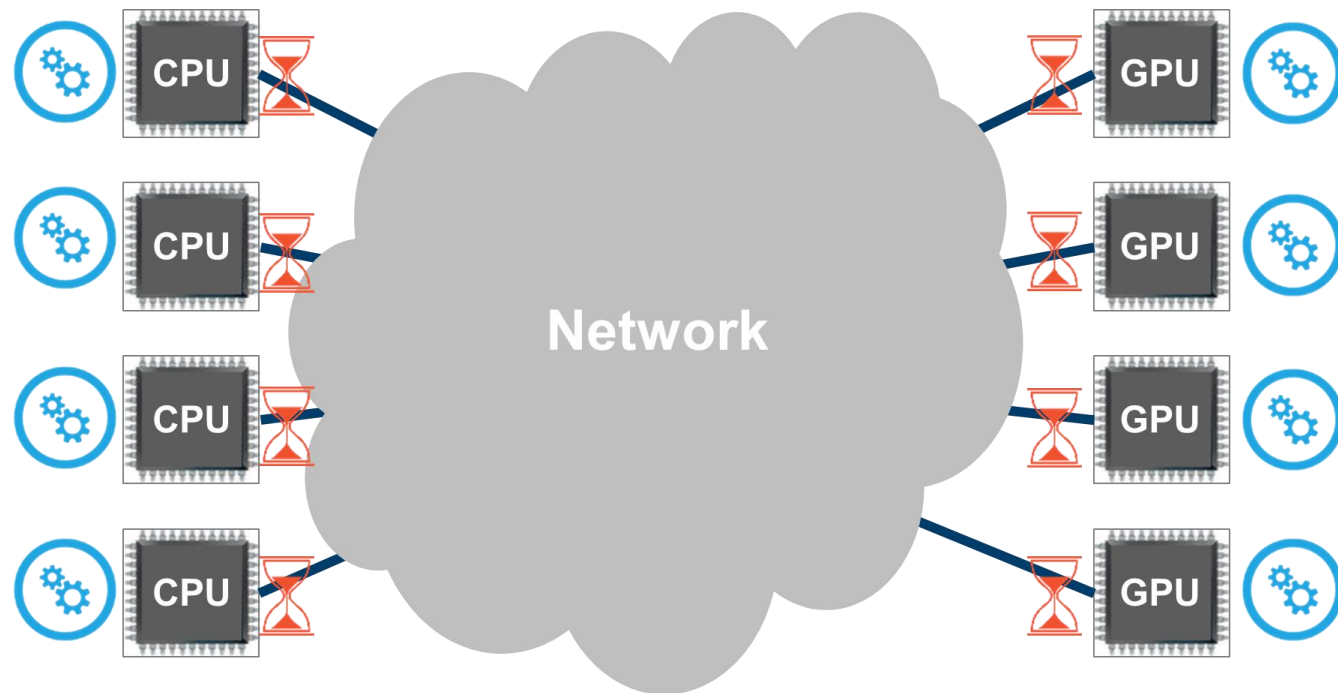
Single-Core to Many-Core



Application  
Software  
Hardware

Co-Design

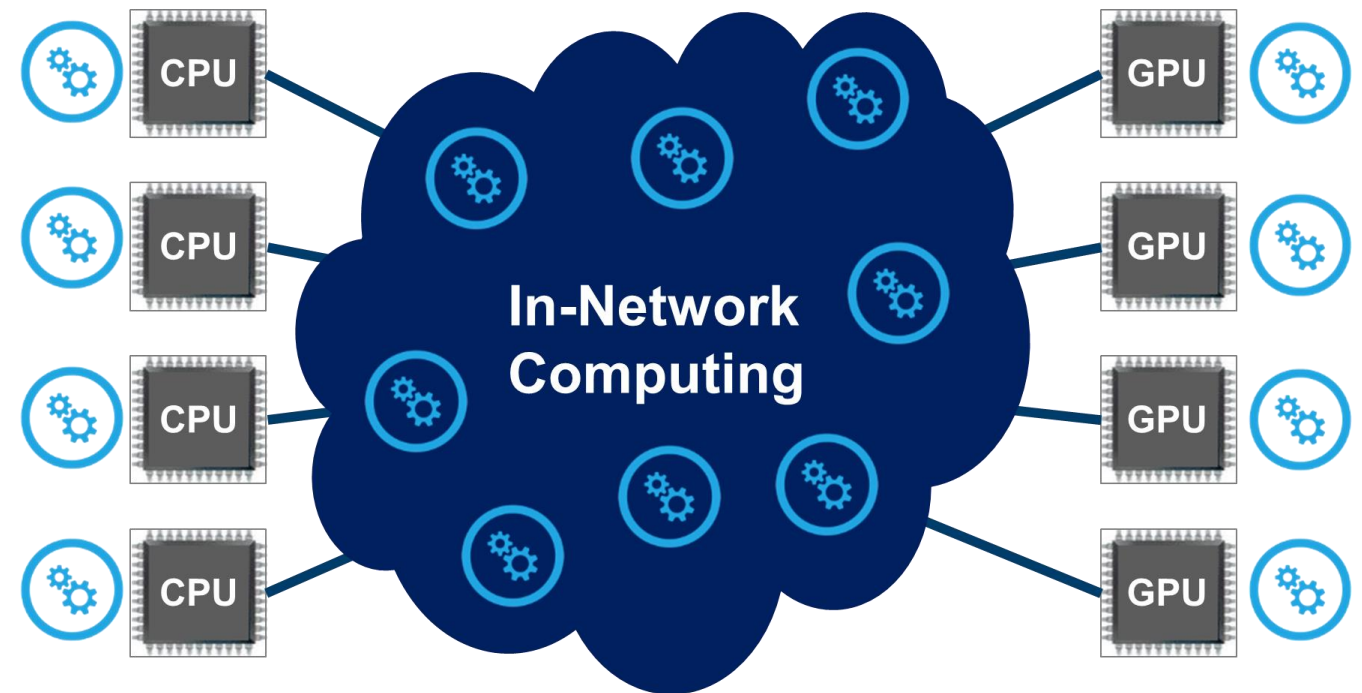
## CPU-Centric



Limited to Main CPU Usage  
Results in Performance Limitation

**Must Wait for the Data  
Creates Performance Bottlenecks**

## Co-Design



Creating Synergies  
Enables Higher Performance and Scale

**Work on The Data as it Moves  
Enables Performance and Scale**

# State of the Smart

a new generation of co-processors emerges

# Mellanox Smart Interconnect

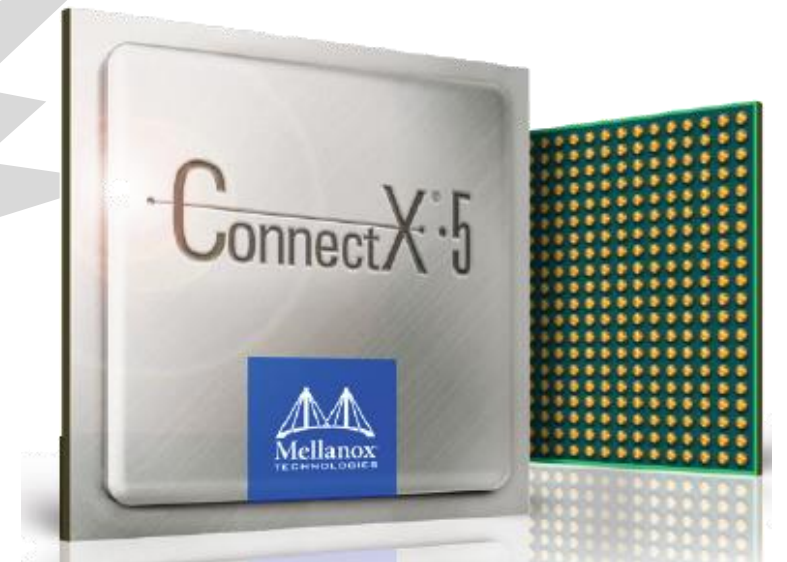
Switch IB™ 2

SHARP

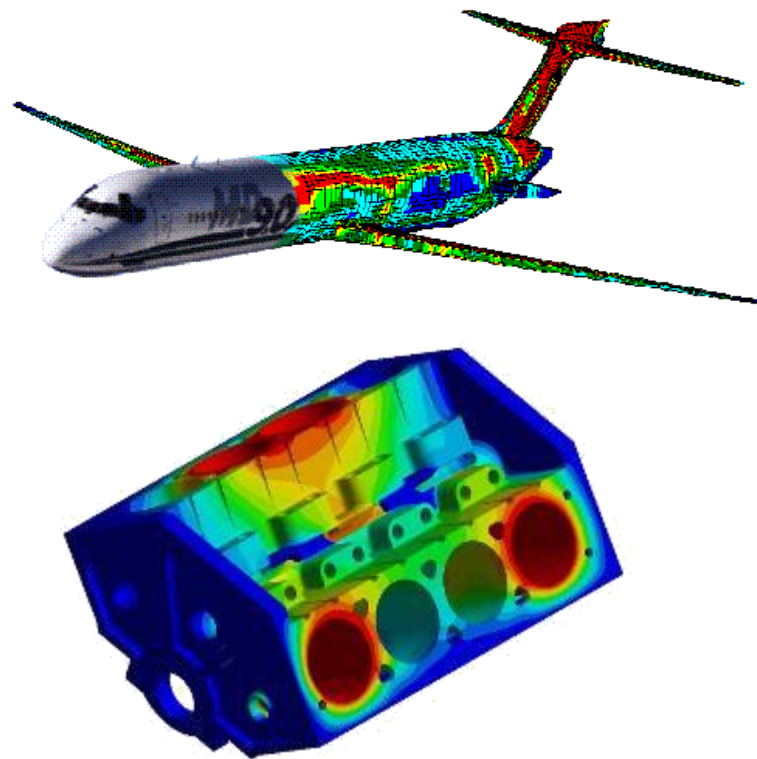
ConnectX® 5



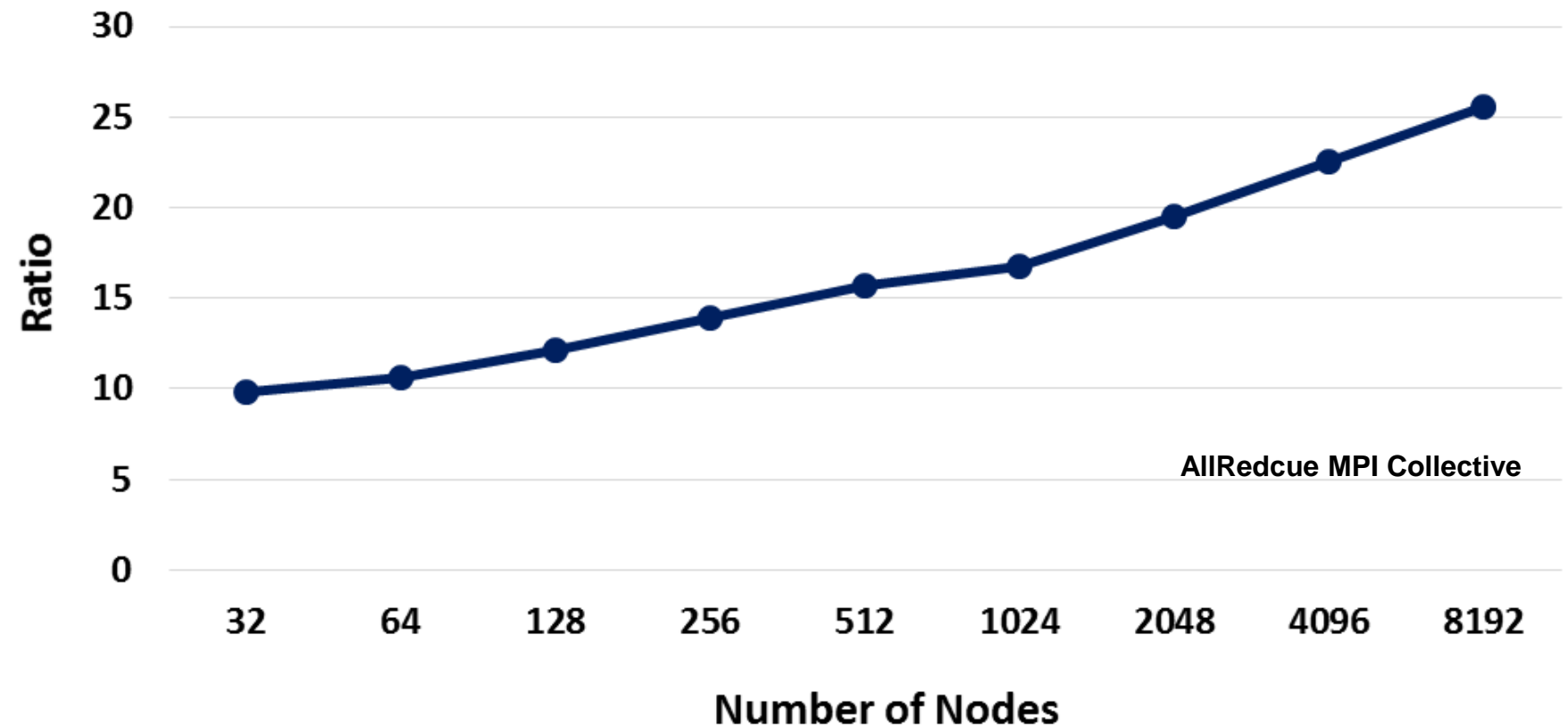
NEW!



