



**MVAPICH**

MPI, PGAS and Hybrid MPI+PGAS Library

# Visualize and Analyze your Network Activities using OSU INAM

Talk at OSU Booth (November '18)

by

**Hari Subramoni**

The Ohio State University

E-mail: [subramon@cse.ohio-state.edu](mailto:subramon@cse.ohio-state.edu)

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

# Outline

- Introduction & Motivation
- Design of OSU INAM
- Impact of Profiling on Application Performance
- Features of OSU INAM & Demo
- Conclusions & Future Work

# Motivation

- IB clusters and the MPI-based applications complex
- Challenging to identify interaction between and impact of underlying IB network on performance of HPC application
- Such knowledge critical to maximize efficiency and performance of HPC applications
- Rely on a plethora of MPI level and IB level tools to analyze and understand an HPC system to answer questions like
  - *Why is my application running slower than usual now?*

# Limitations of Existing IB Fabric Monitoring Tools

- Several tools exist to analyze and inspect the IB fabric
  - e.g.: Nagios, Ganglia, Mellanox Fabric IT, INAM, BoxFish
- **Lack of interaction with & knowledge about MPI library**
  - Cannot classify traffic based on MPI primitives
    - e.g.: Point-to-point, Collective, RMA
  - Cannot correlate of network level and MPI level behavior
- **Lack of interaction with the job scheduler**
  - Cannot classify network traffic as belonging to a particular job
  - Cannot pin point source of conflict at finer granularity

# Limitations of Existing MPI Profiling Tools

- Several tools exist that allow to profile MPI library
  - TAU, HPCToolkit, Intel Vtune, IPM, mpiP
- Lack of interaction with & knowledge about IB fabric
  - Cannot correlate network level and MPI level behavior
- Unable to provide deep insights into MPI library
  - Recently proposed MPI\_T interface enables deep introspection
  - e.g.: MPIAdvisor – No knowledge about the underlying IB fabric

## Broad Challenge

*How can we design a tool that enables in-depth understanding of the communication traffic on the InfiniBand network through tight integration with the MPI runtime?*

# Overview of OSU INAM

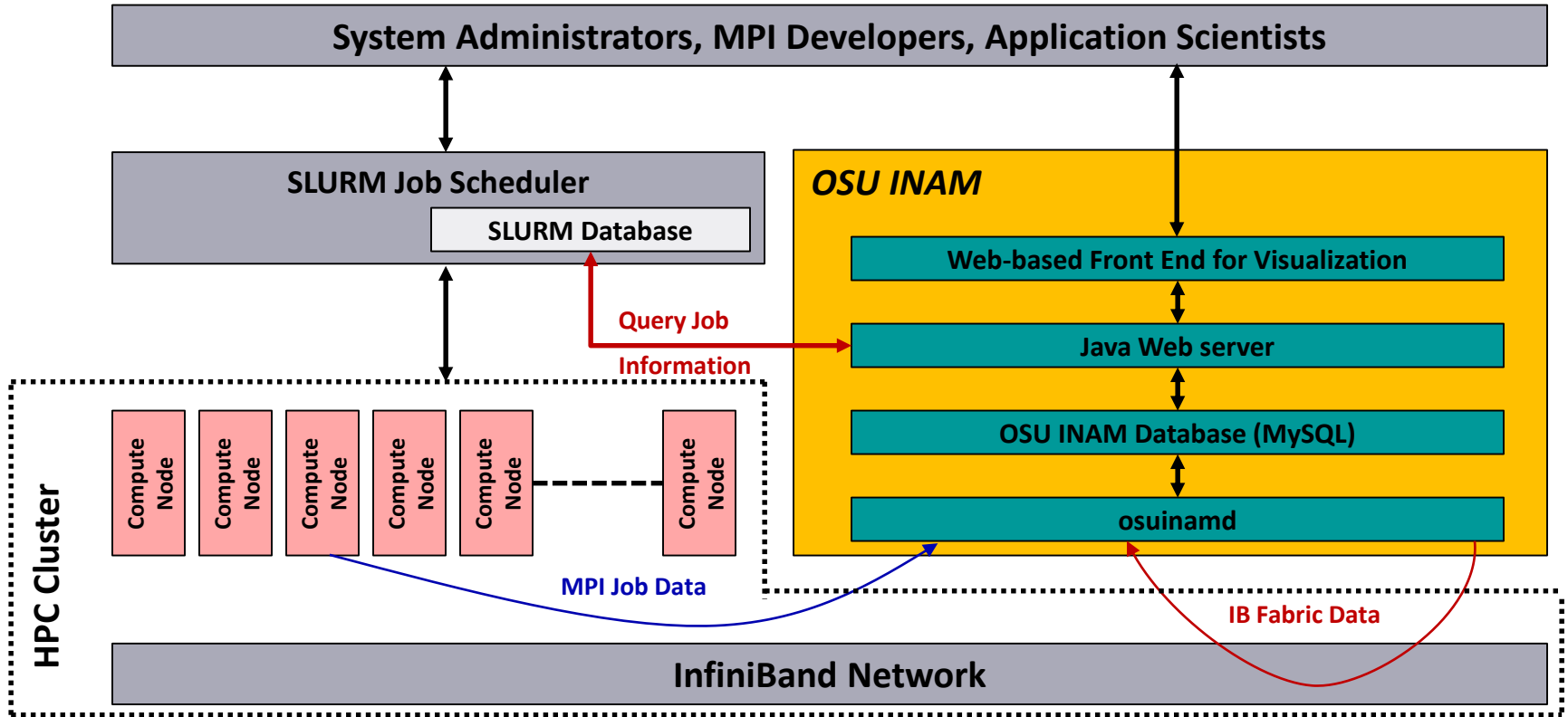
- A network monitoring and analysis tool that is capable of analyzing traffic on the InfiniBand network with inputs from the MPI runtime
  - <http://mvapich.cse.ohio-state.edu/tools/osu-inam/>
- Monitors IB clusters in real time by querying various subnet management entities and gathering input from the MPI runtimes
- Capability to analyze and profile **node-level, job-level and process-level activities** for MPI communication
  - Point-to-Point, Collectives and RMA
- Ability to filter data based on type of counters using “drop down” list
- Remotely monitor various metrics of MPI processes at user specified granularity
- "Job Page" to display jobs in ascending/descending order of various performance metrics in conjunction with MVAPICH2-X
- Visualize the data transfer happening in a “live” or “historical” fashion for entire network, job or set of nodes
- **OSU INAM 0.9.4 released on 11/10/2018**
  - Enhanced performance for fabric discovery using optimized OpenMP-based multi-threaded designs
  - Ability to gather InfiniBand performance counters at sub-second granularity for very large (>2000 nodes) clusters
  - Redesign database layout to reduce database size
  - Enhanced fault tolerance for database operations
    - Thanks to Trey Dockendorf @ OSC for the feedback
  - OpenMP-based multi-threaded designs to handle database purge, read, and insert operations simultaneously
  - Improved database purging time by using bulk deletes
  - Tune database timeouts to handle very long database operations
  - Improved debugging support by introducing several debugging levels

