# Pushing the Limits of AI with In-Network Computing

APNET 2019

Gil  Bloch
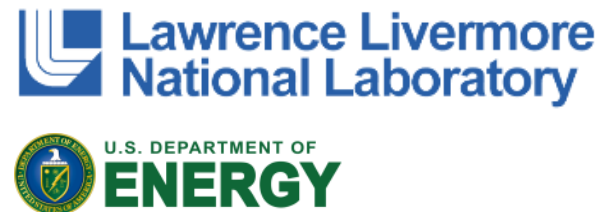
# Mellanox Accelerates Leading HPC and AI Systems

World's Top 3 Supercomputers
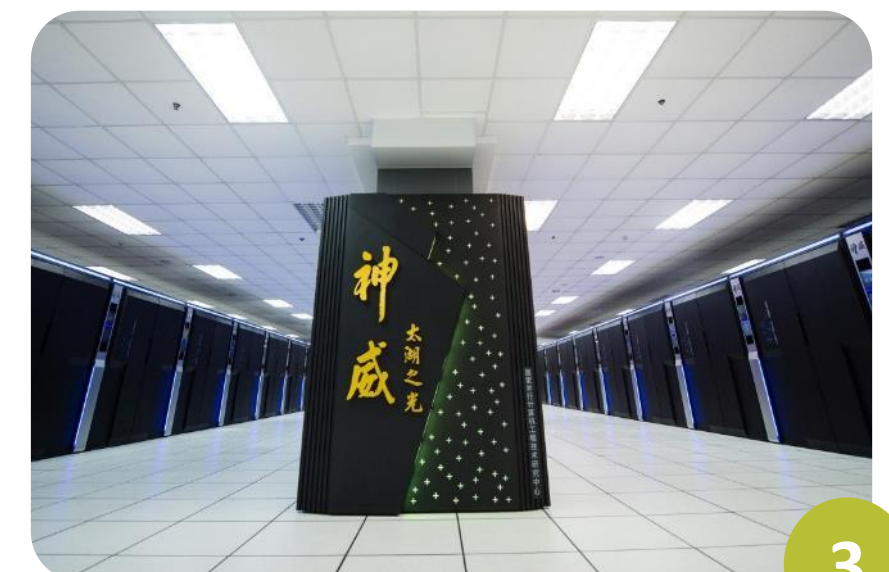


**1**

Summit CORAL System
World's Fastest HPC / AI System
9.2K InfiniBand Nodes



**2**

Sierra CORAL System
#2 USA Supercomputer
8.6K InfiniBand Nodes



**3**

Wuxi Supercomputing Center
Fastest Supercomputer in China
41K InfiniBand Nodes

# Data is Growing Faster Than Ever

Autonomous vehicle generates 4000GByte per day

**CAMERA**
~20-40MB Per/sec
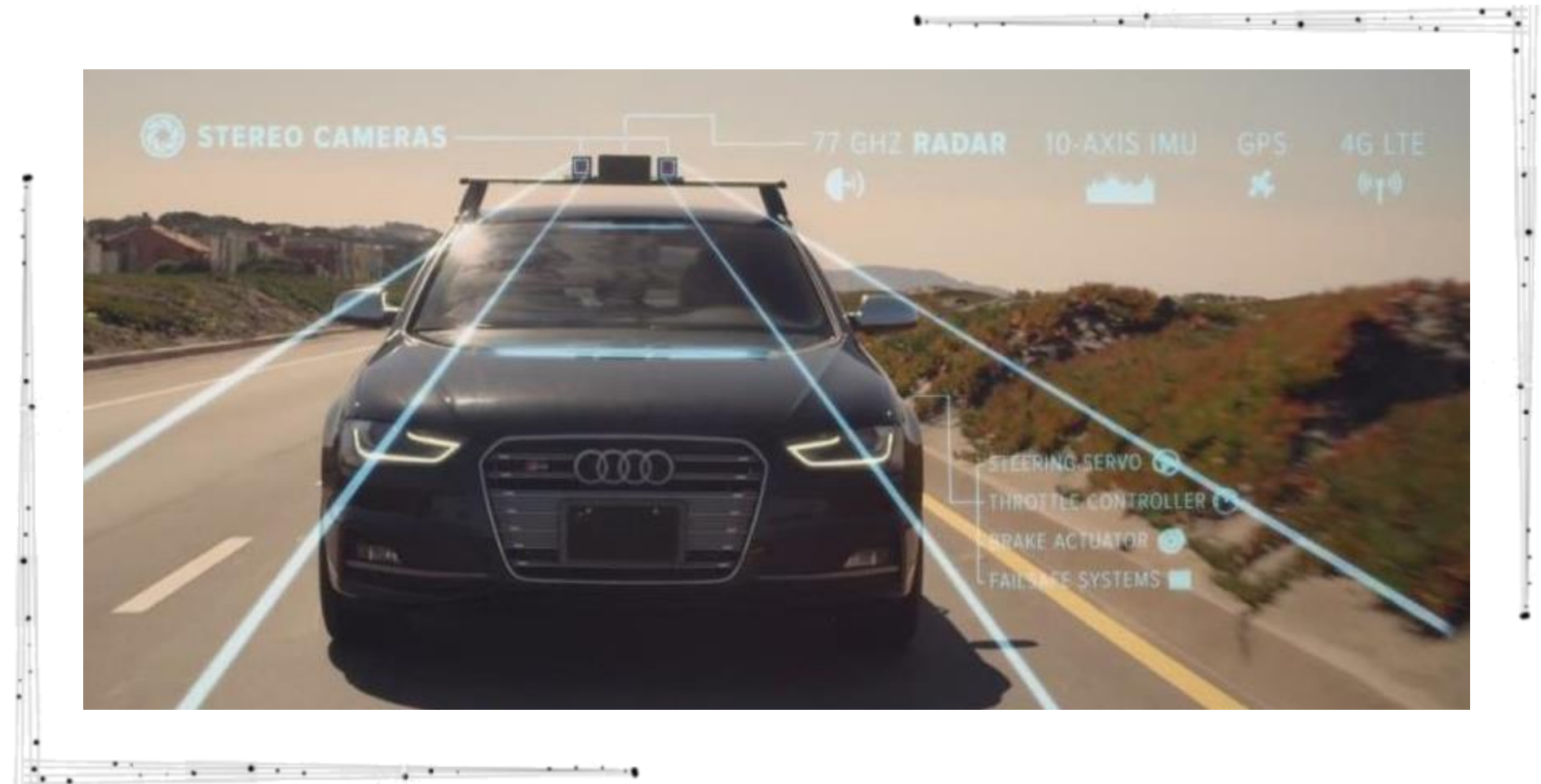
**SONAR**
~10-100KB Per/Sec
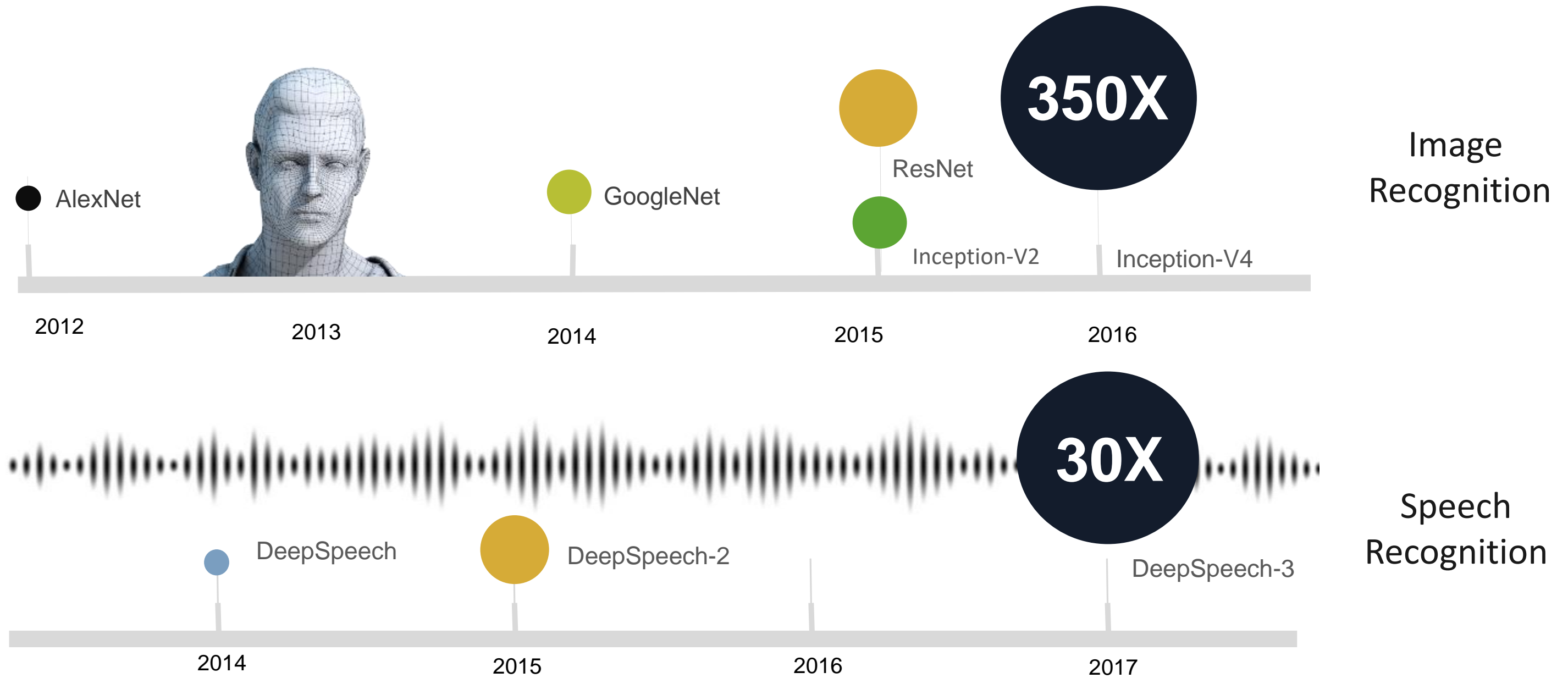
**GPS**
~50KB Per/Sec

**RADAR**
~10-100KB Per/Sec

**Light Detection & Ranging**
~10-70MB Per/Sec



- Data will grow by a factor of 10 over the next decade to 163 Zeta Bytes in 2025 (source: IDC)
- Faster Data processing requires faster Interconnect speeds

