



# Intelligent Data Center Architecture to Enable Next Generation HPC/AI Platforms

Gilad Shainer, HPC-AI Advisory Council

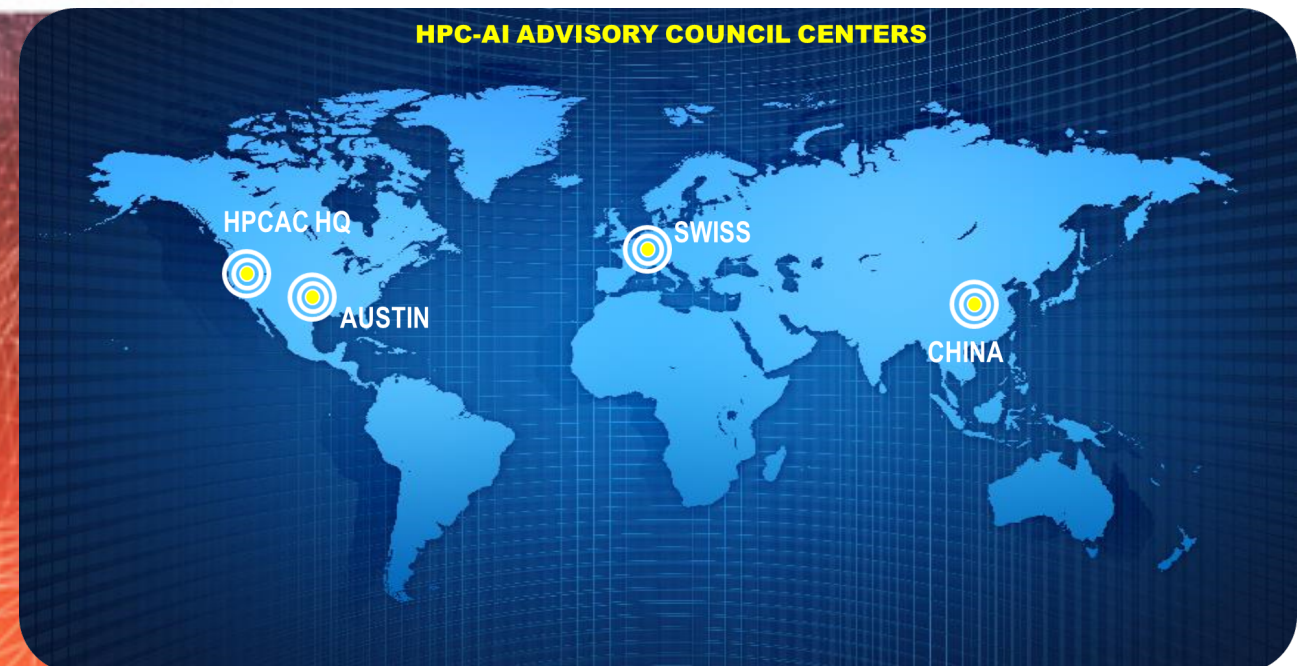
# The HPC-AI Advisory Council

- **World-wide HPC-AI non-profit organization**
- **More than 400 member companies / universities / research centers**
- **Bridges the gap between HPC-AI usage and its potential**
- **Provides best practices and a support/development center**
- **Explores future technologies and future developments**
- **Leading edge solutions and technology demonstrations**

## HPC Advisory Council Objectives

- HPC Technology
- Network of Expertise
- HPC Outreach
- High-Performance Center
- Education
- Best Practices

## HPC-AI ADVISORY COUNCIL CENTERS



# HPC-AI Advisory Council Members

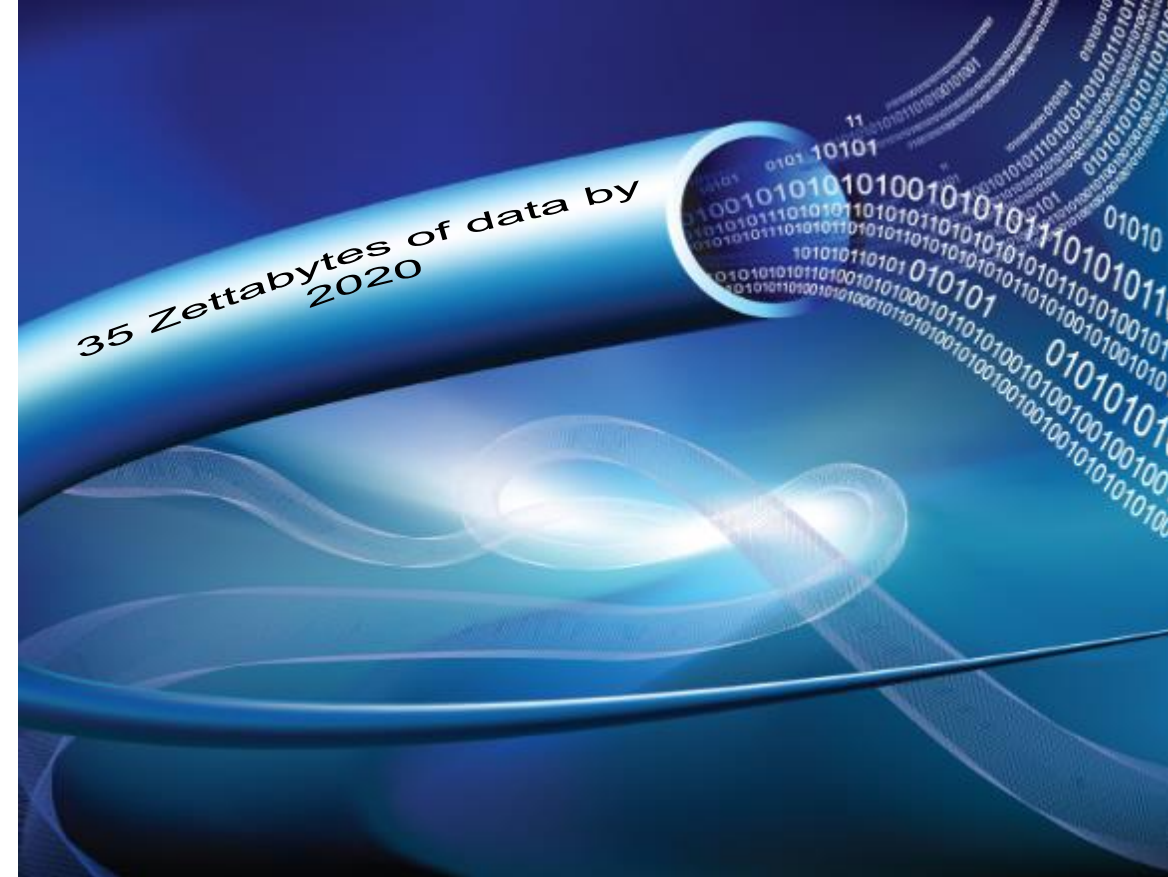


- **Applications Best Practices**
  - Hundreds of cases published
- **Cluster Center and Advanced Technology center**
- **2019 Conferences**
  - USA (Stanford University) – February
  - Switzerland (Swiss Supercomputing Center) – April
  - Australia (Pawsey Supercomputing Center) – August
  - Spain (Barcelona Supercomputing Center) – Sep
  - UK (University of Leicester, DiRAC) – Sep
  - China (HPC China)
- **2019 Competitions**
  - APAC HPC-AI Competition – March
  - China - 7th Annual RDMA Competition – May
  - ISC Germany - 8th Annual Student Cluster Competition – June
- **For more information**
  - [www.hpcadvisorycouncil.com](http://www.hpcadvisorycouncil.com)
  - [info@hpcadvisorycouncil.com](mailto:info@hpcadvisorycouncil.com)