- 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)



**10X to 25X** Performance Improvement!

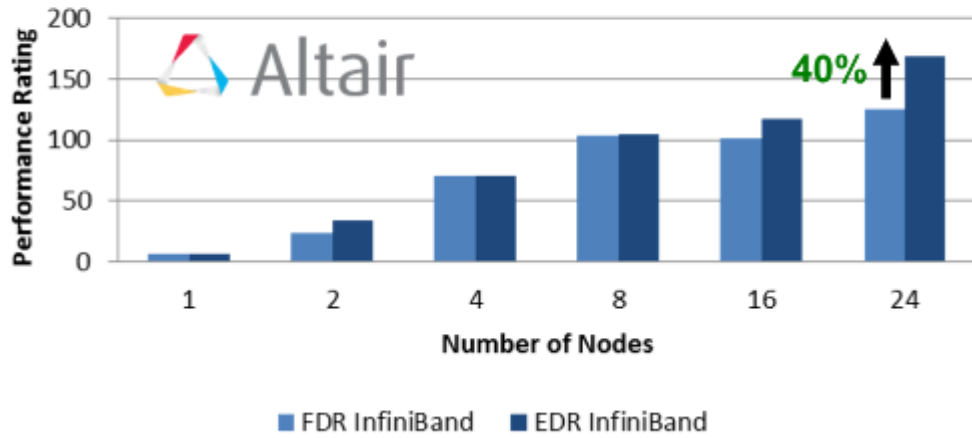
# Highest-Performance 100Gb/s Interconnect Solutions

Adapters		<p>100Gb/s Adapter, 0.6us latency 200 million messages per second (10 / 25 / 40 / 50 / 56 / 100Gb/s)</p>	
Switch		<p>36 EDR (100Gb/s) Ports, &lt;90ns Latency Throughput of 7.2Tb/s 7.02 Billion msg/sec (195M msg/sec/port)</p>	
Switch		<p>32 100GbE Ports, 64 25/50GbE Ports (10 / 25 / 40 / 50 / 100GbE) Throughput of 6.4Tb/s</p>	
Interconnect		<p>Transceivers Active Optical and Copper Cables (10 / 25 / 40 / 50 / 56 / 100Gb/s)</p>	 <p>VCSELs, Silicon Photonics and Copper</p>
Software		<p>MPI, SHMEM/PGAS, UPC For Commercial and Open Source Applications Leverages Hardware Accelerations</p>	

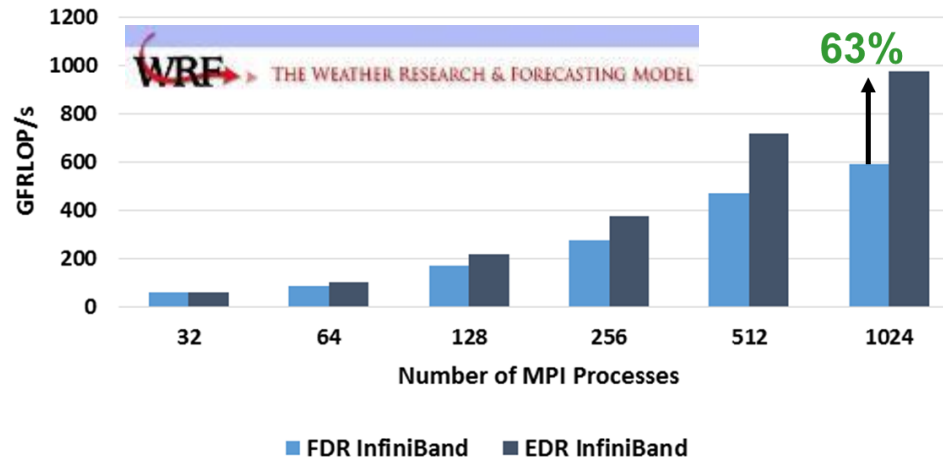
# The Performance Advantage of EDR 100G InfiniBand (28-80%)



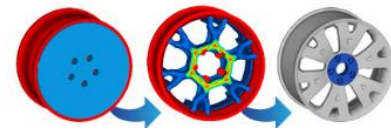
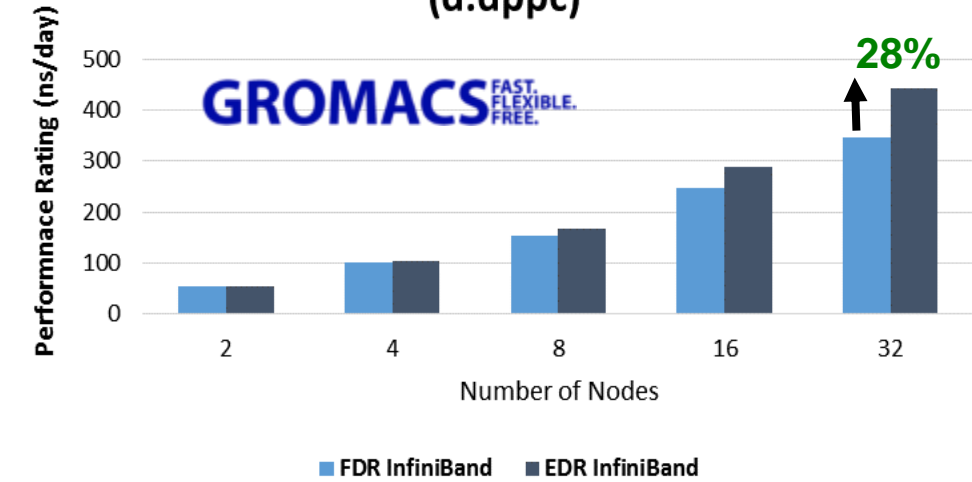
### OptiStruct Performance (Engine\_Assy.fem)



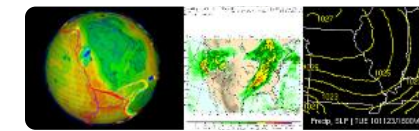
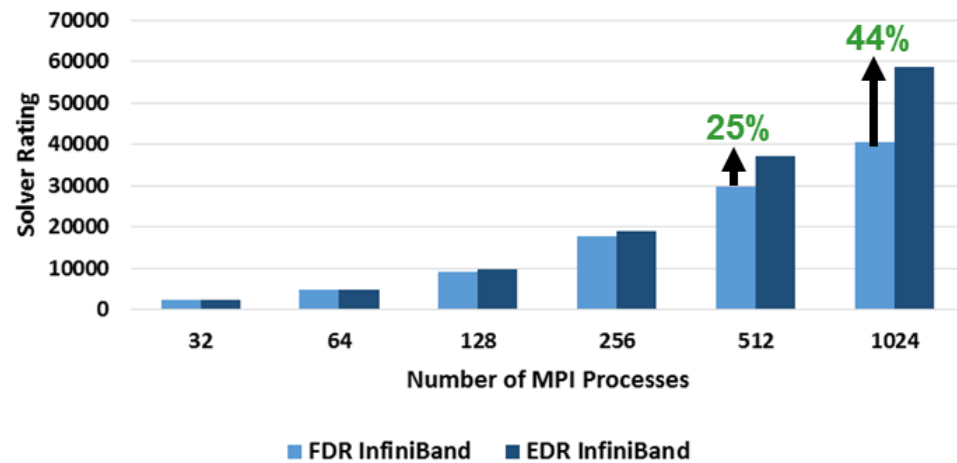
### WRF Performance (conus12km)



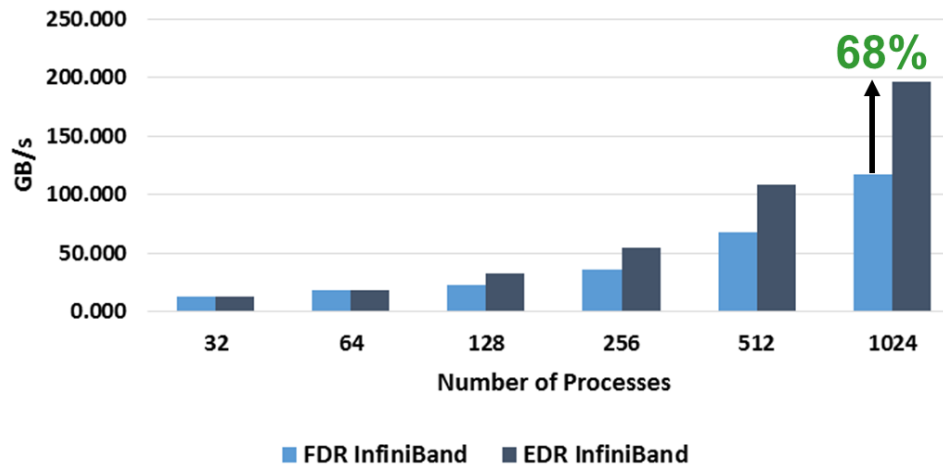
### GROMACS Performance (d.dppc)



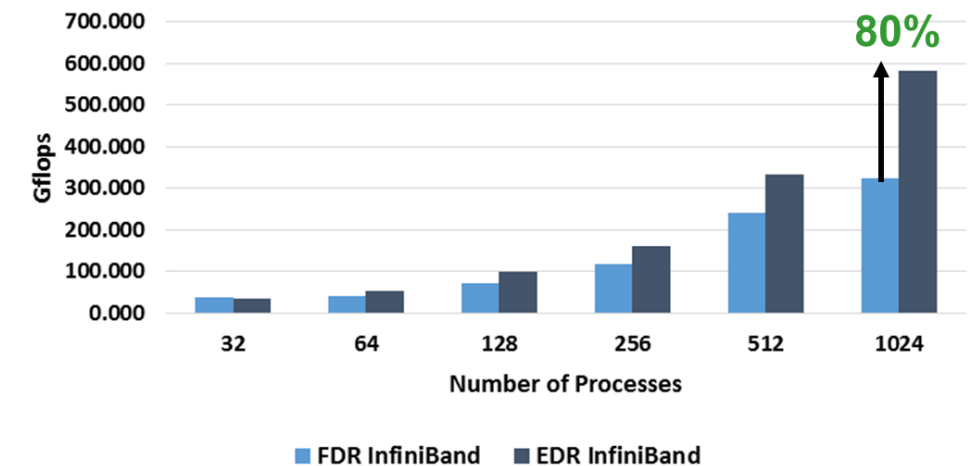
### ANSYS Fluent 16.0 Performance (sedan\_4m)