# Outline

- Introduction & Motivation
- Design of OSU INAM
- Impact of Profiling on Application Performance
- Features of OSU INAM & Demo
- Conclusions & Future Work



# OSU INAM Framework



# MPI Data Collection Thread

- Collect data specific to each MPI process and pushes it to OSU INAM Database
  - Allows analysis and visualization job/node/process level granularities
- Thread is a listener – accepts data from remote MPI processes
  - Avoid bottlenecks that arise where thread actively polls each MPI process
- OSU INAM communication requirements
  - **Small messages; High performance; High scalability, Low latency and no requirement for absolute reliability**
- IB based communication to achieve high performance and low latency
  - **Uses interrupt driven mode in IB**
    - Reduce CPU utilization by eliminating the need to continually poll
- Design choices for IB transport protocol
  - IB supports several transport protocols – RC, XRC, DC, UD
  - UD / DC transport protocols have significant benefits for scalability and memory footprint
  - **UD protocol as the IB transport protocol for the MPI data collection thread**

# Co-design of MPI and OSU INAM

- Enhance MPI\_T based profiling in MVAPICH2-X
  - CPU utilization of each process; Memory utilization of each process; Inter-node and intra-node communication buffer utilization; Intra-node, Inter-node and total bytes sent/received and, Total bytes sent for RMA operations
- MVAPICH2-X collects information via MPI\_T and transmits updates to the MPI data collection thread via UD Queue Pairs (QP) at user specified intervals
  - Default value: 30 seconds
- Each packet sent has some meta data information used later to retrieve the data from the database
- MPI data collection thread dumps the UD QP and Local Identifier (LID) that it is listening on to a file
- This location of this file is passed through environment variables to MPI runtime by the system administrator

# Fabric Discovery Thread & Database Thread

- Fabric Discovery Thread
  - Responsible for discovering the IB fabric and extracting data from selected components
  - Identify the various IB devices present in the network and their current status and stores in DB
  - Computes network path between each pair of hosts and stores in DB
  - Monitor the network for any changes at a user specified interval
  - Queries performance counters from selected components at user specified intervals
  - Queues up the message in FIFO to the database thread for eventual insertion into the database
- Database Thread
  - Responsible for receiving information from the MPI data collection and the FD threads
  - Create the tables in schema that the tool expects
    - Automatically update tables used by earlier versions of tool