20<sup>th</sup> Century



21<sup>st</sup> Century

# World of Data – World of Opportunities

>70% of data  
is generated  
by consumers

80% of data  
stored

3%  
structured

0.5% being  
analyzed

<0.5%  
generates  
information



# The Power of Data – for Everyone

**NOKIA**  
Connecting People

**NAVTEQ**

- ~5M traffic sensors for \$8.1B

 **waze**

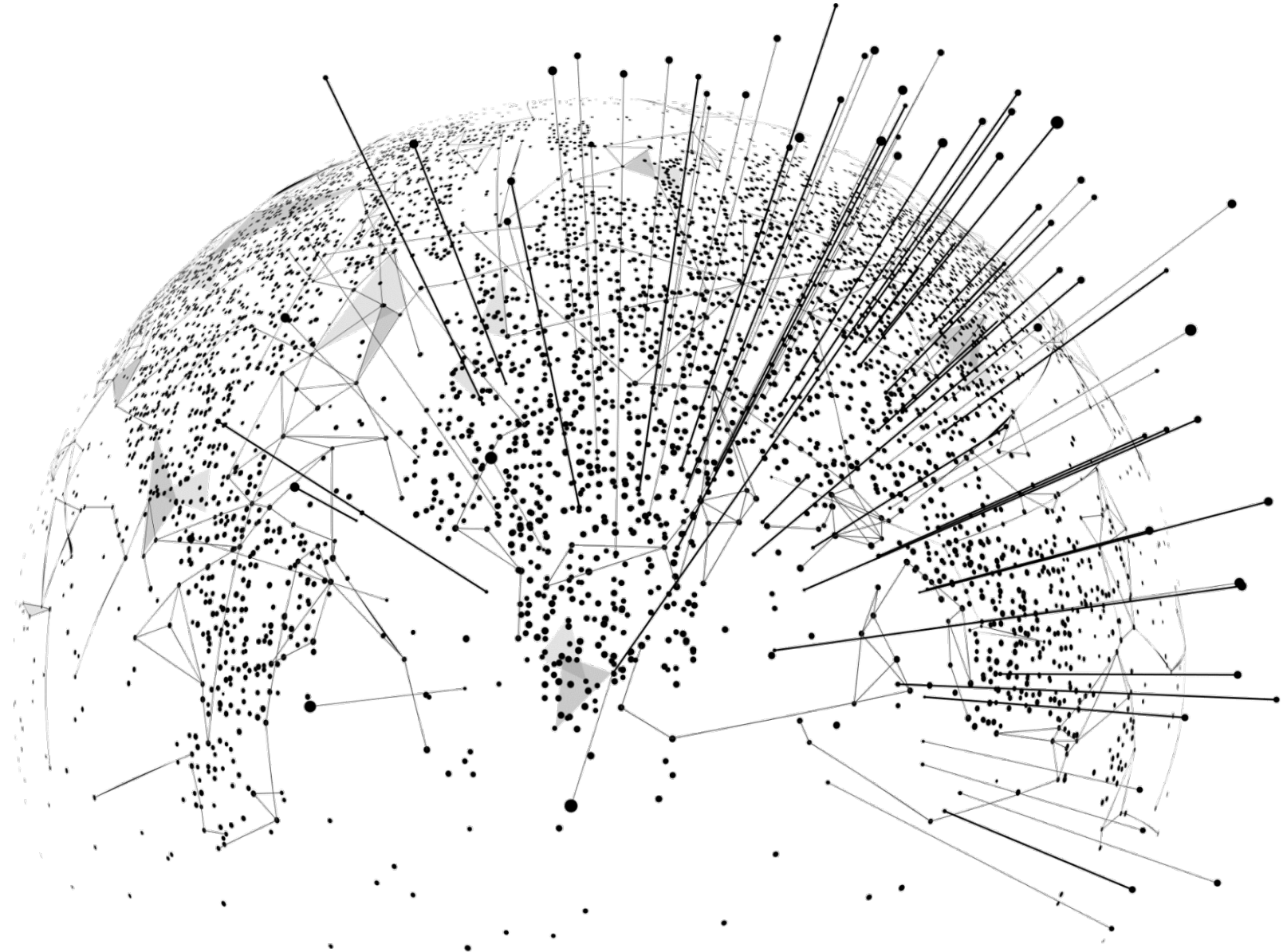


- 50M traffic sensors for free



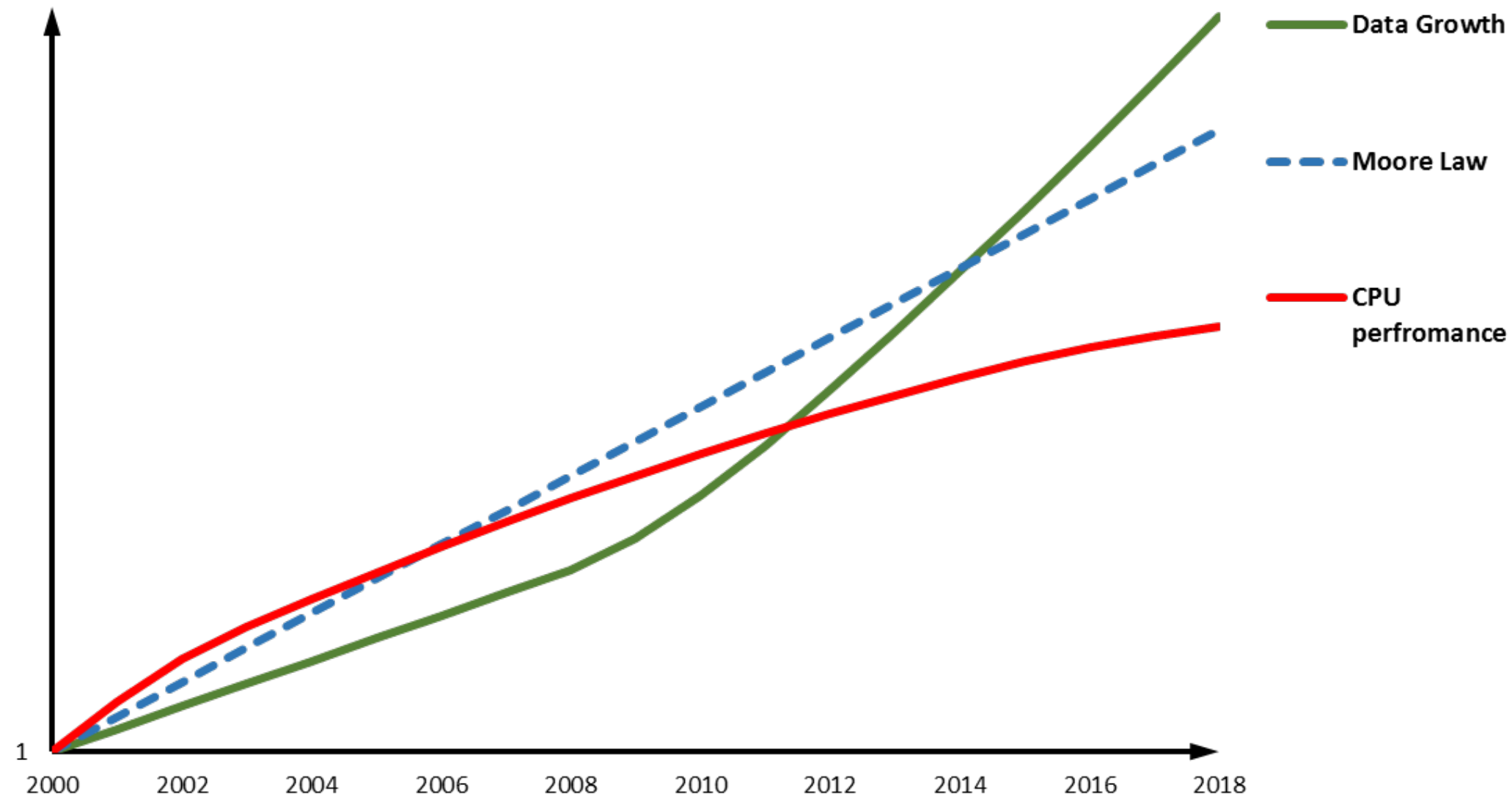
# Data Processing Plant – Convert Data to Information