### HPCC Performance (PTRANS\_GB)



### HPCC Performance (MPIFFT)



# Mellanox Connects the World's Fastest Supercomputer



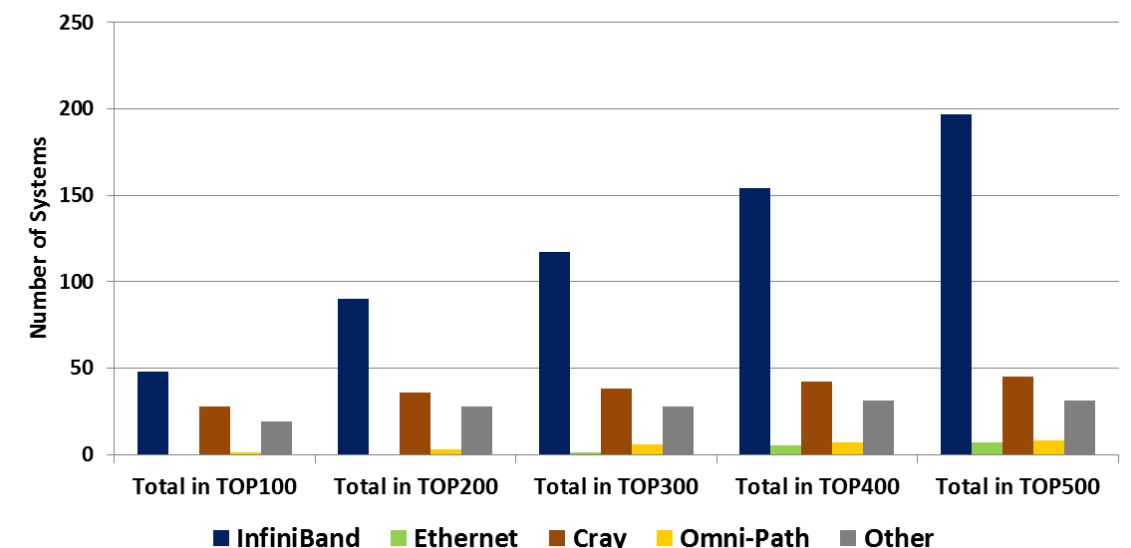
## National Supercomputing Center in Wuxi, China #1 on the TOP500 Supercomputing List

- 93 Petaflop performance, 3X higher versus #2 on the TOP500
- 40K nodes, 10 million cores, 256 cores per CPU
- Mellanox adapter and switch solutions

- The TOP500 list has evolved, includes HPC & Cloud / Web2.0 Hyperscale systems
- Mellanox connects 41.2% of overall TOP500 systems
- Mellanox connects 70.4% of the TOP500 HPC platforms
- Mellanox connects 46 Petascale systems, Nearly 50% of the total Petascale systems

**InfiniBand is the Interconnect of Choice for  
HPC Compute and Storage Infrastructures**

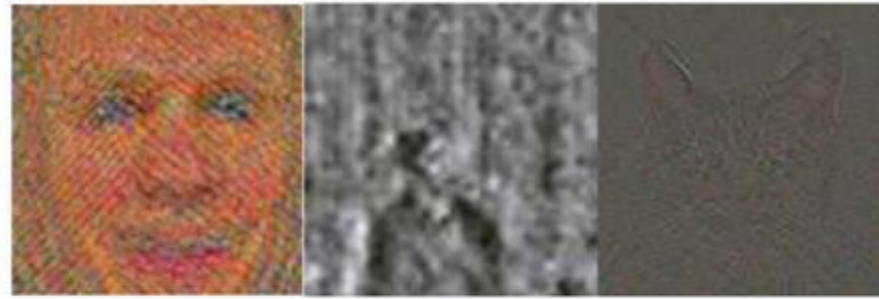
TOP500 - TOP 100, 200, 300, 400, 500 Systems Distribution  
HPC Systems Only







# GPUDirect Enables Efficient Training Platform for Deep Neural Network



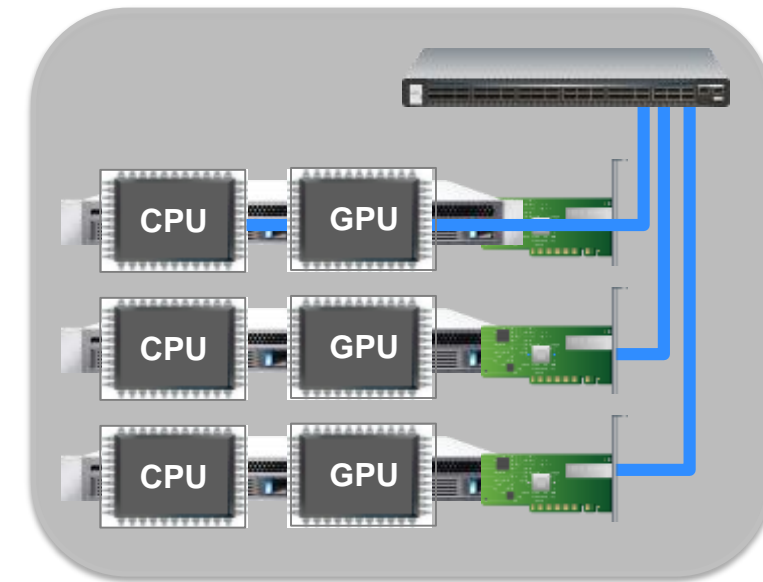
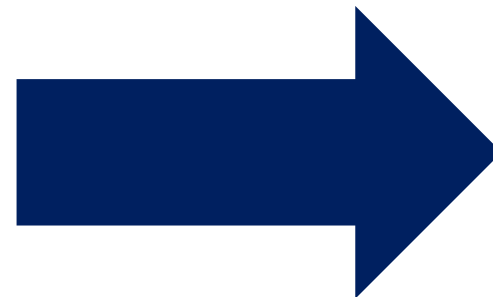
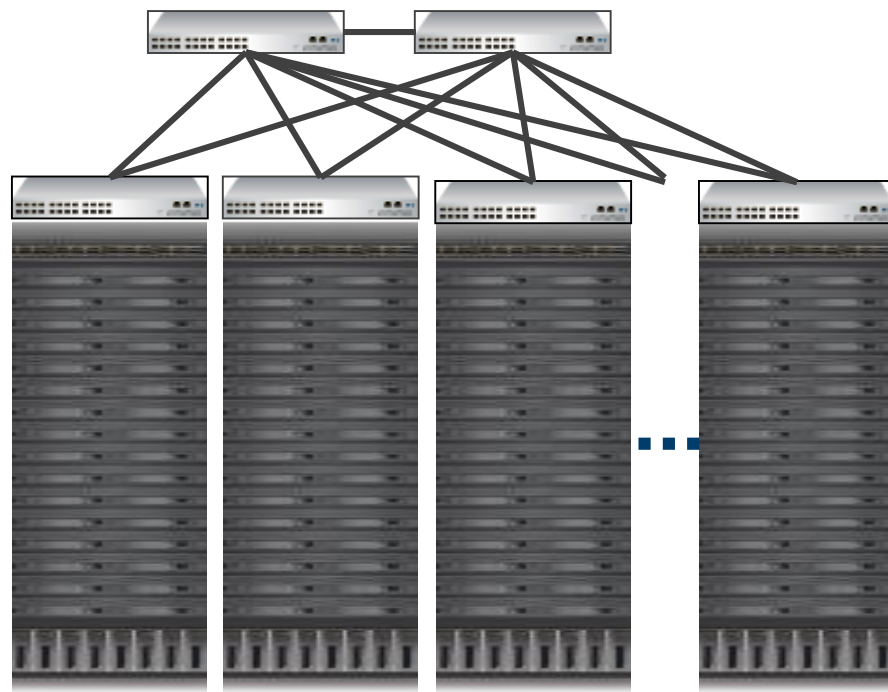
(a) Face

