



THE OHIO STATE UNIVERSITY  
COLLEGE OF ENGINEERING

# Faster Deep Learning Training with MVAPICH2-GDR on NVLink-enabled GPU Clusters

**Ching-Hsiang Chu**

Ph.D. Candidate

Network Based Computing Lab  
Department of Computer Science and Engineering  
The Ohio State University, Columbus, OH

# Outline

- Introduction
- GPU-enabled Allreduce Designs in MVAPICH2-GDR
- Concluding Remarks

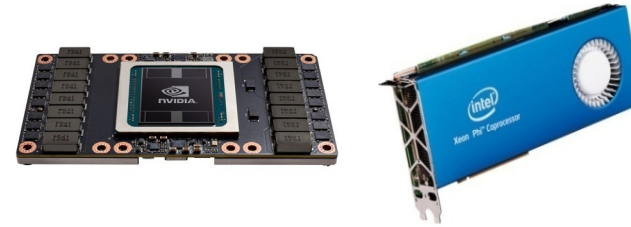
# Trends in Modern HPC Architecture: Heterogeneous



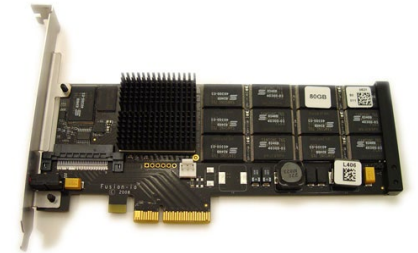
Multi/ Many-core Processors



High Performance Interconnects  
InfiniBand, Omni-Path, EFA  
<1usec latency, 200Gbps Bandwidth



Accelerators / Coprocessors  
high compute density,  
high performance/watt



SSD, NVMe-SSD,  
NVRAM  
Node local storage

- Multi-core/many-core technologies
- High Performance Interconnects

- High Performance Storage and Compute devices
- Variety of programming models (MPI, PGAS, MPI+X)



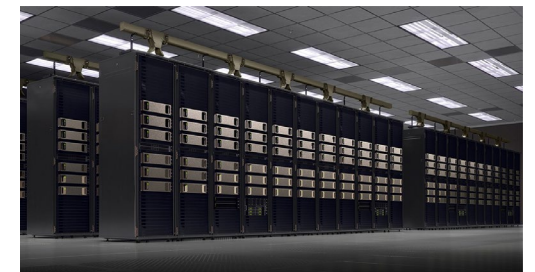
#1 Summit  
(27,648 GPUs)



#2 Sierra (17,280 GPUs)  
#10 Lassen (2,664 GPUs)



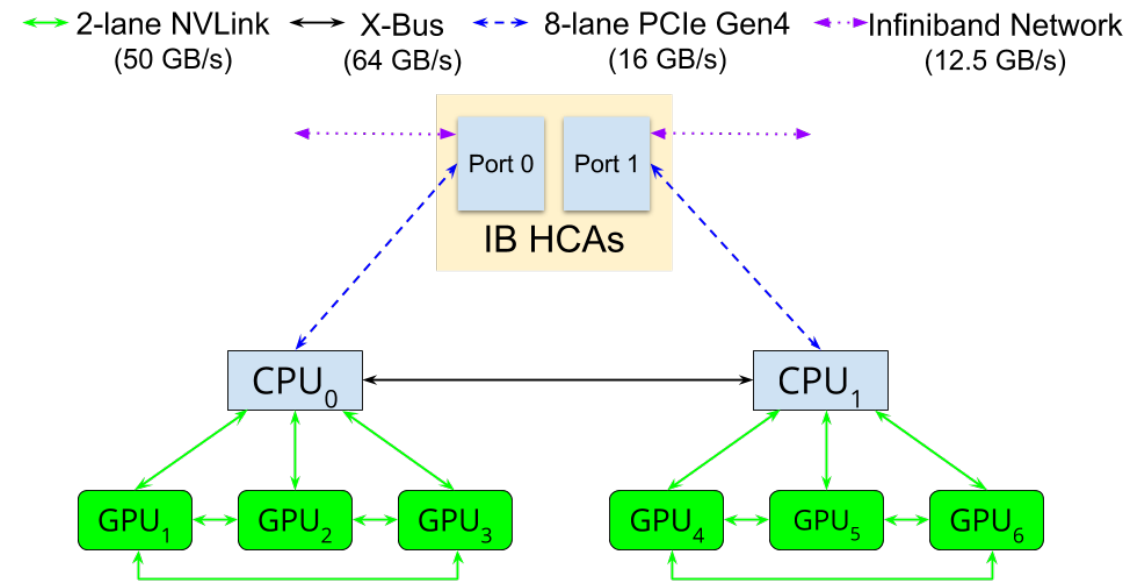
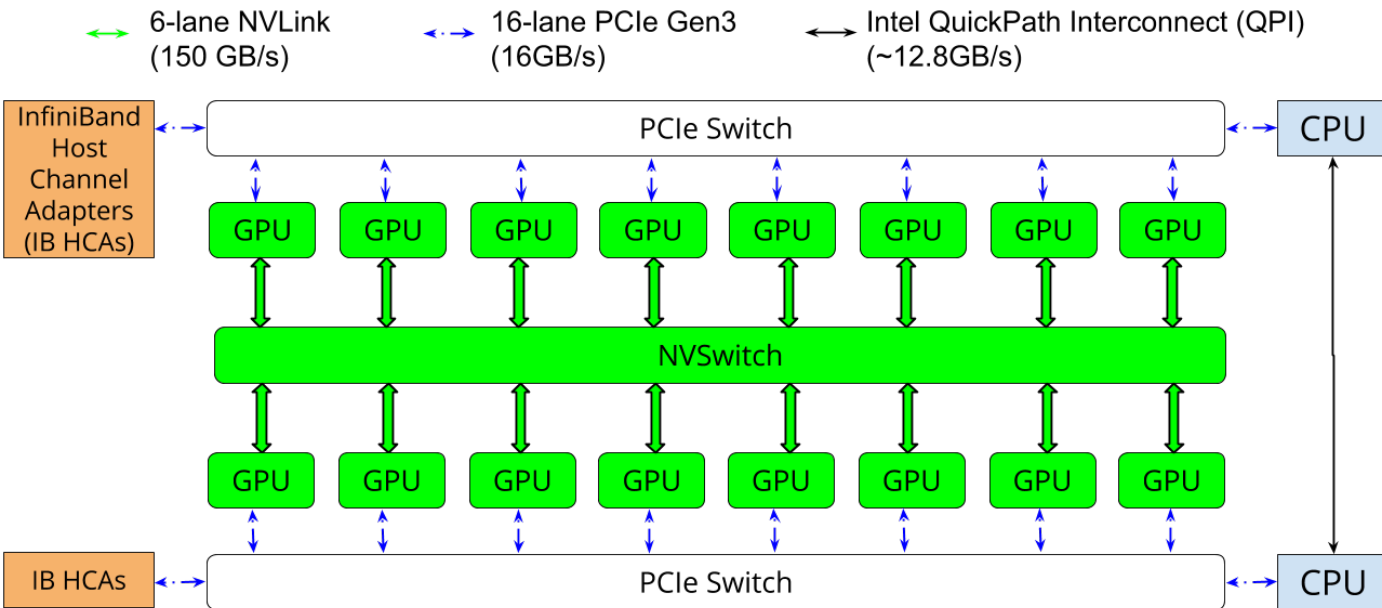
#8 ABCI  
(4,352 GPUs)



#22 DGX SuperPOD  
(1,536 GPUs)

