

Distributed Machine Learning

Alex Deng
Yale University
a.deng@yale.edu

Ryan Jin
Yale University
ryan.jin@yale.edu

Jason Zheng
Yale University
jason.zheng@yale.edu

Abstract

Efficient collective communication is a fundamental requirement for distributed systems, especially in machine learning and scientific computing workflows. This project implements a GPU simulation framework that supports collective operations, focusing on the AllReduce algorithm. Two variants are implemented: a Naive AllReduce, which uses a centralized gather-reduce-broadcast approach, and a Ring AllReduce, which distributes communication and computation efficiently across GPUs. The system employs a GPU Coordinator to manage simulated devices, orchestrating data transfers and synchronization through gRPC-based communication.

The design highlights tradeoffs between simplicity and scalability: while the naive approach is straightforward to implement, it suffers from bottlenecks as the number of GPUs increases. In contrast, the ring-based algorithm achieves communication efficiency and scalability by reducing the communication cost. Comprehensive testing validates correctness and performance, offering insights into algorithm behavior across different workloads.

This framework serves as a versatile platform for exploring, implementing, and analyzing collective communication strategies, enabling research and education in distributed GPU systems.

1 Introduction

In distributed systems, collective communication operations such as AllReduce play a critical role in accelerating computation, particularly in machine learning and scientific computing. These operations aggregate data across multiple nodes and distribute the results efficiently. While hardware platforms such as NVIDIA NCCL offer robust implementations for GPUs, research and teaching contexts often require custom implementations to simulate and

analyze collective communication algorithms on a group of GPUs.

This project implements a simulated GPU coordination framework with two variants of the AllReduce operation:

- Naive AllReduce: A centralized gather-broadcast approach.
- Ring AllReduce: An optimized, communication-efficient algorithm.

The system consists of a GPU Coordinator that manages GPU devices and synchronizes their communication. By enabling fine-grained control over data transfers, the system allows for the exploration of design options, tradeoffs, and evaluation of different algorithms in a simulated environment.

2 Design and Implementation

2.1 System Overview

The system is built using **gRPC-based communication** to emulate interactions between multiple GPU devices. Each simulated GPU:

- Exposes an RPC interface for memory transfers (Memcpy), reductions, and status queries.
- Manages a contiguous memory space to store input and output data for each device.

The **GPU Coordinator** orchestrates GPU devices through the following key components:

1. **Communicator Management:** Initializes communicators to group GPUs for operations.
2. **Collective Operations:** Implements **AllReduce** variants:

- **Naive AllReduce** (centralized gather-reduce-broadcast).
- **Ring AllReduce** (distributed scatter-reduce and all-gather phases).

3. **Concurrency Control:** Operations are parallelized using Goroutines with synchronization (e.g., `sync.WaitGroup`).

2.2 Implementation Details

1. Scatter-Reduce Phase (Ring AllReduce):

Each GPU sends a chunk of its data to the next GPU while receiving and reducing data from the previous GPU. Synchronization is managed using Goroutines and a `WaitGroup`.

2. All-Gather Phase (Ring AllReduce):

GPUs redistribute the reduced data chunks sequentially around the ring. Each GPU sends its chunk to the next device until all GPUs receive the full reduced data.

3. Memory Transfers:

The `Memcopy` RPC interface supports both **Host-to-Device** and **Device-to-Host** transfers, enabling verification of results after collective operations.

4. Testing and Validation:

The system includes comprehensive tests to validate data consistency and correctness:

- Initializing communicators.
- Copying input vectors to GPUs.
- Performing AllReduce operations and verifying outputs.

3 Design Options and Tradeoffs

Option 1: Naive AllReduce

- **Design:** All GPUs send their data to a single "root" GPU, which performs the reduction. The root GPU then broadcasts the reduced data back to all devices.

- **Pros:**

- Simple to implement and debug.
- Suitable for small-scale systems with limited communication overhead.

- **Cons:**

- Bottlenecks on the root GPU, making it non-scalable.

- Communication cost grows with the number of GPUs

- **Time Complexity:**

- **Gather Phase:** N GPUs each send S bytes to the root GPU. This results in a total communication cost of $O(N \cdot S)$.
- **Broadcast Phase:** The root GPU sends S bytes of reduced data back to N GPUs, again incurring $O(N \cdot S)$.

