

Achieving High Performance on Microsoft Azure HPC Cloud using **MVAPICH2**

Microsoft Azure Booth Talk at SC '19

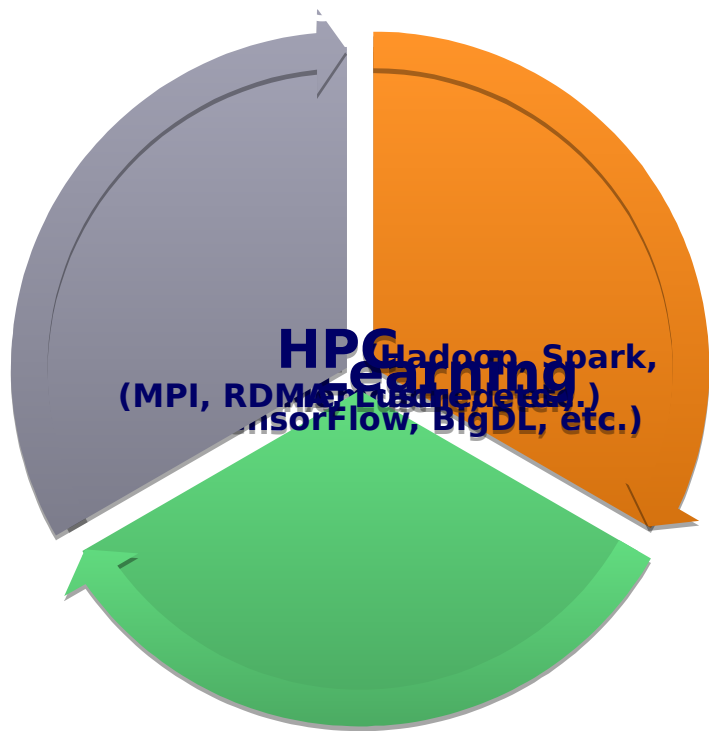
by

Dhabaleswar K. (DK) Panda

The Ohio State University

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

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



**Convergence of HPC,
Big Data, and Deep
Learning!**

**Increasing Need to
Run these
applications on the
Cloud!!**

- **Overview of the MVAPICH2 Project**
- Overview of Azure HPC Environments
- Designing MVAPICH2 for Azure
- Performance Results
- Future Plans

- 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 (Supercomputing '02)
 - MVAPICH2-X (MPI + PGAS), Available since 2011
 - Support for GPGPUs (MVAPICH2-GDR) and MIC (MVAPICH2-MIC), Available since 2011
 - 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) s
 - **Used by more than 3,050 organizations in 89 countries**
 - **More than 615,000 (> 0.6 million) downloads from the OSU site directly**
 - Empowering many TOP500 clusters (Nov '19 ranking)
 - 3rd, 10,649,600-core (Sunway TaihuLight) at National Supercomputing Center in Wuxi, China
 - 5th, 448, 448 cores (Frontera) at TACC
 - 8th, 391,680 cores (ABCI) in Japan
 - 14th, 570,020 cores (Nurion) in South Korea and many others
 - Available with software stacks of many vendors and Linux Distros (RedHat, SuSE, and OpenHPC)



Partner in the #5th TACC Frontera System

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

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 Technologies
(InfiniBand, iWARP, RoCE, Omni-Path, Elastic Fabric)

Transport Protocols

RC

SRD

UD

DC

Modern Features

UMR

ODP

SR-IOV

Multi-Rail

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

Transport Mechanisms

Shared Memory

CMA

IVSHMEM

XPME

Modern Features

Optane*

NVLink

CAP I*

* Upcoming

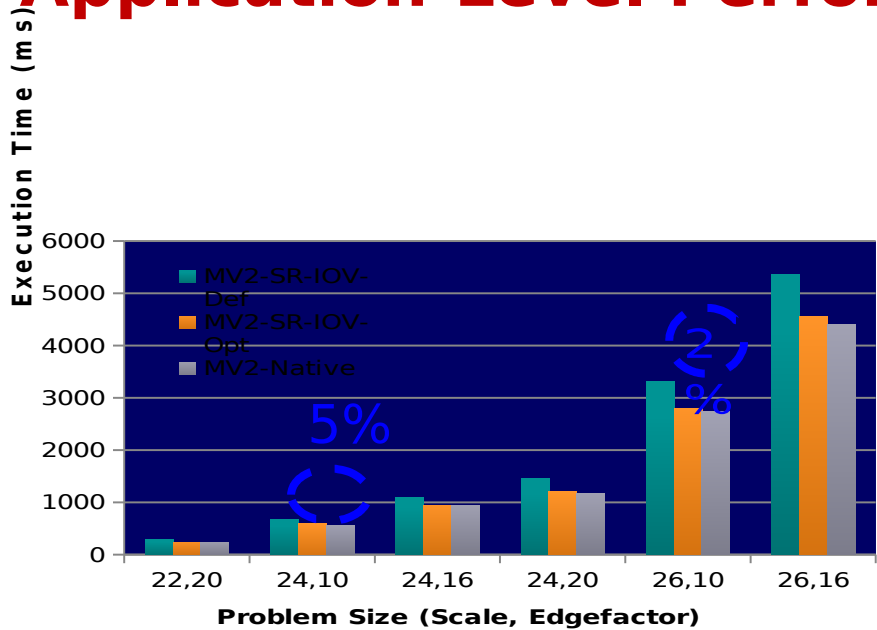
- Virtualization has many benefits
 - Fault-tolerance
 - Job migration
 - Compaction
- Have not been very popular in HPC due to overhead associated with Virtualization
- New SR-IOV (Single Root - IO Virtualization) support available with Mellanox InfiniBand adapters changes the field
- Enhanced MVAPICH2 support for SR-IOV
- MVAPICH2-Virt 2.2 supports:

