

# MPI4AI: Integrating AMD’s RCCL into MPICH for GPU-Accelerated Collective Communication

Grace Li<sup>1</sup>*Mathematics and Computer Science Division, Argonne National Laboratory<sup>a)</sup>*

(\*Electronic mail: grace.li@anl.gov)

(Dated: 1 August 2025)

Message Passing Interface (MPI) remains the de-facto standard for communication in high-performance computing. However, as modern AI workloads increasingly rely on GPU acceleration, MPI’s CPU-bound collectives become a performance bottleneck—particularly for operations like Allreduce that dominate training workloads. This project addresses that gap by integrating AMD’s GPU collective library, RCCL, into MPICH, the premier open-source MPI implementation. Our integration enables flexible, runtime-selectable switching between CPU- and GPU-based collectives based on message size, using tuning profiles. Benchmarking on AMD GPUs shows over 60× lower latency for large messages compared to standard MPICH. This enhancement significantly boosts MPICH’s performance on GPU-rich systems like Aurora and Frontier—Department of Energy exascale supercomputers built for AI and scientific computing.

## I. INTRODUCTION

In modern high-performance computing (HPC) and large-scale artificial intelligence (AI) workloads, communication overhead increasingly limits scalability. The Message Passing Interface (MPI) standard is central to efficient collective communication across distributed systems.

However, the shift to GPU-accelerated workloads has exposed traditional CPU-based MPI collectives as a key bottleneck. To address this, vendors like NVIDIA and AMD have developed GPU-native collective libraries (e.g., NCCL, RCCL) that execute communication kernels directly on the GPU. While these libraries deliver strong performance for large message sizes, they lack MPI compliance, limiting portability.

To bridge this gap, hybrid models—such as the one proposed by Chen *et al.*<sup>1</sup>—that switch between CPU and GPU backends based on workload characteristics have emerged. Building on this idea, our work integrates AMD’s RCCL into MPICH<sup>2</sup>, a widely used, high-performance MPI implementation developed at Argonne National Laboratory, to enable standards-compliant, GPU-accelerated collectives. Specifically, we aim to:

- Implement a HIP-based RCCL backend for MPICH
- Enable backend switching via environment variables
- Support dynamic selection using JSON tuning files
- Benchmark performance on JLSE nodes at Argonne
- Develop a stream-parallel Allreduce implementation

This paper presents the design, implementation, and evaluation of the MPICH-RCCL integration, and discusses its impact on performance portability for GPU-dense HPC and AI workloads.

## II. BACKGROUND

### A. Overview of MPI and Collective Operations

MPI enables distributed processes to communicate in parallel applications and serves as the backbone of HPC workloads, from climate modeling to AI training. At its core are collective operations, which involve coordinated communication among groups of processes—unlike point-to-point communication, where data is exchanged between just two processes. Examples include:

- **Broadcast:** Sends data from one process to all others.
- **Gather/Scatter:** Collects or distributes data among all processes.
- **Allreduce:** Aggregates values (e.g., sum, max, average) across all processes and returns the result to each. Figure 1 illustrates this pattern across four processes.

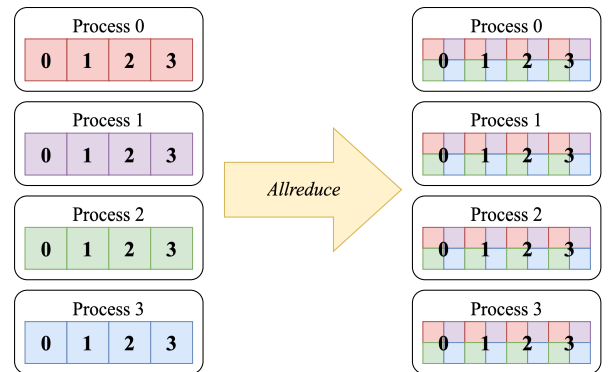


FIG. 1: Visualization of Allreduce across four processes.

Allreduce is especially critical in AI, enabling synchronous data-parallel training by averaging gradients across nodes during backpropagation.

<sup>a)</sup> Also at Columbia University; alternate email: grace.li2@columbia.edu

## B. CPU- vs. GPU-Driven Collectives

Graphics processing units (GPUs) have evolved into highly parallel processors that accelerate compute-intensive workloads, making them ideal for large-scale AI and HPC applications. In contrast, central processing units (CPUs) are responsible for managing execution and orchestrating communication. While MPICH is CUDA-aware and can use GPU-allocated memory as input through transports like FI\_HMEM (a libfabric interface for exposing GPU memory to high-performance communication frameworks), the actual reduction still occurs on the CPU, creating a performance bottleneck.