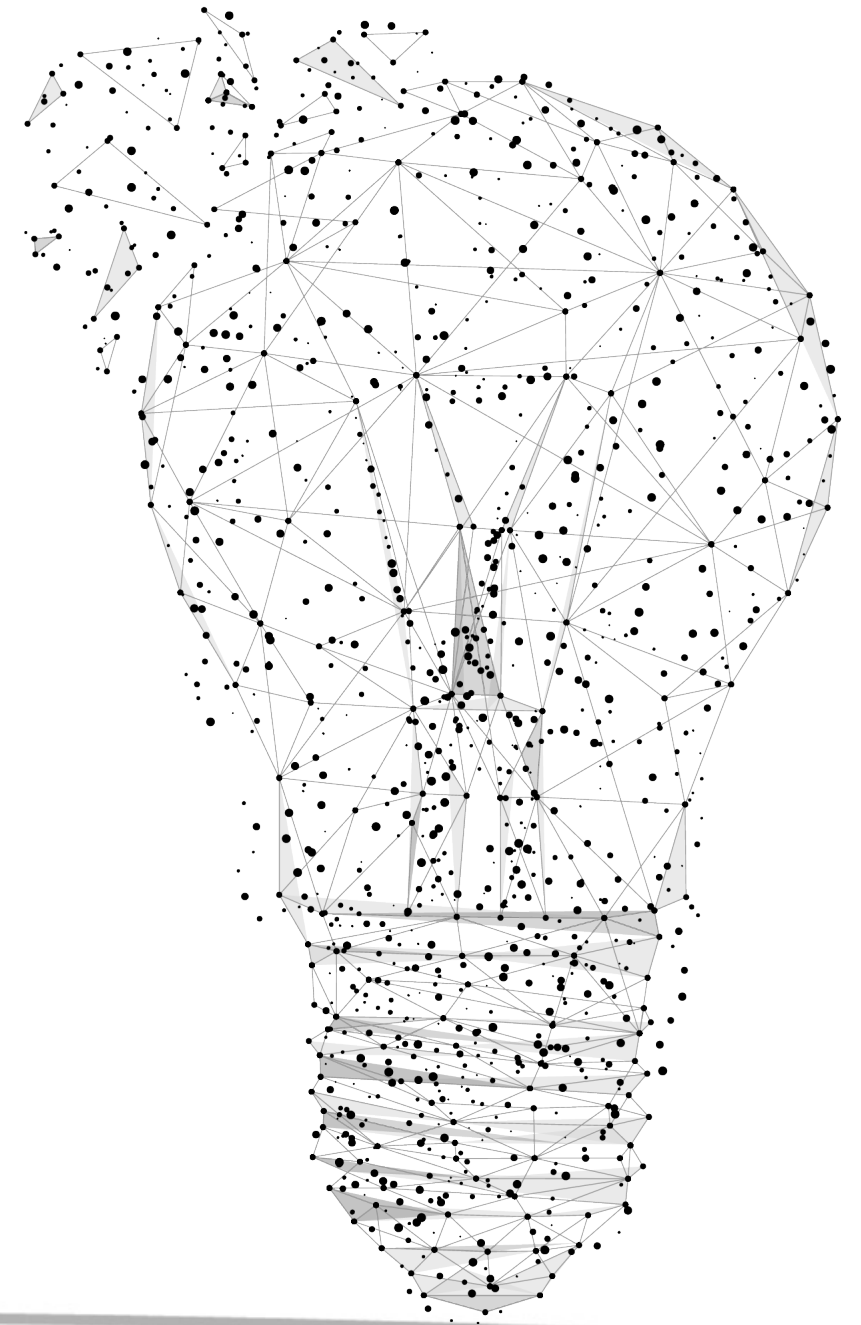




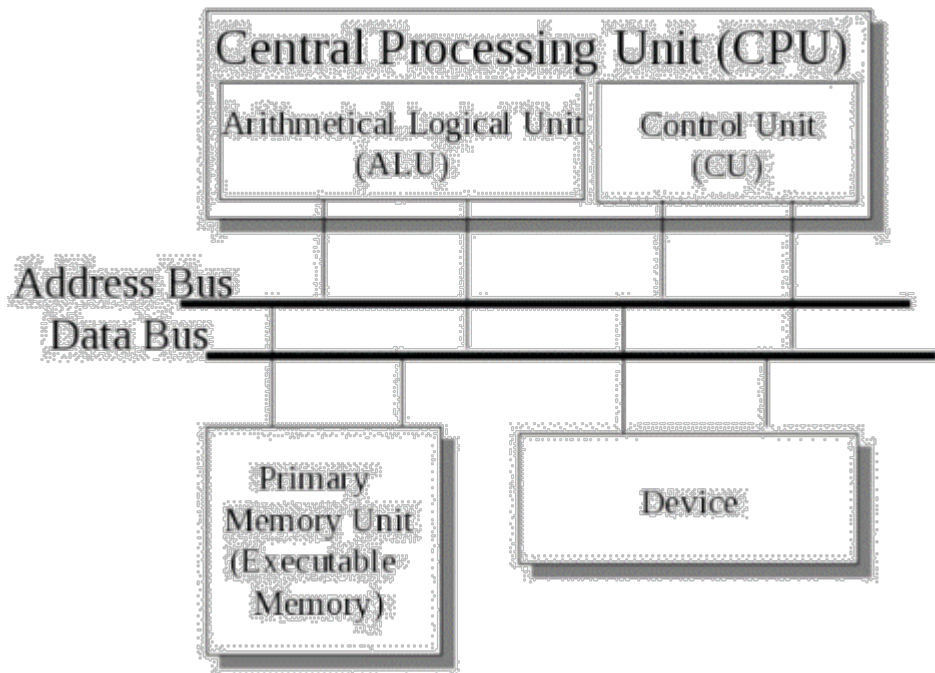
# Computing Challenge



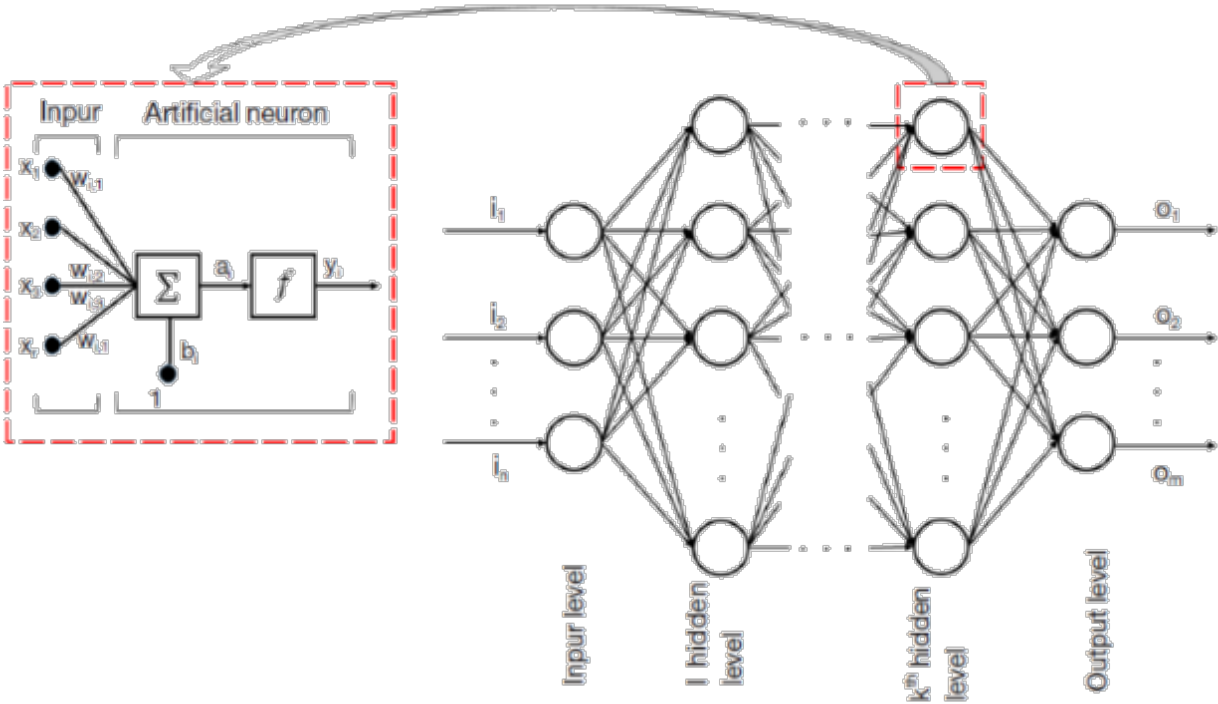
“The Electric Light  
Did Not Come From  
Continuous Improvement  
Of Candles”



