



MVA PICH

MPI, PGAS and Hybrid MPI+PGAS Library



High-Performance
Big Data



High-Performance
Deep Learning

How to Design Convergent HPC, Deep Learning and Big Data Analytics Software Stacks for Exascale Systems?

Keynote Talk at SCAsia (Mar '19)

by

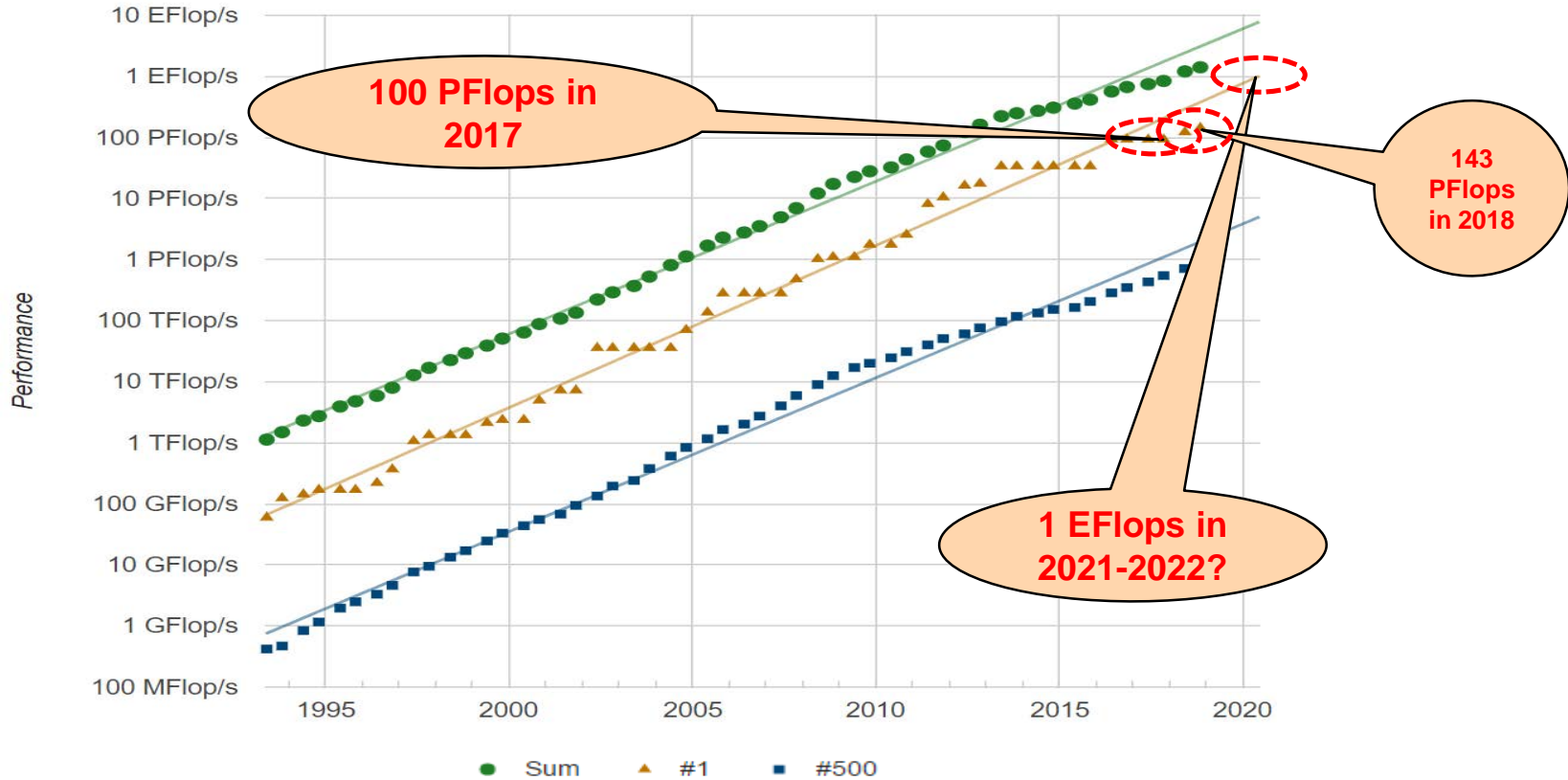
Dhabaleswar K. (DK) Panda

The Ohio State University

E-mail: panda@cse.ohio-state.edu

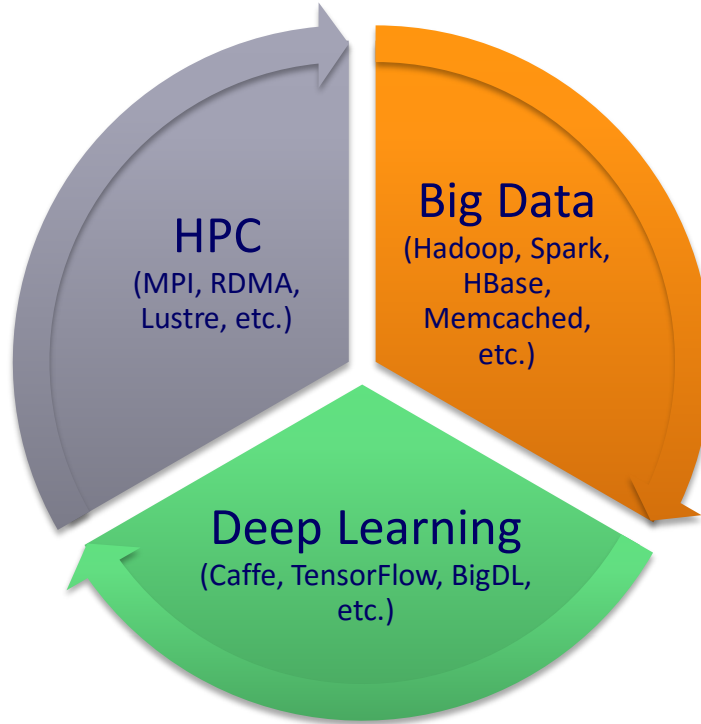
<http://www.cse.ohio-state.edu/~panda>

High-End Computing (HEC): PetaFlop to ExaFlop



Expected to have an ExaFlop system in 2021-2022!

Increasing Usage of HPC, Big Data and Deep Learning



Convergence of HPC, Big Data, and Deep Learning!

Increasing Need to Run these applications on the Cloud!!

Can We Run HPC, Big Data and Deep Learning Jobs on Existing HPC Infrastructure?



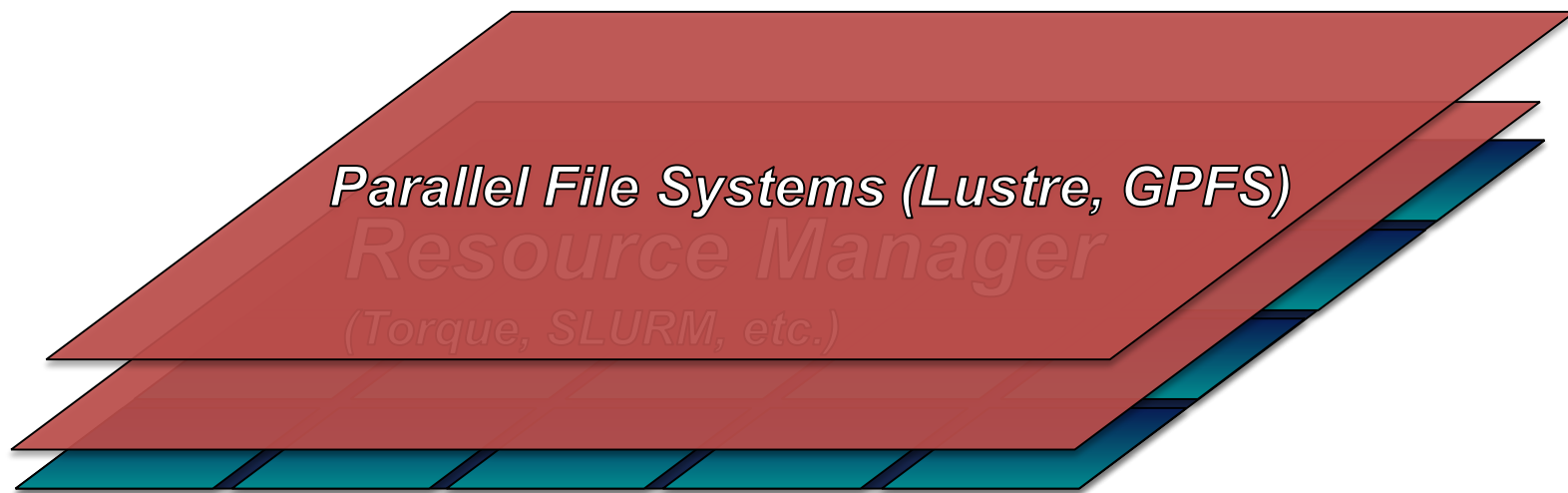
Physical Compute

Can We Run HPC, Big Data and Deep Learning Jobs on Existing HPC Infrastructure?

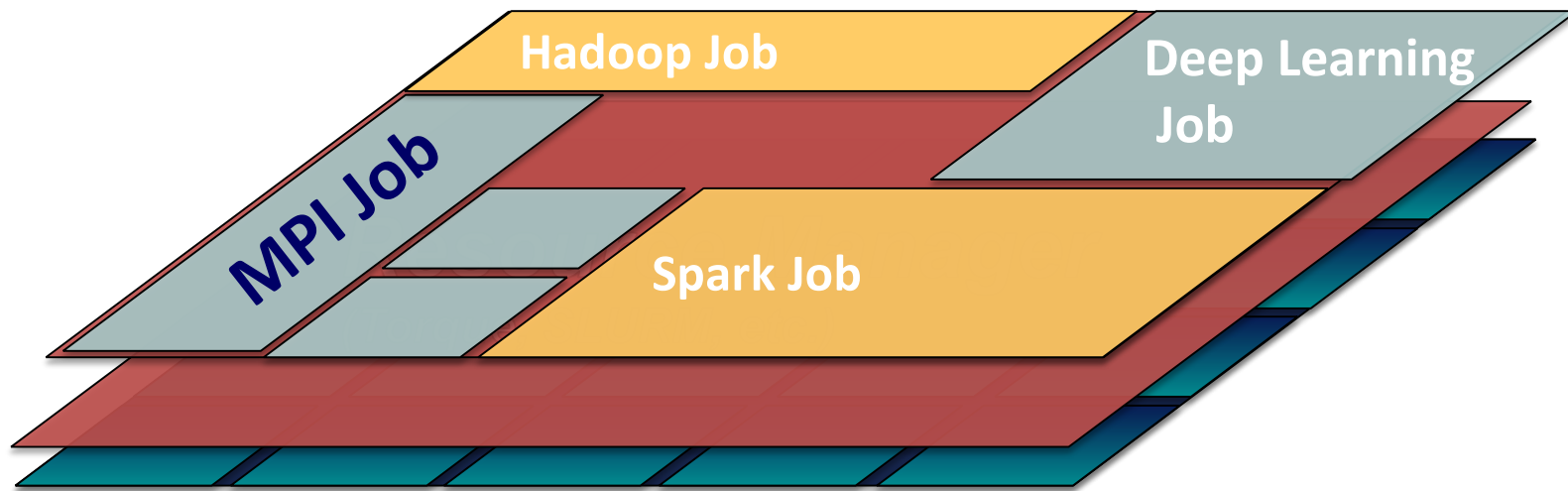


Resource Manager
(Torque, SLURM, etc.)

Can We Run HPC, Big Data and Deep Learning Jobs on Existing HPC Infrastructure?



Can We Run HPC, Big Data and Deep Learning Jobs on Existing HPC Infrastructure?



Designing Communication Middleware for Multi-Petaflop and Exaflop Systems: Challenges

Application Kernels/Applications

Middleware

Programming Models

MPI, PGAS (UPC, Global Arrays, OpenSHMEM), CUDA, OpenMP, OpenACC, Cilk, Hadoop (MapReduce), Spark (RDD, DAG), etc.

Communication Library or Runtime for Programming Models

Point-to-point
Communication

Collective
Communication

Energy-
Awareness

Synchronization
and Locks

I/O and
File Systems

Fault
Tolerance

Networking Technologies

(InfiniBand, 40/100GigE,
Aries, and OmniPath)

**Multi/Many-core
Architectures**

**Accelerators
(NVIDIA and FPGA)**

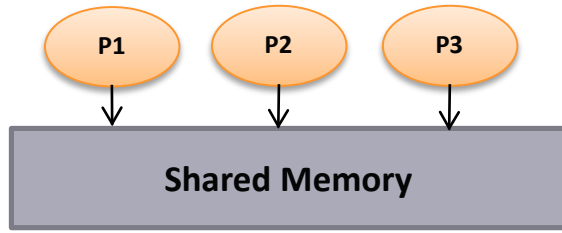
Co-Design
Opportunities
and
Challenges
across Various
Layers

Performance
Scalability
Fault-
Resilience

Presentation Overview

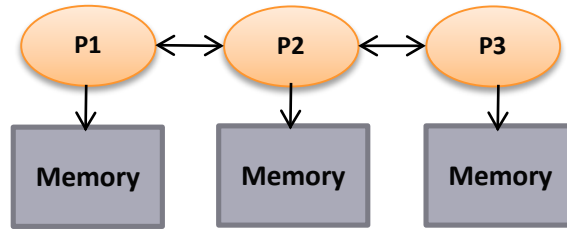
- **MVAPICH Project – MPI and PGAS (MVAPICH) Library with CUDA-Awareness**
- HiDL Project – High-Performance Deep Learning
- HiBD Project – High-Performance Big Data Analytics Library
- Commercial Support from X-ScaleSolutions
- Conclusions and Q&A

Parallel Programming Models Overview



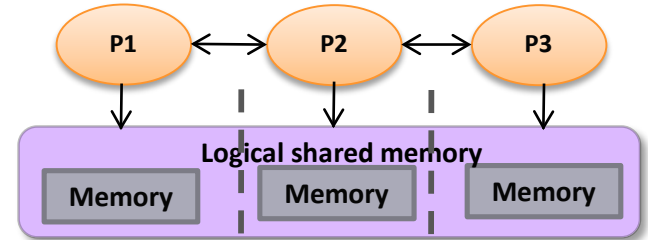
Shared Memory Model

SHMEM, DSM



Distributed Memory Model

MPI (Message Passing Interface)



Partitioned Global Address Space (PGAS)

OpenSHMEM, UPC, Chapel, X10, CAF, ...

- Programming models provide abstract machine models
- Models can be mapped on different types of systems
 - e.g. Distributed Shared Memory (DSM), MPI within a node, etc.
- PGAS models and Hybrid MPI+PGAS models are gradually receiving importance

Overview of the MVAPICH2 Project

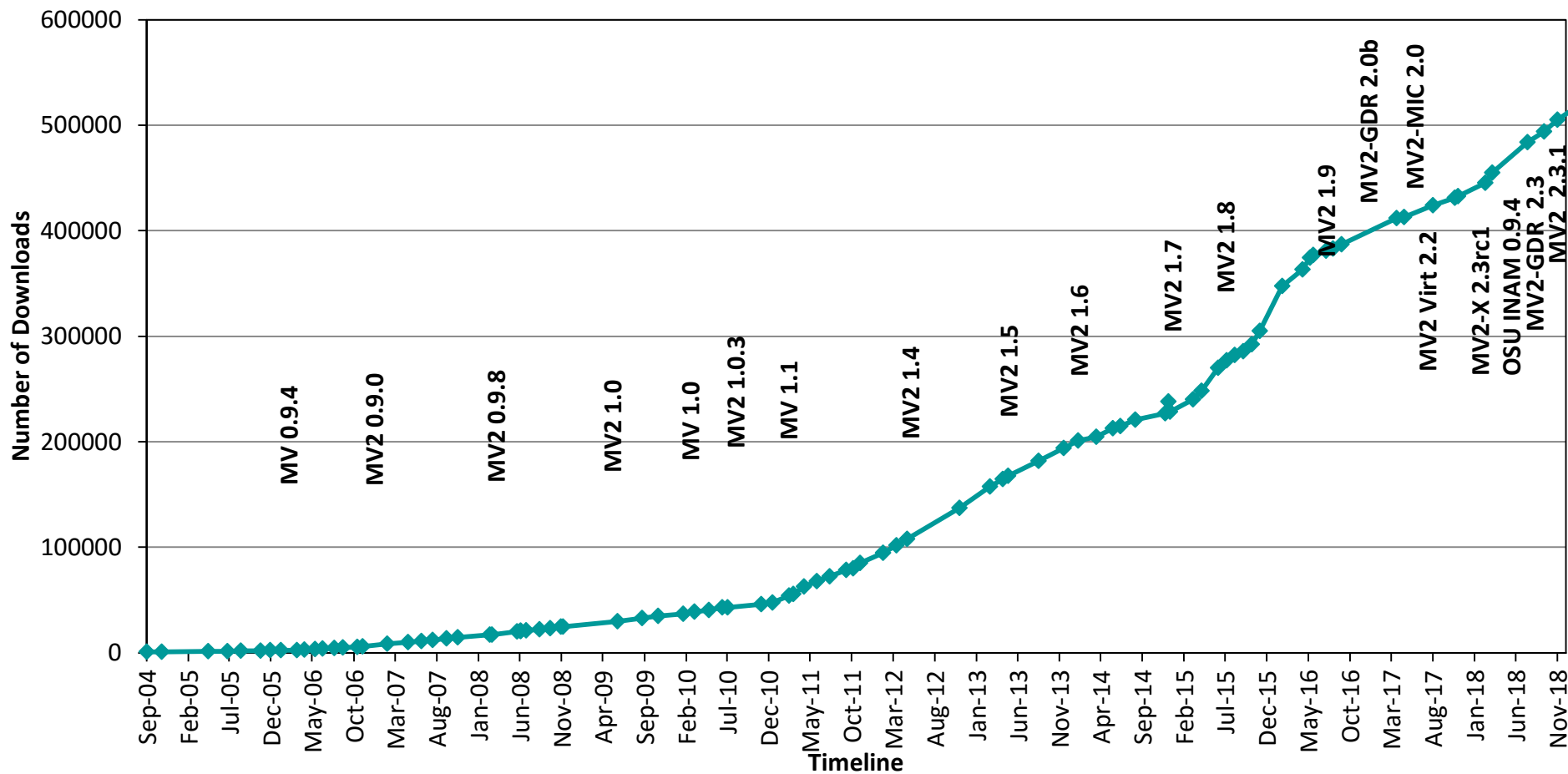
- High Performance open-source MPI Library for InfiniBand, Omni-Path, Ethernet/iWARP, and RDMA over Converged Ethernet (RoCE)
 - MVAPICH (MPI-1), MVAPICH2 (MPI-2.2 and MPI-3.1), Started in 2001, First version available in 2002
 - MVAPICH2-X (MPI + PGAS), Available since 2011
 - Support for GPGPUs (MVAPICH2-GDR) 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 2,950 organizations in 86 countries**
 - **More than 527,000 (> 0.5 million) downloads from the OSU site directly**
 - Empowering many TOP500 clusters (Nov '18 ranking)
 - 3rd ranked 10,649,640-core cluster (Sunway TaihuLight) at NSC, Wuxi, China
 - 14th, 556,104 cores (Oakforest-PACS) in Japan
 - 17th, 367,024 cores (Stampede2) at TACC
 - 27th, 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>



