2024 OFA Virtual Workshop

# DESIGNING IN-NETWORK COMPUTING AWARE REDUCTION COLLECTIVES IN MPI

Bharath Ramesh and Dhabaleswar K. (DK) Panda
Network Based Computing Laboratory
The Ohio State University
http://nowlab.cse.ohio-state.edu/

# Outline

- **Introduction**

- Background

- Motivation

- Problem Statement and Contributions

- Design

  - Overview

  - Registration cache design

  - Proposed Allreduce design

- Results

- Conclusion and Future work

# Introduction: Drivers of Modern HPC Cluster Architectures

**Multi-/Many-core Processors**

**High Performance Interconnects – InfiniBand**
**<1usec latency, 200-400Gbps Bandwidth>**

**Accelerators high compute density, high performance/watt**
**>9.7 TFlop DP on a chip**

**SSD, NVMe-SSD, NVRAM**

- Multi-core/many-core technologies

- Remote Direct Memory Access (RDMA)-enabled networking (InfiniBand, RoCE, Slingshot)

- Solid State Drives (SSDs), Non-Volatile Random-Access Memory (NVRAM), NVMe-SSD

- Accelerators (NVIDIA GPGPUs)

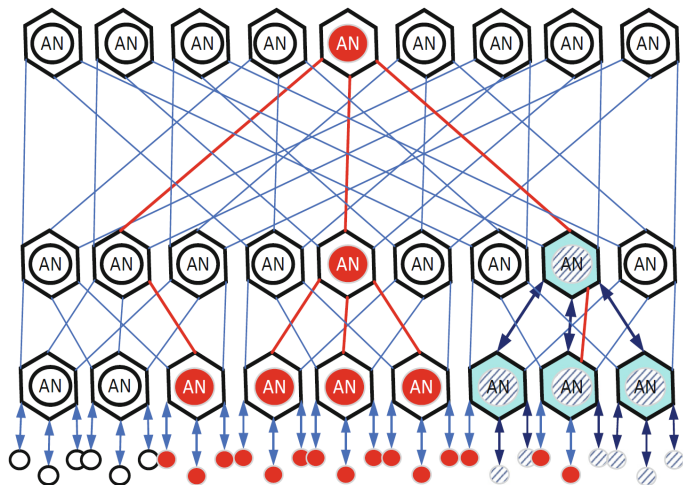*Frontier*　　　　*Fugaku*　　　　*Summit*　　　　*Lumi*

# MPI Reduction collectives and In-network Computing

- Reduction collectives (such as MPI_Allreduce) are important for HPC and AI

  - Involve both compute and communication

- Using CPUs everywhere leads to sub-optimal scale-up and scale-out efficiency

  - Motivates the need for offloading common operations away from the CPU to allow the CPU to perform other useful tasks

- In-network compute allows offloading operations to network devices

  - Switches are a good candidate due to high bandwidth and ability to reduce data on-the-fly eliminating redundancy

  - High scale-out efficiency and network topology awareness
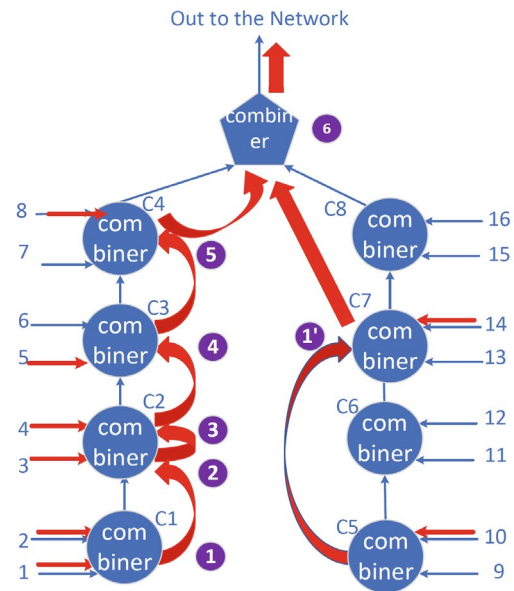
  - Frees up CPU cycles for other operations

