

2020 OFA Virtual Workshop

VISUALIZE AND ANALYZE YOUR NETWORK ACTIVITIES USING OSU INAM Hari Subramoni, Pouya Kousha, Kamal Raj Ganesh, Dhabaleswar K. (DK) Panda

The Ohio State University

OUTLINE

- Introduction & Motivation
- Design of OSU INAM
- Impact of Profiling on Application Performance
- Usage Scenarios & Demonstration
- Conclusions & Future Work

PROFILING TOOLS PERSPECTIVE AND 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?



Unified and

holistic view

for all users

PROFILING TOOLS PERSPECTIVE AND CHALLENGES (CONT.)

- Understanding the interaction between applications, MPI libraries, 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



BROAD CHALLENGE

There are tools to give insight into each layerThere is a gap though!

How can we design a tool that enables in-depth understanding of the communication traffic on the interconnect and GPU through tight integration with the MPI runtime and job scheduler?



OSU INAM - OFAW20 2020

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), and AWS EFA
- Support for multiple platforms
 - x86, OpenPOWER, ARM, Xeon-Phi, GPGPUs
- Started in 2001, first open-source version demonstrated at SC '02
- Supports the latest MPI-3.1 standard
- http://mvapich.cse.ohio-state.edu
- Additional optimized versions for different systems/environments:
 - MVAPICH2-X (Advanced MPI + PGAS), since 2011
 - MVAPICH2-GDR with support for NVIDIA GPGPUs, since 2014
 - 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 3,090 organizations in 89 countries
- More than 759,000 (> 0.75 million) downloads from the OSU site directly
- Empowering many TOP500 clusters (Nov '19 ranking)
 - 3rd, 10,649,600-core (Sunway TaihuLight) at NSC, Wuxi, China
 - 5th, 448, 448 cores (Frontera) at TACC
 - 8th, 391,680 cores (ABCI) in Japan
 - 14th, 570,020 cores (Nurion) in South Korea and many others
- Available with software stacks of many vendors and Linux Distros (RedHat, SuSE, OpenHPC, and Spack)
- Partner in the 5th ranked TACC Frontera system
- Empowering Top500 systems for more than 15 years

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
 - <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 processlevel 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 0.9.6 released

- 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

Introduction & Motivation

- Design of OSU INAM
 Profiler Interface for MPI+CUDA Communication
 Profiler Interface for MPI Data collection
- Impact of Profiling on Application Performance
- Usage Scenarios & Demonstration
- Conclusions & Future Work

OSU INAM FRAMEWORK



OSU INAM - OFAW20 2020

PROFILER INTERFACE FOR MPI+CUDA COMMUNICATION

Low overhead GPU profiler module consist of intra-node topology and metrics inquiry

Each node aggregates and sends the metrics to the OSU INAM daemon at user-defined interval

Startup Phase

- Each rank discovers the topology and updates shared region with rank and its device info.
- Local rank zero will setup and start a profiler thread on CPU to profile all GPUs on node
 - Happens when the list of GPUs that will be used are known
- Query Phase
 - Profiler thread periodically profiles all selected GPUs
- Exit Phase
 - Once the ranks stop using device, profiler thread will perform one last read and send data then exit.



Information Flow for Intra-node GPU Profiling

PROFILER INTERFACE FOR MPI DATA COLLECTION

Extending MPI_T Performance Variable (PVAR)

- 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, Most used MPI primitive, Total bytes sent for RMA operations
- For each collective and point-to-point operation every rank
 - Stores the total bytes sent to and received from every other rank
 - An array of start and end time-stamps
 - Selected algorithm for the communication
 - The number of times a particular algorithm/function was called
- PVAR information is only sent if the aggregated bytes sent for a particular MPI operation exceeds a user-specified threshold

OUTLINE

Introduction & Motivation

- Design of OSU INAM
- Impact of Profiling on Application Performance
- Usage Scenarios & Demonstration
- Conclusions & Future Work

OVERHEAD ANALYSIS

• Overhead of GPU Profiling – reading metrics per node

TIMING OF THE GPU PROFILER THREAD PHASES FOR EACH NODE. EACH NODE HAS FOUR GPUS

Metrics	Average	Min	Max	STDEV.p	
Startup phase	1.632 s	1.561 s	1.672 s	0.035 s	
CUDA context create	1.624 s	1.548 s	1.663 s	0.035 s	
Query phase <	2.33 ms	1.63 ms	208.03 ms	4.43 ms	
Exit phase	88 us	85 us	93 us	28 us	