Partner in the upcoming TACC Frontera System

- Empowering Top500 systems for over a decade

MVAPICH2 Release Timeline and Downloads



Architecture of MVAPICH2 Software Family

High Performance Parallel Programming Models

Message Passing Interface
(MPI)

PGAS
(UPC, OpenSHMEM, CAF, UPC++)

Hybrid --- MPI + X
(MPI + PGAS + OpenMP/Cilk)

High Performance and Scalable Communication Runtime

Diverse APIs and Mechanisms

Point-to-point
Primitives

Collectives
Algorithms

Job Startup

Energy-Awareness

Remote
Memory
Access

I/O and
File Systems

Fault
Tolerance

Virtualization

Active
Messages

Introspection
& Analysis

Support for Modern Networking Technology (InfiniBand, iWARP, RoCE, Omni-Path)

Transport Protocols

RC

XRC

UD

DC

Modern Features

UMR

ODP

SR-IOV

Multi
Rail

Support for Modern Multi-/Many-core Architectures (Intel-Xeon, OpenPOWER, Xeon-Phi, ARM, NVIDIA GPGPU)

Transport Mechanisms

Shared
Memory

CMA

IVSHMEM

XPMEM

Modern Features

MCDRAM*

NVLink

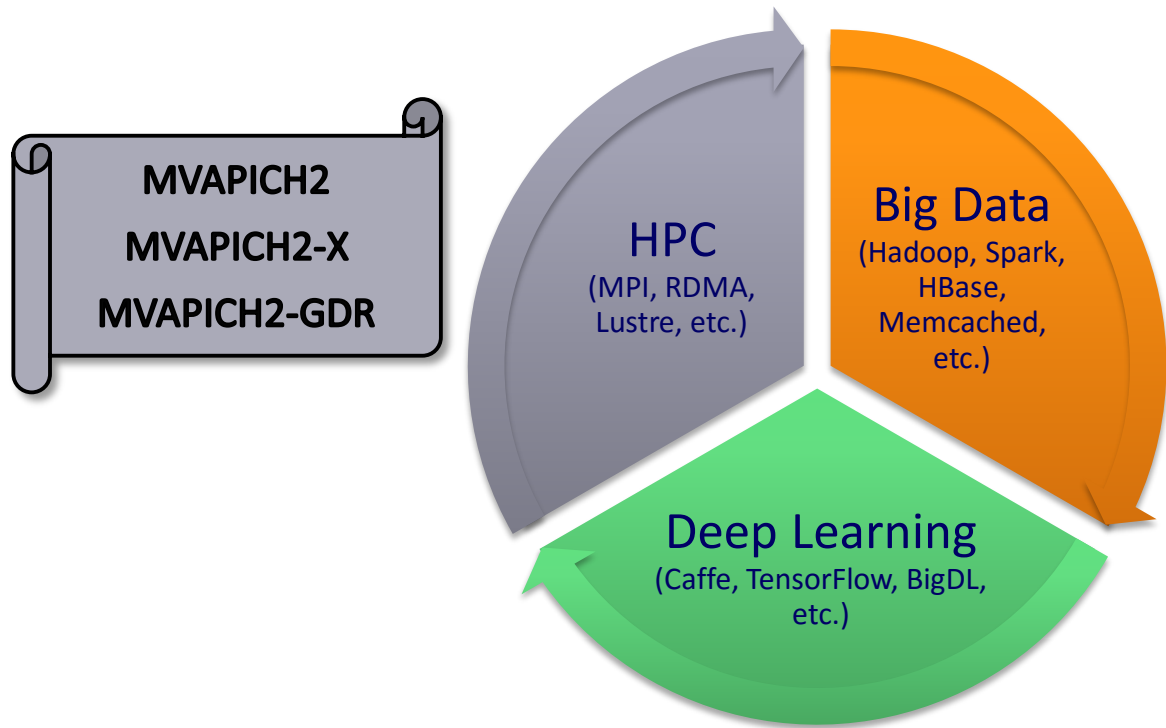
CAPI*

* Upcoming

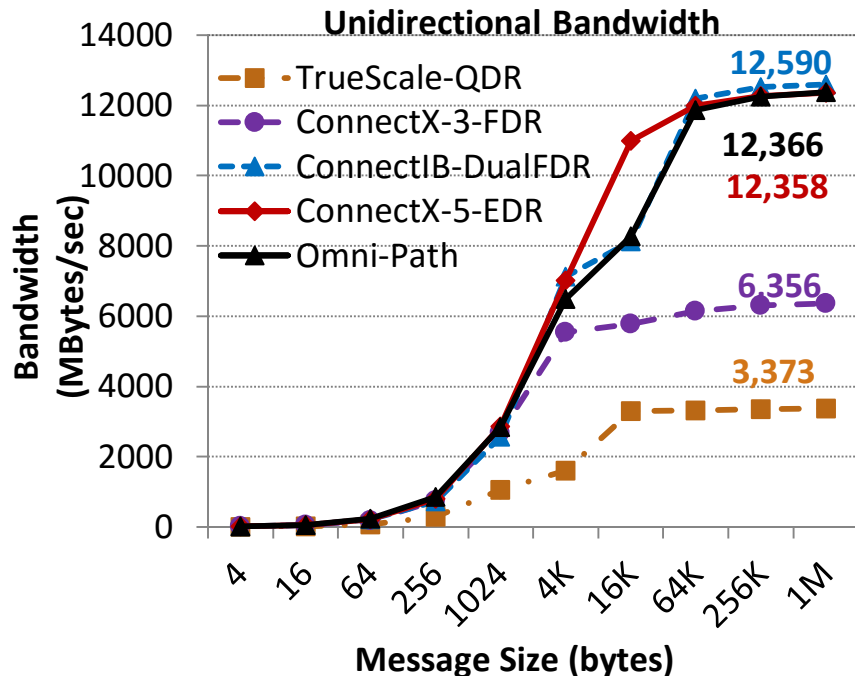
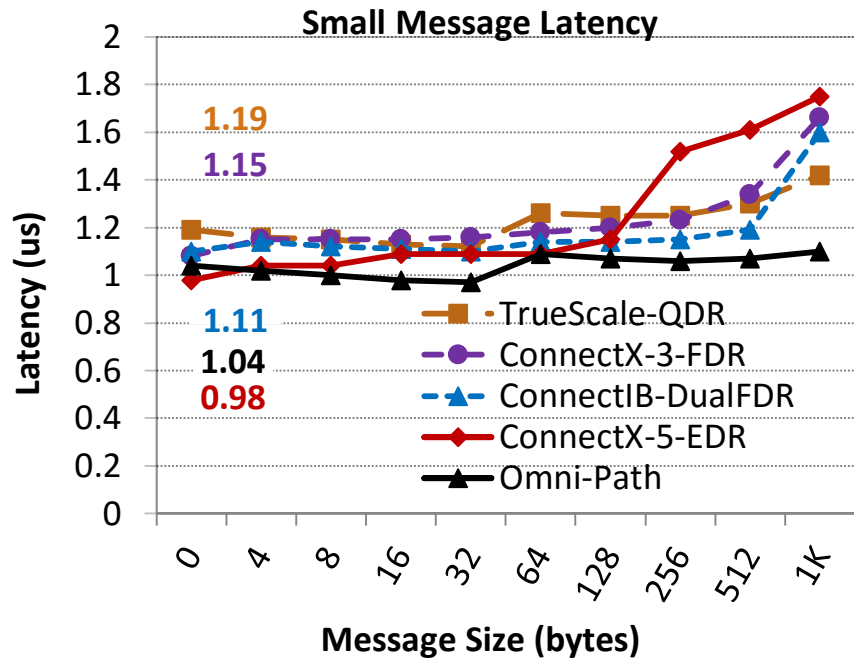
MVAPICH2 Software Family

Requirements	Library
MPI with IB, iWARP, Omni-Path, and RoCE	MVAPICH2
Advanced MPI Features/Support, OSU INAM, PGAS and MPI+PGAS with IB, Omni-Path, and RoCE	MVAPICH2-X
MPI with IB, RoCE & GPU and Support for Deep Learning	MVAPICH2-GDR
HPC Cloud with MPI & IB	MVAPICH2-Virt
Energy-aware MPI with IB, iWARP and RoCE	MVAPICH2-EA
MPI Energy Monitoring Tool	OEMT
InfiniBand Network Analysis and Monitoring	OSU INAM
Microbenchmarks for Measuring MPI and PGAS Performance	OMB

Convergent Software Stacks for HPC, Big Data and Deep Learning

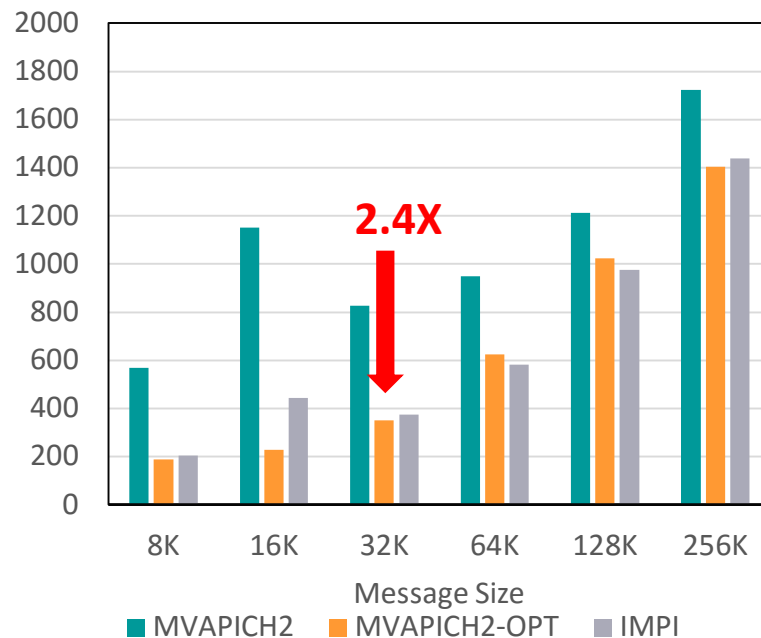
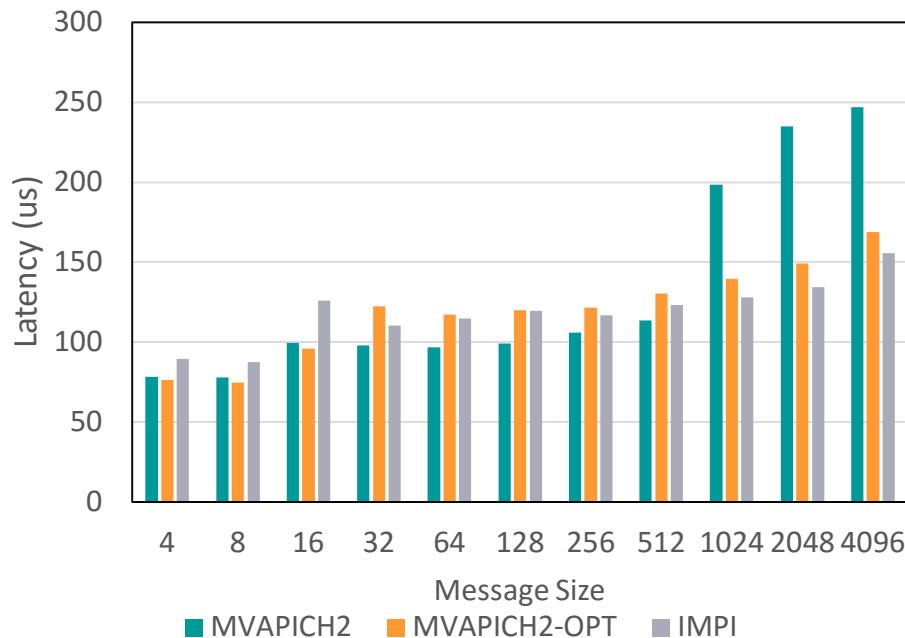


One-way Latency and Bandwidth: MPI over IB with MVAPICH2



- TrueScale-QDR - 3.1 GHz Deca-core (Haswell) Intel PCI Gen3 with IB switch
- ConnectX-3-FDR - 2.8 GHz Deca-core (IvyBridge) Intel PCI Gen3 with IB switch
- ConnectIB-Dual FDR - 3.1 GHz Deca-core (Haswell) Intel PCI Gen3 with IB switch
- ConnectX-5-EDR - 3.1 GHz Deca-core (Haswell) Intel PCI Gen3 with IB Switch
- Omni-Path - 3.1 GHz Deca-core (Haswell) Intel PCI Gen3 with Omni-Path switch

MPI_Allreduce on KNL + Omni-Path (10,240 Processes)



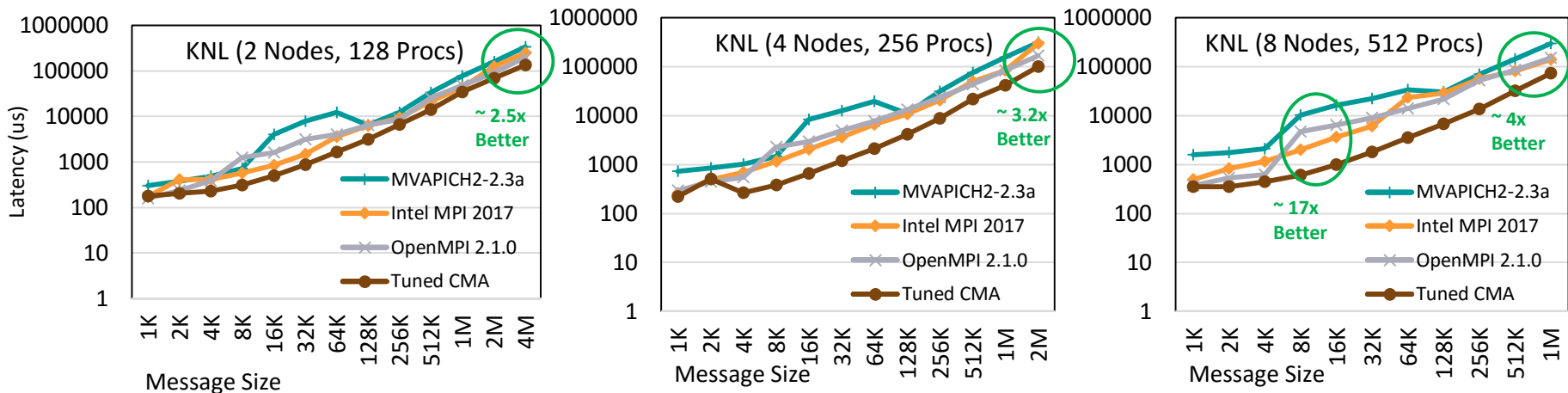
OSU Micro Benchmark 64 PPN

- For MPI_Allreduce latency with 32K bytes, MVAPICH2-OPT can reduce the latency by **2.4X**

M. Bayatpour, S. Chakraborty, H. Subramoni, X. Lu, and D. K. Panda, Scalable Reduction Collectives with Data Partitioning-based Multi-Leader Design, SuperComputing '17.

