

CO-LOCATION OF ML JOBS IN GPU CLUSTERS

by

Iman Tabrizian

A thesis submitted in conformity with the requirements
for the degree of Master of Applied Science
Graduate Department of Electrical and Computer Engineering
University of Toronto

© Copyright 2020 by Iman Tabrizian

Abstract

Co-location of ML Jobs in GPU Clusters

Iman Tabrizian

Master of Applied Science

Graduate Department of Electrical and Computer Engineering

University of Toronto

2020

I would like to thank my parents and my dear brother for their amazing support. Without their support, I would not been able to do anything. I want to sincerely thank my advisor Professor Alberto Leon-Garcia for giving me the ability to work independently. His insightful comments have greatly impacted this project. I would like to thank Kaveh Mahdavian and Professor Bianca Schroeder for their time and insightful discussions.

Contents

1	Introduction	1
2	Motivation and Related Work	3
2.1	GPUs	3
2.1.1	GPU Memory and Compute Capabilities	4
2.1.2	CUDA	5
2.1.3	Thread Block Scheduling	5
2.1.4	Life Cycle of a GPU accelerated Application	7
2.1.5	CUDA Streams	7
2.1.6	Multi-Process Service (MPS)	8
2.2	Related Work	8
3	Scheduler Design	12
3.1	Co-location Effect	12
3.2	Why Some Jobs are not Compatible with each other?	12
3.3	Scheduling algorithm	12
4	Evaluation	13
5	Future Steps	14
	Bibliography	15

List of Tables

List of Figures

2.1	GPU Architecture	3
2.2	GPU Capabilities over time	4
2.3	Grid Block vs Thread Block vs Thread	6
2.4	Lifecycle of a GPU accelerated application	7
2.5	Interaction between CUDA contexts and Work Queues on a GPU	8
2.6	Memory Allocation in Forward and Backward Pass	9

Chapter 1

Introduction

Machine Learning has transformed the world. Nowadays, there are various machine learning models used in production systems. Models that provide image classification, movie recommendations, and even photography. These advancements have been made possible thanks to the specialized hardwares such as GPUs, TPUs[7], and FPGAs. The process of training these models is a trial-and-error approach combined with intuition. It requires tuning a large number of hyperparameters to achieve the desired accuracy. To achieve the best accuracy in the shortest amount of time, usually many similar jobs are dispatched that only change a single hyperparameter. Research institutes and companies have access to compute clusters that provide 100s to 1000s of GPUs. It will require a scheduler to map the jobs to the resources. The users of the jobs need to specify the resource requirements of the jobs too. The scheduler will take in these specifications and map it to the the available resources. GPUs with their unique characteristics add a new dimension to this problem. The main problem with the GPUs is that they were not originally designed for multi-tasking. Unlike conventional resources like CPUs and main memory, that have hardware support for resource sharing, GPUs were originally designed for single process use. The assumption was that the application is able to effectively use all the resources on a GPU and adding another application to the GPU

will hinder the performance. However, starting from the Volta[3] architecture GPUs now include hardware support for execution of multiple processes together.

In this dissertation, we explore the effect of co-locating multiple deep learning training jobs together. We devise metrics that guide the scheduler on which workloads can benefit from co-location and which workloads should not be placed together as it will decrease the training speed for these jobs. Additionally, we study the potential reasons behind the incompatibility of the workloads and how they can be used to predict the incompatible jobs.

The next chapter gives an overview of the scheduling algorithms in other domains and an introduction on how GPUs are used for general processing. Chapter 2 discusses the motivation for co-location of the DL jobs. Chapter 3 discusses the design of our proposed scheduler that employs co-location to improve the job completion time (JCT). Chapter 4 provides the experimental results for this design and how it compares to the schedulers that do not employ co-location. Finally, Chapter 5 concludes the thesis and provides future directions on how this work can be extended.

Chapter 2

Motivation and Related Work

2.1 GPUs

GPUs were among the first accelerators to be introduced alongside with CPUs. GPUs are throughput optimized in contrast to the CPUs which are latency optimized. They consist of many cheap cores that can perform large number of operations in parallel. They also have a large memory bandwidth that helps them bring data into these cores in efficiently.