J. Zhang, X. Lu, J. Jose, R. Shi and D. K. Panda, Can Inter-VM Shmem Benefit MPI Applications on SR-IOV based Virtualized InfiniBand Clusters? EuroPar'14

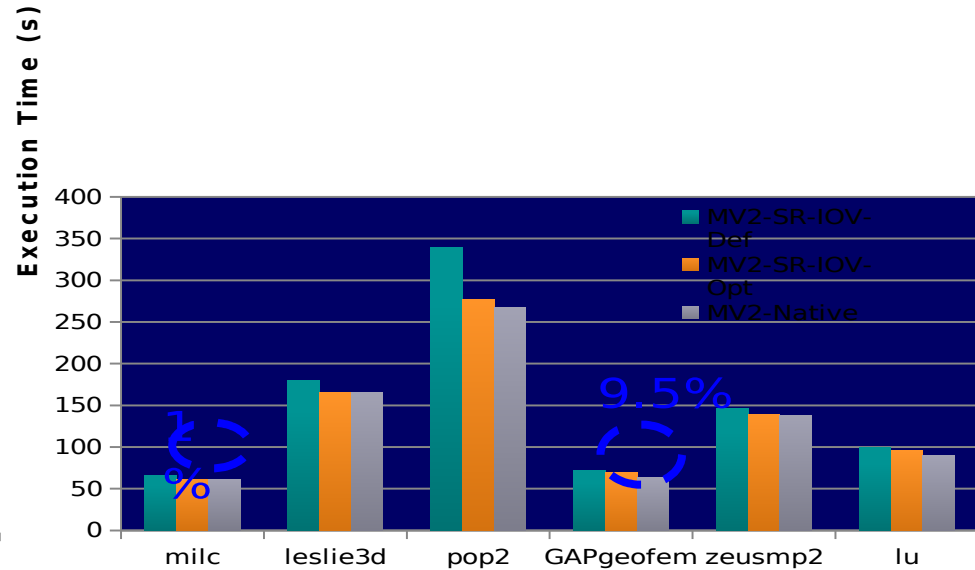
J. Zhang, X. Lu, J. Jose, M. Li, R. Shi and D.K. Panda, High Performance MPI Library over SR-IOV enabled InfiniBand Clusters, HiPC'14

J. Zhang, X. Lu, M. Arnold and D. K. Panda, MVAPICH2 Over OpenStack with SR-IOV: an Efficient Approach to build HPC Clouds, CCGrid'15

Application-Level Performance on Chameleon



Graph500



SPEC

MPI2007

- 32 VMs, 6 Core/VM
- Compared to Native, 2-5% overhead for Graph500 with 128 Procs
- Compared to Native, 1-9.5% overhead for SPEC MPI2007 with 128

- Overview of the MVAPICH2 Project
- **Overview of Azure HPC Environments**
- Designing MVAPICH2 for Azure
- Performance Results
- Future Plans

- Has been using RDMA-enabled network and software stacks for the last several years
- Moved to native InfiniBand support with Mellanox OFED for new instances (HB, HC, and upcoming ones)
- Uses SR-IOV support for virtualization
- Uses one VM per node

- Overview of the MVAPICH2 Project
- Overview of Azure HPC Environments
- **Designing MVAPICH2 for Azure**
- Performance Results
- Future Plans