Available in MVAPICH2-X 2.3b

Optimized CMA-based Collectives for Large Messages



Performance of MPI_Gather on KNL nodes (64PPN)

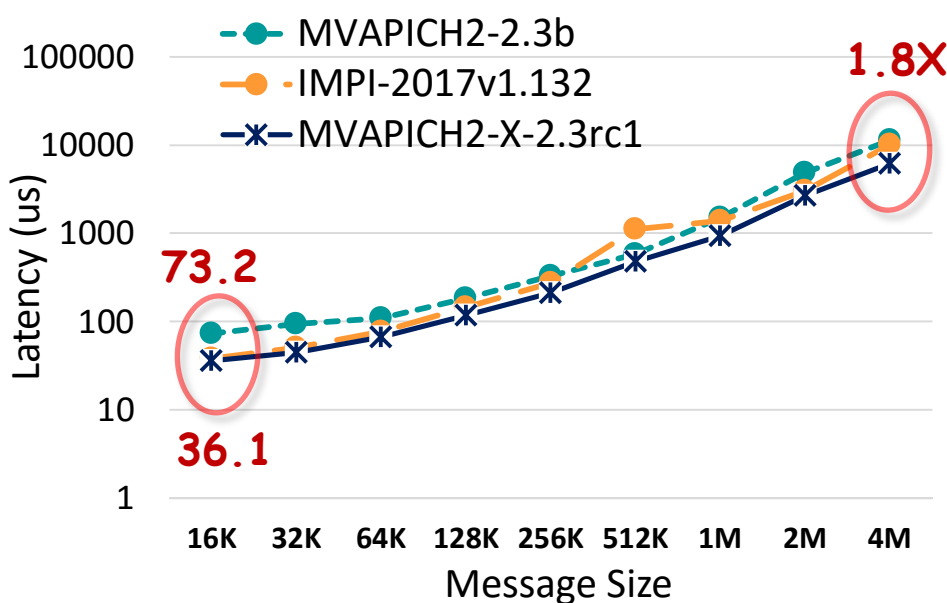
- Significant improvement over existing implementation for Scatter/Gather with 1MB messages (up to 4x on KNL, 2x on Broadwell, 14x on OpenPower)
- New two-level algorithms for better scalability
- Improved performance for other collectives (Bcast, Allgather, and Alltoall)

S. Chakraborty, H. Subramoni, and D. K. Panda, Contention Aware Kernel-Assisted MPI Collectives for Multi/Many-core Systems, IEEE Cluster '17, BEST Paper Finalist

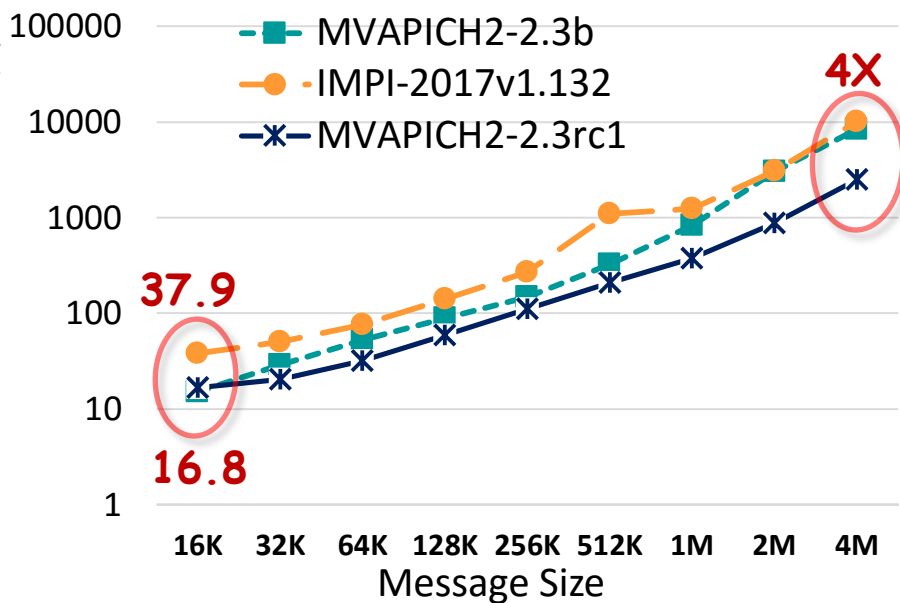
Available in MVAPICH2-X 2.3b

Shared Address Space (XPMEM)-based Collectives Design

OSU_Allreduce (Broadwell 256 procs)



OSU_Reduce (Broadwell 256 procs)

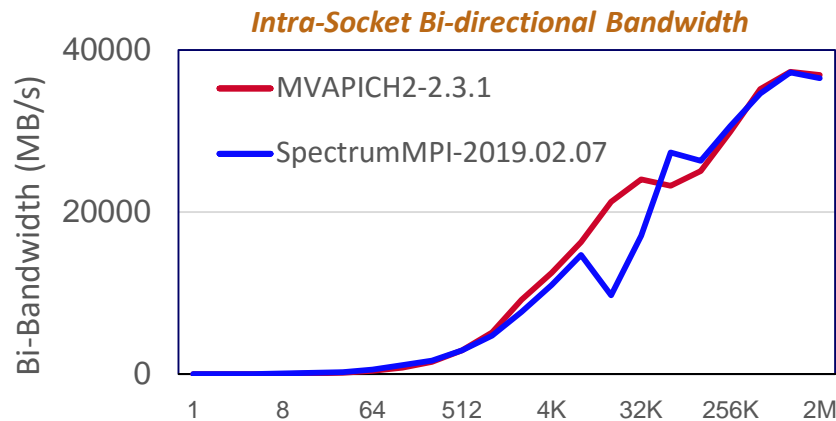
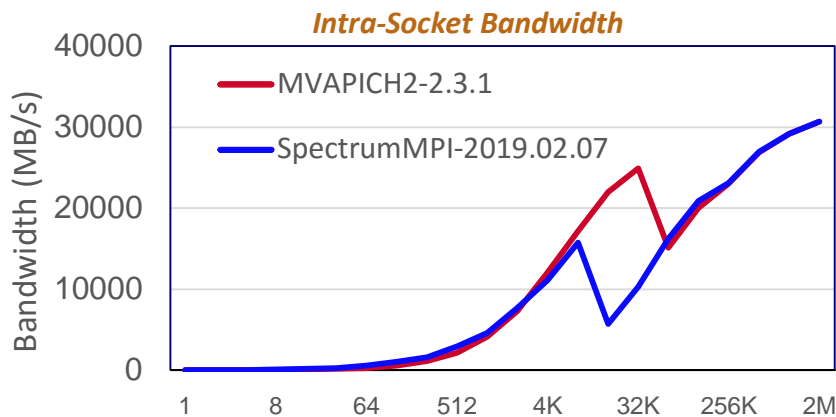
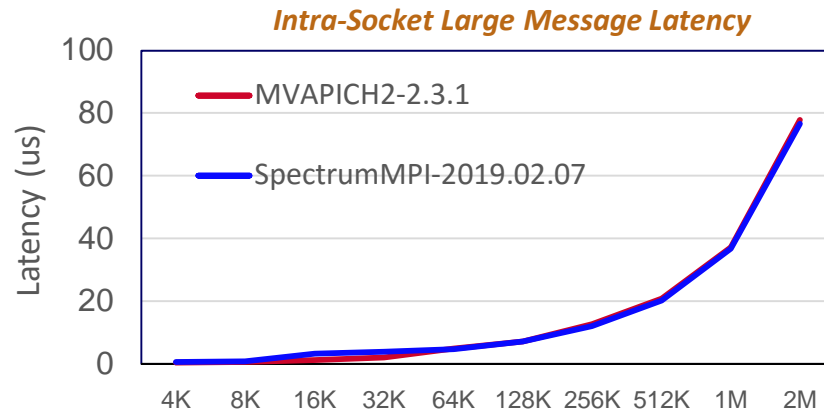
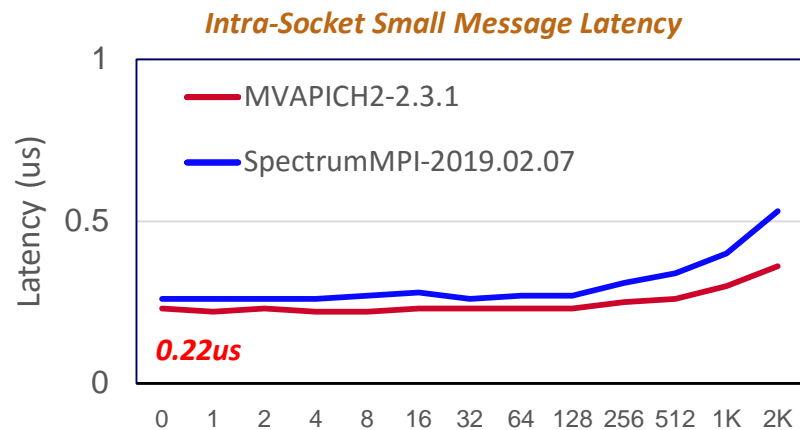


- “Shared Address Space”-based true zero-copy Reduction collective designs in MVAPICH2
- Offloaded computation/communication to peers ranks in reduction collective operation
- Up to **4X** improvement for 4MB Reduce and up to **1.8X** improvement for 4M AllReduce

J. Hashmi, S. Chakraborty, M. Bayatpour, H. Subramoni, and D. Panda, *Designing Efficient Shared Address Space Reduction Collectives for Multi-/Many-cores*, International Parallel & Distributed Processing Symposium (IPDPS '18), May 2018.

Available in MVAPICH2-X 2.3rc1

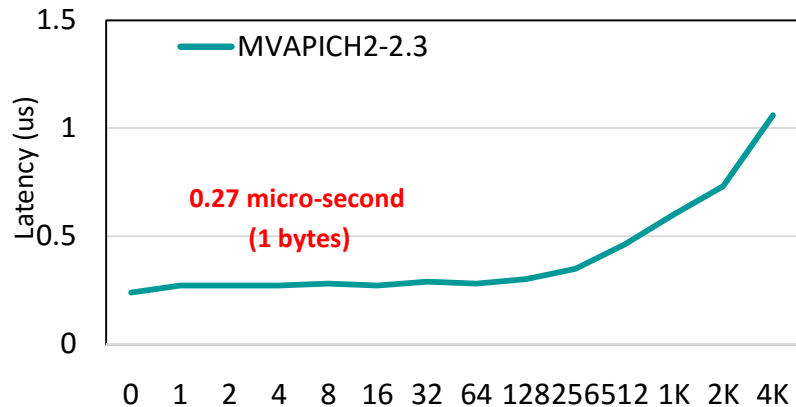
Intra-node Point-to-Point Performance on OpenPower



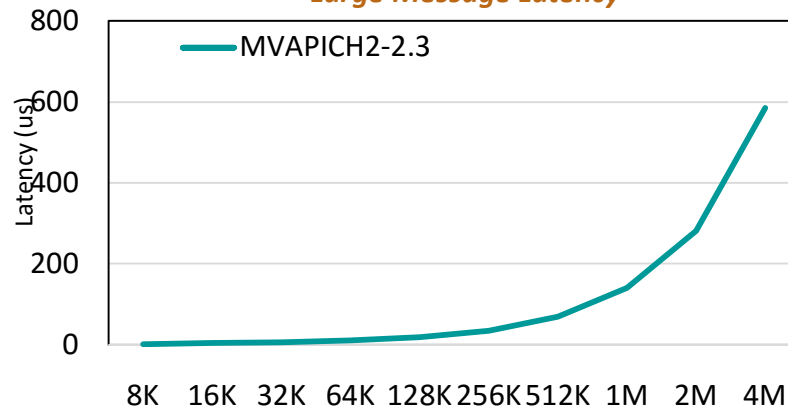
Platform: Two nodes of OpenPOWER (POWER9-ppc64le) CPU using Mellanox EDR (MT4121) HCA

Intra-node Point-to-point Performance on ARM Cortex-A72

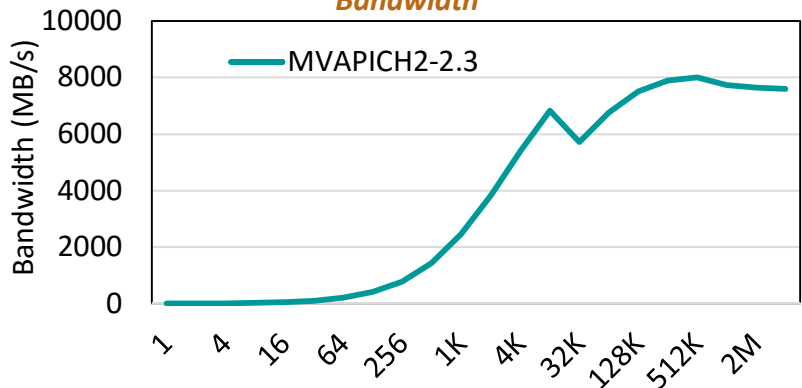
Small Message Latency



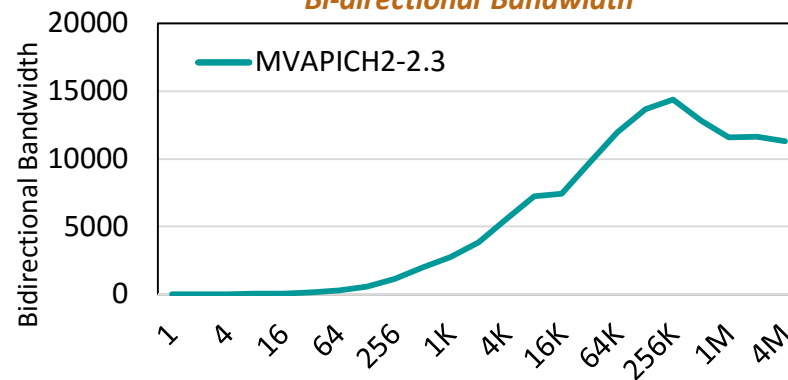
Large Message Latency



Bandwidth

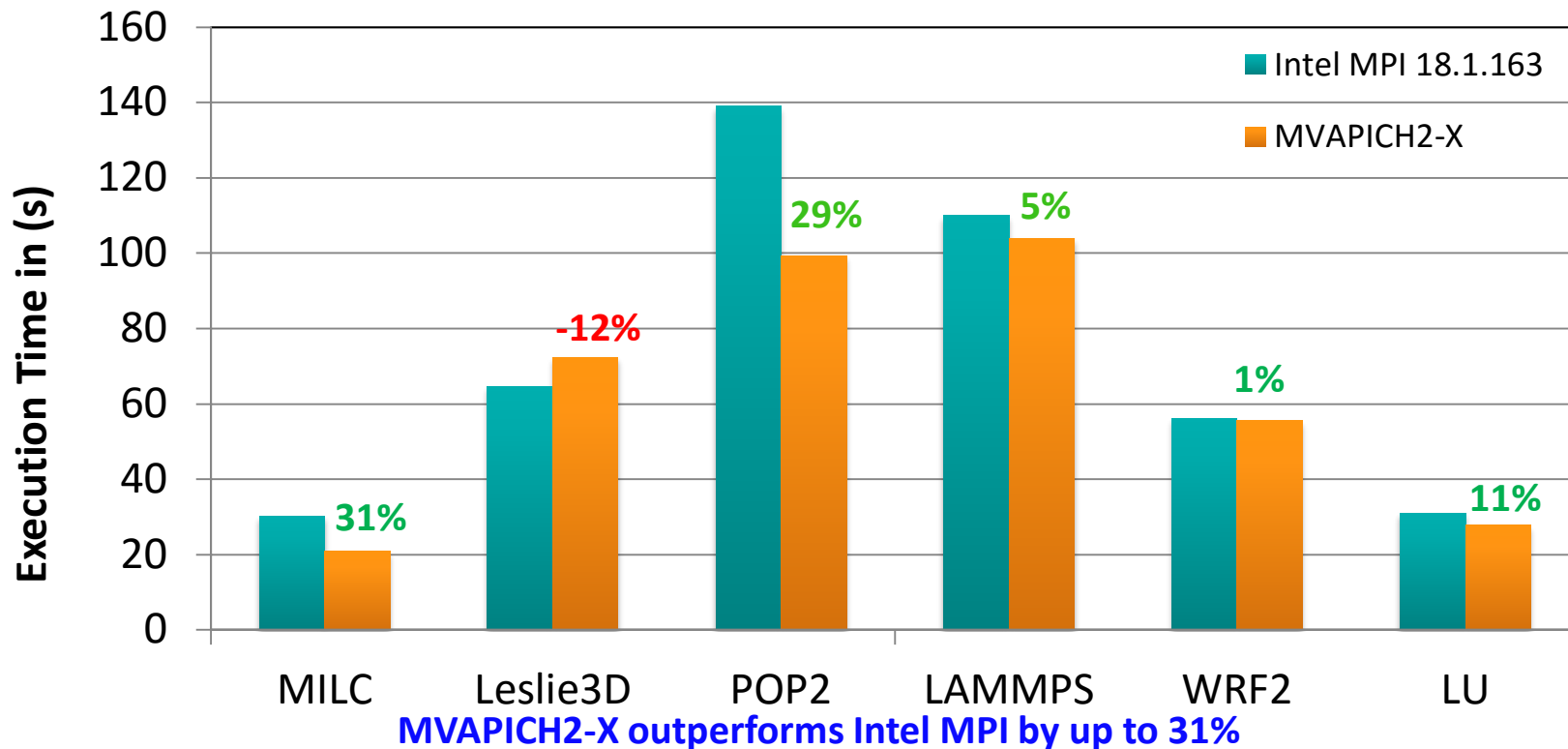


Bi-directional Bandwidth



Platform: ARM Cortex A72 (aarch64) processor with 64 cores dual-socket CPU. Each socket contains 32 cores.

