Ruofan Wu (ruofanw), Jeff Ma (jeffjma), Runyu Lu(runyulu)

# Summary of "FastFlow: Accelerating Deep Learning Model Training with Smart Offloading of Input Data Pipeline"

## Problem and Motivation

In modern deep learning (DL) pipelines, one of the critical challenges is the CPU bottleneck during data preprocessing. As DL models become more complex and data-centric (requiring operations like decoding, augmenting, and batching), and the performance gap between CPU and GPU cycles increases, CPUs often struggle to keep up with the fast computational capabilities of GPUs. This leads to preprocessing stalls, where GPUs are left idle, waiting for the CPUs to prepare the data. These stalls result in underutilized GPUs and longer training times, making the training process inefficient.

Existing solutions, such as DALI or FPGA-based offloading, rely on manual configuration or require translating CPU operations to hardware-specific operations, making them inflexible and time-consuming to implement. With manual offloading using tools like `tf.data.service`, the users must define how to distribute tasks across remote CPUs. This leads to suboptimal offloading and inefficiencies due to a lack of automation and precision in distributing preprocessing tasks. FastFlow is motivated by the need to solve these inefficiencies in DL pipelines by proposing automatic offloading of input pipeline to remote CPUs

## Related Works

1. Alternative Compute

DALI is a GPU-based solution for offloading preprocessing tasks from the CPU to GPU. DLBooster and TrainBox are FPGA-based solutions that use FPGAs to handle preprocessing tasks. They require manual effort to adapt CPU-based operations to the specialized hardware, which is time-consuming and inflexible.

2. Auto-tuning on local node

TensorFlow's `tf.data` API provides a way to optimize data loading pipelines on local nodes by using techniques like prefetching, parallelism, and caching. Plumber is an analysis and tuning tool for optimizing the parallelism and buffer sizes of input pipelines. Although these techniques

improve local performance, they are constrained by the local node's resources and cannot scale to remote or distributed systems for the large-scale datasets or complex preprocessing tasks.

3. Remote CPU Offload

TensorFlow's `tf.data.service` employs a dispatcher/worker model for manual offloading of preprocessing tasks to remote CPUs. However, the users must decide which tasks to offload and how to distribute the load. This often results in suboptimal offloading , as users may not have the expertise to make optimal offloading decisions. The lack of automation and precision limits its performance.

# Solution Overview

FastFlow is designed to be auto-tuning and smart offloading input pipelines to remote CPUs to improve GPU utilization and reduce training time, therefore addressing CPU bottlenecks during input data preprocessing in DL pipelines. It profiles the application and resource environment to decide when, which operations, and how much data to offload, ensuring a balance between both local and remote resources. The system is designed on top of TensorFlow and allows users to leverage these optimizations without modifying their main training logic. FastFlow's approach prevents inefficiencies that arise from naive offloading by dynamically adjusting based on the performance metrics through profiling, ultimately optimizing DL training throughput. The evaluations demonstrate great generalization ability for matching the state-of-art performance under various workloads.

More specifically, FastFlow measures throughput metrics before starting the actual training with the whole dataset to determine the existence of prep stalls. If a prep stall exists, the best offloading pipeline is selected from several candidate pipelines based on estimated training throughput and prep stalls. Finally, an optimal data offloading ratio is decided in consideration of the local and remote computing/network capacity, threading overheads, and offloading overheads.

# Limitations

The limitations of the FastFlow system include the following:

- The performance of smart offloading is very unlikely to be optimal as the operations are coarse-grained.
- FastFlow's design assumes a fixed allocation of resources per job, which might not be suitable for dynamic environments.
- The profiling metrics of 50 ~ 100 steps are from empirical experiments, which may be workload dependent. That is, this value may not be the best for training LLMs or other models. Also, this may introduce high overhead for smaller models little preprocessing stalls. (Although this inflexibility is probably not considered by the authors, the decision to use smart offloading is not "smart.")

# Future Research Directions

1. **More Fine-Grained Control**:
   - Develop techniques to optimize specific preprocessing operations with more fine-grained control, allowing for better customization and efficiency in handling diverse data types and preprocessing tasks.
2. **More Optimal Offloading Strategies**:
   - Investigate and implement more optimal offloading strategies that dynamically adjust based on real-time performance metrics and resource availability, ensuring efficient utilization of both local and remote resources.
3. **Workload-Specific Pipelines**:
   - Design and evaluate workload-specific pipelines tailored to different types of deep learning tasks, such as NLP versus image recognition, to maximize preprocessing efficiency and overall training performance.
4. **Cross-Node Communication and Data Synchronization**:
   - Address challenges related to cross-node communication, data synchronization, and network bottlenecks to ensure seamless and efficient data transfer between local and remote CPUs, minimizing latency and maximizing throughput.