Figure 2.1 shows a simple illustration of the GPU architecture. GPUs consist of several *Streaming Multiprocessors (SMs)*. SMs are the processing units in GPUs. Each

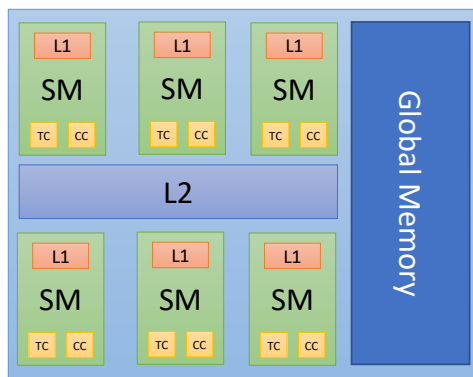


Figure 2.1: GPU Architecture

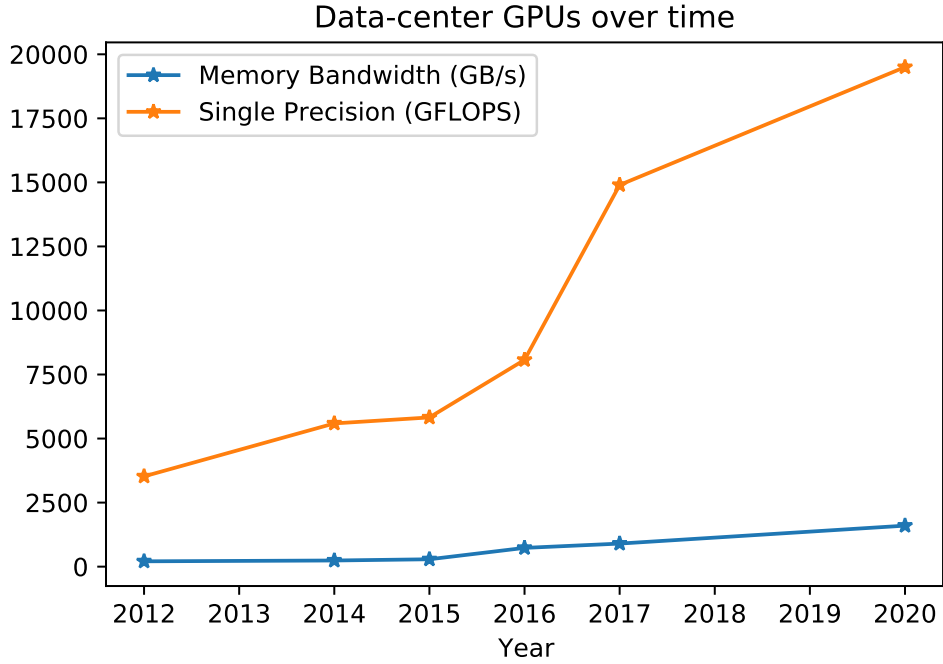


Figure 2.2: GPU Capabilities over time

of the SMs contains various cores designed for operations on different data types. For example, V100 GPU contains 84 SMs where each of them has 64 FP32 cores, 64 INT32 cores, 32 FP64 cores, and 8 tensor cores[3]. In fig. 2.1, *CC* refers to the CUDA cores which consists of all the cores present in each SM except the tensor cores. Tensor Cores are abbreviated using *TC*. They provide accelerated performance for the reduced precision operations which are present in the Deep Learning workloads. They were introduced in the Volta[3] microarchitecture in 2017.

2.1.1 GPU Memory and Compute Capabilities

GPUs memory bandwidth and compute capabilities have increased by a tremendous amount over time. Figure 2.2 shows this trend in the Tesla GPU class. Tesla GPUs are NVIDIA’s data center GPU class. The most recent data center GPU class as of writing this thesis is the A100 GPU with the ability to perform 312 TFLOPS half-precision operations. These increases in the compute and memory capabilities have made it harder

for application developers to fully saturate the GPU resources. In this thesis, we use co-location to better utilize GPUs even if they do not fully utilize the GPU individually.

2.1.2 CUDA

CUDA is a set of extensions to C/C++ to enable easier application development in GPUs. CUDA also introduces a set of language abstractions that make it easier to think about GPU programs. GPU accelerated programs use many threads to perform the computation. CUDA groups the threads into *thread blocks* and *grid blocks*. A *thread block* is a group of threads which are guaranteed to be running on the same SM. A *grid block* is a group of thread blocks which contain all the processing necessary for the computation of a given *kernel*. A *kernel* is a function that runs on a GPU. Both thread blocks and grid blocks can be represented using three dimensions. Figure 2.3 shows the difference between a thread block, grid block, and a thread.