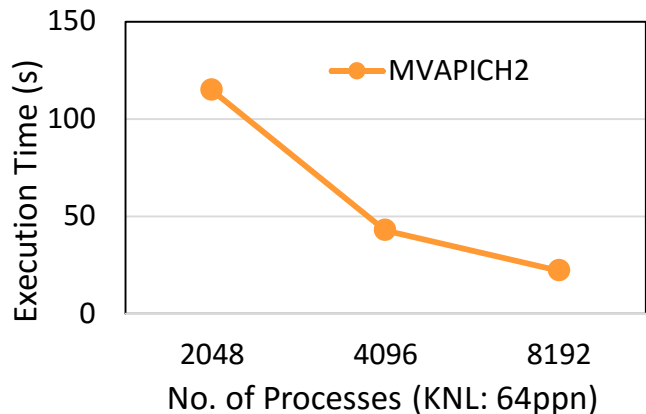
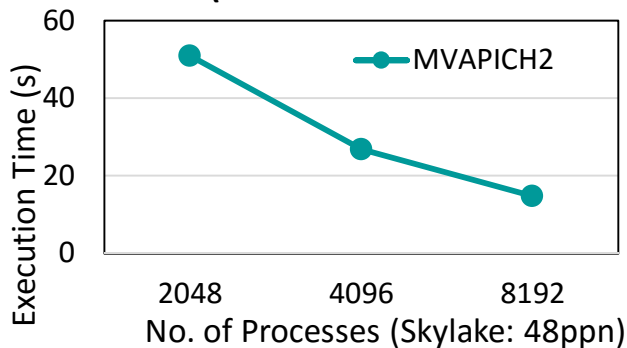
SPEC MPI 2007 Benchmarks: Broadwell + InfiniBand



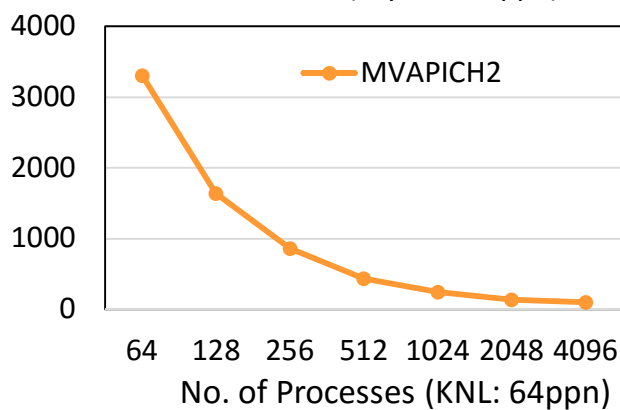
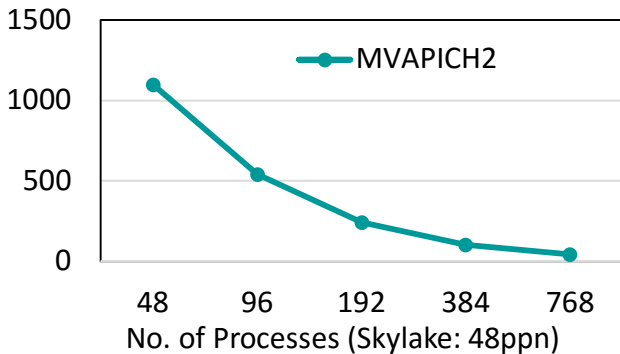
Configuration: 448 processes on 16 Intel E5-2680v4 (Broadwell) nodes having 28 PPN and interconnected with 100Gbps Mellanox MT4115 EDR ConnectX-4 HCA

Application Scalability on Skylake and KNL (Stampeede2)

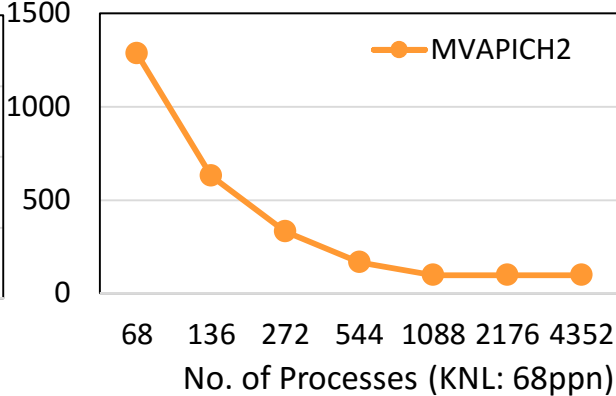
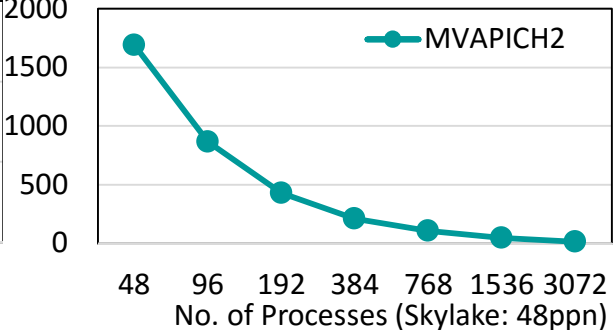
MiniFE (1300x1300x1300 ~ 910 GB)



NEURON (YuEtAl2012)



Cloverleaf (bm64) MPI+OpenMP,
NUM_OMP_THREADS = 2



Courtesy: Mahidhar Tatineni @SDSC, Dong Ju (DJ) Choi@SDSC, and Samuel Khuviz@OSC ---- Testbed: TACC Stampede2 using MVAPICH2-2.3b

Runtime parameters: MV2_SMPI_LENGTH_QUEUE=524288 PSM2_MQ_RNDV_SHM_THRESH=128K PSM2_MQ_RNDV_HFI_THRESH=128K

GPU-Aware (CUDA-Aware) MPI Library: MVAPICH2-GPU

- Standard MPI interfaces used for unified data movement
- Takes advantage of Unified Virtual Addressing (\geq CUDA 4.0)
- Overlaps data movement from GPU with RDMA transfers

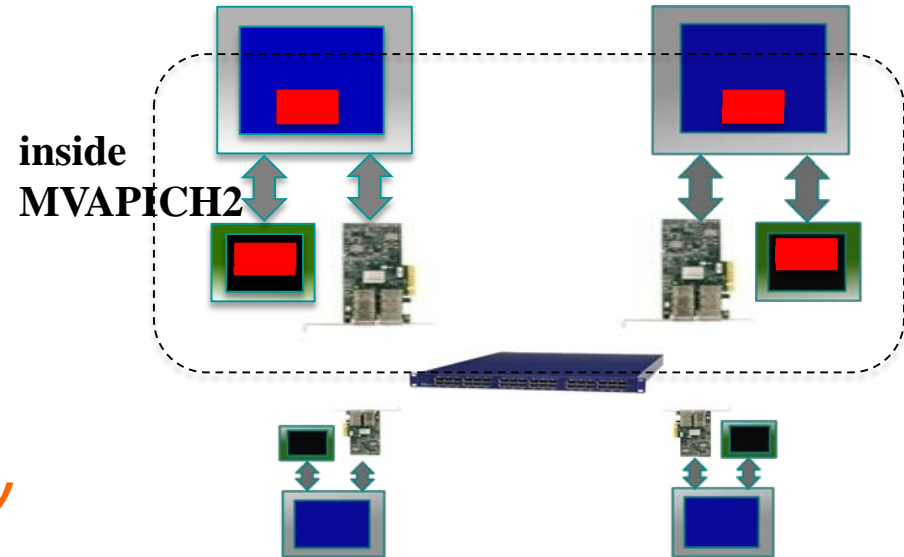
At Sender:

```
MPI_Send(s_devbuf, size, ...);
```

At Receiver:

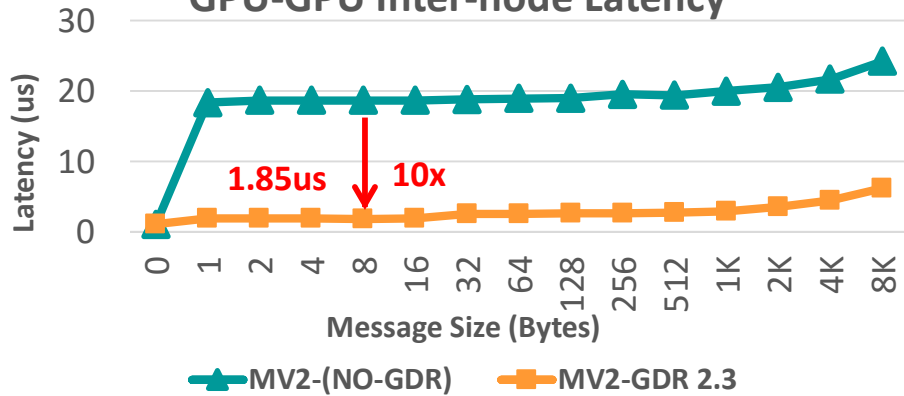
```
MPI_Recv(r_devbuf, size, ...);
```

High Performance and High Productivity

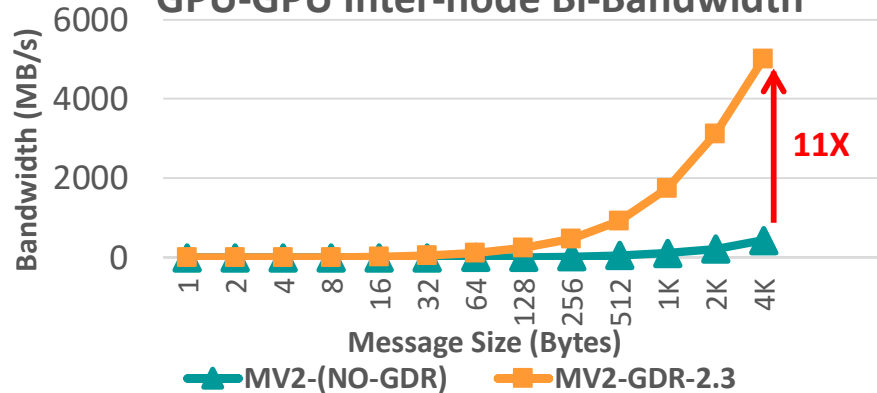


Optimized MVAPICH2-GDR Design

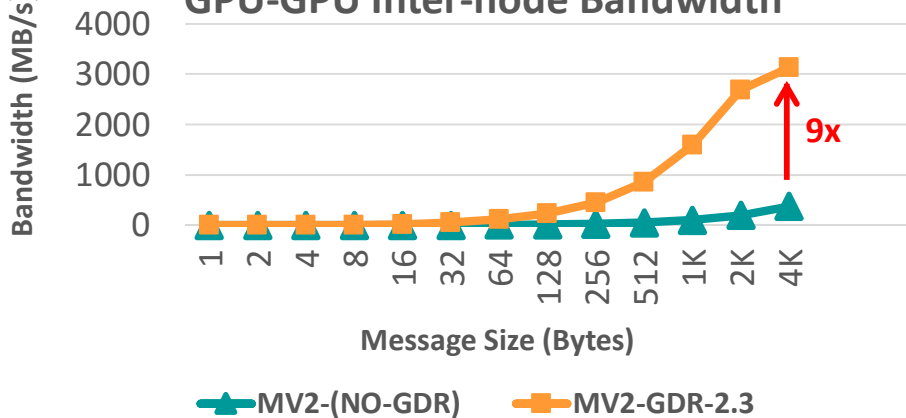
GPU-GPU Inter-node Latency



GPU-GPU Inter-node Bi-Bandwidth



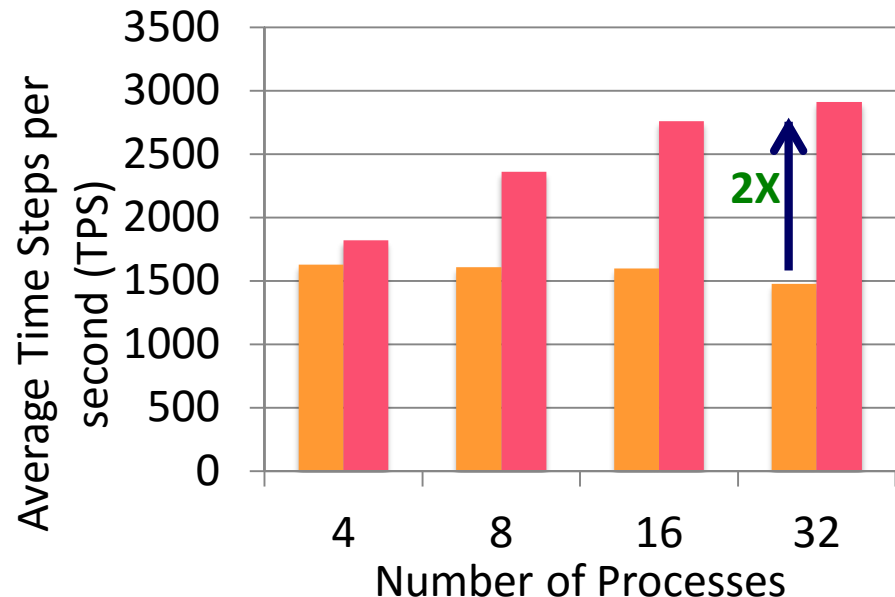
GPU-GPU Inter-node Bandwidth



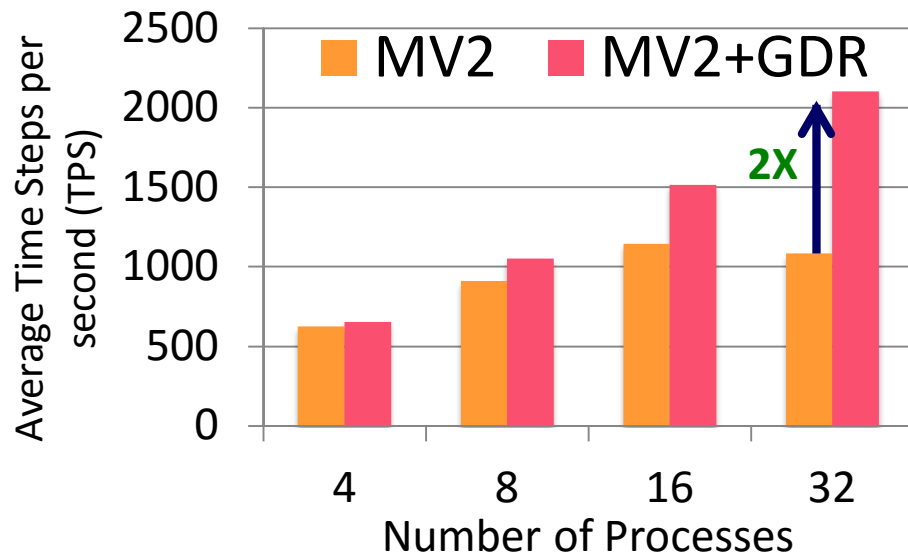
MVAPICH2-GDR-2.3
Intel Haswell (E5-2687W @ 3.10 GHz) node - 20 cores
NVIDIA Volta V100 GPU
Mellanox Connect-X4 EDR HCA
CUDA 9.0
Mellanox OFED 4.0 with GPU-Direct-RDMA

Application-Level Evaluation (HOOMD-blue)

64K Particles



256K Particles



- Platform: Wilkes (Intel Ivy Bridge + NVIDIA Tesla K20c + Mellanox Connect-IB)
- **HoomdBlue Version 1.0.5**
 - GDRCOPY enabled: MV2_USE_CUDA=1 MV2_IBA_HCA=mlx5_0 MV2_IBA_EAGER_THRESHOLD=32768 MV2_VBUF_TOTAL_SIZE=32768 MV2_USE_GPUDIRECT_LOOPBACK_LIMIT=32768 MV2_USE_GPUDIRECT_GDRCOPY=1 MV2_USE_GPUDIRECT_GDRCOPY_LIMIT=16384

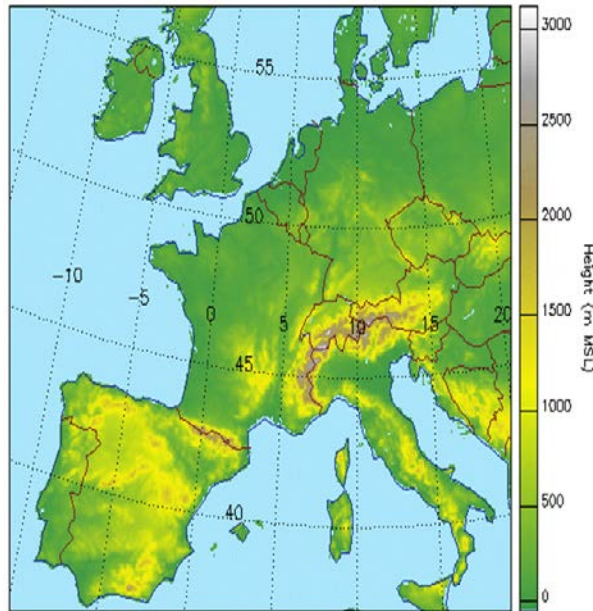
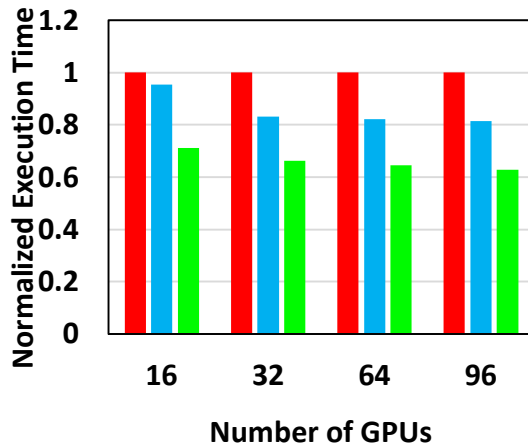
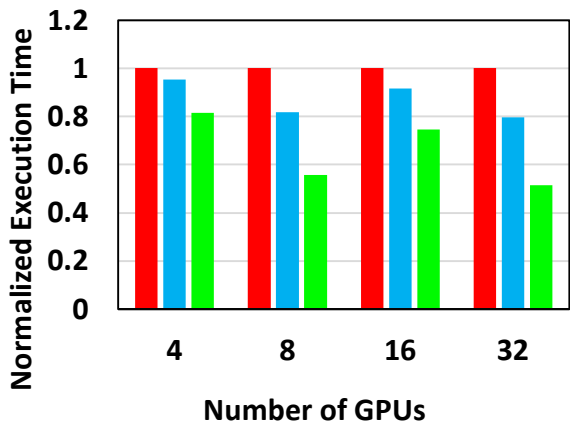
Application-Level Evaluation (Cosmo) and Weather Forecasting in Switzerland

Wilkes GPU Cluster

CSCS GPU cluster

■ Default ■ Callback-based ■ Event-based

■ Default ■ Callback-based ■ Event-based



- 2X improvement on 32 GPUs nodes
- 30% improvement on 96 GPU nodes (8 GPUs/node)

