# Performance of MPI Communications on Nurion Utilizing MVAPICH2-X with XPMEM

**Minsik Kim, Ph.D.**
Supercomputing Infrastructure Center, KISTI

The International Conference for
High Performance Computing, Networking,
Storage, and Analysis (SC'19), OSU Booth

KiSTi Korea Institute of
Science and Technology Information
www.kisti.re.kr

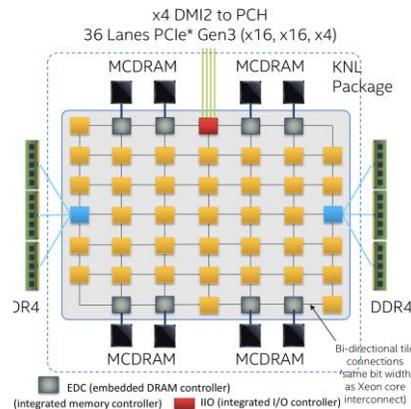# Introduction to KISTI-5 Supercomputer, Nurion

# KISTI-5 Compute Nodes

⚙️ The Largest KNL/OPA based commodity cluster System
Rpeak 25.7PFlops, Rmax 13.9PFlops

**Compute nodes** — 8,305 KNL Computing modules, 116 Racks, 25.3PF
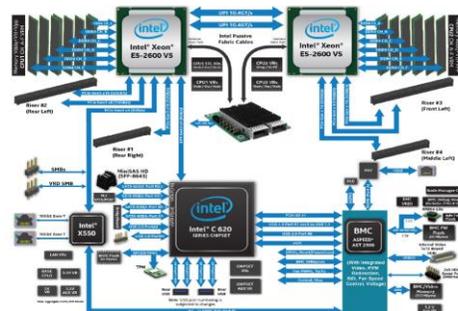
- ➢ 1x Xeon Phi KNL 7250, 68Cores 1.4GHz, AVX512
- ➢ 3TFlops Peak, ~0.2 Bytes/Flops,
- ➢ 96GB (6x16GB) DDR4-2400 6 channel RAM,
- ➢ 16GB HBM (460GB/s)
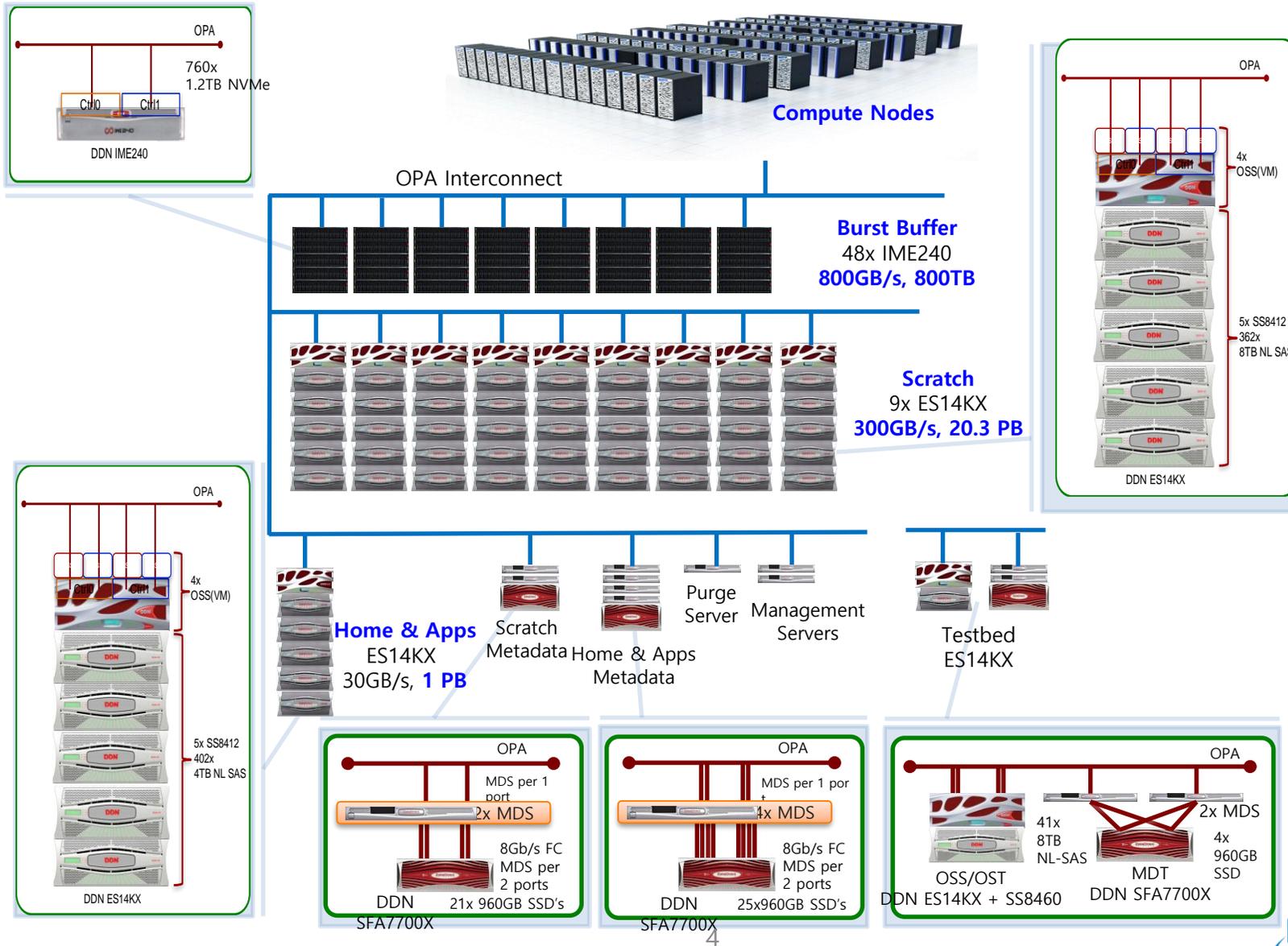- ➢ 1x 100Gbps OPA HFI, 1x On-board GigE Port