(b) Body

(c) Cat



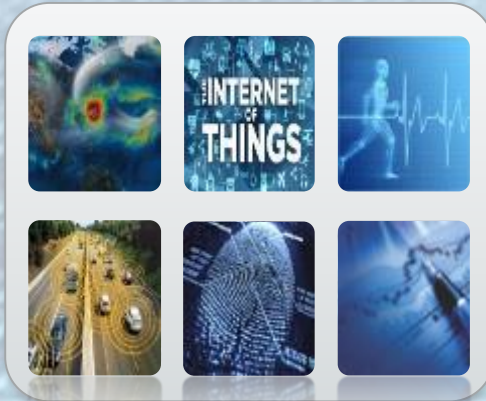
THE OHIO STATE UNIVERSITY



1K nodes (16K cores) for 1 week

3 Nodes with 3 GPUs for 3 days  
Mellanox InfiniBand and GPU-Direct

## More Data



## Better Models



## Faster Compute

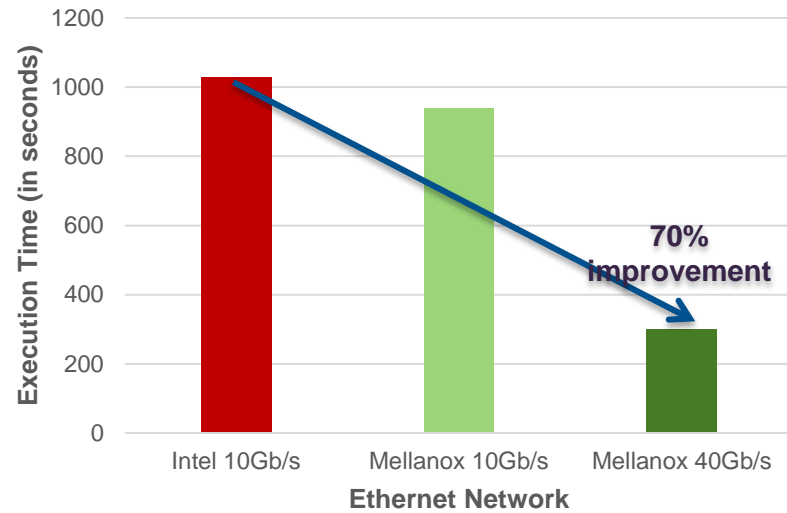


Smart Interconnect Required to Unleash The Power of Data

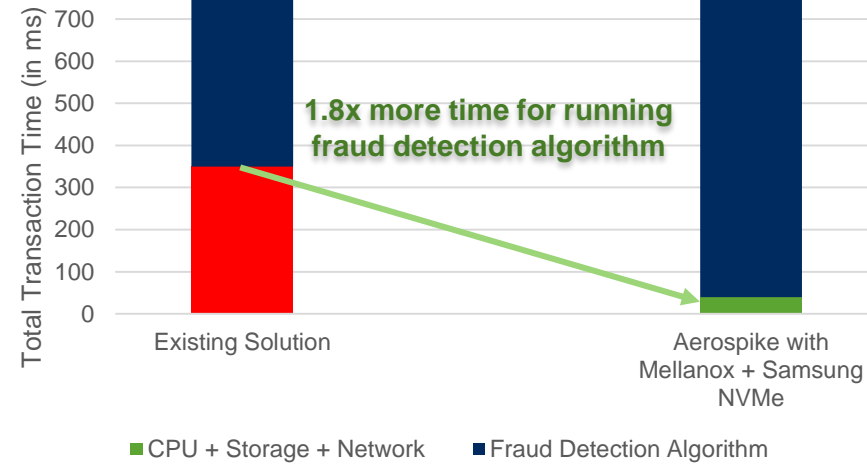


# Enable Real-time Decision Making

## TeraSort

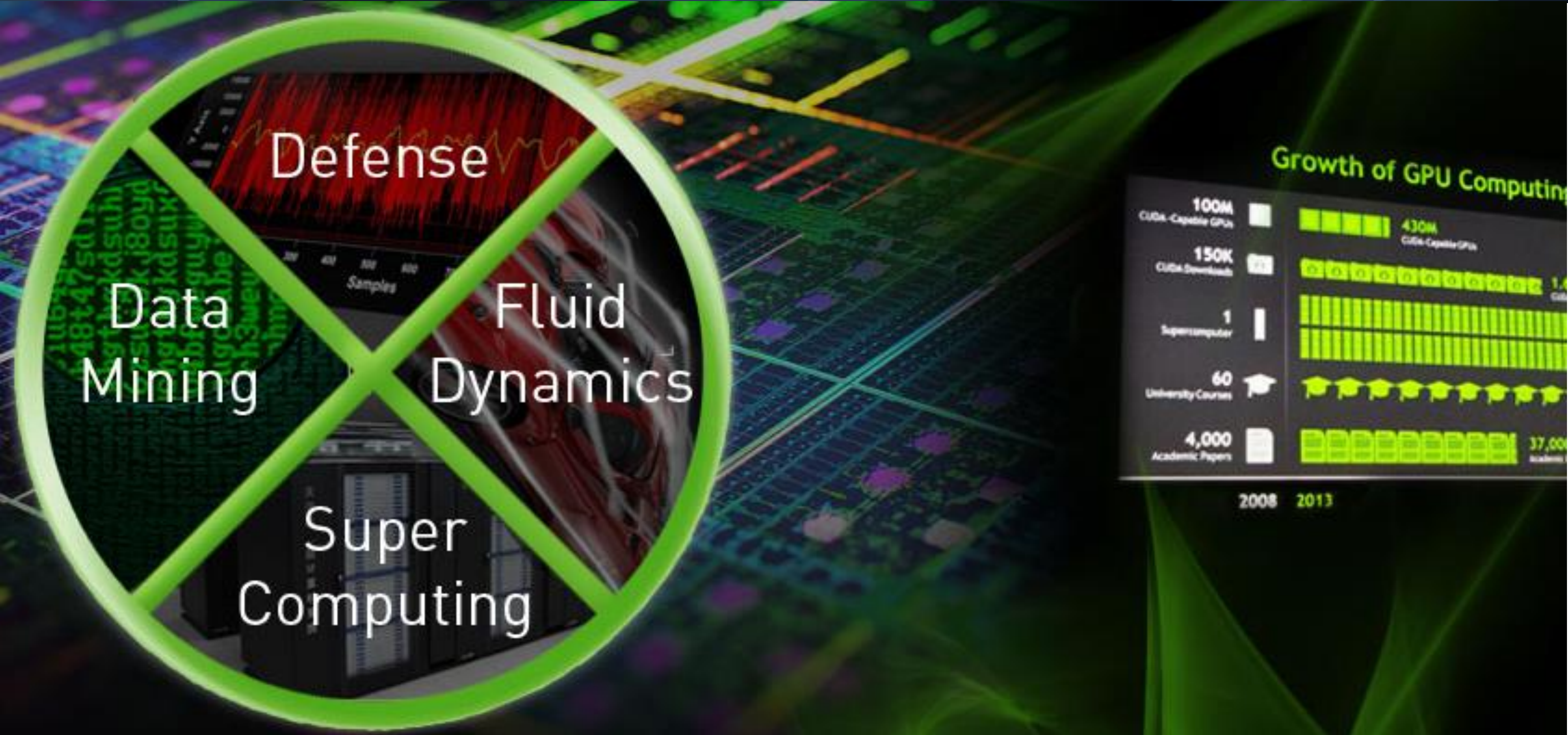


## Fraud Detection workload

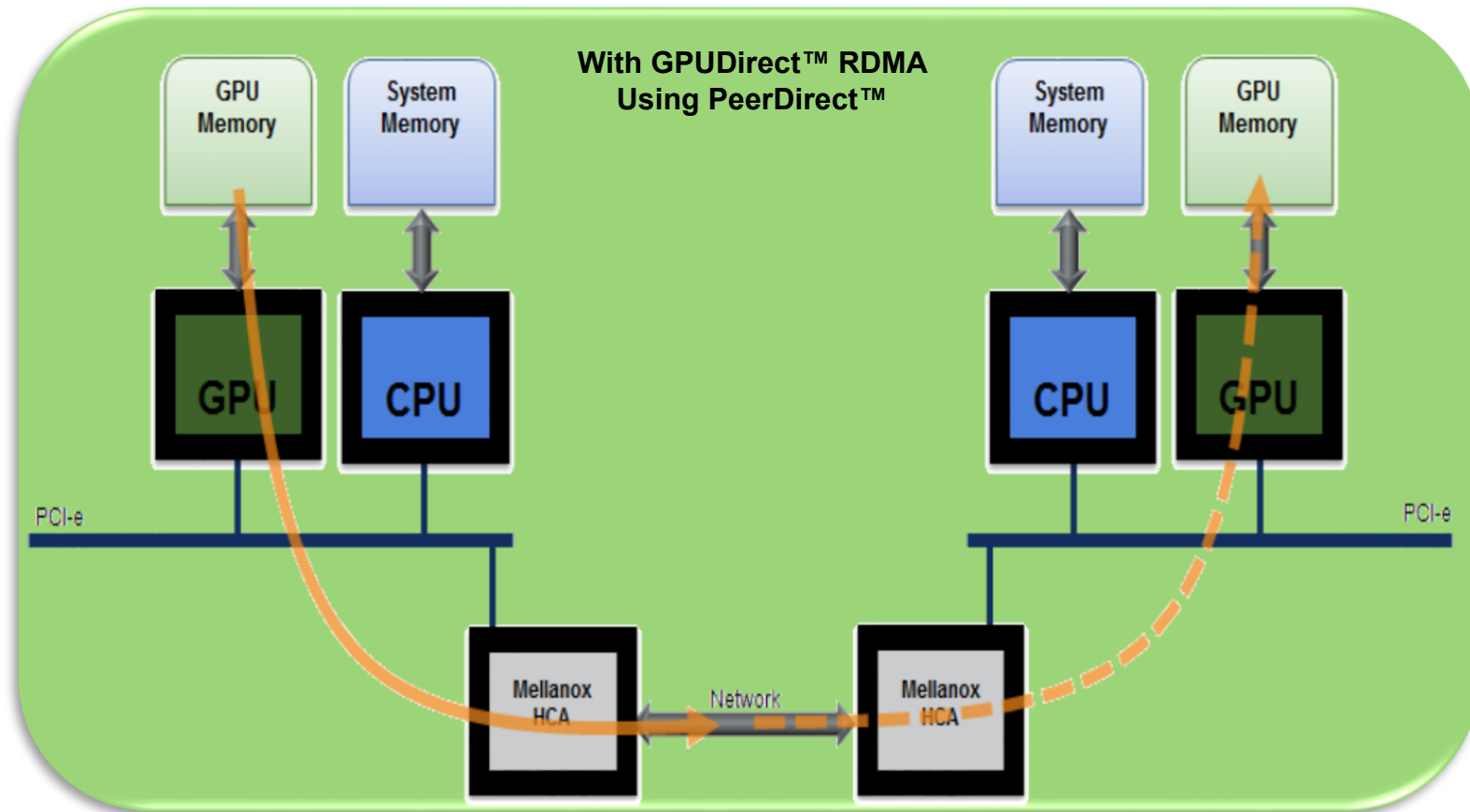


## Big Sur Machine Learning Platform





# GPUDirect™ RDMA Ecosystem



Mission Systems and Training



筑波大学  
University of Tsukuba



NORTHROP GRUMMAN



HZDR



SPACEX

map-D

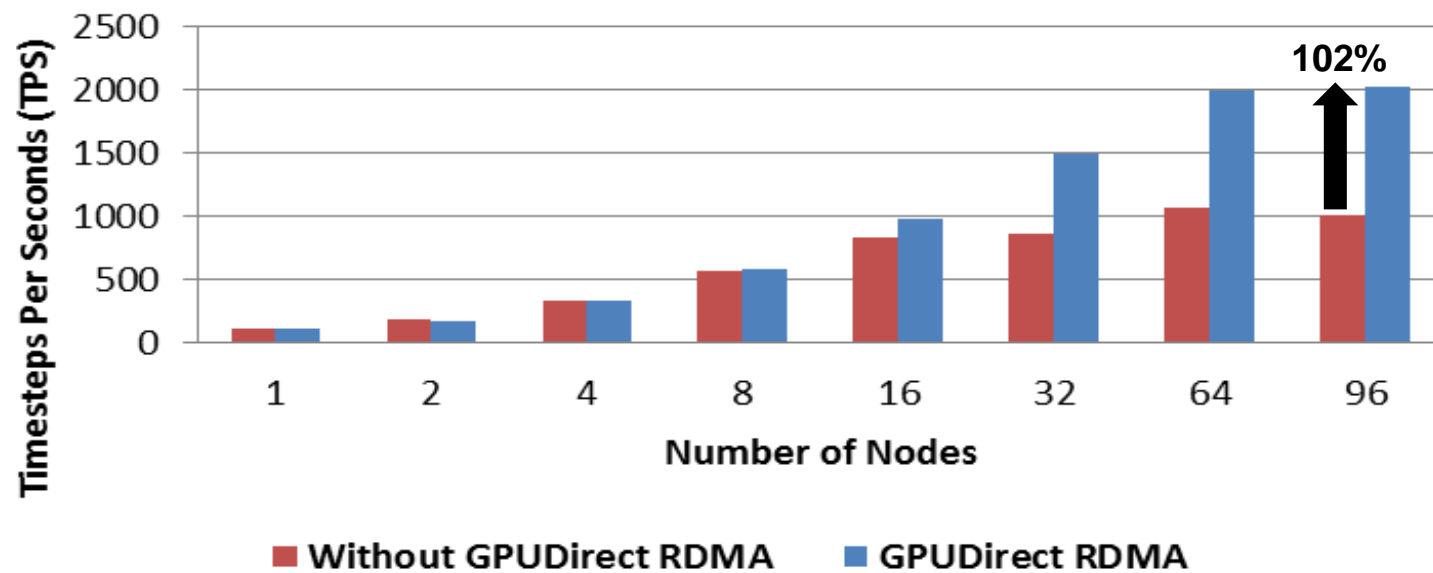
SQREAM TECHNOLOGIES



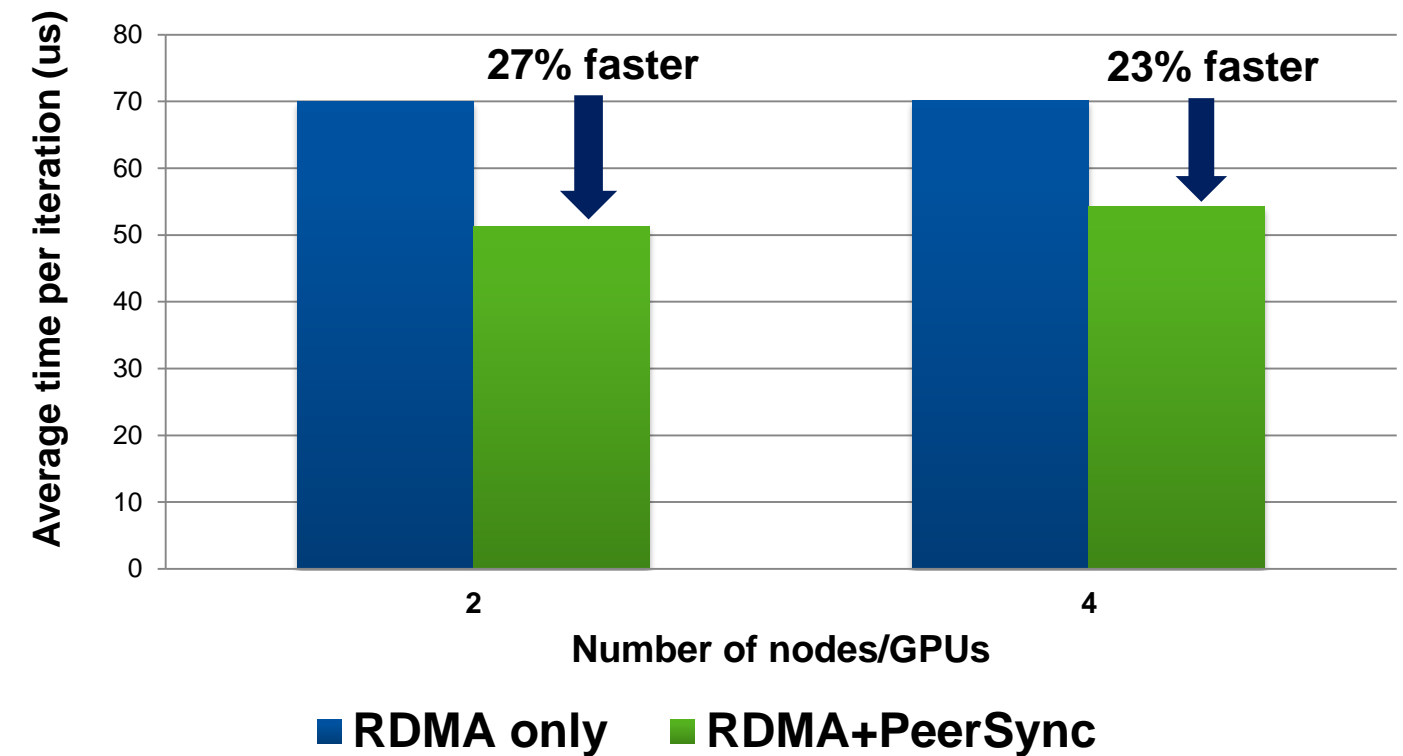
HELMHOLTZ ZENTRUM DRESDEN ROSSENDORF



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

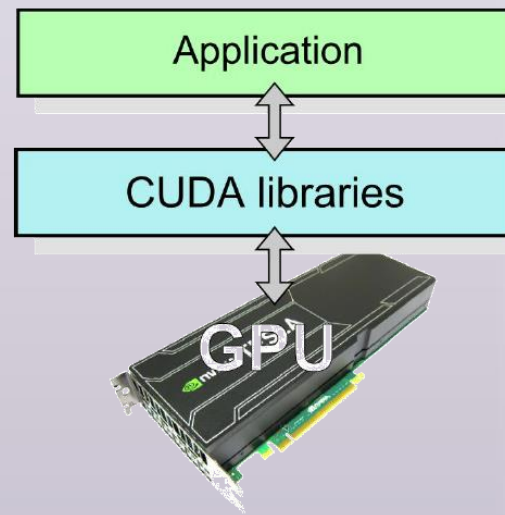
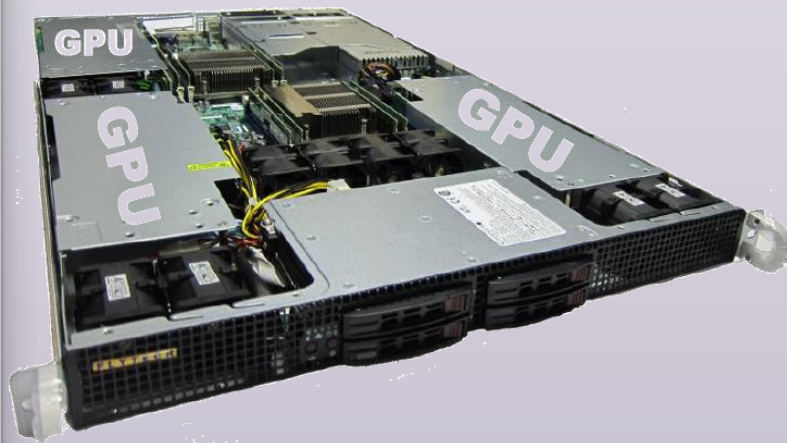