# Neural Networks Complexity Growth



**Image Recognition**

AlexNet
GoogleNet
ResNet
Inception-V2
Inception-V4

**350X**

2012    2013    2014    2015    2016

**30X**

DeepSpeech
DeepSpeech-2
DeepSpeech-3

**Speech Recognition**

2014    2015    2016    2017

Complexity = GOPS X Bandwidth

# Enabling World-Leading Artificial Intelligence Solutions
Mellanox Unleashes the Power of Artificial Intelligence

**More
Data**

$+$

**Better
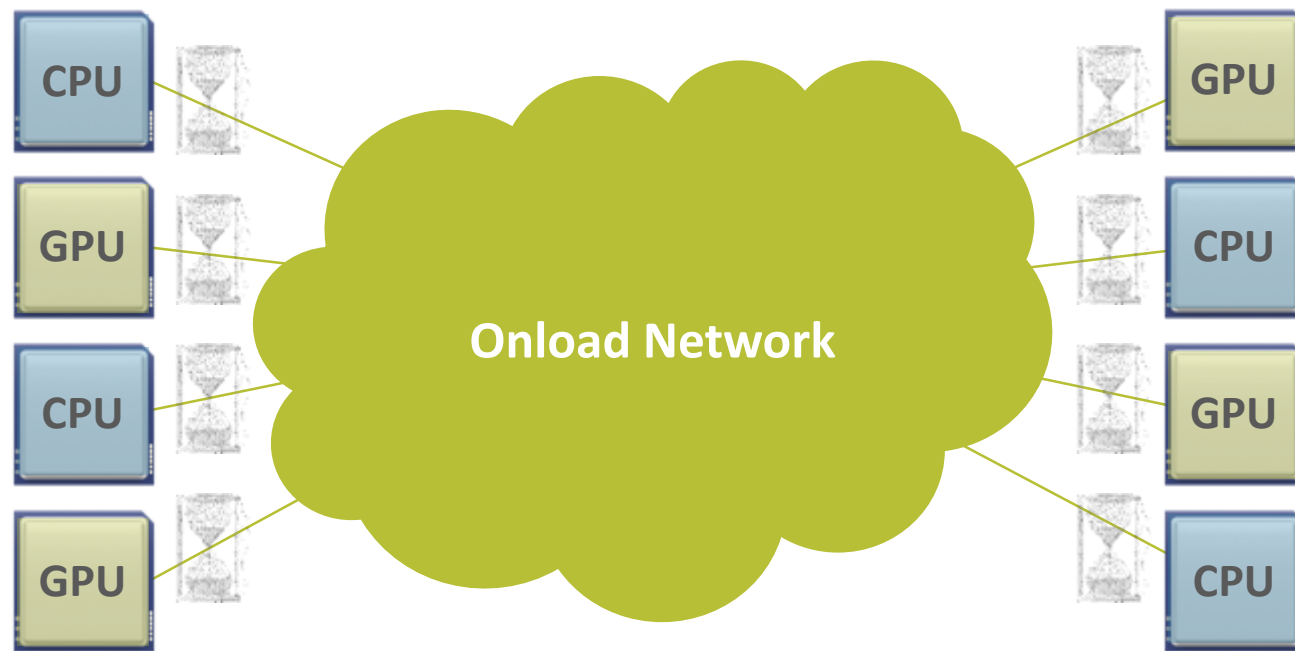Models**

$+$

**Faster
Interconnect**

GPUs

CPUs

ASIC
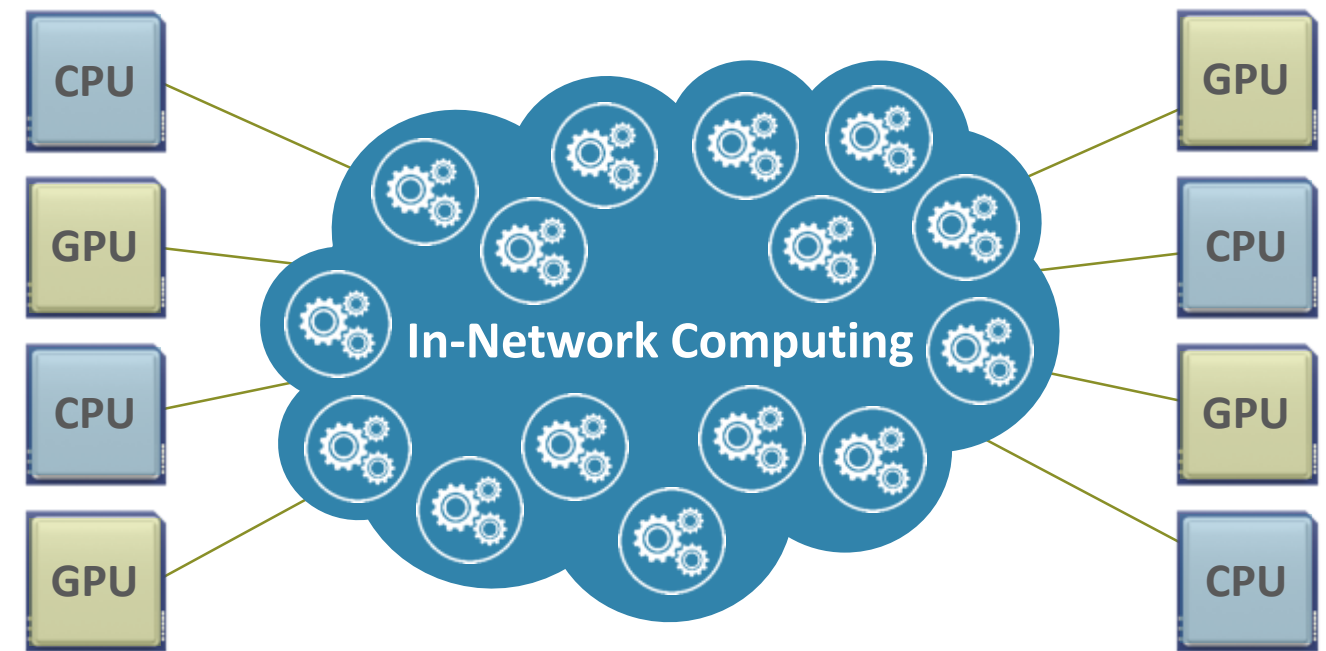
FPGAs

Storage

# The Need for Intelligent and Faster Interconnect

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

**CPU-Centric (Onload)**

**Data-Centric (Offload)**



Must Wait for the Data
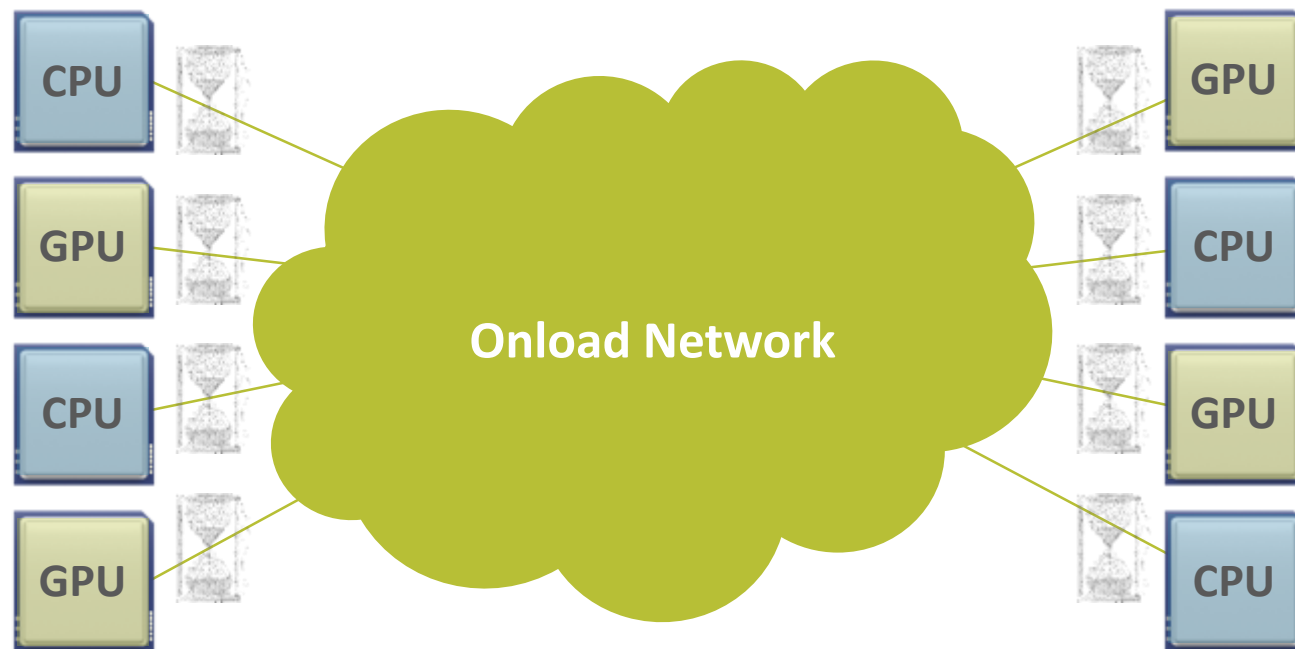Creates Performance Bottlenecks

