#### Tiresias: A GPU Cluster Manager for Distributed Deep Learning

NSDI'19

Guanbin Xu 2019-03-27, Reading Group



•

#### Tiresias

#### A GPU Cluster Manager

# for **Distributed Deep Learning**

NSDI'19

### Outline

#### **1. Deep Learning**

2. GPU Cluster Manager & Weakness of them

3. Tiresias:

1.Scheduler

2. Placement

4. Tiresias: Evaluation

#### Growing Use of Deep Learning at Google

Number of directories containing model description files



Ż А







Large-Scale Deep Learning for Intelligent Computer Systems, Google Research





f (w, x)









#### SIMD



### Outline

1. Deep Learning

#### 2. GPU Cluster Manager & Weakness of them

- 3. Tiresias:
  - 1.Scheduler
  - 2. Placement
- 4. Tiresias: Evaluation

# GPU Cluster Manager

#### Design Objectives

 Minimize Cluster-Wide Average Job Completion Time (JCT)

#### • Achieve

 High Resource (GPU) Utilization



GPU Cluster

#### Challenge 1: Unpredictable Training Time

Unknown execution time of DL training jobs

- Job execution time is useful when minimizing JCT
- Predict job execution time
  - Use the smooth loss curve of DL training jobs (Optimus [1])



[1]. Optimus: An Efficient Dynamic Resource Scheduler for Deep Learning Clusters, EuroSys'18

#### Challenge 2: Over-Aggressive Job Consolidation

- Network overhead in DDL training
  - Consolidated placement for good training performance
    - Fragmented free GPUs in the cluster
    - Longer queuing delay



# **Prior Solutions**

	I. Unpredictable Training Time ( <mark>Scheduling</mark> )	II. Over-Aggressive Job Consolidation (Job Placement)	
<b>Optimus</b> [1]	None	None	
YARN-CS	FIFO	None	
Gandiva <sub>[2]</sub>	Time-sharing	Trial-and-error	

[1]. Optimus: An Efficient Dynamic Resource Scheduler for Deep Learning Clusters, EuroSys'18[2]. Gandiva: Introspective Cluster Scheduling for Deep Learning, OSDI'18

#### Outline

- 1. Deep Learning
- 2. GPU Cluster Manager & Weakness of them
- 3. Tiresias:
  - **1.Scheduler**
  - 2. Placement
- 4. Tiresias: Evaluation

#### **Tiresias:**

A GPU cluster manager for Distributed Deep Learning Without Complete Knowledge

Age-Based Scheduler



 Minimize JCT without complete knowledge of jobs

 Model Profile-Based Placement



 Place jobs without additional information from users

Variations in both temporal and spatial aspects



Variations in both temporal and spatial aspects

128-

#### Scheduler should consider both temporal and spatial aspects of DL training jobs



- Spatial: number of GPUs
- Temporal: executed time
- distribution of job execution time(maybe)



# Tiresias: Scheduler

Short Remaining-Time First(SRTF)

- Least-Attained Service<sup>[1]</sup> (LAS) Short Job First(SJF)
  - Prioritize job that has the shortest executed time
- Gittins Index policy<sup>[2]</sup>
  - Need the distribution of job execution time
  - Prioritize job that has the highest probability to complete in the near future



- [1]. Feedback queueing models for time-shared systems. JACM, 1968
- [2]. Multi-armed bandit allocation indices. Wiley, Chichester, 1989

# Tiresias: Scheduler

- Least-Attained Service<sup>[1]</sup> (LAS)
  - Prioritize job that has the shortest executed time
- Gittins Index policy<sup>[2]</sup>
  - Need the distribution of job execution time
  - Prioritize job that has the highest probability to complete in the near future



- [1]. Feedback queueing models for time-shared systems. JACM, 1968
- [2]. Multi-armed bandit allocation indices. Wiley, Chichester, 1989

$$GI_{J} = \sup_{\Delta > 0} \frac{P(S - a_{J} \le \Delta | S > a_{J})}{E[\min\{S - a_{J}, \Delta\} | S > a_{J}]}$$

- P is the probability that J can complete with in  $\Delta$
- **E** is the expected service (cost) of **J** to be complete with in  $\Delta$
- $\boldsymbol{\Delta}$  is the next service quantum
- **P** and **E** are calculated from the distribution of job GPU time

	# of GPUs	Execution time		# of GPUs	Distribution
$J_{1}$	2	2	$J_{\rm f}$	2	2
J <sub>2</sub>	I	8	J <sub>2</sub>	Ι	(4, 8, 12)
J <sub>3</sub>	2	6	Ja	2	6



Higher probability to complete (Gittins Index), higher priority

	# of GPUs	Distribution	Attained Service	Gittins Index
$\mathbf{J}_{1}$	2	2	0	0.25
J <sub>2</sub>	I	(4, 8, 12)	0	0.25
$J_3$	2	6	0	0.25



$$GI_{J} = \sup_{\Delta > 0} \frac{P(S - a_{J} \le \Delta | S > a_{J})}{E[\min\{S - a_{J}, \Delta\} | S > a_{J}]}$$

Δ=4

 $G_{j_1} = \frac{P_{s=4}}{\min(4-0,\delta)*1/3 + \min(8-0,\delta)*1/3 + \min(12-0,\delta)*1/3} = \frac{1/3}{1/3*4} = 0.25$ 