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

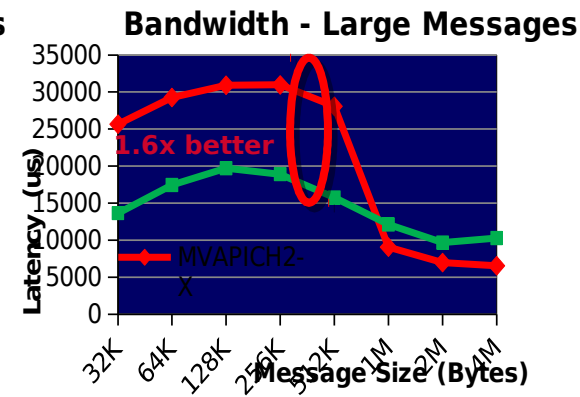
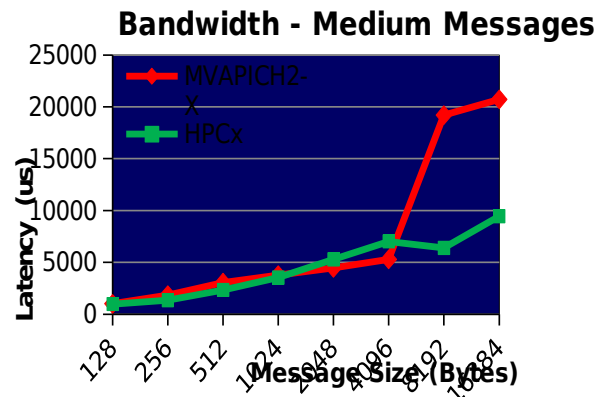
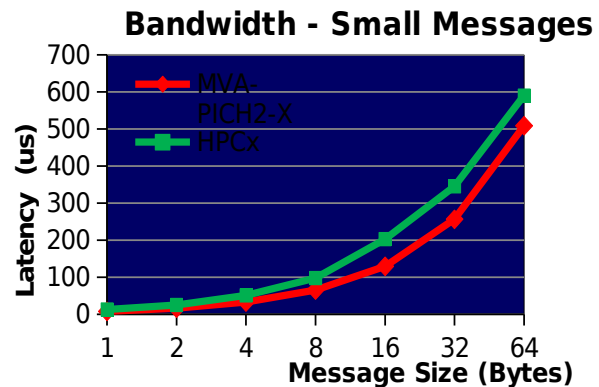
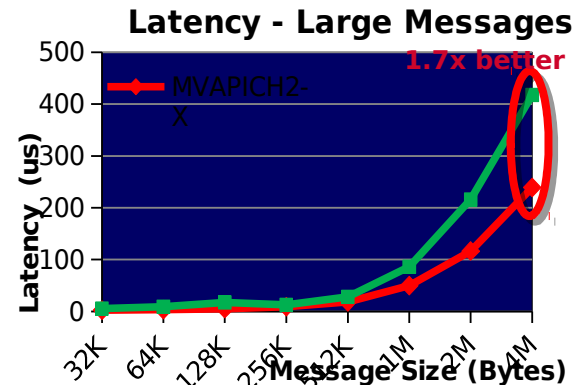
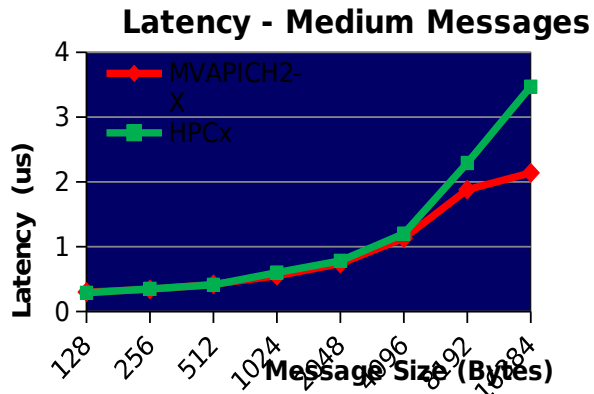
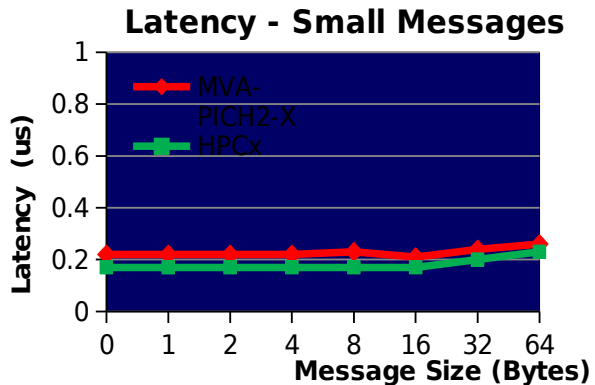
- **Released on 08/16/2019**
- Major Features and Enhancements
 - **Based on MVAPICH2-2.3.2**
 - **Enhanced tuning for point-to-point and collective operations**
 - **Targeted for Azure HB & HC virtual machines**
 - **Flexibility for 'one-click' deployment**
 - **Tested with Azure HB & HC VM instances**
- Available for download from <http://mvapich.cse.ohio-state.edu/downloads/>
- Detailed User Guide: <http://mvapich.cse.ohio-state.edu/userguide/mv2-azure/>
- On-going Work
 - Optimizing MVAPICH2-X features on Azure

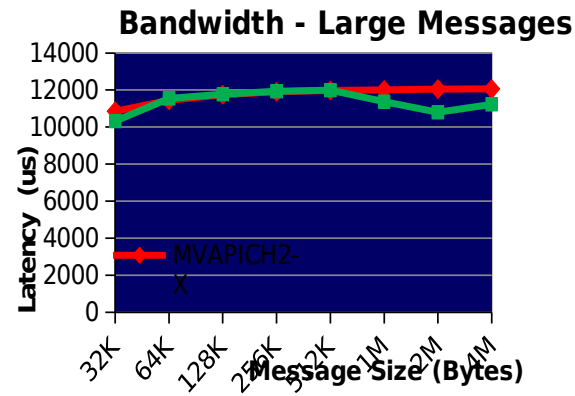
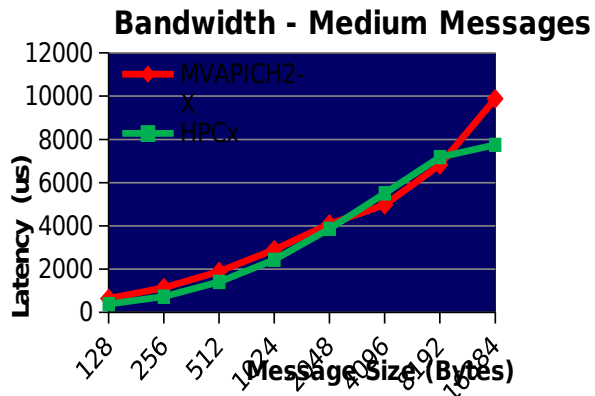
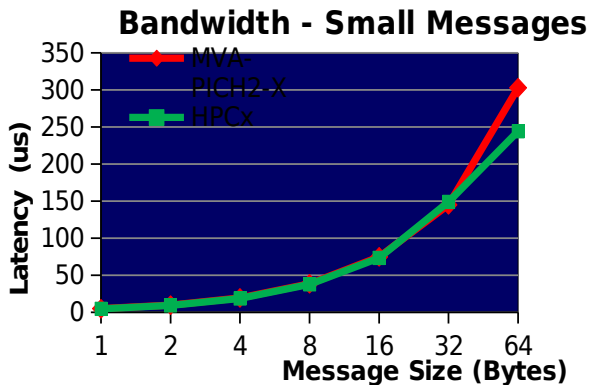
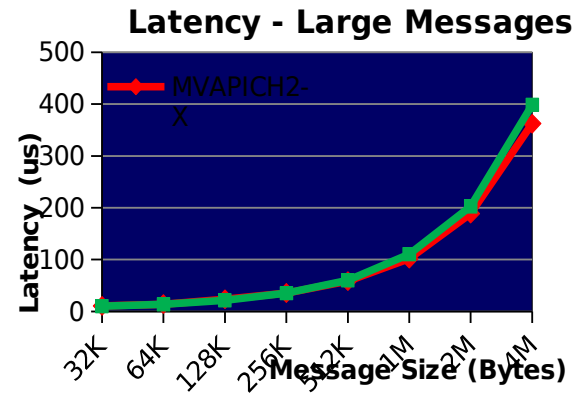
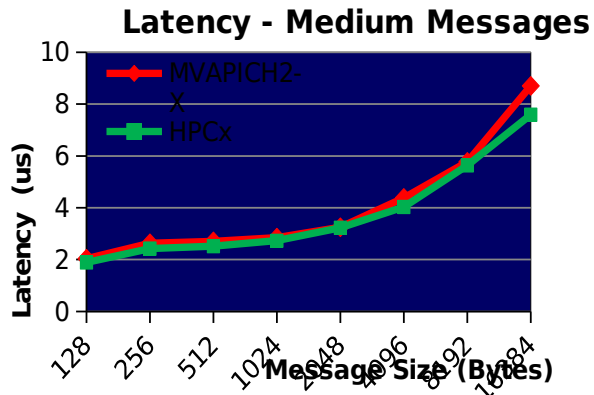
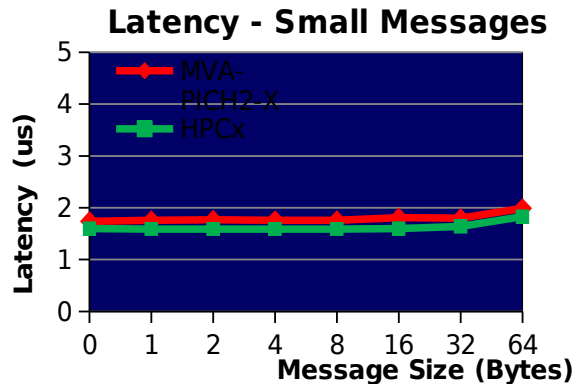