Analyze Data as it Moves!
Higher Performance and Scale

# An Application Example – Pizza Processing

CPU 1 – Pizza Generation

CPU 2 – Pizza Consumption

- Order Pizza
  - Call (or use Pizza application)
- CPU 1 – prepare Pizza
  - Tomato sauce, Cheese, Peperoni…
- CPU 1 – Put in the oven
  - And now we wait…
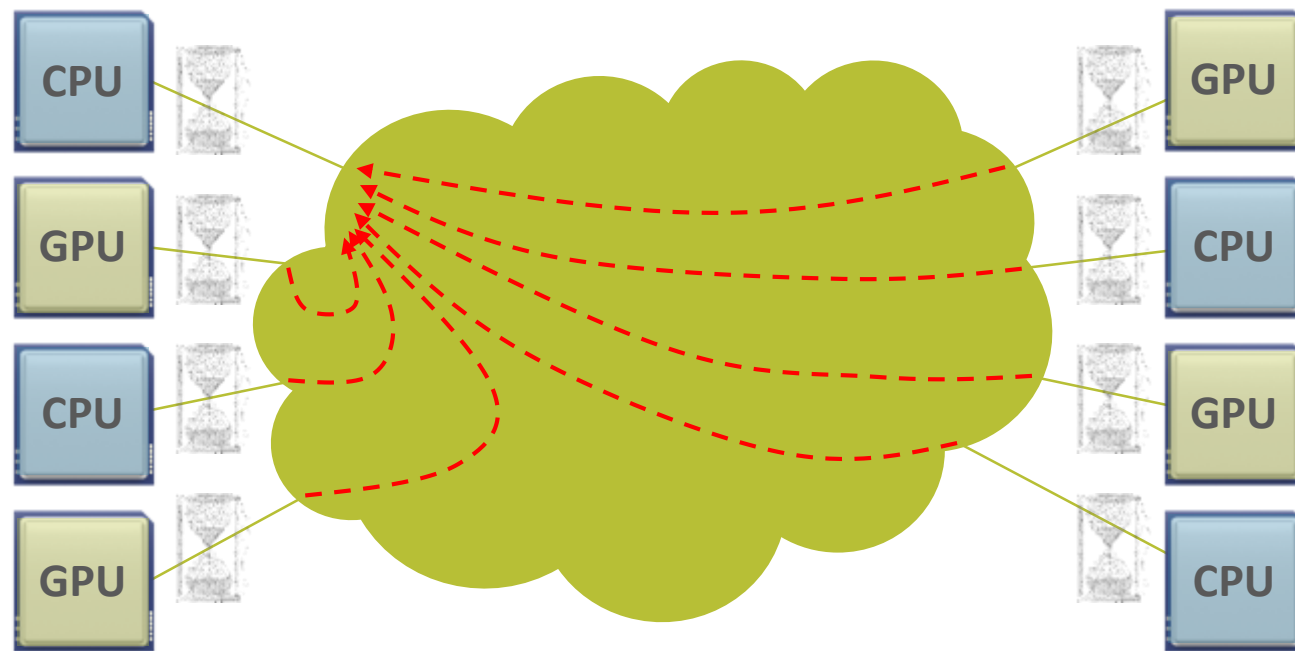- CPU 1 – Pack and send
- Network (Pizza Delivery)

**CPU-Centric (Onload)**



**Onload Network**

CPU  GPU  CPU  GPU

GPU  CPU  GPU  CPU

Must Wait for the Data
Creates Performance Bottlenecks

# What if…

# Data Centric Architecture to Overcome Latency Bottlenecks

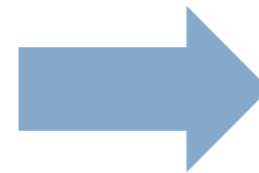Intelligent Interconnect Paves the Road to Exascale Performance

**CPU-Centric (Onload)**
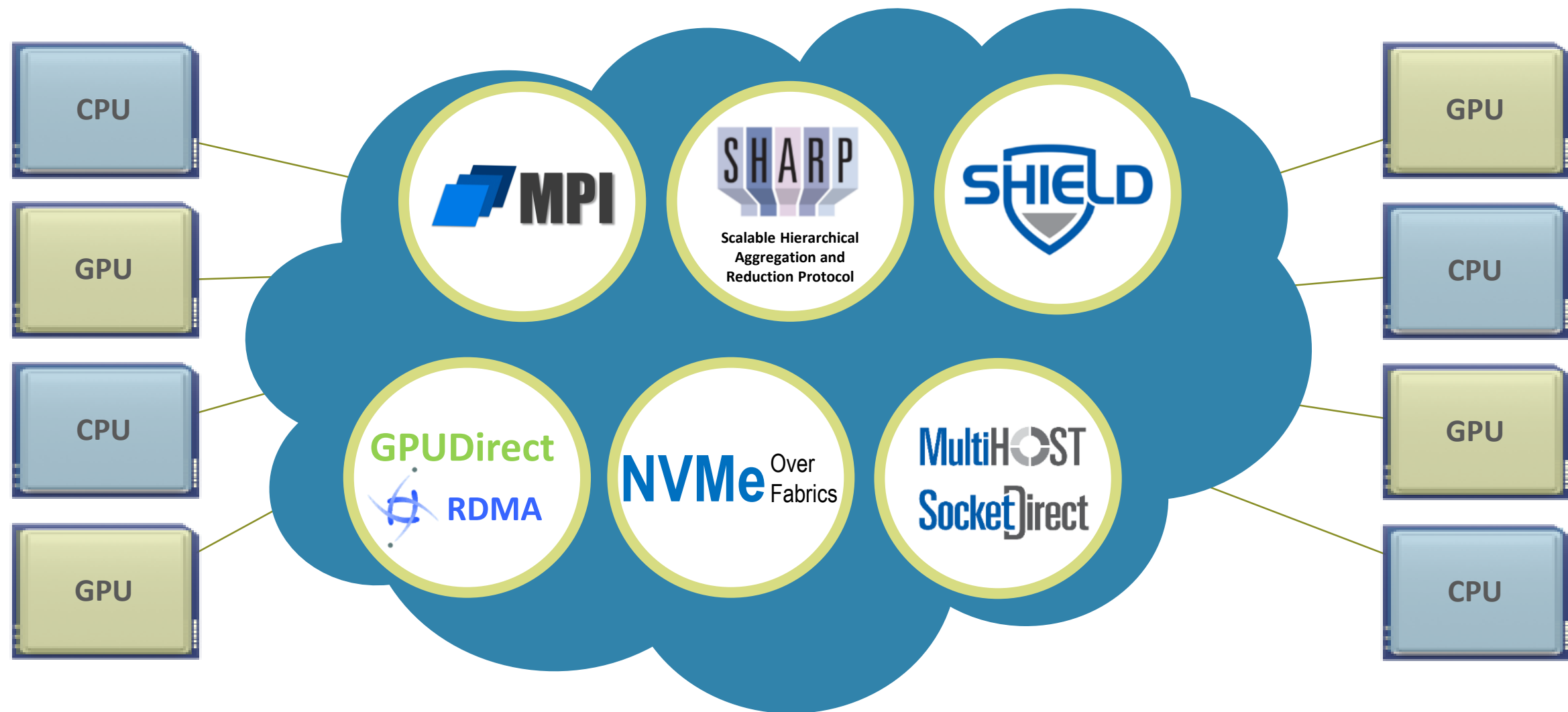
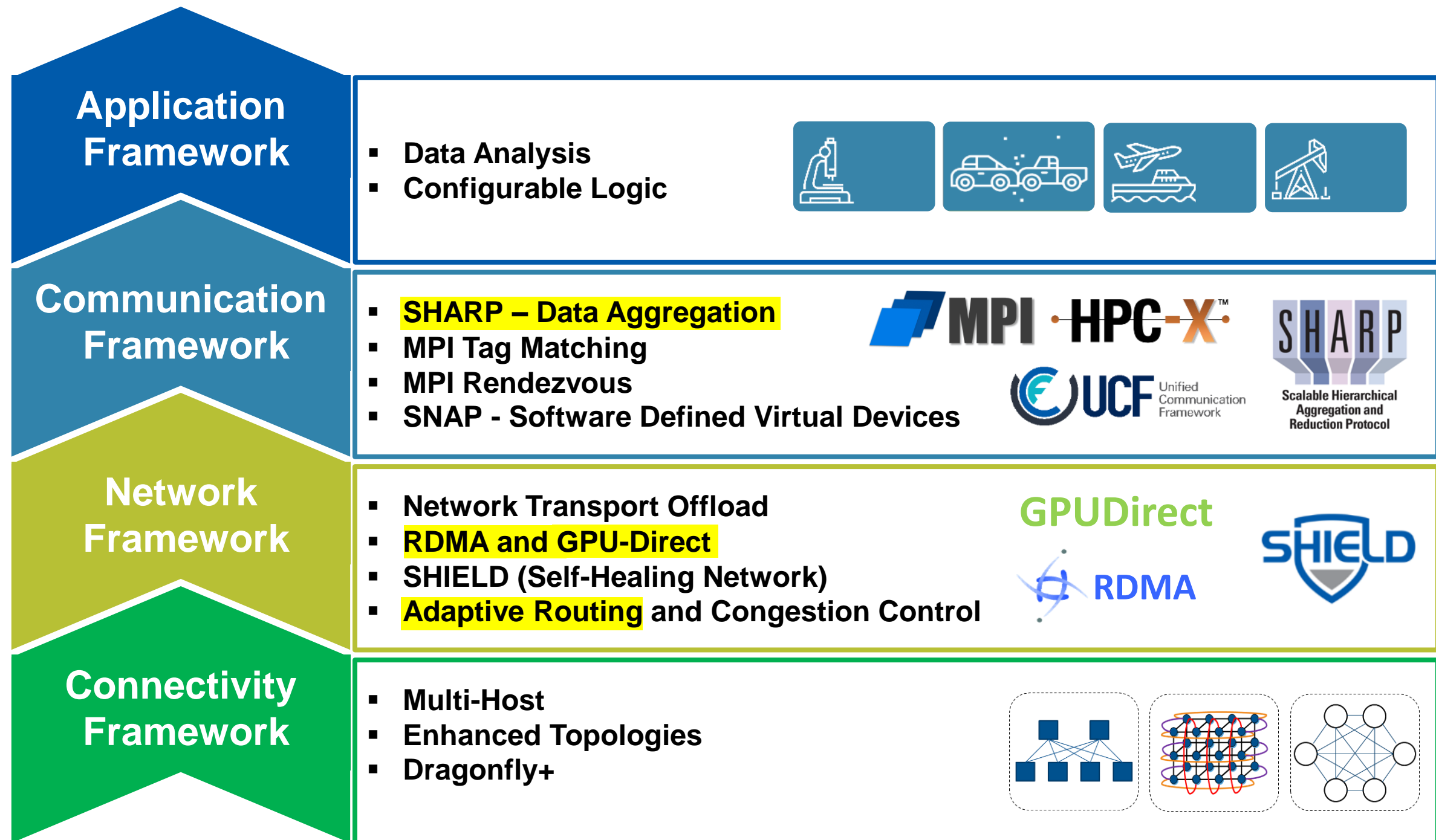**Data-Centric (Offload)**