# Outline

- Introduction

- **Background**

- Motivation

- Problem Statement and Contributions

- Design

  – Overview

  – Registration cache design

  – Proposed Allreduce design

- Results

- Conclusion and Future work

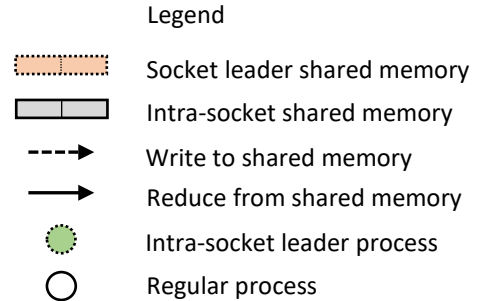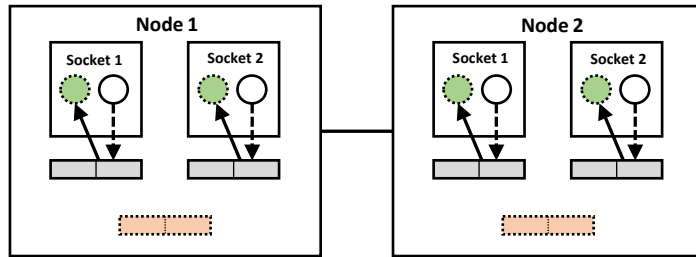# SHARP Reduction trees and Streaming Aggregation (SAT)



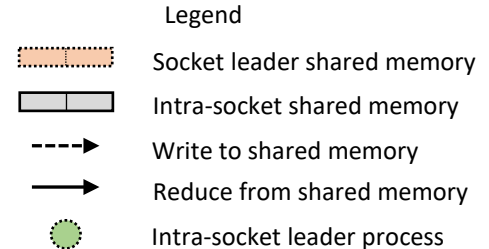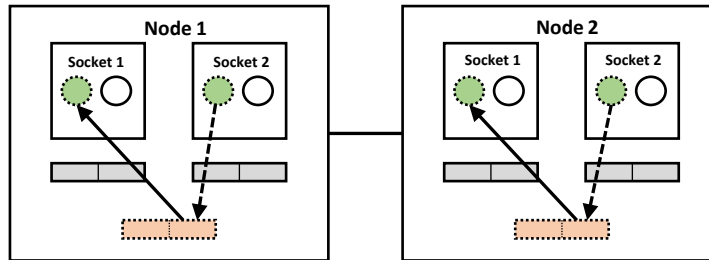Aggregation Tree

Switch-level reduction (radix 16)

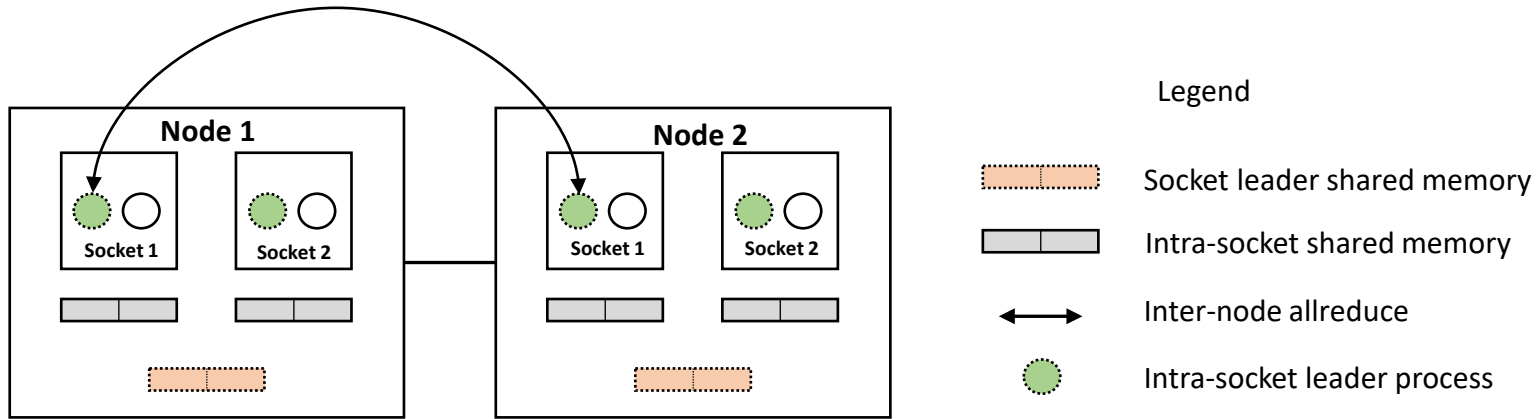# Hierarchical design for small message MPI_Allreduce

