

# CS598AIE Paper Review:

## GPipe: Easy Scaling with Micro-Batch Pipeline Parallelism

Jiaqi Lou (Netid: jiaqil6)

September 11, 2024

### **1 The problem the paper is trying to tackle.**

The paper tries to address the challenges associated with the model size scaling of deep neural networks, out of which increased model size demonstrates promising quality and accuracy improvements. However, as the model size increases, hardware constraints such as memory limitations and communication bandwidth of the accelerators like GPUs and TPUs and existing model partition approaches are task-specific. GPipe tackles these challenges by providing a general and flexible pipeline parallelism library that can efficiently divide and distribute large models to multiple accelerators.

### **2 What's the impact of the work, e.g., why is it an important problem to solve?**

The work is impactful because it addresses the critical challenges (discussed in [section 1](#)) in efficiently scaling deep neural networks. Since larger model capacity improves the model accuracy, the flexible pipeline parallelism library can not only improve the training quality but also achieve high training efficiency and speedup with several design choices and optimizations.

### **3 The main proposed idea(s). + A summary of your understanding of different components of the proposed technique, e.g., the purpose of critical design choices.**

- Pipeline parallelism: GPipe partitions deep neural networks into sequences of layers and each group of layers can be further partitioned into smaller cells such that the model can be distributed across multiple accelerators to tackle the memory limitation on a single accelerator.
- Micro-batches: GPipe splits mini-batches into micro-batches to allow deeper pipelining and better resource utilization because of its finer-grained execution.
- Synchronous gradient update and communication overhead: GPipe pipelines the execution of micro-batches before single synchronous gradient update synchronization of the entire mini-batch. This approach not only minimizes the bubbles but also reduces the communication overhead. This is because the inter-accelerator communication only happens at the partition boundaries for each micro-batch and the latency overhead may be hidden by overlapping with computational tasks.
- Generality and flexibility: GPipe is not task-specific and can be applied to all layered models. Although there is considerable overhead for model partitioning and splitting the micro-batches, its flexibility and generality can offer "easy-to-use" benefits for users.

## 4 Your perceived strengths and weaknesses of the work, e.g., novelty, significance of improvements, quality of the evaluation, easy-to-use.

### 4.1 Strengths

- Generality and flexibility. As is mentioned in [section 3](#), GPipe as an open-sourced library, is easy-to-use for users.
- Scalability: Gpipe also shows the linear scalability to larger model sizes with the number of accelerators used.
- GPipe demonstrates promising speedup and efficiency improvements with all its design choices. The tradeoffs are discussed in detail which are insightful.

### 4.2 Weaknesses

- Although communication overhead is discussed several times in the paper, it is still not very clear to me how they measure the communication overhead. I think they are training the model on a single host, I wonder how GPipe works with multiple hosts connected with high-speed interconnections.
- GPipe assumes a single layer should fit into the memory of a single accelerator. This assumption can limit the use cases of the framework. And spreading a single layer to multiple layers may introduce further overhead.

## 5 Is there room for improvement? If so, which directions you may want to explore or idea you have for improving the techniques?

I would like to see more breakdowns in turns of computation time, communication time, etc. I think the paper assumes that all accelerators are identical (i.e., the same GPU model). I wonder if GPipe can be extended to run on a heterogeneous platform where GPUs may have different memory capacities and communication bandwidth or other specialized accelerators have significant speedup for some operations in the pipeline.