Communications Latencies
of 30-40us

Communications Latencies
of 3-4us

# In-Network Computing to Enable Data-Centric Data Centers
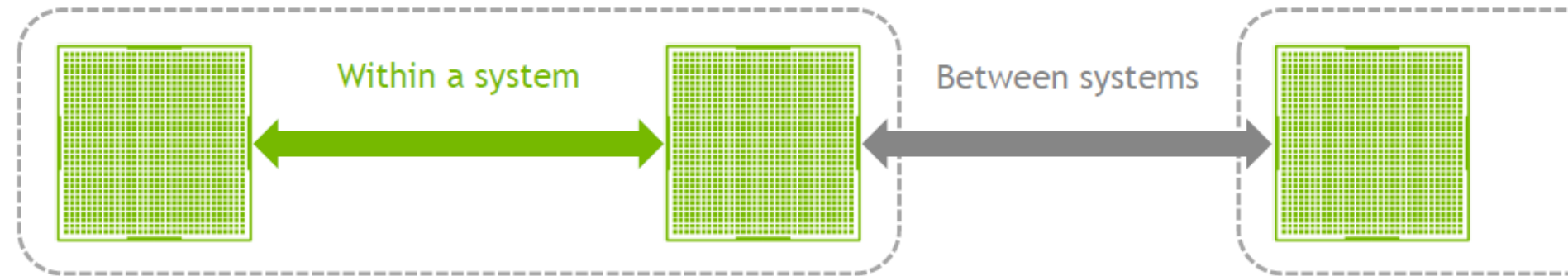
# Accelerating All Levels of HPC/AI Frameworks

**Application Framework**
- Data Analysis
- Configurable Logic

**Communication Framework**
- SHARP – Data Aggregation
- MPI Tag Matching
- MPI Rendezvous
- SNAP - Software Defined Virtual Devices

MPI · HPC-X™

UCF Unified Communication Framework

SHARP
Scalable Hierarchical Aggregation and Reduction Protocol

**Network Framework**
- Network Transport Offload
- RDMA and GPU-Direct
- SHIELD (Self-Healing Network)
- Adaptive Routing and Congestion Control

GPUDirect

RDMA

SHIELD

**Connectivity Framework**
- Multi-Host
- Enhanced Topologies
- Dragonfly+

# The Need for Speed

# Matching Inter and Intra Node Bandwidth

## INTER-GPU COMMUNICATION
### Intra-node and Inter-node

|        | Within a system | | Between systems | |
| --- | --- | --- | --- | --- |

| | |
| --- | --- |
| 6-9 | QPI (shared memory) |
| 9-12 | PCI Express Gen3 x16 (P2P) |
| 62 | NVLink, P100 (P2P) |
| 132 | NVLink, V100 (P2P) |

| | |
| --- | --- |
| 1.2 | 10GbE, TCP/IP Sockets |
| 12 | 100Gb IB or RoCE, RDMA (IB verbs) |
| 47 | 4x 100Gb (DGX1) |
| 82 | 8x 100Gb (DGX2) |

*Effective bandwidth in GB/s*

6  NVIDIA

**NVIDIA**

## S9656 - DISTRIBUTED TRAINING AND FAST INTER-GPU COMMUNICATION WITH NCCL
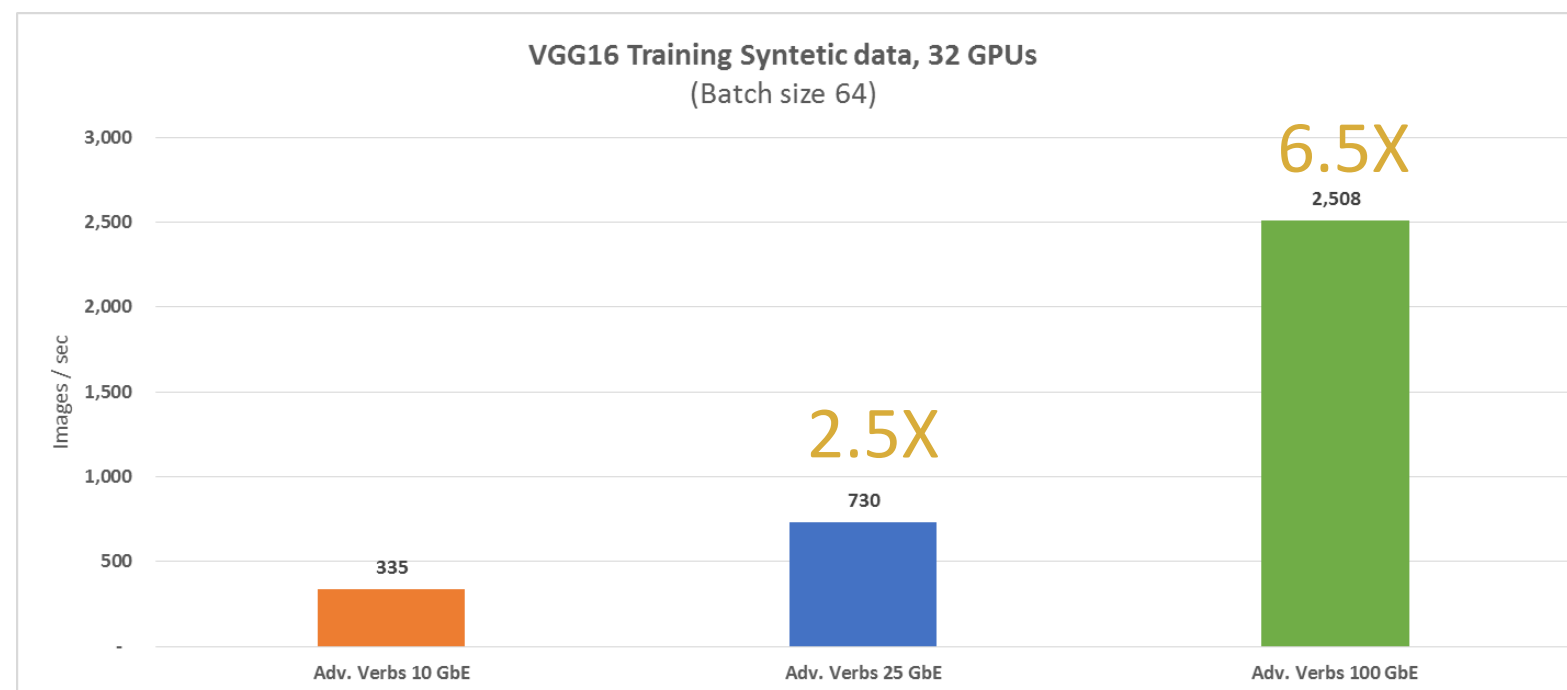Sylvain Jeaugey, NVIDIA

# Mellanox Accelerates TensorFlow 1.5

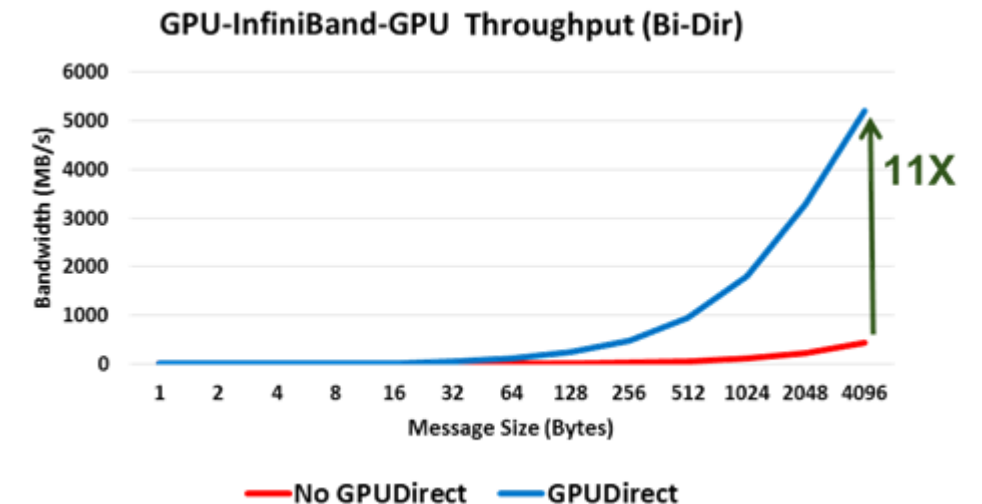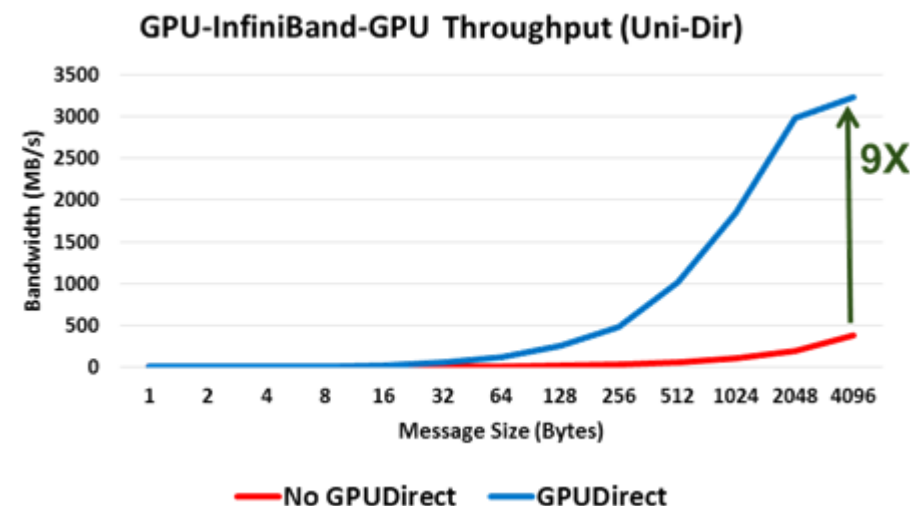## 100G is a Must For Large Scale Models

