





# Visualize, Analyze, and Correlate Networking Activities for Your Parallel Programs on InfiniBand HPC Clusters using the OSU INAM Tool

# **Tutorial at MUG'22**

# by

Hari Subramoni

The Ohio State University

subramon@cse.ohio-state.edu

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

Pouya Kousha

The Ohio State University

kousha.2@buckeyemail.osu.edu

https://www.linkedin.com/in/pouya-kousha2/

Ayyappa Kolli

The Ohio State University

kolli.28@buckeyemail.osu.edu

https://www.linkedin.com/in/ayyappa-kolli/

**Network Based Computing Laboratory** 

**MUG'22** 

#### Outline

#### • Overview

- Motivation
- Profiling Tools Perspective and Broad Challenges
- OSU INAM Components
- Overhead Analysis
- Downloading, Installing, and Configuring OSU INAM
- Demos
- Future Work and Concluding Remarks

#### Introduction

- With the recent advances in HPC, providing efficient data movement is critical to have end-to-end high throughput solutions
- A detailed and real-time insight of network level and the communication is needed to identify and alleviate performance bottlenecks
- A newer set of challenges arise as HPC systems are becoming larger, users expect to have better capabilities like real-time profiling at fine granularity and scalability concerns
- The problem of identifying any performance issues with HPC and DL applications is akin to finding "a needle in a haystack" in HPC communication layers



# **Profiling Tools Perspective and Broad Challenges**

- There are 30+ profiling tools for HPC systems
- System level vs User level
  - User level novelty
- Different set of users have different needs
  - HPC administrators
  - HPC Software developers
  - Domain scientists
- Different HPC layers to profile
  - How to correlate them and pinpoint the problem source?





# Summary of existing profiling tools and their capabilities

| Tools              |              | MPI Runtime    |                       |
|--------------------|--------------|----------------|-----------------------|
| TOOIS              | Applications | Network Fabric | Job scheduler         |
| INAM*              | 1            | 1              | ✓                     |
| TAU                | 1            | 1              | ×                     |
| HPCToolkit         | 1            | ×              | ×                     |
| Intel Vtune        | 1            | ×              | ×                     |
| IPM                | 1            | ×              | ×                     |
| mpiP               | 1            | ×              | ×                     |
| Intel ITAC         | 1            | ×              | ×                     |
| ARM MAP            | 1            | ×              | ×                     |
| HVProf             | 1            | ×              | ×                     |
| PCP(used by XDMOD) | ×            | 1              | <ul> <li>✓</li> </ul> |
| Prometheus         | ×            | 1              | <ul> <li>✓</li> </ul> |
| Mellanox FabricIT  | ×            | 1              | ×                     |
| BoxFish            | ×            | 1              | ×                     |
| LDMS               | ×            | 1              | ×                     |

\* This design has been publicly released on 06/08/2020 and is available for free here https://mvapich.cse.ohio-state.edu/tools/osu-inam/

# **Profiling Tools Perspective and Broad Challenges**

- Understanding the interaction between applications, MPI libraries, I/O and the communication fabric is challenging
  - Find root causes for performance degradation
  - Identify which layer is causing the possible issue
  - Understand the internal interaction and interplay of MPI library components and network level
  - Online profiling

**HPC** Applications Job Scheduler **MPI Library** Rank 0 Rank 1 ••• Rank K MPI T **HPC** Network Communication Fabric I/O File System

6

How can we design a tool that enables holistic, real-time, scalable and in-depth understanding of communication traffic through tight integration with the MPI runtime and job scheduler?

# **Overview of the MVAPICH2 Project**

- High Performance open-source MPI Library
- Support for multiple interconnects
  - InfiniBand, Omni-Path, Ethernet/iWARP, RDMA over Converged Ethernet (RoCE), AWS EFA, Rockport Networks, and Slingshot10/11, Broadcom, Cornelis Networks OPX
- Support for multiple platforms
  - x86, OpenPOWER, ARM, Xeon-Phi, GPGPUs (NVIDIA and AMD)
- Started in 2001, first open-source version demonstrated at SC '02
- Supports the latest MPI-3.1 standard
- <u>http://mvapich.cse.ohio-state.edu</u>
- Additional optimized versions for different systems/environments:
  - MVAPICH2-X (Advanced MPI + PGAS), since 2011
  - MVAPICH2-GDR with support for NVIDIA (since 2014) and AMD (since 2020) GPUs
  - MVAPICH2-MIC with support for Intel Xeon-Phi, since 2014
  - MVAPICH2-Virt with virtualization support, since 2015
  - MVAPICH2-EA with support for Energy-Awareness, since 2015
  - MVAPICH2-Azure for Azure HPC IB instances, since 2019
  - MVAPICH2-X-AWS for AWS HPC+EFA instances, since 2019
- Tools:
  - OSU MPI Micro-Benchmarks (OMB), since 2003
  - OSU InfiniBand Network Analysis and Monitoring (INAM), since 2015



- Used by more than 3,275 organizations in 90 countries
- More than 1.61 Million downloads from the OSU site directly
- Empowering many TOP500 clusters (June '22 ranking)
  - 6<sup>th</sup>, 10,649,600-core (Sunway TaihuLight) at NSC, China
  - **16**<sup>th</sup>, 448, 448 cores (Frontera) at TACC
  - 30<sup>th</sup>, 288,288 cores (Lassen) at LLNL
  - 42<sup>nd</sup>, 570,020 cores (Nurion) in South Korea and many more
- available with software stacks of many vendors and Linux Distros (RedHat, SuSE, OpenHPC, and Spack)
- Partner in the 16<sup>th</sup> ranked TACC Frontera system
- Empowering Top500 systems for more than 20 years

# **Overview of OSU InfiniBand Network Analysis and Monitoring (INAM) Tool**

- A network monitoring and analysis tool that is capable of analyzing traffic on the InfiniBand network with inputs from the MPI runtime
  - <u>http://mvapich.cse.ohio-state.edu/tools/osu-inam/</u>
- 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
- Sub-second port query and fabric discovery in less than 10 mins for ~2,000 nodes

- OSU INAM 0.9.8 released (08/11/2022)
  - Support for MySQL and InfluxDB as database backends
  - Enhanced database insertion using InfluxDB
  - Support for continuous queries to improve visualization performance
  - Support for SLURM multi-cluster configuration
  - Significantly improved database query performance when using InfluxDB
  - Support for automatic data retention policy when using InfluxDB
  - Support for PBS and SLURM job scheduler as config time
  - Ability to gather and display Lustre I/O for MPI jobs
  - Enable emulation mode to allow users to test OSU INAM tool in a sandbox environment without actual deployment
  - · Generate email notifications to alert users when user defined events occur
  - Support to display node-/job-level CPU, Virtual Memory, and Communication Buffer utilization information for historical jobs
  - Support to handle multiple job schedulers on the same fabric
  - Support to collect and visualize MPI\_T based performance data
  - Support for MOFED 4.5, 4.6, 4.7, and 5.0
  - Support for adding user-defined labels for switches to allow better readability and usability
  - Support authentication for accessing the OSU INAM webpage
  - Optimized webpage rendering and database fetch/purge capabilities
  - Support to view connection information at port level granularity for each switch
  - Support to search switches with name and lid in historical switches page
  - Support to view information about Non-MPI jobs in live node page

#### Outline

- Overview
- OSU INAM Components
- Overhead Analysis
- Downloading, Installing, and Configuring OSU INAM
- Demos
- Future Work and Concluding Remarks



#### **Flow of Using OSU INAM**



#### Outline

- Overview
- OSU INAM Components
- Overhead Analysis
  - Experiment Setup
  - Fabric Discovery Evaluation
  - Port Inquiry Evaluation
  - Visualization and MPI\_T Evaluation
  - Application-Level Overhead Analysis and Scaling
- Downloading, Installing, and Configuring OSU INAM
- Demos
- Future Work and Concluding Remarks

