

Optimus: An Efficient Dynamic Resource Scheduler for Deep Learning Clusters

Yanghua Peng

The University of Hong Kong
yhpeng@cs.hku.hk

Yixin Bao

The University of Hong Kong
yxbao@cs.hku.hk

Yangrui Chen

The University of Hong Kong
yrchen@cs.hku.hk

Chuan Wu

The University of Hong Kong
cwu@cs.hku.hk

Chuanxiong Guo

Bytedance Inc.
guochuanxiong@bytedance.com

ABSTRACT

Deep learning workloads are common in today's production clusters due to the proliferation of deep learning driven AI services (e.g., speech recognition, machine translation). A deep learning training job is resource-intensive and time-consuming. Efficient resource scheduling is the key to the maximal performance of a deep learning cluster. Existing cluster schedulers are largely not tailored to deep learning jobs, and typically specifying a fixed amount of resources for each job, prohibiting high resource efficiency and job performance. This paper proposes *Optimus*, a customized job scheduler for deep learning clusters, which minimizes job training time based on online resource-performance models. *Optimus* uses online fitting to predict model convergence during training, and sets up performance models to accurately estimate training speed as a function of allocated resources in each job. Based on the models, a simple yet effective method is designed and used for dynamically allocating resources and placing deep learning tasks to minimize job completion time. We implement *Optimus* on top of Kubernetes, a cluster manager for container orchestration, and experiment on a deep learning cluster with 7 CPU servers and 6 GPU servers, running 9 training jobs using the MXNet framework. Results show that *Optimus* outperforms representative cluster schedulers by about 139% and 63% in terms of job completion time and makespan, respectively.

CCS CONCEPTS

- Computing methodologies → Machine learning;
- Computer systems organization → Cloud computing;

KEYWORDS

Resource management; deep learning

ACM Reference Format:

Yanghua Peng, Yixin Bao, Yangrui Chen, Chuan Wu, and Chuanxiong Guo. 2018. *Optimus: An Efficient Dynamic Resource Scheduler for Deep Learning Clusters*. In *EuroSys '18: Thirteenth EuroSys Conference 2018, April 23–26, 2018, Porto, Portugal*. ACM, New York, NY, USA, 14 pages. <https://doi.org/10.1145/3190508.3190517>

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for components of this work owned by others than ACM must be honored. Abstracting with credit is permitted. To copy otherwise, or republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Request permissions from permissions@acm.org.

EuroSys '18, April 23–26, 2018, Porto, Portugal

© 2018 Association for Computing Machinery.

ACM ISBN 978-1-4503-5584-1/18/04...\$15.00

<https://doi.org/10.1145/3190508.3190517>

Clusters. In *EuroSys '18: Thirteenth EuroSys Conference 2018, April 23–26, 2018, Porto, Portugal*. ACM, New York, NY, USA, 14 pages. <https://doi.org/10.1145/3190508.3190517>

1 INTRODUCTION

The recent five years have witnessed substantial progress and successful applications of deep learning in various domains of AI, such as computer vision [39], natural language processing [65] and speech recognition [68]. The rising amount of data and increasing scale of training models (e.g., deep neural networks) significantly improve the learning accuracy, as well as remarkably extend the training time. Distributed machine (deep) learning frameworks have been designed and deployed to expedite model convergence using parallel training with multiple machines, e.g., TensorFlow [23], MXNet [59]. Most leading IT companies have been operating distributed machine learning (ML)/deep learning (DL) clusters, with hundreds or thousands of (GPU) servers, to train various ML models over large datasets for their AI-driven services.

Even with parallel training, a deep learning job is resource intensive and time consuming. For example, to train the DeepSpeech2 model [25] on the LibriSpeech dataset (1000 hours of speech) [9], it takes 3–5 days to achieve the state-of-the-art accuracy when training on 16 GPUs [25]. In a shared deep learning cluster with various training jobs submitted over time, efficient resource scheduling is the key to maximize utilization of expensive resources (e.g., GPUs and RDMA networks) for expedited training completion. However, achieving high training performance and resource efficiency in deep learning clusters is challenging with existing cluster schedulers.

First, schedulers used in existing ML/DL clusters (e.g., Google uses **Borg** [63], Microsoft, Tencent and Baidu use YARN-like schedulers [61]) allocate a fixed amount of resources to each job upon its submission, according to resource requirements specified by the job owner. Jobs already running in the cluster cannot benefit from extra resources when they become available (e.g., during night time when there are lower workloads), unless the cluster operator manually reconfigures their resource composition or a job owner resubmits the job as new. This may well lead to low resource utilization efficiency.

Second, existing schedulers are designed for different workloads but deep learning. For example, Mesos, Yarn and Borg are for general-purpose cluster resource management, *Corral* [42] is designed for periodic data-parallel jobs, and *Tetrisched* [60] handles reservation-based workloads. There is room for improving resource utilization in deep learning clusters with a tailor-made resource scheduler that leverages structures of deep learning frameworks

(e.g., the parameter server architecture) and characteristics of deep learning jobs (e.g., iterativeness, convergence properties) for maximal training efficiency.

This paper proposes *Optimus*, a customized cluster scheduler for deep learning jobs in production clusters, which minimizes job training time and improves resource efficiency as a result. We focus on data-parallel DL training jobs using the parameter server framework (§2). *Optimus* builds resource-performance models for each job on the go, and dynamically schedules resources to jobs based on job progress and the cluster load to minimize average job completion time and makespan. Specifically, we make the following contributions in developing *Optimus*.

- » We build accurate performance models for deep learning jobs (§3). Through execution of a training job, we track the training progress on the go and use online fitting to predict the number of steps/epochs required to achieve model convergence (§3.1). We further build a resource-performance model by exploiting communication patterns in the parameter server architecture and iterativeness of the training process (§3.2). Different from existing detailed modeling of a distributed deep learning job (such as in [69]), our resource-performance model requires no knowledge about internals of the ML model and hardware configuration of the cluster. The basis is an online learning idea: we run a job for a few steps with different resource configurations, learn the training speed as a function of resource configurations using data collected from these steps, and then keep tuning our model on the go.

- » Based on the performance models, we design a simple yet effective method for dynamically allocating resources to minimize average job completion time (§4.1). We also propose a task placement scheme for deploying parallel tasks in a job onto the servers, given the job's resource allocation (§4.2). The scheme further optimizes training speed by mitigating communication overhead during training.

- » We discover a load imbalance issue on parameter servers with the existing parameter server framework (as in MXNet [59]), which significantly lowers the training efficiency. We resolve the issue by reducing communication cost and assigning model slices to parameter servers evenly (§5.3). We integrate our scheduler *Optimus* with Kubernetes [14], an open-source cluster manager for production-grade container orchestration. We build a deep learning cluster consisting of 7 CPU servers and 6 GPU servers, and run 9 representative DL jobs from different application domains (see Table 1). Evaluation results show that *Optimus* achieves high job performance and resource efficiency, and outperforms widely adopted cluster schedulers by 139% and 63% in job completion time and makespan, respectively (§6).

2 BACKGROUND AND MOTIVATION

2.1 DL Model Training

A deep learning job trains a DL model, such as a deep neural network (DNN), using a large number of training examples, to minimize a loss function (typically) [48].

Iterativeness. The model training is usually carried out in an iterative fashion, due to the complexity of DNNs (*i.e.*, no closed-form solution) and the large size of training dataset (*e.g.*, 14 million images in the full Imagenet dataset [12]). The dataset is commonly divided

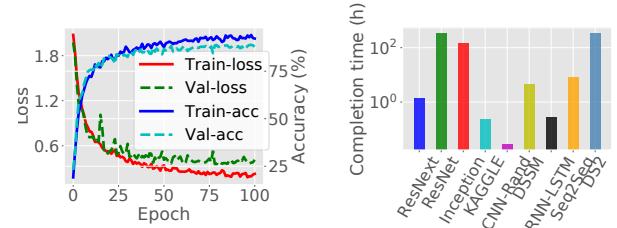


Figure 1: Training curves of ResNext-110 on the CIFAR10 dataset

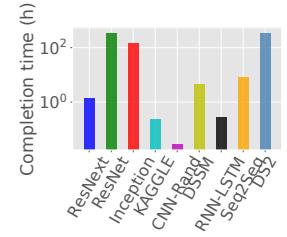


Figure 2: Training time of deep learning models in Table 1

into equal-sized data chunks, and each data chunk is further divided into equal-sized *mini-batches*. In each training *step*, we process one mini-batch by computing what changes to be made to the parameters in the DL model to approach their optimal values (typically expressed as *gradients*, *i.e.*, directions of changes), using examples in the mini-batch, and then update parameters using a formula like `new_parameter = old_parameter - learning_rate × gradient`. A training performance metric is also computed for each mini-batch, *e.g.*, training loss (the sum of the errors made for each example in the mini-batch) or accuracy (the percentage of correct predictions compared to the labels), validation loss or accuracy (computed on validation dataset for model evaluation). After all mini-batches in training dataset have been processed once, one training *epoch* is done.

Convergence. The dataset is usually trained for multiple epochs (tens to hundreds) until the model converges, *i.e.*, the decrease or increase in the performance metric's value between consecutive epochs becomes very small. An illustration of the training curves, the variation of training/validation loss and accuracy vs. the number of training epochs, is given in Fig. 1, with the example of training ResNext-110 [66] on the CIFAR10 dataset [2]. DNN models are usually non-convex and we can not always expect convergence [29]. However, different from experimental models, production models are mature and can typically converge to the global/local optimum very well since all hyper-parameters (*e.g.*, learning rate – how quickly a DNN adjusts itself, mini-batch size) have been well-tuned during the experimental phase. In this work, we focus on such production models, and leverage their convergence property to estimate a training job's progress towards convergence.

Especially, we use the convergence of training loss to decide the completion of a DL job. The DL model converges if the decrease of training loss between two consecutive epochs has consistently fallen below a threshold that the job owner specified, for several epochs. Training loss based training convergence is common in practice [48, 71] and the convergence of training loss often implies the convergence of other metrics (*e.g.*, accuracy) for production models (*i.e.*, no overfitting) [5]. Training/validation accuracy is difficult to be defined in some scenarios where there is no “right answer”, *e.g.*, language modeling [6]. Validation loss is usually used to prevent model overfitting, and evaluation on validation dataset is performed only when necessary (*e.g.*, at the end of each epoch), while we can obtain training loss after each step for more accurate curve fitting (§3.1).