2.1.3 Thread Block Scheduling

There are various constraints that limit the scheduling of thread blocks into the SMs. Amount of the registers, shared memory, and number of threads inside a thread block are among the factors that limit the number of blocks that can be scheduled into a single SM. Since there is a limited amount of these resources available in each SM, the number of thread blocks will be limited to the available resources. Apart from that, different GPU architectures have hard limits on the number of thread blocks and threads that can be scheduled on a given SM. All these factors lead to a metric called *Theoretical Occupancy*[1] of a kernel. Theoretical Occupancy is a metric in percent which determines the percentage of active warps in comparison with the total warps that could be scheduled on a given GPU. There is another concept called *Achieved Occupancy*. Achieved Occupancy measures the scheduled number of warps when the kernel is actually running on the GPU. This can be different from the Theoretical Occupancy because a given thread

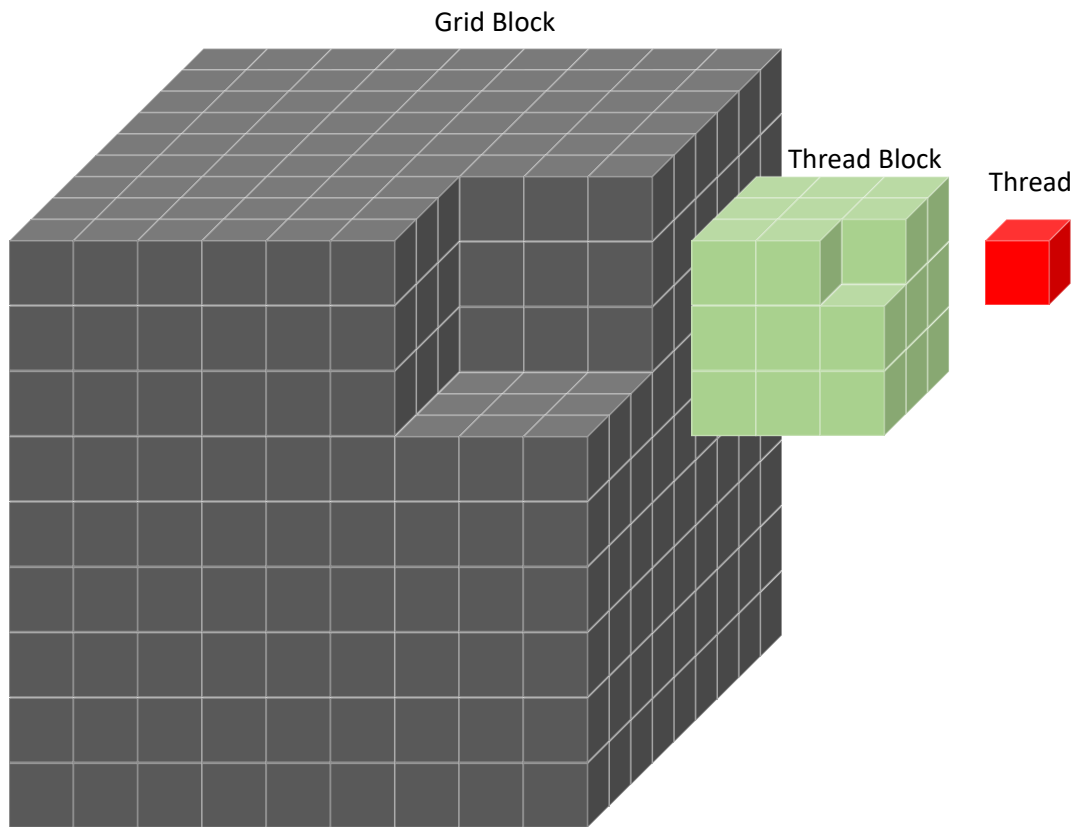


Figure 2.3: Grid Block vs Thread Block vs Thread

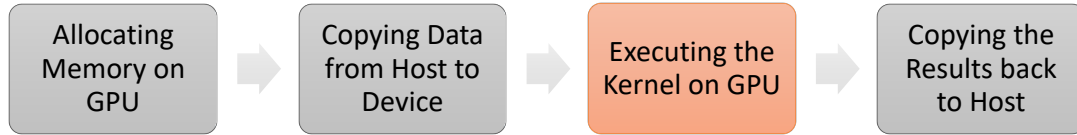


Figure 2.4: Lifecycle of a GPU accelerated application

in warp might be stalled on a memory load and is not yet ready to be scheduled. If there is not enough warps in flight ready to be scheduled instead of the stalled thread block, the achieved occupancy will be lower than the Theoretical Occupancy. Theoretical Occupancy serves as the upper bound for the Achieved Occupancy.