• Overhead of PVAR Collection – per MPI rank

OVERHEAD OF COLLECTING PVAR DATA AT NANOSECOND GRANULARITY

Metrics	Average	Min	Max	STDDEV.p
Collecting PVAR	517.63 ns	140 ns	16,204 ns	305.91 ns

Designing a Profiling and Visualization Tool for Scalable and In-Depth Analysis of High-Performance GPU Clusters, P. Kousha, B. Ramesh, K. Kandadi Suresh, C. Chu, A. Jain, N. Sarkauskas, and DK Panda HiPC'19, Dec 2019

Very Low Overhead for Query Phase

PROFILING VARIATION AND SCALABILITY



Overall Time for Gathering GPU Metrics

- Scales linearly
- Time proportional to number of GPUs
- Metrics are gathered per node

Designing a Profiling and Visualization Tool for Scalable and In-Depth Analysis of High-Performance GPU Clusters, P. Kousha, B. Ramesh, K. Kandadi Suresh, C. Chu, A. Jain, N. Sarkauskas, and DK Panda HiPC'19, Dec 2019



Time/Query for Gathering GPU Metrics

- Average time per query is 2ms
- Timing is stable across thousands of queries

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
- Usage Scenarios & Demonstration
- Conclusions & Future Work

USAGE CASES

Absolute

279280

MF0;ibswitch:MTS3610/L11/U1 -- > node008 HCA-1





k Com	munication (Grid						
	0	1	2	3	4	5	6	7
0	758.78 MB	759.04 MB	758.78 MB	758.78 MB	758.87 MB	759.04 MB	758 78 MB	0.00 bytes
1	759.04 MB	758.78 MB	0.00 bytes	759.04 MB	758.78 MB	758.91 MB	759.04 MB	759.04 MB
2	0.00 bytes	758.78 MB	758.78 MB	759.04 MB	758.78 MB	758.78 MB	759.04 MB	759.01 MB
3	758.78 MB	758.97 MB	758.78 MB	0.00 bytes	758 78 MB	759.04 MB	758 78 MB	759.04 MB
4	758.78 MB	759.04 MB	758.81 MB	758.78 MB	759.04 MB	0.00 bytes	758 78 MB	758,78 MB
5	759.04 MB	758.78 MB	759.04 MB	758.78 MB	0.00 bytes	758.78 MB	758.91 MB	758.78 MB
6	758.78 MB	758 78 MB	758.78 MB	759.04 MB	758.78 MB	759.04 MB	0.00 bytes	759.04 MB
7	758 78 MB	0.00 bytes	758.84 MB	758 78 MB	759.04 MB	758.78 MB	758.78 MB	758.78 MB

Live MPI level communication for each rank on a link Intra-node Topology showing physical links between CPUs and GPUs

Rank level communication grid for an MPI_Allreduce Operation

Each element (i,j) in the grid represents amount data transferred from rank i to rank j



Network View with expanded and hidden modes showing Ohio Supercomputer Center (OSC) with 3 heterogeneous clusters all connected to the same InfiniBand Fabric (114 switches and 1,428 compute nodes connected through 3,402 links)

Network and Live Jobs View Generation Timing on OSC with 1K Jobs

View	Average	Min	Max	STDEV.p
Network View	196.15 ms	187 ms	206.09 ms	5.75 ms
Live Jobs View	18.17 ms	16 ms	20 ms	1 ms

Monitoring Jobs Based on Various Metrics

Job ID 🕴	CPU User Usage 🌖	Virtual Memory Size	Total Communication	Total Inter Node	Total Intra Node	Total Collective	RMA Sent
270747	99	8.19 Mb	92.35 Gb	36.69 Gb	55.66 Gb	64.46 Gb	0.00 bytes
270748	99	15.12 Mb	149.98 Gb	58.23 Gb	91.76 Gb	102.78 Gb	0.00 bytes
270749	99	30.39 Mb	151.23 Gb	58.35 Gb	92.88 Gb	100.34 Gb	0.00 bytes
270759	99	17.99 Mb	58.71 Gb	37.29 Gb	21.43 Gb	303.73 Kb	0.00 bytes
270765	99	9.42 Mb	32.52 Gb	23.19 Gb	9.33 Gb	0.00 bytes	0.00 bytes