2.2 The Parameter Server Architecture

Most distributed ML/DL frameworks (*e.g.*, MXNet [59], TensorFlow [23], PaddlePaddle [17], Angel [43], Petuum [67]) employ the parameter server (PS) architecture [48] (Fig. 3). In this architecture, the model (*i.e.*, a DNN) is partitioned among multiple parameter servers and the training data are split among workers. Each worker computes parameter updates (*i.e.*, gradients) locally using its data partition and pushes them to parameter servers maintaining the respective model parameters. After receiving gradients, parameter servers update the model parameters using some optimization algorithm, *e.g.*, Stochastic Gradient Descent (SGD) [20]. Updated parameters are sent back to the workers, which then start the next training step, using the updated parameters.

There are two training modes: *asynchronous training*, where the training progress (*i.e.*, number of steps) at different workers in a job is not synchronized and a parameter server updates its model partition each time upon receiving gradients from a worker; *synchronous training*, where training progress at all workers is synchronized in each step and a parameter server updates parameters after it has collected gradients from all workers.

2.3 Existing Cluster Schedulers

Static resource allocation. Parameter servers and workers typically run in containers or virtual machines in a DL cluster, and a cluster scheduler manages the resource allocation to training jobs, *e.g.*, Mesos [40] in a TensorFlow cluster [23], Yarn [61] for clusters running MXNet [59] or Angel [43]. With these schedulers, the owner of a training job specifies resource requirements, *e.g.*, the numbers of parameter servers and workers, which remain unchanged throughout the training process.

The numbers of workers and parameter servers used to run a training job influence the training speed (*i.e.*, the average number of training steps completed per second), and hence the training completion time significantly. Fig. 4 shows the training speed varies with different numbers of workers and parameter servers deployed, when we synchronously train a ResNet-50 model [39], one of the state-of-the-art DNNs for image classification (details in Table 1), on the ImageNet dataset [12]. Each container is configured with 5 CPU cores and 10GB memory, and can run 1 worker or 1 parameter server. In Fig. 4(a), we fix the total number of containers to be 20, *i.e.*, if the number of workers is x , then the number of parameter servers used is $20 - x$. We can see the maximal training speed is achieved when there are 8 workers and 12 parameter servers. In

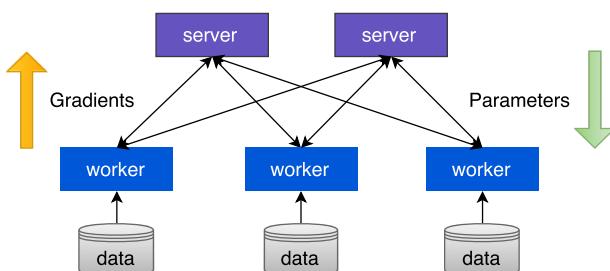


Figure 3: Parameter server architecture

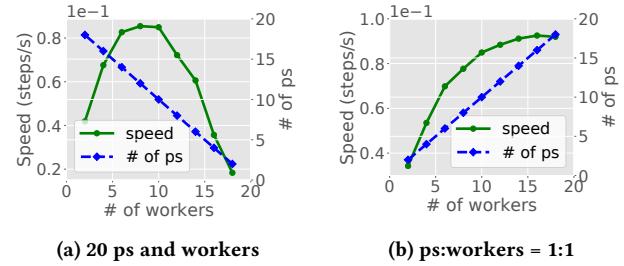


Figure 4: Varying training speeds with different resource configurations

Fig. 4(b), we fix the ratio of the number of parameter servers to the number of workers to be 1:1. We see that increasing resources do not lead to linear training speed improvement, and can even slow down model training.

In a production cluster, job training speed is further influenced by many runtime factors, such as available bandwidth at the time. Configuring a fixed number of workers/parameter servers upon job submission is hence unfavorable. In *Optimus*, we maximally exploit varying runtime resource availability by adjusting numbers and placement of workers and parameter servers, aiming to pursue the best resource efficiency and training speed at each time. Note that the resource composition of each worker or parameter server is still specified by the job owner.

Job size unawareness. Existing schedulers largely adopt FIFO (as in Spark [70]), Dominant Resource Fairness (DRF) [34] (as in Mesos [40] and Yarn [61]) or their variants as default scheduling strategies, which are ignorant of job sizes (represented by input data size, model complexity, or time taken to complete the job). It has been shown that job performance can be improved by considering job sizes when making scheduling decisions [30, 31]. For example, a long job may block a series of short jobs with an FIFO scheduler that is oblivious to the job sizes, causing starvation or long completion time for short jobs.

Training completion time varies significantly among DL jobs. Fig. 2 shows the training time of several representative DL models on respective datasets, as given in Table 1, on a TITAN X Pascal GPU. The training time varies from minutes (for simple models on small datasets, *e.g.*, CNN-rand [46]) to weeks (for complex models on large datasets, *e.g.*, ResNet-50 [39]). *Optimus* takes into account projected job completion time for different DL jobs when dynamically adjusting their resource allocation, to minimize average job completion time.

3 PERFORMANCE MODELING OF DL JOBS

To make good resource scheduling decisions, we would like to know the relation between resource configuration and the time a training job takes to achieve model convergence. We derive this relation by estimating online how many more training epochs a job needs to run for convergence (§3.1), and how much time a job needs to complete one training epoch given a certain amount of resources (§3.2).

Table 1: Deep learning jobs used for tests and experiments

| Model | # of parameters (Million) | Network type | Application domain | Dataset | Dataset size (# of examples) |
|---------------------------|---------------------------|--------------|-------------------------|--------------------------|------------------------------|
| ResNext-110 [66] | 1.7 | CNN | image classification | CIFAR10 [2] | 60,000 |
| ResNet-50 [39] | 25 | CNN | image classification | ILSVRC2012-ImageNet [12] | 1,313,788 |
| Inception-BN [58] | 11.3 | CNN | image classification | Caltech [1] | 30,607 |
| KAGGLE [13] | 1.4 | CNN | image classification | Kaggle-NDSB1 [4] | 37,920 |
| CNN-rand [46] | 6 | CNN | sentence classification | MR [26] | 10,662 |
| DSSM [54] | 1.5 | RNN | word representation | text8 [49] | 214,288 |
| RNN-LSTM-Dropout [22] | 4.7 | RNN | language modeling | PTB [18] | 1,002,000 |
| Sequence-to-Sequence [33] | 9.1 | RNN | machine translation | WMT17 [21] | 1,000,000 |
| DeepSpeech2 [25] | 38 | RNN | speech recognition | LibriSpeech [9] | 45,000 |

3.1 Learning the Convergence Curve

We draw the training loss curve with the training progress of each DL job, and do online model fitting, in order to predict how far the model is from convergence.

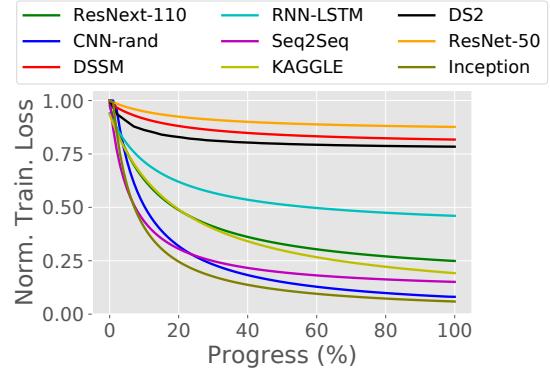
Data preprocessing. For better model fitting, we carry out outlier removal as follows: if a loss data point does not fall within a certain range of its neighbours (*e.g.*, between the minimum loss in subsequent 5 epochs and the maximum loss in previous 5 epochs), we consider the data point as **an outlier**, and use the average value of its neighbours to replace this point when doing model fitting. We also **normalize the loss values**, by dividing each raw value by the maximum loss value collected so far (typically the first loss value). In this way, loss values in different DL jobs are all between 0 and 1. Fig. 5 shows example loss curves collected by running example DL jobs in Table 1 (which are DL examples from official MXNet tutorials [15, 16]), using the MXNet framework on a server with 1 E5-1650 v4 CPU and 2 NVIDIA TITAN X GPUs. The learning rate of each job is set to be fixed. The training progress is the ratio of the number of epochs a model has been trained over the total number of epochs needed for convergence.

Online fitting. We observe that most DL jobs use SGD to update parameters and approximate the optimal parameter values. Since SGD converges at a rate of $O(1/k)$ in terms of the number of steps k , we use the following model to fit the training loss curve:

$$l = \frac{1}{\beta_0 \cdot k + \beta_1} + \beta_2 \quad (1)$$

where l denotes the training loss, and β_0, β_1 and β_2 are nonnegative coefficients. Our online model fitting is carried out as follows: after each training step, we collect a training loss data point (k, l) ; we then preprocess the data as described above and use a non-negative least squares (NNLS) solver [7] to find the best coefficients that fit the loss points collected so far. In some cases hundreds of thousands of steps are needed to achieve model convergence; in such a case we can sample loss data every few steps, or average the values of several data points (*e.g.*, all losses in an epoch) as a single data point, to reduce the number of data points fed into the solver. Since we can collect more and more loss data as the job runs, the fitted model improves continuously. An example of model fitting when training the Seq2Seq model [33] in Table 1 is given in Fig. 7.

At each step, using the fitted loss model and the predefined convergence threshold δ , we can easily calculate the total number of

**Figure 5: Training loss curves for different DL jobs**

steps/epochs a job needs to achieve convergence, as well as the number of steps/epochs left from now until convergence. With more and more data points collected for model fitting, the estimation of the total number of steps/epochs a job needs improves gradually, as illustrated in Fig. 6. Here the prediction error is the difference between the estimated total number of epochs for the model training to converge and the actual total number of epochs needed, divided by the actual number.

3.2 Resource-Speed Modeling

We next build a resource-to-speed model based on computation and communication patterns in the parameter server architecture.