2.1.4 Life Cycle of a GPU accelerated Application

Figure 2.4 shows lifecycle of a typical CUDA application. The application starts with allocating memory on the GPU. Then, it will copy data from the host memory into the GPU memory. After that, the kernel required to run the computation is executed. When the computation is complete, all the results are copied back into the host memory. All these operations must be executed inside a *CUDA context*. Usually there is one CUDA context associated with each process. In the *Exclusive mode*, GPUs give exclusive access to a single CUDA context but in the *Default mode* work submitted from multiple CUDA contexts to the GPU will be scheduled in a time-sharing manner. MPS [2] allows multiple CUDA contexts to run applications on the GPU concurrently. This is explained in more details in section 2.1.6.

2.1.5 CUDA Streams

CUDA Stream is a software construct containing a series of commands that must be executed in order. Work from different streams can be executed concurrently. Recent GPUs are capable of executing work concurrently from different CUDA streams belonging to the same CUDA context. Without MPS [2], it is not possible to run commands from another CUDA context unless the work from the current CUDA context has finished. A

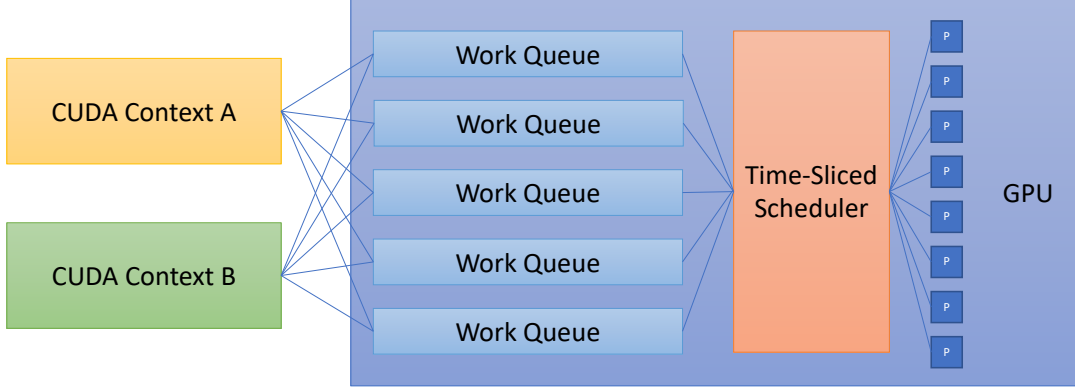


Figure 2.5: Interaction between CUDA contexts and Work Queues on a GPU

common use case for CUDA streams is overlapping computation and communication to speedup the Kernel execution.

2.1.6 Multi-Process Service (MPS)

MPS [2] is a mechanism that enables packing multiple processes together without having to time-share the GPU. MPS achieves this by using a client-server architecture. All the processes that want to run on the GPU are submitted to the the MPS server. MPS is useful when an individual job is not able to saturate all the GPU resources. Before Volta, MPS could not isolate the memory address of different CUDA contexts running on the same GPU. After Volta MPS has improved the address space isolation along with improved performance through hardware support for MPS.

Figure 2.5 shows how CUDA contexts interact with the GPU to schedule work. GPUs have a hardware construct named *Work Queue*. Different CUDA contexts

2.2 Related Work

In the area of GPU schedulers for deep learning workloads there are various related works to the work presented in this thesis. Optimus [11] is one of the earlier works in this area. The main goal of this work is which job should be given more workers so