Collective communication libraries (CCLs), like NVIDIA’s NCCL<sup>3</sup> and AMD’s RCCL<sup>4</sup>, improve efficiency by launching GPU kernels that execute communication operations directly on the device. This reduces CPU-GPU synchronization and host-device data transfers. However, launching GPU kernels introduces latency, which can outweigh the benefits for smaller messages.

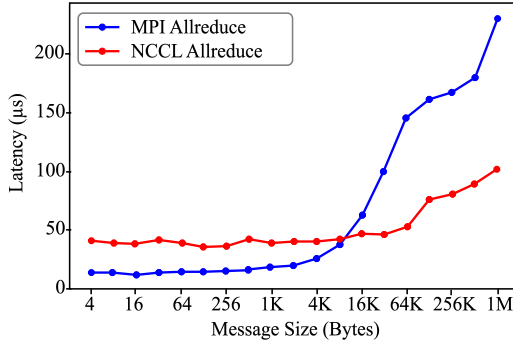


FIG. 2: Allreduce performance comparison between NCCL and MPI. Adapted from Chen *et al.*<sup>1</sup>.

Figure 2 illustrates this trade-off: NCCL Allreduce significantly outperforms MPI for large messages, while MPI retains an edge at smaller sizes. This motivates a hybrid approach of integrating CCLs into MPI.

## C. RCCL and the ROCm Ecosystem

The ROCm ecosystem<sup>5</sup> offers tools for AMD GPU development but lacks full MPI compliance. Our work integrates RCCL—built on HIP (Heterogeneous-computing Interface for Portability)—into MPICH, enabling legacy MPI applications to benefit from efficient GPU-based collectives.

While other CCLs like NCCL<sup>3</sup>, oneCCL<sup>6</sup>, and UCC<sup>7</sup> offer similar capabilities, integrating RCCL is essential given AMD’s growing presence in HPC. U.S. Department of Energy supercomputers such as Frontier at Oak Ridge National Laboratory<sup>8</sup> and El Capitan at Lawrence Livermore National Laboratory<sup>9</sup> are built on AMD GPUs and the ROCm stack, requiring seamless MPI-RCCL support.

## III. IMPLEMENTATION

### A. Enabling RCCL Support in MPICH

Our integration builds on the MPI-xCCL project<sup>1</sup>, which uses tuning tables to switch between MPI-native and CCL backends based on message size. We also draw from the NCCL backend in MPICH<sup>10</sup>, which enables GPU collectives on NVIDIA systems.

To enable RCCL support, we leverage MPICH’s modular device-layer architecture. We modify the build configuration to detect RCCL headers and link against the RCCL library, while also checking for HIP support. These adjustments ensure that RCCL support is enabled only on compatible systems.

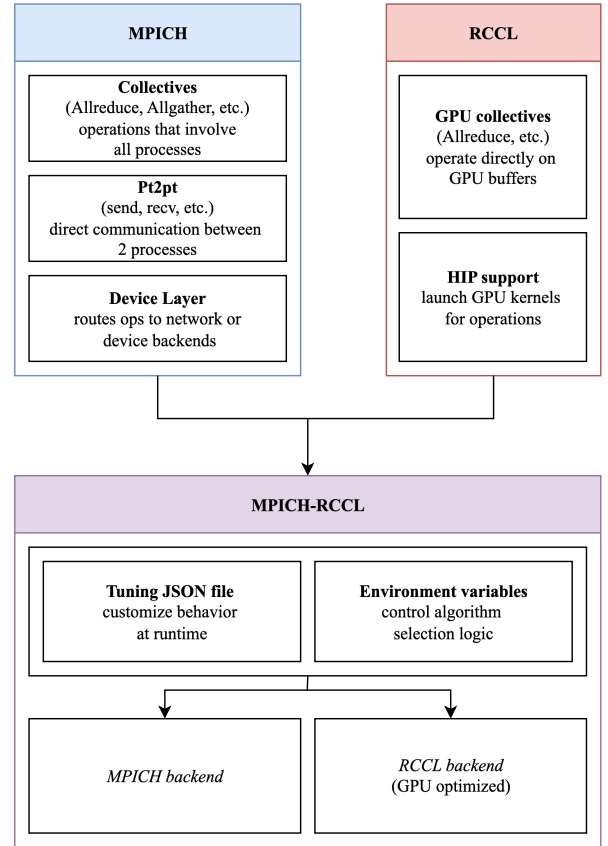
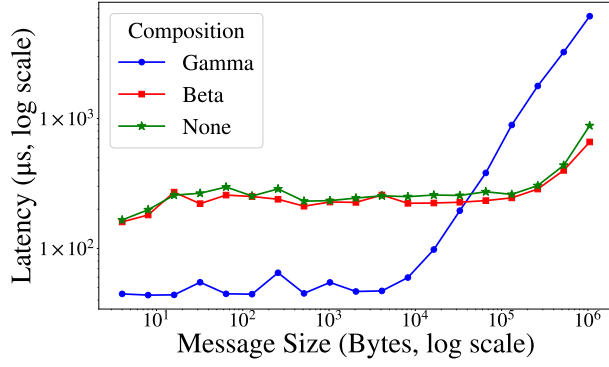