System models. In a typical DL job, the time taken to complete one training step on a worker includes the time for doing forward propagation (*i.e.*, loss computation) and backward propagation (*i.e.*, gradients computation) at the worker, the worker pushing gradients to parameter servers, parameter servers updating parameters, and the worker pulling updated parameters from parameter servers, plus extra communication overhead. Suppose there are p parameter servers and w workers in the job. The bandwidth capacity of each parameter server is B , and the model size (*i.e.*, total bytes of parameters) is S . The forward propagation time when a worker trains a minibatch is $m \cdot T_{forward}$ (the size of a mini-batch times the average processing time of one example). The backward propagation time T_{back} is not related to m and is typically fixed. The size of gradients

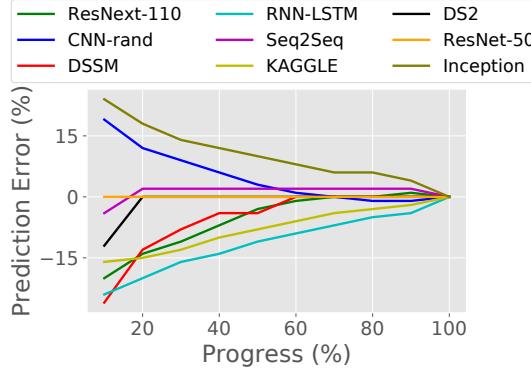


Figure 6: Prediction errors in different DL jobs

is the same as the model size S . If the parameters are evenly distributed on parameter servers (load balanced to be achieved in §5.3), the size of gradients sent between a worker and a parameter server is $\frac{S}{p}$. In practice, the bandwidth bottleneck between a worker and a parameter server usually lies at the parameter server side. Let w' denote the number of workers that send gradients to a parameter server at the same time. Then the bandwidth between the worker and parameter server ρ is $\frac{B}{w'\rho}$. Pushing gradients and pulling updated parameters are symmetric processes, so the data transfer time of each worker is $2\frac{S/p}{B/w'\rho}$. The parameter update time on a parameter server is $\frac{T_{update} \cdot w'}{p}$ on average, where T_{update} is the time to update parameters with size S . In addition, the communication overhead (e.g., handling TCP connections and control messages between parameter servers and workers) increases linearly with the number of parameter servers and the number of workers. It is represented by $\delta \cdot w + \delta' \cdot p$ where δ and δ' are coefficients.

Therefore, the duration of one training step on a worker can be modeled as

$$T = \max_p [m \cdot T_{forward} + T_{back} + 2\frac{S/p}{B/w'\rho} + \frac{T_{update} \cdot w'}{p} + \delta \cdot w + \delta' \cdot p] \quad (2)$$

According to Eqn. 2, the workers should have similar processing speeds and parameter servers should be load-balanced, in order to achieve minimal time per training step. We will discuss how to handle slow workers in §5.2 and achieve load balancing among parameter servers in §5.3.

We next derive the training speed in a job based on Eqn. 2, which is the number of training steps completed per unit time. We divide our models in two cases.

Asynchronous training, where the workers process mini-batches at their own pace. The overall number of training steps completed by all workers per unit time is $w \cdot T^{-1}$. Suppose w' is linear with w , since more workers may concurrently communicate with one parameter server if the total number of workers is larger. Then the training speed achieved with p parameter servers and w workers can be modeled as

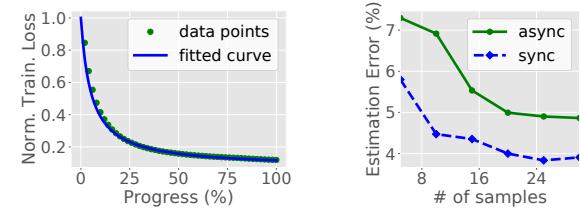


Figure 7: Online model fitting for training

Seq2Seq: $\beta_0 = 0.21$, $\beta_1 = 1.07$, $\beta_2 = 0.07$

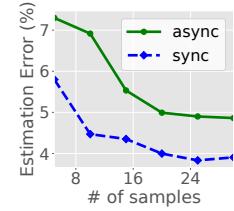


Figure 8: Estimation errors of training speeds

$$f(p, w) = w \cdot (\theta_0 + \theta_1 \cdot \frac{w}{p} + \theta_2 \cdot w + \theta_3 \cdot p)^{-1} \quad (3)$$

where θ are positive coefficients, corresponding to respective terms in Eqn. 2. For example, θ_0 corresponds to the term $m \cdot T_{forward} + T_{back}$ in Eqn. 2. Instead of measuring each term (e.g., $T_{forward}$), we seek to learn the coefficients by fitting the model with runtime data collected for each job.

Synchronous training, where all workers progress from one step to the next at a synchronized pace. The training speed is T^{-1} . w' equals w since all workers are synchronized. For synchronous training, the batch size, i.e., the overall size of all mini-batches trained by all workers in each step, needs to remain the same, no matter how we adjust the number of concurrent workers over time. This guarantees that the same training result (model) can be achieved while varying the number of workers [36]. Let M denote the batch size which is typically specified in the training job when the owner submits it. Then the mini-batch size on each worker is $m = \frac{M}{w}$. The training speed function can be modeled as

$$f(p, w) = (\theta_0 \cdot \frac{M}{w} + \theta_1 + \theta_2 \cdot \frac{w}{p} + \theta_3 \cdot w + \theta_4 \cdot p)^{-1} \quad (4)$$

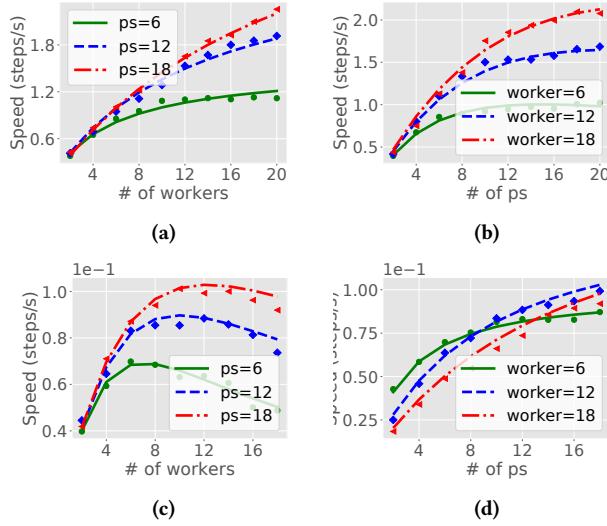
where θ are positive coefficients, to be learned for each job.

Model fitting. To learn the values of θ 's and build the training speed functions in Eqn. 3 and Eqn. 4, we need to collect data points $(p, w, f(p, w))$. Before we run each training job, we train its model on a small sample set of training data for several steps, with possible combinations of p and w . Each run takes about tens of seconds. In each run, we derive the average training speed under (p, w) . Due to the iterative nature of DL model training, training for several steps is enough to give us a good idea of the training speed $f(p, w)$. Then we use NLLS to find θ 's that best fit the collected data points $(p, w, f(p, w))$. This initial training speed function constructed is used for resource scheduling decisions when we start running the actual job. Over the training process, we keep collecting data points $(p, w, f(p, w))$ and use them to calibrate coefficients in our training speed models.

Fig. 9 shows the collected data points and the fitted training speed function curves, when we run the ResNet-50 job in a cluster of 40 containers using synchronous training and asynchronous training, respectively. We make three important observations: (a)

Table 2: Coefficients in speed functions

| | θ_1 | θ_2 | θ_3 | θ_4 | θ_5 | Residual sum of squares for fitting |
|-------|------------|------------|------------|------------|------------|-------------------------------------|
| Async | 2.83 | 3.92 | 0.00 | 0.11 | - | 0.10 |
| Sync | 1.02 | 2.78 | 4.92 | 0.00 | 0.02 | 0.00 |

**Figure 9: Data points and fitted curves of speed functions for asynchronous training (a)(b) and synchronous training (c)(d)**

our speed function can closely describe the relationship between the training speed and resource configurations; (b) due to communication overhead, there is a trend of diminishing return where adding more parameter servers or workers does not improve the training speed much; (c) for synchronous training, more workers may lead to lower training speed. This is because more workers lead to smaller mini-batch size $\frac{M}{w}$ (*i.e.*, a lower workload on each worker), which may cause CPU/GPU under-utilization. Meanwhile, a larger number of workers lead to higher synchronization cost and communication overhead.

Table 2 lists derived coefficients in the speed functions for asynchronous and synchronous training, respectively. We find that forward propagation, backward propagation and data transfer make up most of the training time in one step, since coefficients of these quantities are relatively large.

The reason why we produce the initial training speed function under possible combinations of p and w before running the actual job, is the following: (p, w) pairs used in actual resource configuration when running each job are limited; training speed functions learned using the limited data points may be biased, diverting resource allocation decisions away from the optimum. One question is how many possible (p, w) pairs we should try out to initialize the speed function, to achieve high model fitting accuracy. For the above ResNet-50 example, there are 780 possible (p, w) pairs. Fig. 8 shows the estimation errors of training speeds when we randomly select a number of samples, *i.e.*, (p, w) pairs, to produce the training

speed function. The estimation error is the ratio of the gap (between the measured speed and the estimated speed) over the measured speed. We observe that: (a) we can get a less than 10% error even when we only use 10 (p, w) pairs to learn the speed function; (b) using more (p, w) pairs leads to smaller error, but with a diminishing return.

4 DYNAMIC SCHEDULING

In our DL cluster, jobs arrive in an online manner. *Optimus* periodically allocates resources to the active jobs (new jobs submitted in the previous scheduling interval and unfinished jobs submitted earlier), by adjusting the numbers and placement of parameter servers/workers in each job in the shared DL cluster. Its scheduling algorithm consists of two parts: resource allocation and task placement.