Total Communication Cost:

$$O(N \cdot S) + O(N \cdot S) = O(N \cdot S)$$

Option 2: Ring AllReduce

- **Design:** A more communication-efficient approach where:

1. GPUs send and receive **partial sums** in a "ring" topology during the **scatter-reduce phase**.
2. Each GPU redistributes its chunk of the reduced result during the **all-gather phase**.

- **Pros:**

- Each GPU communicates only with its immediate neighbors, and the workload is evenly distributed across all GPUs. This avoids a central bottleneck.
- Scales efficiently as the number of GPUs increases, since each device only sends its data to its neighboring device.

- **Cons:**

- More complex to implement due to synchronization and chunk management.
- Potential latency in small-scale systems where naive reduction suffices.

- **Time Complexity:**

- **Scatter-Reduce Phase:** Each GPU sends and receives $\frac{S}{N}$ data $N - 1$ times (one iteration per neighbor). Total data communicated per GPU: $O(S)$.
- **All-Gather Phase:** Similarly, each GPU sends and receives $\frac{S}{N}$ data $N - 1$ times so incurring $O(S)$ cost.

Total Communication Cost: $O(S) + O(S) = O(S)$

3.1 Reasoning for Choices

We implemented both **Naive** and **Ring AllReduce** to provide a baseline comparison and highlight the performance benefits of the optimized algorithm under different workloads. We wanted to highlight the performance benefits of Ring AllReduce by comparing it to a naive, baseline version of the AllReduce Algorithm. The system’s modular design allows easy extension to other collective algorithms, such as **Tree-based reductions** or pipeline approaches.

Choosing gRPC

- **Modularity:** gRPC provides a well-defined and language-agnostic protocol, making it easier to extend or integrate with other components in the future.
- **Ease of Debugging:** The structured nature of gRPC simplifies debugging communication between simulated devices.
- **Scalability:** While not the most optimized for high-performance systems, gRPC allows for rapid prototyping and can simulate large-scale systems effectively.

Choosing a Contiguous Memory Array

- Using a contiguous memory array for each device simplifies the simulation and enables direct indexing for memory operations.
- Mutex locks ensure thread safety during memory updates in a multi-threaded environment.

Comparisons Between Our Implementations of Naive and Ring AllReduce

We implemented test cases to test the time complexity between our Naive AllReduce and Ring AllReduce implementations. In our findings, we saw that with smaller vector sizes, Ring AllReduce didn’t show that much, if any, improvements in runtime over Naive AllReduce. However, when we increased the size of the vectors in our test, the benefit of Ring AllReduce was clear. When testing with vector sizes of 100,000, 200,000, and 450,000 elements, Ring AllReduce saw performance runtime improvements of **3.56x**, **3.54x**, and **3.81x** over Naive AllReduce, respectively. Ring AllReduce truly showed its performance optimizations when we had large amounts of data that needed to be communicated between different GPU Devices, as well as when we had increased the total number of GPU Devices themselves.

4 Related Work

4.1 Collective Communication Libraries

1. **NVIDIA NCCL** [1]:

NCCL provides high-performance, hardware-optimized collective operations, including Ring AllReduce. While our system is a simulation, we adopt NCCL’s ring-based algorithm for its scalability.

2. **Horovod** [2]:

Horovod is a deep learning framework that uses Ring AllReduce to efficiently train models across distributed GPUs. Our work mirrors Horovod’s algorithmic design but focuses on simulation and modularity for experimentation.

3. **MPI (Message Passing Interface):**

MPI’s collective operations, such as `MPI_Allreduce`, are widely used in HPC systems. The implementation explores similar designs while leveraging modern gRPC for inter-device communication.

5 Conclusion

This project implements a **GPU simulation framework** that supports collective communication operations, focusing on two AllReduce algorithms: Naive AllReduce and Ring AllReduce. The system highlights the tradeoffs between simplicity and scalability, offering insights into communication-efficient algorithms. Future work includes extending the system to support other collective operations (e.g., AllGather, Scatter) and further optimizing concurrency and synchronization. Additionally, we would like to use our framework to actually train a toy model on sample datasets, such as the MNIST dataset.

References

- [1] NVIDIA NCCL: “NVIDIA Collective Communications Library,” <https://developer.nvidia.com/nccl>.
- [2] A. Sergeev and M. Del Balso. *Horovod: fast and easy distributed deep learning in TensorFlow*. arXiv preprint arXiv:1802.05799, 2018.
- [3] Message Passing Interface Forum. *MPI: A Message-Passing Interface Standard*. 1994.
- [4] D. G. Murray, et al. *Ciel: a universal execution engine for distributed data-flow computing*. NSDI, 2011.