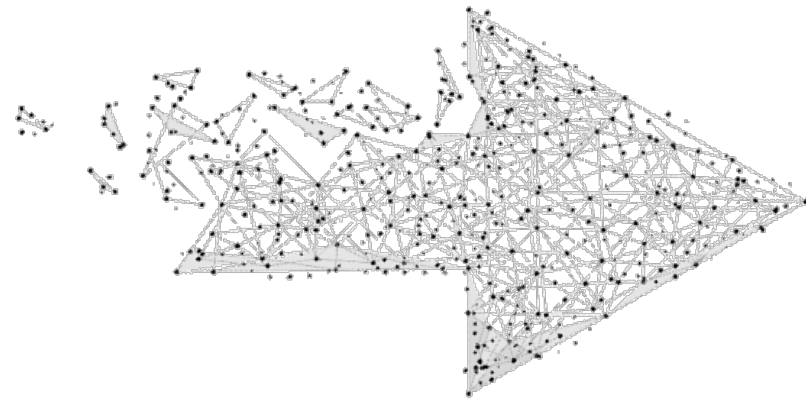
## Compute-Centric



## Data-Centric

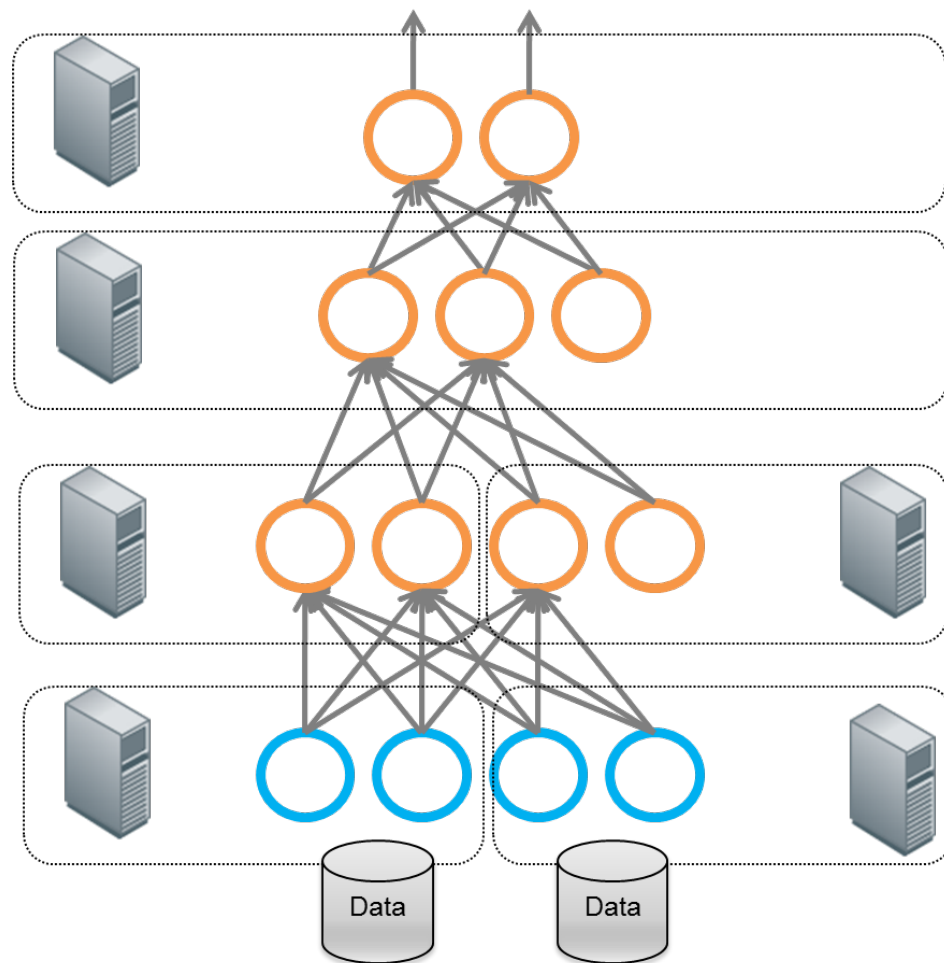


## Von Neumann Machine

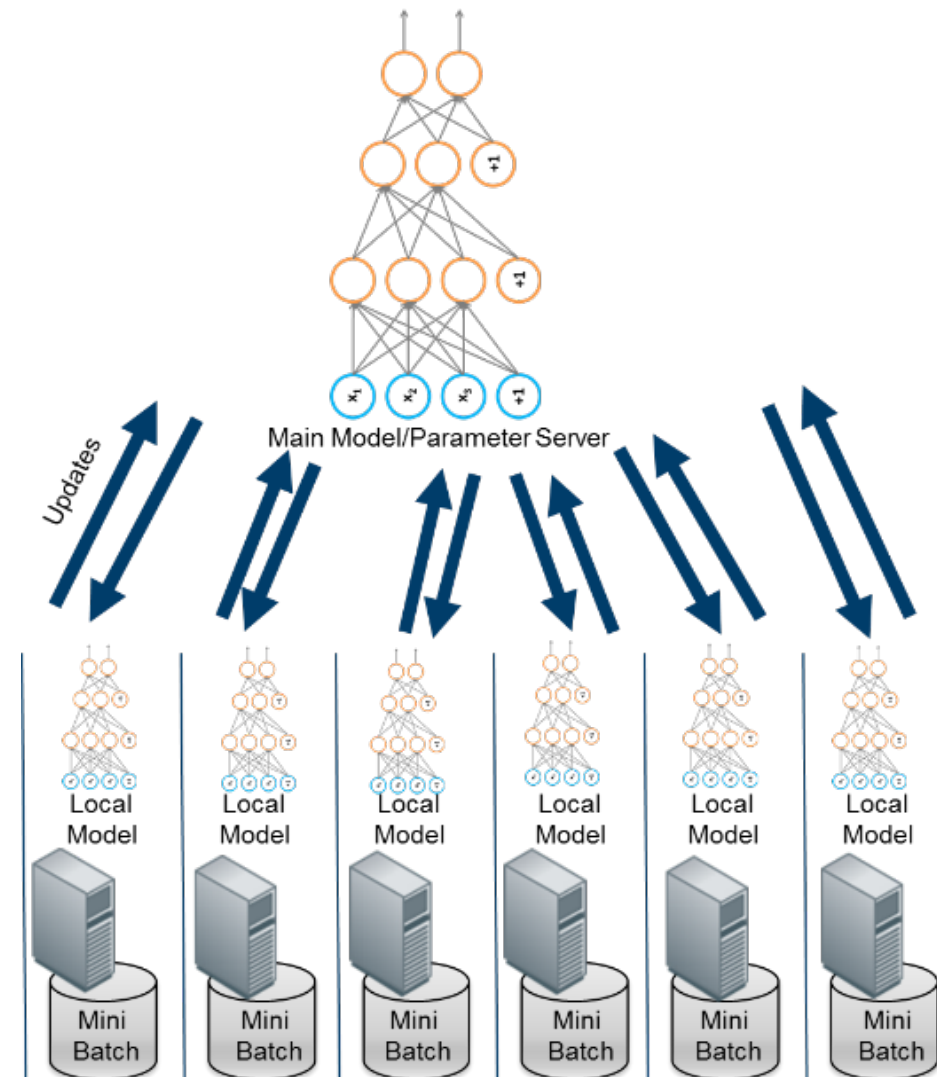


## DataFlow Machine

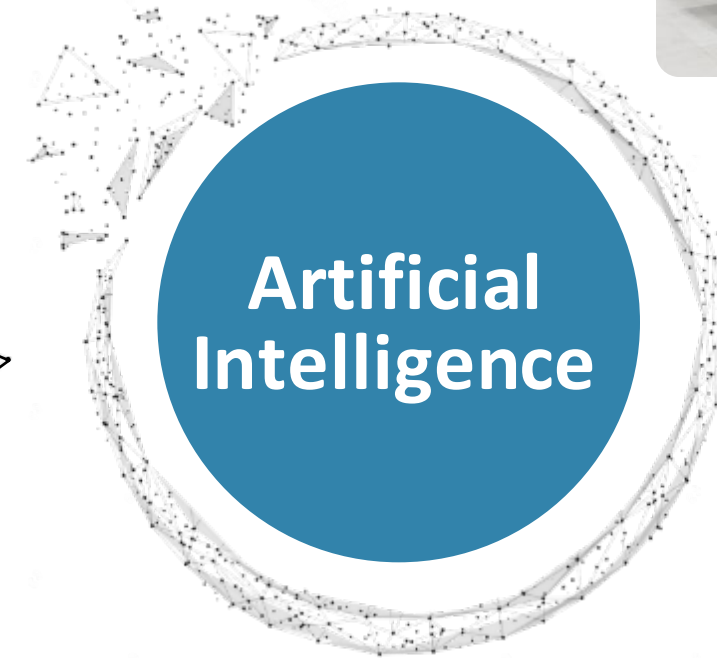
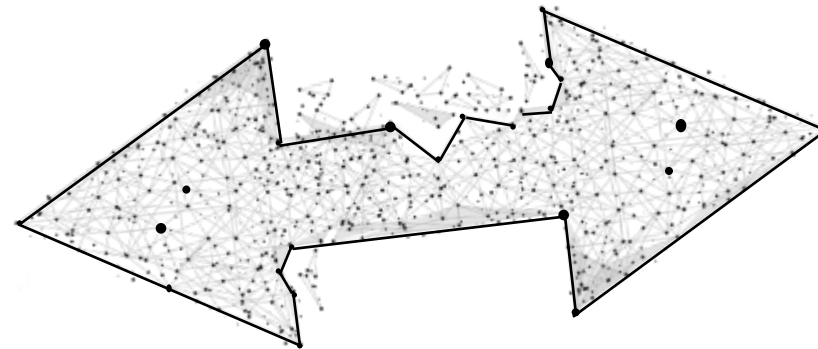
## Model Parallelism



## Distributed Training



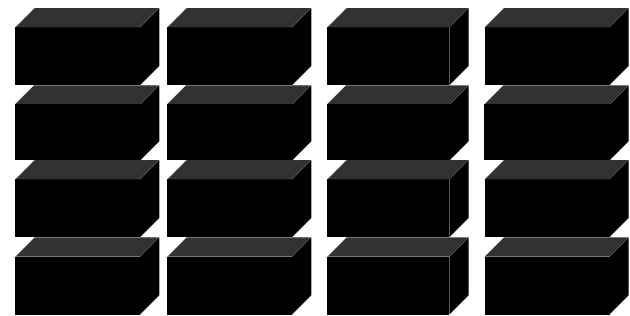
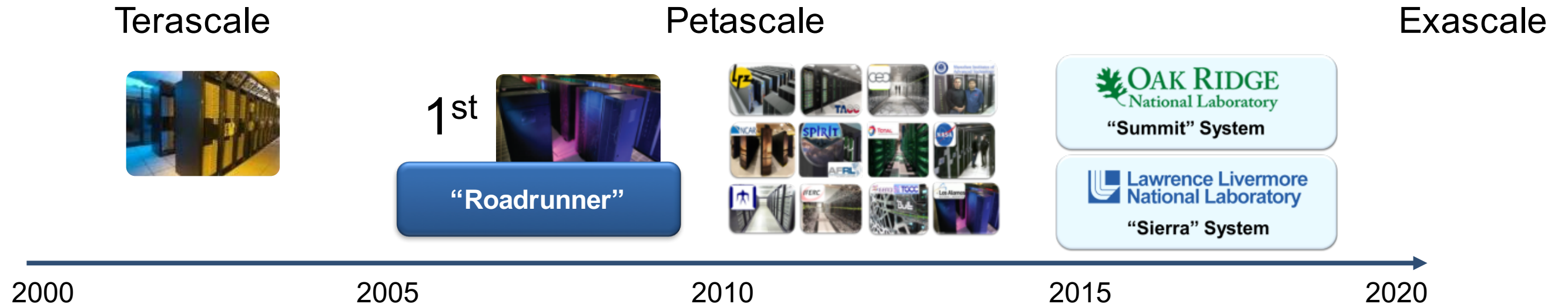
# High Performance Computing and Artificial Intelligence



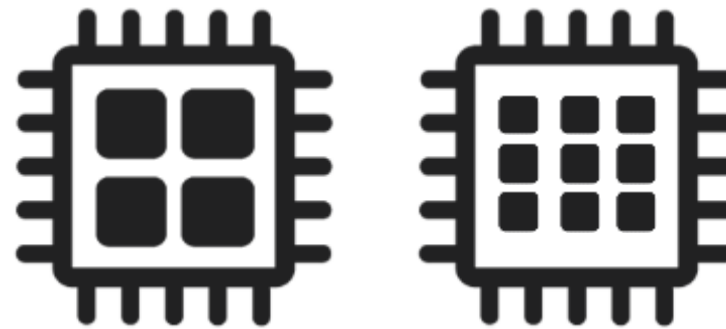
Same Technology



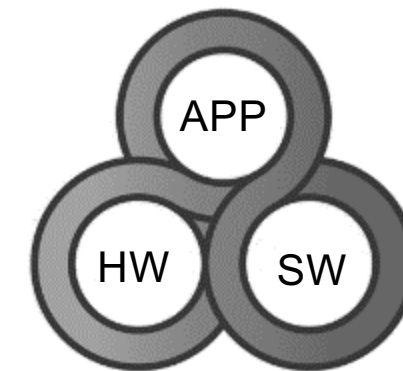
# The Ever Growing Demand for Higher Performance



