

Thunderfield

Abhiram Anil

*School of Computer Science
RV University
abhiramabtech24@rvu.edu.in*

Bhyresh B.S.

*School of Computer Science
RV University
bhyreshbsbtech24@rvu.edu.in*

Busetty Sugnesh

*School of Computer Science
RV University
busettsugneshbtech24@rvu.edu.in*

Adithiyaa Kala Kandan

*School of Computer Science
RV University*

adithiyaakalakandanbtech24@rvu.edu.in

Abstract—Large-scale distributed deep learning workloads running on GPU clusters often suffer from energy inefficiencies caused by bursty inter-GPU communication. These short communication bursts increase network congestion, synchronization delay, and power consumption, leading to a higher Energy Delay Product (EDP).

In this work, we propose DL-AMCC, a learning based runtime controller that dynamically forms temporary GPU micro clusters during communication bursts. The system observes traffic and system telemetry and adjusts cluster behavior to reduce cross rack communication and improve energy efficiency.

We evaluate DL-AMCC in a large-scale GPU cluster simulator under realistic distributed training workloads. Results show that our approach reduces Energy Delay Product and communication overhead while maintaining overall training performance.

Index Terms—Deep Learning, Temporal Graph Neural Networks, Reinforcement Learning, GPU Clusters, Distributed Deep Learning, Energy Efficiency, Energy-Delay Product, Inter-GPU Communication, Runtime Optimization, Traffic-Aware Micro-Clustering

I. INTRODUCTION

Modern deep learning models are trained using *large scale* GPU clusters, where multiple GPUs work together to process different parts of the same model. In distributed training, each GPU computes gradients locally and then communicates with other GPUs to synchronize these updates. This communication is required to keep the model parameters consistent across all devices. Operations such as all reduce are commonly used for this synchronization process.

Although distributed training improves computational speed, it introduces heavy inter-GPU communication. These synchronization steps occur repeatedly during training and often create short but intense communication bursts. During such bursts, network traffic increases sharply, GPUs may wait for synchronization to complete, and power consumption can temporarily spike. As cluster sizes grow, these effects reduce energy efficiency and increase the Energy Delay Product.

Most existing cluster management systems make decisions at job launch time and do not adapt to runtime communication dynamics. As a result, systems cannot respond effectively to *burst driven* inefficiencies that arise during training.

To address this limitation, we present **Thunderfield**, an autonomous runtime framework for *burst aware* GPU clus-

ter control. Thunderfield is built around DL-AMCC, a deep learning based control mechanism that dynamically forms temporary GPU micro clusters during communication bursts. By observing traffic and system telemetry, Thunderfield adjusts cluster behavior in real time to reduce *cross rack* communication and improve energy efficiency without degrading performance.

The main contributions of this work are:

- Thunderfield, a *burst aware* runtime micro clustering framework for GPU clusters.
- DL-AMCC, a deep learning based controller that predicts and responds to communication dynamics.
- A *large scale* simulation study demonstrating improvements in Energy Delay Product while maintaining throughput.

II. LITERATURE SURVEY

Efficient management of GPU clusters for distributed deep learning workloads has been a topic of significant research. Prior works have explored cluster scheduling, energy aware optimization, communication primitive tuning, and system telemetry. However, these efforts rarely address the combined challenges of runtime communication dynamics, burst driven inefficiencies, and energy optimization within a unified control framework. Below we discuss major research trends and position Thunderfield against them.

A. Static and Placement Driven Scheduling

Early work such as Tiresias [1] and HiveD [2] focus on optimizing job placement and resource allocation across GPU clusters. Tiresias uses trace based heuristics to improve job completion time through initial placement decisions, while HiveD introduces hierarchical placement cells to balance fairness and utilization.

These systems advance placement quality at the batch (job) level, but they share a common limitation: they do not adjust decisions once training begins. Communication patterns within a job can vary significantly during execution, especially in large models with varying layer sizes or dynamic batch strategies. Static placement cannot address transient inefficiencies that arise during runtime bursts.