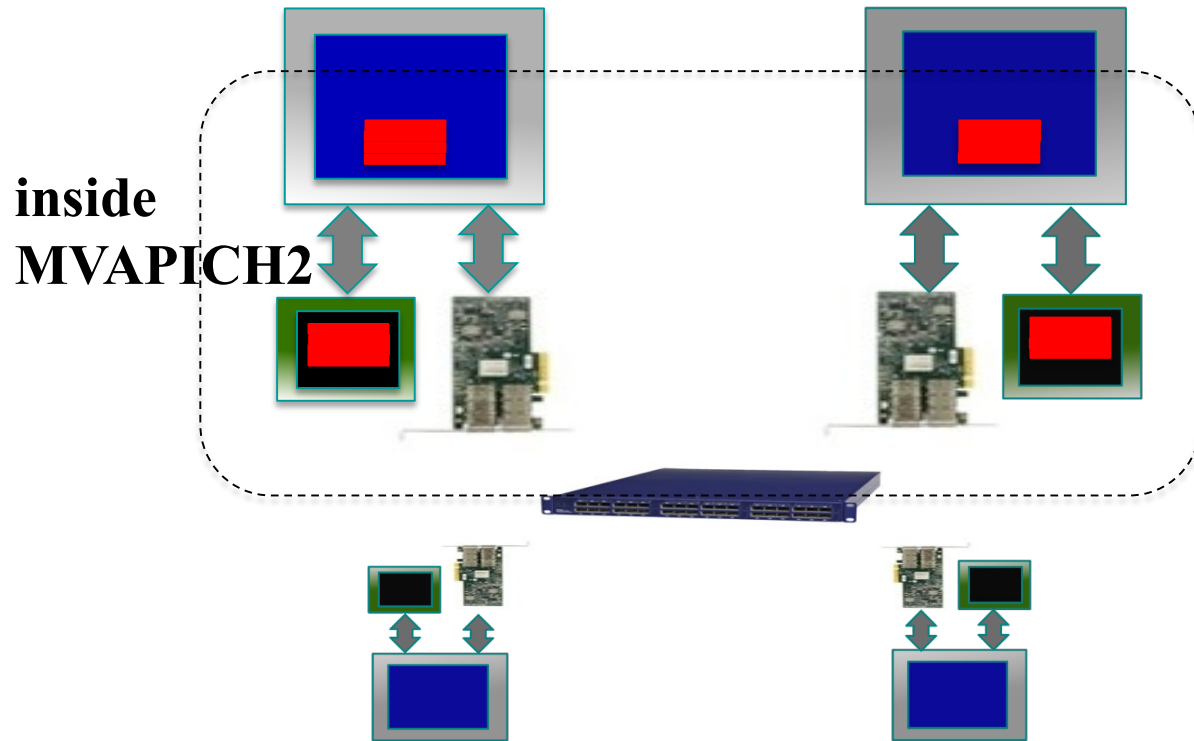
# Trends in Modern Large-scale Dense-GPU Systems

- **Scale-up** (up to 150 GB/s)
  - PCIe, NVLink/NVSwitch
  - Infinity Fabric, Gen-Z, CXL
- **Scale-out** (up to 25 GB/s)
  - InfiniBand, Omni-path, Ethernet
  - Cray Slingshot



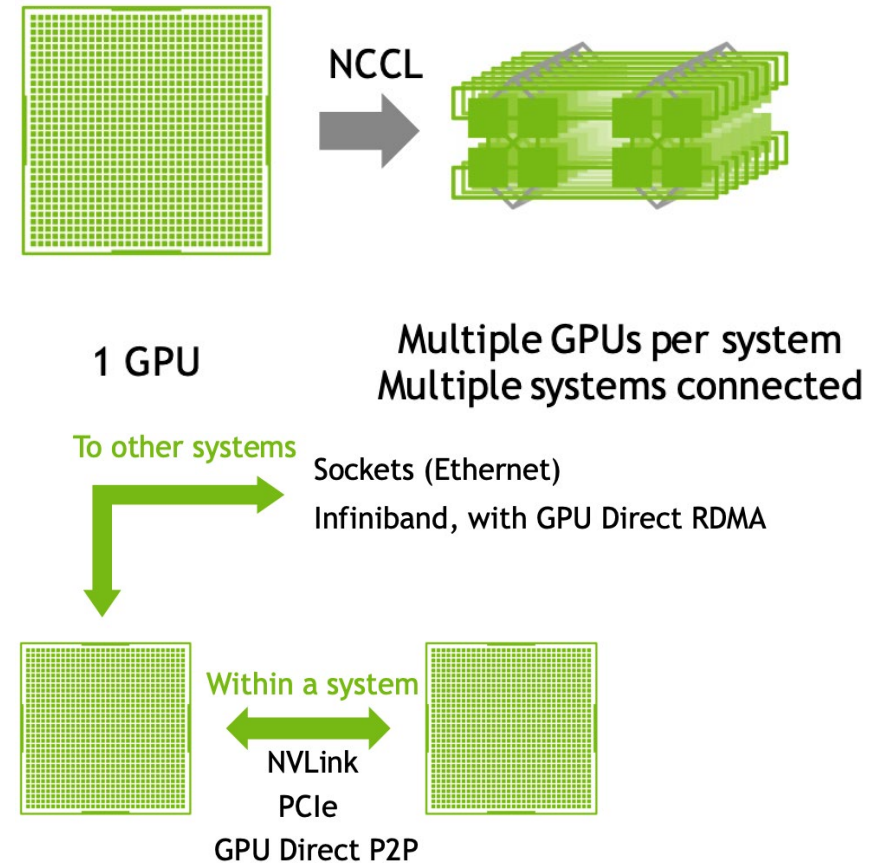
# GPU-Aware (CUDA-Aware) Communication Middleware

## MPI-based Generic Communication Middleware



- Supports and optimizes various communication patterns
- Overlaps data movement from GPU with RDMA transfers

## DL-Specific Communication Middleware



- Ring-based collective operations
- Optimized for DL workloads on GPU systems

# GPU-enabled Emerging Deep Learning Applications

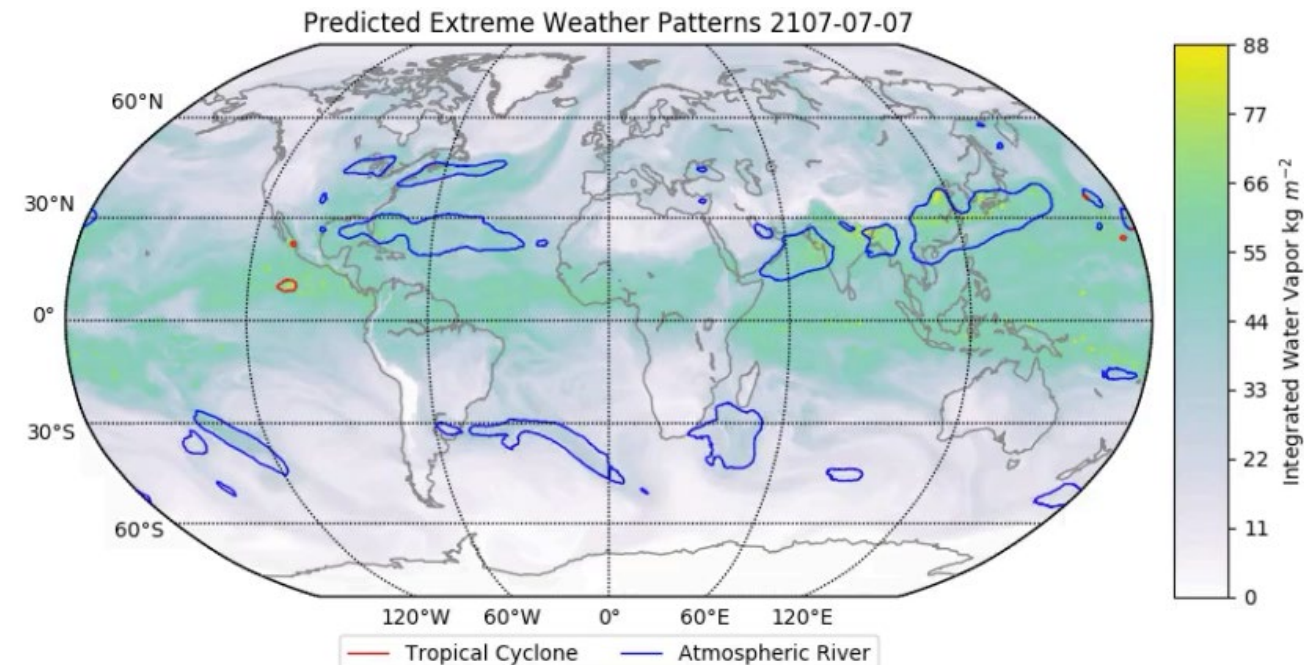
- Easy-to-use and high-performance frameworks



- Wide range of applications

- Image Classification
- Speech Recognition
- Self-driving car
- Healthcare
- Climate Analytic