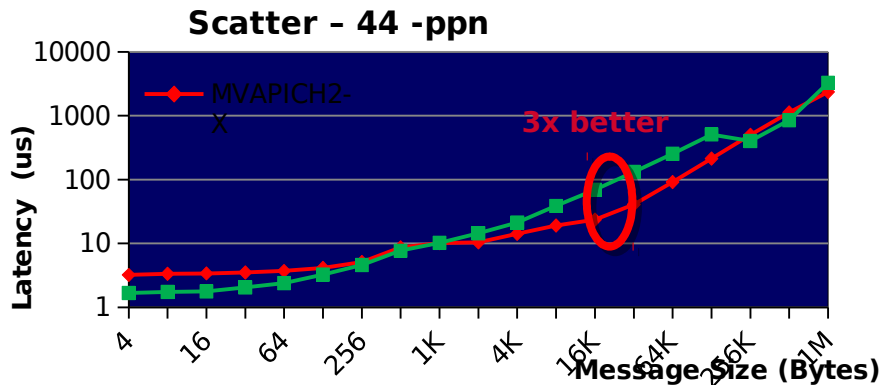
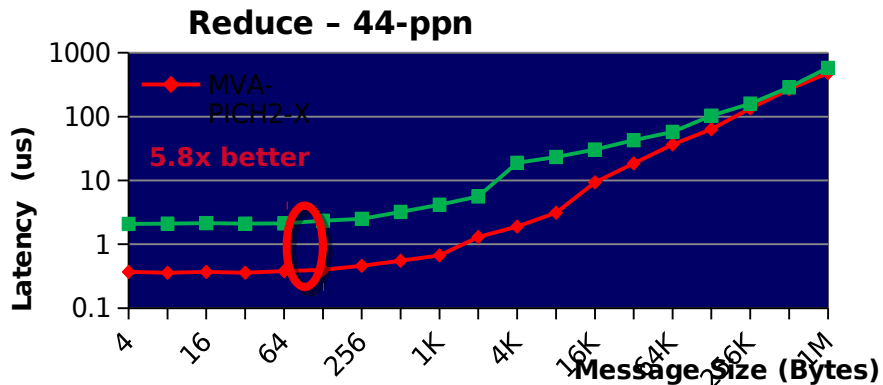
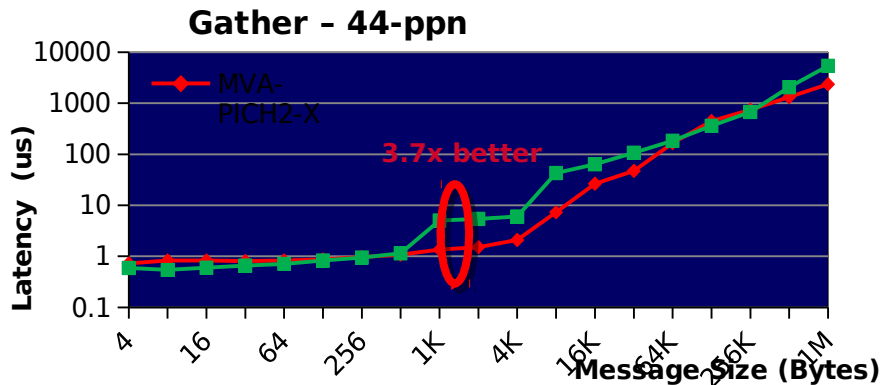
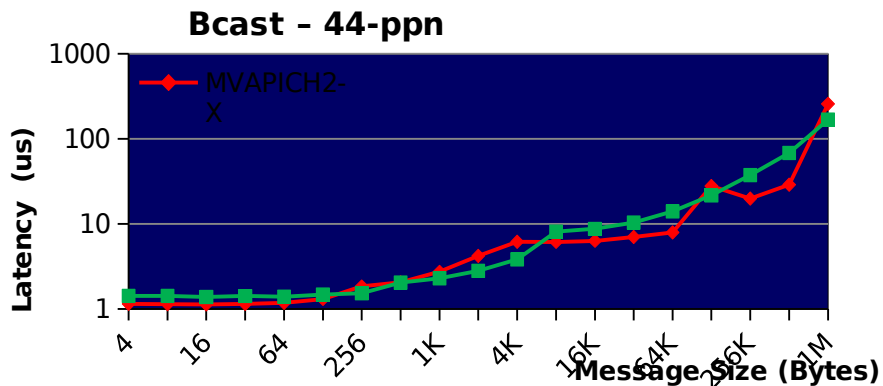
Deploy to Azure