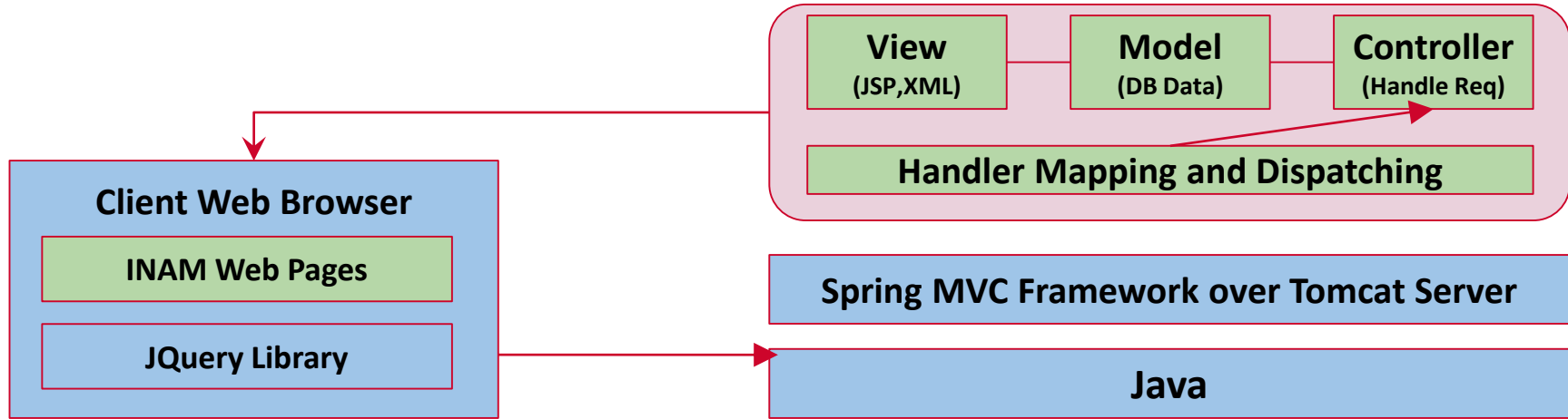
# Design of OSU INAM Database

- Consists of multiple tables to enable various features of OSU INAM
  - Tables to hold InfiniBand network infrastructure related data
    - “route”, “links”, “nodes”, “port data counters”, and “port errors”
    - Hold data for links, nodes, ports and routes
  - Tables to keep track of MPI process communication characteristics
    - “process info”, “process comm main”, and “process comm grid”
- Allows OSU INAM to
  - Analyze and profile node-level, job-level and process-level activities for MPI
  - Profile and report parameters/counters of MPI processes at the node-, job- and process-level
  - Visualize the communication map at process-level and node-level granularities
  - Analyzing and classifying InfiniBand network traffic flows in a physical link

# Design of Java-based Web Server

- Queries OSU INAM and SLURM databases to obtain MPI, Network and Job specific information
  - Users can modify frequency of query
- Validates and correlates results of different queries and presents data to the user in an unified fashion
- Based on the Spring MVC (Model, View and Controller) architecture
- Client side uses light-weight JQuery library to send HTTP requests through AJAX
- OSU INAM can send data to and retrieve responses from the server asynchronously
  - Dramatically improves user experience by hiding data processing / page rendering in the background

# Design of Web-based Front-end Visualization



1. HTTP request by users action sent to server side by Web browser / JQuery library with AJAX
2. Tomcat server receives the request, passes it to Spring framework
3. Spring framework dispatches request to the corresponding controller
4. Selected controller queries the model for some information in database
5. After processing, the Spring framework receives response to build the view through JSP, XML, etc
6. HTTP response will be sent back to the browser at the client side and the Web page will get updated

# Outline

- Introduction & Motivation
- Design of OSU INAM
- Impact of Profiling on Application Performance
- Features of OSU INAM & Demo
- Conclusions & Future Work



# Experimental Setup

- Each node of our 184 node testbed has eight Intel Xeon cores running at 2.53 Ghz with 12 MB L3 cache; 12 GB of memory and Gen2 PCI-Express bus
- Equipped with MT26428 QDR ConnectX-2 HCAs
- Interconnected using Mellanox MTS3610 QDR switch, with 11 leafs, each having 16 ports.
- The operating system used is Red Hat Enterprise Linux Server release 6.5 (Santiago), with the 2.6.32-431.el6.x86\_64 kernel version
- Mellanox OFED version 2.2-1.0.1 is used on all machines.