	# of GPUs	Distribution	Attained Service	Gittins Index
J	2	2	4	0.2
J <sub>2</sub>	I	(4, 8, 12)	0	0.25
J3	2	6	0	0.25





	# of GPUs	Distribution	Attained Service	Gittins Index
Ji	2	2	4	0.2
J <sub>2</sub>	I	(4, 8, 12)	4	0.2
J3	2	6	0	0.25





	# of GPUs	Distribution	Attained Service	Gittins Index
$J_{1}$	2	2	4	0.2
J <sub>2</sub>	I	(4, 8, 12)	4	0.2
J3	2	6	4	0.2





	# of GPUs	Distribution	Attained Service	Gittins Index
$\mathbf{J}_{\mathbf{i}}$	2	2	4	0.2
J2	I	(4, 8, 12)	8	0.125
J <sub>3</sub>	2	6	4	0.2





	# of GPUs	Distribution	Attained Service	Gittins Index
J	2	2	4	0.2
J <sub>2</sub>	I	(4, 8, 12)	8	0.125
J3	2	6	12	N/A





$$GI_{J} = \sup_{\Delta > 0} \frac{P(S - a_{J} \le \Delta | S > a_{J})}{E[\min\{S - a_{J}, \Delta\} | S > a_{J}]}$$

#### Two-Dimensional Age-Based Scheduler (2DAS)

- Age calculated by two-dimensional attained service
  - i.e., a job's total executed GPU time (# of GPUs × executed time)
- No prior information
  - 2D-LAS
- With partial information: distribution of job GPU time
  - 2D-Gittins Index

#### Two-Dimensional Age-Based Scheduler (2DAS)

- Age calculated by two-dimensional attained service
  - i.e., a job's total executed GPU time (# of GPUs × executed time)

#### Fewer job switches: Priority discretization: Discretized-2DAS

- 2D-LAS
- With partial information: distribution of job GPU time
  - 2D-Gittins Index

# **Prior Solutions**

	I. Unpredictable Training Time (Scheduling)		I. Unpredictable Training Time ( <mark>Scheduling</mark> )		II. Over-Aggressive Job Consolidation (Job Placement)
<b>Optimus</b> <sub>[1]</sub>	None		None		None
YARN-CS	FIFO		None		
Gandiva <sub>[2]</sub>	Time-sharing		Trial-and-error		
Tiresias	LAS	Gittins Index	?		

[1]. Optimus: An Efficient Dynamic Resource Scheduler for Deep Learning Clusters, EuroSys'18[2]. Gandiva: Introspective Cluster Scheduling for Deep Learning, OSDI'18

### Outline

- 1. Deep Learning
- 2. GPU Cluster Manager & Weakness of them
- 3. Tiresias:
  - 1.Scheduler

#### 2. Placement

4. Tiresias: Evaluation



4 concurrent 8-worker jobs with different placement schemes.

600 500 400 Size (MB) 300 200 100 0 VGGIP VGG19 Restlect52 Resfletto Pestveriol VGGI Alexiver Inception3 Googletter InceptionA

#### Tensor size in DL models

• Large tensors cause network imbalance and contention

- Tensor size in DL models
  - Large tensors cause network imbalance and contention

#### Consolidated placement: when the model is highly skewed in its tensor size



#### • Tensor size in DL models



### Outline

- 1. Deep Learning
- 2. GPU Cluster Manager & Weakness of them
- 3. Tiresias:
  - 1.Scheduler
  - 2. Placement
- 4. Tiresias: Evaluation

#### Tiresias





https://github.com/SymbioticLab/Tiresias

# **Evaluation - Setup**

- Testbed Experiment(Michigan ConFlux cluster)
  - 100 Gbps EDR Mellanox IB + RDMA protocol
  - 15\*(4-NVIDIA Tesla P100s with NVlink + 256GB DDR4)
  - GPFS(1.2GB/s)
- Large-scale Trace-driven(from Microsoft) Simulation
  - Information: job arrival, completion, demotion, propotion, preemption, #GPU, training time
- #queues: k=2
- threshold= $\Delta = 1$

### **Evaluation - Workload**

- 480 jobs
  - 240 \* 1-GPU jobs
  - 40 \* 2-GPU jobs
  - 80 \* 4-GPU jobs
  - 90 \* 8-GPU jobs
  - 25 \* 16-GPU jobs
  - 5 \* 32-GPU jobs



#### JCT Improvements inTestbed Experiment



• Avg. JCT improvement (w.r.t.YARN-CS): 5.5×

• Comparable performance to SRTF

#### JCT in Testbed Experiment



Bins	l (Small-Short)	2(Small-Long)	3(Large-Short)	4(Large-Long)
% of Jobs	63.5%	12.5%	16.5%	7.5%

#### Queuing Delay inTestbed Experiment

	Average	Median	95th
YARN-CS	8146s	7464s	l 5327s
SRTF	593s	32s	3133s
Tiresias-G	1005s	39s	7933s
Tiresias-L	963s	l 3s	7755s

#### Time overhead of Job switch

Total size (MB)

Largest tensor (MB)

Model

VGG19

VGG16

VGGTT



#### **GPU Utilization in Testbed**



The makespan is improved by 1.21× (w.r.t.YARN-CS)

#### Training Performance in Testbed Experiment



• Training time when Tiresias-L running with and without placement

#### JCT Improvements in Trace-Driven <u>Simulation</u>



#### Sensitivity Analysis of 2D-LAS



Fig14a



Fig14c

4

2







#### Tiresias: A GPU Cluster Manager for Distributed Deep Learning

NSDI'19

#### Q&A



•

#### Backup

• N/A