SMP to Clusters



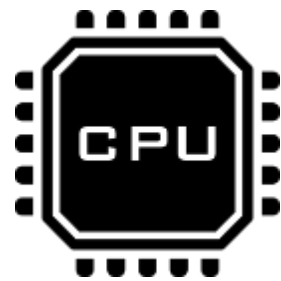
Single-Core to Many-Core



Co-Design

Application  
Software  
Hardware

# From CPU-Centric to Data-Centric Data Centers



**CPU**



**Network**



**Storage**

# From CPU-Centric to Data-Centric Data Centers

Workload

Workload

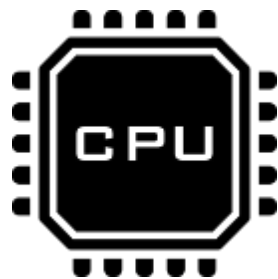
Workload

Communication  
Framework (MPI)

CPU Functions

Network Functions

Storage Functions



**In-CPU Computing**



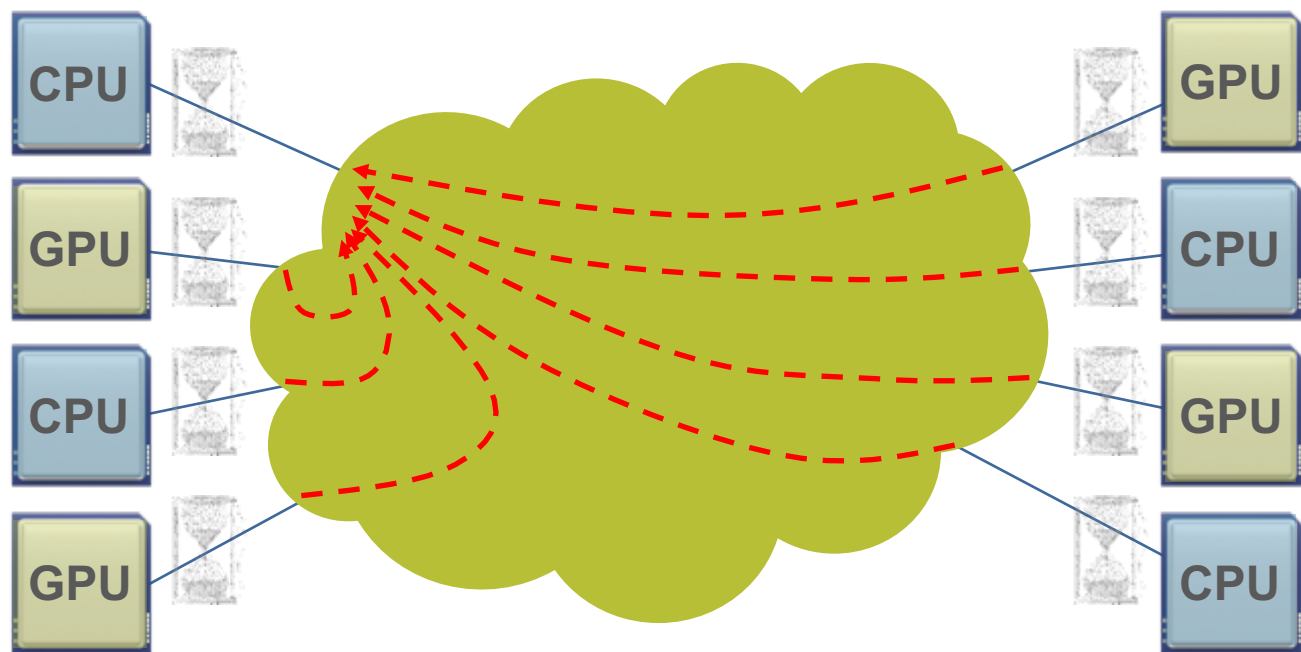
**In-Network Computing**



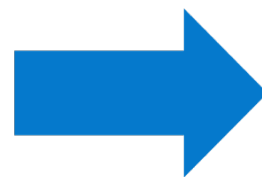
**In-Storage Computing**



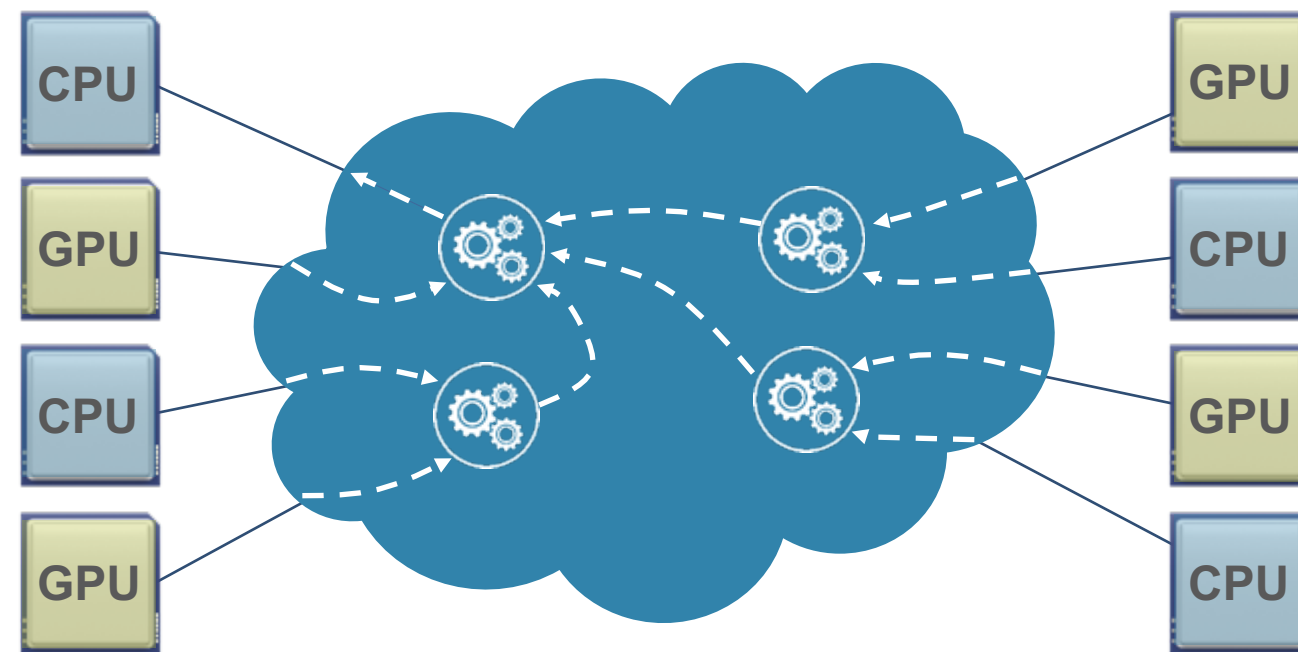
### CPU-Centric (Onload)