# Overview of the MVAPICH2 Project

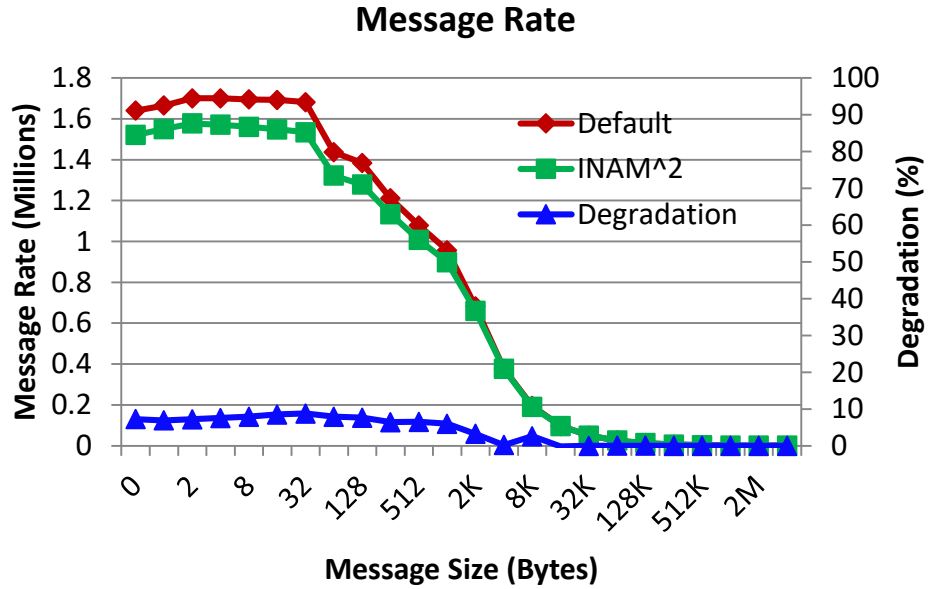
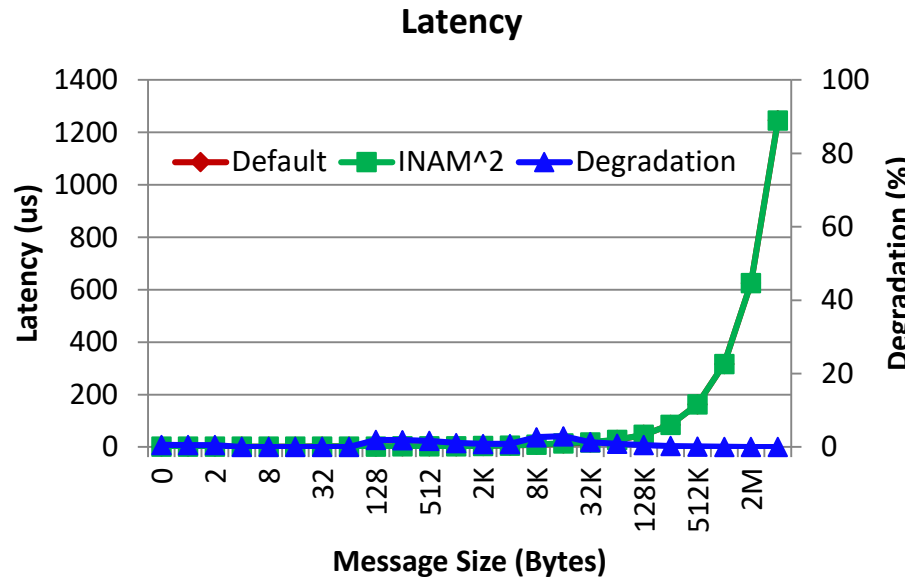
- High Performance open-source MPI Library for InfiniBand, Omni-Path, Ethernet/iWARP, and RDMA over Converged Ethernet (RoCE)
  - MVAPICH (MPI-1), MVAPICH2 (MPI-2.2 and MPI-3.1), Started in 2001, First version available in 2002
  - MVAPICH2-X (MPI + PGAS), Available since 2011
  - Support for GPGPUs (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 505,000 (> 0.5 million) downloads from the OSU site directly**
  - Empowering many TOP500 clusters (Jul '18 ranking)
    - 2<sup>nd</sup> ranked 10,649,640-core cluster (Sunway TaihuLight) at NSC, Wuxi, China
    - 12<sup>th</sup>, 556,104 cores (Oakforest-PACS) in Japan
    - 15<sup>th</sup>, 367,024 cores (Stampede2) at TACC
    - 24<sup>th</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

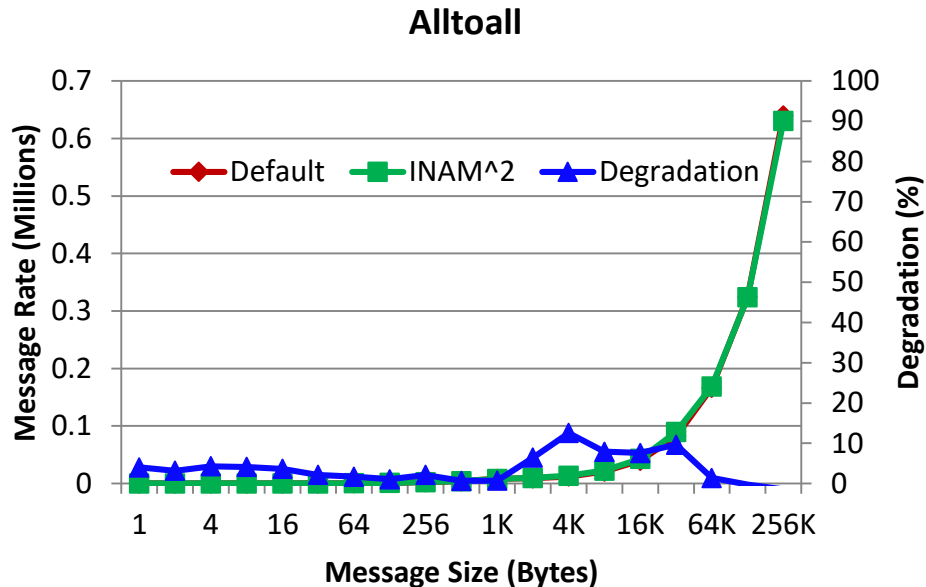
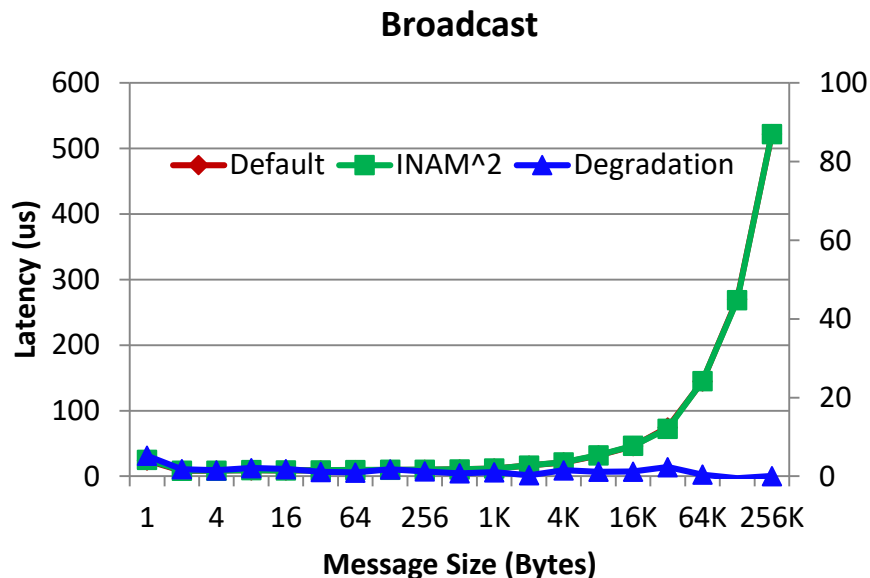
Partner in the upcoming TACC Frontera System