4.1 Resource Allocation

In each scheduling interval, let Q_j denote the remaining number of steps/epochs that a job j needs to run to achieve model convergence (§3.1), and $f(p_j, w_j)$ be the current training speed function for job j (§3.2). We can estimate the remaining running time t_j of job j as $\frac{Q_j}{f(p_j, w_j)}$. Let $O_j^r(N_j^r)$ denote the amount of type- r resource each worker (parameter server) in job j occupies. C_r is the overall capacity of type- r resource in the DL cluster and R is the number of resource types. J is the set of current active jobs. Our scheduler aims to minimize the average completion time of these jobs. We can solve the following optimization problem to decide the numbers of workers/parameter servers for each job $j \in J$, where (7) is the capacity constraint:

$$\text{minimize} \quad \sum_{j \in J} t_j \quad (5)$$

$$\text{subject to: } t_j = \frac{Q_j}{f(p_j, w_j)} \quad \forall j \in J \quad (6)$$

$$\sum_{j \in J} (w_j \cdot O_j^r + p_j \cdot N_j^r) \leq C_r \quad \forall r \in R \quad (7)$$

$$p_j \in Z^+, w_j \in Z^+ \quad \forall j \in J \quad (8)$$

The problem is a non-linear (and even non-convex) integer programming problem since Eqn. 6 is not a linear/convex constraint. It can not be solved using LP/convex solvers and is NP-hard in general, so we design an efficient heuristic to solve it. We define the *marginal gain* in job completion time reduction as follows:

$$\max\left\{\left(\frac{Q_j}{f(p_j, w_j)} - \frac{Q_j}{f(p_j + 1, w_j)}\right)/N_j^D, \left(\frac{Q_j}{f(p_j, w_j)} - \frac{Q_j}{f(p_j, w_j + 1)}\right)/O_j^{D'}\right\} \quad (9)$$

Here D (D') is the dominant resource of workers (parameter servers) in job j . A dominant resource is the type of resource that has the maximal share in the overall capacity of the cluster, among all resources used by a worker (parameter server) [34]. $\frac{Q_j}{f(p_j, w_j)} - \frac{Q_j}{f(p_j, w_j + 1)} \left(\frac{Q_j}{f(p_j, w_j)} - \frac{Q_j}{f(p_j + 1, w_j)}\right)$ is the reduction in job completion time when one worker (parameter server) is added to job j ; dividing it by the amount of dominate resource that a worker (parameter

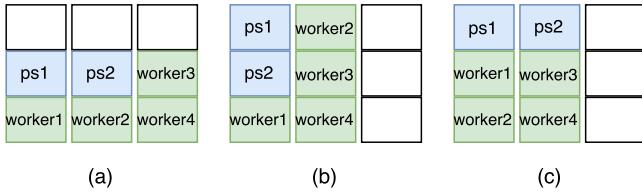


Figure 10: An example of worker/parameter server placement: (c) is the best

server) occupies, we obtain the marginal gain per unit dominant resource consumption.

Our resource allocation algorithm in each scheduling interval works as follows. We first allocate one worker and one parameter server to each active job to avoid starvation, and then sort all jobs in order of their marginal gains computed using (9). Then we iteratively select the job with the largest marginal gain and add one worker or parameter server to the job, according to which of the two terms in (9) is larger (*i.e.*, whether adding a worker or adding a parameter server brings larger marginal gain). Marginal gains of the jobs are updated when their resource allocation changes. The procedure repeats until some resource in the cluster is used up, or marginal gains of all jobs become non-positive.

The algorithm makes use of predictions based on online fitted models in §3. To mitigate its performance degradation due to prediction errors, we can downgrade the priority of a job a bit when it is at the beginning state (*i.e.*, larger prediction errors) by multiplying its marginal gain (*i.e.*, the computed value in (9)) by a factor (e.g., 0.95). Smaller marginal gain of a job means less resources allocated to it, thus mitigating the influence of large prediction errors at the start of training.

4.2 Task Placement

In our model of the training step duration in Eqn. 2, processing time on workers and parameter servers are fixed. We can reduce the time, a.k.a. improve the training speed, by reducing the time spent on parameters/gradients exchange among workers and parameter servers, which is mainly decided by their placement on different servers in the cluster.

To understand how placement affects the training speed, consider a cluster with 3 servers and a synchronous training job using 2 parameter servers and 4 workers. Each server can host 3 parameter servers or workers. The bandwidth at each parameter server or worker is 1. The size of gradients/parameters transferred between a parameter server and a worker in one training step is 1. Fig. 10 illustrates 3 possible ways of placing the workers/parameter servers. With placement (a), the 2 parameter servers and 4 workers need to transfer 3, 3, 1, 1, 2, 2 units of data across servers, respectively. Take *ps1* in Fig. 10 as an example for illustration: it needs to communicate with *worker2*, *worker3* and *worker4* across servers, so it transfers 3 units of data. Note that the bandwidth between a parameter server and a worker is determined by the bandwidth capacity at both ends and the time is decided by the slowest transfer. Therefore, the data transfers of *ps1* and *ps2* are the slowest, and we can obtain that the data transfer time in one training step with placement (a) is 3.

Similarly, the data transfer time with placement (b) is 3 and with placement (c) is 2. Therefore, in this example, placement (c) is the best solution.

THEOREM 1. *Given the numbers of workers and parameter servers in a synchronous training job, the optimal worker/parameter server placement principle to achieve the maximal training speed for the job, in a cluster of homogeneous servers, is to use the smallest number of servers to host the job, such that the same number of parameter servers and the same number of workers are deployed on each of these servers.*

The detailed proof is given in the Appendix. The principles behind the proof are that (a) colocating workers and parameter servers can reduce cross-server data transfers, and (b) packing the same number of workers/parameter servers of a job on each server can minimize the maximal data transfer time in each step of synchronous training. We can also apply these principles to asynchronous training jobs to balance the training speeds of multiple workers.

Based on these principles, we design a placement scheme to minimize the data transfer time during training as follows. We sort all servers in the cluster in descending order of their current resource availability (available CPU capacity is used in our experiments). We place jobs in increasing order of their resource demand (*i.e.*, smallest job first) in order to avoid job starvation (*i.e.*, small jobs do not get any resources). For each job, we check whether the resources on the first k servers are sufficient to host the job (starting with $k = 1$). If so, we place parameter servers and workers in the job evenly on the k servers; otherwise, we check the first $k + 1, k + 2, \dots$ servers until we find enough servers to place the job. We then update available resources on the k servers and sort the server list again. The above procedure repeats until all jobs are placed or no sufficient resources on the servers are left to host more jobs. Note that the number of jobs the servers can accommodate might be smaller than the number of jobs we allocate resource to through the resource allocation algorithm (which considers overall resource capacity in the entire cluster). Jobs which are not placed will be temporarily paused and rescheduled in the next scheduling interval.

5 SYSTEM IMPLEMENTATION

We next present some implementation details of *Optimus*.

5.1 Data Serving

We store training data in Hadoop Distributed File System (HDFS) [3] with a default chunk size of 128MB and a replication factor of 2. At the beginning of a job, we assign a roughly equal number of chunks to each worker in a round-robin manner, so that each worker has a similar workload. When the number of workers changes due to our dynamic scaling, we reassign the data chunks so that the workload on each worker is still balanced.

5.2 Straggler Handling

Stragglers, *i.e.*, slow workers (we will discuss the case of slow parameter servers in §5.3), influences a synchronous training job significantly, due to the need of synchronizing all workers in each training step. For asynchronous training, it is also important to

ensure the workers have similar training speeds so that the parameters on any worker are not too stale; parameter staleness may lead to unstable training progress and hence additional training steps to achieve convergence [27]. In a distributed DL framework, stragglers may happen due to a number of reasons, *e.g.*, resource contention, unbalanced workload.

To detect stragglers in an asynchronous training job, we simply monitor each worker’s training speed: if a worker is too slow (*e.g.*, half speed from the median), we consider it as a straggler. For synchronous training, the training speeds at the workers are the same since they are synchronized. To identify a straggler, we monitor the arrival time of each worker’s gradients on parameter servers and calculate the training speed of each worker as the gap between the arrival time of two steps. We replace a straggler by launching a new worker.

5.3 Load Balancing on Parameter Servers

Our DL jobs are running on the MXNet framework. We identify possible significant load imbalance among parameter servers in MXNet, due to its way of dividing model parameters among parameter servers: for each block of parameters (*i.e.*, parameters of one layer in an NN), if its size (*i.e.*, the number of parameters) is smaller than a threshold (10^6 by default), then it is assigned to one parameter server randomly; otherwise it is sliced evenly among all parameter servers. Setting the threshold is difficult since different models may have different appropriate thresholds, and different threshold values often lead to a big difference in computation workload among parameter servers. Such a load imbalance problem also exists in other distributed ML frameworks such as TensorFlow.

To balance the workload among parameter servers (mainly due to parameter update computation and communication overhead), we seek to minimize (a) the maximal difference of parameter sizes between two parameter servers, (b) the total number of parameter update requests between parameter servers and workers during one training step (each request from a worker asks for one updated parameter block), and (c) the maximal difference of the number of parameter update requests between two parameter servers. We design a parameter assignment algorithm (PAA) as follows.

We sort parameter blocks in decreasing order of size and calculate the average parameter size avg_size , *i.e.*, the overall parameter size divided by the number of parameter servers. For each block, if its size is very small (*e.g.*, less than 1% of avg_size), then we assign it to the parameter server with the least number of update requests. If the block size is between 1% of avg_size and avg_size , we assign the block to the parameter server with the smallest remaining capacity (avg_size minus the size of parameters assigned), that can accommodate it (a best-fit approach). If the block size is larger than avg_size , we further slice it into partitions with size avg_size or less (for the last partition), and assign the sliced partitions to the parameter server with the smallest size of parameters assigned. Once a parameter block (or partition) is assigned to a parameter server, we add the number of parameter update requests on the server by 1.

5.4 Elastic Training on MXNet

To adjust resource allocation to jobs (*i.e.*, numbers of workers and parameter servers) during training, we adopt a checkpoint-based method. When the number of workers or parameter servers assigned to a job changes, we checkpoint the model parameters and save them to HDFS [3]. Then we restart the job from the checkpoint and redeploy parameter servers and workers based on the scheduling decisions. In practical DL clusters, multiple distributed training frameworks may be used. Our approach is simple and general, and can be easily extended for resource scaling in other frameworks with little code modification.

5.5 Scheduler on Kubernetes

We deploy our scheduler *Optimus* as a normal pod (*i.e.*, a unit of deployment that couples one or more containers tightly) on Kubernetes 1.7 [14], which polls the Kubernetes master to obtain cluster information and job states. For fault-tolerance, we use etcd [10] (*i.e.*, a distributed reliable key-value storage) as fault-tolerant storage of job states. Kubernetes will automatically restart the scheduler if it fails.

6 EVALUATION

6.1 Methodology

Testbed. We built a testbed that consists of 7 CPU servers and 6 GPU servers. Each CPU server has two 8-core Intel E5-2650 CPUs, 80GB memory, two 300GB HDDs. Each GPU server has one 8-core Intel E5-1660 CPU, two GeForce 1080Ti GPUs, 48GB memory, one 500GB SSD and one 4TB HDD. They are connected by a 48-port Dell N1548 1GbE switch. We deployed Kubernetes 1.7 [14] and HDFS 2.8 [3] in the cluster.

Simulator. To evaluate *Optimus* at a larger scale of cluster and understand its performance with more parameter choices, we also implemented a discrete-time simulator. The simulator uses the following from the traces collected from our testbed experiments: training losses of each kind of jobs, training speeds under different resource configurations, resource capacities of each server, job configurations (*e.g.*, resource requirements of workers/parameter servers), DL model details (*e.g.*, parameter size).