12

# **Experimental Setup**

| Cluster/Configuration      | RI  | OSC   |
|----------------------------|---|---|
| Overview                   | Equipped with<br>MT26428 QDR<br>ConnectX-2 HCAs<br>with PCI-Express<br>interfaces | 3 heterogeneous<br>clusters connected<br>to the same<br>InfiniBand<br>fabric                    |
| #Nodes                     | 146   | 1,428   |
| #Links                     | 542   | 3,402   |
| #Switches                  | 20 (36 ports/switch)  | 114 (36 ports/switch)   |
| CPU                        | Intel Xeon 2.53Ghz  | Intel(R) Xeon(R)<br>CPU E5-2680 v3<br>@ 2.50GHz   |
| Cores/node<br>for OSU INAM | 8   | 24  |
| L3 cache                   | 12MB  | 30MB  |
| Memory/node                | 12GB  | 128GB   |
| Switch<br>Technology       | Mellanox MTS3610QDR   | Two clusters have<br>Mellanox EDR<br>InfiniBand<br>(100Gbps)<br>and one with<br>FDR/EN (56Gbps) |
| Job Scheduler              | SLURM   | PBS   |

#### **Fabric Discovery Evaluation**

- We evaluated the impact of multithreading and performance variability
- To be fair and mimic the behavior of real deployment, we set some threads for port inquiry while measuring performance of fabric inquiry.
- The timings measured include discovery of fabric and insertion of data into the database



Impact of multi-threading on fabric discovery module on OSC cluster with hourly interval for 1428 nodes Impact of multi-threading on fabric discovery module on RI cluster with 1-minute interval Histogram of the time taken for fabric discovery module using 8 threads on RI cluster. Port inquiry frequency is set to 1 sec

• Enhanced performance for fabric discovery using optimized OpenMP-based multi-threaded designs with **6x speedup compared to serial version** 

#### **Port Inquiry Evaluation**

- We evaluated the impact of multithreading and performance variability
- To be fair and mimic the behavior of real deployment, we set some threads for fabric module while measuring performance of port inquiry.
- The timings measured include reading port metrics and insertion of data into the database



Impact of multi-threading on the latency of port inquiry sweep with 1 second query interval on the OSC clusters. OSC has 3,402 active ports Impact of multi-threading on latency of port inquiry sweep with 1 second querying interval for RI cluster. RI has 542 active port The latency of port inquiry sweep for 35,000 samples for OSC cluster with query interval of 250ms and using 16 thread

• Ability to remotely gather InfiniBand performance counters at sub-second granularity for very large (>2,000 nodes) clusters

# MPI\_T and Visualization Overhead

- Overhead of PVAR Collection per MPI Process
  - Query interval of 1s

Overhead of collecting PVAR data at nanosecond granularity

| Metrics          | Average   | Min    | Max       | STDDEV.p  |  |
|------------------|-----------|--------|-----------|-----------|--|
| Collecting PVARs | 517.63 ns | 140 ns | 16,204 ns | 305.91 ns |  |

- **Overhead of Visualization and**
- Rendering
  - Evaluating the timing of accessing and
  - visualizing on web UI

NETWORK AND LIVE JOBS VIEW GENERATION TIMING ON OSC WITH 1K JOBS

**Very Low** 

**Overhead** 

for Query

Phase

| View                | Average   | Min    | Max       | STDEV.p |
|---------------------|-----------|--------|-----------|---------|
| <b>Network View</b> | 196.15 ms | 187 ms | 206.09 ms | 5.75 ms |
| Live Jobs View      | 18.17 ms  | 16 ms  | 20 ms     | 1 ms    |

## **Application-Level Analysis and Session Scaling**

- Introspect performance of different NAS Class D benchmarks with and without MPI profiling (PVARs) with an interval of 5 seconds for 4,096 processes across 256 nodes
- Numbers are average of 25 integrations
- There is less than 5% performance degradation.



- Network overhead of adding additionall profiling sessions running on 3DStencil application
  - Compared to disabled MPI\_T mode

| 5 | #Sessions | Bytes sent for profiling | Bytes sent by application | Overhead |   |
|---|-----------|--------------------------|---------------------------|----------|---|
|   | 1         | 165 KB                   | 306 MB                    | 0.05%    |   |
|   | 4         | 390 KB                   | 306 MB                    | 0.12%    |   |
|   | 8         | 689 KB                   | 306 MB                    | 0.22%    |   |
|   |           |                          |                           |          | / |

#### Outline

- Overview
- OSU INAM Components
- Overhead Analysis
- Downloading, Installing, and Configuring OSU INAM
  - Downloading and Architecture Support
  - Deciding on Intervals for Query Data from HPC Communication Stack
  - Storage Data Modeling
  - OSU INAM Configuration
  - Best Practices for Resource Allocation among INAM Components
  - Configuration for Profiling MPI Jobs
  - Running Example MVAPICH2-X Jobs with INAM
- Demos
- Future Work and Concluding Remarks

18

# **Downloading and Architecture Support**

• OSU INAM is available at <u>http://mvapich.cse.ohio-state.edu/downloads/</u>

#### **OSU INAM**

OSU InfiniBand Network Analysis and Monitoring (OSU INAM) Tool v0.9.6 (06/08/20)

- The OSU INAM package is distributed under the BSD License .
- A detailed user guide with instructions to build, install and run OSU INAM is available here. This document also contains guidelines for troubleshooting and best practice deployment.
- Please see CHANGES for the full changelog.
- To estimate the expected size of the database on your system, please see the Database Size Calculator section.
- If you need OSU INAM for specific architecture and OFED version please email us at mvapich-help@cse.ohio-state.edu .
- These RPMs contain the OSU-INAM software on the corresponding distro. Please note that the RHEL RPMs are compatible with CentOS as well. For Debian/Ubuntu users, please follow the instructions in the install section in the userguide.

| Environment                  | Download link                                    |
|------------------------------|--|
| x86_64 MOFED 4.5             | Download EL7 MOFED 4.5 RPM (RHEL7/CentOS7)       |
| x86_64 MOFED 4.6 & MOFED 4.7 | Download EL7 MOFED 4.6 & 4.7 RPM (RHEL7/CentOS7) |
| x86_64 MOFED 5.0             | Download EL7 MOFED 5.0 RPM (RHEL7/CentOS7)       |
| ARM MOFED 4.5                | Download ARM EL7 MOFED 4.5 RPM (RHEL7/CentOS7)   |

Please note that, all RPMs were built using GNU (GCC) 4.8.5

#### Modeling of Aggregated InfiniBand Data on the Network and Storage

• How do we identify the amount of traffic added to the network from our tool?

 $Traffic_{(tool)} = Traffic_{(PC\_total)} + Traffic_{(FB\_total)}$ 

 $Traffic_{(PC\_total)} = 200Bytes imes QueryFrequency \ imes Number of Switches imes Number of SwitchPorts$ 

• For fabric discovery, it depends on the topology of the network and the diameter of the network

 $Traffic_{(FB\_total)} = NetworkDiameter imes 20Bytes imes$ 

Number of Switch Ports imes Query Frequency imes Number of Switches

- Considering Coc (114 switches and 1,420 compute noues connected in ough 0,402 links)
  - Port Inquiry traffic will be 1,603 KB/sec with interval of half second
  - fabric discovery traffic will be 801 KB/sec when using interval of 1 second for updating fabric
  - Given that the OSC clusters are capable of 56 Gbps to 100 Gbps, this is an insignificant amount of data.