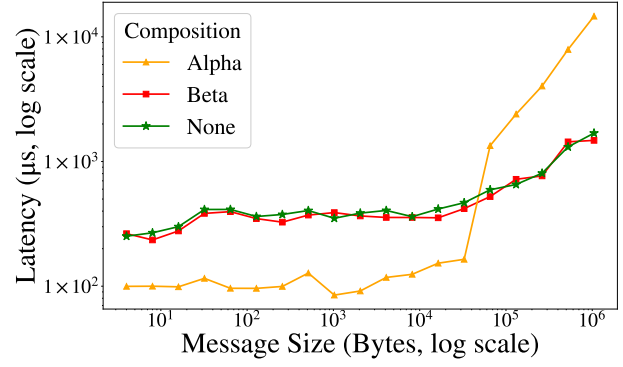
FIG. 3: Architecture overview of MPICH-RCCL.

Figure 3 illustrates the hybrid architecture of our MPICH-RCCL integration. At runtime, MPICH dynamically selects between CPU and GPU backends depending on message size.

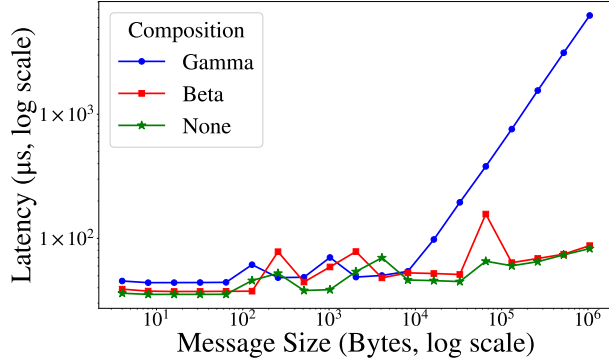
By enabling device collectives, MPICH can now execute operations like Allreduce entirely within GPU memory using RCCL, avoiding costly host-device transfers and reducing latency. If GPUs are unavailable, MPICH seamlessly defaults to CPU-based logic. This hybrid model combines MPI’s portability with GPU performance, allowing runtime backend switching to optimize communication across diverse HPC and AI workloads.



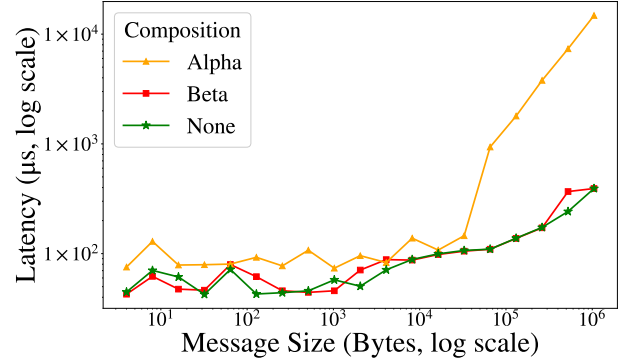
(a) MPICH, 1 node (gamma, beta, DC off)



(b) MPICH, 2 nodes (alpha, beta, DC off)



(c) RCCL, 1 node (gamma, beta, DC off)



(d) RCCL, 2 nodes (alpha, beta, DC off)

FIG. 4: Allreduce composition latency across MPICH (top) and RCCL (bottom) backends for 1- (left) and 2-node (right) runs. Lower latency corresponds to faster performance.

## B. Composition Algorithms

A composition algorithm describes how a collective operation is implemented. MPICH supports several such algorithms for Allreduce:

**Alpha:** A leader-based algorithm optimized for multi-node systems. Each node reduces data to a leader, which then participates in inter-node Allreduce, followed by local broadcast. This reduces cross-node traffic and exploits fast intra-node communication.

**Beta:** A default, topology-agnostic algorithm that treats all processes equally. It doesn't differentiate between intra- and inter-node communication and is used when device collectives are disabled.