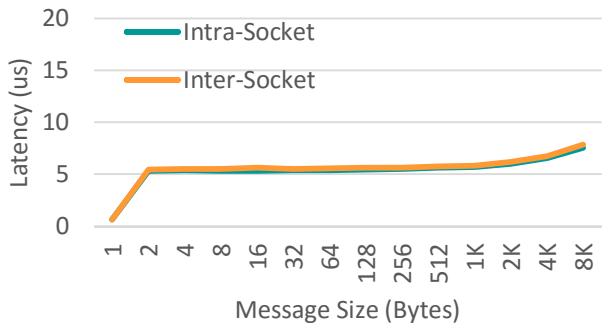
Cosmo model: <http://www2.cosmo-model.org/content/tasks/operational/meteoSwiss/>

On-going collaboration with CSCS and MeteoSwiss (Switzerland) in co-designing MV2-GDR and Cosmo Application

C. Chu, K. Hamidouche, A. Venkatesh, D. Banerjee, H. Subramoni, and D. K. Panda, Exploiting Maximal Overlap for Non-Contiguous Data Movement Processing on Modern GPU-enabled Systems, IPDPS'16

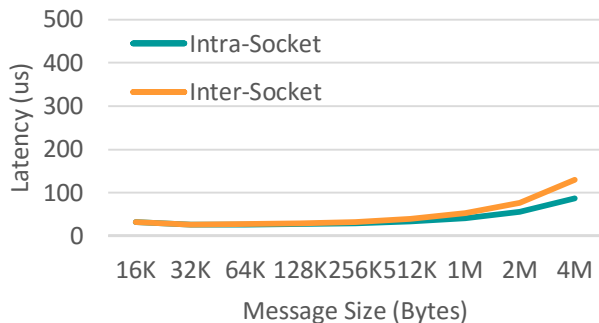
MVAPICH2-GDR: Performance on OpenPOWER (NVLink + Volta)

INTRA-NODE LATENCY (SMALL)

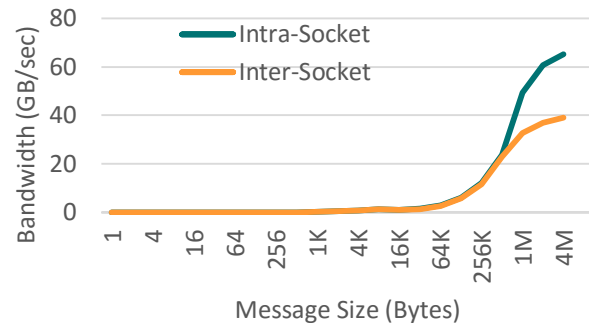


Intra-node Latency: 5.77 us (without GDRCopy)

INTRA-NODE LATENCY (LARGE)

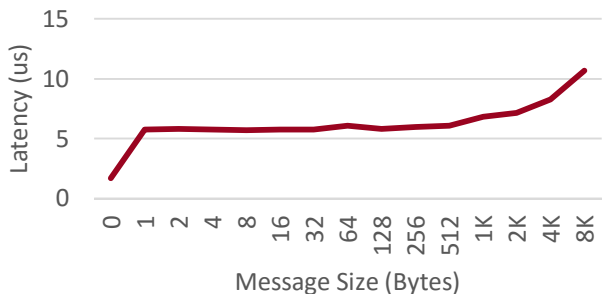


INTRA-NODE BANDWIDTH



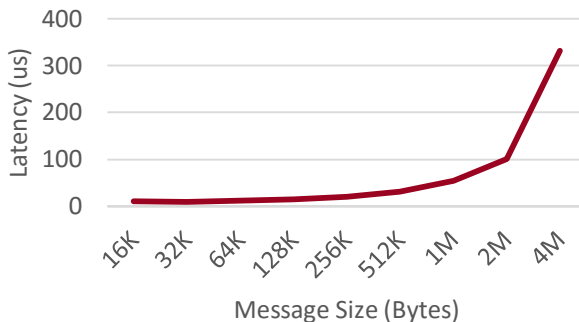
Intra-node Bandwidth: 63.7 GB/sec (NVLINK)

INTER-NODE LATENCY (SMALL)

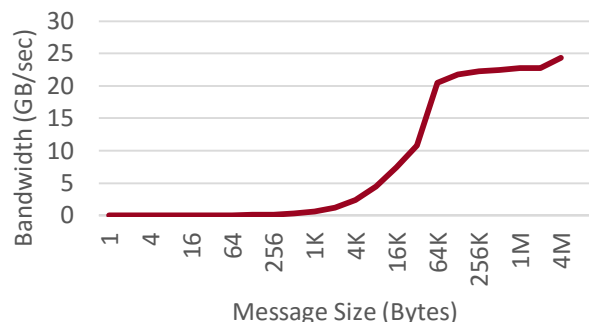


Inter-node Latency: 5.77 us (without GDRCopy)

INTER-NODE LATENCY (LARGE)



INTER-NODE BANDWIDTH

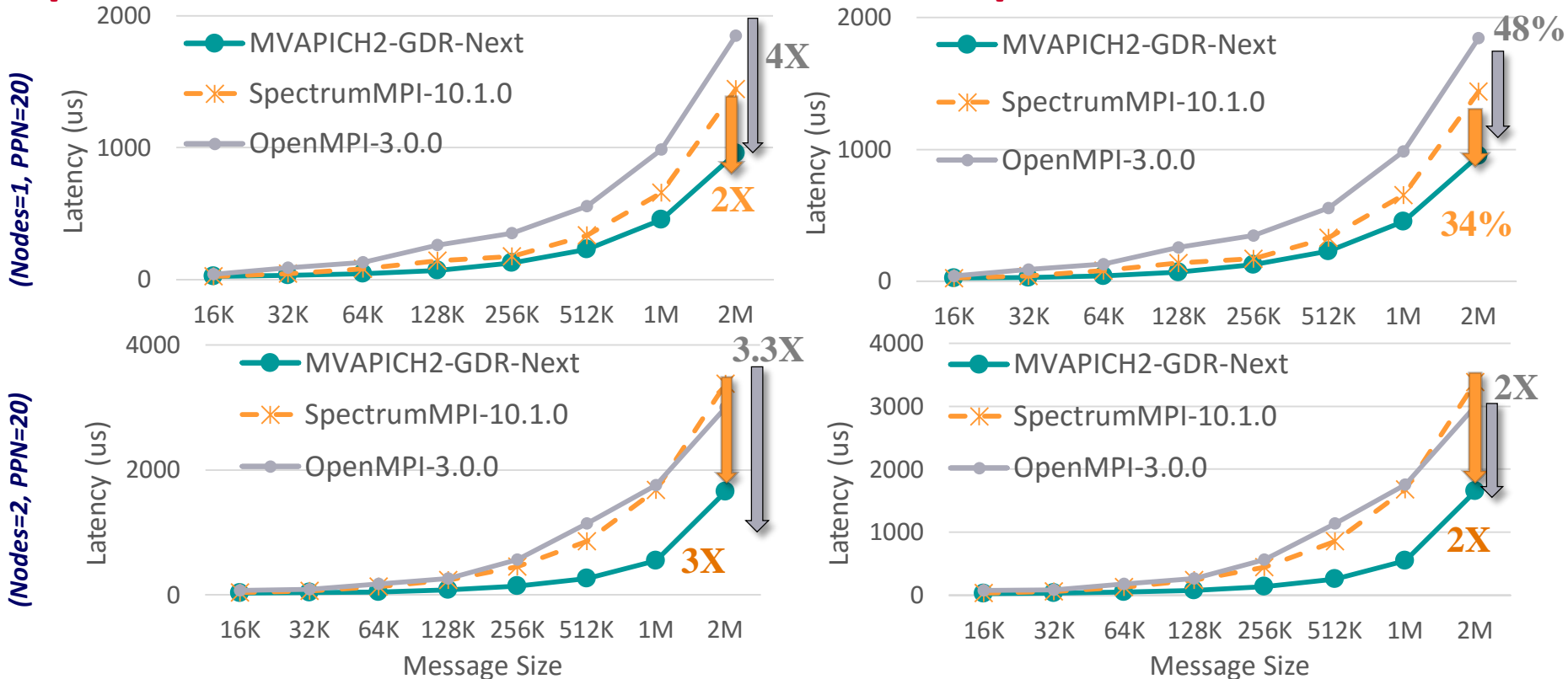


Inter-node Bandwidth: 23.7 GB/sec (2 port EDR)

Available since MVAPICH2-GDR 2.3a

Platform: OpenPOWER (ppc64le) nodes equipped with a dual-socket CPU, 4 Volta V100-SXM2 GPUs, and 2port EDR InfiniBand Interconnect

Optimized All-Reduce with XPMEM on OpenPOWER



- **Optimized MPI All-Reduce Design in MVAPICH2**

- **Up to 2X** performance improvement over Spectrum MPI and **4X** over OpenMPI for intra-node

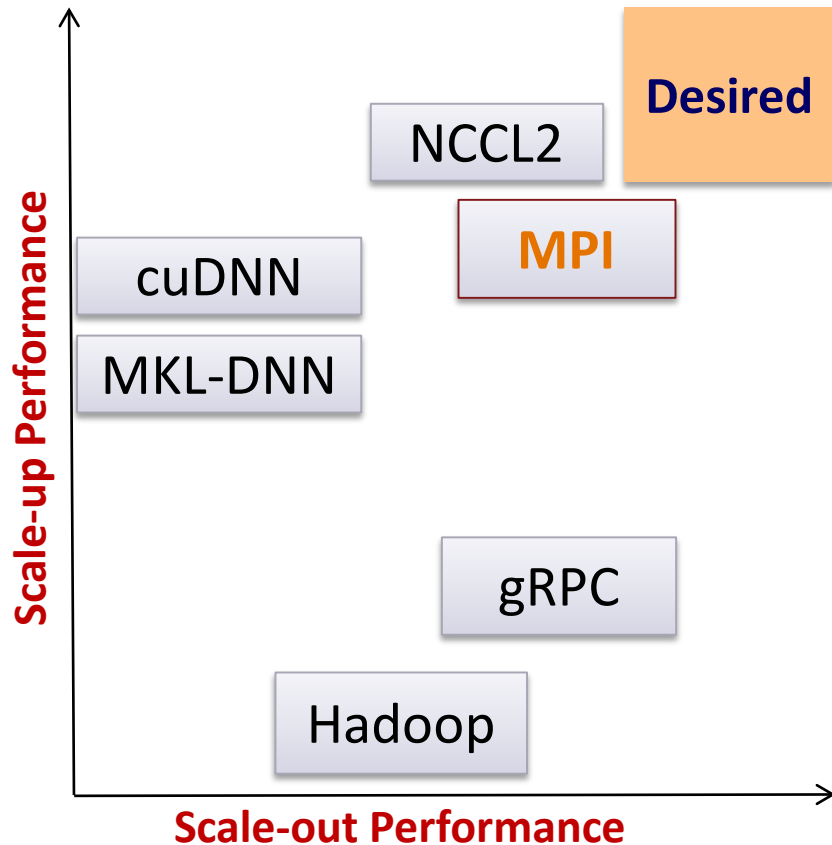
Optimized Runtime Parameters: MV2_CPU_BINDING_POLICY=hybrid MV2_HYBRID_BINDING_POLICY=bunch

Presentation Overview

- MVAPICH Project – MPI and PGAS (MVAPICH) Library with CUDA-Awareness
- **HiDL Project – High-Performance Deep Learning**
- HiBD Project – High-Performance Big Data Analytics Library
- Commercial Support from X-ScaleSolutions
- Conclusions and Q&A

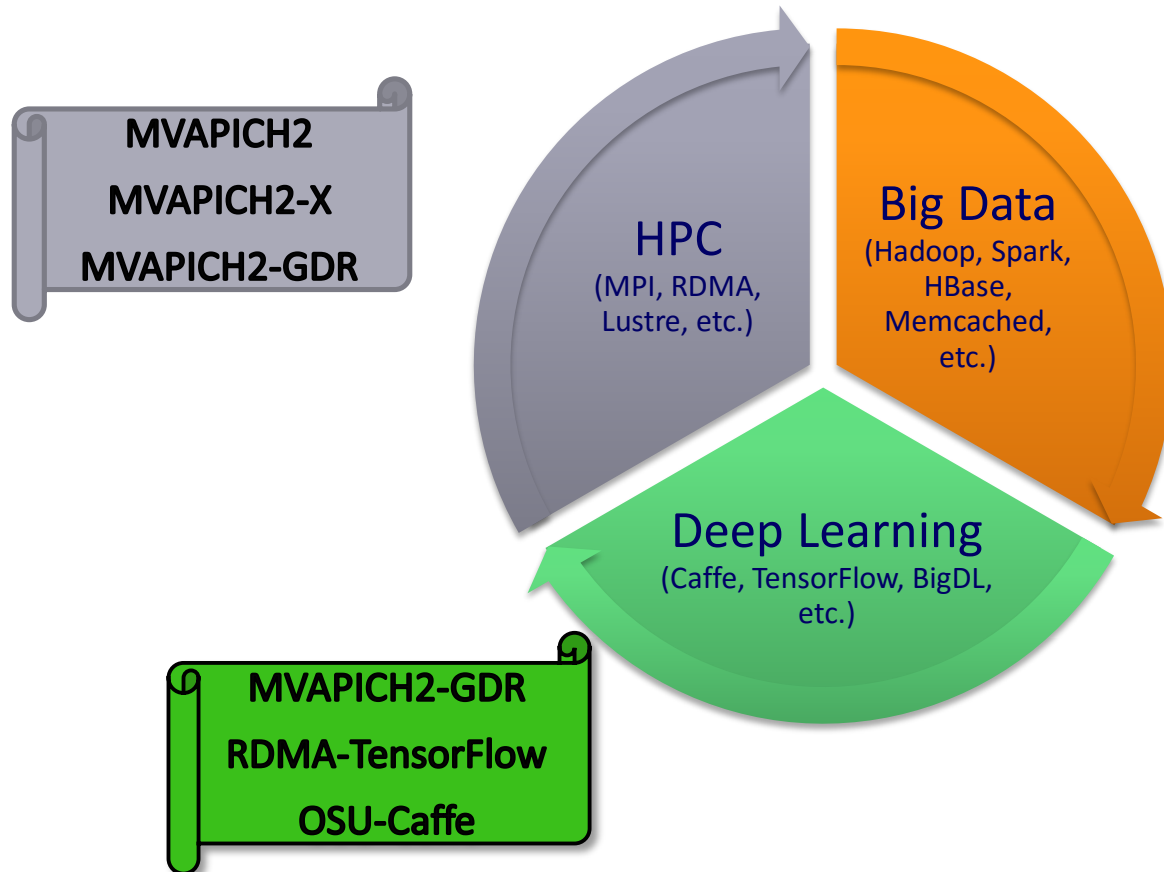
Deep Learning: New Challenges for MPI Runtimes

- Deep Learning frameworks are a different game altogether
 - Unusually large message sizes (order of megabytes)
 - Most communication based on GPU buffers
- Existing State-of-the-art
 - cuDNN, cuBLAS, NCCL --> **scale-up** performance
 - NCCL2, CUDA-Aware MPI --> **scale-out** performance
 - For small and medium message sizes only!
- Proposed: Can we **co-design** the MPI runtime (**MVAPICH2-GDR**) and the DL framework (**Caffe**) to achieve both?
 - Efficient **Overlap** of Computation and Communication
 - Efficient **Large-Message** Communication (Reductions)
 - What **application co-designs** are needed to exploit **communication-runtime co-designs**?