#### Modeling of MPI\_T Data Storage Requirement

• For profiling of MPI Performance Variables (PVARs) in MPI jobs:

 $NumProcs = NumNodes \times ProcessesPerNode \times AverageClusterLoad$ 

 $Traffic_{PVARs} = NumPVAR(300) \times RecSize_{PVAR}(255B) \times$ 

 $NumProcs \times (1 + Sessions) \times Frequency$ 

• For CPU, Memory and basic processor information we have:

 $Traffic_{ProcessInfo} = NumProcs * RecSize(336B) * Frequency$ 

- Considering the full-load for cluster size of 2,000 nodes with 28 processes per node and 1Hz profiling query frequency
  - Assuming somehow one could activate all PVAR records per MPI process
  - Maximum PVARs traffic/storage will be 4.28 GB/s per profiling session on cluster level
  - 18.8 MB/s for MPI process information on cluster level

#### **OSU INAM Configuration**

- OSU INAM Daemon configuration specifies resource allocation for components of Daemon
  - Running only on one node!
  - OSU INAM Daemon configuration file is located at \$OSU\_INAM\_INSTALL\_PREFIX/etc/osuinamd.conf
- MySQL configuration can be optimized to use caching based on cluster size
  - See user guide:

http://mvapich.cse.ohio-state.edu/userguide/osu-inam/#\_mysql\_tuning\_parameters

- MySQL configuration file is the same as default located at /etc/my.conf
- OSU INAM web is configurable for access and visualization setting
  - User Authentication, Read-Interval for updating visualizations, etc.
  - OSU INAM Web configuration file is located at /etc/osu-inam.properties.

#### **Best Practice for OSC INAM and Thread Load balancing**

• What is the proper allocation of number of thread based on number of CPU cores for each module inside OSU INAM Daemon?

| Cluster size    | fabric discovery | performance<br>port inquiry | MPI_T and job<br>thread | Purge thread |
|-----------------|------------------|-----------------------------|-------------------------|--------------|
| < 500           | 2                | 2+                          | 1                       | 2            |
| 500< size <1000 | 4                | 8+                          | 1                       | 2            |
| > 1000          | 8                | 16+                         | 2                       | 2            |

• Load Balancing of Threads in Port Inquiry



(a) Timing of write phase for each (b) Timing of write phase for each (c) Timing of query phase for each (d) Timing of query phase for each thread for fabric discovery thread for port inquiry thread for port inquiry thread for fabric discovery

#### **Configuration for Profiling MPI Jobs**

- OSU INAM Daemon generates a configuration file that users need to pass to MPI jobs to send the information to OSU INAM Daemon
- This file is located in /tmp/mv2\_mpit.conf or \$OSU\_INAM\_INSTALL\_PREFIX/etc/osu-inam.conf
- Query intervals passed to the jobs are user-defined variables and can be edited in this file
- Unless specified they are system-defaults as specified in daemon config file

(base) kousha.2@head:~\$ cat inam.conf MV2\_TOOL\_QPN=190564 MV2\_TOOL\_LID=10 MV2\_TOOL\_COUNTER\_INTERVAL=30 MV2\_TOOL\_REPORT\_CPU\_UTIL=1 MV2\_TOOL\_REPORT\_MEM\_UTIL=1 MV2\_TOOL\_REPORT\_IO\_UTIL=1 MV2\_TOOL\_REPORT\_IO\_UTIL=1 MV2\_TOOL\_REPORT\_COMM\_GRID=1 MV2\_TOOL\_REPORT\_COMM\_GRID=1 MV2\_TOOL\_REPORT\_LUSTRE\_STATS=0 MV2\_TOOL\_REPORT\_PVARS=1 (base) kousha.2@head:~\$

# **Deciding on Intervals for Query Data from HPC Communication Stack**

- Deciding on Query interval is a key factor for profiling jobs
- Ideal Query interval should be the same across stacks
- Selecting Interval has a trade-off
  - Low intervals (1-5 s) of query result in massive amount of profiling data
    - Fine grained profiling
    - Results in slower database
    - Delay in web interface
  - Large (>60s) of query interval results in hiding concurrency bugs across layers or missing information
    - Faster web interface specifically for historical jobs
    - coarse-grained profiling
- You will explore this as part of Demo 0!



#### **Running Example MVAPICH2-X Jobs with INAM**

• MVAPICH2 Running Example without Performance Variables (PVARS)

\$ mpirun\_rsh -np 4 n0 n0 n1 n1 MV2\_ON\_DEMAND\_THRESHOLD=1
MV2\_TOOL\_INFO\_FILE\_PATH=inam.conf MV2\_TWO\_LEVEL\_COMM\_THRESHOLD=1
MV2\_USE\_RDMA\_CM=0 ./test

MVAPICH2 Running Example with Performance Variables (PVARS)

\$ mpirun\_rsh -rsh -np 4 n0 n0 n1 n1 MV2\_ON\_DEMAND\_THRESHOLD=1
MV2\_TOOL\_INFO\_FILE\_PATH=inam.conf MV2\_TWO\_LEVEL\_COMM\_THRESHOLD=1 MV2\_USE\_RDMA\_CM=0
MV2\_TOOL\_REPORT\_PVARS=1 MV2\_ENABLE\_PVAR\_TIMER=1 MV2\_ENABLE\_PVAR\_COUNTER=1
MV2\_ENABLE\_PVAR\_TIMER\_BUCKETS=1 MV2\_ENABLE\_PVAR\_COUNTER\_BUCKETS=1
MV2\_TOOL\_REPORT\_SESSIONS=1 MV2\_TOOL\_SESSIONS\_DEFAULT\_ALL\_HANDLES=1 ./test

#### Outline

- Overview
- OSU INAM Components
- Overhead Analysis
- Downloading, Installing, and Configuring OSU INAM
- Demos
  - Demo 1: Visualizing Entire Network
  - Demo 2: Visualizing Specific Job
  - Demo 3: Understanding Live Jobs on Cluster
  - Demo 4: Setting up Notifications
- Future Work and Concluding Remarks

# **Demo 1 - Visualizing Entire Network**

- Objectives
  - How to visualize network traffic at system level?
  - How to find congested links in a cluster?
  - How to view network traffic flow for historical jobs in a cluster?
- Tasks:
  - System level network traffic visualization
  - Viewing live node information
  - Finding nodes in the cluster
  - Finding most used links
  - Historical replay of network traffic

# **Demo 1: Visualizing Entire Network**

- Visualize the entire network
  - Represents a logical layout of the HPC network
  - Need not show how the compute nodes and switches are physically laid out on the data center
  - Use physics-based visualization package for better layout
  - Use advanced caching schemes to accelerate visualization time for very large clusters



The OSU InfiniBand Network Analysis and Monitoring tool - OSU INAM monitors IB clusters in real time by querying various subnet management entities in the network. It is also capable of interacting with the MVAPICH2-X software stack to gain insights into the communication pattern of the application and classify the data transferred into Point-to-Point, Collective and Remote Memory Access (RMA). INAM can also remotely monitor the CPU utilization of MPI processes in conjunction with MVAPICH2-X.