Phase 1 : Intra-socket reduction

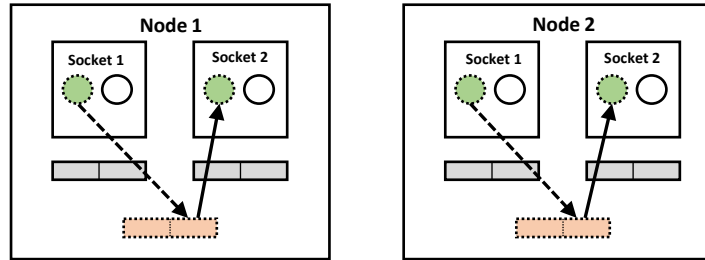Phase 2 : Inter-socket reduction

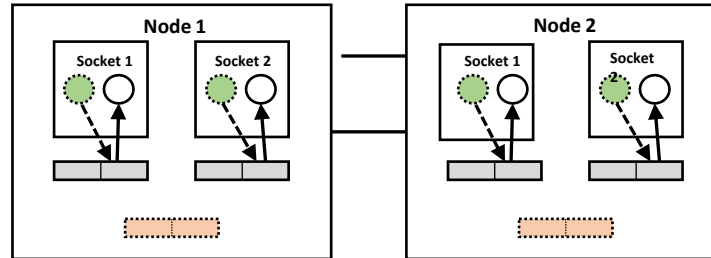# Hierarchical design for small message MPI_Allreduce



Phase 3 : Inter-node allreduce. Uses SHARP for scale-out performance
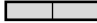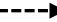
# Hierarchical design for small message MPI_Allreduce

Phase 4 : Inter-socket broadcast



Phase 5 : Intra-socket broadcast

# Overview of the MVAPICH Project

- High Performance open-source MPI Library

- Support for multiple interconnects
    - InfiniBand, Omni-Path, Ethernet/iWARP, RDMA over Converged Ethernet (RoCE),  AWS EFA, OPX, Broadcom RoCE, Intel Ethernet, Rockport Networks, Slingshot 10/11

- 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

- 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 (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

*23 Years & Counting!*

*2001-2024*

- **Used by more than 3,375 organizations in 91 countries**

- **More than 1.77 Million downloads from the OSU site directly**

- Empowering many TOP500 clusters (Nov '23 ranking)
    - 11th , 10,649,600-core (Sunway TaihuLight) at NSC, Wuxi, China
    - 29th , 448, 448 cores (Frontera) at TACC
    - 46th, 288,288 cores (Lassen) at LLNL
    - 61st, 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 29th ranked TACC Frontera system

- **Empowering Top500 systems for more than 18 years**

# Results for small MPI_Allreduce – Varying message sizes

1 ppn, 7861 nodes

16 ppn, 1024 nodes



- Scaling with message size, average latency

- **Close to a flat curve** across message sizes up to 2K

*Available in the MVAPICH 3.0 release*

# Results for small MPI_Allreduce – Varying node counts

- Scaling with increasing node counts, 16 bytes, average latency

- Same as trends with reduce (implementations are almost the same except for the intra-node broadcast phases)

*More information in the following paper*
*B. Ramesh, K. Suresh, N. Sarkauskas, M. Bayatpour, J. Hashmi, H. Subramoni, and DK Panda – "Scalable MPI Collectives using SHARP: Large Scale Performance Evaluation on the TACC Frontera System", ExaMPI'20*

# Outline

- Introduction

- Background

- **Motivation**

- Problem Statement and Contributions

- Design

  - Overview

  - Registration cache design

  - Proposed Allreduce design

- Results

- Conclusion and Future work

# Limitations of state-of-the-art schemes for large message reduction collectives

- Two-copy reduction collectives with SHARP

    - Used leader-based schemes that had a reduction, followed by a SHARP operation and finally a broadcast

    - Not suitable for large message sizes (>=128k)

- Single-copy schemes are very efficient for large message data movement

    - XPMEM allows remote process to have load/store access through address space mapping

- Using Sharp SAT in MPI has a few limitations and bottlenecks that need to be addressed for achieving good scale-out performance

- Motivates the need for large message reduction designs that combine advantages of SHARP and single-copy schemes like XPMEM

# Motivation

- SHARP SAT provides excellent bandwidth with close to point-to-point latency

- Registration involves pinning pages to memory (like InfiniBand registration)

  - Overhead increases significantly with increase in message size

  - Requires a cache that avoids expensive calls to sharp_coll_reg_mr

- Switch resources are limited

  - Causes bottlenecks when scaling up on modern CPUs with hundreds of cores

  - The SHARP runtime places limits to manage resources

- Motivates need for designs that are aware of SHARP runtime capabilities, overcome bottlenecks and scale-up efficiently for many processes per node