## 2D stencil benchmark



## Basic GPU computing

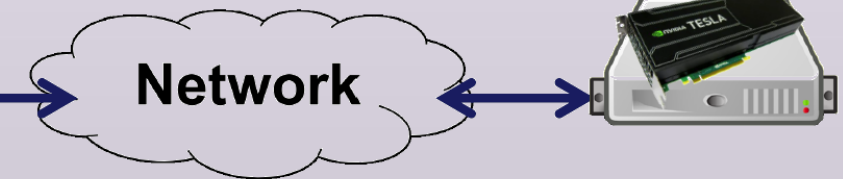
Basic behavior of CUDA



## rCUDA – remote CUDA

A software technology that enables a more flexible use of GPUs in computing facilities

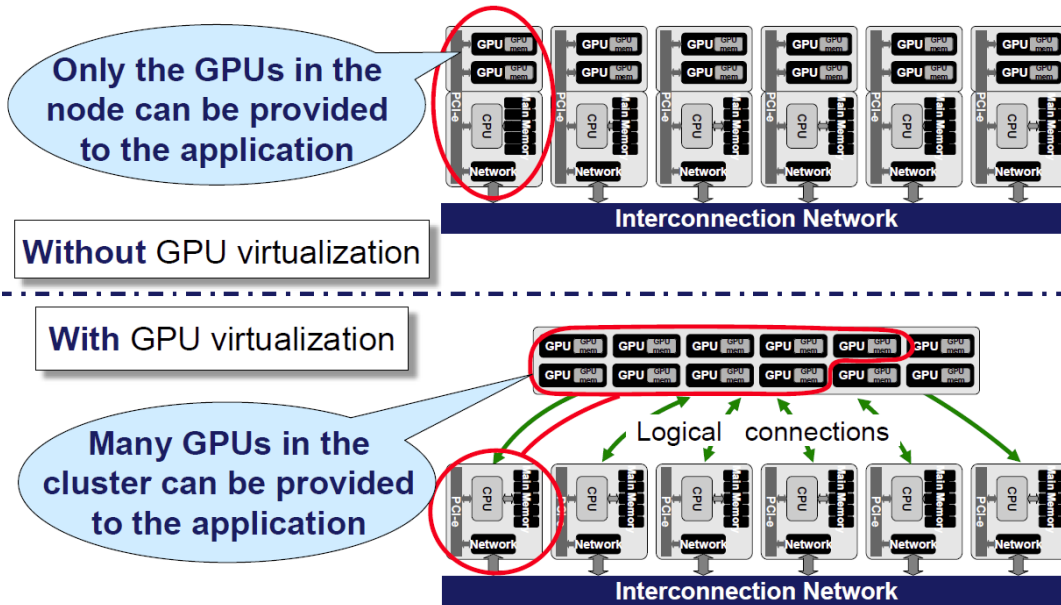
**No GPU**



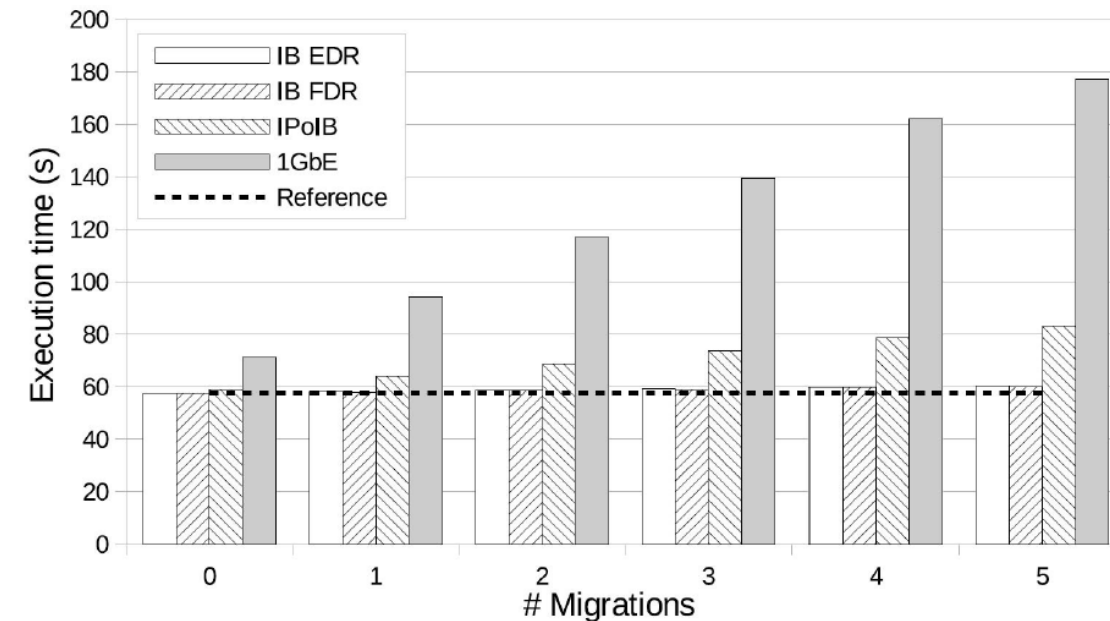
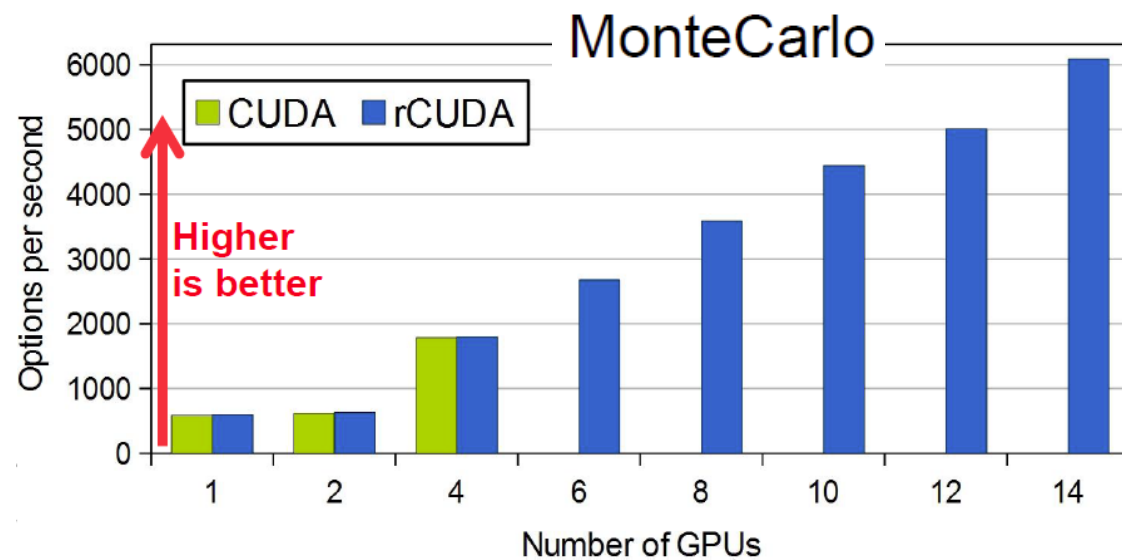
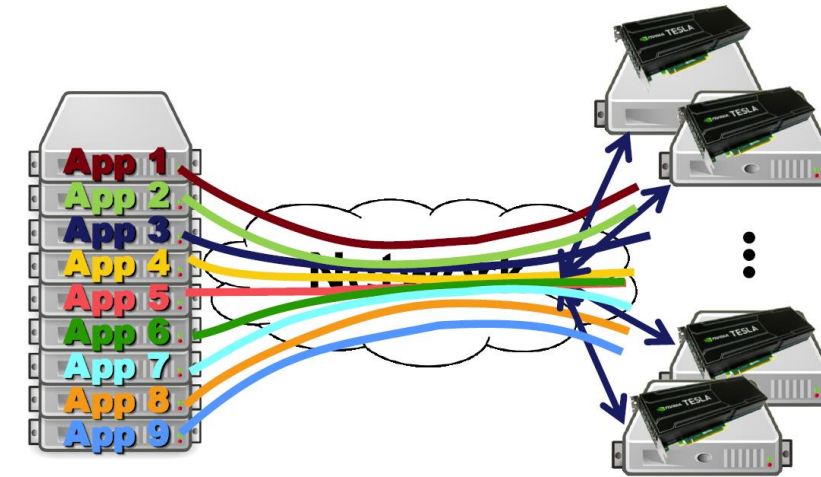
rCUDA supports CUDA 7.5