# **Demo 1: Visualizing Entire Network – System Level Visualization**

- Visualize the entire network
  - Contains two types of nodes
    - Switches (Red)
    - Compute Nodes (Blue)
  - Links represent physical IB connectivity between switches and compute nodes
  - Color on the links represent link usage
    - Grey (0% 5%)
    - Light Green (5% 25%)
    - Dark Green (25% 50%)
    - Orange (50% 75%)
    - Red (75% 100%)
  - Link utilization information is obtained in a live fashion from the fabric



# **Demo 1: Visualizing Entire Network - View Node Information**

- Users can interact with the different elements on the network view
- Hovering over the switches and nodes gives more information about specific element
- Provides information gathered and combined from job scheduler (e.g., Job ID) and IB fabric (e.g., link utilization)



# **Demo 1: Visualizing Entire Network – Finding Nodes**

- To find specific nodes, users can use search bar at bottom
- Start typing the node names, the interface provides suggestions on possible options using data from IB fabric
- Highlights desired node in live network view



# **Demo 1: Visualizing Entire Network – Finding Most-used Links**

- The IB links can be filtered by selecting/deselecting appropriate values in the "Link Usage" legend
- Provides an easy way to find out/visualize specific categories of links in the user interface

| Filter By Complete Network Vew Careford                                       | ng (i) Live View () Historical View | Link Usage<br>0% - 5% □ 5% - 25% □<br>25% - 50% □ 50% - 75% □<br>75% - 100% ♥  |
|---|-------------------------------------|--|
| € Usage Hints   |                                     | 95   |
| B   |                                     | 85   |
|   |                                     | Yei         55           10         50           10         45           40         40           33         35           30         25 |
| (f)<br>(d) (d) (d)  |                                     | (□) 10<br>(□) 10<br>(□) (□) 10<br>(□) 10<br>5  |
| You must click on Find Node to get the righ LID      Enter LID      Find Node | ıt result                           |  |

# **Demo 1: Visualizing Entire Network - Historical Replay**

- OSU INAM provides a live view as well as historical replay feature
- Users can select desired date and time ranges from the user interface
- Data is retrieved from the database and displayed as a moving timeline
- Time taken to render/fetch the data depends on the size of the cluster and amount of data requested
  - Start with a small interval like one hour to get an idea



## **Demo 2 - Visualizing Specific Job**

- Objectives
  - How to visualize a specific job's traffic in the cluster?
  - How to find most congested links in a job with respect to system level traffic?
  - How to understand job's topology to explain performance jitters?
- Tasks:
  - Job level network traffic visualization
  - Understanding job topology
  - Finding participating nodes in the job and accessing their information
  - Finding most used links in the job
  - Historical replay of job level network traffic

# **Demo 2: Visualizing Specific Job – Understanding Job Topology**

- Live network can be viewed at multiple granularities and in different ways
  - Entire network
  - Nodes and switches/links participating in a specific job
  - User specified nodes and the switches/links interconnecting them
- Job-level view can be obtained by filtering the network elements using the job ID
  - If users do not know the ID of the job, it can be obtained from the "Live Jobs" page
  - More details about the "Live Jobs" page and associated details will be provided in Demo 3



# **Demo 2: Visualizing Specific Job – Finding Node Information**

- Users can interact with different elements in the Job-Level view of the live network page
- Hovering over the node provides details of the node using data collected from the job scheduler (e.g., Job ID) and the IB fabric (e.g., LID, GID, HCA name etc)

| Filter By                    |                | O Live View | <ul> <li>Historical View</li> </ul>   | L'ink Usa | age            |  |         |
|------------------------------|----------------|-------------|---------------------------------------|-----------|----------------|--|---------|
| Job Id ~ 729984              | C Rendering    |             |                                       |           | ■ 0% - 5% 🗹    | 5%   | - 25% 🔽 |
| Network Metrics Max [Xmit Da | ta/Rcv Data] V |             |                                       | j         | ■ 75% - 100% 🗹 | _ 30%  |         |
| O Usage Hints                |                |             |                                       |           |                |  |         |
|                              |                |             |                                       |           |                | 95 -   |         |
|                              |                |             | GUID: 0x0002c903000a                  | 8fdd      |                | 90 -   |         |
|                              |                |             |                                       |           |                |  |         |
|                              | 154            | 90          | NAME: node143 HCA-1<br>Job ID: 729984 |           |                | 85 -   |         |
|                              | 34 154         |             | NAME: node143 HCA-1<br>Job ID: 729984 |           |                | 85   |         |
|                              | 39<br>32       |             | NAME: node143 HCA-1<br>Job ID: 729984 |           |                | 85   |         |
|                              | 8              |             | NAME: node143 HCA-1<br>Job ID: 729984 |           |                | 85   |         |
|                              | 8              |             | NAME: node143 HCA-1<br>Job ID: 729984 |           |                | 85   |         |
|                              | 8              |             | 102                                   |           |                | 85   |         |
|                              | 8              |             | NAME: node143 HCA-1<br>Job ID: 729984 |           |                | 25   |         |
|                              |                |             | NAME: node143 HCA-1<br>Job ID: 729984 |           |                | 85   |         |
|                              |                |             | NAME: node143 HCA-1<br>Job ID: 729984 |           |                | 85   |         |
|                              |                |             | NAME: node143 HCA-1<br>Job ID: 729984 |           |                | 85   |         |
|                              |                |             | NAME: node143 HCA-1<br>Job ID: 729984 |           |                | 85   |         |
|                              |                |             | NAME: node143 HCA-1<br>Job ID: 729984 |           |                | 85 -<br>80 -<br>75 -<br>70 -<br>70 -<br>65 -<br>60 -<br>55 -<br>50 -<br>50 -<br>50 -<br>45 -<br>30 -<br>25 -<br>20 - |         |

# **Demo 2: Visualizing Jobs- Historical Replay**

- OSU INAM provides historical replay feature at the job-level filter of the live network view
- Users can select desired date and time ranges from the user interface
- Data is retrieved from the database and displayed as a moving timeline
- Time taken to render/fetch the data depends on the size of the job and amount of data requested
  - Start with a small interval like one hour to get an idea



# **Demo 2: Visualizing Specific Job – Live Node View**



- From the Live Network or Jobs page, users can interact with different elements (switches and nodes) to get more information about them by double clicking on the element
- Double clicking on the node takes one to the live "Node-Level" view
- Uses information from job scheduler, MPI library, and IB fabric as available

# **Demo 3 - Understanding Live Job on Cluster**

- Objectives
  - How to access running jobs statistics in an easy manner?
  - How to sort live jobs on a cluster based on different metrics?
  - How much traffic does each job generate at system-level?
  - How to view CPU usage, memory usage, MPI information usage and IB level counters for live jobs running in the cluster?
- Tasks:
  - Viewing the list of running jobs
  - Filtering jobs and sorting jobs based on different metrics

40

# **Demo 3: Viewing the list of running jobs**

| OSU INAM         | Home Networ  | k View Historical Graph                                     | Live Jo                           | bbs D bug - Notifica  | tions User Guide   | DB Size Calculate | or                 |                  |         |
|------------------|--|---|-----------------------------------|-----------------------|--------------------|-------------------|--------------------|------------------|---------|
| • Note<br>• Plea | e that all counter valu<br>ase refer to the OSU IN | es in this page are instanta<br>IAM userguide for more deta | neous values.<br>ails about indiv | ridual counters       |                    |                   |                    |                  |         |
|                  |  |   |                                   |                       |                    |                   | Search             | n E              | C III - |
| Job ID 🍦         | CPU User Usage 🍦                                   | Virtual Memory Size 👙                                       | Total IO 👙                        | Total Communication 👙 | Total Inter Node 👙 | Node Count 👙      | Total Intra Node 👙 | Total Collective | RMA Se  |
| 729984           | 0  | 59.99 MB  | 0.00 bytes                        | 4.26 TB               | 4.26 TB            | 16                | 63.09 MB           | 11.36 GB         | 31.31 A |
| 729985           | 0  | 59.51 MB  | 14.27 MB                          | 706.79 GB             | 706.75 GB          | 16                | 41.78 MB           | 1.88 GB          | 20.79 ٨ |
| 729986           | 0  | 60.00 MB  | 0.00 bytes                        | 3.40 TB               | 3.40 TB            | 16                | 50.32 MB           | 9.06 GB          | 24.98 / |
| -                |  |   |                                   |                       |                    |                   |                    |                  |         |

Showing 1 to 3 of 3 rows

- Users can click on the "Live Jobs" tab to obtain list of all the jobs running currently on the cluster
- Uses information from job scheduler, MPI library, and IB fabric as available