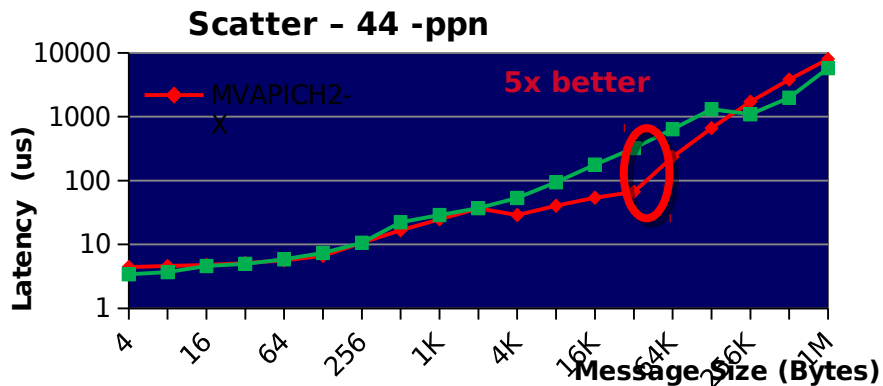
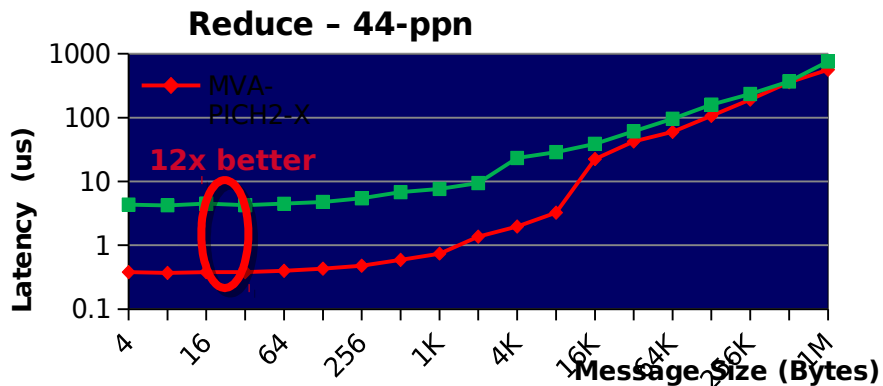
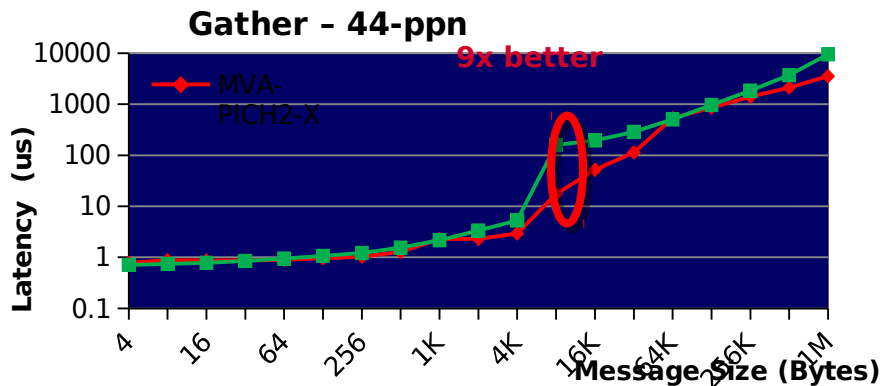
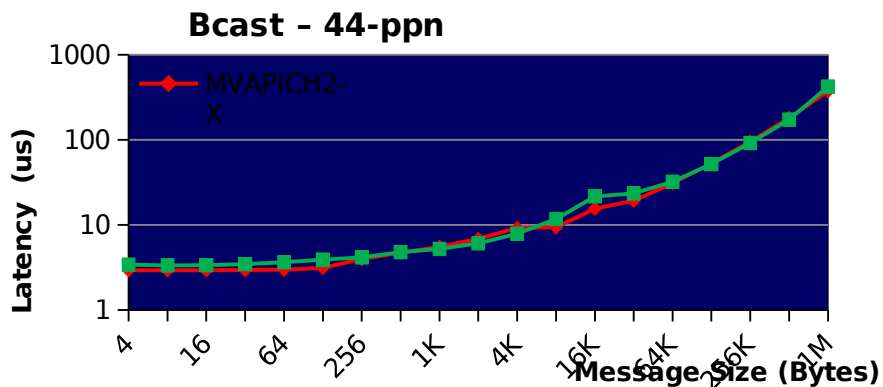
- Overview of the MVAPICH2 Project
- Overview of Azure HPC Environments
- Designing MVAPICH2 for Azure
- **Performance Results**
- Future Plans