Workload. Job arrival happens randomly between [0,12000] seconds. Upon an arrival event, we randomly choose the job among the examples in Table 1 and decide to run it using asynchronous training or synchronous training randomly. We vary the convergence threshold of jobs between 1% and 5%. For jobs training large models, *e.g.*, the ResNet-50 model or the DeepSpeech2 model, we downscale their dataset sizes so that the experiment can be finished in a reasonable amount of time, as otherwise each experiment run would last for weeks. We verified that the models still converge with the small datasets. After downscaling, one experiment run takes about 6 hours and we repeat each experiment for 3 times to obtain the average results.

Baselines. We compare *Optimus* with two representative schedulers, implemented on Kubernetes as well: (i) A fairness-based scheduler adopted in many resource managers such as Hadoop [11], Yarn [61] and Mesos [40], which uses Dominant Resource Fairness (DRF) [34] to allocate resources to jobs and dynamically reschedules

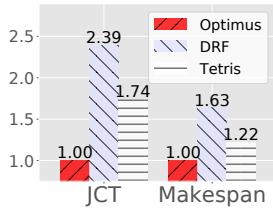


Figure 11: Performance comparison

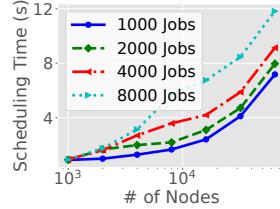


Figure 12: Scalability test

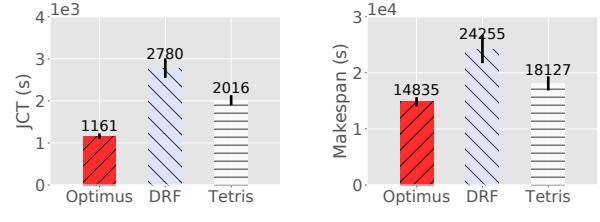


Figure 13: Average value and standard deviation

jobs in each scheduling interval. The workers/parameter servers are placed in a load balancing way, according to the default behavior of Kubernetes. (ii) Tetris [37], which preferentially allocates resources to jobs with low duration or small resource consumption and packs jobs to servers to minimize resource fragmentation. Since Tetris does not have its own mechanism to estimate the remaining time of a deep learning job, we use our speed function and convergence estimation to provide Tetris with such information. We set the ratio of the number of parameter servers to the number of workers to 1:1 [19] in both schedulers.

Metrics. We use the average job completion time (JCT) as an indicator of system performance. In addition, we evaluate the makespan as an indicator of resource efficiency, which is the total time elapsed from the arrival of the first job to the completion of all jobs. Minimizing makespan is equivalent to maximizing resource efficiency [37].

To initialize the training speed function for each job, we pre-run a job on a small dataset with 5 different combinations of (p, w) . Each scheduling interval is 10 minutes long. We set the priority factor in §4.1 to 1 and set the very small parameter block size in §5.3 to 1% of *avg_size* by default.

6.2 Performance

Comparison with baselines. Fig. 11 shows that *Optimus* can reduce the average completion time and makespan by 2.39x and 1.63x respectively in comparison to the DRF-based fairness scheduler. The average value and standard deviation of JCT and makespan are further presented in Fig. 13. We see that *Optimus*, DRF and Tetris use 4.1, 6.7 and 5.0 hours to finish all jobs, respectively. To see more details, Fig. 14 shows the number of running tasks and normalized CPU utilization of tasks in each time slot (*i.e.*, CPU utilization divided by overall allocated CPU capacity on a parameter server or a worker) during the whole experiment run. *Optimus* does not run a large number of tasks as compared to DRF. The reason is that DRF is work-conserving and allocates as many resources to a job as possible, but more resources do not mean higher training speed, as demonstrated in §3.2. Further, the normalized CPU utilization of workers and parameter servers in *Optimus* is larger than that of DRF and Tetris. It shows that *Optimus* can utilize allocated resources more efficiently.

Resource adjustment overhead. The overhead of changing from one (p, w) configuration to another in a job is measured by the percentage of time spent on adjusting resources for the job. In our

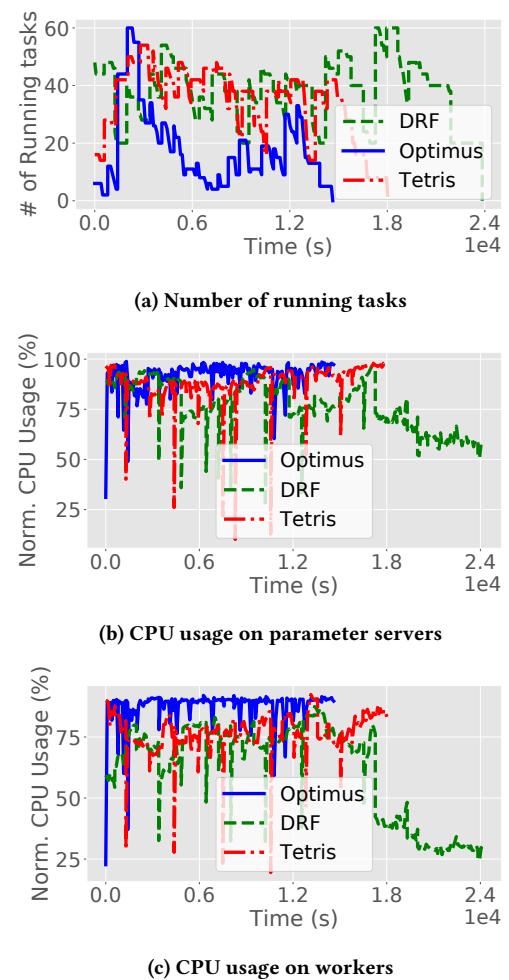
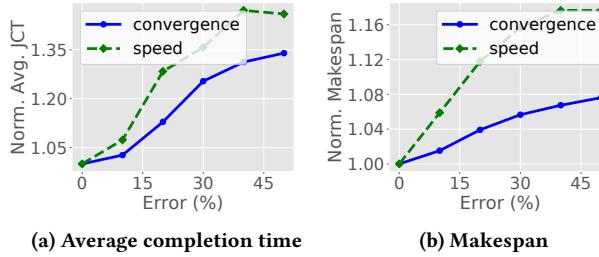


Figure 14: Number of running tasks and CPU usage during an experiment

experiments, the overall scaling overhead is 2.54% of the makespan, which is acceptable compared to the performance gain.

Scalability. To evaluate whether *Optimus* is sufficiently fast and scalable to large-scale clusters, we emulate submitting and scheduling a large number of jobs in a cluster with thousands of nodes.



(a) Average completion time

(b) Makespan

Figure 15: Sensitivity to prediction errors

Figure 12 shows the scheduling time when *Optimus* runs on one core of Intel E5-1620 v4 CPU. *Optimus* can schedule 4,000 jobs (about 100,000 tasks) within 5 seconds on a cluster of 16,000 nodes. This is comparable to the performance of Kubernetes’ default scheduler, *i.e.*, 150,000 tasks in 5,000 nodes within 5 seconds [64]. Besides, since *Optimus* makes scheduling decisions at each scheduling interval (*e.g.*, 10 minutes), the scheduling overhead is very small.

6.3 Sensitivity analysis

Prediction error. We examine to what extent *Optimus* is affected by the prediction errors of convergence time and training speed. We carry out simulation under different error levels: suppose the true number of epochs for convergence (training speed) is v and the error is e ; we use $v \cdot (1+e)$ or $v \cdot (1-e)$ as the initial input to our scheduler, decreasing with job progress. We run each simulation for 100 times to obtain average results.

In Fig. 15, the convergence (speed) curve plots the resulting average JCT/makespan when we add errors of different levels in convergence epoch (training speed) prediction. When the error is larger, JCT and makespan both increase, but with a diminishing speed. If the error of convergence estimation is 20% and the error of training speed estimation is 10%, there is about 15% performance gap compared to the case where the estimation errors are 0. Compared to the error of convergence estimation, the error of speed estimation affects the performance more. Fortunately, we can estimate training speed much more accurate (10% error) than training convergence (20% error).

To see the effect of our technique at the end of §4.1, we have also done evaluation with the priority factor set to 0.95. In this case, the average JCT and makespan are 2.66% and 1.88% smaller, respectively, which validates the effectiveness of our technique in improving the overall scheduling performance.

Varying workloads. We examine how training modes affect the performance. Instead of randomly selecting between asynchronous and synchronous training (§6.1), we either train all jobs in asynchronous mode or synchronous mode. Fig. 16 shows that *Optimus* outperforms the other two schedulers in both cases, and the performance gain is larger when all jobs use synchronous training. This is because all workers have the most updated parameters with synchronous training, such that model convergence is more stable and convergence estimation error is smaller. The training speed of all workers are the same in synchronous training and the speed estimation error is smaller, as verified in §3.2.

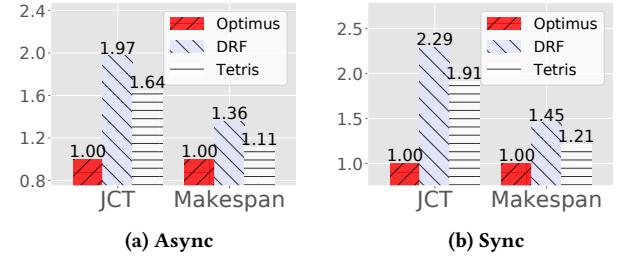


Figure 16: Sensitivity to workloads: training modes

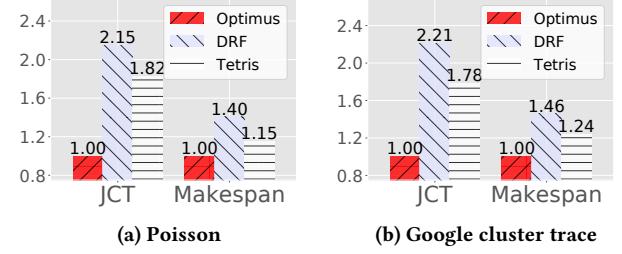


Figure 17: Sensitivity to workloads: job arrival processes

We further investigate *Optimus*’s performance under two other job arrival processes. The first is a Poisson process with 3 arrivals per scheduling interval. The second is extracted from Google cluster workload traces over a 7 hour period [8]. Fig. 17 shows that *Optimus* still outperforms the other two schedulers and the performance gain is larger when using Google cluster traces. There are many job arrival spikes in the traces and *Optimus* can handle them better than DRF and Tetris by efficiently allocating resources.