## Flexible Resources Assignment

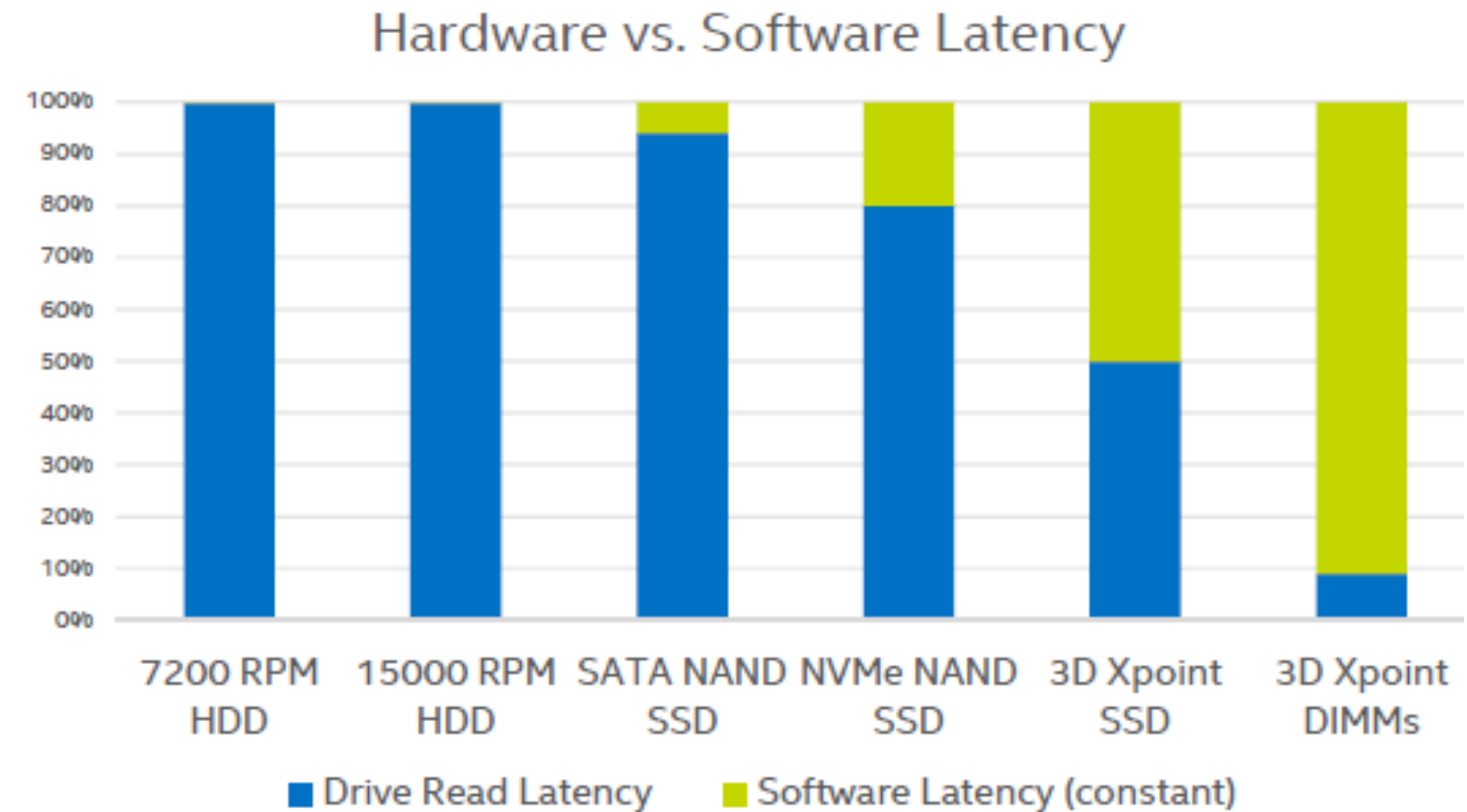
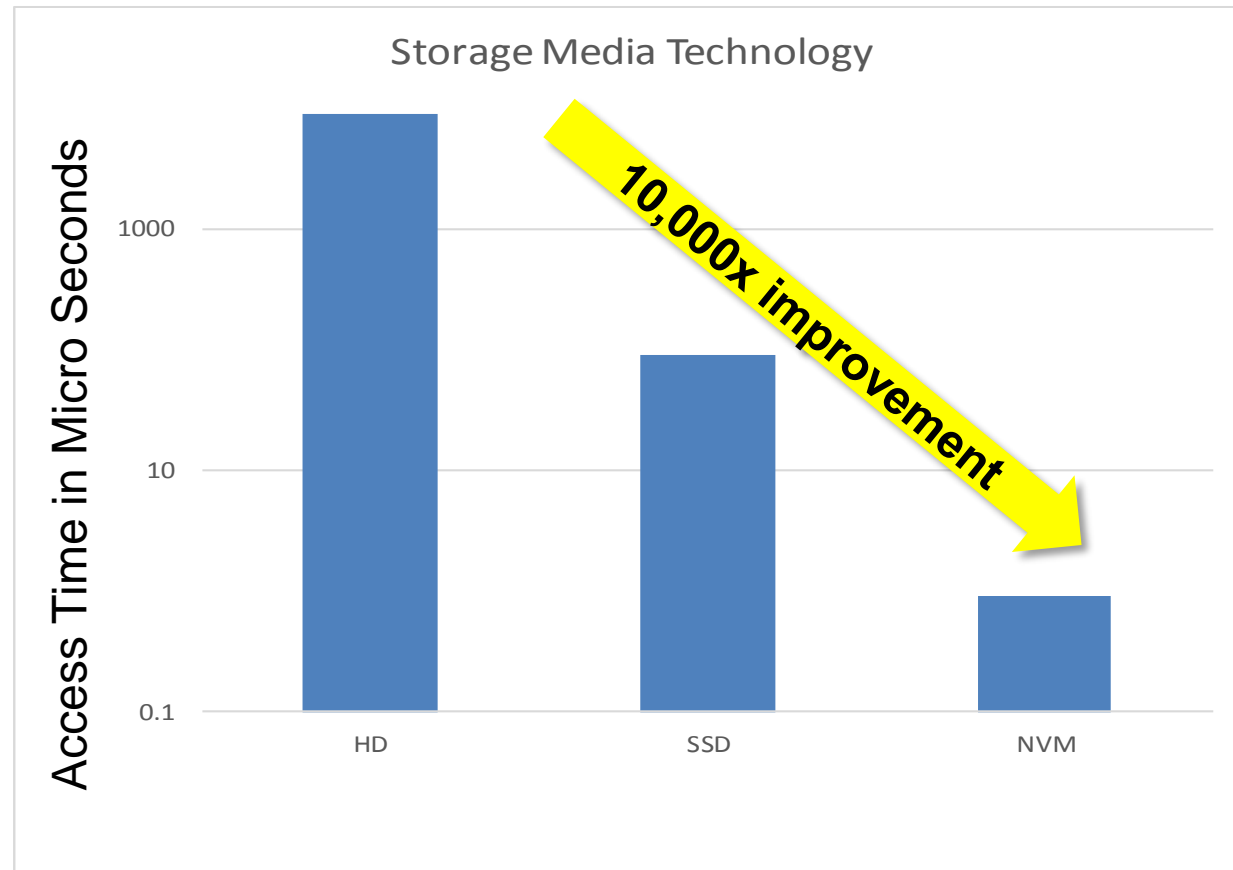


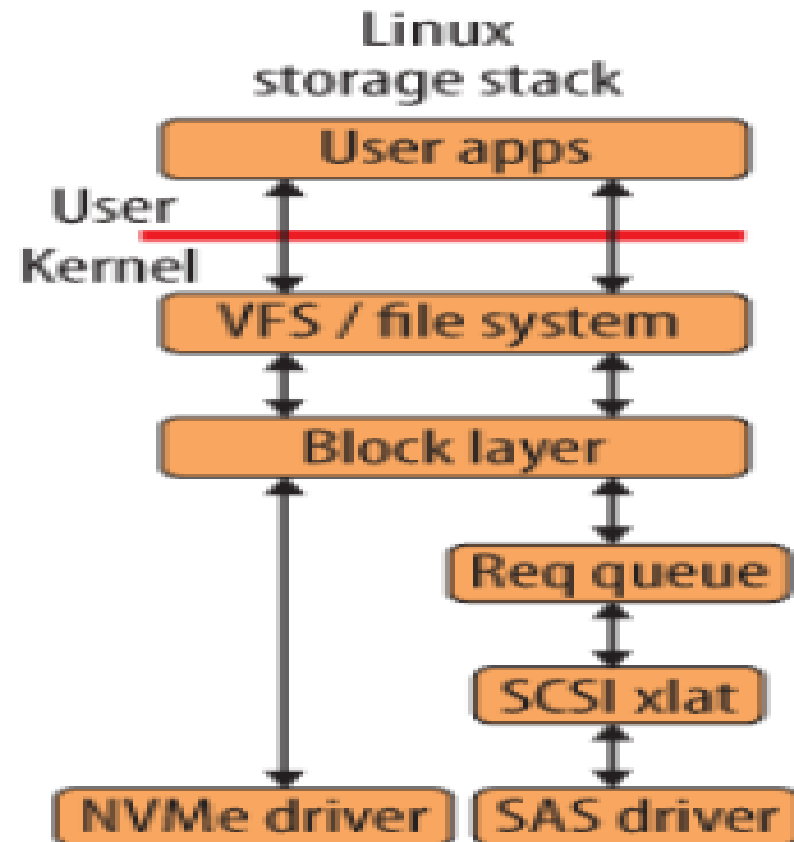
## Resources' Consolidation and Management



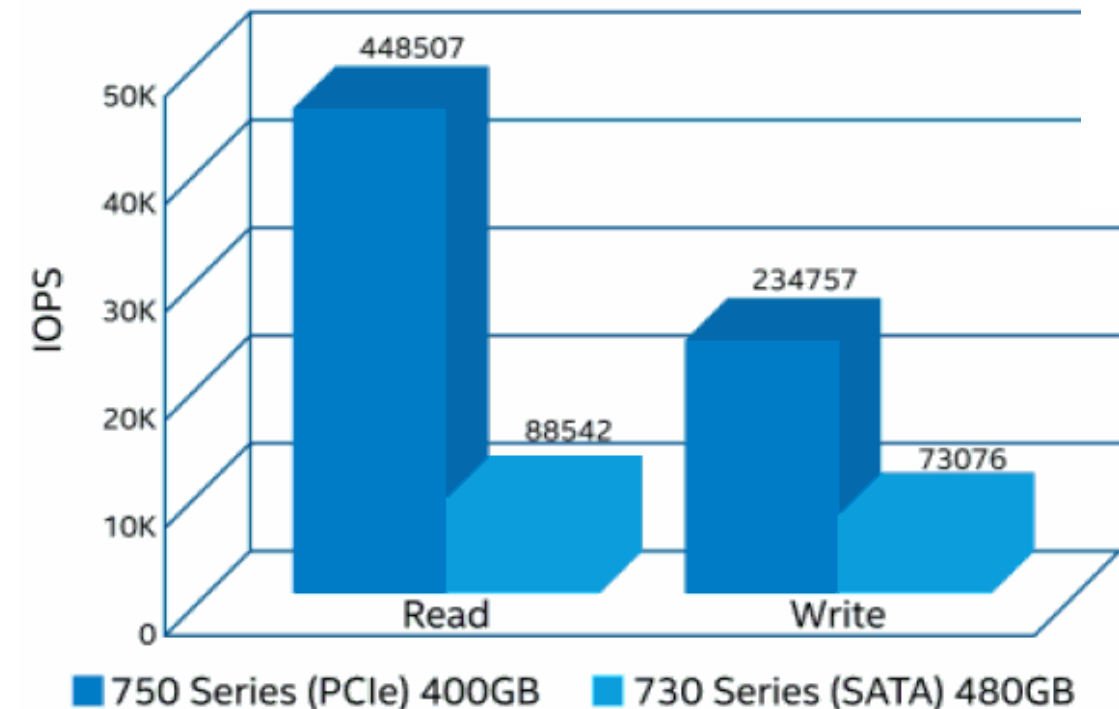


# Storage Technology Evolution





Random Read/Write Performance<sup>†</sup>  
750 Series (PCIe) vs. 730 Series (SATA)



## Specification Strategy and Breakdown of Work

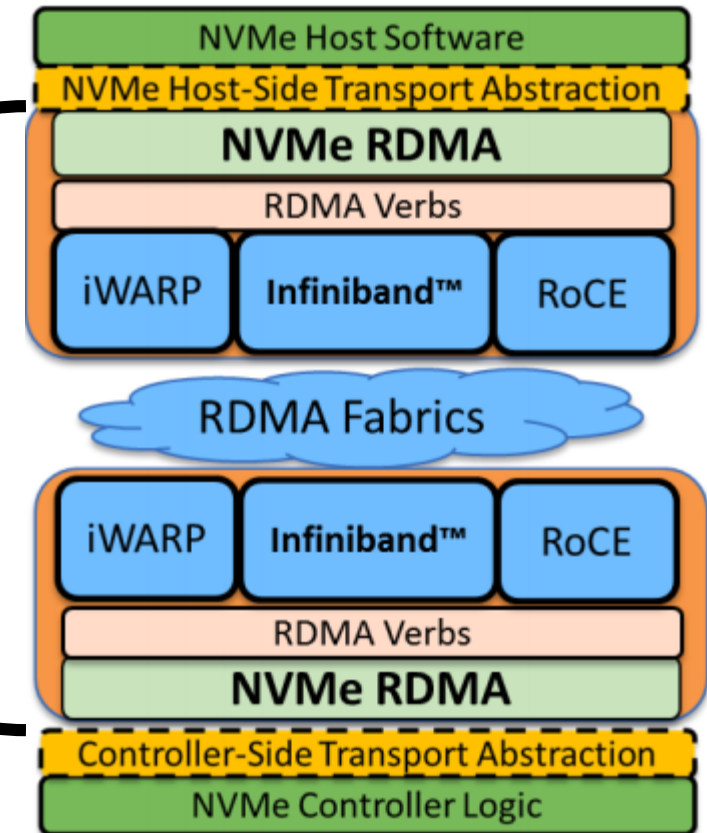
### Do not create a standalone specification

- Initial goal is to minimize changes to existing specification
- Cleanly separate out the non-PCIe NVMe Transport layers through separate chapters/sections
  - Fabrics Core (concepts and RDMA binding)
  - Fabrics Base Differences (SGL changes, etc.)
- Long-term goal is to create a Transport agnostic base spec