**Gamma:** Designed for single-node systems with multiple GPUs. All communication stays within the node, avoiding inter-node overhead to maximize local performance. Not used in multi-node setups.

Figure 4 shows the results of our initial benchmarking. On a single node with 4 AMD GPUs, we compare mpi-gamma

and rccl-beta. For small messages (less than 8 KB), mpi-gamma performs slightly better due to MPICH's optimized shared-memory communication. However, for larger messages, rccl-beta achieves significantly higher throughput by leveraging GPU memory.

We also evaluate performance on two nodes (8 GPUs total), comparing mpi-alpha and rccl-beta. For mid-sized messages, mpi-alpha benefits from its leader-based structure. Yet as message sizes grow, rccl-beta outperforms MPICH, demonstrating the advantages of GPU-driven collectives for high-volume communication.

## C. Custom Tuning

CPU-based collectives often perform better for small messages, while GPU-driven algorithms excel with larger ones. MPICH's composition engine supports dynamic algorithm selection at runtime, which we leverage using environment variables and JSON-based tuning. Users can specify message-size thresholds and preferred algorithms, enabling fine-grained control and easy experimentation.

Our performance results inform the switching thresholds:

RCCL is used for messages over 8 KB on single-node and 32 KB on multi-node jobs. MPICH’s JSON-based framework applies these thresholds at runtime for automatic backend selection. Figure 5 illustrates the tuning pipeline, from environment variables to threshold-driven algorithm choice.

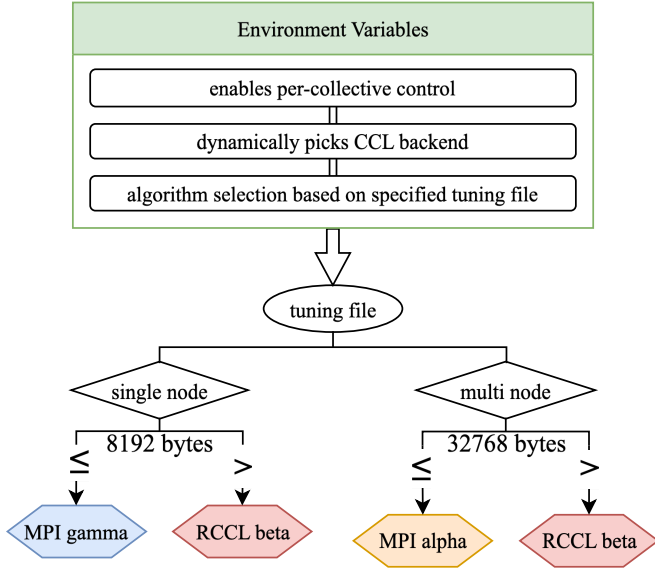


FIG. 5: MPICH-RCCL Allreduce tuning via environment variables and a JSON-based tuning file.

The result is a modular, user-friendly system that dynamically switches between backends based on message size, with no code changes required. This design unifies portability and performance, making it ideal for GPU-intensive applications.

## IV. EVALUATION

### A. Methodology

We evaluate the MPICH-RCCL hybrid using OSU Micro-Benchmarks<sup>11</sup>, measuring Allreduce latency across a range of message sizes. These benchmarks are widely used in HPC for assessing the performance of collective communication.

Latency refers to the time it takes for an operation to complete. Lower latency is critical in HPC because it directly impacts the efficiency of synchronized workloads.

We conduct our experiments on Argonne’s JLSE (Joint Laboratory for System Evaluation) cluster, a testbed designed for HPC. Each node features four AMD Instinct GPUs—MI100, MI250, or MI300X—equipped with 32 GB of HBM2, 128 GB of HBM2e, and 192 GB of HBM3, respectively.

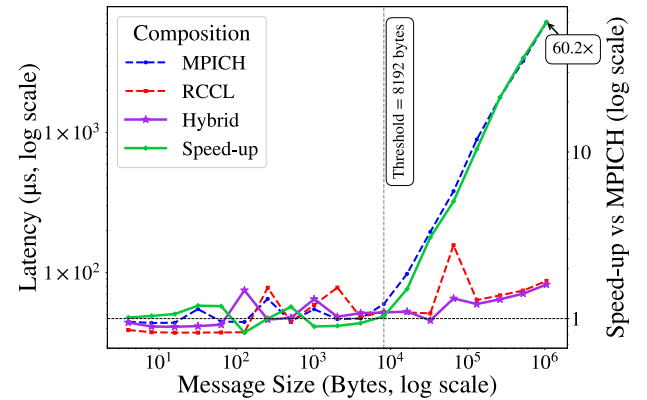
These GPUs are specifically optimized for HPC and rely on High Bandwidth Memory (HBM) to meet the demands of data-intensive workloads. The use of HBM2, HBM2e, and HBM3 across the AMD Instinct lineup provides increasing bandwidth and capacity.