# **Demo 3: Filtering jobs and sorting jobs based on different metrics**

| SU INAM   | Home Netwo             | ork view Historical Grap      | hs - Live Job     | os Debug - Notif    | ications User Guide | DB Size Calculat | or                         |                          |      |
|-----------|------------------------|-------------------------------|-------------------|---------------------|---------------------|------------------|----------------------------|--------------------------|------|
| - Not     | e that all counter val | uss in this page are instants |                   |                     |                     |                  |                            |                          |      |
| • Plea    | ase refer to the OSU   | INAM userguide for more det   | ails about indivi | dual counters       |                     |                  |                            |                          |      |
|           |                        |                               |                   |                     |                     |                  | Search                     | S                        | ₩.   |
|           | CD1111                 |                               |                   | <b>T</b> . 10       |                     | N-1-6            |                            | V Job ID                 | Colu |
| ¢ UIC     | CPU User Usage         | ■ VIITUAI Memory Size =       | Total IO 🏺        | Total Communication | ☐ I otal inter Node | Node Count 🏺     | Iotal Intra Node 🗧   Iotal | CPU User                 |      |
| 729984    | 0                      | 60.05 MB                      | 0.00 bytes        | 4.31 TB             | 4.31 TB             | 16               | 15.96 MB                   | System                   |      |
| 29985     | 0                      | 58.58 MB                      | 2.13 MB           | 1.00 TB             | 1.00 TB             | 16               | 29.70 MB                   | Usage                    |      |
| 729986    | 0                      | 57.71 MB                      | 0.00 bytes        | 3.28 TB             | 3.28 TB             | 16               | 12.09 MB                   | 🗌 Idle Time              |      |
|           |                        |                               |                   |                     |                     |                  |                            | ✓ Virtual<br>Memory Size |      |
| /ing 1 to | 3 of 3 rows            |                               |                   |                     |                     |                  |                            | 🗹 Total IO               |      |
|           |                        |                               |                   |                     |                     |                  |                            | 🗌 IO Read                |      |
|           |                        |                               |                   |                     |                     |                  |                            |                          |      |

- Different metrics are available, and users can select which metrics to be listed
- Ability to sort table in ascending or descending order based on different metrics

#### **Demo 3: Live Job View – Job Info and CPU/Mem Usage**



• Users can go to any level of granularity and visualize elements relevant to that level

#### **Demo 3: Live Job View - Vbuf Usage for UD and RC**

- Users can go to any level of granularity and visualize elements relevant to that level
  - Shows MPI internal communication buffer usage for different transport protocols
    - Supports Reliable Connected (RC) and Unreliable Datagram (UD)



# **Demo 3: Live Job View - Job level Counters**

- Provides counters collated from MPI runtime as well as IB fabrics at different granularities
  - Job Level
  - Node Level
  - Process Level
- Users can select the type of counter from the dropdown box



#### **Demo 3: Live Job View - Node Level Counters**

- Provides counters collated from MPI runtime as well as IB fabrics at different granularities
  - Job Level
  - Node Level
  - Process Level
- Users can select the type of counter from the dropdown box



# **Demo 3: Live Job View - MPI Rank level Counters**

- Provides counters collated from MPI runtime as well as IB fabrics at different granularities
  - Job Level
  - Node Level
  - Process Level
- Users can select the type of counter from the dropdown box



# **Demo 3: Live Job View – InfiniBand Level Counters**

- The counters are applicable to the switches as well as the compute nodes
- Only IB level counters
   provided at
   switch level

| ort Counte        | ers                                   |                            |                    |             |                      |          |   |           |   |       |          |
|-------------------|---------------------------------------|----------------------------|--------------------|-------------|----------------------|----------|---|-----------|---|-------|----------|
| Port Counters     |                                       |                            |                    |             |                      |          |   |           |   |       |          |
|                   |                                       |                            |                    |             |                      |          |   |           |   |       |          |
| lode              |                                       |                            |                    |             |                      |          |   |           |   |       |          |
| √ node008 (0x0002 | 2c90300                               | 0a852d)                    |                    |             |                      |          |   |           |   |       |          |
| node009 (0x0002   | 2c90300                               | 0a8439)                    |                    |             |                      |          |   |           |   |       |          |
| node012 (0x0002   | 2c90300                               | 0a82cd)                    |                    |             |                      |          |   |           |   |       |          |
| noaco 12 (0x000)  | 2_90300                               | (Dat-rad)                  |                    |             |                      |          |   |           |   |       |          |
| node014 (0x0002   | 2c90300                               | 0a84b5)                    |                    |             |                      |          |   |           |   |       |          |
| node015 (0x0002   | 2090300                               | 0a8549)                    |                    |             |                      |          |   |           |   |       |          |
| node016 (0x0002   | 200300                                | 0a82c5)                    |                    |             |                      |          |   |           |   |       |          |
| node133 (0x0002   | 2 < 90300                             | 0a83cd)                    |                    |             |                      |          |   |           |   |       |          |
| node135 (0x0002   | 2c90300                               | 0a9259)                    |                    |             |                      |          |   |           |   |       |          |
| node136 (0x0002   | 2c90300                               | 0a9389)                    |                    |             |                      |          |   |           |   |       |          |
| node139 (0x0002   | 2c90300                               | 0a8dc1)                    |                    |             |                      |          |   |           |   |       |          |
| node143 (0x0002   | 2c90300                               | 0a8fdd)                    |                    |             |                      |          |   |           |   |       |          |
| node144 (0x0002   | 2c90300                               | 0a8f2d)                    |                    |             |                      |          |   |           |   |       |          |
| node147 (0x0002   | 2c90300                               | 0a927d)                    |                    |             |                      |          |   |           |   |       |          |
| node148 (0x0002   | 2c90300                               | 0a8c19)<br>8 HCA-11        |                    |             |                      |          |   |           |   |       |          |
|                   | Jueoo                                 | oncarij                    |                    |             |                      |          |   |           |   |       |          |
| D. in the set     | o collor                              | stad from the quitch Con-  | d and Poor hore as | ra from the | a parapactive of the | owitch   |   |           |   |       |          |
| Metric:           | e coller                              | .ted from the switch. Sent | and Nety here a    | e nom die   | s perspective of the | Metric:  |   |           |   |       |          |
| Xmit Data         |                                       | ~                          |                    |             |                      | Rcv Data |   | ~         |   |       |          |
| 750 GB            |                                       | Aggregate                  |                    | Delta 🚽     | 150 GB               | 750 GB   | - | Aggregate |   | Delta | 150 GB   |
| 700 GD            |                                       |                            | l                  | 1           | - 100 CD             | 700 GB   |   |           |   |       | 140 CD   |
| 700 GB            |                                       |                            |                    | /           | 140 GB               | 700 GB   |   |           |   | /     | 140 GB   |
| 650 GB            |                                       |                            | 1                  | 1           | 130 GB               | 650 GB - |   |           |   | 1     | 130 GB   |
| 600 GB -          |                                       |                            | /                  |             | - 120 GB             | 600 GB - |   |           |   | /     | - 120 GB |
| 550 GB            | · · · · · · · · · · · · · · · · · · · |                            |                    |             | 110 GB               | 550 GB   |   |           |   |       | 110 GB   |
|                   |                                       |                            |                    |             | 100.00               | 500.00   |   |           | / |       | 100.00   |

### **Demo 4 - Performance Notifications**

- Objectives:
  - Understand how to set up customized notifications
  - Understand how to manage and delete notifications
  - Understand how to setup recurring and non-recurring notifications
- Tasks:
  - Accessing Notifications
  - Creating Notifications
  - Managing notifications
  - Notifications occurrence view