# Impact of Profiling on Performance of Point-to-point Operations



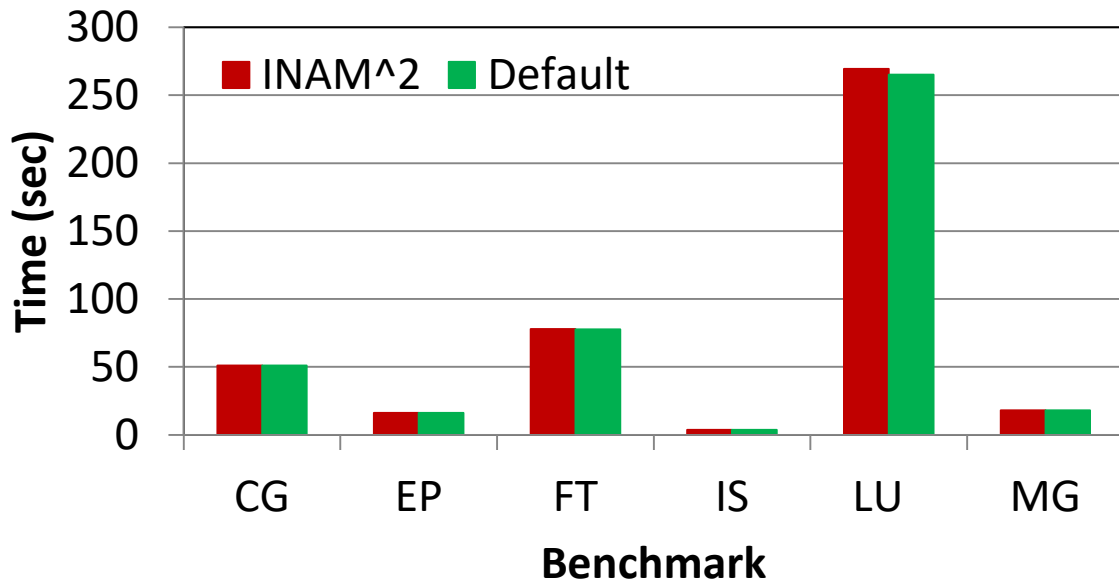
- Data collection adds very less than degradation when compared to the native performance

# Impact of Profiling on Performance of Collectives



- Performance of Broadcast and Alltoall at 512 processes
- Data collection adds very less degradation when compared to the native performance

# Impact of Profiling on Performance of NAS Parallel Benchmarks



- Performance of NAS parallel benchmarks at 512 processes
- Little to no impact on the performance due to the addition of the data collection and reporting