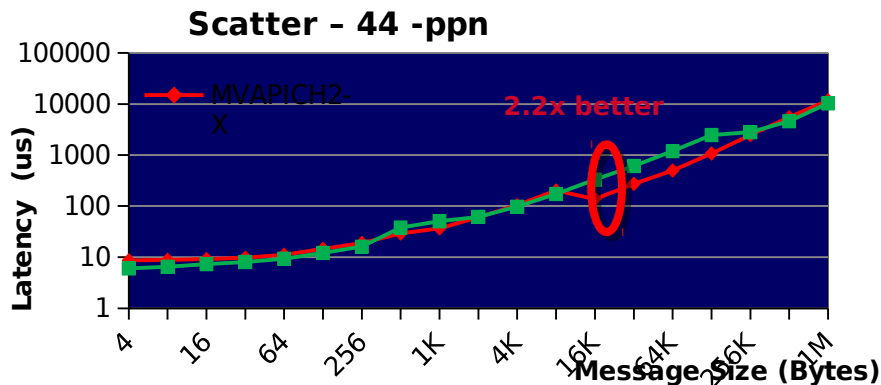
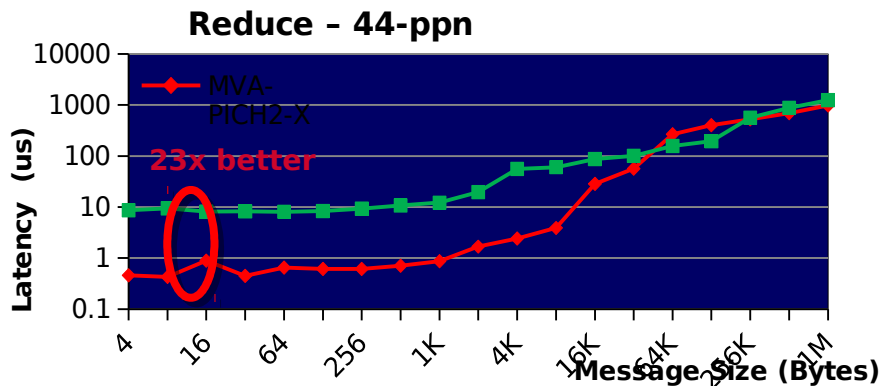
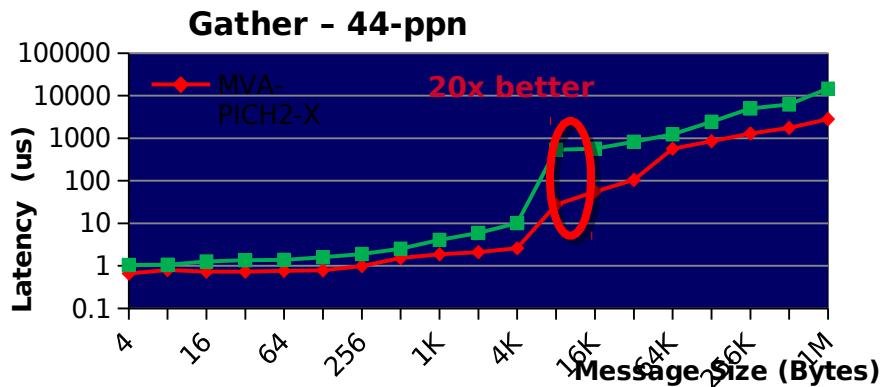
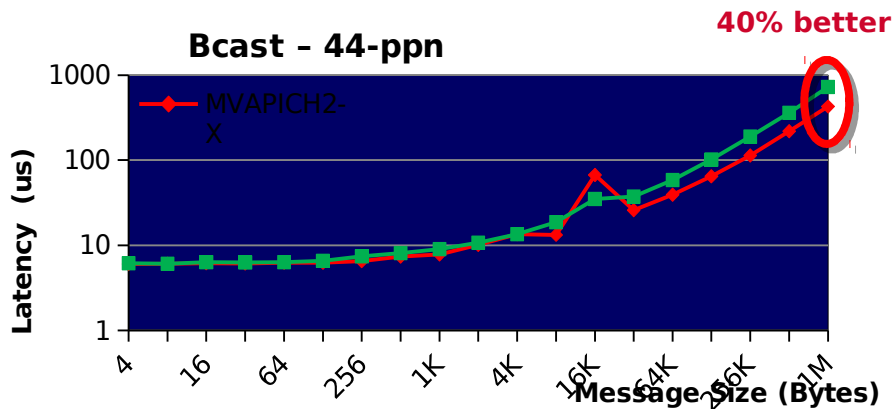
- VM type: Azure HC
- CPU: Intel(R) Xeon(R) Platinum 8168 CPU @ 2.70GHz
- Cores: 2 socket, 44 cores
- MVAPICH2 Version: Latest MVAPICH2-X w/ XPMEM support
- HPCx Version: Built-in HPCx-v2.5.0-gcc-MLNX_OFED_LINUX-4.7-1.0.0.1-redhat7.6-x86_64
- OMB Version: OSU-MicroBenchmars-5.6.2