**Comparison of SHARP protocols on 8 nodes**



SHARP-LLT    SHARP-SAT

**Allreduce runtime registration overhead**



■ SHARP-allreduce-without-registration ■ SHARP-registration
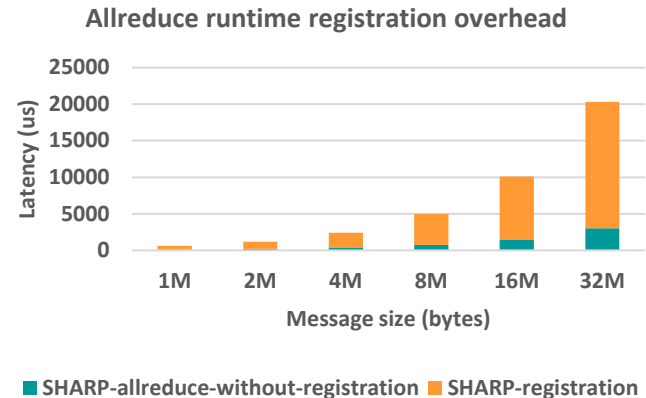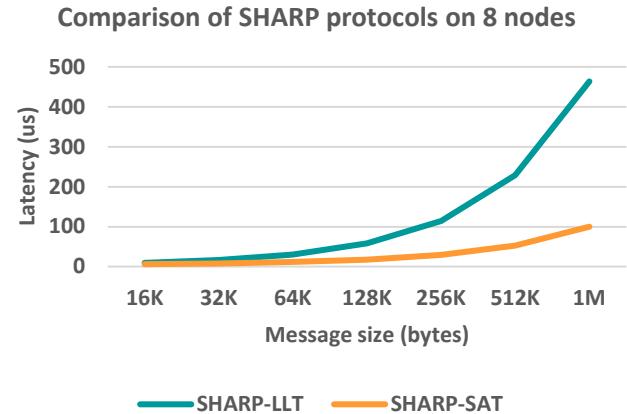
# Outline

- Introduction

- Background

- Motivation

- **Problem Statement and Contributions**

- Design

  - Overview

  - Registration cache design

  - Proposed Allreduce design

- Results

- Conclusion and Future work

# Problem Statement and Contributions

- **Problem Statement - Can we propose an algorithm for large message AllReduce that overcomes bottlenecks and resource constraints in the SHARP runtime by making efficient use of node and network level resources?**

- Contributions

  - Identify registration overheads involved in the use of SHARP streaming aggregation for large messages and propose solutions to address them

  - Analyze the impact of chunking reductions when using streaming aggregation for different message sizes to empirically determine ways to overlap intra-node reductions with SHARP-based reductions

  - Propose an algorithm for large AllReduce that utilizes SAT and CPUs efficiently

  - Evaluate the proposed design by comparing it against state-of-the-art MPI libraries

# Outline

- Introduction

- Background

- Motivation

- Problem Statement and Contributions

- **Design**

  – **Overview**

  – Registration cache design

  – Proposed Allreduce design

- Results

- Conclusion and Future work
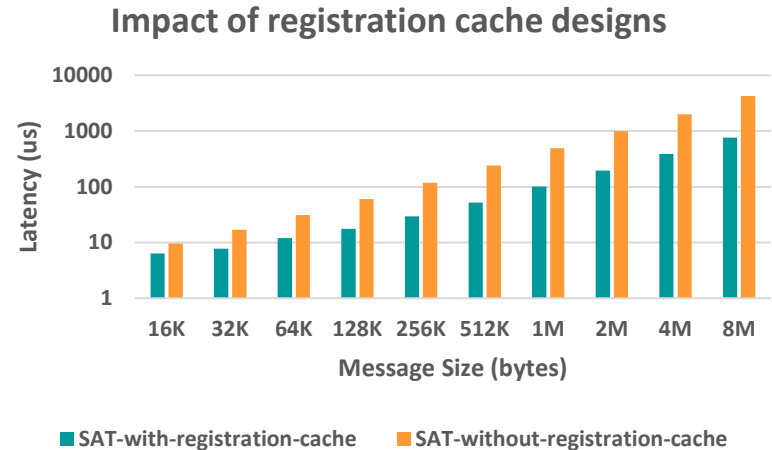
# Proposed Design Overview

- Use a registration cache to amortize registration costs in the SHARP runtime

- Designate a "leader" process on each node to interact with SHARP

- Chunk buffer into PPN (number of processes per node) chunks and reduce to a single buffer belonging to the leader process
  - Uses XPMEM for load/store access
  - All processes perform local reductions, but only leader process calls the SHARP runtime
  - Once local reductions are complete, leader calls a non-blocking MPI_Allreduce
    - Perfect overlap of intra-node and inter-node steps
  - Local reduction happens in batches for good network utilization
  - Final result broadcast within the node

# Outline

- Introduction

- Background

- Motivation

- Problem Statement and Contributions

- **Design**

  – Overview

  – **Registration cache design**

  – Proposed Allreduce design

- Results

- Conclusion and Future work

# Registration cache design

- Use an AVL tree or similar, to store buffer addresses
  - O(log n) insertion/query time
  - If buffer entry exists, directly get registration information from cache
- Up to 5.6X reduction in latency