We evaluate using an MPICH build configured with a ch4:ucx device layer and both RCCL and HIP support. Uni-

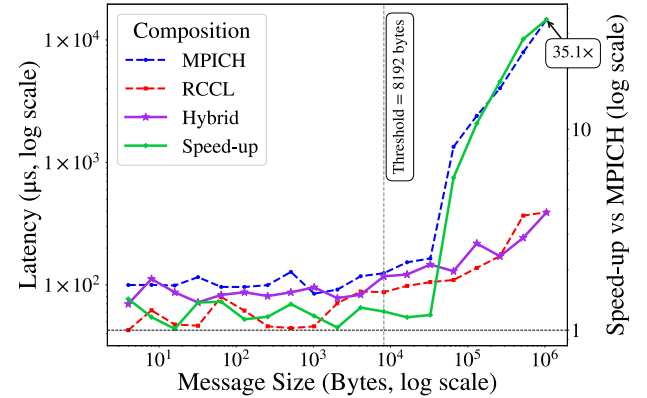
fied Communication X (UCX) enables GPU-aware communication over both shared-memory and network transports, including ROCm-specific transport layers. In single-node experiments, UCX uses the rocm, sm (shared memory), and self (loopback) transports for intra-node communication. For multi-node experiments, UCX also uses the tcp transport (socket-based communication over a network).

We test both single-node and two-node configurations to isolate intra-node and inter-node behaviors. Message sizes range from 4 bytes to 1 megabyte, covering both latency- and bandwidth-dominated regimes.

### B. Results



(a) 1-node: switch from MPICH to RCCL at 8 KB.



(b) 2-nodes: switch from MPICH to RCCL at 32 KB.

FIG. 6: Allreduce latency. The hybrid achieves 60.2× and 35.1× speedup for 1 MB messages. Lower latency (purple) and higher speedup (green) indicate better performance.

We compare native MPICH collectives, RCCL device collectives, and a JSON-tuned hybrid configuration. The hybrid results from our custom tuning described earlier. To quantify its performance gains, we plot speedup relative to native MPICH. The results are shown in Figure 6.

In a single-node configuration, `mpi-gamma` algorithm performs best at low message sizes. However, as the message size increases beyond 8 KB, `rccl-beta` significantly outperforms MPICH. At 1 MB, RCCL demonstrates  $60.2\times$  lower latency than native MPICH, showcasing the power of GPU-resident data paths and optimized hardware interconnects.

In a two-node configuration, we raise the switching threshold to 32 KB to account for inter-node overhead. Here, `mpi-alpha` initially performs well, but beyond 32 KB, the overhead outweighs the benefits. RCCL again becomes dominant, achieving  $35.1\times$  speedup over native MPICH.

These results validate our design choices and tuning heuristics. They also highlight the need for hybrid approaches in heterogeneous environments, where no single algorithm or backend is universally optimal.

## V. PIPELINED ALLREDUCE

### A. Motivation

To reduce collective latency for large messages, we explore a stream-parallel version of Allreduce. Single-stream execution tends to underutilize available GPU concurrency, especially on modern architectures equipped with multiple compute engines and extensive support for asynchronous kernel execution and data transfers.

The broader objective is to establish a foundation for more sophisticated collective strategies. These are particularly relevant in heterogeneous computing environments where overlapping communication and computation can significantly reduce latency.

### B. Implementation

We develop this pipelined approach with two strategies:

1. **Multi-Stream Extension:** We begin by extending the RCCL backend to support multiple HIP streams within a single communicator. For each Allreduce call, the message is evenly divided into four chunks. Each chunk is assigned to a different HIP stream, concurrently launching four RCCL Allreduce operations. Although these are dispatched to separate streams, they share the same communicator context—meaning the underlying collective operations remain logically serialized. This approach improves kernel launch concurrency but is ultimately limited by intra-communicator synchronization semantics.
2. **Multi-Communicator Support:** To address the above limitation, we implement support for multiple RCCL communicators. At startup, we initialize four communicators, each bound to a unique HIP stream and reused across runs. Each chunk of the message is assigned to a distinct communicator and stream, enabling four truly concurrent Allreduce operations. This design

eliminates inter-chunk serialization caused by single-communicator constraints and better exploits GPU concurrency by decoupling collective execution.

### C. Evaluation

Experiments were conducted on JLSE using 2 nodes, each equipped with 4 AMD MI100 GPUs (32 GB of HBM2 memory per GPU). Communication is managed using the UCX transport layer.