# 6.5X
### Faster Training with 100G

**VGG16 Training Syntetic data, 32 GPUs**
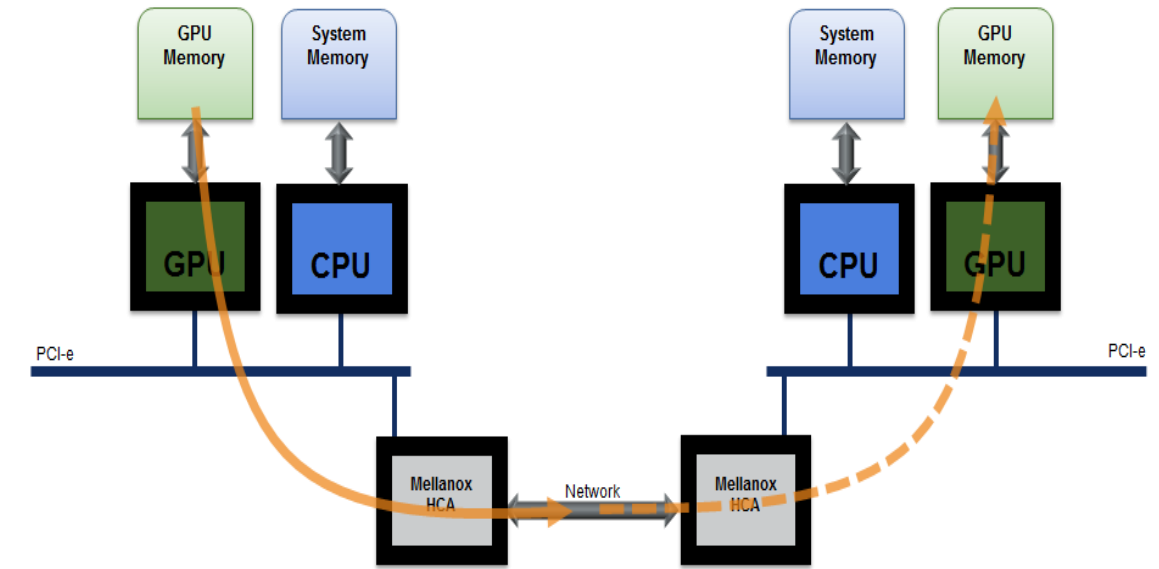(Batch size 64)

6.5X
2,508

2.5X
730

335

Adv. Verbs 10 GbE          Adv. Verbs 25 GbE          Adv. Verbs 100 GbE

Images / sec

3,000
2,500
2,000
1,500
1,000
500
-

# RDMA and GPUDirect

# 10X Higher Performance with GPUDirect™ RDMA

- Accelerates HPC and Deep Learning performance
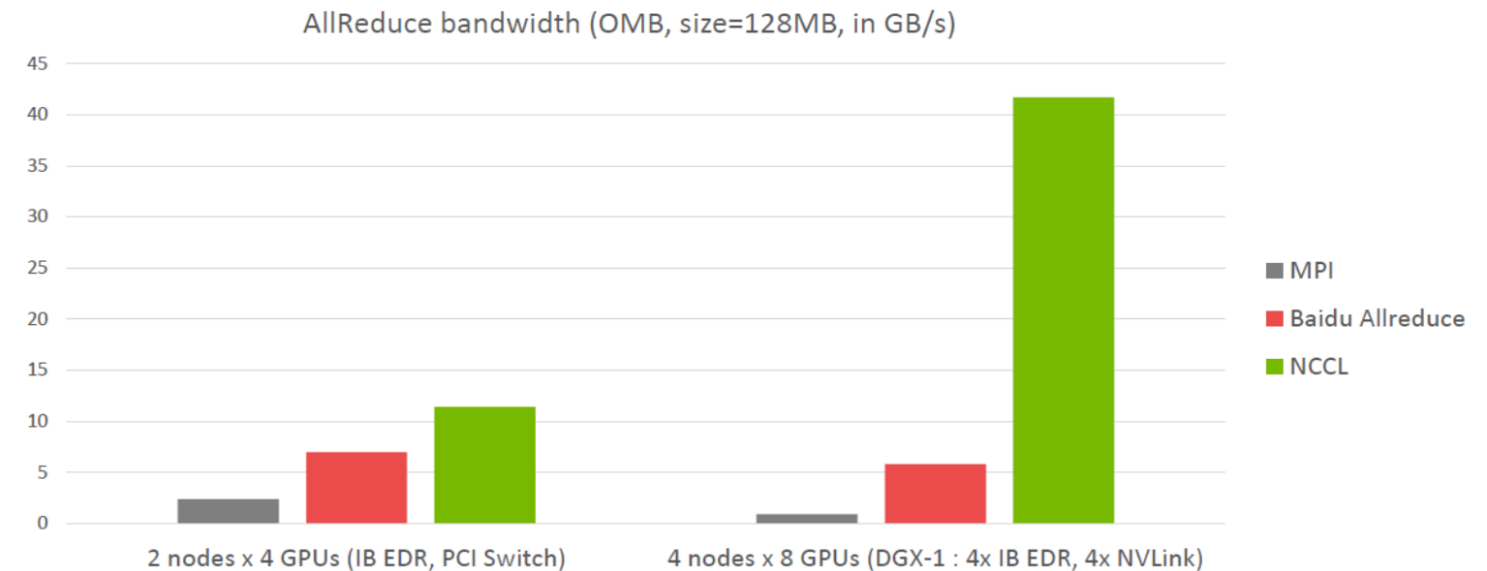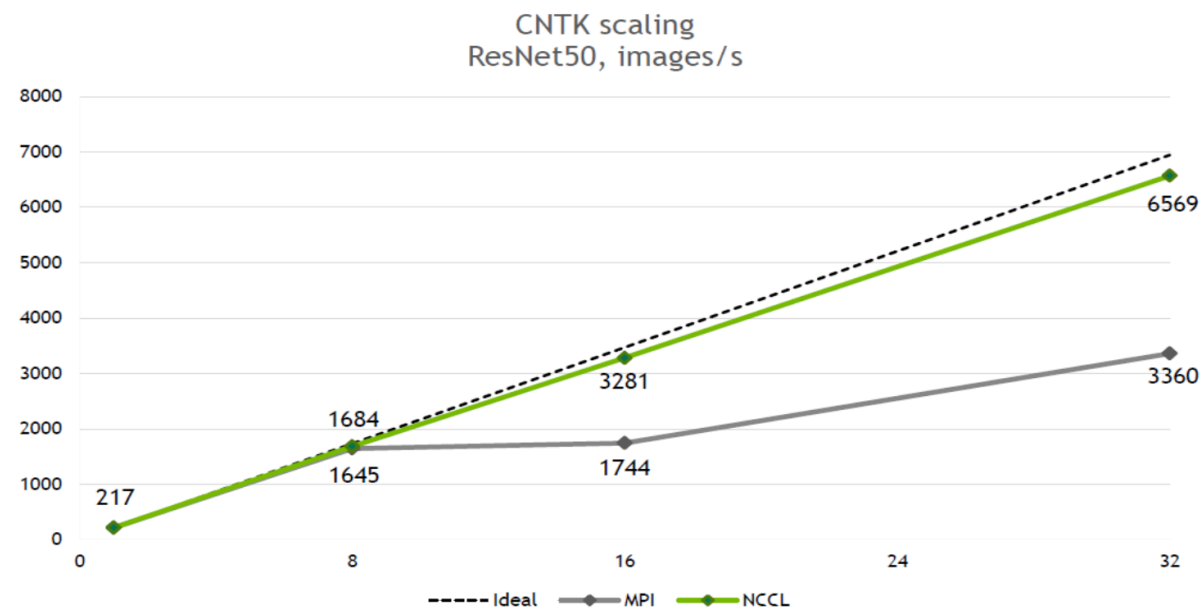
- Lowest communication latency for GPUs

GPUDirect™ RDMA

GPU-InfiniBand-GPU Latency

10X

1.88 usec

No GPUDirect — GPUDirect

GPU-InfiniBand-GPU Throughput (Uni-Dir)

9X

No GPUDirect — GPUDirect

GPU-InfiniBand-GPU Throughput (Bi-Dir)

11X

No GPUDirect — GPUDirect

# Mellanox Accelerates NVIDIA NCCL 2.0

## 50% Performance Improvement

with NVIDIA® DGX-1 across 32 NVIDIA Tesla V100 GPUs Using InfiniBand RDMA and GPUDirect™ RDMA



CNTK scaling
ResNet50, images/s

217
1645
1684
1744
3281
3360
6569

Ideal — MPI — NCCL



AllReduce bandwidth (OMB, size=128MB, in GB/s)