Communications Latencies  
of 30-40us

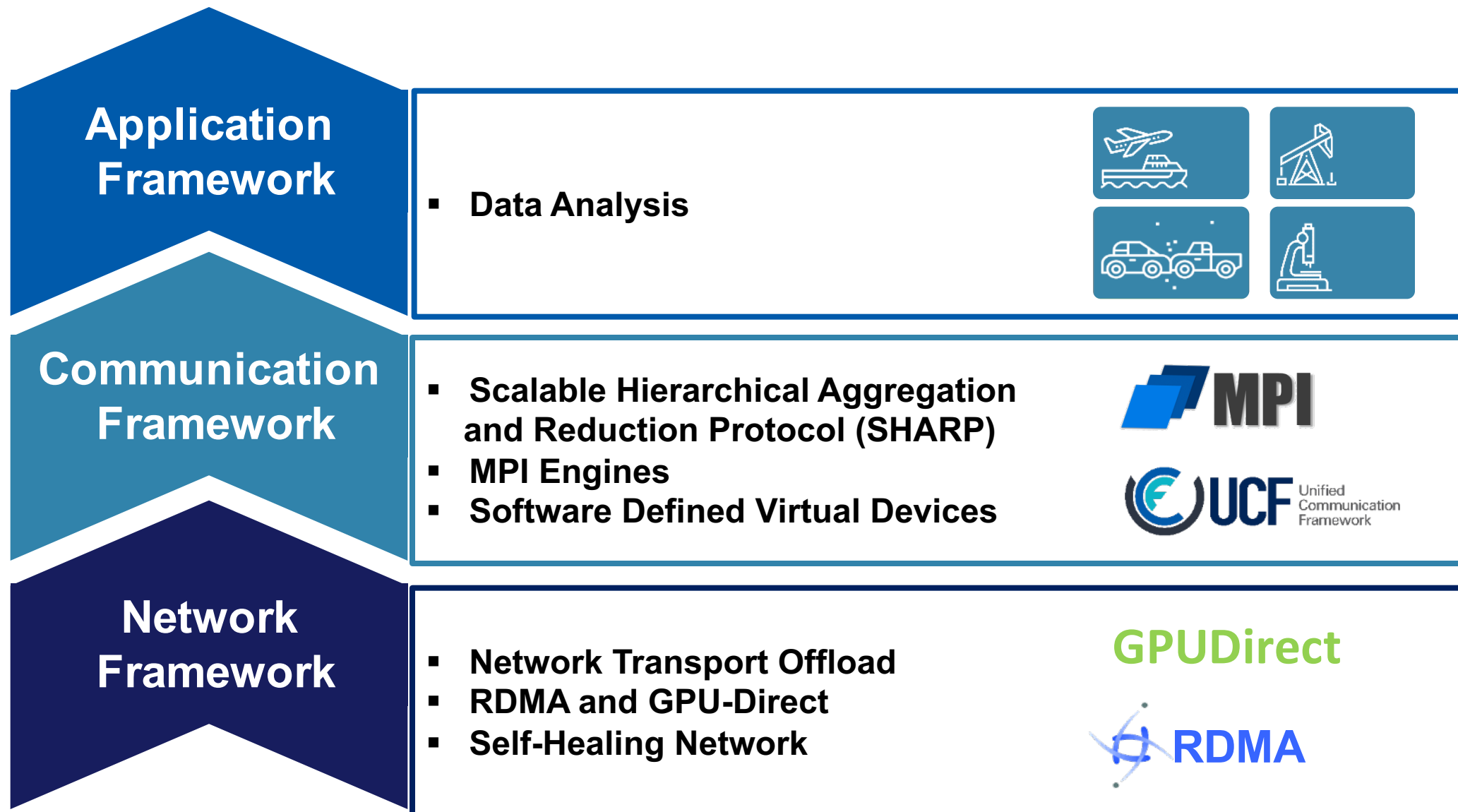


### Data-Centric (Offload)

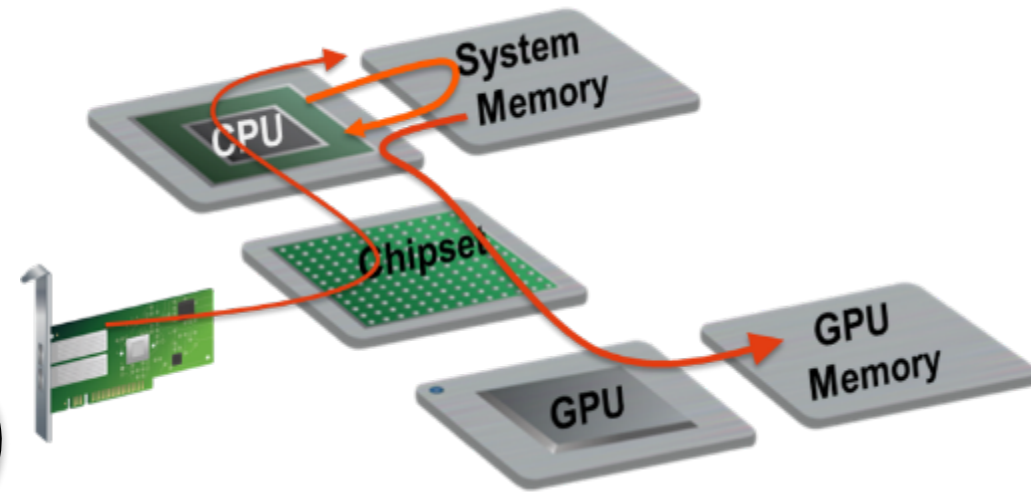
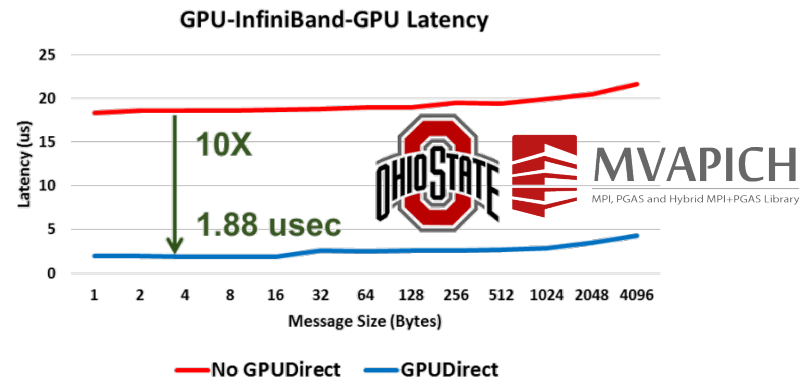


Communications Latencies  
of 3-4us

# The Roadmap of In-Network Computing

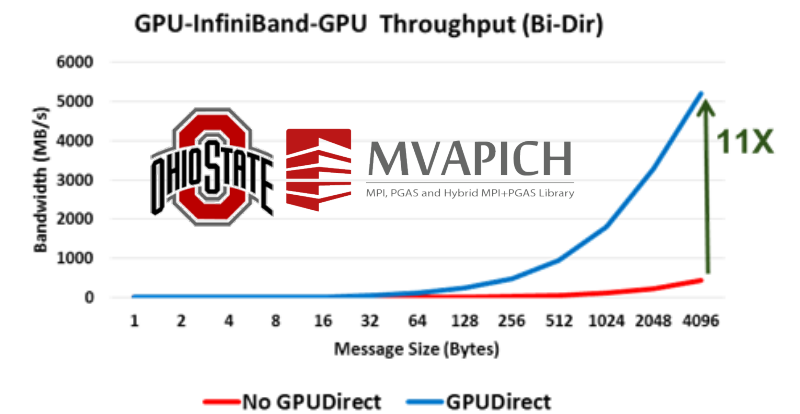
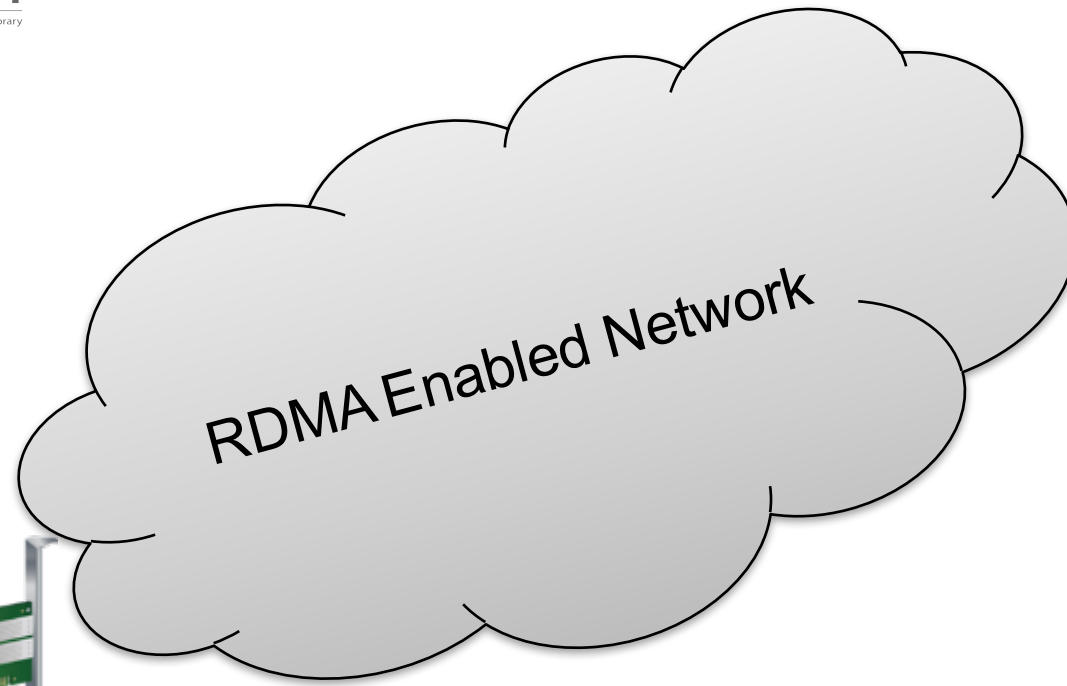
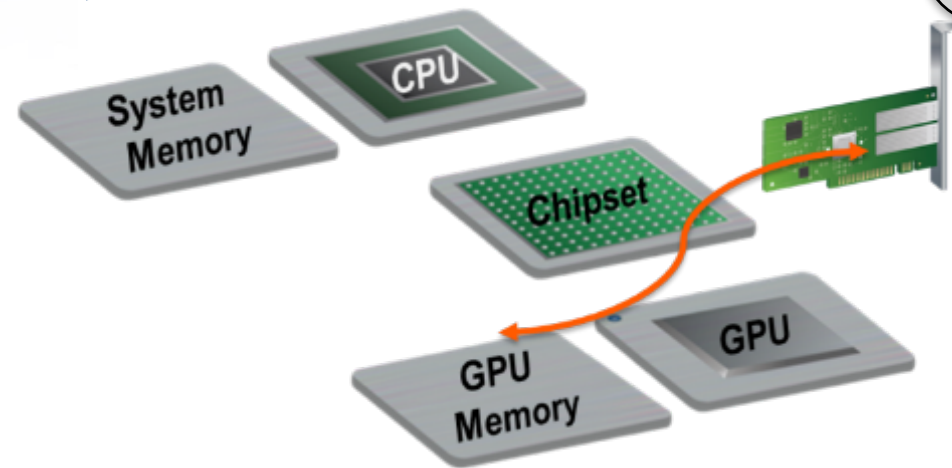