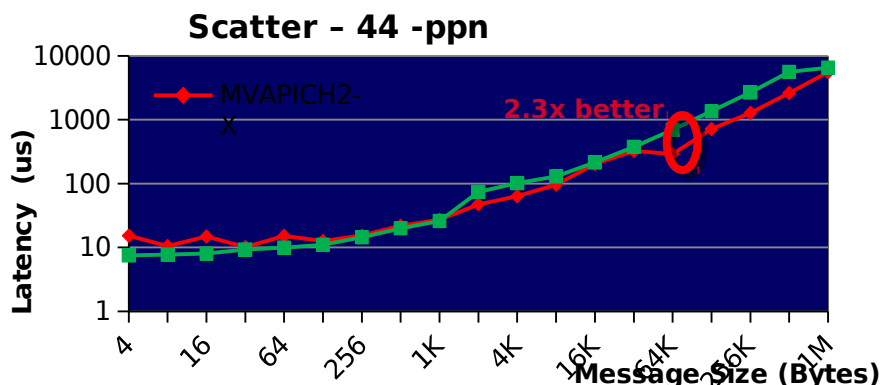
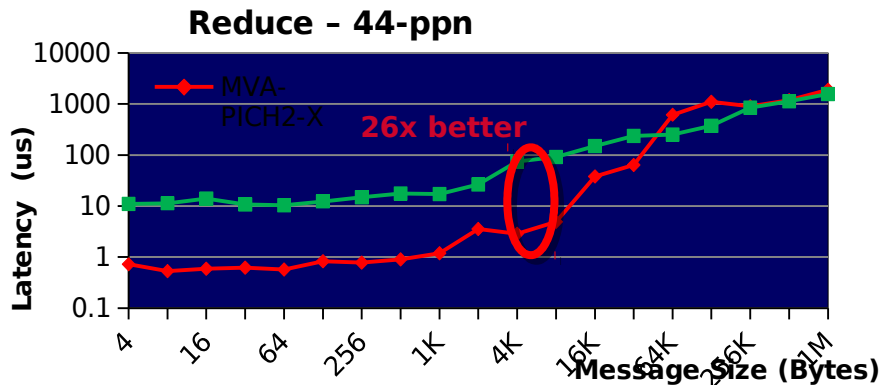
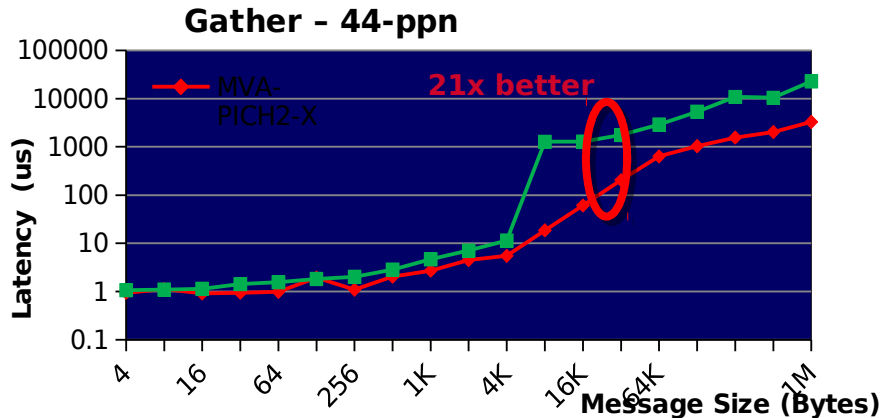
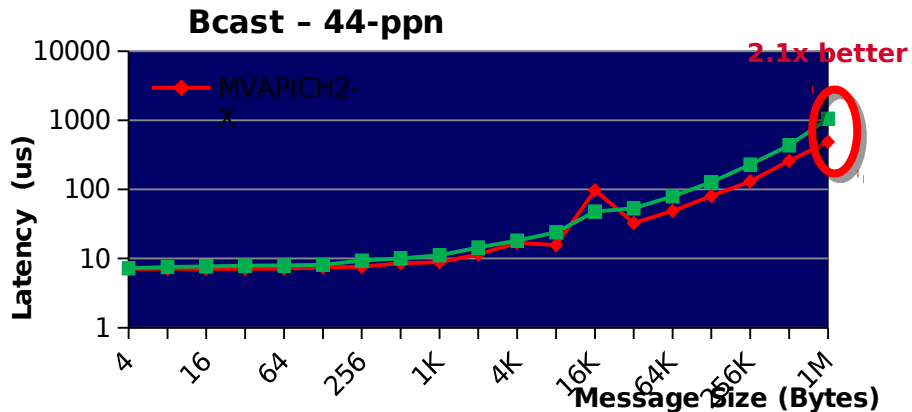