**Impact of registration cache designs**



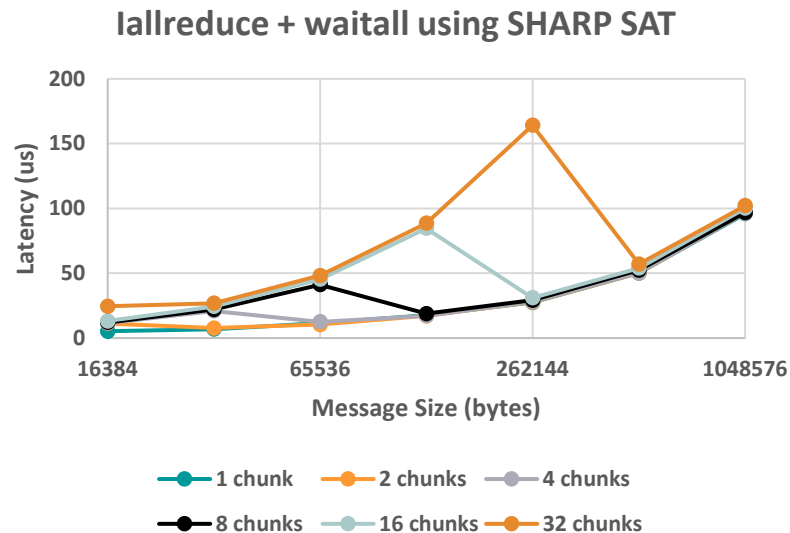Legend: ■ SAT-with-registration-cache   ■ SAT-without-registration-cache

# Outline

- Introduction

- Background

- Motivation

- Problem Statement and Contributions

- **Design**

  - Overview

  - Registration cache design

  - **Proposed Allreduce design**

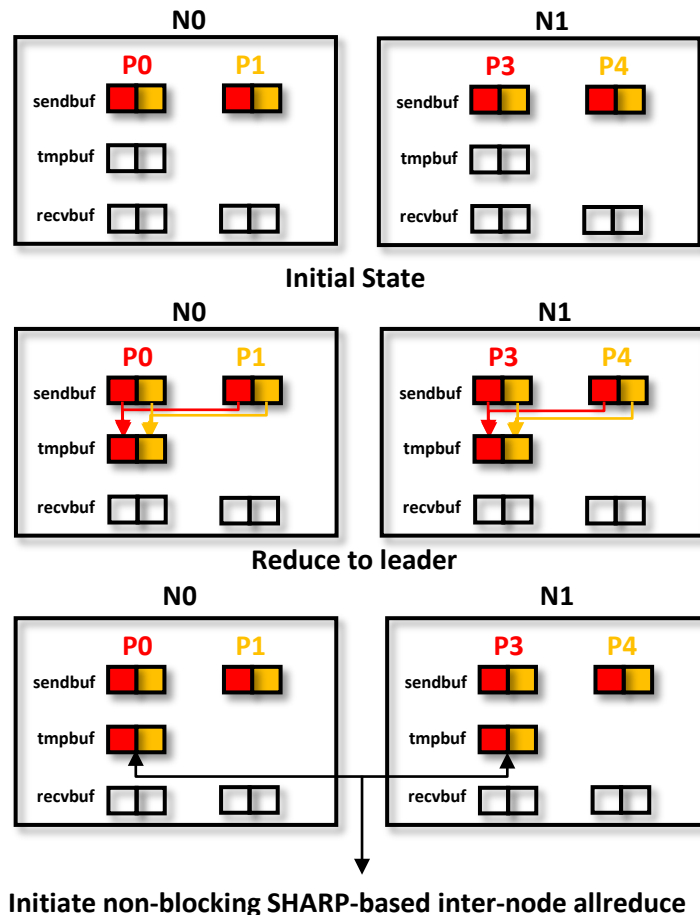- Results

- Conclusion and Future work

# Analyzing impact of chunking iallreduce operations

- Measure impact of a message sent using one call to the SHARP library vs multiple calls

- Given a message size M and number of chunks C, call non-blocking SHARP allreduce C times (of size M/C each) followed by waitall