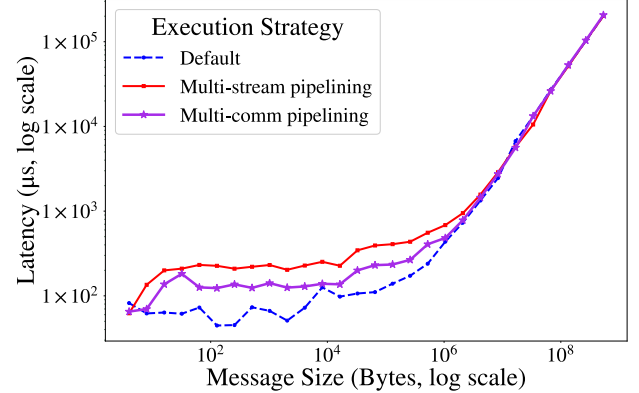


FIG. 7: Stream-parallel Allreduce variants. For smaller message sizes, our variants underperform compared to the default implementation. At larger message sizes (above 8 MB), performance converges with the default as execution becomes bandwidth-bound.

As shown in Figure 7, our stream-parallel implementations do not currently yield significant speedups. For messages under 8 MB, the multi-stream and multi-communicator strategies underperform due to kernel launch overhead and limited parallelism benefits. At larger message sizes (8 MB to 537 MB), performance converges to the default single-stream baseline as execution becomes bandwidth-bound, diminishing the advantage of stream-level parallelism.

While current gains are modest, the approach remains promising. Ongoing work focuses on improving stream concurrency and coordination, with the goal of reducing latency.

## VI. RELATED AND FUTURE WORK

Our work builds on several lines of research in MPI collectives, GPU communication, autotuning, and large-scale HPC systems.

### A. Foundations of GPU Integration and Portability

The foundation for this project draws from MPI-xCCL by Chen *et al.*<sup>1</sup>, which introduced a portable MPI layer capable of integrating GPU collective libraries with dynamic backend

selection. We extend this design into the AMD ROCm ecosystem, enabling transparent GPU communication via RCCL within MPICH.

We mirror the MPICH-NCCL integration<sup>10</sup>, adapting its modular backend structure and tuning interface for RCCL. This minimized development friction and ensured consistency.

Unified Communication X (UCX)<sup>12</sup> and Unified Collective Communication (UCC)<sup>7</sup> offer accelerator-aware frameworks, but diverge from standard MPI semantics. Nihar Kodkani<sup>13</sup> explores accelerating GPU collectives in MPICH using UCC in a concurrent effort. Our approach maintains strict MPI compliance, ensuring legacy compatibility while improving performance.

## B. Collective Algorithm Optimization

Our composition logic is informed by prior work on optimizing MPICH’s collectives. Thakur and Gropp<sup>14</sup> developed message-size-dependent algorithm selection strategies that directly influence our switching thresholds for GPU vs. CPU backends.

Additionally, Balaji *et al.*<sup>15</sup> examined MPI scalability to petascale systems, identifying memory and algorithmic scalability concerns. Our work aligns with these goals by reducing communication bottlenecks through GPU collectives.

Studies show that Allreduce accounts for a significant share of MPI core-hours in production systems, highlighting its critical role<sup>16,17</sup>.

## C. Autotuning and Machine Learning for MPI

A growing body of work investigates replacing hardcoded algorithm selection with intelligent autotuning. ACCLaIM<sup>18</sup> and FACT<sup>19</sup> use machine learning to adaptively choose optimal collective algorithms. Our work complements these efforts by exposing GPU-accelerated paths as tunable options.

Other tuners like PGMPITuneLib<sup>20</sup> and OMPICollTune<sup>21</sup> intercept MPI calls to optimize algorithm selection at runtime. SCCL<sup>22</sup> offers a systematic method for synthesizing collective algorithms along a latency-bandwidth Pareto frontier. Alternative techniques like quadtree encoding<sup>23</sup> further highlight fast runtime decision logic.

## D. GPU-Aware MPI Communication

We also acknowledge extensive work optimizing MPI on GPU-enabled clusters. Studies on GPUDirect RDMA<sup>24</sup> and CUDA IPC<sup>25</sup> demonstrate latency reductions for intra/inter-node GPU communication.

Ayala *et al.*<sup>26</sup> show challenges in large-scale FFT workloads using collectives, motivating libraries like HEFFTE. BlueConnect<sup>27</sup> decomposes Allreduce across network hierarchy, a model we may explore for future scaling on Frontier.