6.4 Inspecting Detailed Designs in *Optimus*

Resource allocation. To see how effective our marginal gain-based resource allocation algorithm is, we replace it with the resource allocation schemes in the fairness scheduler or Tetris, while still adopting the same task placement algorithm in *Optimus*. Fig. 18 shows that the average completion time and makespan are reduced by 62% and 31% respectively when using *Optimus*, as compared to the fairness scheduler. That is, the resource allocation algorithm in *Optimus* is critical for high job performance and resource efficiency.

Task placement. We further examine the task placement algorithm in *Optimus* to see to what extent it contributes to job performance and resource efficiency. For comparison, we place tasks using the placement algorithm in the fairness scheduler (*i.e.*, in a load-balancing way) and Tetris (*i.e.*, minimizing resource fragmentation), but still use the resource allocation algorithm in *Optimus*. Fig. 19 shows that our algorithm reduces average completion time and makespan by about 10% compared to Tetris and 15% compared to DRF.

Parameter server load balancing. The difference of parameter sizes among parameter servers, the difference of the number of parameter update requests among parameter servers and the total

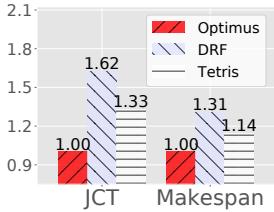


Figure 18: Effectiveness of resource allocation algorithm

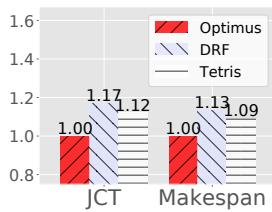


Figure 19: Effectiveness of task placement algorithm

Table 3: Comparison of parameter distribution

| Algorithm | Difference of parameter sizes | Difference of # of requests | Total # of requests |
|-----------|-------------------------------|-----------------------------|---------------------|
| MXNet | 3.6M | 43 | 247 |
| PAA | 0.1M | 1 | 157 |

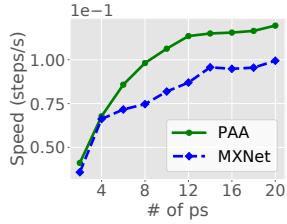


Figure 20: Training speed comparison of ResNet-50 by varying # of ps

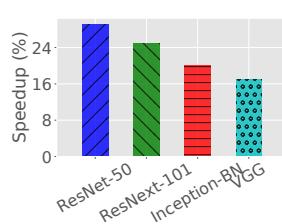


Figure 21: Training speed improvements on different models

number of update requests between parameter servers and workers are three main factors that represent load imbalance or overhead on parameter servers. Table 3 shows values of the three factors achieved with our PAA algorithm (in §5.3) and with the default parameter distribution algorithm in MXNet, using the ResNet-50 model [39] with 25 million parameters formed into 157 parameter blocks. Our algorithm does not split any parameter block further (since the total number of parameter update requests is 157, the minimal for 157 parameter blocks) while keeping minimal the difference of parameter sizes (*i.e.*, 0.1M) and the difference of the number of requests (*i.e.*, 1).

To see the effectiveness of balanced parameter distribution on the training speed, we train ResNet-50 on the ILSVRC2012 ImageNet dataset [12] by fixing the number of workers to 10 and varying the number of parameter servers, using synchronous training. Fig. 20 shows the training speed with and without our load balancing algorithm. We can see PAA improves the training speed especially when the number of parameter servers is large. Fig. 21 further shows the improvement when more models are trained using synchronous training, 10 workers and 10 parameter servers: PAA achieves up to 29% speedup compared to the MXNet algorithm. We observed similar results with asynchronous training.

In summary, the highlights of our evaluation results are as follows.

- (1) Testbed experiments show that *Optimus* improves average job completion time and makespan by 139% and 63% compared to the fairness scheduler. Further, *Optimus* can scale to schedule 100,000 tasks on 16,000 nodes in 5 seconds, and its resource adjustment overhead is small, *i.e.*, 2.54%.
- (2) Further improvement of estimation accuracy will not increase *Optimus*'s performance much (15%) and *Optimus* performs better than DRF and Tetris under various workloads.
- (3) The resource allocation algorithm, the task placement scheme and the parameter server load balancing algorithm contribute to *Optimus*'s performance improvement by about 62%, 17%, 20% respectively.

7 DISCUSSIONS

We now discuss extensions and future work on *Optimus*.

Various workloads. While *Optimus* targets scheduling of deep learning jobs, it can be used in DL clusters with mixed workloads (*e.g.*, data analytics, online services). For example, in a Kubernetes cluster, we can plug in multiple schedulers and each scheduler is responsible for one kind of workloads. In such case, *Optimus* may ask for resources from a central cluster resource manager and schedule deep learning jobs on a varying portion of cluster resources.

Convergence estimation. For some DL models (*e.g.*, ResNet [39]), the learning rate may be reduced significantly (*e.g.*, by a factor of 0.1) when training reaches a predefined condition, in order to minimize loss further (*e.g.*, as with SGD). In such a case, we can treat the model training after learning rate adjustment as a new training job and restart online fitting. In addition, the training loss curves of some models (*e.g.*, A3C [51]) cannot be described or can only be partly described using our fitting function in Eqn. (1), but they may be fitted using other functions based on the convergence speed of optimization algorithm [71]. One possible solution is to let the job owner provide the functions, based on the previous running experience of such jobs [71].

Scaling overhead. We use a checkpoint-based method to adjust the resource configuration of a job due to its simplicity and general implementability. This approach may bring quite large overhead if the job has hundreds of workers/parameter servers. To reduce checkpointing overhead, we may set a threshold of checkpointing times for each job to limit the restarting frequency. For long or large jobs, the threshold can be smaller to avoid frequent resource reallocation.

8 RELATED WORK

Performance modeling. Jockey [32] and Morpheus [44] use historical traces of periodic jobs and dynamically adjust resource allocations to meet deadlines, while *Optimus* does not depend on the previous run of the same job since production training data often change (*e.g.*, daily). PerfOrator [53] builds a resource-to-performance model of big data queries by estimating query size and profiling hardware, while we use high-level system modeling approach without the knowledge about hardware or the internal details of a job. To map resources to training speed, we build and fit a parametric performance model based on sample runs. This approach has been applied in other work, such as job execution

time estimation [56] and data size estimation of SQL queries [53]. Ernest [62] also adopts a similar approach to estimate the completion time of data analytics. It designs an experimental theory to minimize sampling overhead. In comparison, *Optimus* has a relatively small configuration space (*i.e.*, the number of tasks) and 5-10 sample runs are enough for fitting the performance model quite accurately. PREDiT [52] uses sample runs for capturing the convergence trend of a graph algorithm, which is infeasible for deep learning training since the size of dataset affects convergence. Yan *et al.* [69] model the training of deep learning neural networks at a very fine granularity (*e.g.*, the computation time of each operator on a specific CPU, neural network structures, etc.) while our models capture high-level computation and communication patterns. Bayesian Optimization is a parameter-free approach used in many other works (*e.g.*, FABOLAS [47], BOAT [28], CherryPICK [24]) to search best hyperparameters/resource configuration for a model/job. This approach is not applicable to our problem since we need a parametric performance model to describe the relation between the number of tasks and training speed, so that the scheduler can exploit this relation to optimize global scheduling of all concurrent jobs.

Job scheduling. There have been many efforts on cluster/cloud resource allocation to achieve different objectives. Corral [42] and Morpheus [44] focus on periodic or predictable workloads. Borg [63], Fuxi [72], Firmament [35] are designed for heterogeneous workloads in a large-scale cluster and support policy-based scheduling (*e.g.*, fairness, data locality, job priority). Instead, our work focuses on deep learning workload. Mesos [40] and Yarn [61] use DRF [34] to allocate resources while we focus on resource efficiency and job performance. TetriSched [60] and Morpheus [44] also dynamically allocate resources in a global way, but they focus on reservation-based or periodic jobs with specified deadlines. Eagle [30] is a hybrid scheduler designed to solve the head-of-line problem: it dynamically divides the cluster resources into two partitions for short jobs and long jobs. Optimus instead focuses on dynamic resource configurations of jobs. There are several studies [41, 57, 71] on resource allocation of classical machine learning jobs (*e.g.*, clustering, logistic regression) on Spark MLLib [70]. Huang *et al.* [41] propose a memory optimizer for Spark master and workers given a machine learning program. SLAQ [71] targets the training quality of experimental ML models instead of models in production. It adopts similar online fitting technique to estimate the training loss of convex algorithms. Dorm [57] uses a utilization-fairness optimizer to schedule jobs. The main difference is that our work focuses on deep learning jobs running on parameter server architecture. We leverage the characteristics of the jobs to design resource allocation algorithm and task placement scheme, and demonstrate significant performance improvement. STRAD [45] proposes a programming approach to improve model convergence by scheduling parameter updates for model-parallel machine learning, while we did not delve into modifying the underlying ML frameworks. Azalia *et al.* [50] use a model-free deep reinforcement learning method to achieve model parallelism that maximizes training speed of a given model in a single machine. Such an approach is yet to be general and efficient for resource allocation in a deep learning cluster. Proteus [38] exploits transient virtual machines in EC2 to complete ML jobs in

an efficient and cheap way. They use a simpler performance model and focus more on the expected cost due to dynamic bidding prices.

Distributed machine learning frameworks. The parameter server architecture was first introduced in [55] and improved with update primitives, fault tolerance and communication optimization in [29, 48, 67]. Most distributed machine learning frameworks (*e.g.*, MXNet [59], Petuum [67], TensorFlow [23], Angel [43]) are implicitly or explicitly built upon this architecture. Our work targets scheduling jobs running on these frameworks. We find that the load imbalance problem is common in these distributed frameworks and we propose and implement the PAA algorithm in one of the frameworks, MXNet.

9 CONCLUSION

Optimus is a customized cluster scheduler targeting high training performance and resource efficiency in deep learning clusters. At its core is an accurate performance model for deep learning workloads, built by exploiting the characteristics of DL model training (*e.g.*, convergence property, iterativeness) and communication patterns of the parameter server architecture. Based on the performance model, we design a marginal gain-based resource allocation algorithm and a training speed-maximizing task placement scheme. Our experiments on a Kubernetes cluster show that *Optimus* outperforms representative cluster schedulers significantly.