# GPUDirect RDMA Technology and Advantages



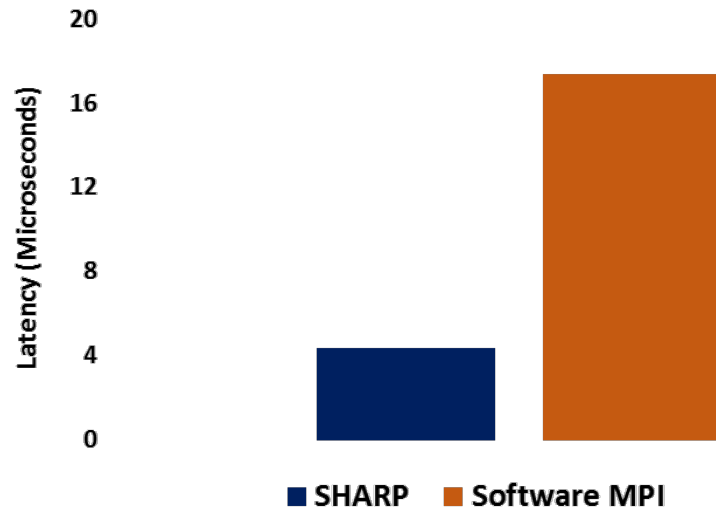
**Without GPU Direct  
- Same Data Copied 3 Times**

**With GPUDirect**

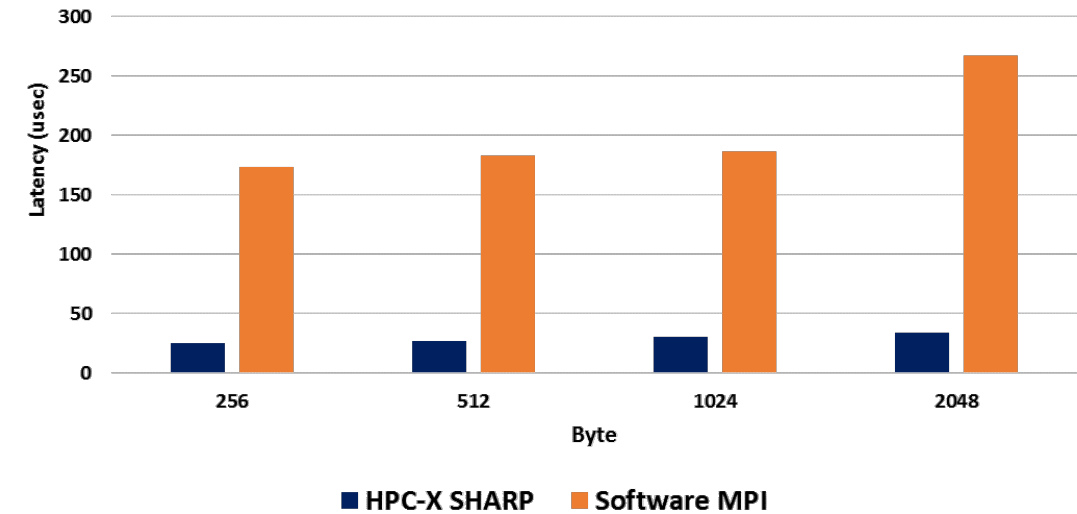


# Scalable Hierarchical Aggregation and Reduction Protocol (SHARP™) Performance Advantages

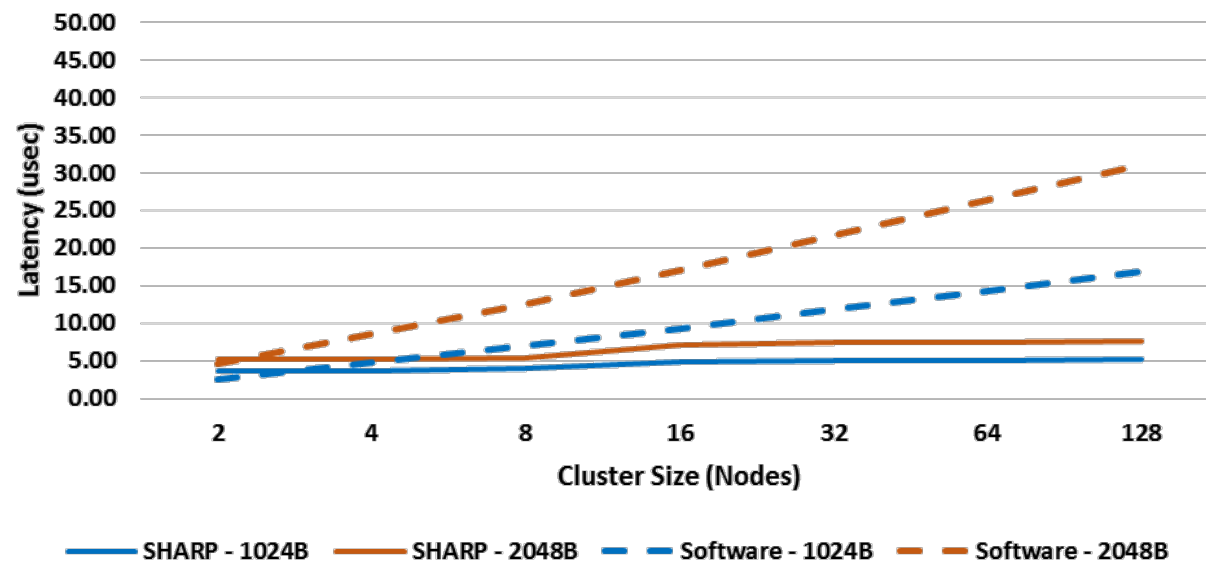
**MPI Barrier Latency**  
Summit Supercomputer, 512 Nodes



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



**Allreduce Latency**

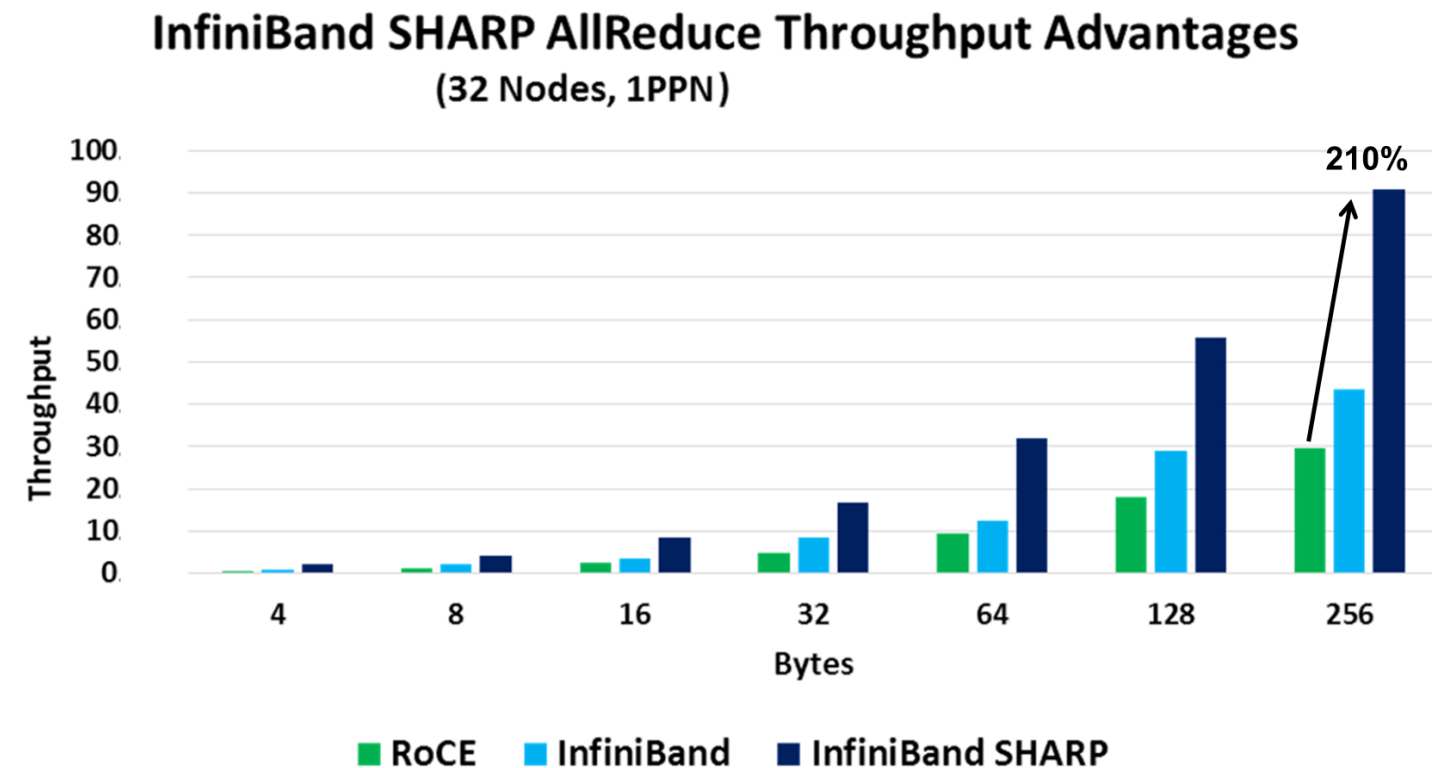
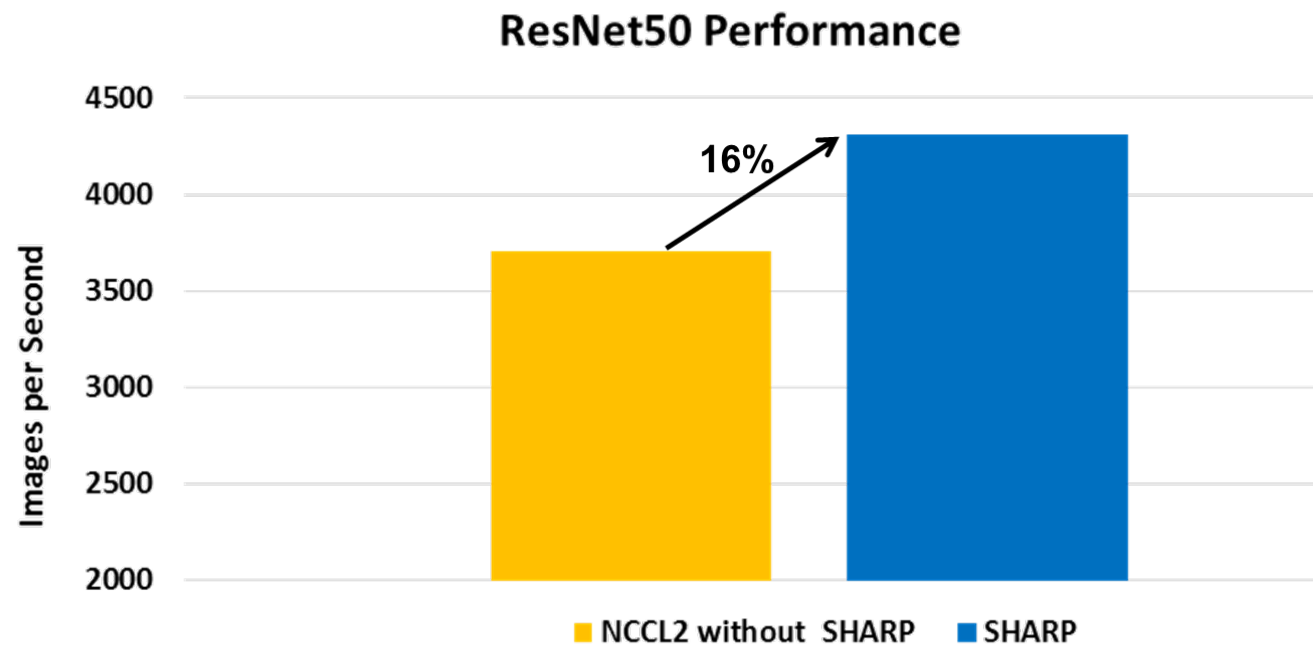