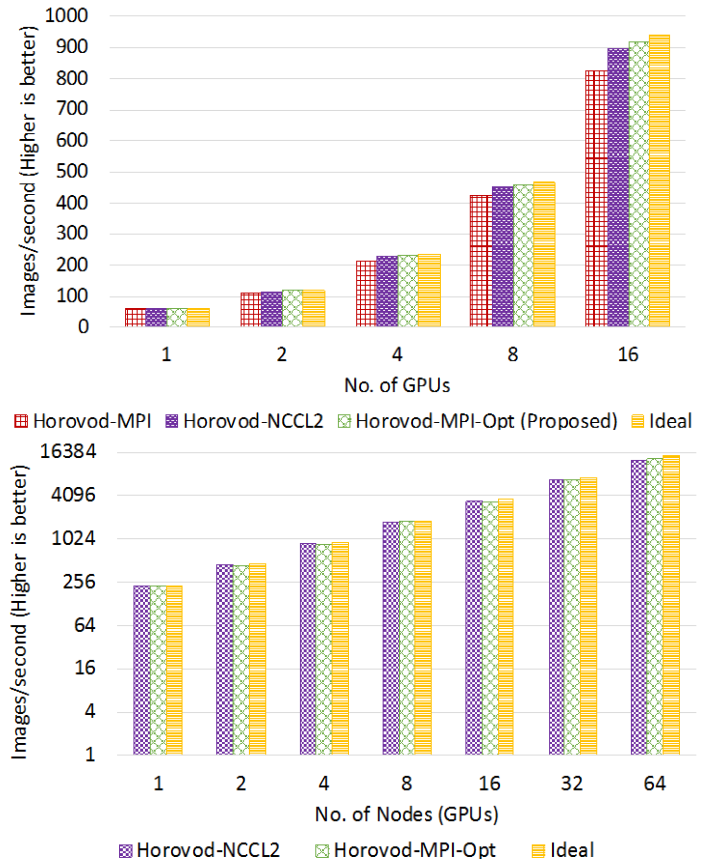
A. A. Awan, K. Hamidouche, J. M. Hashmi, and D. K. Panda, S-Caffe: Co-designing MPI Runtimes and Caffe for Scalable Deep Learning on Modern GPU Clusters. In *Proceedings of the 22nd ACM SIGPLAN Symposium on Principles and Practice of Parallel Programming (PPoPP '17)*

Convergent Software Stacks for HPC, Big Data and Deep Learning



Exploiting CUDA-Aware MPI for TensorFlow (Horovod)

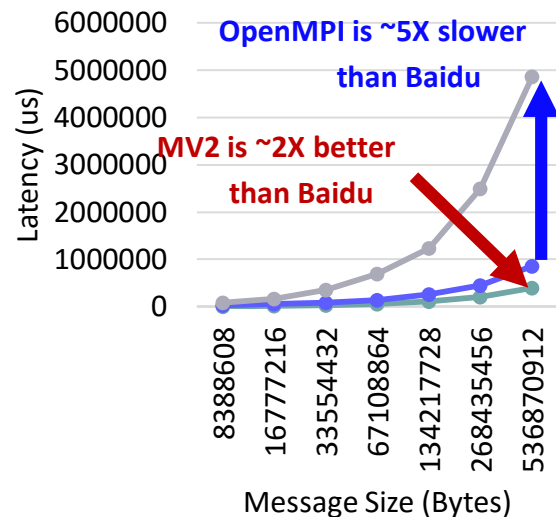
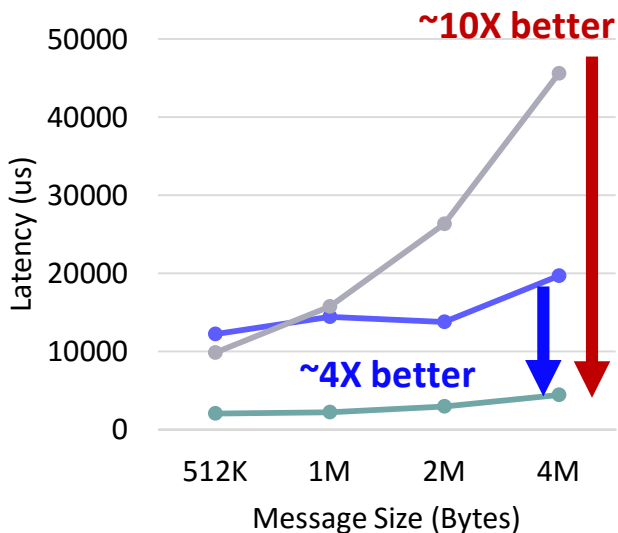
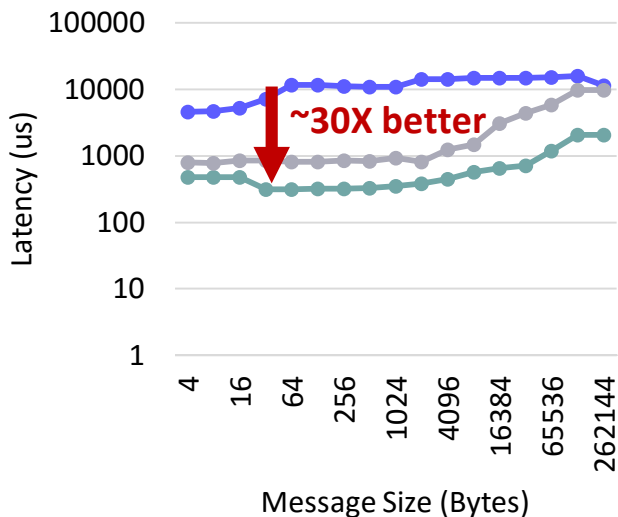
- Efficient Allreduce is crucial for Horovod's overall training performance
 - Both MPI and NCCL designs are available
- We have evaluated Horovod extensively and compared across a wide range of designs using gRPC and gRPC extensions
- MVAPICH2-GDR achieved up to **90%** scaling efficiency for ResNet-50 Training on 64 Pascal GPUs



Awan et al., "Scalable Distributed DNN Training using TensorFlow and CUDA-Aware MPI: Characterization, Designs, and Performance Evaluation", (To be presented) CCGrid '19.
<https://arxiv.org/abs/1810.11112>

MVAPICH2-GDR: Allreduce Comparison with Baidu and OpenMPI

- 16 GPUs (4 nodes) MVAPICH2-GDR vs. Baidu-Allreduce and OpenMPI 3.0



Legend: MVAPICH2 (green), BAIDU (blue), OPENMPI (grey)

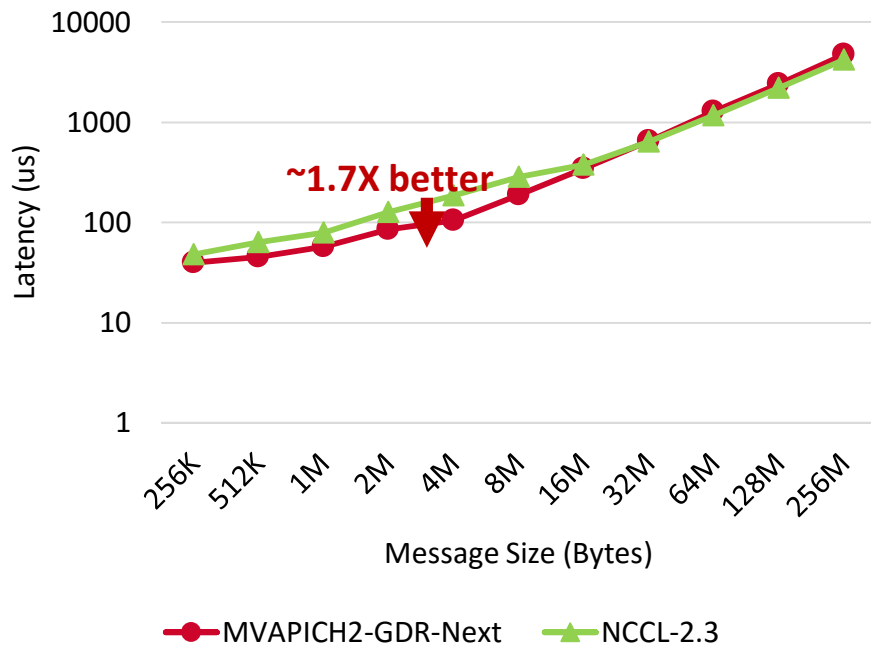
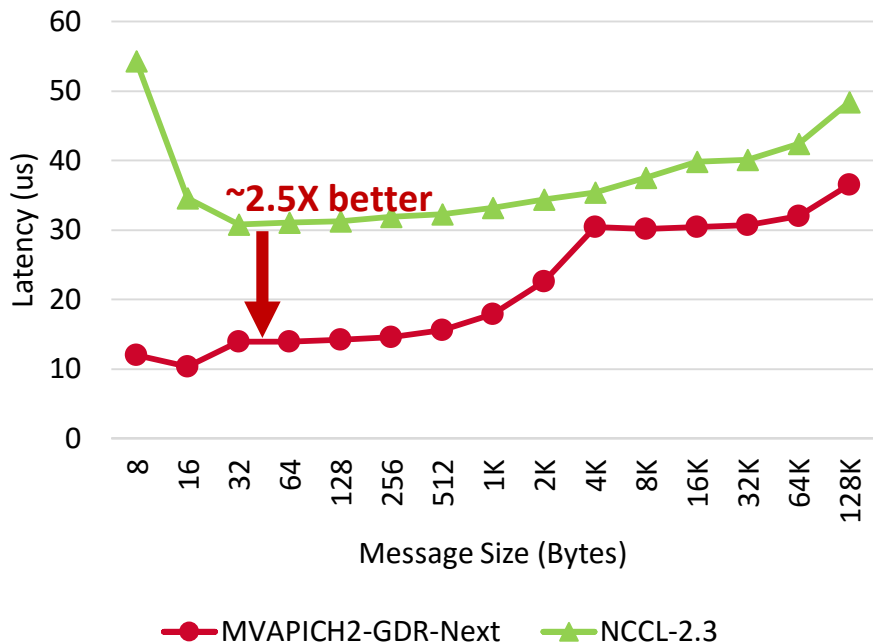
Legend: MVAPICH2 (green), BAIDU (blue), OPENMPI (grey)

Legend: MVAPICH2 (green), BAIDU (blue), OPENMPI (grey)

*Available since MVAPICH2-GDR 2.3a

MVAPICH2-GDR vs. NCCL2 – Allreduce Operation (DGX-2)

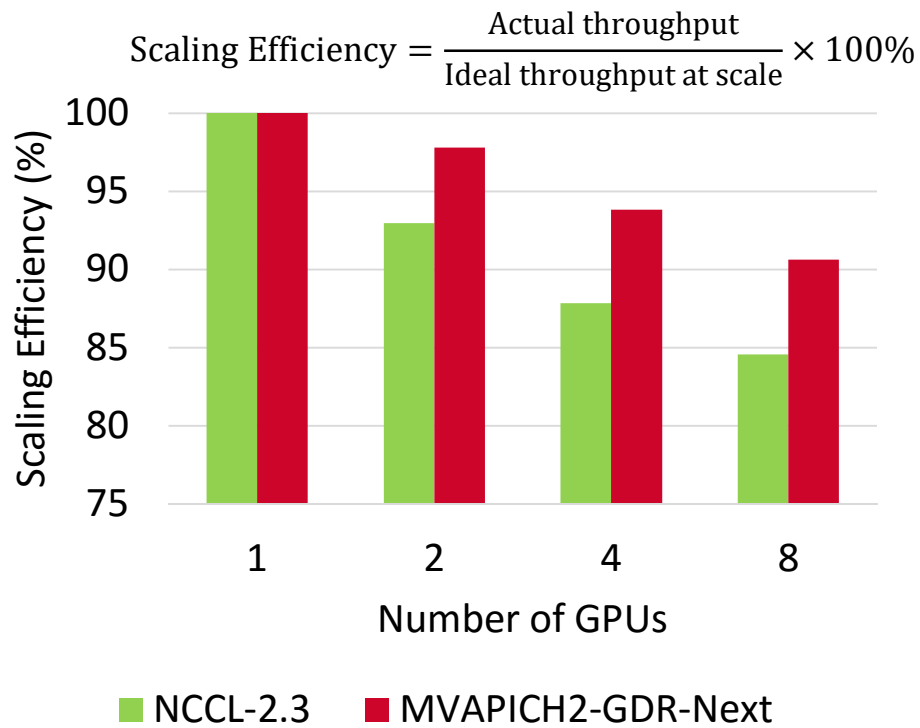
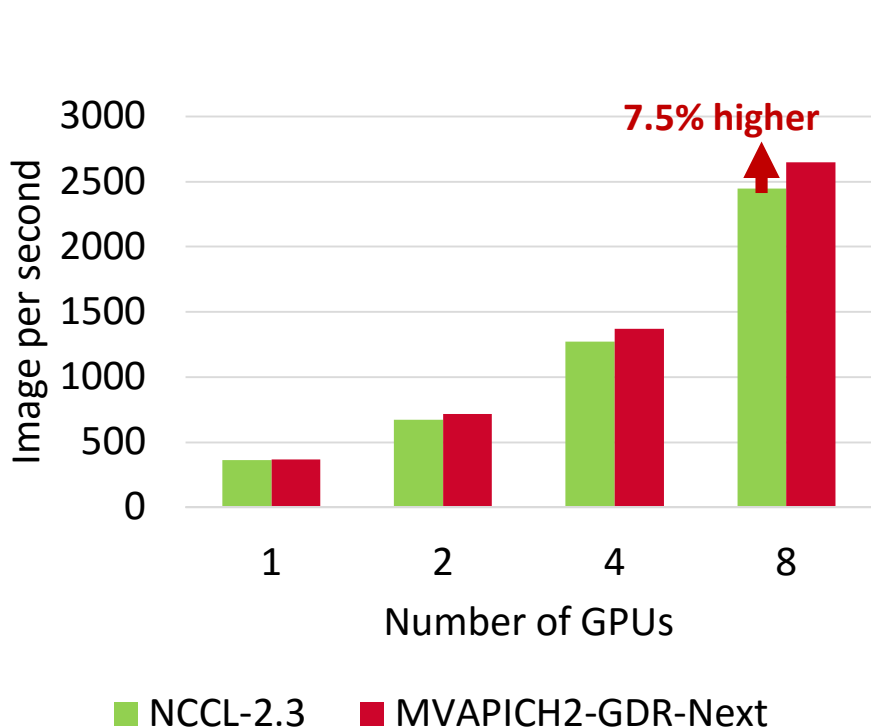
- Optimized designs in upcoming MVAPICH2-GDR offer better/comparable performance for most cases
- MPI_Allreduce (MVAPICH2-GDR) vs. ncclAllreduce (NCCL2) on 1 DGX-2 node (16 Volta GPUs)



Platform: Nvidia DGX-2 system (16 Nvidia Volta GPUs connected with NVSwitch), CUDA 9.2

MVAPICH2-GDR vs. NCCL2 – ResNet-50 Training

- ResNet-50 Training using TensorFlow benchmark on 1 DGX-2 node (8 Volta GPUs)



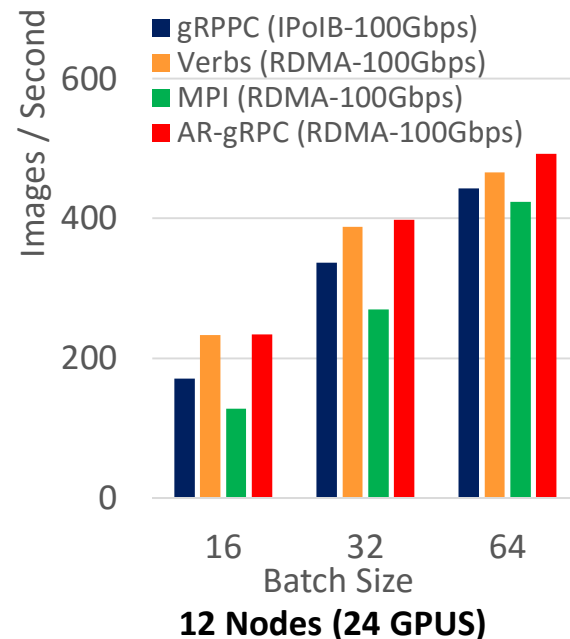
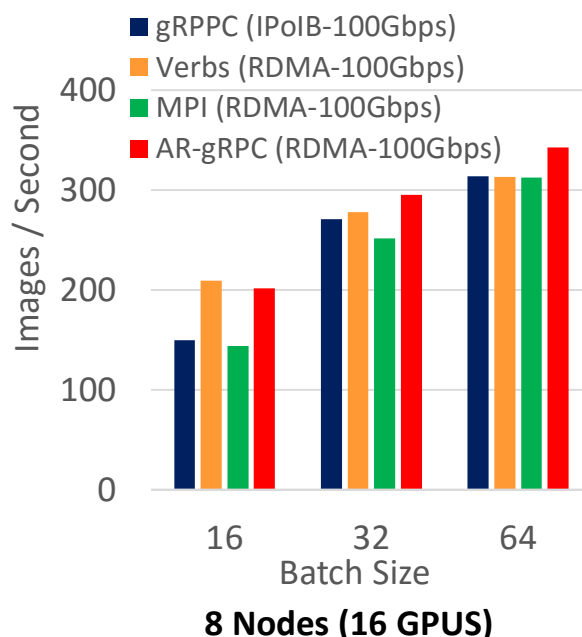
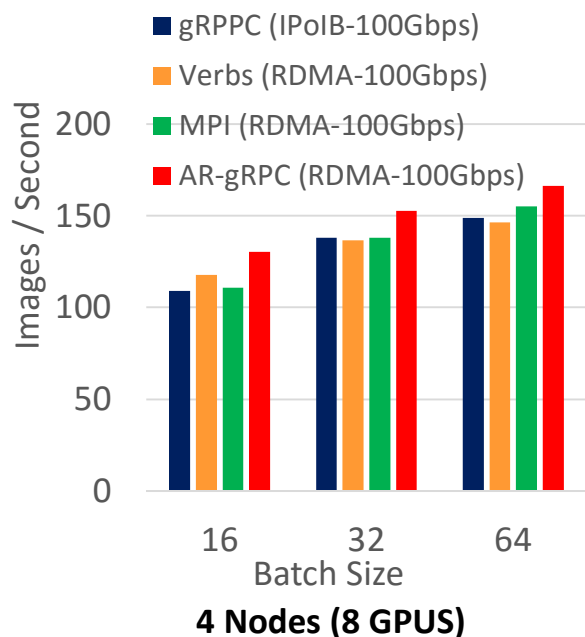
Platform: Nvidia DGX-2 system (16 Nvidia Volta GPUs connected with NVSwitch), CUDA 9.2

RDMA-TensorFlow Distribution

- High-Performance Design of TensorFlow over RDMA-enabled Interconnects
 - High performance RDMA-enhanced design with native InfiniBand support at the verbs-level for gRPC and TensorFlow
 - RDMA-based data communication
 - Adaptive communication protocols
 - Dynamic message chunking and accumulation
 - Support for RDMA device selection
 - Easily configurable for different protocols (native InfiniBand and IPoIB)
- Current release: **0.9.1**
 - Based on Google TensorFlow **1.3.0**
 - Tested with
 - Mellanox InfiniBand adapters (e.g., EDR)
 - NVIDIA GPGPU K80
 - Tested with CUDA 8.0 and CUDNN 5.0
 - <http://hidl.cse.ohio-state.edu>

OpenPOWER support
will be coming soon

Performance Benefit for RDMA-TensorFlow (Inception3)



- TensorFlow **Inception3** performance evaluation on an IB EDR cluster
 - Up to 20% performance speedup over Default gRPC (IPoIB) for 8 GPUs
 - Up to 34% performance speedup over Default gRPC (IPoIB) for 16 GPUs
 - Up to 37% performance speedup over Default gRPC (IPoIB) for 24 GPUs