49

# **Demo 4: Setting up Notifications – Notifications**

- Choose Notifications tab from tabs
- This tab shows the notifications system
- The top half shows the notifications that have already been created/is in effect
- The bottom half allows to create a new notification
- Choose to add a new notification criteria

| otifications      |             |         |             |                  |                 |        |                        |                    |
|-------------------|-------------|---------|-------------|------------------|-----------------|--------|------------------------|--------------------|
|                   |             |         |             |                  |                 |        | Search                 | <i>C</i> <b>II</b> |
| Notification ID   | Criteria ID | ÷ Categ | ory 🔶 Metri | c 🔶 Condition    | Threshold Value | † Time | Additional Information | Action             |
|                   |             |         |             | No matching reco | rds found       |        |                        |                    |
|                   |             |         |             |                  |                 |        |                        |                    |
| otification Crite | ria         |         |             |                  |                 |        |                        |                    |
|                   |             |         |             |                  |                 |        | Add Notific            |                    |
|                   |             |         |             |                  |                 |        |                        | ation Criteria     |
|                   |             |         |             |                  |                 |        |                        | ation Criteria     |
|                   |             |         |             |                  |                 | Searc  | h O 2                  | ation Criteria     |
| Criteria ID       | Category    | 🔶 Metr  | ic 🔶 Co     | mparison         | Threshold Value | Searc  | h 🖸 💭                  | ation Criteria     |

# **Demo 4: Setting up Notifications – Adding Notifications**

- Upon Click a new window should appear for adding notification
- Explore the available options as this part is under development
- Request for feedback: What type of notifications are you interested to have support here?
- Future Work: Ability to create multi-variable notifications.

| otifications & Criteria |                 |   |     |                 |          |                 |
|-------------------------|-----------------|---|-----|-----------------|----------|-----------------|
| Notifications           | Category        | Process Counters  | •   |                 |          |                 |
|                         | Metric          | Bytes sent  | •   | Search          |          | a               |
|                         | Comparison      | Greater than  | •   | Search          |          |                 |
| Notification ID         | Threshold value |   | 5   | Additional Info | ormation | Action          |
|                         |                 | For percentage, enter values between 0 and 100  | - 1 |                 |          |                 |
| Notification Criteria   | Is Recurring    | ○ Yes ○ No  | - 1 |                 |          |                 |
|                         | Email List      |   |     |                 |          | den Celheria    |
|                         |                 |   |     |                 |          | tion Criteria 4 |
|                         |                 |   | c   |                 | 0 0      |                 |
| Criteria ID Category    |                 | Enter email recipients separated by comma.<br>E.g. mvapich-help@cse.ohio-state.edu,panda@cse.ohio-<br>state.edu | 1   | s Recurring     |          | Action          |
|                         | Email Subject   | OSU INAM Notification   |     |                 |          |                 |
|                         | Email Prologue  | Hello,  |     |                 |          |                 |
|                         |                 | Text that appears at the start of the notification email  | li  |                 |          |                 |
|                         | Email Epilogue  | -OSU INAM   |     |                 |          |                 |
|                         |                 | Text that appears at the end of the polification email  | lie |                 |          |                 |
|                         |                 | ross ons appears as one end of one normcation email   |     |                 |          |                 |
|                         |                 |   |     |                 |          |                 |

# **Demo 4: Setting up Notifications – Setting Up Notifications**

- Try to add a new notification by choosing your favorite values
- Two modes support Recurring and nonrecurring
- The text for the message and subject could be customized
- Click save after you finish

| OSU INAM Home Network View Histori | ical Graphs - Live Jobs | Debug - Notifications User Guide DB Size Ca              | alculator |                               |
|------------------------------------|-------------------------|--|-----------|-------------------------------|
|                                    | Notification Criteria   |  | ×         |                               |
| Notifications & Criteria           |                         |  |           |                               |
|                                    |                         |  |           |                               |
| Notifications                      | Category                | Process Counters   | •         |                               |
|                                    | Metric                  | Packets sent   | •         | Search C III -                |
| Notification ID                    | Comparison              | Less than  | •         | Additional Information Action |
|                                    | Threshold value         | 15000  |           |                               |
|                                    |                         | For percentage, enter values between 0 and 100           |           |                               |
| Notification Criteria              | Is Recurring            | Yes ○ No   | - 1       |                               |
|                                    | Email List              | Kousha.2@osu.edu   |           | Add Natification Critoria     |
|                                    |                         |  |           |                               |
|                                    |                         |  |           | h 🖸 🗭 🔳 🏭 -                   |
| Criteria ID Category               |                         | Enter email recipients separated by comma.               | 11.       | Is Recurring                  |
|                                    |                         | E.g. mvapich-help@cse.ohio-state.edu,panda@cse.ohio-     | - 1       |                               |
|                                    |                         | state.edu  |           |                               |
|                                    | Email Subject           | OSU INAM Notification                                    |           |                               |
|                                    | Email Prologue          | Hello,   |           |                               |
|                                    |                         |  |           |                               |
|                                    |                         |  |           |                               |
|                                    |                         | Text that appears at the start of the notification email |           |                               |
|                                    | Email Epilogue          | -OSU INAM  |           |                               |
|                                    |                         |  |           |                               |
|                                    |                         |  |           |                               |
|                                    |                         |  | 11.       |                               |
|                                    |                         | Text that appears at the end of the notification email   |           |                               |
|                                    |                         |  |           |                               |
|                                    |                         |  |           |                               |
|                                    |                         | Save   | lose      |                               |
|                                    |                         |  |           |                               |

# **Demo 4: Setting up Notifications – Managing Notifications**

- Your saved notification should show in Notification page
- You can edit or delete the notification and view the details of it
- Request for feedback: What would you suggest to improve managing notifications section?

|                                      |  |  |   |  | Search  | C III.                  |
|--------------------------------------|--|--|---|--|---|-------------------------|
| Notification ID                      | ♦ Criteria ID  | Category 🔶 Metric  | Condition     Three                               | shold Value 🔶 Time                               | Additional Informat                                   | ion Action              |
|                                      |  |  | No matching records foun                          | d  |   |                         |
|                                      |  |  |   |  |   |                         |
|                                      |  |  |   |  |   |                         |
| otification Crit                     | eria   |  |   |  | Add   | Notification Criteria 🕇 |
| otification Crit<br>Criteria ID      | eria   | 🔶 Metric 🔤   | Comparison  | Se Threshold Value                               | Add arch  | Notification Criteria + |
| otification Crit<br>Criteria ID<br>1 | eria Category Process Counters   | <ul> <li>Metric</li> <li>Packets sent</li> </ul>                         | <b>Comparison</b><br>Less than                    | Se<br>Threshold Value<br>15000                   | Add arch C  | Notification Criteria + |
| Criteria ID<br>1                     | eria Category Process Counters Process Counters                              | <ul> <li>Metric</li> <li>Packets sent</li> <li>Bytes received</li> </ul> | Comparison<br>Less than<br>Less than              | Se<br>Threshold Value<br>15000<br>15000          | Add<br>arch C<br>b Is Recurring<br>c<br>x             | Notification Criteria + |
| Criteria ID<br>1<br>2<br>3           | eria<br>Category<br>Process Counters<br>Process Counters<br>Process Counters | Metric       Packets sent       Bytes received       Bytes received      | Comparison<br>Less than<br>Less than<br>Less than | Se<br>Threshold Value<br>15000<br>15000<br>15000 | Add<br>arch C<br>b Is Recurring<br>Add<br>X<br>X<br>X | Notification Criteria + |