# Summary of Class Discussion

1. **Models and Preprocessing Stalls**:
   - **Question**: How do preprocessing stalls differ for LLMs compared to other models?
   - **Answer**: Different models have varying preprocessing requirements, but they all need operations like batching tensors and prefetching.
2. **User-Defined Processing**:
   - **Question**: Is user-defined processing more flexible?
   - **Answer**: Generalization, accessibility, and trade-offs need to be considered to balance flexibility and performance.
3. **Remote CPU Job Handling**:
   - **Question**: Does the remote CPU wait until it gets offloaded jobs?
   - **Answer**: Prefetching cannot occur until the remote CPU finishes its job, requiring parallelism between batches.
4. **Workload-Specific Pipeline Example**:
   - **Question**: Can you provide an example of a workload-specific pipeline?
   - **Answer**: NLP tasks compared to image recognition tasks require customized pipelines based on the form of data.
5. **Profiling vs. Analytical Component**:
   - **Question**: Is profiling better than an analytical component?
   - **Answer**: Profiling has empirically negligible overhead and provides real-time performance insights.
6. **Remote Cluster GPU Availability**:

- ○ **Question**: Does the remote cluster have GPUs?
- ○ **Answer**: No, profiling happens on both local and remote CPUs.
7. **Using the Best Case**:
   - ○ **Question**: Why wouldn't we just use the best case scenario?
   - ○ **Answer**: The contribution lies in performance-matching generalization, ensuring the system adapts to various workloads and conditions.

# Summary of "Rail-only: A Low-Cost High-Performance Network for Training LLMs with Trillion Parameters"

## Problem and Motivation

As Large Language Models (LLMs) continue to grow in size, the demand for computational resources has exceeded the advancements in hardware accelerators, leading to the need for hyper-scale GPU data centers. However, the current design of these data centers, particularly the network infrastructure, is inefficient for LLM training workloads.

Traditional GPU clusters rely on Clos networks with spine switches to provide any-to-any connectivity. While this architecture supports full-bisection bandwidth, it is excessive for the sparse traffic patterns exhibited by LLMs, where most communication is within high-bandwidth (HB) domains. Traffic analysis reveals that a very small percentage of GPU pairs are involved in significant data transfers, and most traffic remains confined within local domains or "rails". Only minimal traffic crosses different HB domains, so a complex network infrastructure is unnecessary.

The Rail-only network is motivated by the need to design a more efficient and cost-effective network architecture tailored to the specific communication requirements of LLM training.

## Related Works

1. Clos Networks

These multistage switching networks provide any-to-any connectivity between servers through spine and leaf switches, offering full-bisection bandwidth.

2. High-bandwidth Interconnects (HBIs) and Intra-domain Communication

Recent advances in HBIs, such as Nvidia NVLink or AMD Infinity Fabric, have enabled high-speed intra-domain communication within small groups of GPUs, called HB domains.

3. Parallelism Strategies in Large-Scale Neural Networks

Tensor Parallelism (TP), Data Parallelism (DP), and Pipeline Parallelism (PP) are common parallelization strategies used to distribute LLM workloads across multiple GPU.

4. Mixture-of-Experts (MoE) Models and All-to-All Traffic

Mixture-of-Experts (MoE) models introduce a more complex all-to-all communication pattern, where each expert is distributed across different GPUs, creating more frequent inter-domain communication.

5. Hierarchical Collective Communication Algorithms

Hierarchical collective communication strategies, such as hierarchical AllReduce and AllGather, have been explored to reduce inter-node traffic by creating intermediate aggregation points. These techniques are designed to optimize traffic in distributed systems and to minimize latency.

## Solution Overview

This paper presents a new network architecture: Rail-only. From the insight above in the motivation section, LLM training does not require full any-to-any connectivity across all GPUs in a datacenter because of sparse communication patterns. The Rail-only architecture eliminates the spine layer of traditional Clos networks, connecting only those GPUs that engage in significant communication. When training dense LLMs, TP/DP communications are achieved Hierarchical Collective Communication Algorithms that allow networks traffic to stay within a rail. For sparse models, they also prove that the Rail-only network has a very small slow-down factor for all-to-all traffic, which is also a small portion of total traffic in training MoE models.