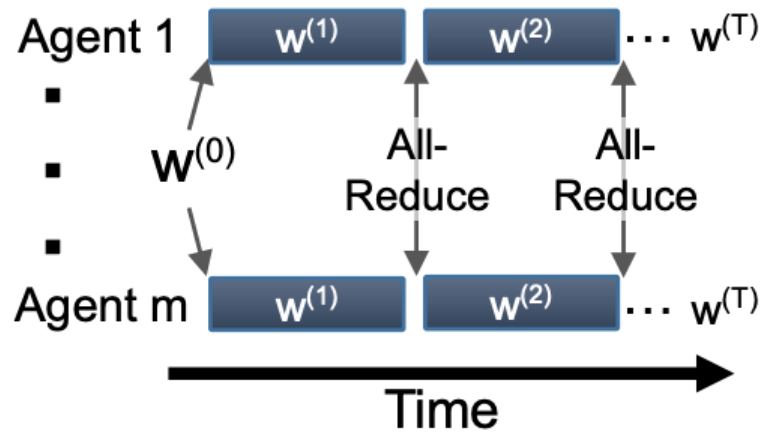
999 PetaFlop/s sustained, and 1.13 ExaFlop/s peak FP 16 performance over 4560 nodes (27,360 GPU)



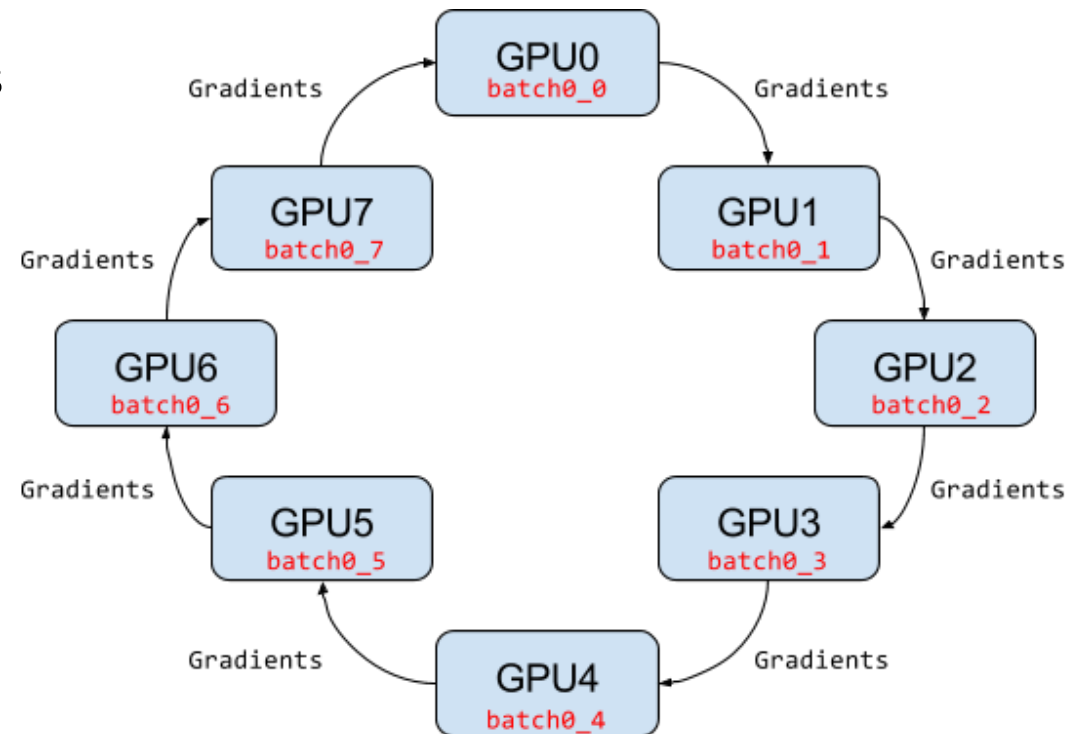
Kurth T, Treichler S, Romero J, Mudigonda M, Luehr N, Phillips E, Mahesh A, Matheson M, Deslippe J, Fatica M, Houston M. Exascale deep learning for climate analytics. SC 2018 Nov 11 (p. 51). (Golden Bell Prize)

# Motivated Example – Reduction Op. for DL Training

- Can GPU resources help improving compute-intensive communications?
  - E.g., MPI\_Reduce, MPI\_Allreduce, MPI\_Scan
  - **Emerging distributed deep learning training**
    - Exchange and update weights
  - Requires **fast and high-bandwidth** solutions



Ben-Nun T, Hoefler T. Demystifying parallel and distributed deep learning: An in-depth concurrency analysis. arXiv preprint arXiv:1802.09941. 2018 Feb 26.



<https://www.oreilly.com/ideas/distributed-tensorflow>

# How to leverage GPUs for MPI Reduction Operations?

## Existing designs

1. Explicit copy the data from GPU to host memory
2. Host-to-Host communication to remote processes
3. Perform computation on CPU
4. Explicit copy the data from host to GPU memory

Expensive!

Fast

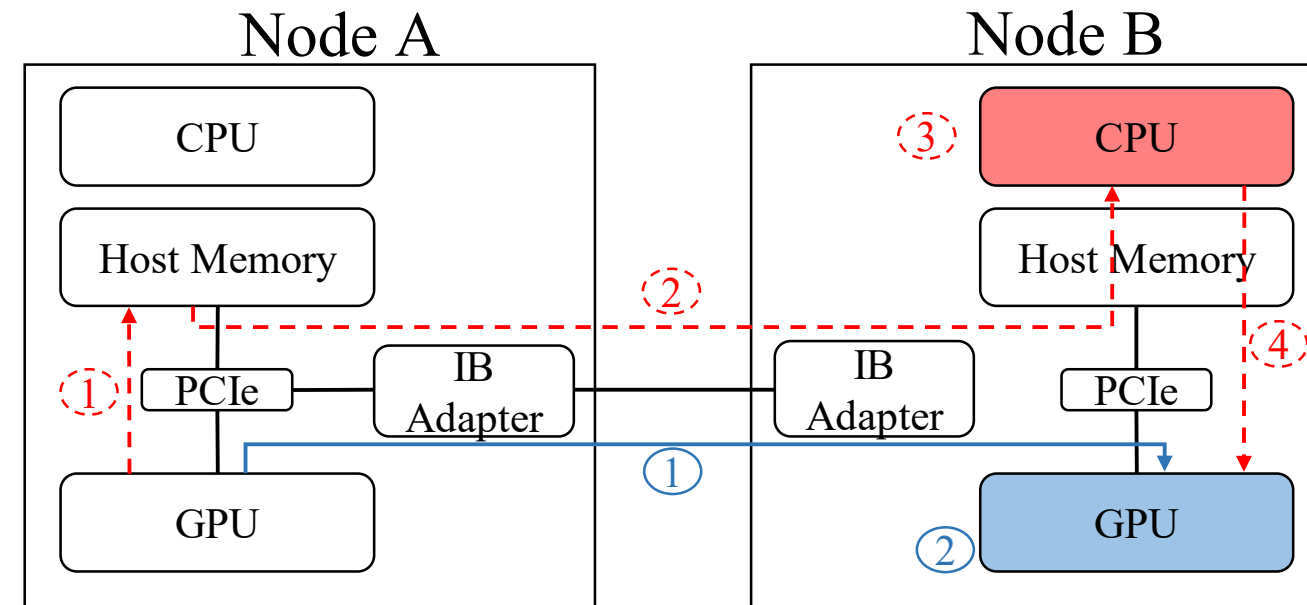
Relative slow for large data

Good for small data

Expensive!

## Proposed designs

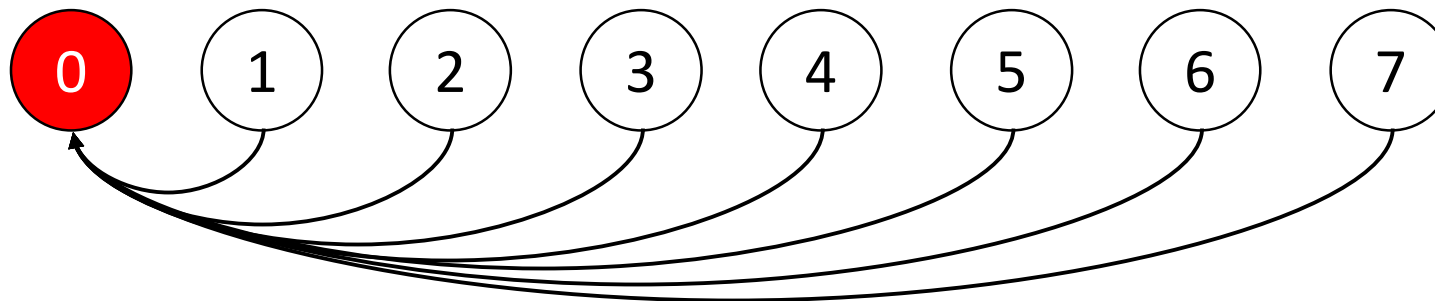
1. GPU-to-GPU communication
  - NVIDIA GPUDirect RDMA (GDR)
  - Pipeline through host for large msg
2. Perform computation on GPU
  - Efficient CUDA kernels





