

---

# Week 7 Reading Summaries: GPipe + Alpa

---

**Amogh Patankar**

University of California, San Diego  
3869 Miramar Street Box #3037, La Jolla, CA- 92092  
apatankar@ucsd.edu

## Abstract

With these readings, we dive deep into GPipe and Alpa, two methods to use parallelism in order to scale the sizes of workloads executed. GPipe uses a strategy called micro-batch parallelism, resulting in an almost linear speedup. Alpa uses inter-parallelism and intra-parallelism for distributed deep learning, outperforming hand-tuned systems, and generalizing to models as well.

## 1 GPipe

### 1.1 Introduction

Deep learning has improved due to the scale of neural networks, and is especially true for computer vision tasks. However, scaling these networks brings up challenges such as memory constraints, and bandwidth issues. The solutions cannot generalize well; GPipe allows to efficiently train a neural network, by creating microbatches. Complementing GPipe with data parallelism is the best way of scaling training.

### 1.2 GPipe Library

The GPipe library has many design features.

#### 1.2.1 Interface

GPipe's interface allows user specification for a cost estimation function for a neural network. That is, the user only specifies model partitions, # of micro-batches, and sequences/definitions of their layers.

#### 1.2.2 Algorithm

GPipe's algorithm partitions the network using the number of partitions specified. It then divides mini-batches into micro-batches during forward pass, and computes gradients for each during the backward pass. Finally, gradients are all accumulated and update the model parameters.

#### 1.2.3 Performance Optimization

GPipe allows only output activations to be stored in forward passes, and recomputations of the forward function in backward passes. GPipe has a low communication overhead, and scaling performance can be boosted on any kind of accelerators due to passing activation tensors at boundaries.

## 1.3 Performance Analyses

With respect to performance analysis of GPipe, they evaluate based on AmoebaNet CNN and Transformer Seq2Seq model. Notably, they show that with a higher number of micro-batches, the overhead is negligible, and that with a small M value, the overhead is no longer negligible.

## 30 1.4 Image Classification

31 The authors scale AmoebaNet using GPipe, and show that scaling through GPipe results in higher  
32 accuracy in six of the eight tests they run on various image classification datasets.

## 33 1.5 Massive Massively Multilingual Machine Translation

34 The authors use GPipe to scale neural machine translation models for NLP tasks, using a multilingual  
35 dataset of 25 billion training examples across 102 languages and English. They show the ability  
36 to handle both low-resource and high-resource languages, demonstrating that a single NMT model  
37 can learn mappings and outperform bilingual models. Comparing a 1.3B wide model and a 1.3B  
38 deep model, the deeper model significantly outperforms the wide model on low-resource languages,  
39 meaning that increasing depth improves generalization. Training deep models has its challenges  
40 including sharp activations and dataset noise, causing training instability.

## 41 1.6 Design Features and Tradeoffs

42 Single Program Multiple Data (SPMD) methods like Mesh-Tensorflow can enable linear scaling  
43 of model parameters but have high communication overhead and limited applicability. In con-  
44 trast, pipeline parallelism approaches like PipeDream struggle with weight staleness and memory  
45 inefficiency. GPipe minimizes communication overhead and bubble inefficiency by pipelining micro-  
46 batches and applying synchronous gradient updates. The only assumption is single-layer memory  
47 constraints.

# 48 2 Alpa

## 49 2.1 Introduction

50 Deep learning advances recently, have been simply to increase the model size. The issue here is that  
51 it requires manual tuning and scaling models efficiently is very complex. The authors automate the  
52 parallelization of large-scale models, to speed up ML research. Moreover, parallel processing can  
53 occur hierarchically, through intra- and inter-operator parallelism.

## 54 2.2 Background: Distributed Deep Learning

### 55 2.2.1 Conventional View of ML Parallelism

56 There are four approaches: data parallelism, operator parallelism, pipeline parallelism, and a com-  
57 bination of parallelisms. In essence, all these parallelisms are used to compute workloads in some  
58 parallel schema, and a combination of parallelisms means you have a configuration to follow.