UNIVERSITY OF  
**TORONTO**



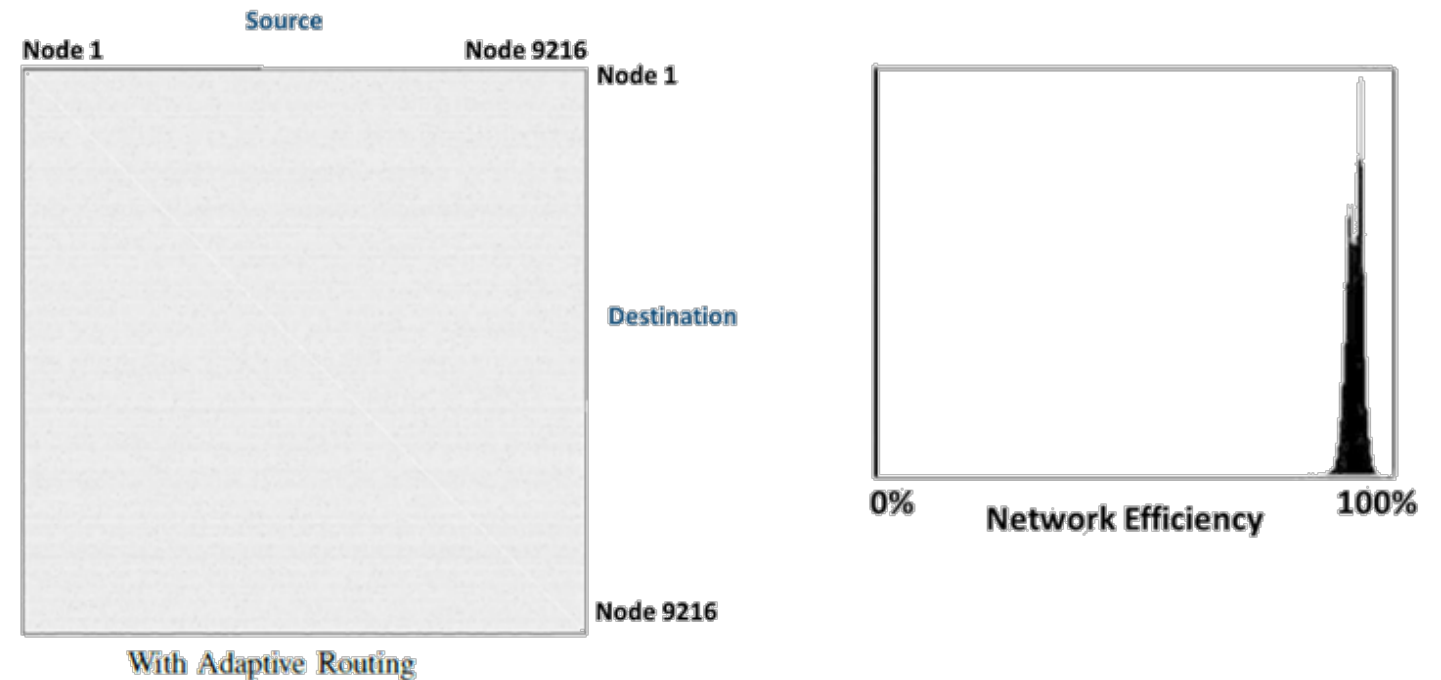
# SHARP Performance Advantage for Deep Learning

- SHARP provides 16% Performance Increase for deep learning, initial results
- TensorFlow with Horovod running ResNet50 benchmark, HDR InfiniBand



- Oak Ridge National Laboratory – Coral Summit supercomputer
- Bisection bandwidth benchmark, based on mpiGraph
- InfiniBand demonstrates an average performance of 96% network utilization

mpiGraph explores the bandwidth between possible MPI process pairs. In the histograms, the single cluster with adaptive routing indicates that all pairs achieve nearly maximum bandwidth while single-path static routing has nine clusters as congestion limits bandwidth, negatively impacting overall application performance.

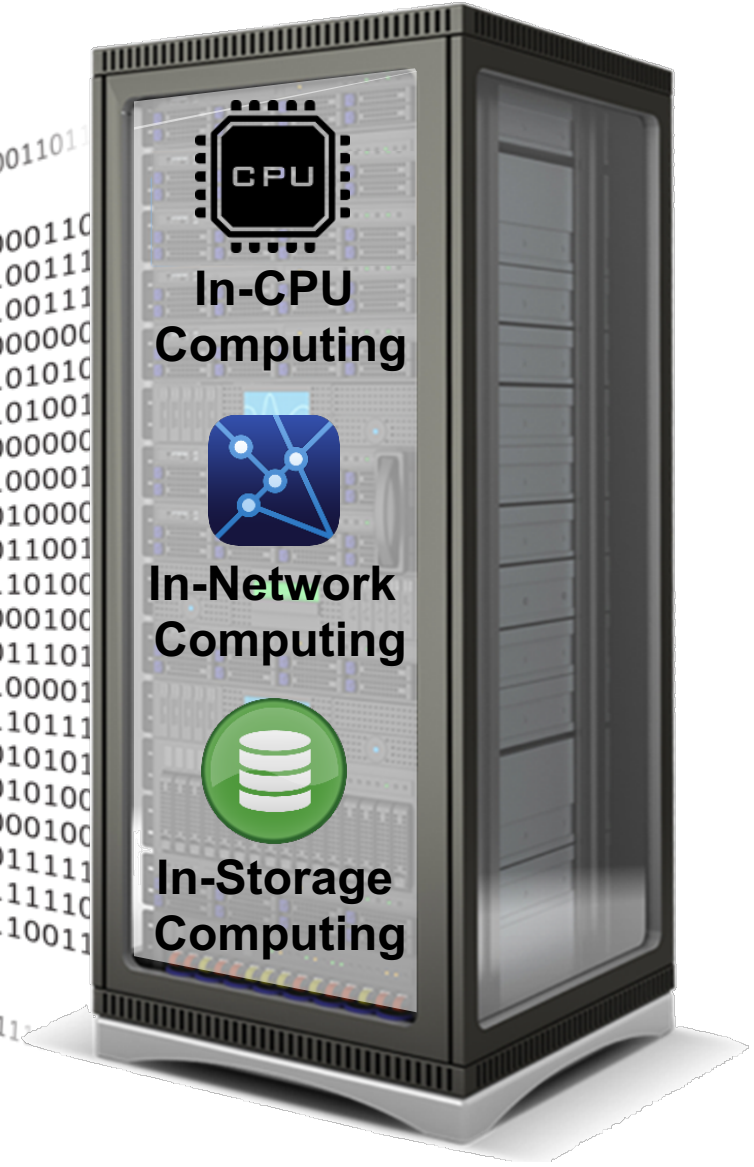


Adaptive Routing

Summit's MpiGraph Output

Source: "The Design, Deployment, and Evaluation of the CORAL Pre-Exascale Systems" paper

# Unleashing the Power of Data



# Thank You