## Limitations

First, the paper lacks empirical evaluation, making it challenging to convince the audience of its effectiveness. Additionally, the Rail-only architecture is less suitable for dynamic workloads that require reconfiguration or flexibility, such as multi-job training scenarios. Moreover, the design includes redundant connections between GPUs, which could be optimized further. Lastly, although the overhead during sparse model training may be tolerated, the Rail-only network will be constrained if the portion of the all-to-all communications increases.

## Future Research Directions

1. **Empirical Evaluation and Benchmarking**:
   - Conduct comprehensive empirical evaluations of the Rail-only architecture using real-world LLM training workloads, which includes measuring performance

metrics such as training time, latency, and throughput across various model sizes and configurations.

2. **Elasticity and Scalability**:
   ○ Investigate the elasticity of the Rail-only network, particularly how it handles dynamic changes in workload and GPU availability. Develop strategies to dynamically reconfigure the network to maintain performance as the number of GPUs scales up or down.

3. **Redundancy and Fault Tolerance**:
   ○ Explore methods to introduce redundancy in the Rail-only network to enhance fault tolerance. This includes designing mechanisms to handle failures of individual GPUs or network links without significantly impacting overall performance.

4. **Optimization for Different Parallelism Strategies**:
   ○ Study the impact of the Rail-only architecture on various parallelism strategies beyond Tensor Parallelism (TP), Data Parallelism (DP), and Pipeline Parallelism (PP). This includes exploring its effectiveness for emerging parallelism techniques and hybrid models.

5. **Application to Other Model Architectures**:
   ○ Extend the Rail-only network design to other types of models beyond LLMs, such as convolutional neural networks (CNNs) and graph neural networks (GNNs). Evaluate its performance and cost-effectiveness for these different workloads.

6. **Integration with Hierarchical Collective Communication**:
   ○ Investigate the integration of hierarchical collective communication algorithms with the Rail-only architecture. This includes optimizing hierarchical AllReduce and AllGather operations to further reduce inter-node traffic and latency.

7. **Energy Efficiency and Sustainability**:
   ○ Conduct studies on the energy efficiency and environmental impact of the Rail-only network. Develop techniques to further reduce power consumption and improve the sustainability of hyperscale GPU data centers.

8. **Advanced Routing and Scheduling Algorithms**:
   ○ Develop advanced routing and scheduling algorithms tailored to the Rail-only architecture. These algorithms should optimize data transfer paths and minimize communication overhead, especially for complex all-to-all traffic patterns in Mixture-of-Experts (MoE) models.

# Summary of Class Discussion

1. **Actual Experiments**:
   ○ **Question**: Are these actual experiments?
   ○ **Answer**: Yes, they used 384 DGX-A100 nodes, each with 8 GPUs.
2. **Parallelism Strategies (TP, DP, PP)**:
   ○ **Question**: How do they implement Tensor Parallelism (TP), Data Parallelism (DP), and Pipeline Parallelism (PP)?

- ○ **Answer**: TP is implemented within a single node, while DP and PP are used across different nodes, which is considered the optimal strategy.
3. **Network Topology**:
    - ○ **Question**: Did they change the topology?
    - ○ **Answer**: No, they did not change the topology.
4. **Rate Calculation**:
    - ○ **Question**: What does the equation y(x-1) = x(y-1) represent?
    - ○ **Answer**: For x = 8 in A100, increasing y results in a maximum rate of $\frac{7y}{8(y-1)}$, which is acceptable.
5. **Model Fitting in a Single GPU Node**:
    - ○ **Question**: How do you fit your model in a single GPU node?
    - ○ **Answer**: Pipeline parallelism is used across GPU nodes.
6. **Rail-Only and Rail-Optimization Evaluation**:
    - ○ **Question**: What about the evaluation of training time for rail-only and rail-optimization?
    - ○ **Answer**: The training time is not known as they did not actually perform the training.
7. **Application to Other Models**:
    - ○ **Question**: Can this be applied to models other than LLMs?
    - ○ **Answer**: Yes, it can be applied as long as all-to-all operations are ensured. Fully Sharded Data Parallel (FSDP) could also work if there are some collective operations, providing gains.