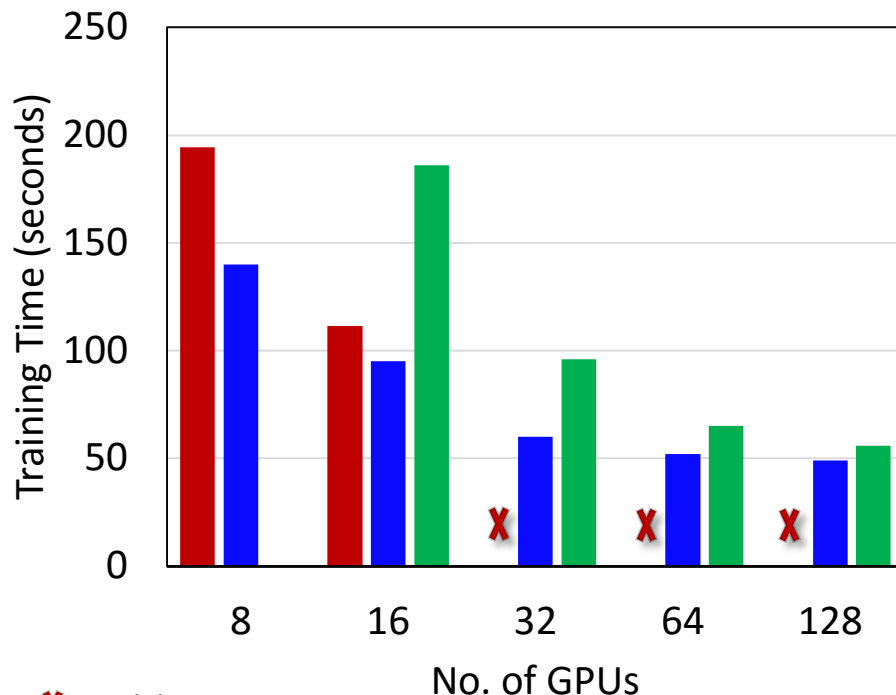
OSU-Caffe: Scalable Deep Learning

- Caffe : A flexible and layered Deep Learning framework.
- Benefits and Weaknesses
 - Multi-GPU Training within a single node
 - Performance degradation for GPUs across different sockets
 - Limited Scale-out
- OSU-Caffe: MPI-based Parallel Training
 - Enable Scale-up (within a node) and Scale-out (across multi-GPU nodes)
 - Scale-out on 64 GPUs for training CIFAR-10 network on CIFAR-10 dataset
 - Scale-out on 128 GPUs for training GoogLeNet network on ImageNet dataset

OSU-Caffe publicly available from

<http://hidl.cse.ohio-state.edu/>

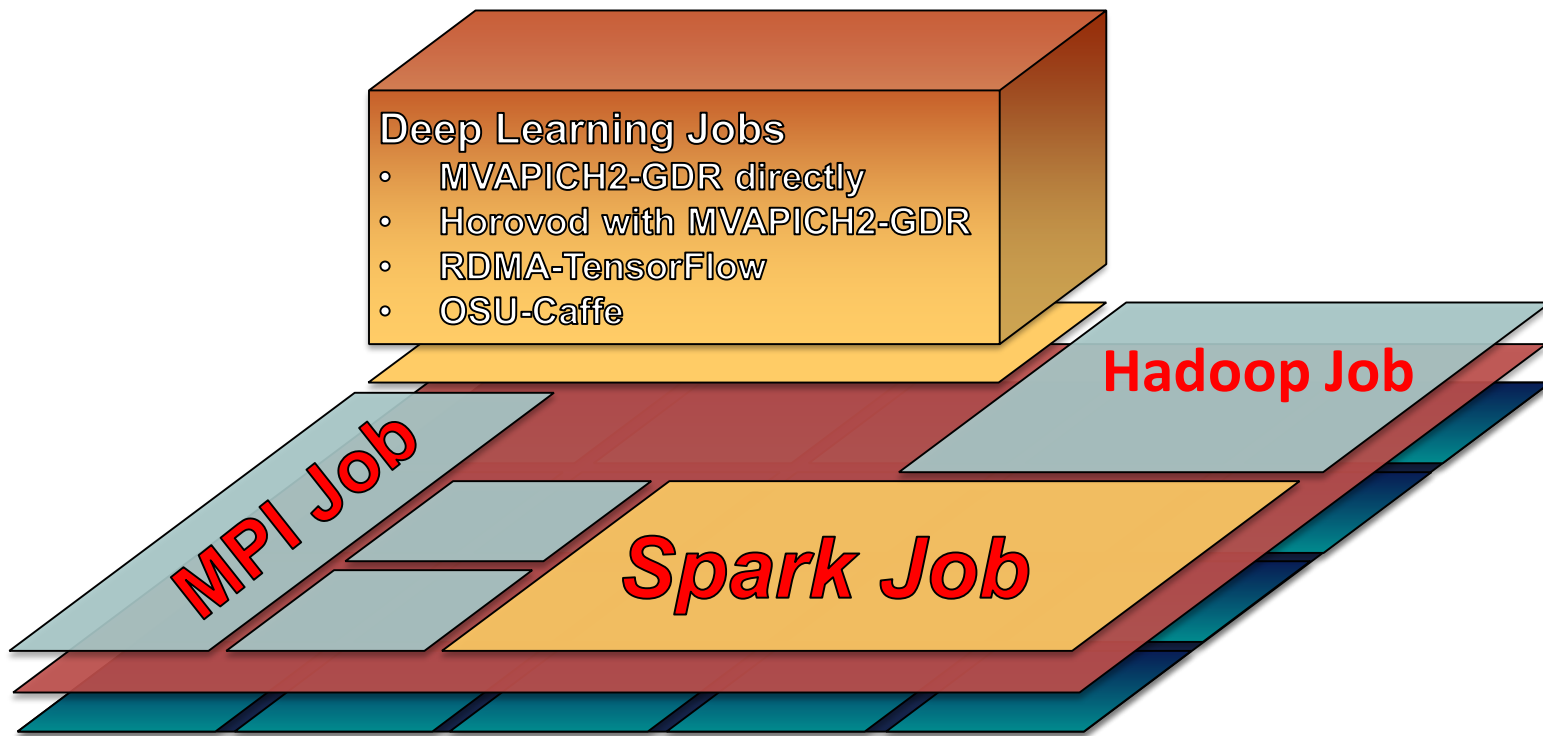
GoogLeNet (ImageNet) on 128 GPUs



X Invalid use case

■ Caffe ■ OSU-Caffe (1024) ■ OSU-Caffe (2048)

Using HiDL Packages for Deep Learning on Existing HPC Infrastructure

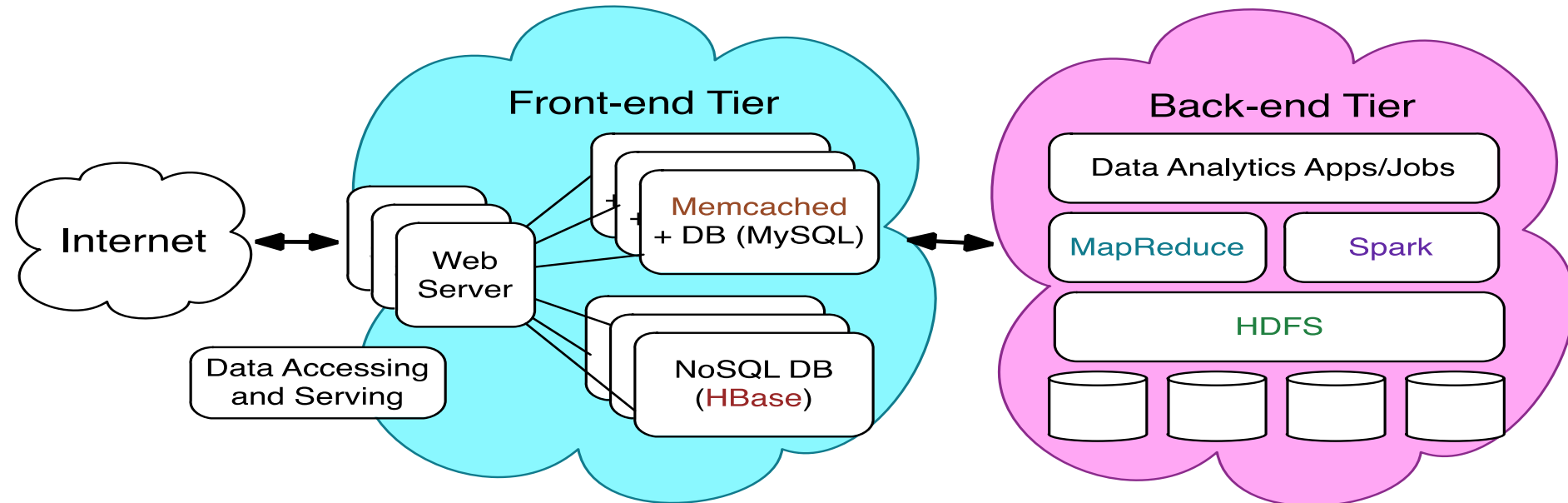


Presentation Overview

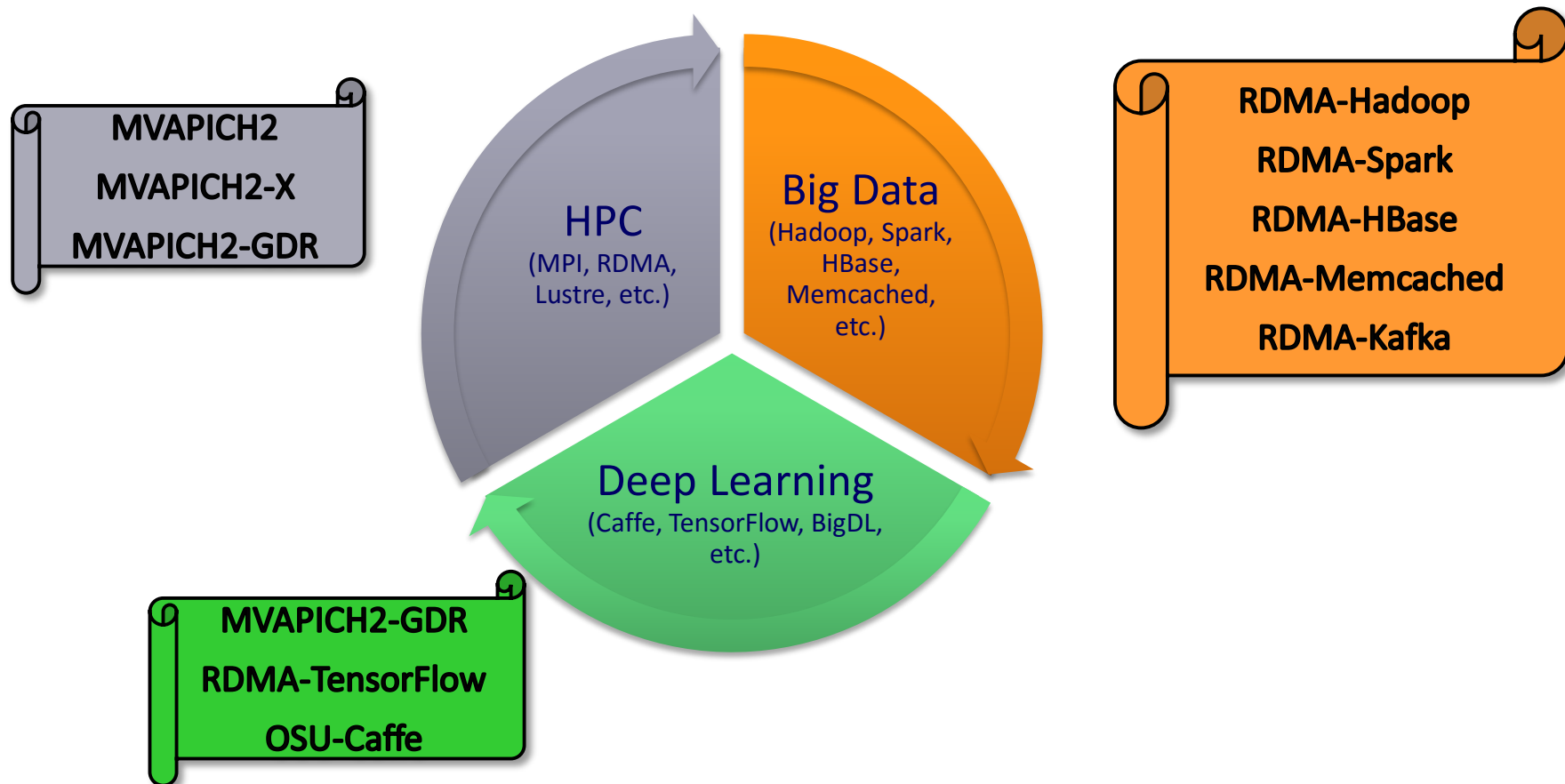
- MVAPICH Project – MPI and PGAS (MVAPICH) Library with CUDA-Awareness
- HiDL Project – High-Performance Deep Learning
- **HiBD Project – High-Performance Big Data Analytics Library**
- Commercial Support from X-ScaleSolutions
- Conclusions and Q&A

Data Management and Processing on Modern Datacenters

- Substantial impact on designing and utilizing data management and processing systems in multiple tiers
 - Front-end data accessing and serving (Online)
 - Memcached + DB (e.g. MySQL), HBase
 - Back-end data analytics (Offline)
 - HDFS, MapReduce, Spark



Convergent Software Stacks for HPC, Big Data and Deep Learning



The High-Performance Big Data (HiBD) Project

- RDMA for Apache Spark
- RDMA for Apache Hadoop 3.x (RDMA-Hadoop-3.x)
- RDMA for Apache Hadoop 2.x (RDMA-Hadoop-2.x)
 - Plugins for Apache, Hortonworks (HDP) and Cloudera (CDH) Hadoop distributions
- RDMA for Apache Kafka
- RDMA for Apache HBase
- RDMA for Memcached (RDMA-Memcached)
- RDMA for Apache Hadoop 1.x (RDMA-Hadoop)
- OSU HiBD-Benchmarks (OHB)
 - HDFS, Memcached, HBase, and Spark Micro-benchmarks
- <http://hibd.cse.ohio-state.edu>
- Users Base: 300 organizations from 35 countries
- More than 29,300 downloads from the project site

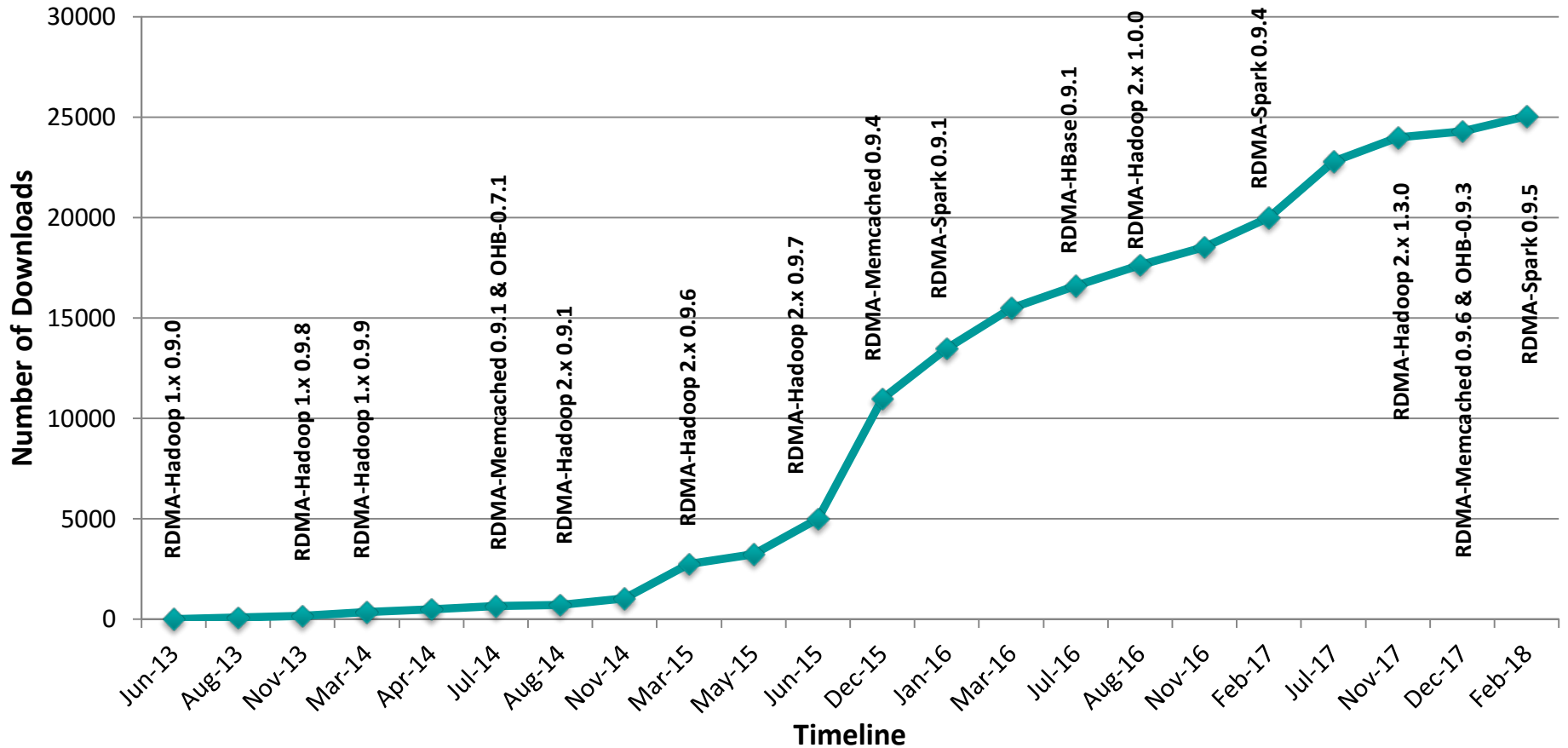
Available for InfiniBand and RoCE
Also run on Ethernet