Showing 1 to 5 of 5 rows

Profiling and Reporting Performance Metrics at Different Granularities





Each MPI Primitive & List Most used MPI Primitives inside the MPI job

MPI Primitives Information							
Session name: global +							
MPI Primitives: usage over time			0	MPI Primitives: most used	0		
MPI Primitives: XTop3 Metric: Bytes sent *				Granularity: Job • Metric: Bytes sent • Top: 5 •			
34 T MPI_Allgather(agg)		MPI_Allgather(delta) - 3.	8.4 T	# MPI Primitive Node PVAR Value			
32 T		MPI_Allreduce(delta) - 3. MPI_Alltoall(delta) - 3.	3.2 T	MPI_Alltoall Job level MV2_COLL_ALLTOALL_BYTES_SEND 33.881T			
28 T		2.	2.8 T	2 MPI_Allreduce Job level MV2_COLL_ALLREDUCE_BYTES_SEND 109.736K			
26 T 24 T		2.	2.6 T 2.4 T	3 MPI_Allgather Job level MV2_COLL_ALLGATHER_BYTES_SEND 6.152K			
22 T 22 O T		2.	2.2 T	4 MPI_Reduce Job level MV2_COLL_REDUCE_BYTES_SEND 56			
18 T 16 T 14 T 12 T 10 T 8 T			.8 T Detta .6 T .4 T .2 T .7				
6 T 4 T		60	500 G 100 G				
2 T		- 20	200 G				
10:55 11:00 Fri 15 May	11:05 11:10	11:15		Showing 1 to 4 of 4 rows			
				Legend: K - Kilo (10 ³) M - Mega (10 ⁶) G - Giga (10 ⁹) T - Tera (10 ¹²) P - Peta (10 ¹⁵)			

Lustre read/write traffic with an OSS

Lustre Communication Grid





Proposed Chart Displaying MPI_T PVAR

X-axis: Current time Y-axis: Number of bytes sent over the network to the OSU INAM daemon

Designing a Profiling and Visualization Tool for Scalable and In-Depth Analysis of High-Performance GPU Clusters, P. Kousha, B. Ramesh, K. Kandadi Suresh, C. Chu, A. Jain, N. Sarkauskas, and DK Panda HiPC'19, Dec 2019



Physical and Logical NVLink Metrics

X-axis: Time Y-axis: Bandwidth utilization for the link

USAGE SCENARIOS & DEMONSTRATION: ALLREDUCE RING-BASED ALGORITHM

- The links in the clockwise direction have relatively lower link utilization compared to the links in the counter-clockwise direction
- Ineffective use of bi-directional bandwidth
- For developer: the ring should use both directions



NVLink utilization GPU0 to GPU1

Designing a Profiling and Visualization Tool for Scalable and In-Depth Analysis of High-Performance GPU Clusters, P. Kousha, B. Ramesh, K. Kandadi Suresh, C. Chu, A. Jain, N. Sarkauskas, and DK Panda HiPC'19, Dec 2019





NVLink utilization GPU0 to GPU1

USAGE SCENARIO: TENSORFLOW

Running TensorFlow v1.12 with MVAPICH2 using Horovod with Resnet50 Model

- What is the impact of batch size on GPU communication?
 - Usually users are interested in images/sec
 - Useful to understand lower layer communication efficiency
- The smaller batch size result in lower link utilization and is less communication efficient



NVLink Metrics chart with a batch size of 2



NVLink Metrics chart with a batch size of 32

Designing a Profiling and Visualization Tool for Scalable and In-Depth Analysis of High-Performance GPU Clusters, P. Kousha, B. Ramesh, K. Kandadi Suresh, C. Chu, A. Jain, N. Sarkauskas, and DK Panda HiPC'19, Dec 2019

USAGE SCENARIO: TENSORFLOW (CONT.)

- What is the impact of batch size on MPI level?
 - The peak data (message size) transferred between ranks are the same, but showing different patterns between ranks
 - Horovod uses certain message sizes in the Allreduce operations and it depends on the Deep Learning model, batch size, GPU architecture, and other Deep learning parameters





PVAR Metrics chart with a batch size of 32

Designing a Profiling and Visualization Tool for Scalable and In-Depth Analysis of High-Performance GPU Clusters, P. Kousha, B. Ramesh, K. Kandadi Suresh, C. Chu, A. Jain, N. Sarkauskas, and DK Panda HiPC'19, Dec 2019

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 (v0.9.6) 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 Deep Learning frameworks

THANK YOU!

subramon@cse.ohio-state.edu, kousha.2@osu.edu, sankarapandiandayalaganeshr.1@osu.edu,

panda@cse.ohio-state.edu





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



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