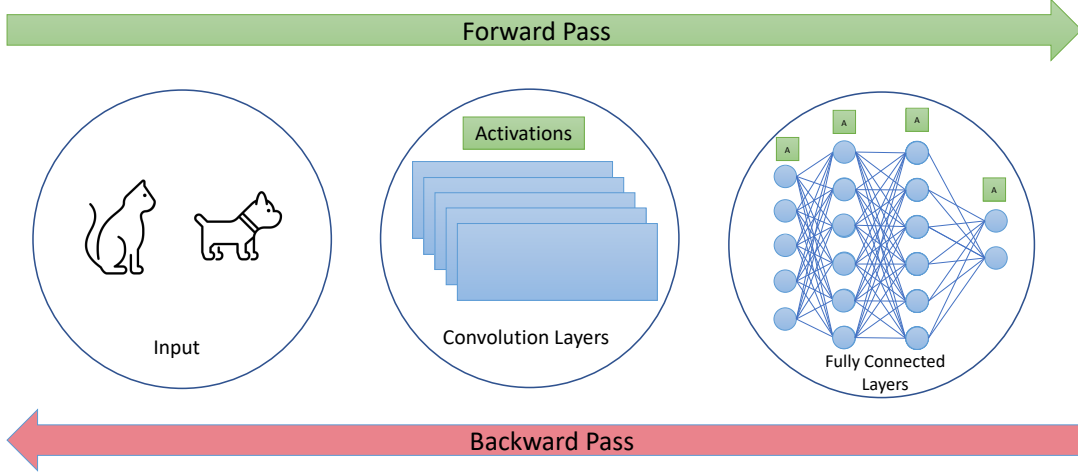


Figure 2.6: Memory Allocation in Forward and Backward Pass

that the total job completion time for a set of jobs is minimized. It is assumed that the jobs use the Parameter Server [9] architecture. They learn the convergence curve for a given job to predict the number of required iterations until the completion. Using this information, they can measure how long the job is going to last. Then, they create a heuristic algorithm that increases the number of workers or parameter servers for a job that gains more speed up compared to the other jobs. This work is complementary to our work. We can augment their techniques with our co-location algorithms to better utilize the GPUs.

Tiresias [6] presented a more realistic scheduling algorithm. Tiresias is able to use the historical data of the jobs to minimize the job completion time. They do not adaptively change the number of workers for a given job. This is a more realistic approach since increasing the number of workers affects the accuracy and may require retuning all the hyperparameters. They do not consider the packing of multiple jobs on the same GPU.

$$w(t+1) = w(t) - \eta * \nabla Q(w(t)) \quad (2.1)$$

Gandiva [12] introduced a fast context switch mechanism to avoid starvation of the jobs in deep learning clusters. They observed that the GPU memory usage of a job is not

constant during the training and is minimum between the iterations. While this is not true for PyTorch and Tensorflow frameworks, some frameworks like MXNet deallocate the memory when it is no longer needed. Figure 2.6 shows how memory allocation is performed during the training. In this figure, we are assuming the training of a convolutional neural network [8]. The weights associated with each layer is always present in the GPU. As the input traverses different layers of the neural network, it creates *activations*. Activations must be stored for each layer. The activations are required for calculating the gradients in eq. (2.1). These gradients will be calculated during the backward pass. In the backward pass, the activation values can be discarded and as a consequence the memory allocated for each of the layers may be freed. Gandiva leverages this pattern, and does not interrupt the job during the forward or backward passes to reduce the amount of data that needs to be copied during the checkpointing process. Gandiva also included a mechanism for "packing" the jobs to reduce the queuing time and improve JCT. They employed random packing to find the matching job pairs. As mentioned in [10], this strategy is not sufficient for finding beneficial co-locations.

Chic [5] introduced the idea of adaptively changing the number of workers using a reinforcement learning algorithm. They showed that using reinforcement learning will lead to better results compared to Optimus [11]. However, they still didn't consider the packing of multiple jobs into a single GPU.

Gavel [10] was the first scheduler to design an scheduling mechanism that can use many different scheduling objectives. Gavel has support for hierarchical and multi-domain scheduling policies. Their modeling of the scheduling problem is able to take into account packing of multiple jobs into a single GPU if the appropriate profiling information is available. They also include a throughput estimator that is able to estimate the co-location throughput of unseen jobs.

Although interference of co-located jobs in GPU clusters has not been very explored, there is a significant body of work on the interference effect of co-located jobs in CPU

clusters. Quasar [4] employs PQ-reconstruction with Stochastic Gradient Descent, to fill the unspecified elements of a matrix. In this matrix, the rows are jobs and the columns are different platforms. Quasar first profiles a couple of jobs extensively on a number of platforms. For unseen jobs, it profiles the job on a limited set of platforms and then uses the PQ-reconstruction to predict the performance on unseen platforms. Gavel [10] used this technique in the "Throughput Estimator" to predict the performance of co-location. They treated the co-located jobs as a new job and tried to fill in the matrix appropriately.