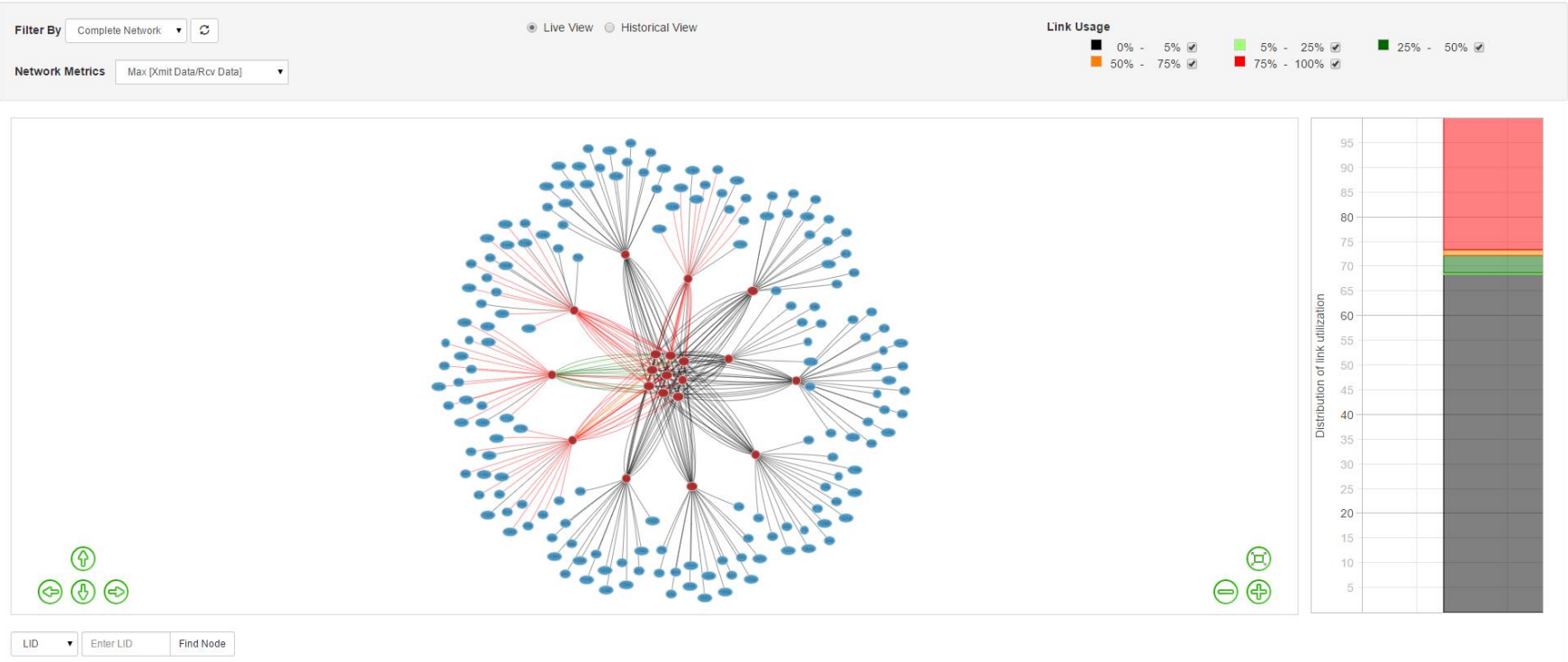
# Outline

- Introduction & Motivation
- Design of OSU INAM
- Impact of Profiling on Application Performance
- Features of OSU INAM & Demo
- Conclusions & Future Work

## Discussion on Features of OSU INAM

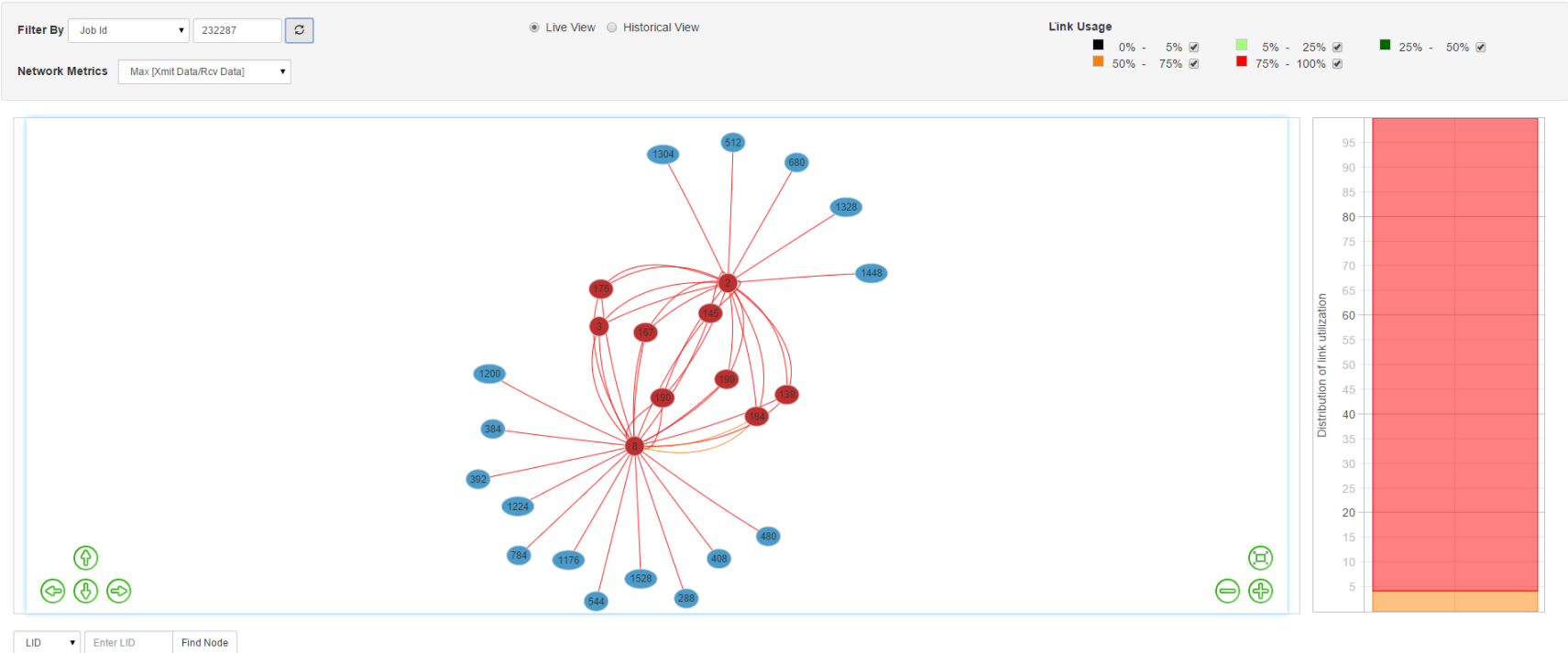
- Analyzing and Understanding Inter-node Communication Buffer Allocation and Use
- Identifying and Analyzing Sources of Link Congestion
- Monitoring Jobs Based on Various Metrics
- Capability to Profile and Report Several Metrics of MPI Processes at Different Granularities

