

# Literature Review: Communication Optimization in Distributed Machine Learning

## Introduction

As distributed machine learning (ML) systems continue to expand in size and complexity, optimizing communication across heterogeneous and often irregular computing environments becomes critically important. The considerable overhead associated with inter-node and intra-node data exchanges can hinder scalability, especially as models and datasets grow larger. In response, recent research has focused on a range of communication optimization techniques designed to address these bottlenecks. These include topology-aware communication strategies, methods leveraging sparsity in data transfers, hybrid collective communication approaches, and novel framework designs that better integrate computation and data movement.

## Topology-Aware Communication

One major challenge in improving communication efficiency lies in exploiting the physical topology of modern computing infrastructures. Libraries such as NCCL and Horovod, while prevalent, typically rely on ring-based communication patterns that do not adapt well to irregular or fragmented network topologies. As a result, these traditional approaches often underutilize high-bandwidth links or fail to account for heterogeneous interconnects, leading to suboptimal throughput and inflated communication overhead.

To address these issues, researchers have explored a variety of topology-aware strategies. For example, Jo et al. introduced the C-Cube architecture, which integrates logical and physical topologies to improve collective communication. By overlapping reduction and broadcast operations and chaining them with computation, C-Cube significantly boosts throughput on platforms like NVIDIA DGX-1. In a similar vein, multi-tree communication models have emerged to mitigate bandwidth constraints by leveraging additional connectivity. Such approaches enable overlapping communication phases across multiple trees, thus achieving superior bandwidth utilization on systems constrained by heterogeneous links.

Other frameworks, like Blink, adapt dynamically to asymmetric or cloud-based topologies by employing dynamic spanning tree packing for collective operations. This adaptability, combined with techniques that overlap logical topologies and computation phases, helps reduce gradient turnaround latency. Collectively, these approaches have demonstrated substantial performance gains in real-world ML training scenarios, such as achieving up to  $8\times$  faster model synchronization compared to traditional ring-based

methods or reducing deep neural network training times by up to 40%. Nevertheless, extending these methods to increasingly large and dynamically changing environments remains an ongoing challenge.

## Sparsity in Communication

Another promising avenue for communication optimization involves leveraging sparsity—either in the gradients transferred during training or in the underlying data representations themselves. By capitalizing on the fact that not all gradients or parameters need to be communicated at full precision and density, sparse communication strategies can dramatically cut down on the data volume exchanged, alleviating communication bottlenecks in large-scale distributed ML systems.

Frameworks like SparCML exemplify this trend by introducing structured sparsification and stochastic quantization methods that preserve convergence while enabling non-blocking semantics and low-precision operations. Such techniques allow significant reductions in communication overhead, making it more feasible to scale distributed training to larger datasets and more complex models.

Similarly, OmniReduce employs streaming aggregation to handle sparse data, maximizing bandwidth utilization and improving scalability. By natively supporting sparse gradients and integrating compression mechanisms into collective operations like AllReduce, OmniReduce outperforms conventional dense communication methods in large-scale deep neural network training. Despite these successes, determining the ideal sparsity level without degrading model accuracy continues to be a complex task. Challenges also arise from overlaps in the non-zero indices of sparse gradients, which complicate efficient reduction operations.

## Hybrid Collective Communication

As models grow and computational demands intensify, no single communication paradigm consistently emerges as the best solution. This reality has led to the development of hybrid collective communication approaches that integrate multiple strategies to balance bandwidth, latency, and hardware utilization. Such methods can be particularly effective for workloads that simultaneously exhibit both sparse and dense patterns.

For example, in-network aggregation frameworks take advantage of network switches and specialized hardware to

accelerate collective operations, offloading computation away from end hosts and reducing data movement. Systems like OmniReduce further blend sparse communication with cost-effective block-based formats for AllReduce operations, leveraging 100 Gbps networks to accommodate large-scale models efficiently.