### Break the work into functional sub-sections

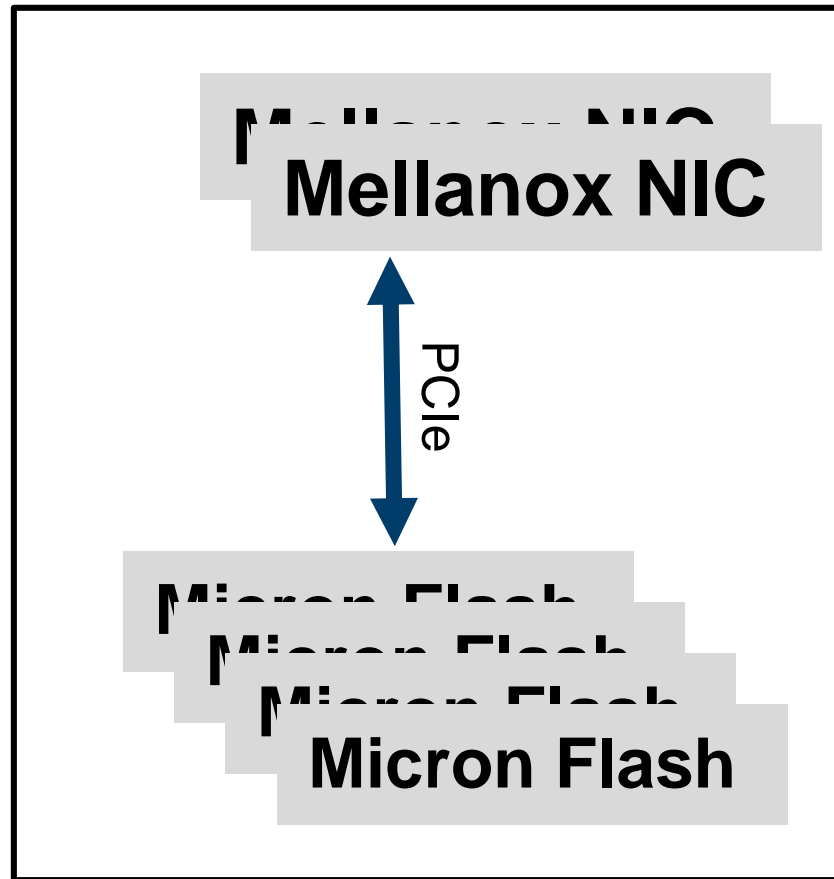
- Capsules
  - Discovery
  - Connections
  - Flow Control
  - Naming
  - Binding
  - Error Handling



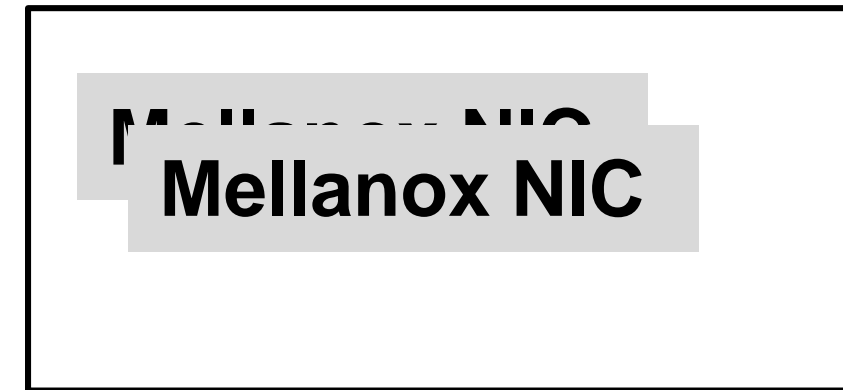
# RDMA-based Remote NVME Access (NVME over Fabrics)