MPI
Baidu Allreduce
NCCL

2 nodes x 4 GPUs (IB EDR, PCI Switch)          4 nodes x 8 GPUs (DGX-1 : 4x IB EDR, 4x NVLink)
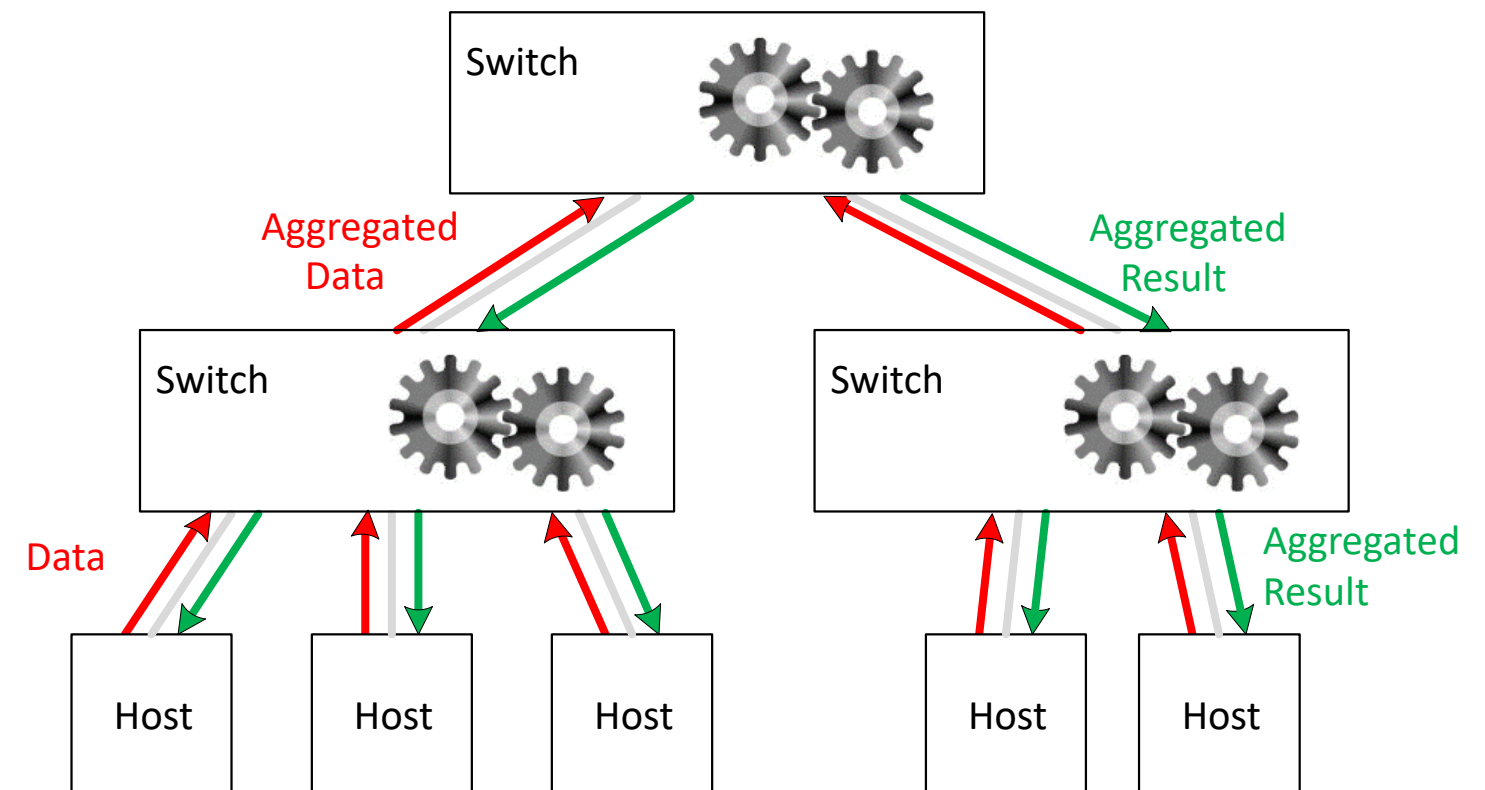
# Scalable Hierarchical Aggregation and Reduction Protocol (SHARP)
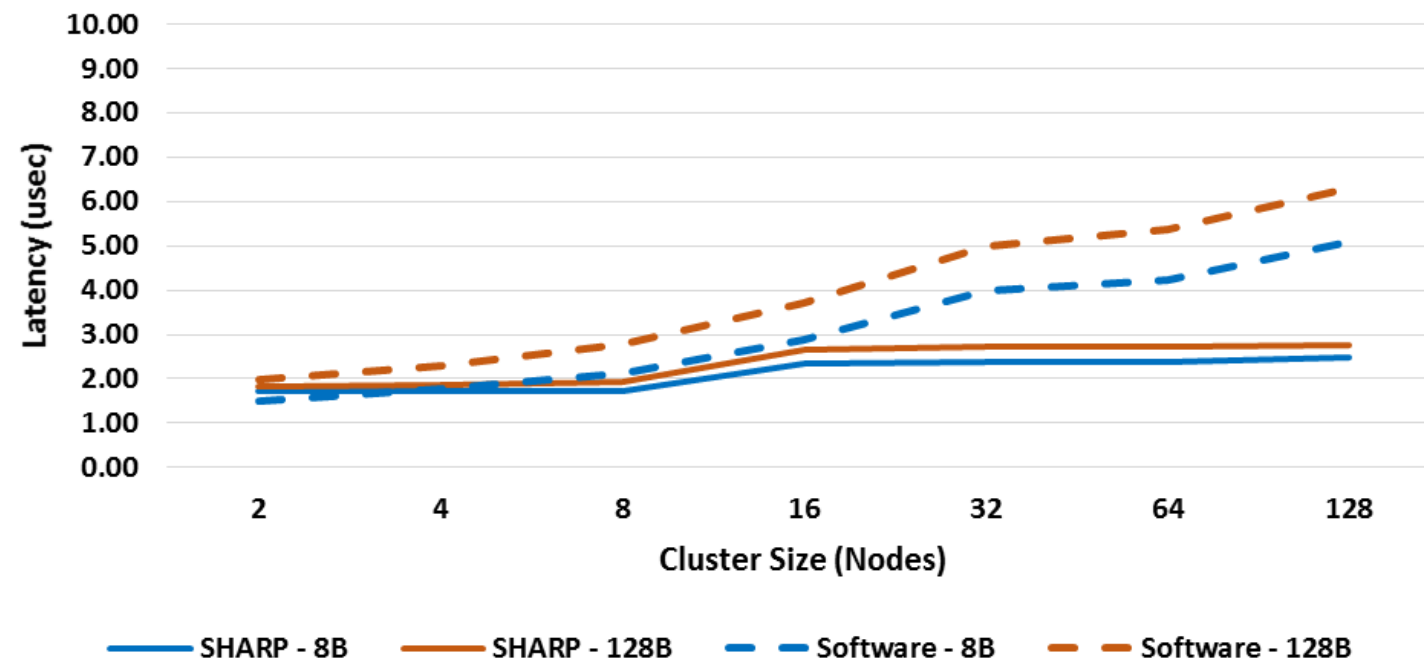
# Scalable Hierarchical Aggregation and Reduction Protocol (SHARP)

- Reliable Scalable General Purpose Primitive

- Applicable to Multiple Use-cases in ML/HPC

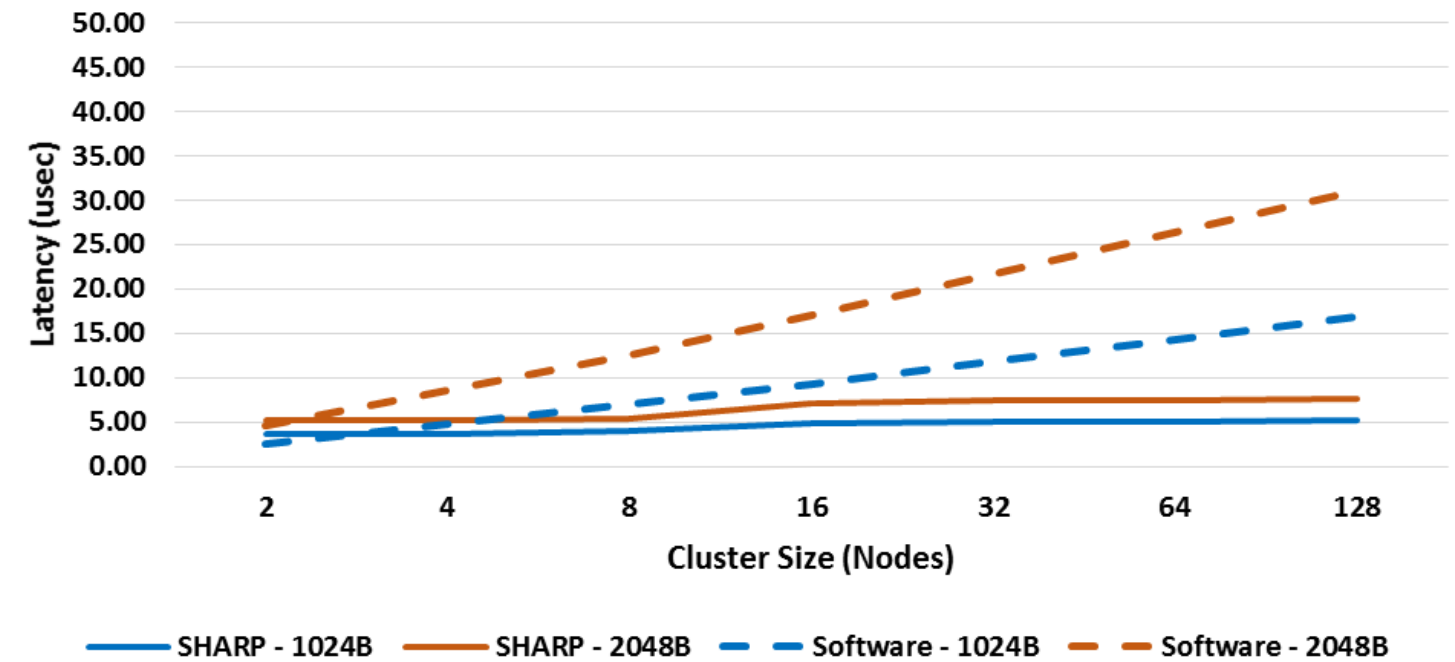- Scalable High Performance Collective Offload

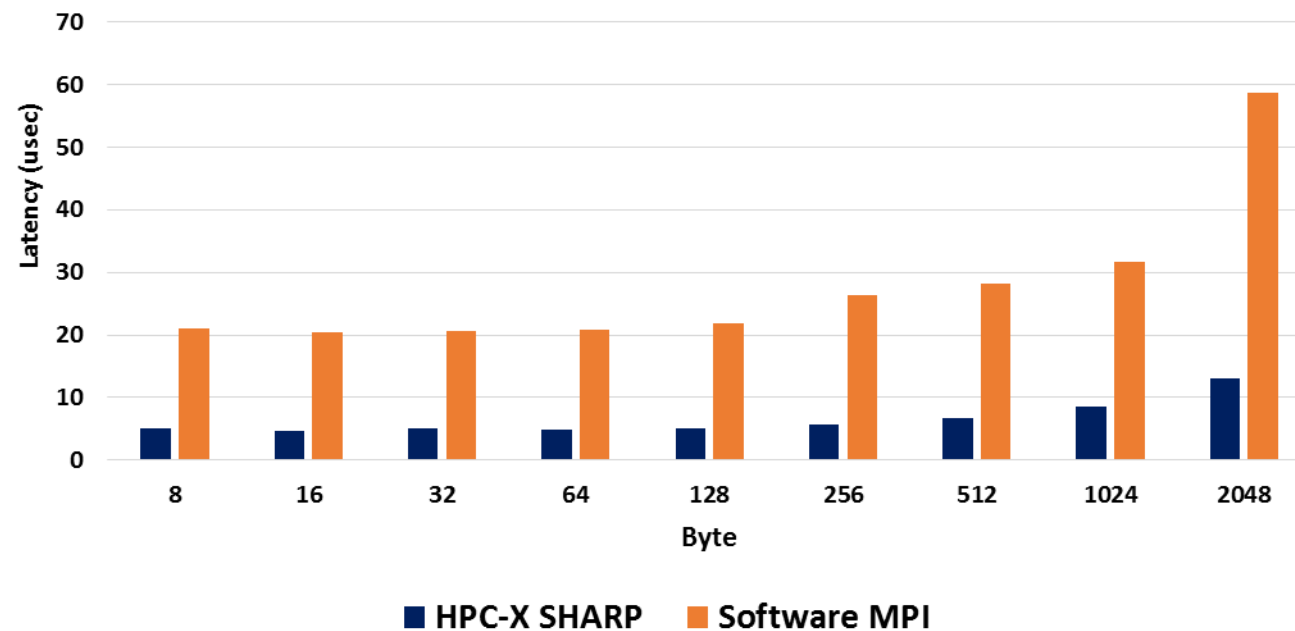# SHARP AllReduce Performance Advantages (128 Nodes)