## E. Future Directions

We plan to:

- Extend RCCL support to other collectives (e.g., Broadcast, Allgather, Scatter)
- Contribute upstream to the MPICH repository
- Perform scalability experiments on Frontier
- Evaluate full AI training workloads using our MPICH-RCCL backend
- Refine pipelined collective strategies to improve stream overlap and reduce latency

Our goal is to unify portability, performance, and scalability for AI and HPC, aligned with DOE’s mission of accelerating scientific discovery on exascale systems<sup>28</sup>.

## VII. CONCLUSION

This work presents an integration of AMD’s RCCL into MPICH, enabling hybrid CPU-GPU Allreduce collectives that dynamically switch between backends based on message size and system characteristics. Our results demonstrate substantial performance gains—over 60×—compared to native MPICH for large messages, without sacrificing MPI portability or runtime flexibility. The hybrid execution model enables users to dynamically balance latency and throughput, enabling them to extract maximum performance from GPU-dense systems. By embedding GPU collectives into a standards-compliant MPI stack, our approach advances the goals of MPI4AI: making MPI efficient and portable in emerging AI workflows.

## ACKNOWLEDGMENTS

This project was made possible thanks to the mentorship of Dr. Mike Wilkins and the support of collaborators Nihar and Sean. We thank the MPICH/PMRS team for their foundational NCCL backend work and guidance throughout this process. Special thanks to Argonne National Laboratory for the research opportunity, and to Oak Ridge National Laboratory for access to Frontier.

This work was supported in part by the U.S. Department of Energy, Office of Science, Office of Workforce Development for Teachers and Scientists (WDTS) under the Science Undergraduate Laboratory Internships (SULI) program<sup>29,30</sup>.

## REFERENCES

- <sup>1</sup>C.-C. Chen, K. S. Khorassani, P. Kousha, Q. Zhou, J. Yao, H. Subramoni, and D. K. Panda, in *Proceedings of the SC '23 Workshops* (ACM, New York, NY, USA, 2023) pp. 847–854.