# Proposed Gather-first MPI\_Reduce / MPI\_Scan

- Gather-first algorithm
  - Root gathers all the data and perform the computation
  - Low computation overhead
  - Poor scalability

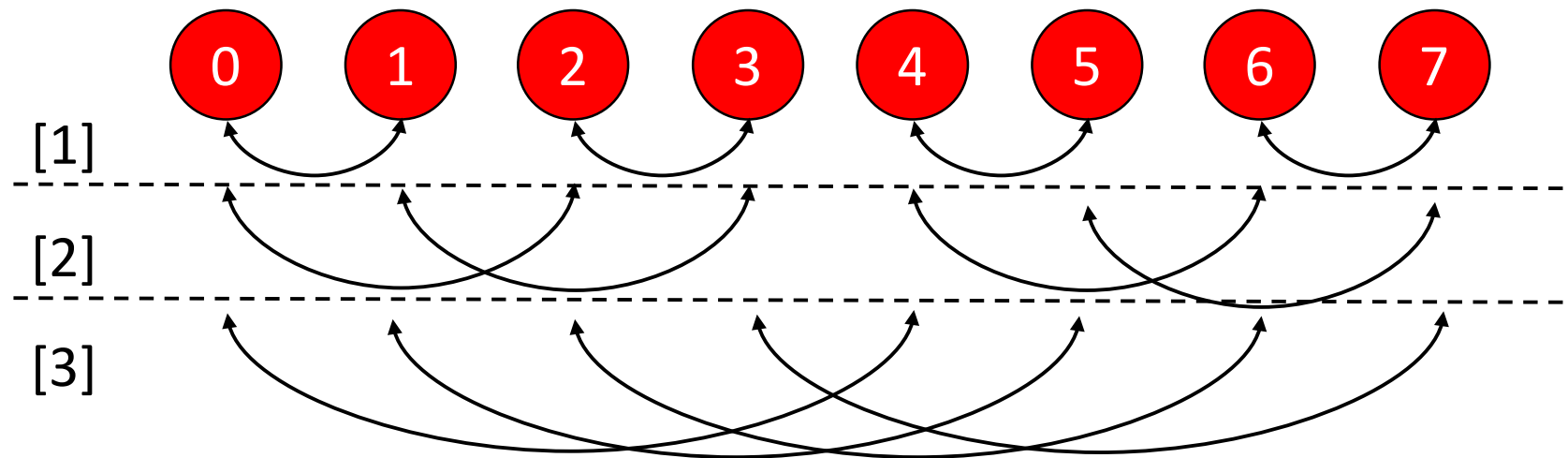


$$\underline{(n - 1) \times (Comm_{Host}(M) + Comp_{Host}(M))} + 2 \times Copy(M)$$

*Good for small messages and small scale*

# Proposed GPU-enabled MPI\_Allreduce / MPI\_Scan

- GPU-enabled Recursive doubling algorithm
  - Every processor needs to perform computation
  - Load balance, Efficient/scalable communication
  - Higher average latency



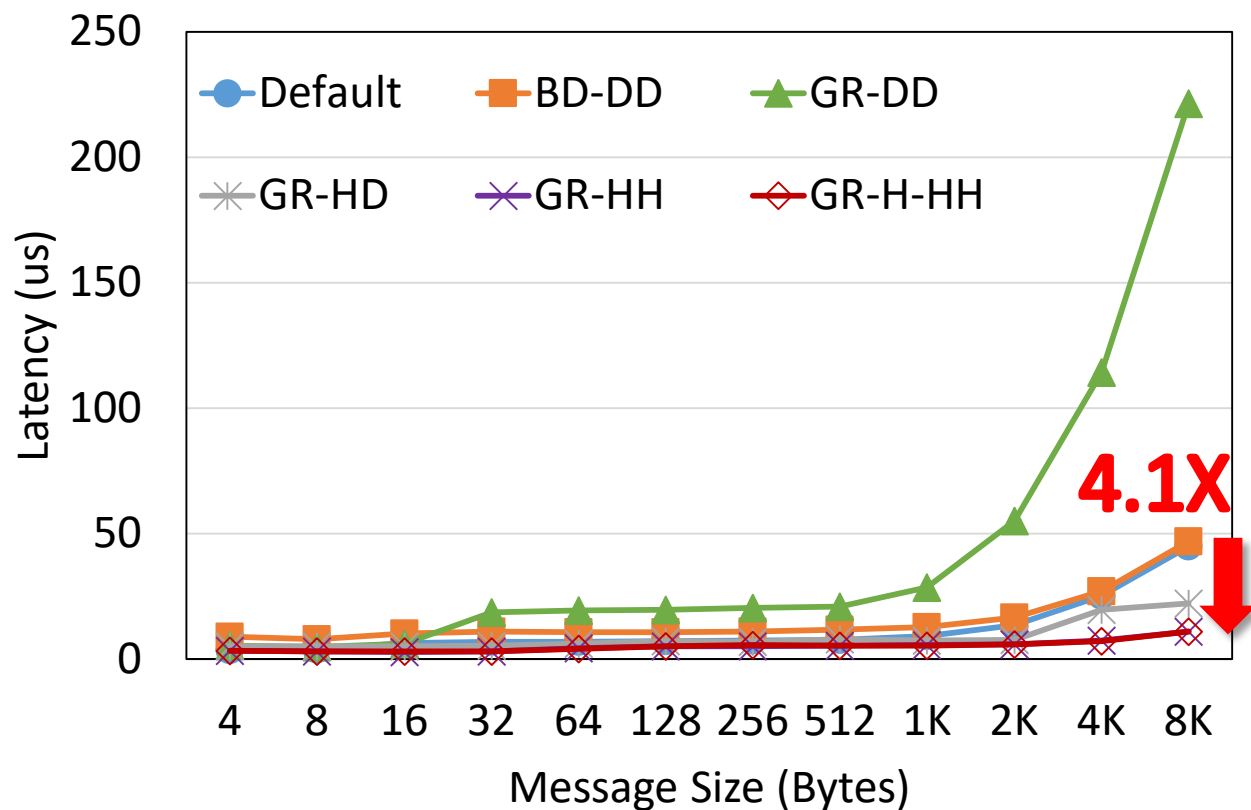
$$[\log_2 n] \times (\epsilon \times Comm_{GDR}(M) + Overhead_{GPU}(M) + Comp_{GPU}(M))$$

# Alternative and Extended Designs

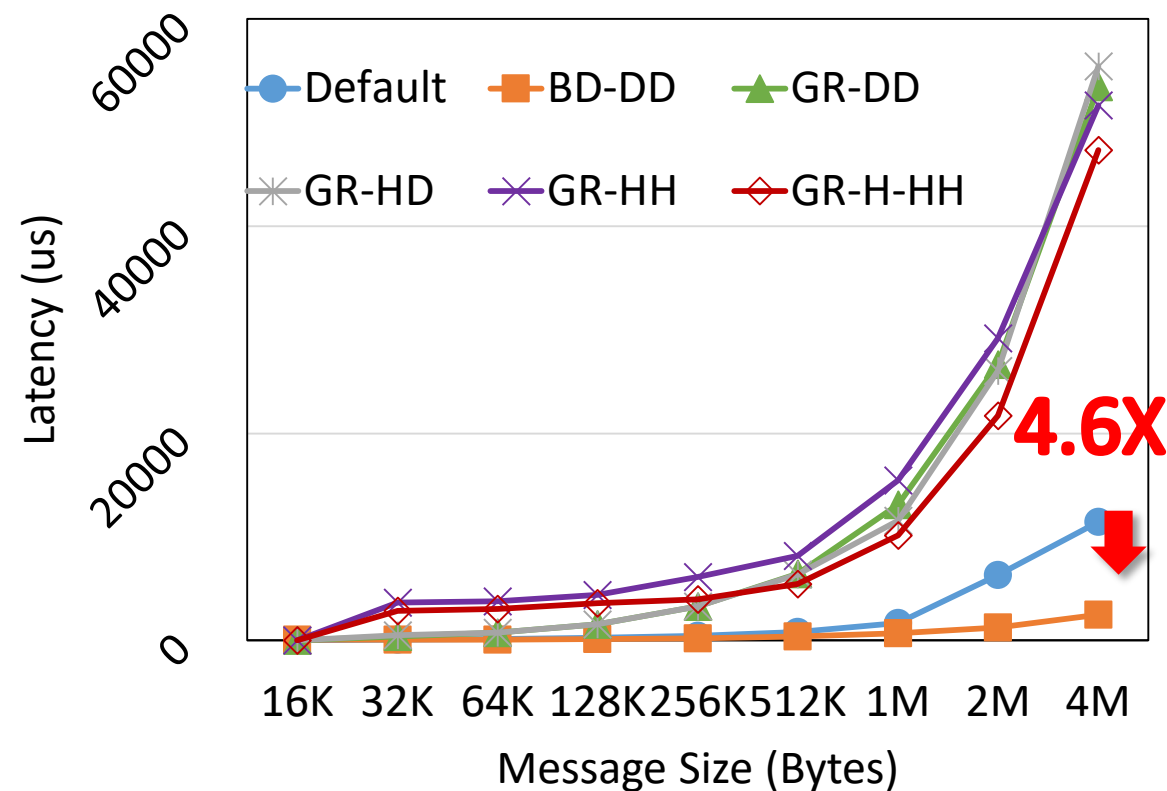
Communication	Computation	Design	Algorithm	Benefit
Host<->Host	CPU	<b>BR-H-HH (Default)</b>	Binomial-Reduce	<i>Large scale, small messages</i>
		<b>RD-H-HH (Default)</b>	Recursive doubling	
		<b>GR-H-HH</b>	Gather-Reduce	<i>Small scale, small messages</i>
<b>GR-HH</b>				
<b>GR-HD / GR-DH</b>				
Host<->Device (GDR)				
Device<->Device (GDR)	GPU	<b>GR-DD</b>	Binomial-Reduce	<i>Large messages for any scale</i>
		<b>BR-DD</b>		
		<b>BRB-DD</b>	Binomial-Reduce-Bcast	
		<b>RD-DD</b>	Recursive doubling	
Host<->Device (GDR)	<b>RD-HD/RD-DH</b>			