Notification

# **Demo 4: Setting up Notifications - Viewing Occurrences**

- Upon occurring of event that was specified in the notification you should see a history of them here
- You can click on additional information to see the details

| otifications   |  |   |  |           |  |                |                                  |  |  |                                     |         |                     |          |   |
|--|--|---|--|-----------|--|----------------|----------------------------------|--|--|-------------------------------------|---------|---------------------|----------|---|
|  |  |   |  |           |  |                |                                  |  |  | Search                              |         |                     | C        |   |
| Notification ID  | Criteria ID 🔶  | Category  | ≜ Metric ≜   | Condition | Threshold Value  | Time           | .≜<br>▼                          | Additiona                              | l Information                                  |                                     |         |                     |          | Actio   |
| 1  | 4  | Port<br>Counters  | Link recovers  | Less than | 15000  | 18 Jul 22:15:2 | 2021                             | (0x0002c9)<br>(0x0002c9)<br>(0x0002c9) | tch:MTS3610/I<br>0200424408):1<br>03000a857d), | .09/U1<br>8 - node027<br>MF0;ibswit | HCA-1   | >                   |          | 面   |
| 2  | 5  | Port  | Link utilization   | Greater   | 90   | 18 Jul 22:16:0 | 2021<br>07                       | MF0;ibswit                             | tch:MTS3610/I<br>020040e518):2                 | .05/U1<br>4 -                       | 1026902 | 004247              | 0        | 曲   |
| nowing 1 to 2 of 2   | rows   | Counters  | (percentage)   | (Fair     |  |                |                                  | MFU; IDSWII                            |  |                                     | Add Not | tificatio           | on Crit  | teria •   |
| nowing 1 to 2 of 2   | rows   | Counters  | (percentage)   |           |  |                |                                  | MFU; IDSWI                             | Search   |                                     | Add Not | cificatio<br>S      | on Crit  | teria ·   |
| howing 1 to 2 of 2<br>otification Crite<br>Criteria ID   | rows<br>rria<br>Category   | ¢ A   | Aetric   |           | Comparison   | \$             | Thres                            | hold Value                             | Search   | Is Recurri                          | Add Not | tificatio<br>G<br>¢ | on Crit  | teria •<br>III •                                      |
| howing 1 to 2 of 2<br>otification Crite<br>Criteria ID   | rows eria Category Process Counters  | ⇒ A<br>s F  | Aetric<br>Vackets sent   |           | Comparison     Less than   | \$             | Thres<br>15000                   | hold Value                             | Search   | Is Recurri                          | Add Not | tificatio<br>G<br>¢ | on Crit  | teria •<br>III →<br>ction<br>; m                      |
| <ul> <li>howing 1 to 2 of 2</li> <li>otification Crite</li> <li>Criteria ID</li> <li>1</li> <li>2</li> </ul>                       | rows eria Category Process Counters Process Counters                                   | ↓ A<br>s F F F F F F F F F F F F F F F F F F F  | Aetric<br>Vackets sent<br>bytes received   |           | <ul> <li>♦ Comparison</li> <li>Less than</li> <li>Less than</li> </ul>                                     | \$             | Thres<br>15000<br>15000          | hold Value                             | Search   | Is Recurri                          | Add Not | cificatio           | on Crit  | teria •<br>III •<br>ction<br>i m<br>i m               |
| <ul> <li>howing 1 to 2 of 2</li> <li>otification Crite</li> <li>Criteria ID</li> <li>1</li> <li>2</li> <li>3</li> </ul>            | rows eria Category Process Counters Process Counters Process Counters                  | <ul> <li>Counters</li> <li>A</li> <li>S</li> <li>F</li> <li>S</li> <li>E</li> <li>S</li> <li>E</li> </ul> | Aetric<br>Packets sent<br>hytes received<br>hytes received   |           | Comparison       Less than       Less than       Less than       Less than                                 | \$             | Thres<br>15000<br>15000          | hold Value                             | Search   | Is Recurri                          | Add Not | tificatio<br>♀      | n Crit   | teria •<br>III •<br>ction<br>i m<br>i m<br>i m<br>i m |
| <ul> <li>howing 1 to 2 of 2</li> <li>otification Crite</li> <li>Criteria ID</li> <li>1</li> <li>2</li> <li>3</li> <li>4</li> </ul> | rows eria Category Process Counters Process Counters Process Counters Process Counters | ⇔ A<br>s F<br>s E<br>s E  | Aetric<br>Aetric<br>Arackets sent<br>Attric Arackets sent<br>Arackets received<br>Arackets received<br>Arackets received<br>Arackets received<br>Arackets sent |           | <ul> <li>Comparison</li> <li>Less than</li> <li>Less than</li> <li>Less than</li> <li>Less than</li> </ul> | \$             | Thres<br>15000<br>15000<br>15000 | hold Value                             | Search   | Is Recurri<br>×<br>×<br>×           | Add Nor | tificatic<br>C<br>¢ | Acc<br>C | teria -<br>III -<br>ction<br>; m<br>; m<br>; m        |

# **Demo 4: Setting up Notifications – Viewing Occurrences' Details**

- After clicking you can see the details of node GUID and names with their ports for which the notification has happened
- Request for feedback: What would you suggest to improve this notifications section?

|                   | ne Network View H          | iistorical Graphs – Live Jobs Debug    | - Notifications Use         | r Guide DB Size Calcul                | ator      |                |            |
|-------------------|----------------------------|--|-----------------------------|---------------------------------------|-----------|----------------|------------|
|                   | Additional II              | nformation                             |                             |                                       | ×         |                |            |
| tifications & (   | Criteria<br>Notification I | D: 1                                   |                             |                                       |           |                |            |
| Notifications     |                            |  |                             |                                       |           |                |            |
|                   | Additional Inf             | formation                              |                             |                                       |           |                | ;<br>₩-    |
|                   | MF0;ibswitch:              | MTS3610/L09/U1 (0x0002c90200424408):1  | 8 - node027 HCA-1 (0x0002c  | 903000a857d)                          |           |                |            |
| Notification II   | MF0;ibswitch:              | MTS3610/L01/U1 (0x0002c90200424420):2- | 4 - MF0;ibswitch:MTS3610/S  | 03/U1 (0x0002c902004247(              | 00):2     |                | Action     |
| 1                 | 4 MF0;ibswitch:            | MTS3610/L02/U1 (0x0002c9020040d970):1  | 0 - storage02 HCA-1 (0x0002 | c903000aae37)                         |           | -1             | Ê          |
|                   | MF0;ibswitch:              | MTS3610/L10/U1 (0x0002c90200424410):3  | 0 - MF0;ibswitch:MTS3610/S  | 06/U1 (0x0002c902004246               | f8):20    |                |            |
| 2                 | 5 MF0;ibswitch:            | MTS3610/L05/U1 (0x0002c9020040e518):3- | 4 - MF0;ibswitch:MTS3610/S  | 08/U1 (0x0002c902004246               | d0):9     |                | 龠          |
| -                 | MF0;ibswitch:              | MTS3610/L09/U1 (0x0002c90200424408):2  | 5 - MF0;ibswitch:MTS3610/S  | 04/U1 (0x0002c9020042466              | :0):17    | 020042470      |            |
| Showing 1 to 2 of | of 2 rows MF0;ibswitch:    | MTS3610/S07/U1 (0x0002c902004246f0):16 | 5 - MF0;ibswitch:MTS3610/L0 | 08/U1 (0x0002c9020042441              | 8):32     |                |            |
|                   | MF0;ibswitch:              | MTS3610/L06/U1 (0x0002c9020040d990):3  | 0 - MF0;ibswitch:MTS3610/S  | 06/U1 (0x0002c902004246               | f8):12    |                |            |
| Notification C    | MF0;ibswitch:              | MTS3610/L03/U1 (0x0002c9020040e3d8):8  | - storage16 HCA-1 (0x0002c  | 903000a8cf9)                          |           |                |            |
| Notification c    | MF0;ibswitch:              | MTS3610/L01/U1 (0x0002c90200424420):3  | 2 - MF0;ibswitch:MTS3610/S  | 07/U1 (0x0002c <sup>001</sup> 004z40) | 0):2      |                | _          |
|                   | Showing 1 to 10            | 0 of 329 rows 10 × records per page    |                             | « 、 1 2                               | 3 4 5 , » | Notification C | criteria 🕂 |
| Criteria ID       | Cat                        |  |                             |                                       | Close     |                | Action     |
| 3                 | Process Counters           | Bytes received                         | Less than                   | 15000                                 | *         |                | c 💼        |
| 4                 | Port Counters              | Link recovers                          | Less than                   | 15000                                 | ×         |                | c 🖻        |
| 5                 | Port Counters              | Link utilization (percentage)          | Greater than                | 90                                    | ×         |                | <b>区</b> 前 |