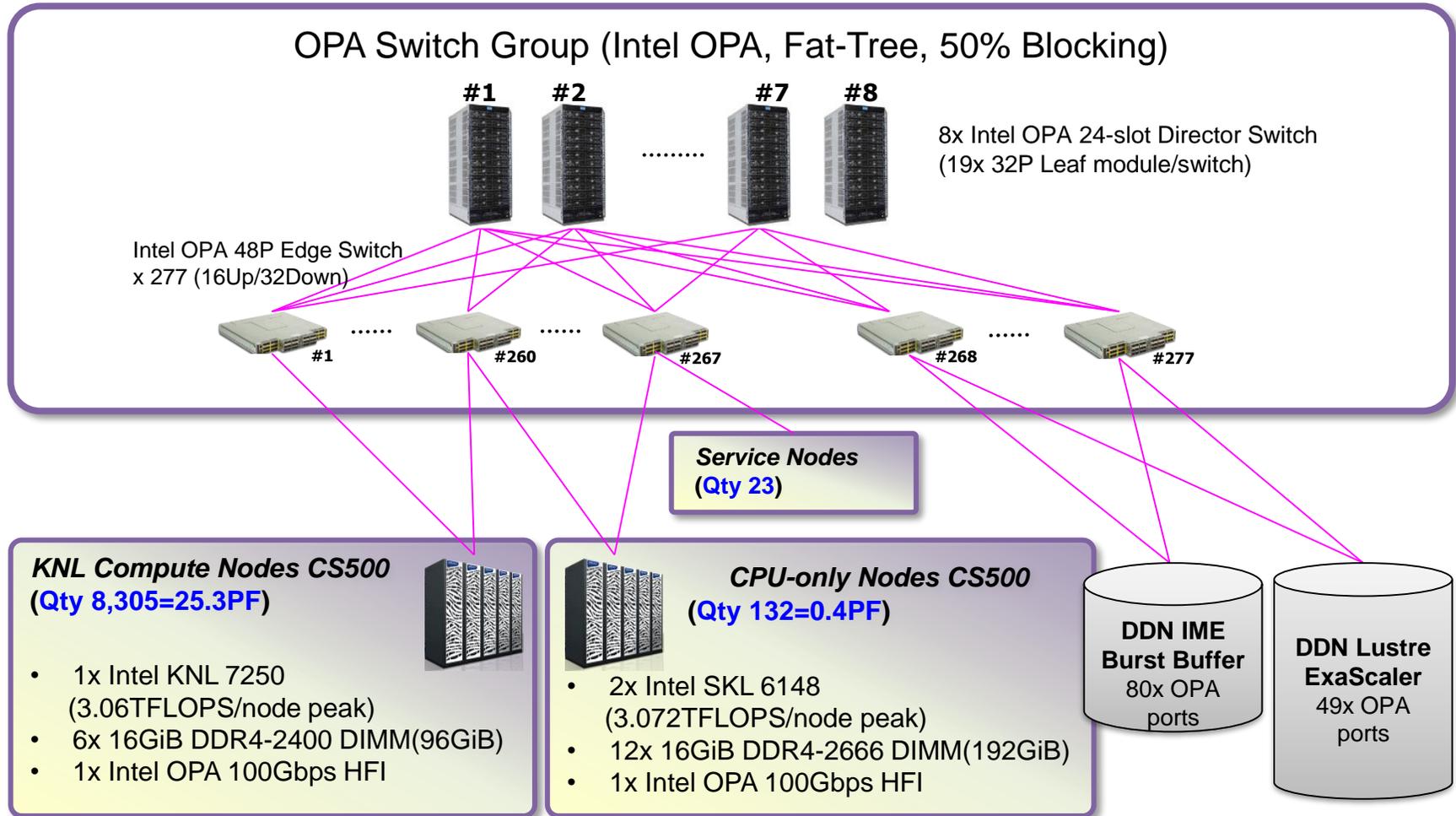**CPU-only nodes** — 132 Skylake Computing modules, 4 Racks, 0.4PF