B. Energy Aware Scheduling

Energy aware scheduling techniques such as PowerFlow [3] explore the trade off between energy consumption and job completion time under an energy budget. PowerFlow’s contribution lies in introducing energy awareness into the scheduler’s cost model, but it remains focused on scheduler level allocation rather than runtime adaptation.

In particular, PowerFlow treats energy as a constraint in planning rather than as a dynamic signal to optimize throughout execution. Burst driven spikes in communication and power consumption occur within iterations and are invisible to schedulers that operate at coarse temporal scales.

C. Communication and Collective Optimization

Efforts to optimize collective operations and network utilization include NV-Group [4] and Libra [5]. NV-Group tailors reduction kernels to leverage NVLink effectively, reducing traffic volume across slower links. Libra proposes topology aware collective scheduling to reduce communication cost in multi dimensional network configurations.

While these works reduce fundamental communication costs, they do so by improving the efficiency of primitives and leveraging static topology knowledge. They do not provide a mechanism for runtime reshaping of communication groups or adjustment of cluster configuration based on observed traffic patterns.

D. Telemetry and Monitoring Tools

Production clusters can use NVIDIA Data Center GPU Manager (DCGM) [6] to access real time metrics for power, utilization, and thermal state. Such telemetry is crucial for understanding cluster behavior, but existing tools do not translate observed signals into automated control actions. They function primarily as monitoring or alerting systems rather than as closed loop controllers.

E. Gap and Opportunity

The above works offer valuable insights into placement, energy, and communication efficiency, yet they fall short in key areas that Thunderfield targets:

- **Temporal adaptation:** Most systems make decisions at job start time or operate at coarse scheduler timescales, missing short lived communication bursts that occur during iteration level execution.
- **Runtime control:** Energy aware schedulers do not react dynamically to telemetry, whereas communication optimizations do not reshape cluster behavior at runtime.
- **Unified control:** No existing work jointly considers burst prediction, dynamic micro grouping, and energy performance optimization inside a learned control policy.

Thunderfield addresses these gaps by using DL-AMCC, a learned runtime controller that observes traffic and power telemetry to predict and respond to communication driven inefficiencies, forming temporary micro clusters that minimize cross rack traffic and improve the Energy Delay Product without reducing throughput.

REFERENCES

- [1] J. Gu, M. Chowdhury, K. G. Shin, Y. Zhu, M. Jeon, J. Qian, H. Liu, and C. Guo, “Tiresias: A GPU Cluster Manager for Distributed Deep Learning,” in *Proceedings of the 16th USENIX Symposium on Networked Systems Design and Implementation (NSDI)*, 2019.
- [2] H. Zhao, Z. Zhang, C. Delimitrou, and others, “HiveD: Sharing a GPU Cluster for Deep Learning with Guarantees,” in *Proceedings of the 14th USENIX Symposium on Operating Systems Design and Implementation (OSDI)*, 2020.
- [3] D. Gu, X. Xie, G. Huang, X. Jin, and X. Liu, “Energy Efficient GPU Clusters Scheduling for Deep Learning (PowerFlow),” arXiv preprint arXiv:2304.06381, 2023.
- [4] C.-H. Chu, P. Kousha, A. A. Awan, K. S. Khorassani, H. Subramoni, and D. K. Panda, “NV-Group: Link Efficient Reduction for Distributed Deep Learning on Modern Dense GPU Systems,” in *Proceedings of the 34th International Conference on Supercomputing (ICS)*, 2020.
- [5] W. Won, S. Rashidi, S. Srinivasan, and T. Krishna, “LIBRA: Enabling Workload Aware Multi Dimensional Network Topology Optimization for Distributed Training of Large AI Models,” arXiv preprint arXiv:2109.11762, 2021.
- [6] NVIDIA Corporation, “NVIDIA Data Center GPU Manager (DCGM) User Guide,” NVIDIA Documentation. [Online]. Available: <https://docs.nvidia.com/datacenter/dcgm/>