# Evaluation - MPI\_Reduce @ CSCS (96 GPUs)

Gather-first approaches  
win for small messages

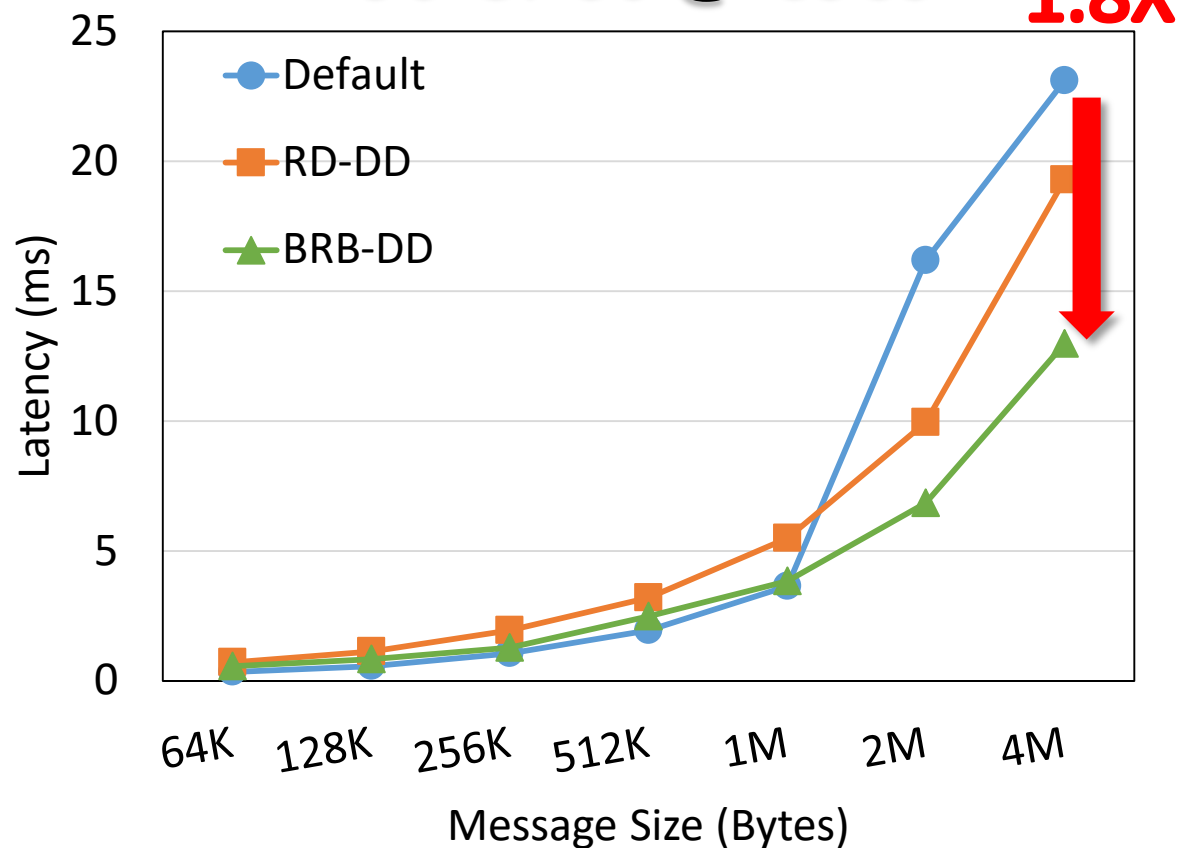


K-nomial GPU-based approach win  
for large messages



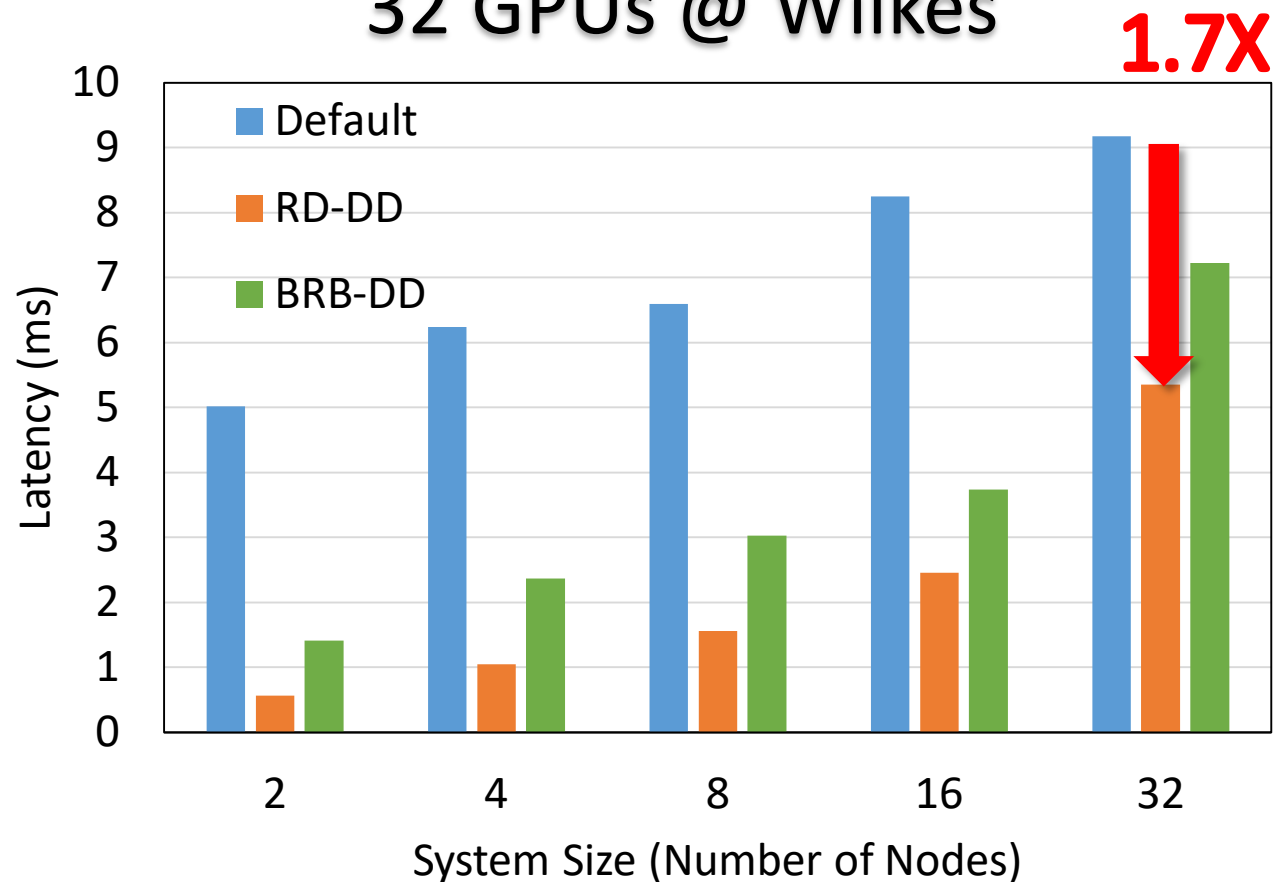