- ➢ 2x Xeon SKX 6148 CPUs, 2.4GHz, AVX512
- ➢ 192GB (12x 16GB) DDR4-2666 RAM
- ➢ 1x Single-port 100Gbps OPA HFI card
- ➢ 1x On-board GigE (RJ45) port

# KISTI-5 Storage System



**Compute Nodes**

OPA Interconnect

**Burst Buffer**
48x IME240
**800GB/s, 800TB**

**Scratch**
9x ES14KX
**300GB/s, 20.3 PB**

OPA
760x
1.2TB NVMe

Ctrl0  Ctrl1

DDN IME240

OPA
4x
OSS(VM)

5x SS8412
362x
8TB NL SAS

DDN ES14KX

OPA
4x
OSS(VM)

5x SS8412
402x
4TB NL SAS

DDN ES14KX

**Home & Apps**
ES14KX
30GB/s, **1 PB**

Scratch
Metadata

Purge
Server

Management
Servers

Home & Apps
Metadata

Testbed
ES14KX

OPA
MDS per 1 port
2x MDS
8Gb/s FC
MDS per
2 ports
21x 960GB SSD's
DDN
SFA7700X

OPA
MDS per 1 por
+
4x MDS
8Gb/s FC
MDS per
2 ports
25x960GB SSD's
DDN
SFA7700X

OPA
41x
8TB
NL-SAS
2x MDS
4x
960GB
SSD
OSS/OST
DDN ES14KX + SS8460
MDT
DDN SFA7700X

4

# KISTI-5 OPA Interconnect

OPA Switch Group (Intel OPA, Fat-Tree, 50% Blocking)

**#1** **#2** **#7** **#8**

.........

8x Intel OPA 24-slot Director Switch
(19x 32P Leaf module/switch)

Intel OPA 48P Edge Switch
x 277 (16Up/32Down)

...... ...... ......

**#1** **#260** **#267** **#268** **#277**

*Service Nodes*
*(Qty 23)*

*KNL Compute Nodes CS500*
*(Qty 8,305=25.3PF)*

- 1x Intel KNL 7250
  (3.06TFLOPS/node peak)
- 6x 16GiB DDR4-2400 DIMM(96GiB)
- 1x Intel OPA 100Gbps HFI

*CPU-only Nodes CS500*
*(Qty 132=0.4PF)*

- 2x Intel SKL 6148
  (3.072TFLOPS/node peak)
- 12x 16GiB DDR4-2666 DIMM(192GiB)
- 1x Intel OPA 100Gbps HFI