### 59 2.2.2 Intra- and Inter-Operator Parallelisms

60 Intra-operator parallelism refers to partitioning tensors along a dimension, assigning the computations  
61 to multiple devices, and execute those portions simultaneously. Inter-operator parallelism means that  
62 different operators execute on different distributed devices. Both have differing granularities, but  
63 both fit within compute cluster structure.

## 64 2.3 Overview

65 Alpa generates model-parallel optimization plans at both parallelism levels. Alpa minimizes interop  
66 parallelization latency, and minimizes the execution cost for intraop parallelization.

## 67 2.4 Intra-Operator Parallelism

68 Alpa optimizes parallelism within a device mesh, namely, the SPMD intraop parallelism.

#### 69 **2.4.1 Space of Intra-Operator Parallelism**

70 The authors note that a 2D device mesh with different parallelization strategies, and HW tradeoffs  
71 occur. IN addition, sharding defines partitioning of axes.

#### 72 **2.4.2 ILP Formulation**

73 Integer Linear Programming minimizes the total execution cost of a computational graph including  
74 compute, communication, and resharding. Communication and compute are both estimated based  
75 on bytes and bandwidth with lightweight operators merged to simplify the graph. After ILP, post-  
76 optimization techniques like reduce-scatter and all-gather as opposed to all-reduce are used to reduce  
77 communication and computation overhead.

### 78 **2.5 Inter-Operator Parallelism**

#### 79 **2.5.1 Space for Inter-Operator Parallelism**

80 Using a computational graph, we define the latency as the minimized ILP reported by the intra-op  
81 pass. We try to solve this by adding a constraint for the forward pass, and backward latency.

#### 82 **2.6 Parallelism Orchestration**

83 Parallel orchestration is used to address cross-mesh communication; cross-mesh sharding and pipeline  
84 execution instruction generation are used. Cross-mesh sharding is responsible for sending/receiving,  
85 and gathering instructions in cross-mesh, while pipeline execution instructions are generated by Alpa.

#### 86 **2.7 Limitations and Discussion**

87 The author talk about how Alpa’s hierarchical view of inter- and intra-operator parallelisms gives  
88 three flexibilities, namely, random # of layers, various shape maps, and non-uniform parallelism  
89 configuration.

### 90 **2.8 Evaluation**

#### 91 **2.8.1 End-To-End Performance**

92 The authors evaluate Alpa by using 16K LoC and 6K LoC in Python and C++, with Jax and XLA.  
93 Namely, they show that with an exceeding # of GPUs, Alpa outperforms Megatron-LM, interop only,  
94 and intraop only- for GPT, MoE, and ResNet architectures, using metrics like PFLOPS.

#### 95 **2.8.2 Intra-Op Parallelism Ablation Study**

96 In this study, the authors compare the ILP solution with alternatives like ZeRO; they note that  
97 vanilla data parallelism runs out of memory, ZeRO-methods don’t optimize, while ILP is the highest  
98 performing.

#### 99 **2.8.3 Inter-Op Parallelism Ablation Study**

100 In this study, the authors compare the DP solution with equal layer and equal operator; they note  
101 that the DP method always outperforms equal operator, while DP vs equal layer depends on model  
102 architecture. This is dependent on model and compute, both simultaneously.

#### 103 **2.8.4 Compilation Time**

104 Alpa scales very well to large models and/or clusters, as compilation time grows linearly with model  
105 size and # of GPUs.

### 106 **2.9 Cross-Mesh Resharding**

107 They enable all-gather optimizations to increase speedup.

## 108 **2.10 Wide-ResNet**

109 In the case of Wide-ResNet, Alpa uses intra-operator parallelism, partitioning along the batch axis  
110 for the few layers. For higher quantities of GPUs, Alpa slices the model into more stage, with more  
111 GPUs assigned to later stages.

## 112 **3 Related Work**

113 The authors discuss systems for DP training like Horovod, systems for MP training like PipeDream,  
114 and more advanced techniques to optimize DL execution through compilers.