# Live Network Level View





# Live Job Level View



# Live Node Level View



# Live Node Level View (Cont.)

**Node Information**

**Node Details**

NAME : **node158 HCA-1**  
LID : **384**  
GUID: **0x0002c903000a9119**

**Job Information**

Job Id : 232287  
Start Time :Wed Sep 09 2015 13:56:37 GMT-0400 (Eastern Daylight Time)  
Nodes : node001 node002 node003 node004 node005 node019 node020 node151 node152 node153 node154 node155 node156 node157 node158 node159

**CPU Usage**

Core Level ▾

**CPU Utilization**

■ User ■ System ■ Other ■ Idle

core #	User (%)	System (%)	Other (%)	Idle (%)
0	98	1	1	0
1	98	1	1	0
2	98	1	1	0
3	98	1	1	0
4	98	1	1	0
5	98	1	1	0
6	98	1	1	0
7	98	1	1	0

Rank112 [ core 0]

Rank113 [ core 1]

# Live Switch Level View

## Switch Information

Switch: MF0:bswitch.MTS3610/L04/U1

Start Time: 09/09/15 13:30:24

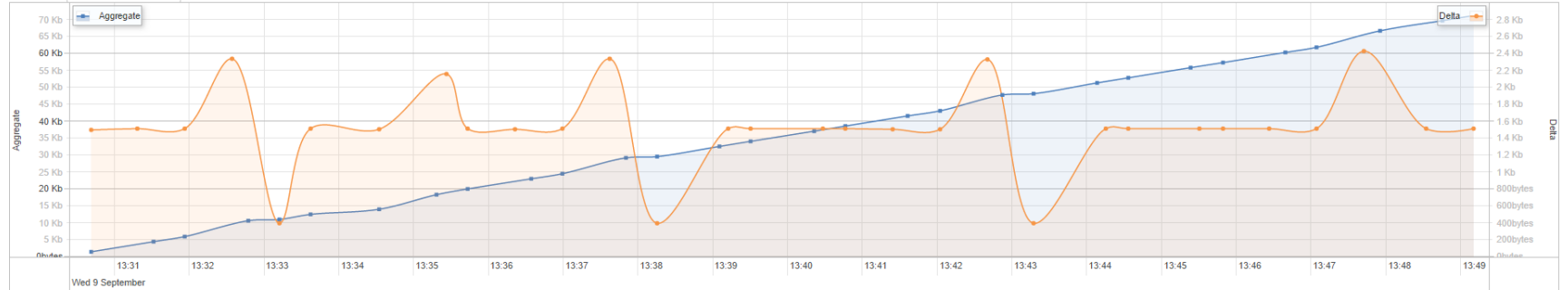
End Time: 09/09/15 13:49:24

- Select switch to monitor
- Enter start time and end time.
- Click on desired port panel

NAME : MF0:bswitch.MTS3610/L04/U1  
LID : 39  
GUID : 0x0002c90200424430

## Port 1 [node124 HCA-1]

Port Counter: Xmit Data



Port 2 [node125 HCA-1]

Port 3 [node126 HCA-1]

Port 4 [node127 HCA-1]

# List of Supported Switch Counters

- The following counters are queried from the InfiniBand Switches
- Xmit Data
  - Total number of data octets, divided by 4, transmitted on all VLs from the port
  - This includes all octets between (and not including) the start of packet delimiter and the VCRC, and may include packets containing errors
  - Excludes all link packets.
- Rcv Data
  - Total number of data octets, divided by 4, received on all VLs from the port
  - This includes all octets between (and not including) the start of packet delimiter and the VCRC, and may include packets containing errors
  - Excludes all link packets.
- Max [Xmit Data/Rcv Data]
  - Maximum of the two values above