ACKNOWLEDGEMENTS

We thank our shepherd Paolo Romano and the anonymous reviewers for their feedback. This work was supported in part by grants from Hong Kong RGC under the contracts HKU 17204715, 17225516, and C7036-15G (CRF). The Titan X Pascal used for this research was donated by the NVIDIA Corporation.

APPENDIX

Proof of Theorem 1

Proof: Assume there are K nodes (physical servers) in the cluster and the number of parameter servers on node k is p_{jk} for job j , the number of workers on node k is w_{jk} for job j . Assume the capacities of the K nodes are sufficient for placing the job. Let B_j denote the bandwidth requirement of each parameter server in job j and b_j be the bandwidth requirement of each worker in job j . Let S_j be the model size of job j . Then the data (gradients/parameters) transmission time in job j for each training step in case of synchronous training is

$$\max_k \left\{ \frac{\frac{S_j}{p_j} (w_j - w_{jk})}{B_j}, \frac{\frac{S_j}{p_j} (p_j - p_{jk})}{b_j} \right\}$$

Then we can formulate the worker/parameter server placement problem for transmission time minimization as follows.

$$\begin{aligned} \text{minimize} \quad & \max_k \left\{ \frac{\frac{S_j}{p_j} (w_j - w_{jk})}{B_j}, \frac{\frac{S_j}{p_j} (p_j - p_{jk})}{b_j} \right\} \\ \text{subject to:} \quad & \sum_k p_{jk} = p_j \\ & \sum_k w_{jk} = w_j \\ & p_{jk} \in Z^+, w_{jk} \in Z^+ \end{aligned}$$

We decompose the above problem to the following two subproblems whose solutions are guaranteed to be the optimal solution of the above problem. Each subproblem is a lexicographical min-max problem whose optimal solution is to place tasks evenly. Combining the optimal solution of the two subproblems, one optimal solution of the original problem is to place parameter servers evenly and place workers evenly on the K nodes.

Subproblem 1:

$$\begin{aligned} \text{minimize} \quad & \max_k \frac{\frac{S_j}{p_j} (w_j - w_{jk})}{B_j} \\ \text{subject to:} \quad & \sum_k w_{jk} = w_j \\ & w_{jk} \in Z^+ \end{aligned}$$

Subproblem 2:

$$\begin{aligned} \text{minimize} \quad & \max_k \frac{\frac{S_j}{p_j} (p_j - p_{jk})}{b_j} \\ \text{subject to:} \quad & \sum_k p_{jk} = p_j \\ & p_{jk} \in Z^+ \end{aligned}$$

The next step is to prove that a smaller K leads to smaller data transmission time. The proof can be done via mathematical induction. The idea is that a smaller K means more parameter servers and workers on each node, so the amount of transferred data via the inter-server network is smaller and hence the communication time decreases.

REFERENCES

- [1] 2006. Caltech 256 Dataset. http://www.vision.caltech.edu/Image_Datasets/Caltech256/. (2006).
- [2] 2009. The CIFAR-10 Dataset. <https://www.cs.toronto.edu/~kriz/cifar.html>. (2009).
- [3] 2014. HDFS. <https://wiki.apache.org/hadoop/HDFS>. (2014).
- [4] 2014. Kaggle NDSB1 Dataset. <https://www.kaggle.com/c/datasciencetutorial/data>. (2014).
- [5] 2014. Overfitting and Regularization. <https://alliance.seas.upenn.edu/~cis520/dynamic/2017/wiki/index.php?n=Lectures/Overfitting>. (2014).
- [6] 2014. Perplexity Versus Error Rate. <https://nlpers.blogspot.hk/2014/05/perplexity-versus-error-rate-for.html>. (2014).
- [7] 2014. SciPy NNLS. <https://docs.scipy.org/doc/scipy-0.14.0/reference/generated/scipy.optimize.nnls.html>. (2014).
- [8] 2015. Google Cluster Workload Traces. <https://github.com/google/cluster-data>. (2015).
- [9] 2015. LibriSpeech ASR Corpus. <http://www.openslr.org/12/>. (2015).
- [10] 2017. etcd. <https://github.com/coreos/etcd>. (2017).
- [11] 2017. Hadoop CapacityScheduler. <https://hadoop.apache.org/docs/r2.7.4/hadoop-yarn/hadoop-yarn-site/CapacityScheduler.html>. (2017).
- [12] 2017. ImageNet Dataset. <http://www.image-net.org>. (2017).
- [13] 2017. KAGGLE-DSB Model. <https://github.com/apache/incubator-mxnet/tree/master/example/kaggle-ndsb1>. (2017).
- [14] 2017. Kubernetes. <https://kubernetes.io>. (2017).
- [15] 2017. MXNet Neural Machine Translation. <https://github.com/awslabs/sockeye>. (2017).
- [16] 2017. MXNet Official Examples. <https://github.com/apache/incubator-mxnet/tree/master/example>. (2017).
- [17] 2017. PaddlePaddle. <http://www.paddlepaddle.org>. (2017).
- [18] 2017. Penn Tree Bank Dataset. <https://catalog.ldc.upenn.edu/ldc99t42>. (2017).
- [19] 2017. Run Deep Learning with PaddlePaddle on Kubernetes. <http://blog.kubernetes.io/2017/02/run-deep-learning-with-paddlepaddle-on-kubernetes.html>. (2017).
- [20] 2017. Stochastic Gradient Descent. https://en.wikipedia.org/wiki/Stochastic_gradient_descent. (2017).
- [21] 2017. WMT17. <http://www.statmt.org/wmt17/>. (2017).
- [22] 2017. Word Language Model. https://github.com/apache/incubator-mxnet/tree/master/example/gluon/word_language_model. (2017).
- [23] Martín Abadi, Paul Barham, Jianmin Chen, Zhifeng Chen, Andy Davis, Jeffrey Dean, Matthieu Devin, Sanjay Ghemawat, Geoffrey Irving, Michael Isard, et al. 2016. TensorFlow: A System for Large-Scale Machine Learning. In *Proc. of the 12th USENIX Symposium on Operating Systems Design and Implementation (OSDI)*.
- [24] Omid Alipourfard, Hongqiang Harry Liu, Jianshu Chen, Shivaram Venkataraman, Minlan Yu, and Ming Zhang. 2017. CherryPick: Adaptively Unearthing the Best Cloud Configurations for Big Data Analytics. In *Proc. of the 14th USENIX Symposium on Networked Systems Design and Implementation (NSDI)*.
- [25] Dario Amodei, Sundaram Ananthanarayanan, Rishita Anubhai, Jingliang Bai, Eric Battenberg, Carl Case, Jared Casper, Bryan Catanzaro, Qiang Cheng, Guoliang Chen, et al. 2016. Deep Speech 2: End-to-end Speech Recognition in English and Mandarin. In *Proc. of the 33th International Conference on Machine Learning (ICML)*.
- [26] Pang Bo and Lee Lillian. 2005. Movie Review Data. <https://www.cs.cornell.edu/people/pabo/movie-review-data/>. (2005).
- [27] Jianmin Chen, Rajat Monga, Samy Bengio, and Rafal Jozefowicz. 2016. Revisiting Distributed Synchronous SGD. *arXiv preprint arXiv:1604.00981* (April 2016).
- [28] Valentin Dalibard, Michael Schaefer, and Eiko Yoneki. 2017. BOAT: Building Auto-tuners with Structured Bayesian Optimization. In *Proc. of the 26th International Conference on World Wide Web (WWW)*.
- [29] Jeffrey Dean, Greg Corrado, Rajat Monga, Kai Chen, Matthieu Devin, Mark Mao, Andrew Senior, Paul Tucker, Ke Yang, Quoc V Le, et al. 2012. Large Scale Distributed Deep Networks. In *Proc. of the 25th Advances in Neural Information Processing Systems (NIPS)*.
- [30] Pamela Delgado, Diego Didona, Florin Dinu, and Willy Zwaenepoel. 2016. Job-aware Scheduling in Eagle: Divide and Stick to Your Probes. In *Proc. of the 7th ACM Symposium on Cloud Computing (SoCC)*.
- [31] Matteo Dell'Amico, Damiano Carra, Mario Pastorelli, and Pietro Michiardi. 2014. Revisiting Size-Based Scheduling with Estimated Job Sizes. In *Proc. of the 22th IEEE International Symposium on Modelling, Analysis & Simulation of Computer and Telecommunication Systems (MASCOTS)*.
- [32] Andrew D Ferguson, Peter Bodik, Srikanth Kandula, Eric Boutin, and Rodrigo Fonseca. 2012. Jockey: Guarantee Job Latency in Data Parallel Clusters. In *Proc. of the 7th ACM European Conference on Computer Systems (Eurosyst)*.
- [33] Jonas Gehring, Michael Auli, David Grangier, Denis Yarats, and Yann N Dauphin. 2017. Convolutional Sequence to Sequence Learning. In *Proc. of the 34th International Conference on Machine Learning (ICML)*.
- [34] Ali Ghodsi, Matei Zaharia, Benjamin Hindman, Andy Konwinski, Scott Shenker, and Ion Stoica. 2011. Dominant Resource Fairness: Fair Allocation of Multiple Resource Types. In *Proc. of the 8th USENIX Symposium on Networked Systems Design and Implementation (NSDI)*.
- [35] Ionel Gog, Malte Schwarzkopf, Adam Gleave, Robert NM Watson, and Steven Hand. 2016. Firmament: Fast, Centralized Cluster Scheduling at Scale. In *Proc. of the 12th USENIX Conference on Operating Systems Design and Implementation (OSDI)*.
- [36] Priya Goyal, Piotr Dollár, Ross Girshick, Pieter Noordhuis, Lukasz Wesolowski, Aapo Kyrola, Andrew Tulloch, Yangqing Jia, and Kaiming He. 2017. Accurate, Large Minibatch SGD: Training ImageNet in 1 Hour. In *arXiv preprint arXiv:1706.02677*.
- [37] Robert Grandl, Ganesh Ananthanarayanan, Srikanth Kandula, Sriram Rao, and Aditya Akella. 2014. Multi-Resource Packing for Cluster Schedulers. In *Proc. of ACM SIGCOMM*.
- [38] Aaron Harlap, Alexey Tumanov, Andrew Chung, Gregory R Ganger, and Phillip B Gibbons. 2017. Proteus: Agile ML Elasticity Through Tiered Reliability in Dynamic Resource Markets. In *Proc. of the 12th ACM European Conference on Computer Systems (EuroSys)*.
- [39] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. 2016. Deep Residual Learning for Image Recognition. In *Proc. of the 29th IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*.
- [40] Benjamin Hindman, Andy Konwinski, Matei Zaharia, Ali Ghodsi, Anthony D Joseph, Randy H Katz, Scott Shenker, and Ion Stoica. 2011. Mesos: A Platform for Fine-Grained Resource Sharing in the Data Center. In *Proc. of the 8th USENIX Symposium on Networked Systems Design and Implementation (NSDI)*.