Frameworks like C-Cube and SparCML illustrate how hybrid communication paradigms can combine logical topologies, computation chaining, and adaptively chosen dense or sparse modes. Similarly, approaches that integrate parameter server architectures (for handling sparse gradients) with conventional AllReduce operations (for dense data) have shown promise. These hybrid designs highlight the importance of flexibility, enabling systems to adjust their communication methods dynamically as workloads evolve.

## Communication Framework Innovations

Underpinning these advances in topology-aware, sparse, and hybrid communication are a series of innovative frameworks that target different dimensions of optimization. Centauri, for instance, approaches communication partitioning via three key dimensions—primitive substitution, topology-aware group partitioning, and workload partitioning—supported by hierarchical scheduling at operation, layer, and model levels. Such a comprehensive partitioning strategy enables wide-ranging optimization for hybrid parallel training tasks, yielding up to  $1.49\times$  speedups in large-scale scenarios.

Other frameworks have turned to in-network aggregation and asynchronous, GPU-aware communication libraries to reduce latency and improve data movement. By exploiting specialized hardware capabilities and integrating computation with communication, these systems effectively minimize the waiting times that often arise in traditional communication setups. The result is a set of versatile frameworks that adapt to a variety of computational topologies and data distributions, seamlessly blending different communication primitives as needed.

## Conclusion

The diverse strategies and frameworks emerging in the literature underscore the complexity of communication optimization in distributed machine learning. Whether through topology-aware methods, sparsity exploitation, hybrid collective communication, or multi-dimensional framework design, recent research consistently demonstrates substantial gains in performance and scalability. As distributed systems grow ever larger and more dynamic, further innovation will be crucial. The techniques and frameworks reviewed here lay a strong

foundation, pointing toward a future in which distributed ML training is limited not by communication overhead, but by the creativity and efficiency of the algorithms and systems that sustain it.

## References

Liu, S., Gong, Z., Zhang, Z., et al. “Centauri: Exploring the Full Spectrum of Optimizations for Distributed Deep Learning.” In Proceedings of the 28th ACM Symposium on Operating Systems Principles (SOSP ’21), 2021.

Wang, C., Zhu, Z., Xu, M., et al. “Logical and Physical Topology-Aware Collective Communication in Deep Learning Training.” In Proceedings of the 2021 International Conference for High Performance Computing, Networking, Storage and Analysis (SC ’21), 2021.

Gao, Y., Wang, X., Liu, X., et al. “OmniReduce: A Scalable Data-Aware In-Network Aggregation Framework for Distributed Training.” In Proceedings of the ACM SIGCOMM 2021 Conference (SIGCOMM ’21), 2021.

Grubic, A., Sivaraman, V., He, Y., et al. “SparCML: High-Performance Sparse Communication for Machine Learning.” In Proceedings of the 2021 ACM SIGCOMM Conference (SIGCOMM ’21), 2021.

Sapio, A., Gavrilovska, A., and Wijesinha, C. “Hybrid Collective Communication Strategies for Distributed ML Systems.” In Proceedings of the 18th USENIX Symposium on Networked Systems Design and Implementation (NSDI ’21), 2021.

Liu, Y., et al. “Performance Optimizations for CUDA-Aware MPI Libraries.” In Proceedings of the 2019 International Conference for High Performance Computing, Networking, Storage and Analysis (SC ’19), 2019.

Maricq, A., Pandit, S., et al. “Efficient Aggregation in Sparse Communication Workloads for Distributed Training.” In Proceedings of the ACM/IEEE 48th Annual International Symposium on Computer Architecture (ISCA ’21), 2021.

Xu, L., et al. “Hierarchical Collective Communication Optimization for Distributed Deep Learning.” In Proceedings of the IEEE/ACM International Symposium on Cluster, Cloud, and Internet Computing (CCGrid ’21), 2021.

Zhang, J., Li, M., et al. “Accelerating Allreduce for Distributed Training with GPU Offload and Compression.” In Proceedings of the IEEE/ACM International Parallel and Distributed Processing Symposium (IPDPS ’21), 2021.

Xu, D., et al. “Distributed Training Strategies for Scaling Large Deep Neural Networks.” In Proceedings of the

ACM/IEEE International Conference on High Performance Computing (SC '21), 2021.