# Evaluation - MPI\_Allreduce

## 96 GPUs @ CSCS



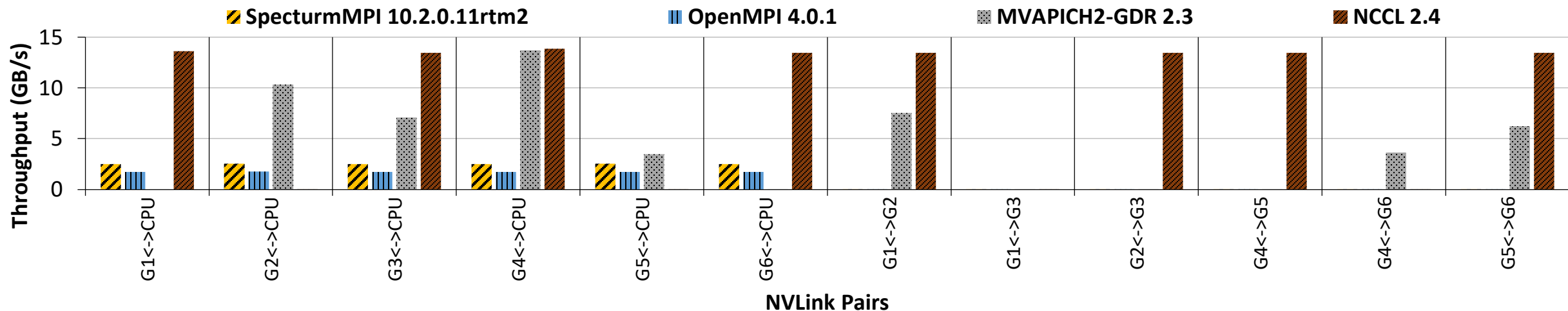
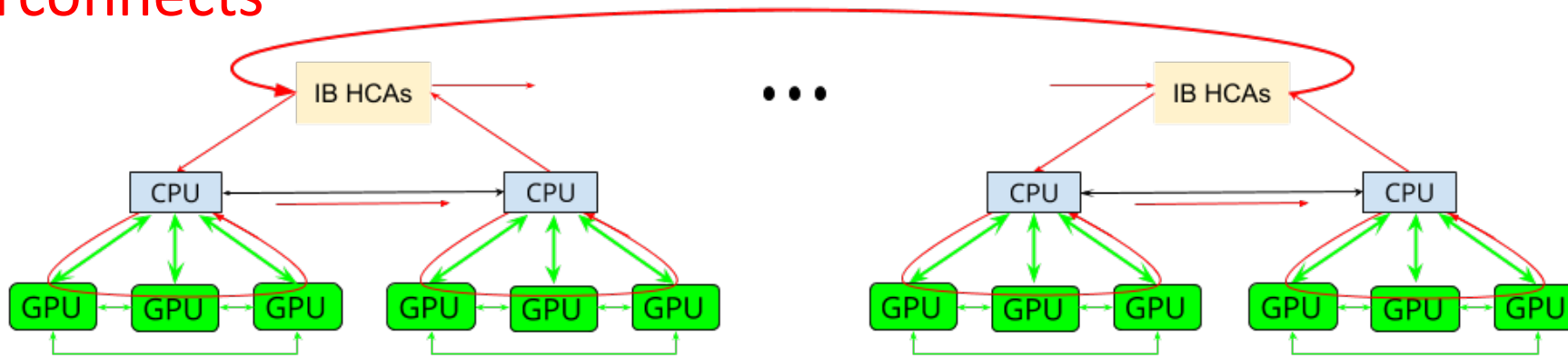
## Good Scalability

## 32 GPUs @ Wilkes



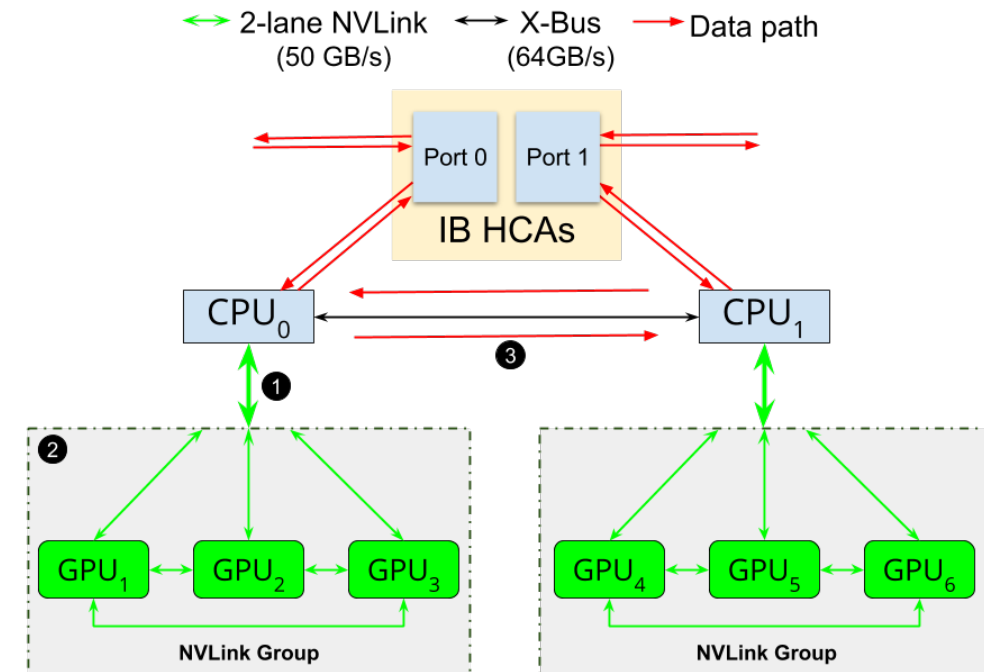
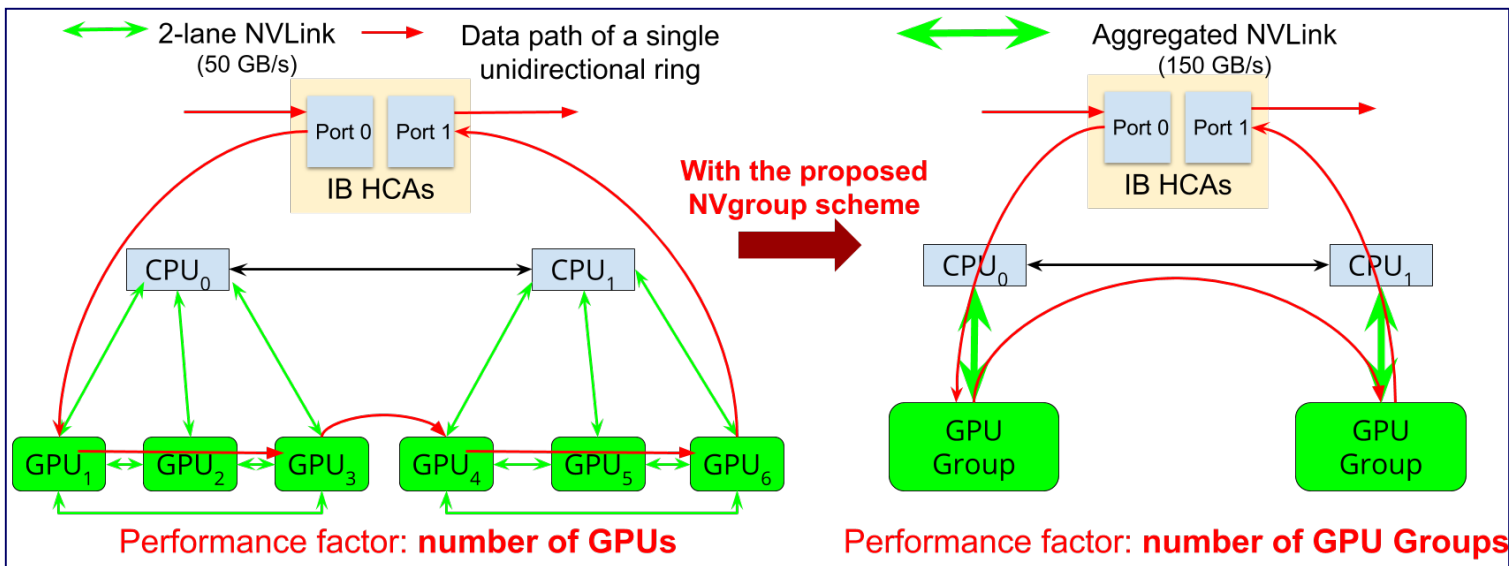
# Allreduce Operations in Modern Dense-GPU System