- [41] Botong Huang, Matthias Boehm, Yuanyuan Tian, Berthold Reinwald, Shirish Tatikonda, and Frederick R Reiss. 2015. Resource Elasticity for Large-Scale Machine Learning. In *Proc. of ACM SIGMOD*.
- [42] Virajith Jalaparti, Peter Bodik, Ishai Menache, Sriram Rao, Konstantin Makarychev, and Matthew Caesar. 2015. Network-Aware Scheduling for Data-Parallel Jobs: Plan When You Can. In *Proc. of ACM SIGCOMM*.
- [43] Jie Jiang, Leli Yu, Jiawei Jiang, Yuhong Liu, and Bin Cui. 2017. Angel: a New Large-Scale Machine Learning System. *National Science Review* (2017), nwx018.
- [44] Sangeetha Abdu Jyothi, Carlo Curino, Ishai Menache, Shravan Narayananmurthy, Alexey Tumanov, Jonathan Yaniv, Íñigo Goiri, Subru Krishnan, Janardhan Kulkarni, and Sriram Rao. 2016. Morphus: Towards Automated SLOs for Enterprise Clusters. In *Proc. of the 12th USENIX Symposium on Operating Systems Design and Implementation (OSDI)*.
- [45] Jin Kyu Kim, Qirong Ho, Seunghak Lee, Xun Zheng, Wei Dai, Garth A Gibson, and Eri P Xing. 2016. STRADS: A Distributed Framework for Scheduled Model Parallel Machine Learning. In *Proc. of the 11th European Conference on Computer Systems (Eurosys)*.
- [46] Yoon Kim. 2014. Convolutional Neural Networks for Sentence Classification. In *Proc. of 19th SIGDAT Conference on Empirical Methods in Natural Language Processing (EMNLP)*.
- [47] Aaron Klein, Stefan Falkner, Simon Bartels, Philipp Hennig, and Frank Hutter. 2017. Fast Bayesian Optimization of Machine Learning Hyperparameters on Large Datasets. In *Proc. of the 20th International Conference on Artificial Intelligence and Statistics (AISTATS)*.
- [48] Mu Li, David G Andersen, Jun Woo Park, Alexander J Smola, Amr Ahmed, Vanja Josifovski, James Long, Eugene J Shekita, and Bor-Yiing Su. 2014. Scaling Distributed Machine Learning with the Parameter Server. In *Proc. of the 11th USENIX Symposium on Operating Systems Design and Implementation (OSDI)*.
- [49] Mahoney Matt. 2017. text8. <http://mattmahoney.net/dc/>. (2017).
- [50] Azalia Mirhoseini, Hieu Pham, Quoc Le, Mohammad Norouzi, Samy Bengio, Benoit Steiner, Yufeng Zhou, Naveen Kumar, Rasmus Larsen, and Jeff Dean. 2017. Device Placement Optimization with Reinforcement Learning. In *Proc. of the 34th International Conference on Machine Learning (ICML)*.
- [51] Volodymyr Mnih, Adria Puigdomenech Badia, Mehdi Mirza, Alex Graves, Timothy Lillicrap, Tim Harley, David Silver, and Koray Kavukcuoglu. 2016. Asynchronous Methods for Deep Reinforcement Learning. In *Proc. of the 33th International Conference on Machine Learning (ICML)*.
- [52] Adrian Daniel Popescu, Andrey Balmin, Vuk Ercegovac, and Anastasia Ailamaki. 2013. PREDICt: Towards Predicting the Runtime of Large-Scale Iterative Analytics. *Proceedings of the VLDB Endowment (PVLDB)* 6, 14 (2013), 1678–1689.
- [53] Kaushik Rajan, Dharmesh Kakadia, Carlo Curino, and Subru Krishnan. 2016. PerfOrator: Eloquent Performance Models for Resource Optimization. In *Proc. of the 7th ACM Symposium on Cloud Computing (SoCC)*.
- [54] Yelong Shen, Xiaodong He, Jianfeng Gao, Li Deng, and Grégoire Mesnil. 2014. A Latent Semantic Model with Convolutional-Pooling Structure for Information Retrieval. In *Proc. of the 23th ACM International Conference on Conference on Information and Knowledge Management (CIKM)*.
- [55] Alexander Smola and Shravan Narayananmurthy. 2010. An Architecture for Parallel Topic Models. *Proceedings of the VLDB Endowment (PVLDB)* 3, 1-2 (2010), 703–710.
- [56] Evan R Sparks, Ameet Talwalkar, Daniel Haas, Michael J Franklin, Michael I Jordan, and Tim Kraska. 2015. Automating Model Search for Large Scale Machine Learning. In *Proc. of the 6th ACM Symposium on Cloud Computing (SoCC)*.
- [57] Peng Sun, Yonggang Wen, Nguyen Binh Duong Ta, and Shengen Yan. 2017. Towards Distributed Machine Learning in Shared Clusters: A Dynamically-Partitioned Approach. In *Proc. of the 3rd IEEE International Conference on Smart Computing (SMARTCOMP)*.
- [58] Christian Szegedy, Vincent Vanhoucke, Sergey Ioffe, Jon Shlens, and Zbigniew Wojna. 2016. Rethinking the Inception Architecture for Computer Vision. In *Proc. of the 29th IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*.
- [59] Chen Tianqi, Li Mu, Li Yutian, Lin Min, Wang Naiyan, Wang Minjie, Xiao Tianjun, Xu Bing, Zhang Chiyuan, and Zhang Zheng. 2016. MXNet: A Flexible and Efficient Machine Learning Library for Heterogeneous Distributed Systems. In *Proc. of NIPS Workshop on Machine Learning Systems (LearningSys)*.
- [60] Alexey Tumanov, Timothy Zhu, Jun Woo Park, Michael A Kozuch, Mor Harchol-Balter, and Gregory R Ganger. 2016. Tetrisched: Global Rescheduling with Adaptive Plan-Ahead in Dynamic Heterogeneous Clusters. In *Proc. of the 11th ACM European Conference on Computer Systems (Eurosys)*.
- [61] Vinod Kumar Vaipaliappalli, Arun C Murthy, Chris Douglas, Sharad Agarwal, Mahadev Konar, Robert Evans, Thomas Graves, Jason Lowe, Hitesh Shah, Siddharth Seth, et al. 2013. Apache Hadoop Yarn: Yet Another Resource Negotiator. In *Proc. of the 4th annual Symposium on Cloud Computing (SoCC)*.
- [62] Shivaram Venkataraman, Zongheng Yang, Michael Franklin, Benjamin Recht, and Ion Stoica. 2016. Ernest: Efficient Performance Prediction for Large-Scale Advanced Analytics. In *Proc. of the 13th USENIX Symposium on Networked Systems Design and Implementation (NSDI)*.
- [63] Abhishek Verma, Luis Pedrosa, Madhukar Korupolu, David Oppenheimer, Eric Tune, and John Wilkes. 2015. Large-Scale Cluster Management at Google with Borg. In *Proc. of the 10th ACM European Conference on Computer Systems (Eurosys)*.
- [64] Tyczynski Wojciech. 2017. Kubernetes Scalability. <http://blog.kubernetes.io/2017/03/scalability-updates-in-kubernetes-1.6.html>. (2017).
- [65] Yonghui Wu, Mike Schuster, Zhifeng Chen, Quoc V Le, Mohammad Norouzi, Wolfgang Macherey, Maxim Krikun, Yuan Cao, Qin Gao, Klaus Macherey, et al. 2016. Google's Neural Machine Translation System: Bridging the Gap Between Human and Machine Translation. *arXiv preprint arXiv:1609.08144* (2016).
- [66] Saining Xie, Ross Girshick, Piotr Dollár, Zhuowen Tu, and Kaiming He. 2017. Aggregated Residual Transformations for Deep Neural Networks. In *Proc. of the 30th IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*.
- [67] Eric P Xing, Qirong Ho, Wei Dai, Jin-Kyu Kim, Jinliang Wei, Seunghak Lee, Xun Zheng, Pengtao Xie, Abhimanyu Kumar, and Yaoliang Yu. 2015. Petuum: A New Platform for Distributed Machine Learning on Big Data. In *Proc. of the 21th ACM International Conference on Knowledge Discovery and Data Mining (SIGKDD)*.
- [68] Wayne Xiong, Jasha Droppo, Xuedong Huang, Frank Seide, Mike Seltzer, Andreas Stolcke, Dong Yu, and Geoffrey Zweig. 2017. The Microsoft 2016 Conversational Speech Recognition System. In *Proc. of the 42th IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*.
- [69] Feng Yan, Olatunji Ruwase, Yuxiong He, and Trishul Chilimbi. 2015. Performance Modeling and Scalability Optimization of Distributed Deep Learning Systems. In *Proc. of the 21th ACM International Conference on Knowledge Discovery and Data Mining (KDD)*.
- [70] Matei Zaharia, Mosharaf Chowdhury, Tathagata Das, Ankur Dave, Justin Ma, Murphy McCauley, Michael J Franklin, Scott Shenker, and Ion Stoica. 2012. Resilient Distributed Datasets: A Fault-Tolerant Abstraction for In-Memory Cluster Computing. In *Proc. of the 9th USENIX conference on Networked Systems Design and Implementation (NSDI)*.
- [71] Haoyu Zhang, Logan Stafman, Andrew Or, and Michael J Freedman. 2017. SLAQ: Quality-Driven Scheduling for Distributed Machine Learning. In *Proc. of the 8th ACM Symposium on Cloud Computing (SoCC)*.
- [72] Zhuo Zhang, Chao Li, Yangyu Tao, Renyu Yang, Hong Tang, and Jie Xu. 2014. Fuxi: a Fault-Tolerant Resource Management and Job Scheduling System at Internet Scale. *Proceedings of the VLDB Endowment (PVLDB)* 7, 13 (2014), 1393–1404.