- Indirect measure of overlap at the network level

- Splitting into chunks of size >= 16384 gives the same latency (independent of num_chunks)
  - Can be overlapped with reductions within the node
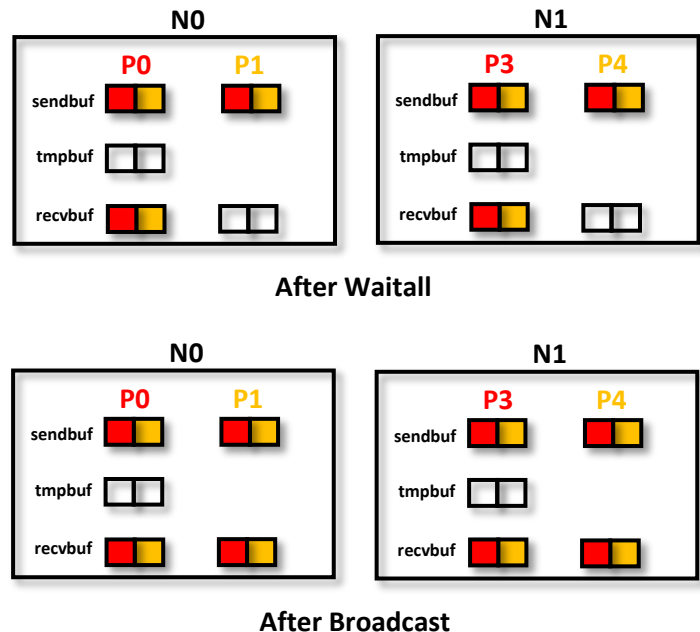
**Iallreduce + waitall using SHARP SAT**

# Proposed Allreduce Design

- First process on each node is designated as leader

- Before reduction, exchange buffer information using shared memory (for XPMEM load/stores)

- Process i reduces the ith chunk from every process and stores to tmpbuf at leader

- At the end of this step, leader on every node has the reduced result for the current phase

- Leader process initiates non-blocking inter-node SHARP allreduce

- Use "request" objects to track progress of SHARP Allreduce operations



Initial State

Reduce to leader

Initiate non-blocking SHARP-based inter-node allreduce

# Proposed Allreduce Design – Continued

- For large buffers, the intra-node reduction and inter-node phases are run multiple times

  - Reduction of large buffers is time consuming

  - Done in multiple phases for good network utilization

  - Chunk size if tuned to get perfect overlap of intra-node and inter-node operations

- Leader waits for non-blocking allreduces to complete after all runs of the first two phases are done

- Perform and intra-node broadcast to get final result

**After Waitall**

**After Broadcast**

*More information in the following paper*

*B. Ramesh, G. Kuncham, K. Suresh, R. Vaidya, N. Alnaasan, M. Abduljabbar, A. Shafi, D. Panda, Designing In-network Computing Aware Reduction Collectives in MPI, Hot Interconnects 2023, Aug 2023.*
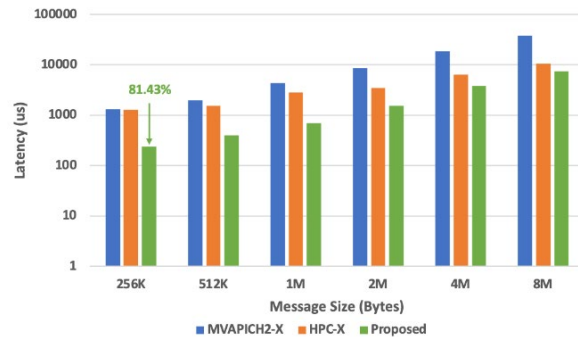
# Outline

- Introduction

- Background

- Motivation

- Problem Statement and Contributions

- Design

  – Overview

  – Registration cache design

  – Proposed Allreduce design

- **Results**

- Conclusion and Future work
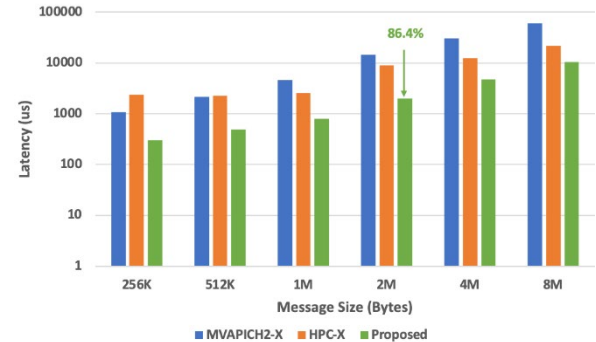
# Experimental setup