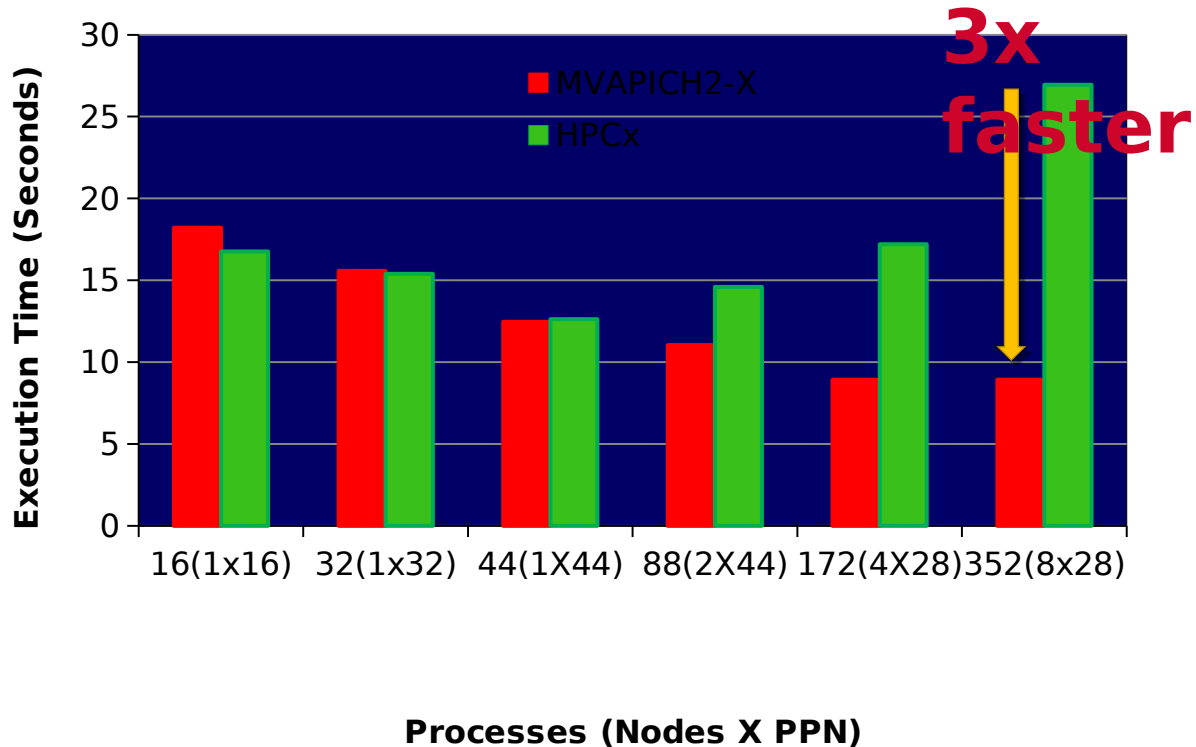


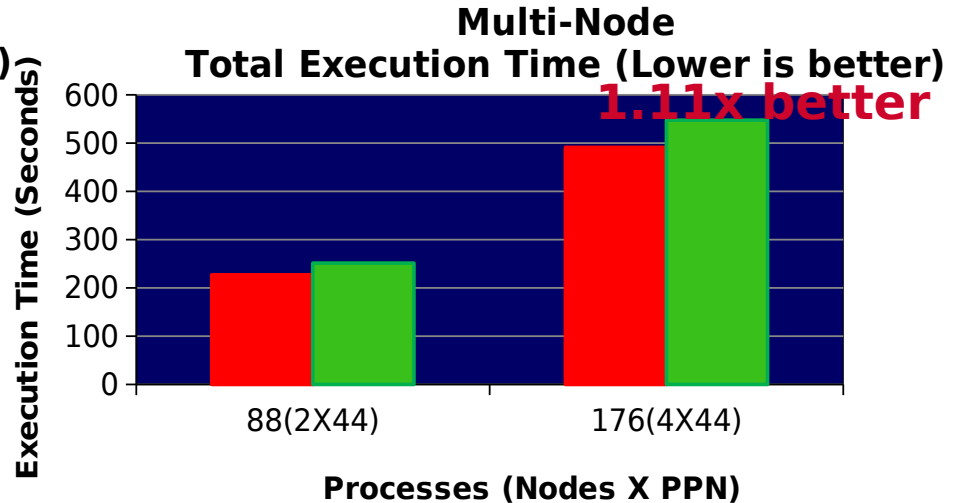
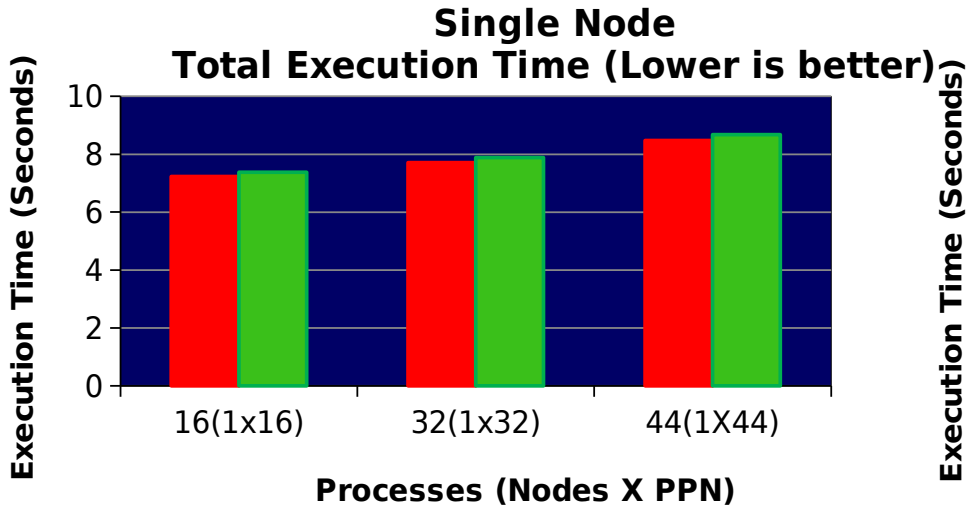






Total Execution Time (Lower is better)





Part of input parameter:

`&MESH IJK=5,5,5, XB=-1.0,0.0,-1.0,0.0,0.0,1.0, MULT_ID='mesh array'`

- Overview of the MVAPICH2 Project
- Overview of Azure HPC Environments
- Designing MVAPICH2 for Azure
- Performance Results
- **Future Plans**

- Additional optimizations and tuning of MVAPICH2-X
- A new version will be released soon
- Making it available in an integrated manner in the Azure portal
- Making it available through Azure Market Place
- **Commercial Support available for End-Users, ISVs, and Organizations**
 - Through X-ScaleSolutions (<http://x-scalesolutions.com>)

Thank You!

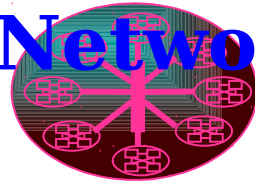
panda@cse.ohio-state.edu



Follow us on

<https://twitter.com/mvapich>

Network Based Computing



Laboratory

Network-Based Computing Laboratory

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



MVAPICH

MPI, PGAS and Hybrid MPI+PGAS Library

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/>



High-Performance
Deep Learning

The High-Performance Deep Learning
Project

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