# List of Supported MPI Process Level Counters

- MVAICH2-X collects additional information about the process's network usage which can be displayed by OSU INAM
- Xmit Data
  - Total number of bytes transmitted as part of the MPI application
- Rcv Data
  - Total number of bytes received as part of the MPI application
- Max [Xmit Data/Rcv Data]
  - Maximum of the two values above
- Point to Point Send
  - Total number of bytes transmitted as part of MPI point-to-point operations
- Point to Point Rcvd
  - Total number of bytes received as part of MPI point-to-point operations
- Max [Point to Point Sent/Rcvd]
  - Maximum of the two values above
- Coll Bytes Sent
  - Total number of bytes transmitted as part of MPI collective operations
- Coll Bytes Rcvd
  - Total number of bytes received as part of MPI collective operations

# List of Supported MPI Process Level Counters (Cont.)

- Max [Coll Bytes Sent/Rcvd]
  - Maximum of the two values above
- RMA Bytes Sent
  - Total number of bytes transmitted as part of MPI RMA operations
  - Note that due to the nature of the RMA operations, bytes received for RMA operations cannot be counted
- RC VBUF
  - The number of internal communication buffers used for reliable connection (RC)
- UD VBUF
  - The number of internal communication buffers used for unreliable datagram (UD)
- VM Size
  - Total number of bytes used by the program for its virtual memory
- VM Peak
  - Maximum number of virtual memory bytes for the program
- VM RSS
  - The number of bytes resident in the memory (Resident set size)
- VM HWM
  - The maximum number of bytes that can be resident in memory (Peak resident set size or High water mark)

# List of Supported Network Error Counters

- The following error counters are available both at switch and process level:
- SymbolErrors
  - Total number of minor link errors detected on one or more physical lanes
- LinkRecovers
  - Total number of times the Port Training state machine has successfully completed the link error recovery process
- LinkDowned
  - Total number of times the Port Training state machine has failed the link error recovery process and downed the link
- RcvErrors
  - Total number of packets containing an error that were received on the port. These errors include:
    - Local physical errors
    - Malformed data packet errors
    - Malformed link packet errors
    - Packets discarded due to buffer overrun
- RcvRemotePhysErrors
  - Total number of packets marked with the EBP delimiter received on the port.
- RcvSwitchRelayErrors
  - Total number of packets received on the port that were discarded because they could not be forwarded by the switch relay



# List of Supported Network Error Counters (Cont.)

- XmtDiscards
  - Total number of outbound packets discarded by the port because the port is down or congested. Reasons for this include:
    - Output port is not in the active state
    - Packet length exceeded NeighborMTU
    - Switch Lifetime Limit exceeded
    - Switch HOQ Lifetime Limit exceeded This may also include packets discarded while in VLStalled State.
- XmtConstraintErrors
  - Total number of packets not transmitted from the switch physical port for the following reasons:
    - FilterRawOutbound is true and packet is raw
    - PartitionEnforcementOutbound is true and packet fails partition key check or IP version check
- RcvConstraintErrors
  - Total number of packets not received from the switch physical port for the following reasons:
    - FilterRawInbound is true and packet is raw
    - PartitionEnforcementInbound is true and packet fails partition key check or IP version check
- LinkIntegrityErrors
  - The number of times that the count of local physical errors exceeded the threshold specified by LocalPhyErrors
- ExcBufOverrunErrors
  - The number of times that OverrunErrors consecutive flow control update periods occurred, each having at least one overrun error
- VL15Dropped
  - Number of incoming VL15 packets dropped due to resource limitations (e.g., lack of buffers) in the port

# Outline

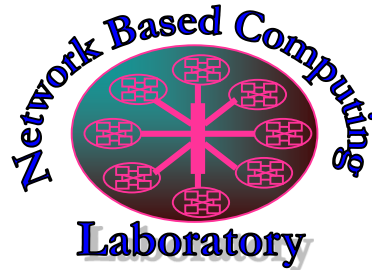
- Introduction & Motivation
- Design of OSU INAM
- Impact of Profiling on Application Performance
- Features of OSU INAM & Demo
- **Conclusions & Future Work**

# Conclusions & Future Work

- Designed OSU INAM capable of analyzing the communication traffic on the InfiniBand network with inputs from the MPI runtime
- Latest version (v0.9.3) available for free download from
  - <http://mvapich.cse.ohio-state.edu/tools/osu-inam/>
- OSU INAM has been downloaded more than 500 times directly from the OSU site
- Provides the following major features
  - Analyze and profile network-level activities with many parameters (data and errors) at user specified granularity
  - Capability to analyze and profile node-level, job-level and process-level activities for MPI communication (Point-to-Point, Collectives and RMA)
  - Remotely monitor CPU utilization of MPI processes at user specified granularity
  - Visualize the data transfer happening in a "live" or historical fashion for Entire Network, Particular Job One or multiple Nodes, One or multiple Switches
- Future Work
  - Add support to profile and analyze GPU-based communication
  - Capability to profile various PGAS programming languages

# Thank You!

[subramon@cse.ohio-state.edu](mailto:subramon@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/>