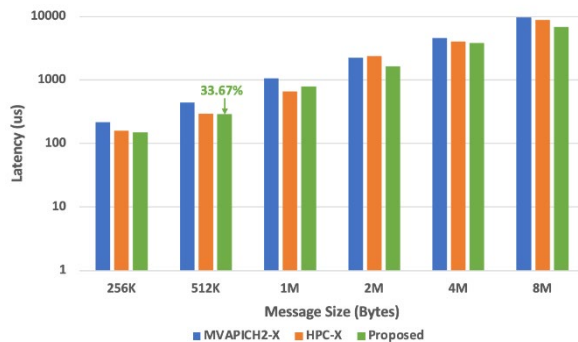| Cluster | MRI | HPCAC |
|---|---|---|
| Processor model | AMD EPYC 7713 | Intel(R) Xeon(R) Gold 6138 |
| Max Clock speed | 3.72GHz | 2GHz |
| Number of sockets | 2 | 2 |
| Cores per socket | 64 | 20 |
| RAM | 256GB | 196GB |
| Interconnect | NVIDIA HDR-200 with Quantum 2 switches | NVIDIA HDR-200 with Quantum 2 switches |
| MPI libraries | MVAPICH2-X, HPC-X | MVAPICH2-X, HPC-X |

# Results for large MPI_Allreduce – 2 nodes

- Increased parallelism by using multiple processes and SHARP for reduction

- Up to 81.43% over state-of-the-art for 32PPN and 86.4% for 64PPN on MRI

- Up to 33.67% over state-of-the-art for 32PPN and 60% for 64PPN on HPCAC

- Increased number of page faults leads to decreased benefits at 1M (Needs to be investigated further)
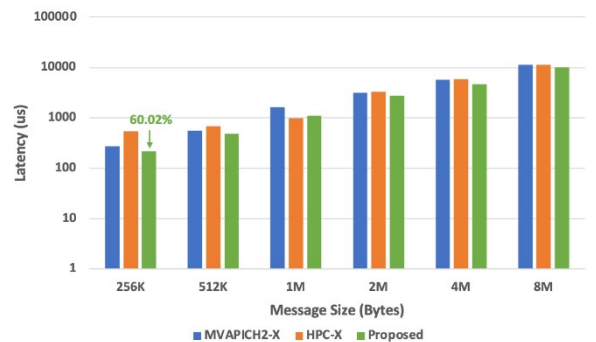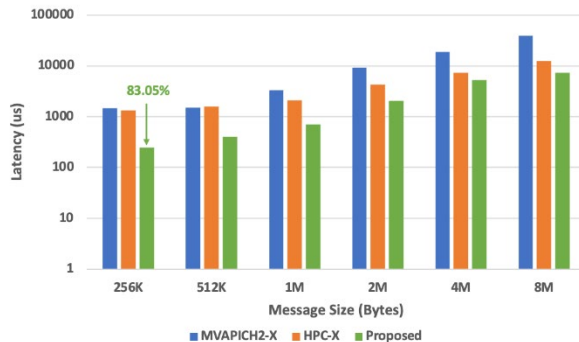


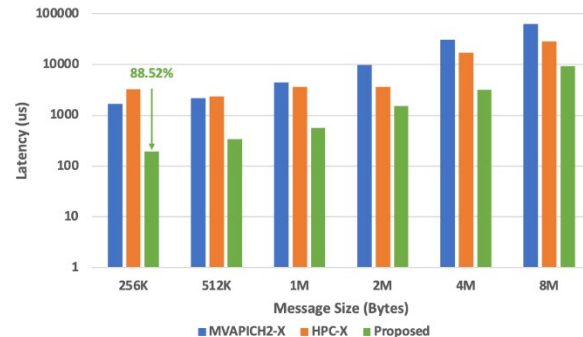MRI - 32PPN



MRI - 64PPN



HPCAC - 32PPN



HPCAC - 64PPN

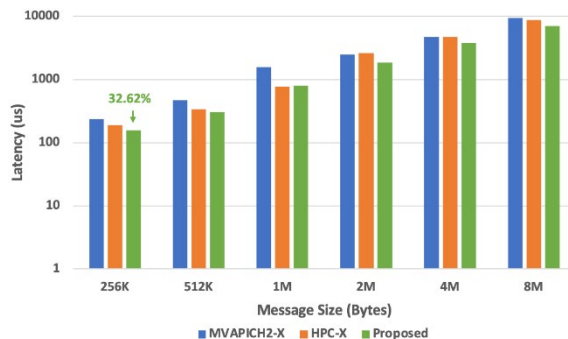# Results for large MPI_Allreduce – 4 nodes