- <sup>2</sup>Argonne National Laboratory, “MPICH,” <https://www.mpich.org> (2025).
- <sup>3</sup>NVIDIA, “NVIDIA Collective Communications Library (NCCL),” <https://developer.nvidia.com/nccl> (2025).
- <sup>4</sup>AMD, “RCCL: ROCm Collective Communication Library,” <https://github.com/ROCm/rccl> (2025).
- <sup>5</sup>AMD, “ROCm Developer Resources,” <https://rocm.docs.amd.com> (2016).
- <sup>6</sup>Intel, “Intel® oneAPI Collective Communications Library,” <https://www.intel.com/content/www/us/en/developer/tools/oneapi/oneccl.html> (2025).
- <sup>7</sup>UCF Consortium, “UCC: Unified Collective Communication,” <https://ucfconsortium.org/projects/ucc/> (2025).
- <sup>8</sup>U.S. Department of Energy, Oak Ridge National Laboratory, Oak Ridge Leadership Computing Facility, “Frontier,” <https://www.olcf.ornl.gov/frontier/> (2022).
- <sup>9</sup>U.S. Department of Energy, Lawrence Livermore National Laboratory, Advanced Simulation and Computing Program, “El Capitan: NNSA’s First Exascale Supercomputer,” <https://asc.llnl.gov/exascale/el-capitan> (2023).
- <sup>10</sup>MPICH Development Team, “MPICH: Initial NCCL Allreduce Backend,” <https://github.com/pmodels/mpich/pull/7298> (2024).
- <sup>11</sup>Ohio State University, “OSU Micro-Benchmarks,” <https://mvapich.cse.ohio-state.edu/benchmarks/> (2024).
- <sup>12</sup>OpenUCX Project, “OpenUCX,” <https://openucx.org/> (2025).
- <sup>13</sup>Nihar Kodkani, “Accelerating AI Workloads: MPI Collectives with NVIDIA’s UCC in MPICH,” Publication Pending (2025).
- <sup>14</sup>R. Thakur and W. D. Gropp, in *Proceedings of Recent Advances in Parallel Virtual Machine and Message Passing Interface*, edited by J. Dongarra, D. Laforenza, and S. Orlando (Springer Berlin Heidelberg, Berlin, Heidelberg, 2003) pp. 257–267.
- <sup>15</sup>P. Balaji, D. Buntinas, D. Goodell, W. Gropp, S. Kumar, E. Lusk, R. Thakur, and J. L. Träff, in *Proceedings of Recent Advances in Parallel Virtual Machine and Message Passing Interface* (Springer Berlin Heidelberg, Berlin, Heidelberg, 2009) pp. 20–30.
- <sup>16</sup>S. Chunduri, S. Parker, P. Balaji, K. Harms, and K. Kumaran, in *SC 2018* (2018) pp. 386–400.
- <sup>17</sup>I. Laguna, R. Marshall, K. Mohror, M. Ruefenacht, A. Skjellum, and N. Sultana, in *SC ’19* (2019).
- <sup>18</sup>M. Wilkins, Y. Guo, R. Thakur, P. Dinda, and N. Hardavellas, in *Proceedings of the 2022 IEEE International Conference on Cluster Computing (CLUSTER)* (2022) pp. 161–171.
- <sup>19</sup>M. Wilkins, Y. Guo, R. Thakur, N. Hardavellas, P. Dinda, and M. Si, in *ExaMPI 2021* (2021) pp. 36–45.
- <sup>20</sup>S. Hunold and A. Carpen-Amarié, in *Proceedings of HPCAsia ’18* (ACM, New York, NY, USA, 2018) pp. 64–74.
- <sup>21</sup>S. Hunold and S. Steiner, in *PMBS 2022* (2022) pp. 123–128.
- <sup>22</sup>Z. Cai, Z. Liu, S. Maleki, M. Musuvathi, T. Mytkowicz, and O. Saarikivi, in *PPoPP ’21* (2021) pp. 62–75.
- <sup>23</sup>J. Pjesivac-Grbovic, G. Bosilca, G. E. Fagg, T. Angskun, and J. J. Dongarra, *Parallel Computing* **33**, 613 (2007).
- <sup>24</sup>S. Potluri, K. Hamidouche, A. Venkatesh, D. Bureddy, and D. K. Panda, in *ICPP 2013* (2013) pp. 80–89.
- <sup>25</sup>S. Potluri, H. Wang, D. Bureddy, A. Singh, C. Rosales, and D. K. Panda, in *IPDPSW 2012* (2012) pp. 1848–1857.
- <sup>26</sup>A. Ayala, S. Tomov, X. Luo, H. Shaeik, A. Haidar, G. Bosilca, and J. Dongarra, in *ExaMPI 2019* (2019) pp. 12–18.
- <sup>27</sup>M. Cho, U. Finkler, M. Serrano, D. Kung, and H. Hunter, *IBM J. Res. Dev.* **63**, 1:1 (2019).
- <sup>28</sup>U.S. Department of Energy, “National Laboratories,” <https://www.energy.gov/national-laboratories> (2025).
- <sup>29</sup>U.S. Department of Energy, “Science Undergraduate Laboratory Internships (SULI),” <https://science.osti.gov/wdts/suli> (2025).
- <sup>30</sup>U.S. Department of Energy, “Frequently Asked Questions,” <https://science.osti.gov/wdts/suli/Frequently-Asked-Questions>

(2025).

## Appendix A: Acronyms

Acronym	Full Form	Description
AI	Artificial Intelligence	Data-driven computational workloads
CCL	Collective Communication Library	GPU backend for collective operations
CPU	Central Processing Unit	General-purpose processor for control and serial tasks
DOE	Department of Energy	U.S. agency funding energy and computing research
FI_HMEM	Fabric Interface for Heterogeneous Memory	Enables GPU memory access over libfabric in HPC
GPU	Graphics Processing Unit	Processor for parallel computation
HBM	High Bandwidth Memory	Type of high-performance memory used in GPUs
HIP	Heterogeneous-computing Interface for Portability	AMD’s GPU programming API
HPC	High-Performance Computing	Large-scale computing for science and engineering
JLSE	Joint Laboratory for System Evaluation	Argonne’s HPC testbed
JSON	JavaScript Object Notation	Lightweight data format
MPI	Message Passing Interface	Communication standard for HPC applications
MPICH	Argonne’s MPI Implementation	Portable, high-performance MPI library
NCCL	NVIDIA CCL	NVIDIA’s GPU collective library
OSU	Ohio State University	Developer of OSU benchmarks and MPI tools
PMRS	Programming Models and Runtime Systems	ANL research group on parallel computing
RCCL	ROCm CCL	AMD’s GPU collective library for ROCm
ROCm	Radeon Open Compute	AMD’s open-source GPU computing platform
SULI	Science Undergraduate Laboratory Internships	DOE research internship program for undergrads
TCP	Transmission Control Protocol	Standard network protocol for inter-node communication
UCC	Unified Collective Communication	Modular GPU backend for collectives
UCX	Unified Communication X	HPC transport and memory abstraction layer

TABLE I: List of acronyms used in this work.