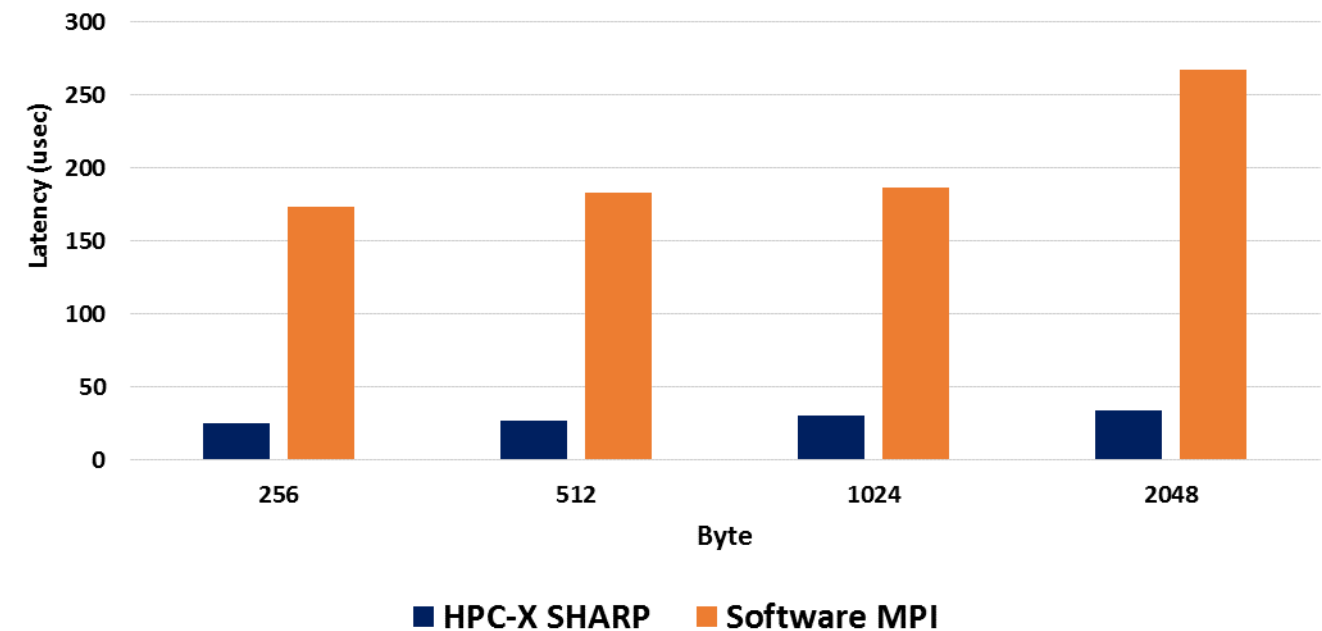
SHARP enables 75% Reduction in Latency
Providing Scalable Flat Latency

# SHARP AllReduce Performance Advantages
## 1500 Nodes, 60K MPI Ranks, Dragonfly+ Topology
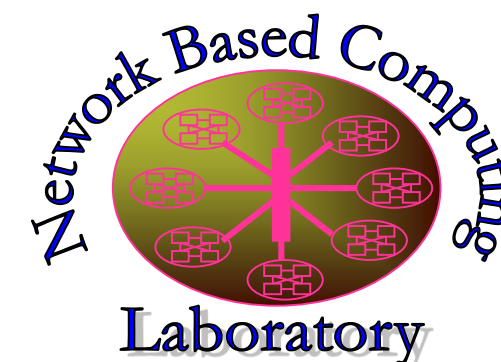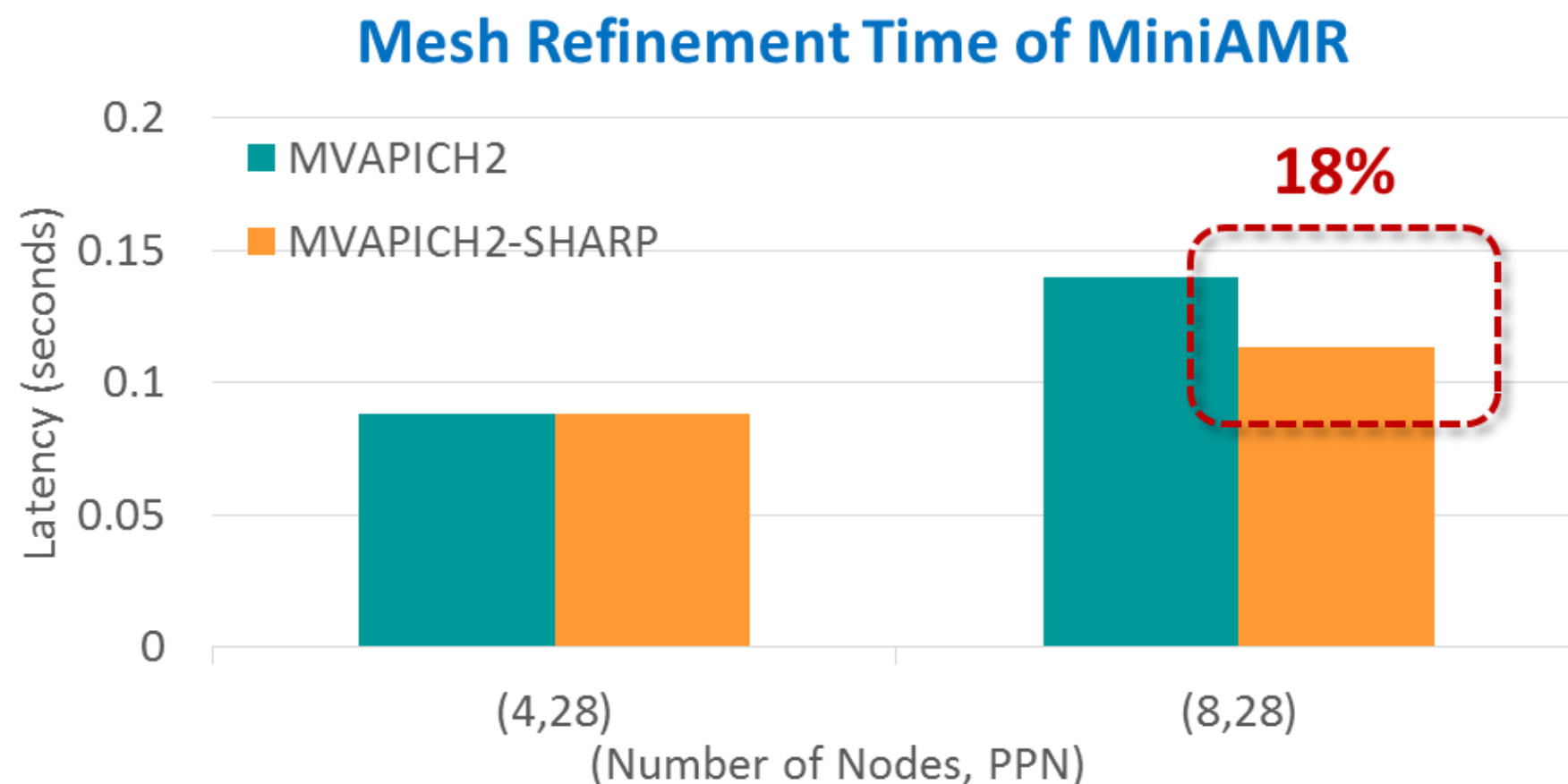
**MPI AllReduce Latency**
**1500 Nodes, 1PPN**



■ HPC-X SHARP  ■ Software MPI

**MPI AllReduce Latency**
**1500 Nodes, 40PPN, 60K MPI Ranks**



■ HPC-X SHARP  ■ Software MPI

**SHARP**
Scalable Hierarchical
Aggregation and
Reduction Protocol

## SHARP Enables Highest Performance

# SHARP Performance – Application (OSU)



Mesh Refinement Time of MiniAMR

- MVAPICH2
- MVAPICH2-SHARP

18%

Latency (seconds)
0.2
0.15
0.1
0.05
0

(4,28)          (8,28)
(Number of Nodes, PPN)

Network-Based Computing Laboratory
http://nowlab.cse.ohio-state.edu/

The MVAPICH2 Project
http://mvapich.cse.ohio-state.edu/

**Source: Prof. DK Panda, Ohio State University**

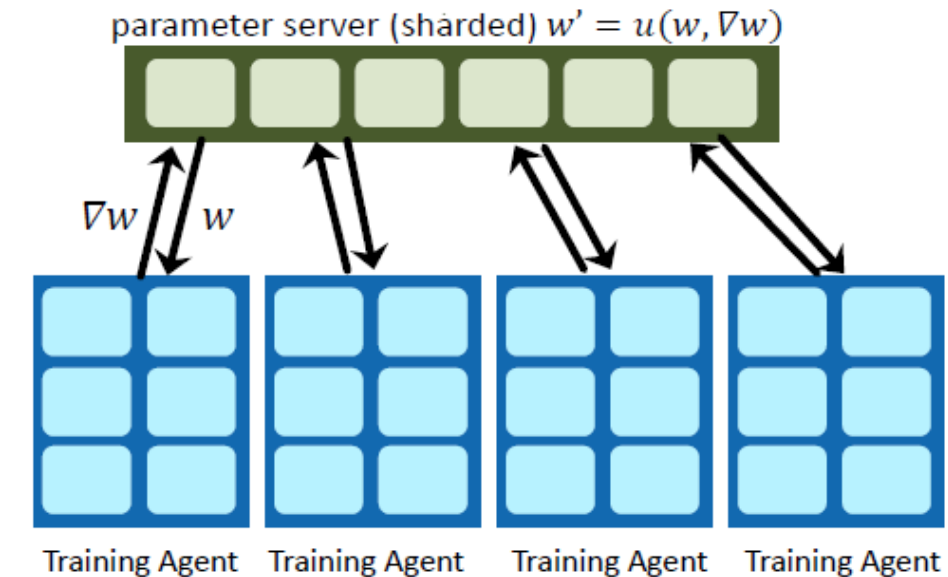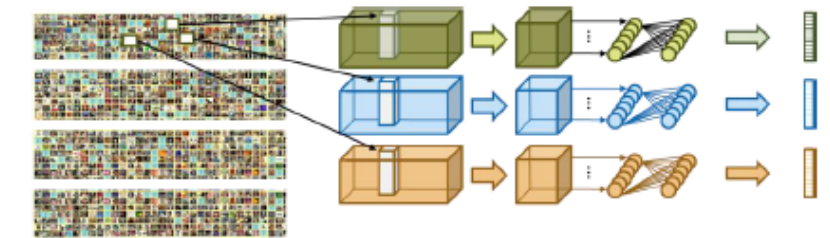# SHARP Accelerates AI Performance