- Ring-based Allreduce for DL workloads **cannot efficiently utilize fast interconnects**

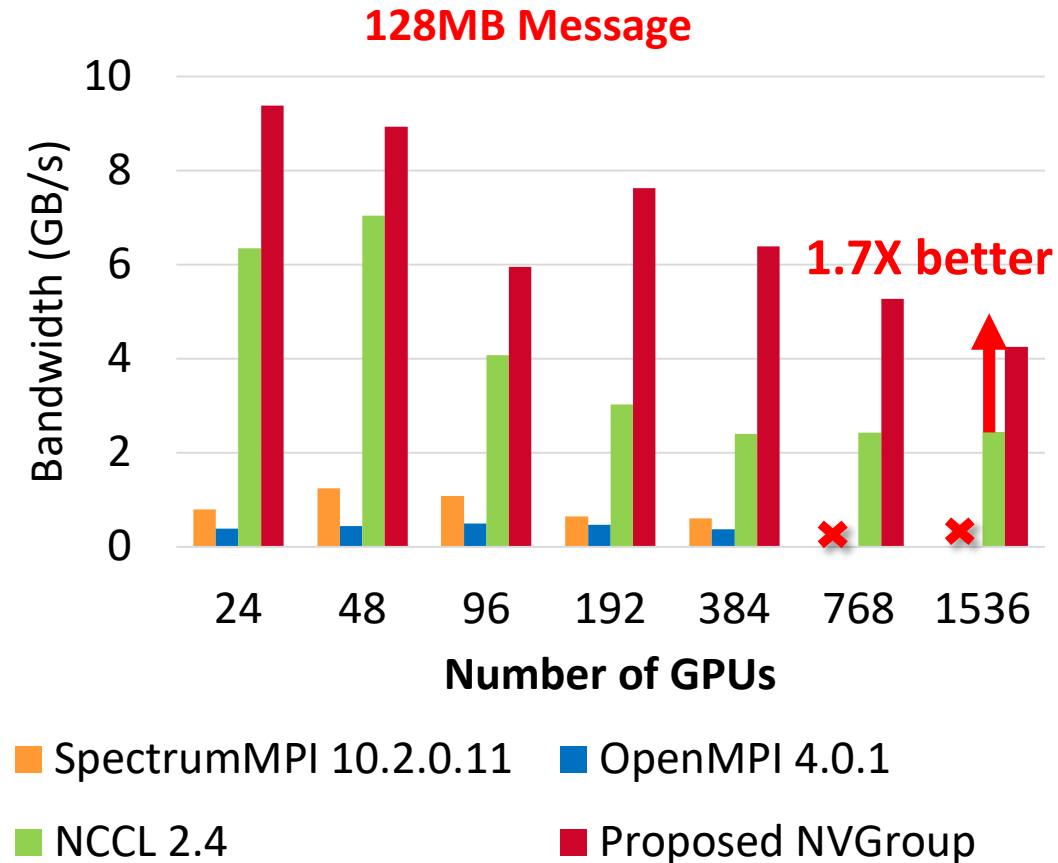
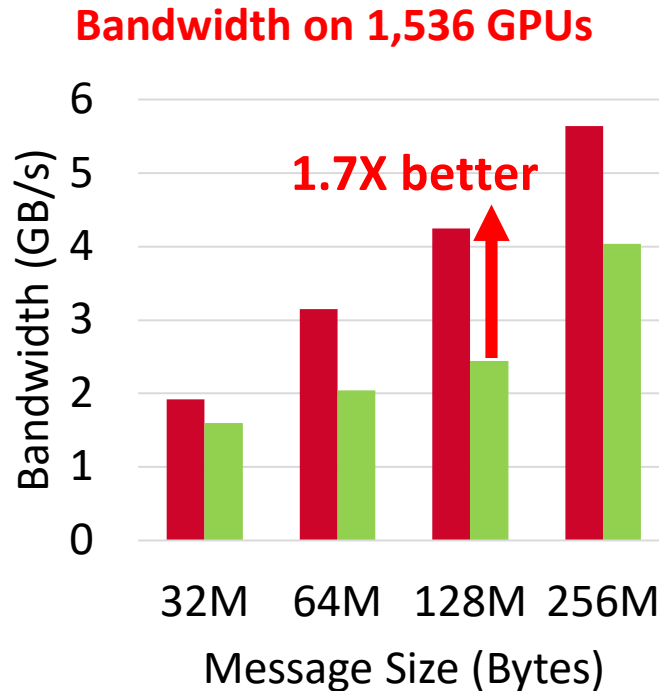
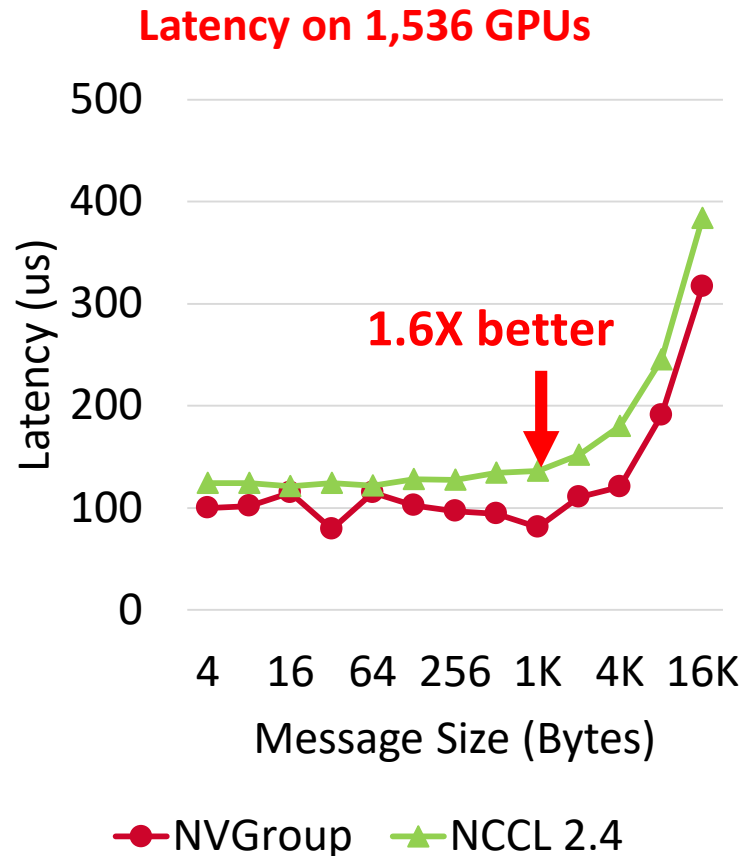


# Topology-aware Allreduce on Dense-GPU Clusters

- Grouping GPUs which are fully connected by NVLinks
  - **Contention-free** communication within the group
- Cooperative Reduction Kernels to exploit **load-compute-store** primitives over NVLinks