Chapter 3

Scheduler Design

We designed SC (Scheduling using Co-location) to utilize the excess amount of compute available in GPUs to reduce the queuing time and job completion time for the deep learning training jobs in GPU clusters. In Section 3.1, we show empirical results on how various jobs respond to co-location. In Section 3.2, we discuss the underlying reason on why some jobs are not compatible with each other. In Section 3.3 we provide the general design and architecture for our scheduling algorithm.

3.1 Co-location Effect

3.2 Why Some Jobs are not Compatible with each other?

3.3 Scheduling algorithm

Chapter 4

Evaluation

Chapter 5

Future Steps

Bibliography

- [1] Cuda warps and occupancy. https://on-demand.gputechconf.com/gtc-express/2011/presentations/cuda_webinars_WarpsAndOccupancy.pdf. (Accessed on 11/07/2020).
- [2] Multi-process service. https://docs.nvidia.com/deploy/pdf/CUDA_Multi_Process_Service_Overview.pdf. (Accessed on 11/16/2020).
- [3] Volta architecture. <https://images.nvidia.com/content/volta-architecture/pdf/volta-architecture-whitepaper.pdf>. (Accessed on 10/20/2020).
- [4] Christina Delimitrou and Christos Kozyrakis. Quasar: Resource-efficient and qos-aware cluster management. *SIGPLAN Not.*, 49(4):127–144, February 2014.
- [5] Yifan Gong, Baochun Li, Ben Liang, and Zheng Zhan. Chic: Experience-driven scheduling in machine learning clusters. In *Proceedings of the International Symposium on Quality of Service, IWQoS '19*, New York, NY, USA, 2019. Association for Computing Machinery.
- [6] Juncheng Gu, Mosharaf Chowdhury, Kang G. Shin, Yibo Zhu, Myeongjae Jeon, Junjie Qian, Hongqiang Liu, and Chuanxiong Guo. Tiresias: A GPU cluster manager for distributed deep learning. In *16th USENIX Symposium on Networked Systems Design and Implementation (NSDI 19)*, pages 485–500, Boston, MA, February 2019. USENIX Association.

- [7] Norm Jouppi. Google supercharges machine learning tasks with tpu custom chip. *Google Blog, May*, 18:1, 2016.
- [8] Yann LeCun, Bernhard Boser, John S Denker, Donnie Henderson, Richard E Howard, Wayne Hubbard, and Lawrence D Jackel. Backpropagation applied to handwritten zip code recognition. *Neural computation*, 1(4):541–551, 1989.
- [9] Mu Li, David G. Andersen, Jun Woo Park, Alexander J. Smola, Amr Ahmed, Vanja Josifovski, James Long, Eugene J. Shekita, and Bor-Yiing Su. Scaling distributed machine learning with the parameter server. In *11th USENIX Symposium on Operating Systems Design and Implementation (OSDI 14)*, pages 583–598, Broomfield, CO, October 2014. USENIX Association.
- [10] Deepak Narayanan, Keshav Santhanam, Fiodar Kazhamiaka, Amar Phanishayee, and Matei Zaharia. Heterogeneity-aware cluster scheduling policies for deep learning workloads. In *14th USENIX Symposium on Operating Systems Design and Implementation (OSDI 20)*, pages 481–498. USENIX Association, November 2020.
- [11] Yanghua Peng, Yixin Bao, Yangrui Chen, Chuan Wu, and Chuanxiong Guo. Optimus: An efficient dynamic resource scheduler for deep learning clusters. In *Proceedings of the Thirteenth EuroSys Conference, EuroSys ’18*, New York, NY, USA, 2018. Association for Computing Machinery.
- [12] Wencong Xiao, Romil Bhardwaj, Ramachandran Ramjee, Muthian Sivathanu, Nipun Kwatra, Zhenhua Han, Pratyush Patel, Xuan Peng, Hanyu Zhao, Quanlu Zhang, Fan Yang, and Lidong Zhou. Gandiva: Introspective cluster scheduling for deep learning. In *13th USENIX Symposium on Operating Systems Design and Implementation (OSDI 18)*, pages 595–610, Carlsbad, CA, October 2018. USENIX Association.