**DDN IME
Burst Buffer**
80x OPA
ports

**DDN Lustre
ExaScaler**
49x OPA
ports

KiSTi
**Korea Institute of
Science and Technology Information**
www.kisti.re.kr

# Benchmark Performance Result

| Category | Features | # of nodes | Score | World Ranking |
|---|---|---|---|---|
| HPL | Large-scale Dense Matrix Computation<br>Used for Top500 | 8,174(KNL)<br>+ 122(SKX) | 13.93PF | 15th (Jun 2019) |
| HPCG | Large-scale Sparse Matrix Computation<br>Similar to normal user applications | 8,250(KNL) | 0.39PF | 8th (Jun 2019) |
| Graph500 | Breadth-First Search,<br>Single-Source Shortest Paths | 1,024(KNL) | 1,456GTEPS<br>337GTEPS | 10Th (Jun 2019)<br>3rd (Jun 2019) |
| IO500 | Various IO Workloads | 2,048(KNL) | 160.67 | 5th (Jun 2019) |

# MVAPICH2-X with XPMEM on Nurion Supercomputer

- Compute node
  - Manycore processor: Intel Xeon Phi KNL 7250 (68 cores)
  - MPI intra-node communication

- XPMEM (Cross-Process Memory Mapping)
  - Linux kernel module
  - Enables a process to map the memory of another process into its virtual address space

- OSU Micro-Benchmarks
  - Collective Communications: Latency
  - Point-to-point communications: Bandwidth

- MVAPICH2-X with XPMEM
  - GCC compiler 4.8.5 (Intel IFS 10.6)
  - Build issues on Intel compiler version

- Experimental environment
  - Intel MPI 19.0.4
  - MVAPICH 2.3
  - MVAPICH2-X 2.3rc2 with XPMEM

KiSTi Korea Institute of Science and Technology Information

# MVAPICH2-X with XPMEM: Single Node

KiSTi Korea Institute of Science and Technology Information
www.kisti.re.kr

# Collective Communication on Single Node: Latency (PPN=64)

# Collective Communication on Single Node: Latency (PPN=64)
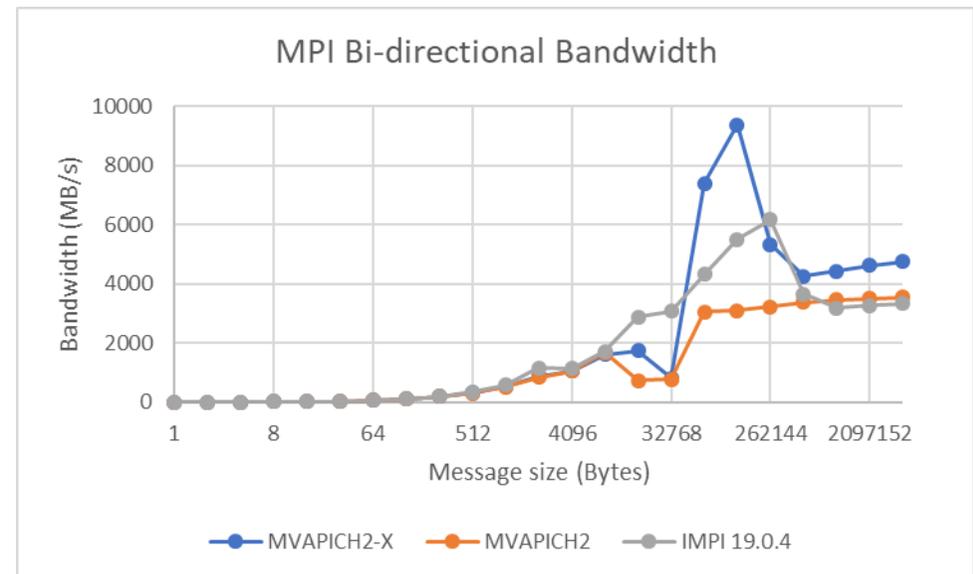


MPI_Allgather



MPI_Allreduce



MPI_Alltoall

- Experimental environment
  - Single KNL node (cache mode)
  - Processes per node = 64

- Performance evaluation
  - MVAPICH2-X better performance on collective communications compared to MVAPICH2
  - Intel MPI 19.0.4 better performance on MPI_Allreduce