- Increased parallelism by using multiple processes and SHARP for reduction

- Up to 83.05% over state-of-the-art for 32PPN and 88.52% for 64PPN on MRI

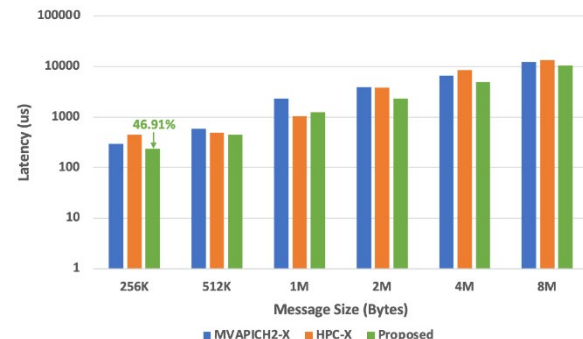- Up to 32.62% over state-of-the-art for 32PPN and 46.91% for 64PPN on HPCAC
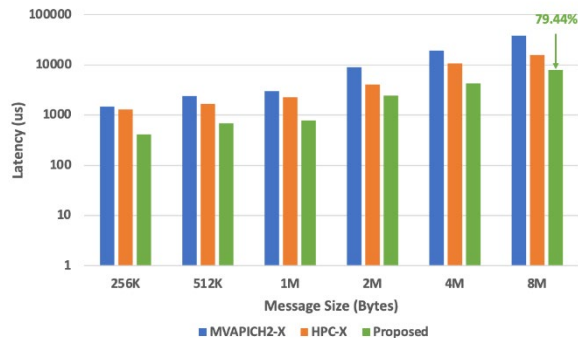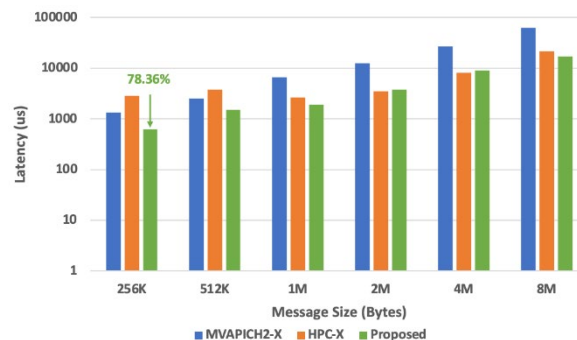


MRI - 32PPN



MRI - 64PPN



HPCAC - 32PPN



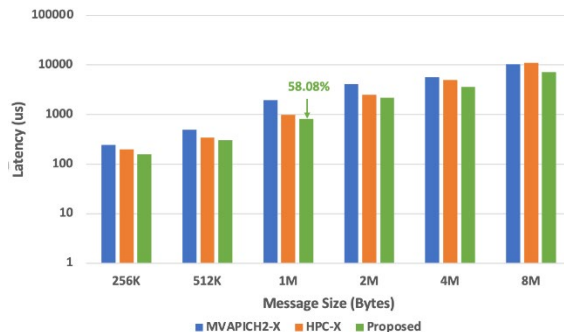HPCAC - 64PPN

# Results for large MPI_Allreduce – 8 nodes

- Increased parallelism by using multiple processes and SHARP for reduction

- Up to 79.44% over state-of-the-art for 32PPN and 78.36% for 64PPN on MRI

- Up to 58.08% over state-of-the-art for 32PPN and 52.13% for 64PPN on HPCAC
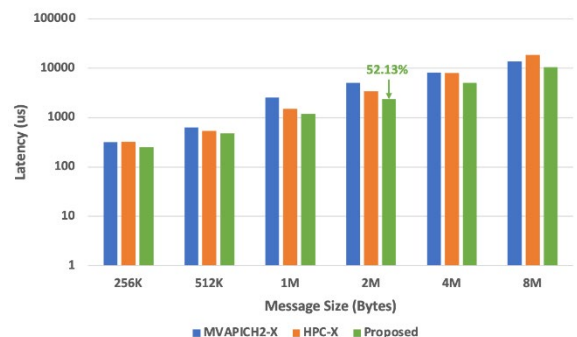


MRI - 32PPN



MRI - 64PPN



HPCAC - 32PPN
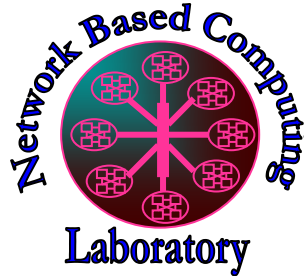


HPCAC - 64PPN

# Outline

- Introduction

- Background

- Motivation

- Problem Statement and Contributions

- Design

  – Overview

  – Registration cache design

  – Proposed Allreduce design

- Results

- **Conclusion and Future work**

# Conclusion and Future Work

- SHARP runtime enables in-network offload with excellent bandwidth utilization

- Proposed designs overcome various bottlenecks by using a leader-based algorithm and streaming aggregation for large message reductions
  - Outperforms state-of-the-art by up to 86%

- Will be available in a future release of MVAPICH-plus

- Future work
  - Comprehensive application evaluation
  - Evaluating performance at larger scales
  - Exploring NUMA-awareness

# THANK YOU!



**Network-Based Computing Laboratory**
**http://nowlab.cse.ohio-state.edu/**







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

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

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