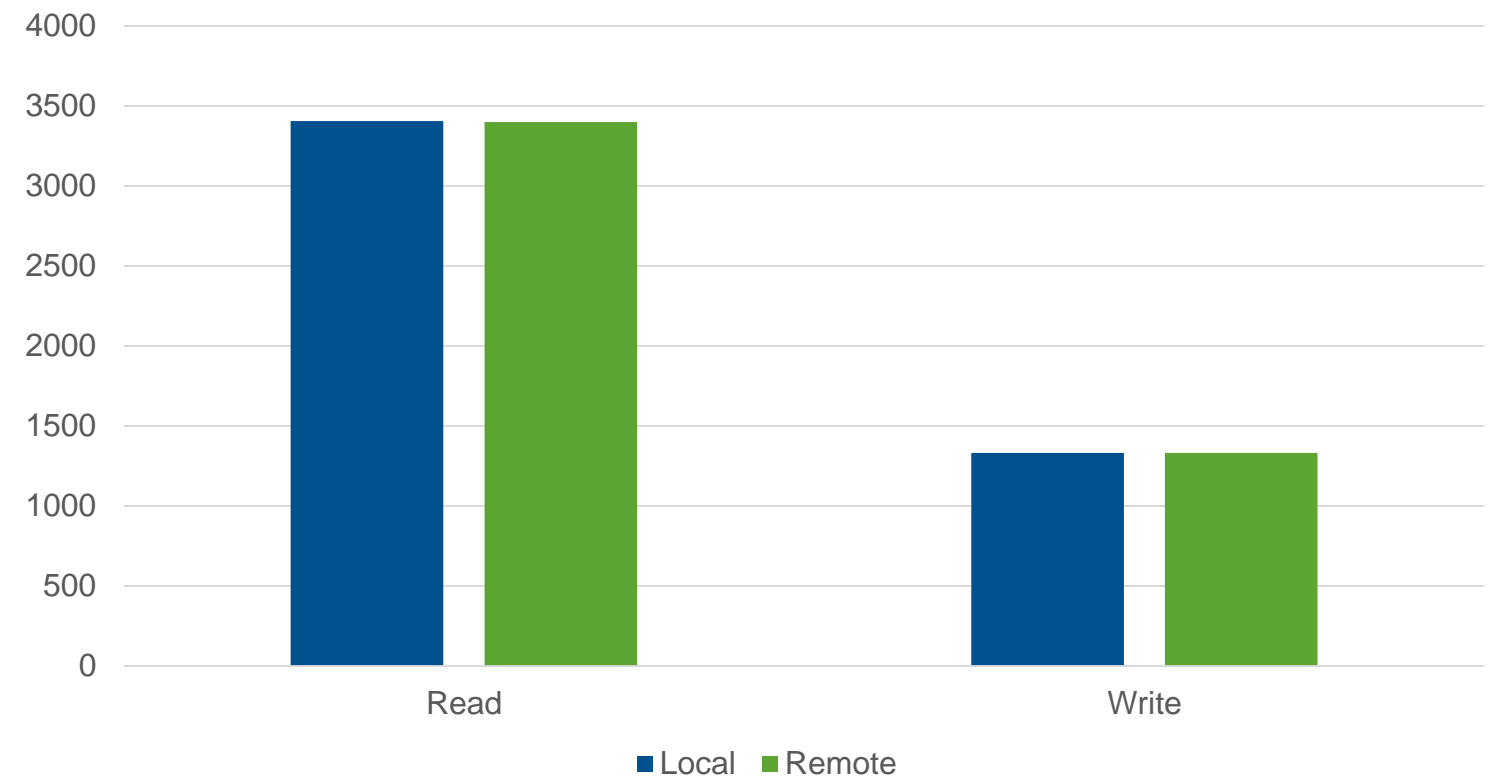
Target Server



Initiator Server



Operations/sec – local vs. remote



# BlueField System-on-a-Chip (SoC) Solution



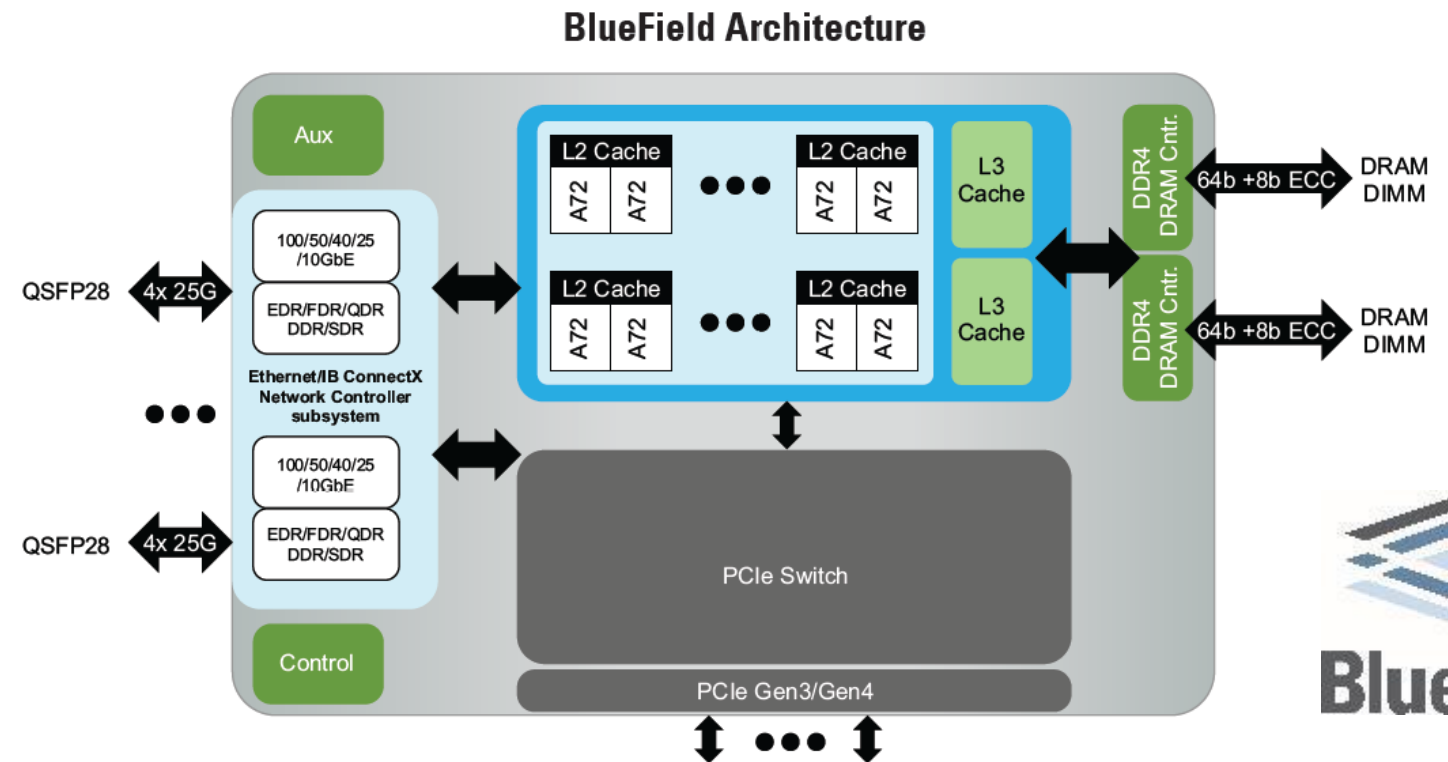
**Storage**

**NVMe Flash Storage Arrays**  
**Scale-Out Storage (NVMe over Fabric)**

**NFV**

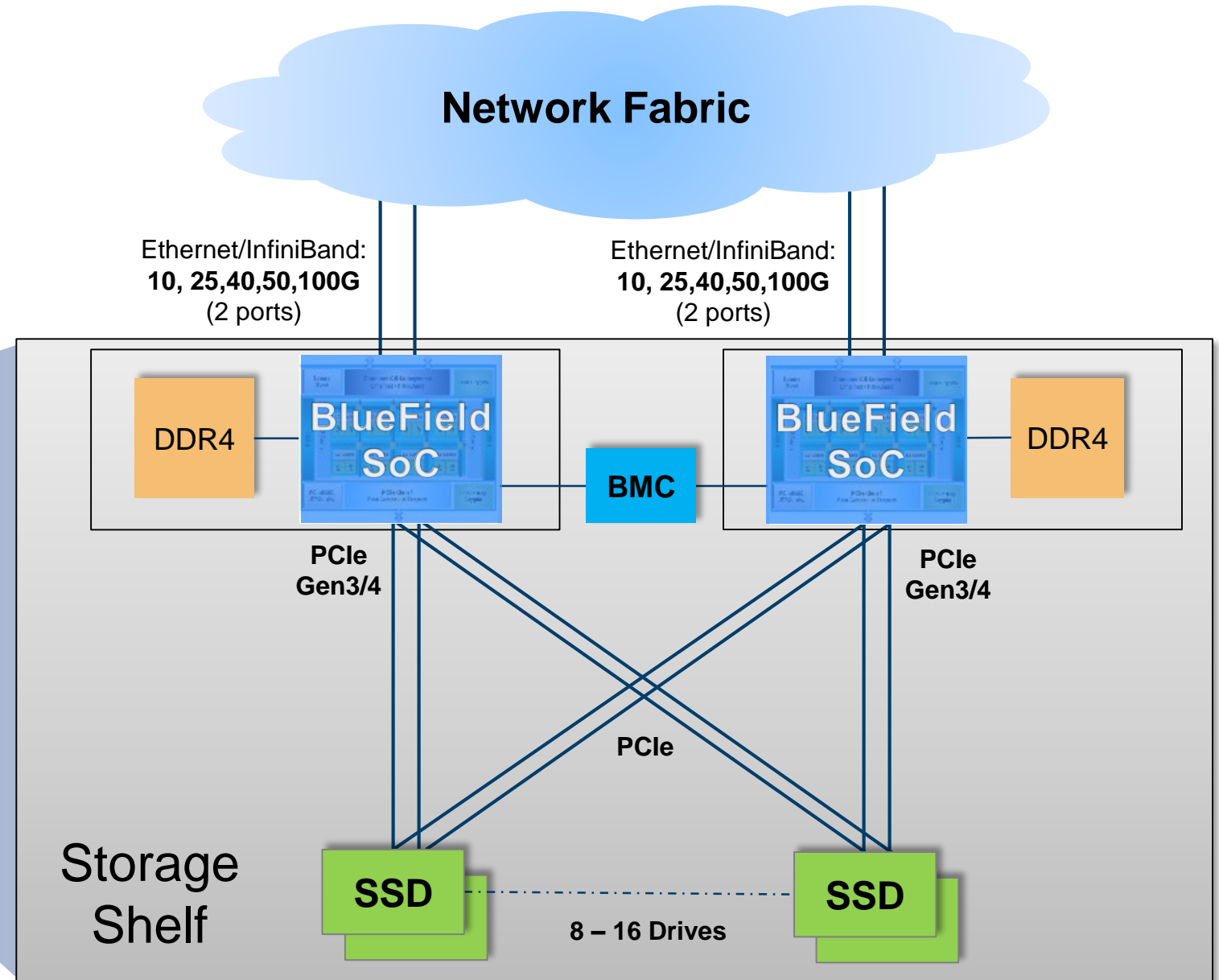
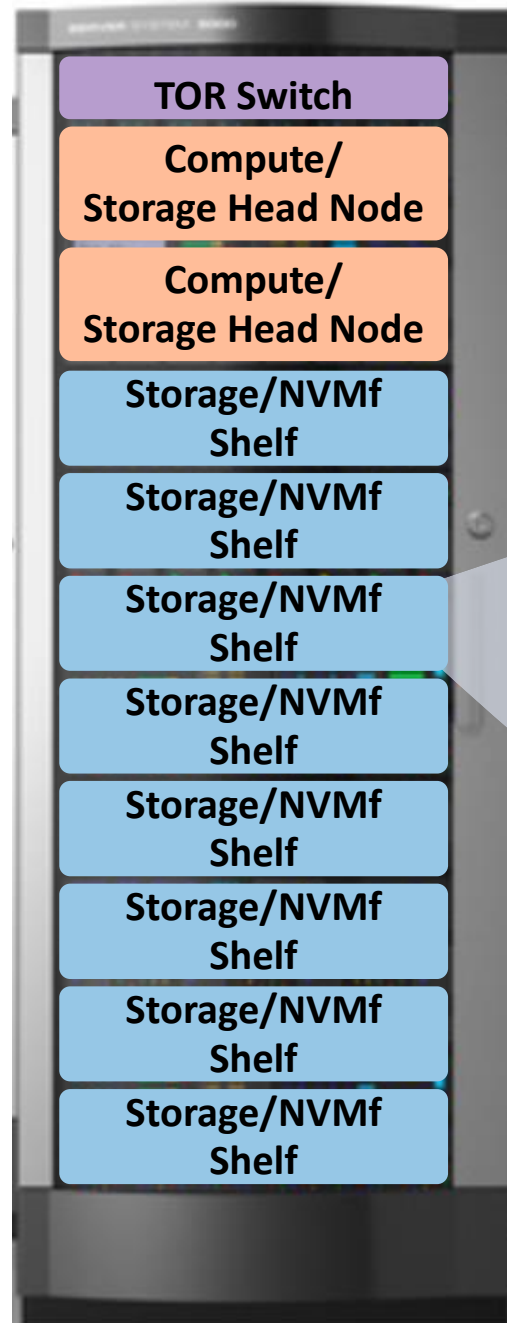
**Accelerating & Virtualizing VNFs**  
**Open vSwitch (OVS), SDN**  
**Overlay networking offloads**

- Integration of ConnectX5 + Multicore ARM
- State of the art capabilities
  - 10 / 25 / 40 / 50 / 100G Ethernet & InfiniBand
  - PCIe Gen3/Gen4
  - Hardware acceleration offload
    - RDMA, RoCE, NVMeF, RAID
- Family of products
  - Range of ARM core counts and I/O ports/speeds
  - Price/Performance points



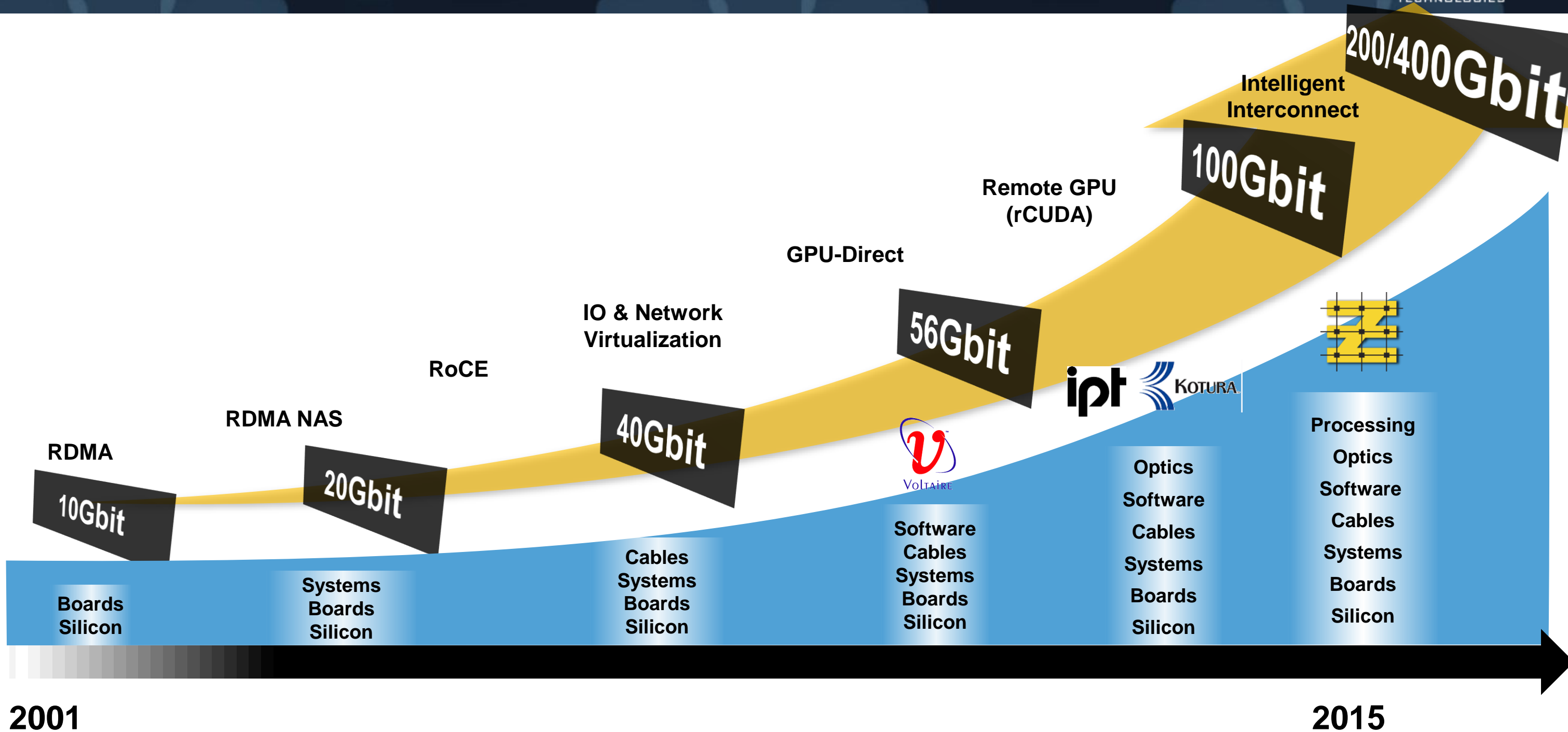
# Scale-Out NVMe Storage System with BlueField

Rack view



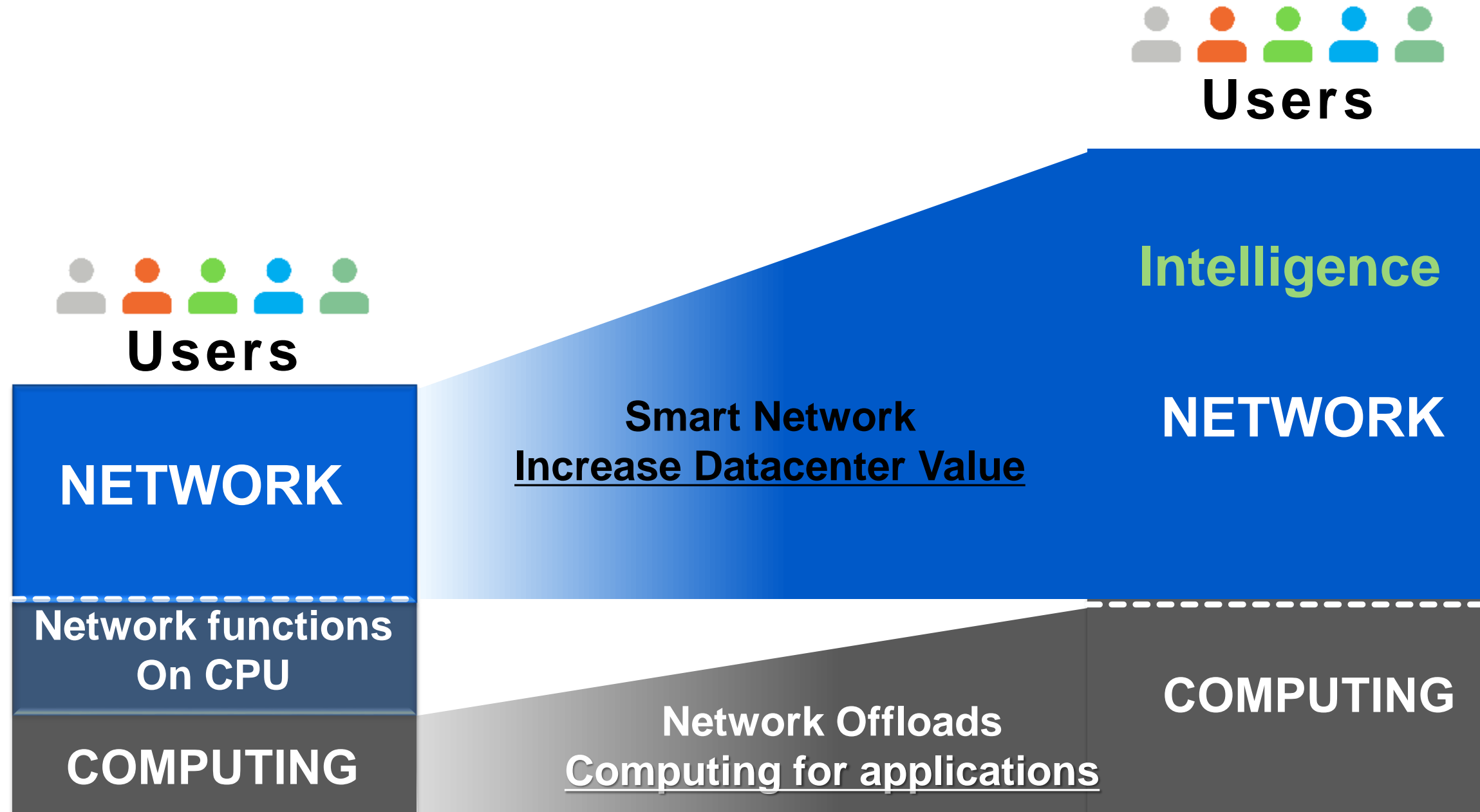


# Technology Leadership



2001

2015





**Bare Metal  
Programming**



**Scale-out  
Services**



Thank You