KiSTi Korea Institute of Science and Technology Information
www.kisti.re.kr

# Point-to-Point Communication on Single Node: Bandwidth



MPI Latency



MPI Bandwidth
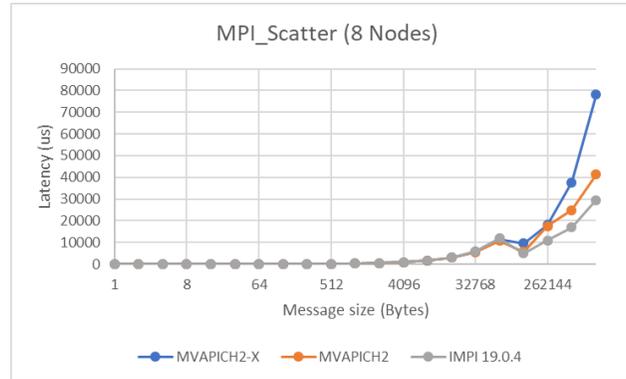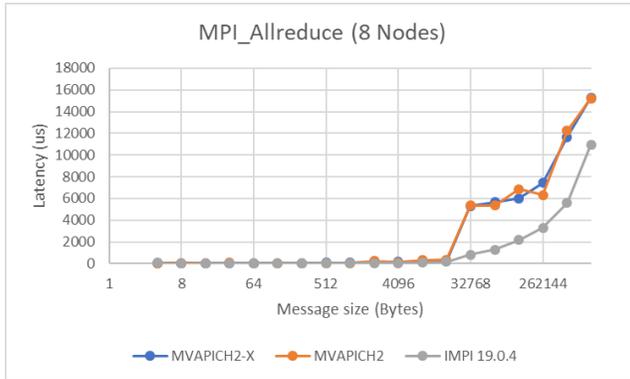


MPI Bi-directional Bandwidth

- Experimental environment
  - Single KNL node (cache mode)
- MPI Latency & Bandwidth
  - MVAPICH2-X better performance on 1B-256KB message size
  - MVAPICH2 better performance on 512KB-4MB message size
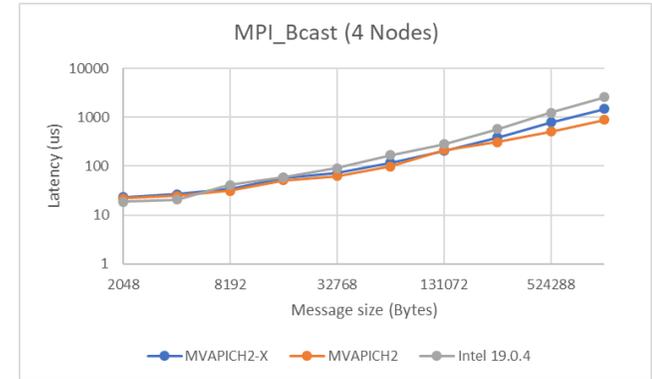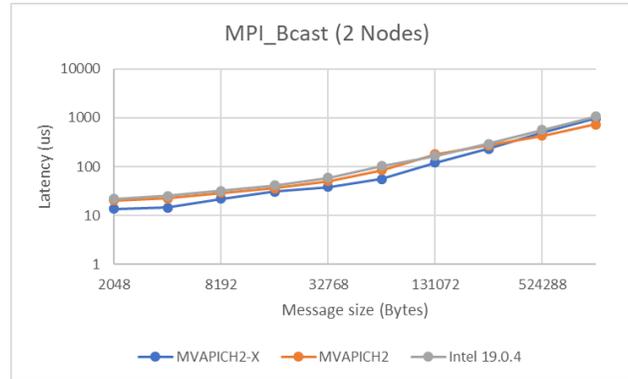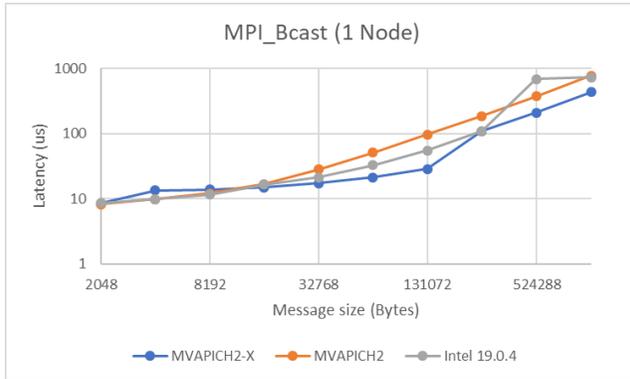- MPI Bi-directional BW
  - MVAPICH2 better performance