The CPU in a parameter server
becomes the bottleneck

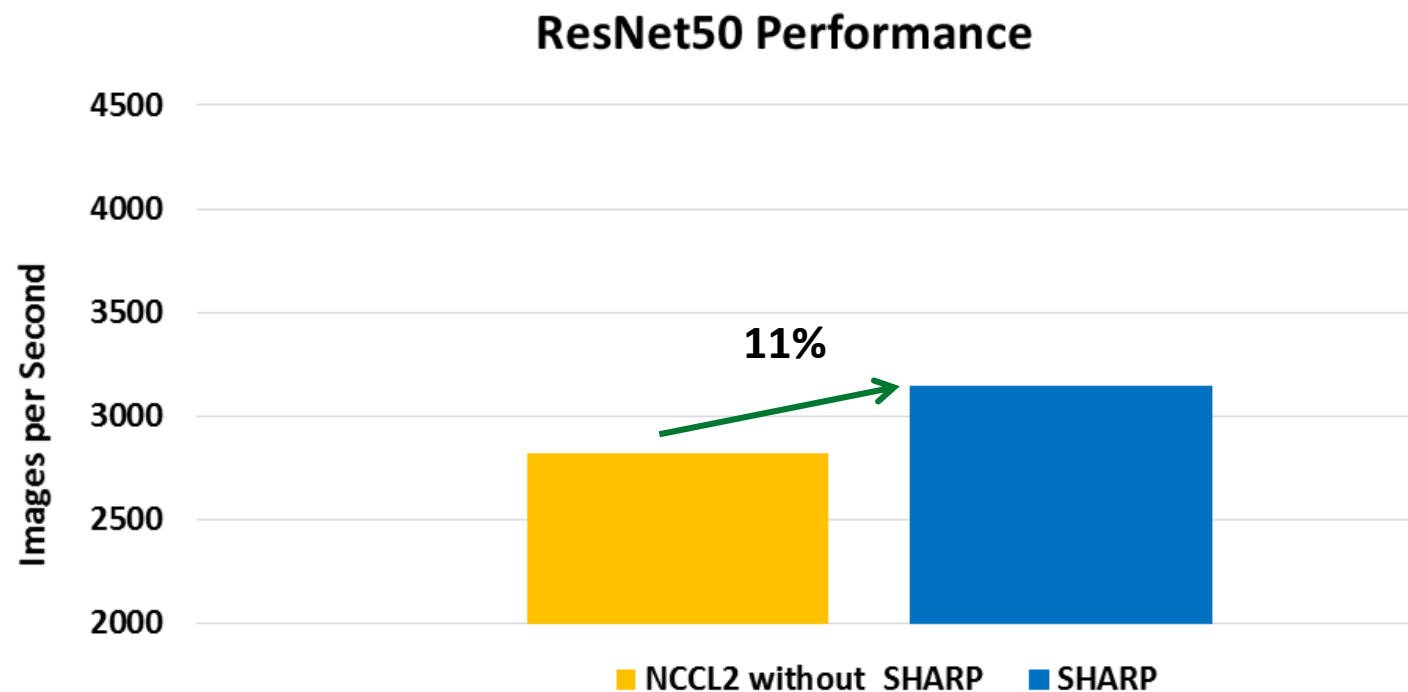**SHARP** Scalable Hierarchical Aggregation and Reduction Protocol

Performs the Gradient Averaging
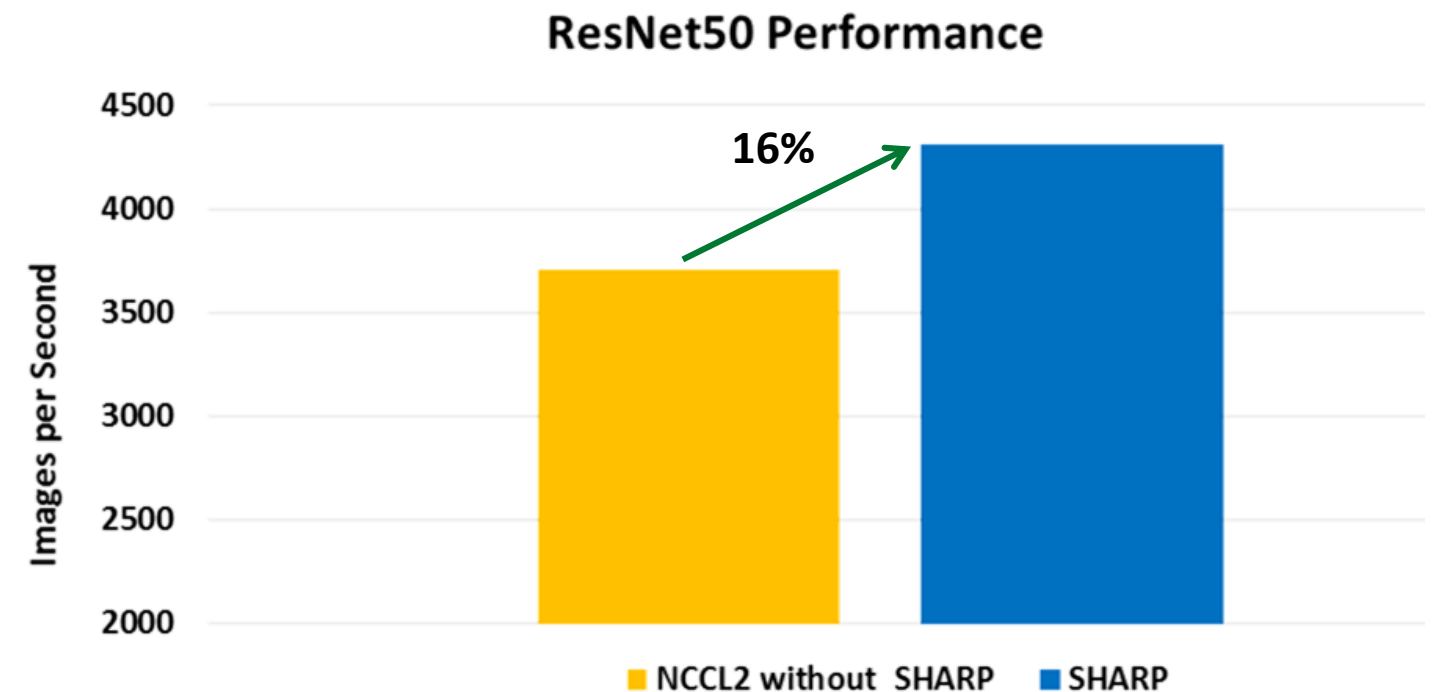Replaces all physical parameter servers
Accelerate AI Performance

parameter server (sharded) $w' = u(w, \nabla w)$

$\nabla w$ $w$

Training Agent    Training Agent    Training Agent    Training Agent

# InfiniBand SHARP Advantage for Deep Learning

- Increase System Performance
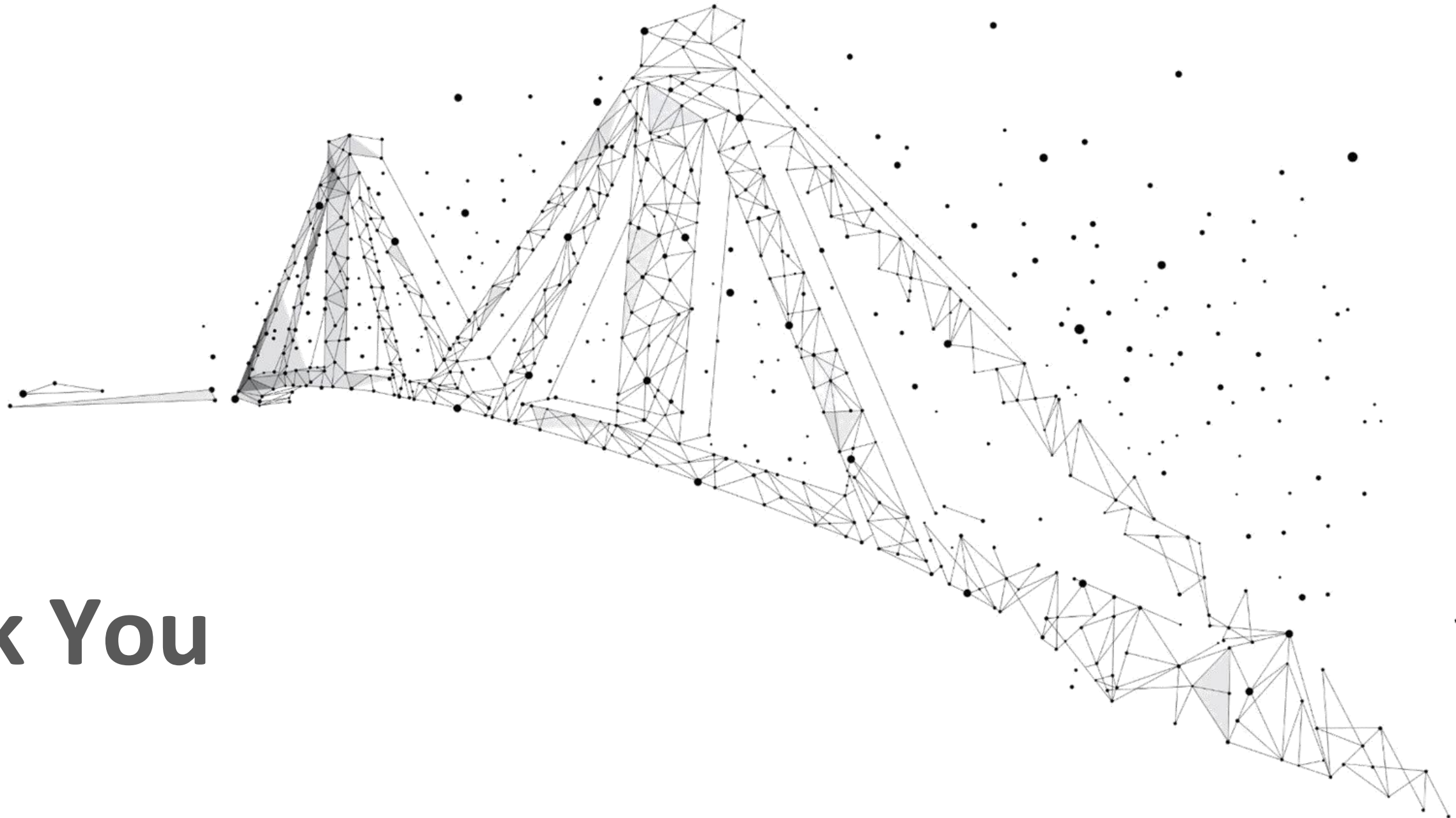- Better Scalability
- Reduces amount of data traversing the network

### ResNet50 Performance

**11%**

*8 Nodes, 16 GPUs, InfiniBand*

Legend: NCCL2 without SHARP | SHARP

### ResNet50 Performance

**16%**

*8 Nodes, 22 GPUs, InfiniBand*

Legend: NCCL2 without SHARP | SHARP

## Scalable Performance for Distributed AI

# Thank You