Available for x86 and OpenPOWER

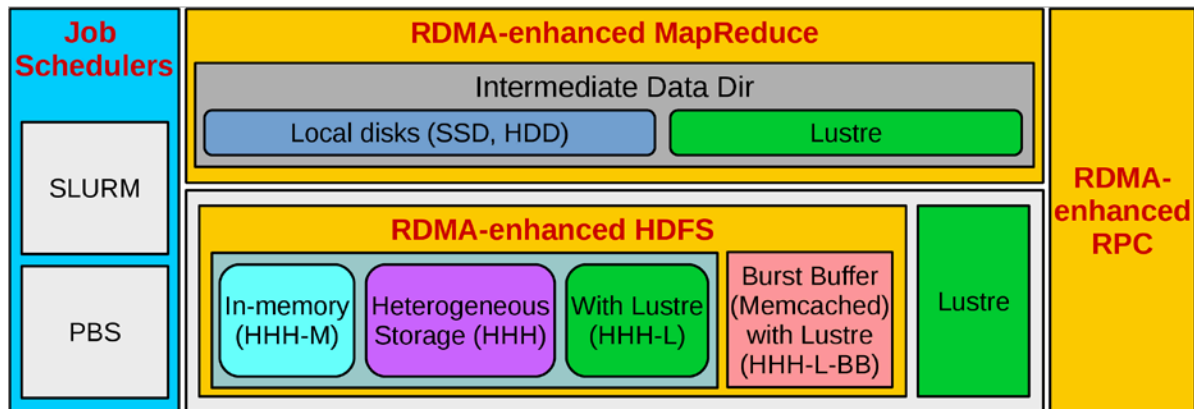
Support for Singularity and Docker



HiBD Release Timeline and Downloads

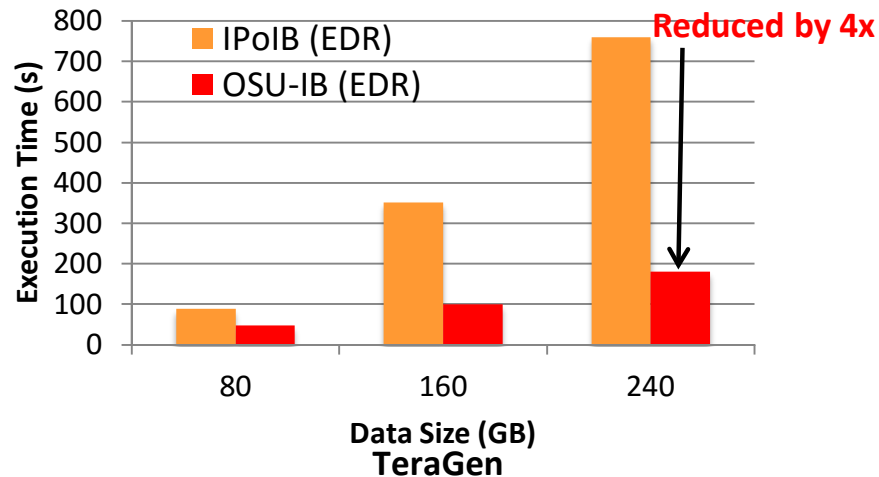
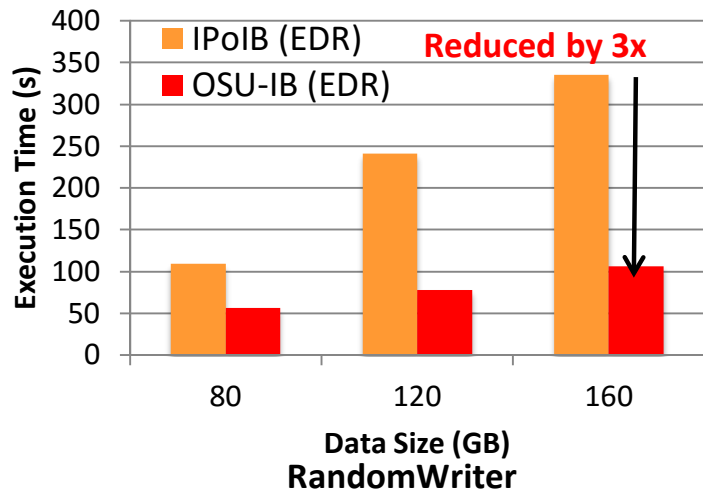


Different Modes of RDMA for Apache Hadoop 2.x



- **HHH:** Heterogeneous storage devices with hybrid replication schemes are supported in this mode of operation to have better fault-tolerance as well as performance. This mode is enabled by **default** in the package.
- **HHH-M:** A high-performance in-memory based setup has been introduced in this package that can be utilized to perform all I/O operations in-memory and obtain as much performance benefit as possible.
- **HHH-L:** With parallel file systems integrated, HHH-L mode can take advantage of the Lustre available in the cluster.
- **HHH-L-BB:** This mode deploys a Memcached-based burst buffer system to reduce the bandwidth bottleneck of shared file system access. The burst buffer design is hosted by Memcached servers, each of which has a local SSD.
- **MapReduce over Lustre, with/without local disks:** Besides, HDFS based solutions, this package also provides support to run MapReduce jobs on top of Lustre alone. Here, two different modes are introduced: with local disks and without local disks.
- **Running with Slurm and PBS:** Supports deploying RDMA for Apache Hadoop 2.x with Slurm and PBS in different running modes (HHH, HHH-M, HHH-L, and MapReduce over Lustre).

Performance Numbers of RDMA for Apache Hadoop 2.x – RandomWriter & TeraGen in OSU-RI2 (EDR)



Cluster with 8 Nodes with a total of 64 maps

- RandomWriter

- **3x** improvement over IPoIB for 80-160 GB file size

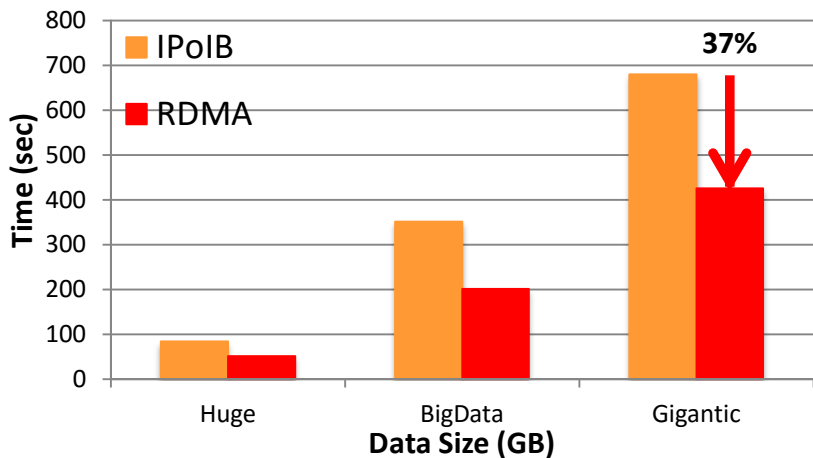
- TeraGen

- **4x** improvement over IPoIB for 80-240 GB file size

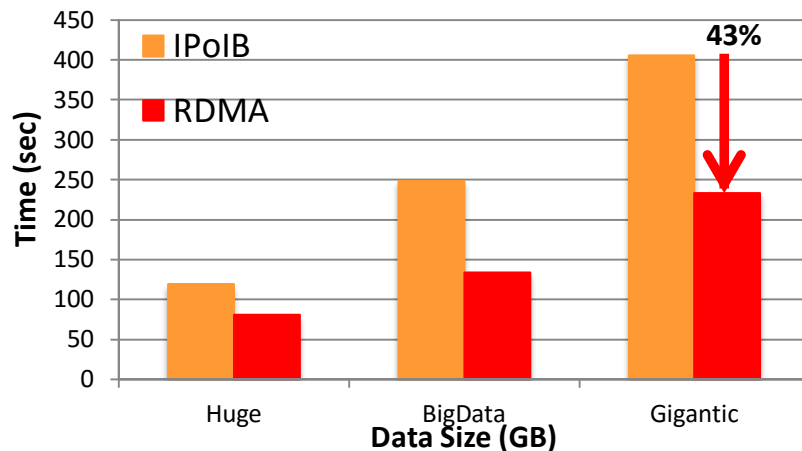
Using HiBD Packages for Big Data Processing on Existing HPC Infrastructure



Performance Evaluation of RDMA-Spark on SDSC Comet – HiBench PageRank



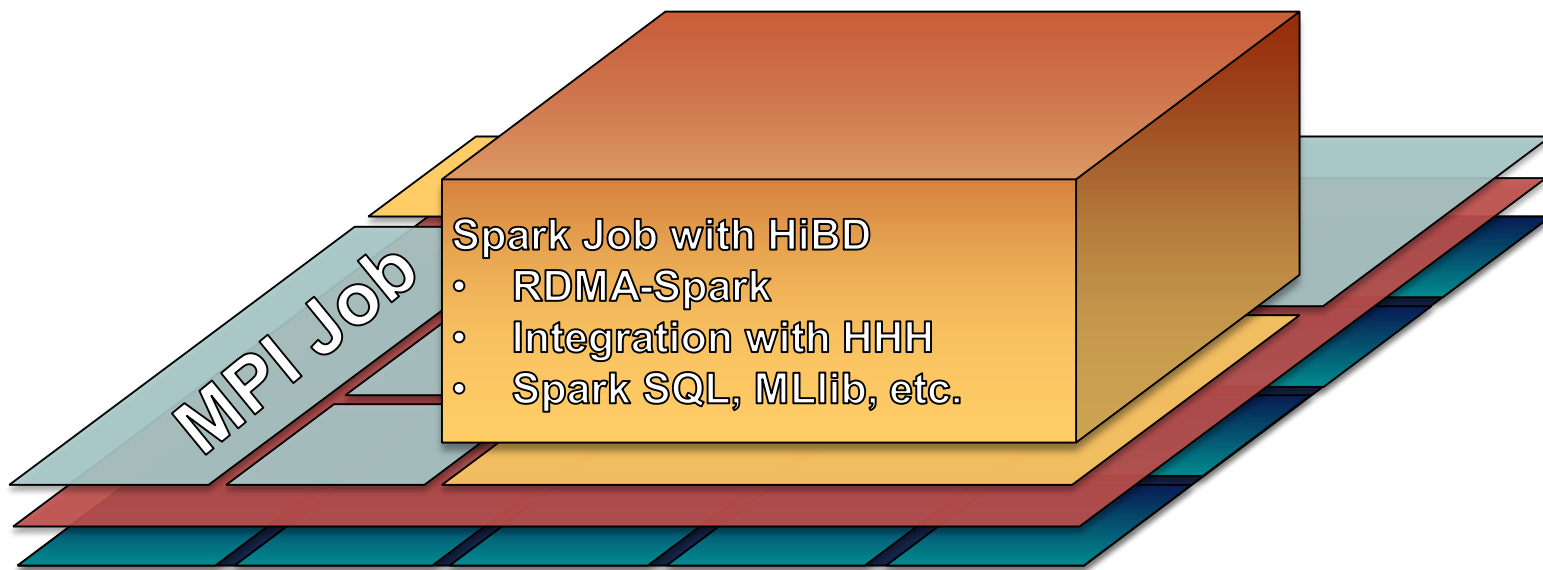
32 Worker Nodes, 768 cores, PageRank Total Time



64 Worker Nodes, 1536 cores, PageRank Total Time

- InfiniBand FDR, SSD, 32/64 Worker Nodes, 768/1536 Cores, (768/1536M 768/1536R)
- RDMA vs. IPoIB with 768/1536 concurrent tasks, single SSD per node.
 - 32 nodes/768 cores: Total time reduced by 37% over IPoIB (56Gbps)
 - 64 nodes/1536 cores: Total time reduced by 43% over IPoIB (56Gbps)

Using HiBD Packages for Big Data Processing on Existing HPC Infrastructure



Presentation Overview

- MVAPICH Project – MPI and PGAS (MVAPICH) Library with CUDA-Awareness
- HiDL Project – High-Performance Deep Learning
- HiBD Project – High-Performance Big Data Analytics Library
- **Commercial Support from X-ScaleSolutions**
- Conclusions and Q&A

Commercial Support for MVAPICH2, HiBD, and HiDL Libraries

- Supported through X-ScaleSolutions (<http://x-scalesolutions.com>)
- Benefits:
 - Help and guidance with installation of the library
 - Platform-specific optimizations and tuning
 - Timely support for operational issues encountered with the library
 - Web portal interface to submit issues and tracking their progress
 - Advanced debugging techniques
 - Application-specific optimizations and tuning
 - Obtaining guidelines on best practices
 - Periodic information on major fixes and updates
 - Information on major releases
 - Help with upgrading to the latest release
 - Flexible Service Level Agreements
- Support provided to Lawrence Livermore National Laboratory (LLNL) for the last two years



Presentation Overview

- MVAPICH Project – MPI and PGAS (MVAPICH) Library with CUDA-Awareness
- HiDL Project – High-Performance Deep Learning
- HiBD Project – High-Performance Big Data Analytics Library
- Commercial Support from X-ScaleSolutions
- **Conclusions and Q&A**

Concluding Remarks

- Upcoming Exascale systems need to be designed with a holistic view of HPC, Big Data, Deep Learning, and Cloud
- Presented an overview of designing convergent software stacks for HPC, Big Data, and Deep Learning
- Presented solutions enable HPC, Big Data, and Deep Learning communities to take advantage of current and next-generation systems

Funding Acknowledgments

Funding Support by



Equipment Support by



Personnel Acknowledgments

Current Students (Graduate)

- A. Awan (Ph.D.)
- M. Bayatpour (Ph.D.)
- S. Chakraborty (Ph.D.)
- C.-H. Chu (Ph.D.)
- S. Guganani (Ph.D.)

Current Students (Undergraduate)

- J. Hashmi (Ph.D.)
- A. Jain (Ph.D.)
- K. S. Khorassani (Ph.D.)
- P. Kousha (Ph.D.)
- D. Shankar (Ph.D.)
- V. Gangal (B.S.)
- M. Haupt (B.S.)
- N. Sarkauskas (B.S.)
- A. Yeretizian (B.S.)

Current Research Asst. Professor

- X. Lu

Current Post-doc

- A. Ruhela
- K. Manian

Current Research Scientist

- H. Subramoni

Current Research Specialist

- J. Smith

Past Students

- A. Augustine (M.S.)
- P. Balaji (Ph.D.)
- R. Biswas (M.S.)
- S. Bhagvat (M.S.)
- A. Bhat (M.S.)
- D. Buntinas (Ph.D.)
- L. Chai (Ph.D.)
- B. Chandrasekharan (M.S.)
- N. Dandapanthula (M.S.)
- V. Dhanraj (M.S.)
- T. Gangadharappa (M.S.)
- K. Gopalakrishnan (M.S.)
- W. Huang (Ph.D.)
- W. Jiang (M.S.)
- J. Jose (Ph.D.)
- S. Kini (M.S.)
- M. Koop (Ph.D.)
- K. Kulkarni (M.S.)
- R. Kumar (M.S.)
- S. Krishnamoorthy (M.S.)
- K. Kandalla (Ph.D.)
- M. Li (Ph.D.)
- P. Lai (M.S.)
- J. Liu (Ph.D.)
- M. Luo (Ph.D.)
- A. Mamidala (Ph.D.)
- G. Marsh (M.S.)
- V. Meshram (M.S.)
- A. Moody (M.S.)
- S. Naravula (Ph.D.)
- R. Noronha (Ph.D.)
- X. Ouyang (Ph.D.)
- S. Pai (M.S.)
- S. Potluri (Ph.D.)

Past Post-Docs

- D. Banerjee
- X. Besseron
- H.-W. Jin
- J. Lin
- M. Luo
- E. Mancini
- S. Marcarelli
- J. Vienne
- H. Wang

- R. Rajachandrasekar (Ph.D.)
- G. Santhanaraman (Ph.D.)
- A. Singh (Ph.D.)
- J. Sridhar (M.S.)
- S. Sur (Ph.D.)
- H. Subramoni (Ph.D.)
- K. Vaidyanathan (Ph.D.)
- A. Vishnu (Ph.D.)
- J. Wu (Ph.D.)
- W. Yu (Ph.D.)
- J. Zhang (Ph.D.)

Past Research Scientist

- K. Hamidouche
- S. Sur

Past Programmers

- D. Bureddy
- J. Perkins

Past Research Specialist

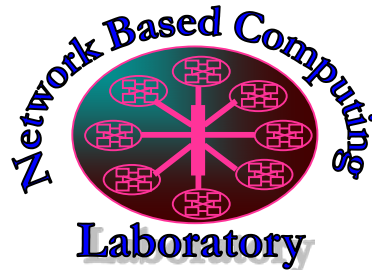
- M. Arnold

Multiple Positions Available in My Group

- Looking for Bright and Enthusiastic Personnel to join as
 - PhD Students
 - Post-Doctoral Researchers
 - MPI Programmer/Software Engineer
 - Hadoop/Big Data Programmer/Software Engineer
 - Deep Learning and Cloud Programmer/Software Engineer
- If interested, please send an e-mail to panda@cse.ohio-state.edu

Thank You!

panda@cse.ohio-state.edu



Network-Based Computing Laboratory

<http://nowlab.cse.ohio-state.edu/>



The High-Performance MPI/PGAS Project
<http://mvapich.cse.ohio-state.edu/>



High-Performance
Big Data

The High-Performance Big Data Project
<http://hibd.cse.ohio-state.edu/>



The High-Performance Deep Learning Project
<http://hidl.cse.ohio-state.edu/>