KiSTi Korea Institute of Science and Technology Information

# MVAPICH2-X with XPMEM: Multi Node

# Collective Communication on Multi Node: Latency



- Experimental environment
  - 8 KNL nodes (Nurion testbed, cache mode)
  - Processes per node = 64

- Performance evaluation
  - MPI_Allreduce: MVAPICH2-X better performance at specific range of the message size (2KB to 64KB)
  - MPI_Scatter: MVAPICH2 better performance
  - MPI_Alltoall: MVAPICH2-X slightly better performance, almost similar
  - In small message case, it could be better performance with threshold option (MV2_XPMEM_COLL_THRESHOLD)

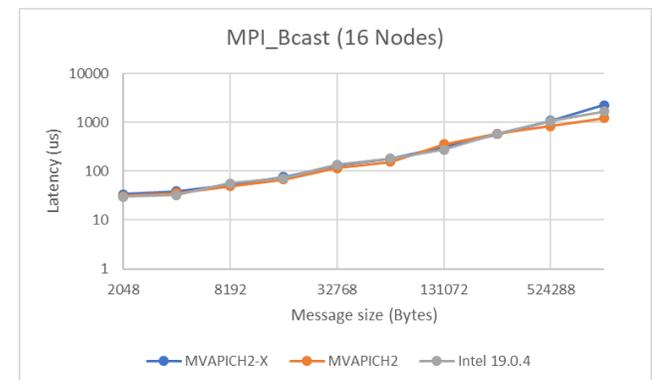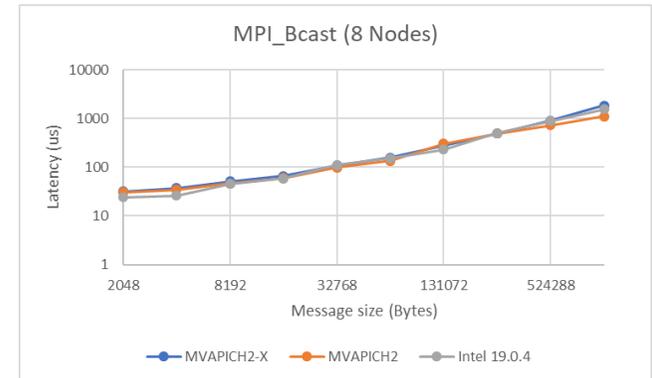# Collective Communication on Multi Node: Latency