# Preliminary Results – Allreduce Benchmark

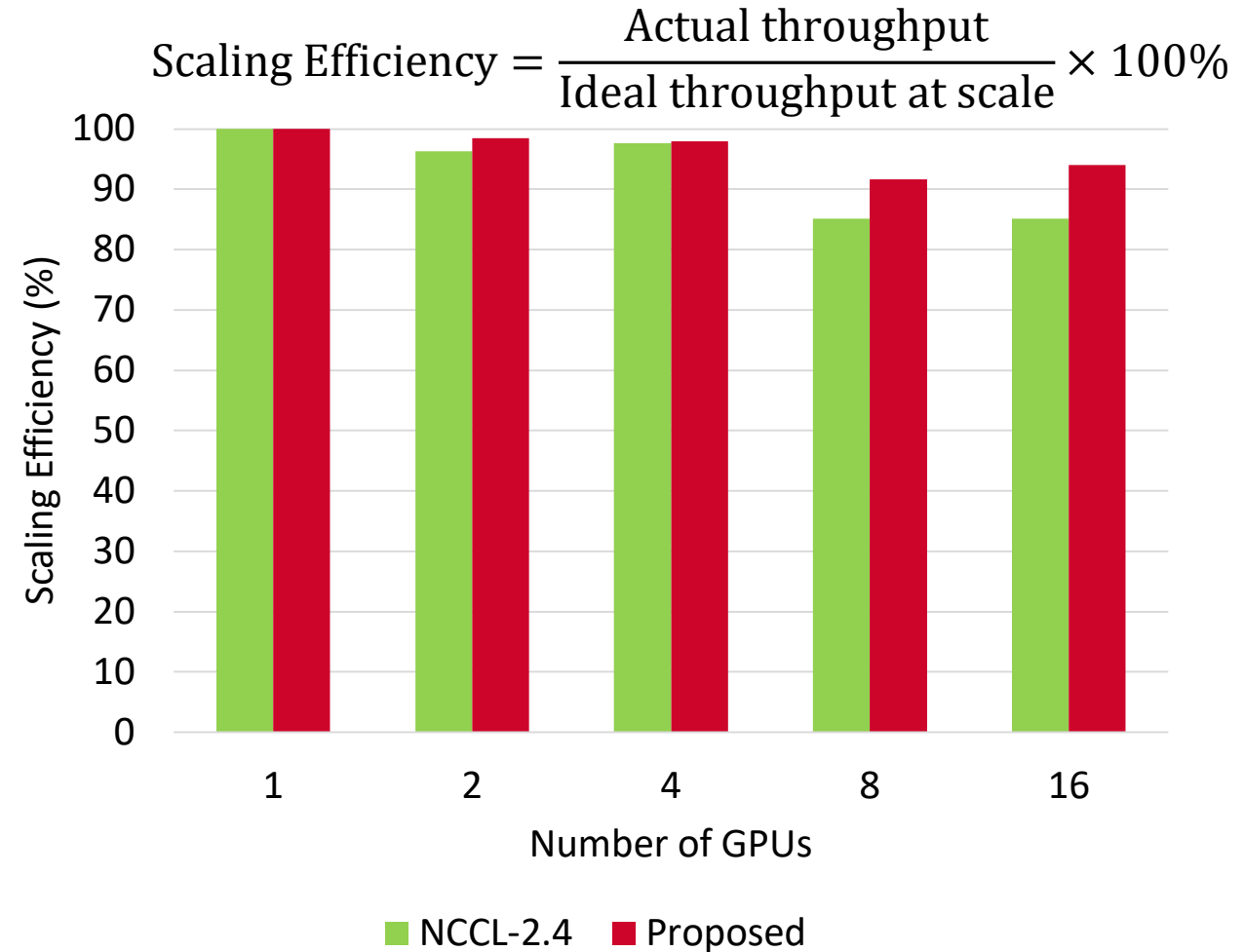
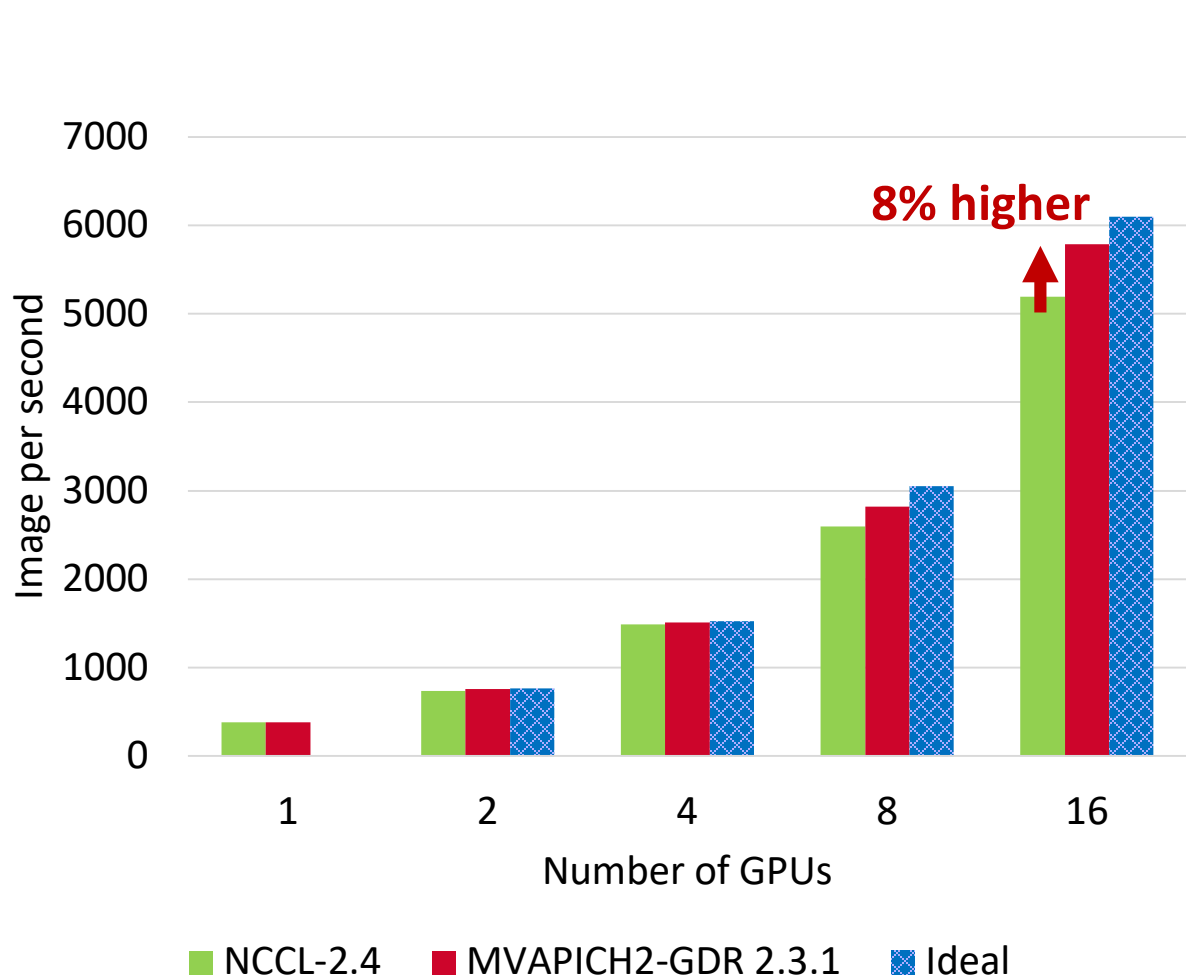


#1 Summit Platform: Dual-socket IBM POWER9 CPU, 6 NVIDIA Volta V100 GPUs, and 2-port InfiniBand EDR Interconnect



# Preliminary Results – Distributed Deep Learning Training

- ResNet-50 Training using TensorFlow benchmark on a DGX-2 machine (16 Volta GPUs)



# MVAPICH2 Project

- High Performance open-source MPI Library for InfiniBand, Omni-Path, Ethernet/iWARP, and RDMA over Converged Ethernet (RoCE)
  - MVAPICH (MPI-1), MVAPICH2 (MPI-2.2 and MPI-3.1), Started in 2001, First version available in 2002
  - **MVAPICH2-X (MPI + PGAS), Available since 2011**
  - **Support for GPGPUs (MVAPICH2-GDR) and MIC (MVAPICH2-MIC), Available since 2014**
  - **Support for Virtualization (MVAPICH2-Virt), Available since 2015**
  - **Support for Energy-Awareness (MVAPICH2-EA), Available since 2015**
  - **Support for InfiniBand Network Analysis and Monitoring (OSU INAM) since 2015**
  - **Used by more than 3,000 organizations in 89 countries**
  - **More than 553,000 (> 0.5 million) downloads from the OSU site directly**
  - Empowering many TOP500 clusters (June '19 ranking)
    - **3<sup>rd</sup> ranked 10,649,640-core cluster (Sunway TaihuLight) at NSC, Wuxi, China**
    - 16<sup>th</sup>, 556,104 cores (Oakforest-PACS) in Japan
    - 19<sup>th</sup>, 367,024 cores (Stampede2) at TACC
    - 31<sup>st</sup>, 241,108-core (Pleiades) at NASA and many others
  - Available with software stacks of many vendors and Linux Distros (RedHat, SuSE, and OpenHPC)
  - <http://mvapich.cse.ohio-state.edu>

*Empowering Top500 systems for over a decade*



**Partner in the 5<sup>th</sup> ranked TACC Frontera System**

# Thank You!

## Questions?

[chu.368@osu.edu](mailto:chu.368@osu.edu)

<http://web.cse.ohio-state.edu/~chu.368>