#### Outline

- Overview
- OSU INAM Components
- Overhead Analysis
- Downloading, Installing, and Configuring OSU INAM
- Demos
- Concluding Remarks and Future Work

## **Concluding Remarks and Future Work**

- Presented details of OSU INAM tool capable of analyzing the communication traffic on the InfiniBand network with inputs from the MPI runtime. Major features include:
  - 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
  - low overhead (5%), scalable (4K processes on 256 nodes), and fined-granularity profiling of HPC communication stack
  - Visualize the data transfer happening in a "live" or historical fashion for
    - Entire Network, Particular Job One or multiple Nodes, One or multiple Switches
- The designs have been publicly released as a part of v0.9.6 and are being used by many users
  - Downloaded over 5000 times form the project website
  - Community engagement with
    - OSC @ USA, NOAA @ USA, U. of Utah @ USA, CAE Services @ Germany, Pratt & Whitney, Ghent University @ Germany, Cyfronet @ Poland, and Georgia Tech Univ @ USA
- As future work we aim to
  - Add support to profile and analyze GPU-based communication
  - Address the storage bottleneck by exploring high-performance database engines
  - Support more intra-node metrics and power consumption profiling

# **Funding Acknowledgments**

#### **Funding Support by**















# Acknowledgments to all the Heroes (Past/Current Students and Staffs)

| Current Studen  | nts (Graduate)   |   |  |  |                                     |   |   |             | Current Research Scientists   | Current Faculty   |
|---|--|---|--|--|-------------------------------------|---|---|-------------|---|---|
| – N. Alna   | aasan (Ph.D.) –  | K. S. Khor                                | assani (Ph.D.) -   | – A. H. Tu (F  | Ph.D.)                              | -   | H. Ahn (Ph.D.)  |             | – M. Abduljabbar  | <ul> <li>H. Subramoni</li> </ul>  |
| – Q. Antł   | hony (Ph.D.) –   | P. Kousha                                 | (Ph.D.) -  | – S. Xu (Ph.I  | D.)                                 | _   | G. Kuncham (Ph.D.)  | )           | – A. Shafi  | Current Software Engineers  |
| <ul> <li>CC. Cl</li> <li>N. Con</li> <li>A. Jain</li> </ul> | hun (Ph.D.) —<br>tini (Ph.D.) —<br>(Ph.D.) —<br><b>Idents</b>  | B. Michal<br>B. Rames<br>K. K. Sure       | sh (Ph.D.) -   | <ul> <li>Q. Zhou (F</li> <li>K. Al Attar</li> <li>L. Xu (Ph.I</li> </ul> | Ph.D.)<br><sup>-</sup> (M.S.<br>D.) |   | R. Vaidya (Ph.D.)<br>J. Yao (Ph.D)<br>M. Han (M.S.)<br>A. Guptha (M.S.)                               |             | Current Students (Undergrads)<br>— V. Shah<br>— T. Chen   | <ul> <li>B. Seeds</li> <li>N. Pavuk</li> <li>N. Shineman</li> <li>M. Lieber</li> </ul>  |
| -<br>-<br>-<br>-<br>-<br>-<br>-<br>-                        | <ul> <li>A. Awan (Ph.D.)</li> <li>A. Augustine (M.S.)</li> <li>P. Balaji (Ph.D.)</li> <li>M. Bayatpour (Ph.D.)</li> <li>R. Biswas (M.S.)</li> <li>S. Bhagvat (M.S.)</li> <li>S. Bhagvat (M.S.)</li> <li>A. Bhat (M.S.)</li> <li>D. Buntinas (Ph.D.)</li> <li>L. Chai (Ph.D.)</li> <li>B. Chandrasekharan (N</li> </ul> | -<br>-<br>-<br>-<br>-<br>-<br>-<br>-<br>- | T. Gangadharapp<br>K. Gopalakrishna<br>J. Hashmi (Ph.D.)<br>W. Huang (Ph.D.)<br>W. Jiang (M.S.)<br>J. Jose (Ph.D.)<br>M. Kedia (M.S.)<br>S. Kini (M.S.)<br>M. Koop (Ph.D.) | oa (M.S.)<br>in (M.S.)   |                                     | P. Lai (M.<br>J. Liu (Ph<br>M. Luo (F<br>A. Mamie<br>G. Marsh<br>V. Meshr<br>A. Mood<br>S. Naravu<br>R. Noron | S.)<br>.D.)<br>Ph.D.)<br>dala (Ph.D.)<br>(M.S.)<br>am (M.S.)<br>y (M.S.)<br>Ila (Ph.D.)<br>ha (Ph.D.) |             | D. Shankar (Ph.D.)<br>G. Santhanaraman (Ph.D.)<br>N. Sarkauskas (B.S. and M.S)<br>N. Senthil Kumar (M.S.)<br>A. Singh (Ph.D.)<br>J. Sridhar (M.S.)<br>S. Srivastava (M.S.)<br>S. Sur (Ph.D.)<br>H. Subramoni (Ph.D.)<br>K. Vaidyanathan (Ph.D.) | Current Research Specialist - R. Motlagh Past Research Scientists - K. Hamidouche - S. Sur - X. Lu Past Senior Research Associate - J. Hashmi |
| _<br>_<br>_<br>_  | S. Chakraborthy (Ph.<br>N. Dandapanthula (M.<br>V. Dhanraj (M.S.)<br>CH. Chu (Ph.D.)   | D.) –<br>.S.) –<br>–                      | R. Kumar (M.S.)<br>S. Krishnamoorth<br>K. Kandalla (Ph.D<br>M. Li (Ph.D.)  | )<br>ny (M.S.)<br>).)  | -<br>-<br>-                         | S. Pai (M<br>S. Potluri<br>K. Raj (M<br>R. Rajach   | g (Ph.D.)<br>(Ph.D.)<br>.S.)<br>andrasekar (Ph.D.)  | -<br>-<br>- | A. Vishnu (Ph.D.)<br>J. Wu (Ph.D.)<br>W. Yu (Ph.D.)<br>J. Zhang (Ph.D.)   | Past Programmers-A. Reifsteck-D. Bureddy-J. Perkins   |
| Past Pos<br>–<br>–<br>–                                     | <b>t-Docs</b><br>D. Banerjee<br>X. Besseron<br>M. S. Ghazimeersaeed  | – H<br>– J.I<br>– M.                      | W. Jin – E.<br>.in – K.<br>Luo – S.  | Mancini<br>Manian<br>Marcarelli  | -<br>-<br>-                         | A. Ruhe<br>J. Vienn<br>H. Wanş  | la<br>e<br>g  |             |   | Past Research Specialist<br>— M. Arnold<br>— J. Smith   |

# **Thank You!**

subramoni.1@osu.edu, kousha.2@buckeyemail.osu.edu, kolli.28@buckeyemail.osu.edu



Network-Based Computing Laboratory <u>http://nowlab.cse.ohio-state.edu/</u>



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



High-Performance Big Data

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



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