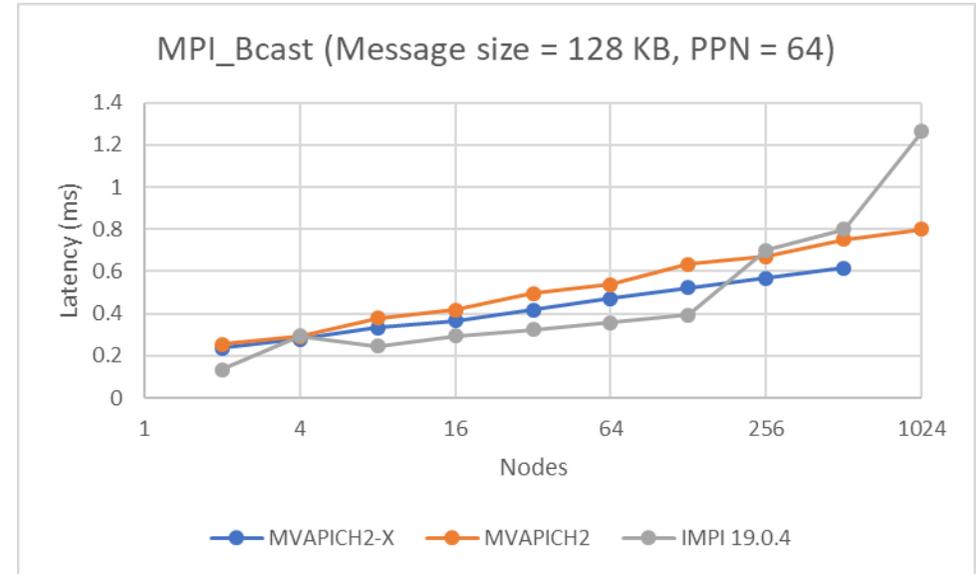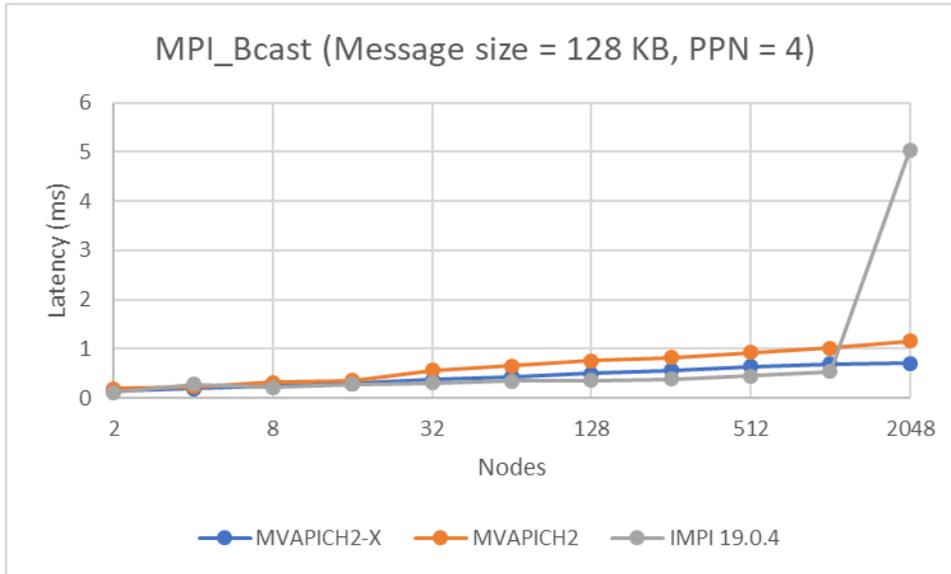

- Experimental environment
  - 1-16 KNL nodes (Nurion testbed, cache mode)
  - MPI_Bcast, processes per node = 64
- Performance evaluation
  - Single-node : MVAPICH2-X with XPMEM better performance if message size is larger than 8KB
  - Multi-node: MVAPICH2-X better performance at specific range of the message size
  - Similar performance results as the number of nodes increases

# Collective Communication on Multi Node: Latency



MPI_Bcast (Message size = 128 KB, PPN = 4)

MPI_Bcast (Message size = 128 KB, PPN = 64)

- Experimental environment
  - 2-2048 KNL nodes (Normal queue, cache mode)
  - MPI_Bcast, message size = 128KB, processes per node = 4, 64

- Performance evaluation
  - MVAPICH2-X with XPMEM is faster than MVAPICH2 (PPN=4, 64)
  - MVAPICH2 and MVAPICH2-X with XPMEM are faster than Intel MPI 19.0.4 on 2048 nodes (PPN=4), 256-1024 nodes (PPN=64)
  - MVAPICH2-X with XPMEM has segmentation fault problem on 1024 nodes
  - Intel IFS version of Nurion is slightly different (Intel IFS 10.8)

# Conclusion & Future Plan

- MVAPICH2-X with XPMEM
  - GCC compiler, Intel IFS 10.6 (Slightly different with version which is installed in Nurion)
  - Better performance at specific range of the message size on Nurion supercomputer
  - Another option for the large size message could improve the performance which is similar to the threshold option for the small size message (XPMEM_COLL_THRESHOLD)
- MVAPICH2-X installation
  - Intel compiler, Intel IFS 10.8
- Benchmark and application test (Presented at MUG'19)
  - Benchmark: HPCG, NPB
  - Application: DNS-TBL, GOTPM, LAMMPS
- Graph500 optimization with MVAPICH2-X
- Application optimization
  - WRF (Weather Research and Forecasting Model) application
  - Deep learning framework (TensorFlow, PyTorch, Caffe)

KiSTi Korea Institute of Science and Technology Information

THANK YOU

KiSTi
Korea Institute of
Science and Technology Information
www.kisti.re.kr