

DL²: A Deep Learning-driven Scheduler for Deep Learning Clusters

Yanghua Peng¹ Yixin Bao¹ Yangrui Chen¹ Chuan Wu¹ Chen Meng² Wei Lin²

¹{yhpeng, yxbao, yrchen, cwu}@cs.hku.hk ²{mc119496, weilin.lw}@alibaba-inc.com

¹*The University of Hong Kong* ²*Alibaba Inc.*

Abstract

More and more companies have deployed machine learning (ML) clusters, where deep learning (DL) models are trained for providing various AI-driven services. Efficient resource scheduling is essential for maximal utilization of expensive DL clusters. Existing cluster schedulers either are agnostic to ML workload characteristics, or use scheduling heuristics based on operators’ understanding of particular ML framework and workload, which are less efficient or not general enough. In this paper, we show that DL techniques can be adopted to design a generic and efficient scheduler.

DL² is a DL-driven scheduler for DL clusters, targeting global training job expedition by dynamically resizing resources allocated to jobs. DL² advocates a joint supervised learning and reinforcement learning approach: a neural network is warmed up via offline supervised learning based on job traces produced by the existing cluster scheduler; then the neural network is plugged into the live DL cluster, fine-tuned by reinforcement learning carried out throughout the training progress of the DL jobs, and used for deciding job resource allocation in an online fashion. By applying past decisions made by the existing cluster scheduler in the preparatory supervised learning phase, our approach enables a smooth transition from existing scheduler, and renders a high-quality scheduler in minimizing average training completion time. We implement DL² on Kubernetes and enable dynamic resource scaling in DL jobs on MXNet. Extensive evaluation shows that DL² outperforms fairness scheduler (*i.e.*, DRF) by 44.1% and expert heuristic scheduler (*i.e.*, Optimus) by 17.5% in terms of average job completion time.

1 Introduction

Recent years have witnessed the breakthrough of DL-based techniques in various domains, such as machine translation [16], image classification [35] and speech recognition [28]. Large companies have deployed ML clusters with tens to thousands of expensive GPU servers, and run distributed training jobs on one or different

distributed ML frameworks (such as TensorFlow [15], MXNet [20], Petuum [65] and PaddlePaddle [11]), to obtain DL models in need for their AI-driven services. Even with parallel training, training a DL model is commonly very time and resource intensive. Efficient resource scheduling is crucial in operating a shared DL cluster with multiple training jobs, for best utilization of expensive resources and expedited training completion.

Two camps of schedulers exist in today’s ML clusters. In the first camp, general-purpose cloud/cluster schedulers are applied, and possibly customized, for distributed ML job scheduling. For example, Google uses Borg [59] as its DL cluster scheduler; Microsoft, Tencent, and Baidu use custom versions of YARN-like schedulers [58] for managing DL jobs. Representative scheduling strategies used include First-In-First-Out (FIFO) and Dominant Resource Fairness (DRF) [24]. These schedulers allocate resources according to user specification and do not adjust resource allocation during training. As we will see in §2.2, setting the right amount of resources for a job is difficult and static resource allocation leads to resource under-utilization in the cluster.

In the second camp, recent studies have proposed white-box heuristics for resource allocation in ML clusters [67][49][18]. Typically they tackle the problem in two steps: set up analytical models for DL/ML workloads, and propose scheduling heuristics accordingly for online resource allocation and adjustment. Designing heuristics requires a deep understanding of ML frameworks and workloads, and the analytical model is tightly coupled with the ML framework implementation (*e.g.*, a new feature or optimization in evolving ML frameworks may invalidate the analytical model) [49]. Further, the modeling typically does not consider interference in a multi-tenant cluster, where in average 27.3% performance variation may happen (§2.2).

In this paper, we pursue a DL cluster scheduler that does not depend on expert heuristics and explicit modeling, resorting to a black-box end-to-end approach enabled by modern learning techniques. We propose DL², a deep learning-driven scheduler for deep learning clusters, that dynamically adjusts resource allocation to train-

ing jobs on the go. DL^2 learns resource allocation policies through experience using deep reinforcement learning (DRL): the policy neural network takes the current system state as input, produces resource allocation decisions for all the current training jobs and gradually improves the decisions based on feedback. However, merely applying off-the-shelf RL algorithms to scheduling does not produce high-quality decisions, and careful design according to the problem nature is in need.

Existing DRL applications in resource scheduling scenarios [39][41] [42] (§8) use simulators to generate training data for offline training, and apply trained models for resource scheduling in a live system. The core of such a simulator is typically an explicit performance model as mentioned above, and hence the inaccuracy of the simulator may lead to low-quality trained model. Instead of extensive offline training over large simulation, DL^2 takes a different approach: we bootstrap the model using minimal offline supervised learning with any available historical job traces and decisions of any existing scheduling strategy employed in the cluster; then we use online training with feedback from ongoing decision making in a live system, with carefully designed techniques to guide model convergence to high-quality decisions, which minimize average job completion time in the cluster.

In summary, we make the following contributions in DL^2 :

- ▷ In contrast to previous DL cluster scheduling approaches that require analytical performance model and job profiling, DL^2 adopts a more generic design, *i.e.*, using DRL to schedule DL workloads. Instead of simulation-driven RL model training, we adopt online training with real feedback from online resource allocation (§2).

- ▷ We identify that direct application of simple RL approaches for our online scheduler training often leads to poor decisions. To avoid poor decisions at the beginning of online RL, we apply past decisions made by an existing scheduler in the DL cluster in a preparatory offline supervised learning stage. Our approach enables a smooth transition from an existing scheduler, and automatically learns a better scheduler beyond the performance level of the existing one (§3). To optimize online RL particularly for DL job scheduling, we propose job-aware exploration for efficient exploration in the action space, and adopt additional training techniques (*e.g.*, actor-critic algorithm, experience replay) for sample-efficient learning (§4).

- ▷ We design and implement elastic scaling in a popular distributed ML framework, MXNet [20], to achieve dynamic worker/parameter server adjustment (§5). We integrate DL^2 with Kubernetes [9], and carefully evaluate DL^2 using testbed experiments and controlled simulations, driven by DL job patterns collected from a production DL cluster. Evaluation results show that DL^2 signifi-

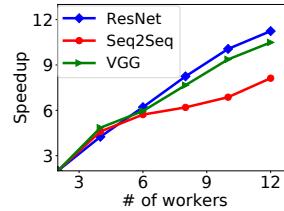


Figure 1: Training speedup with diff. worker/PS numbers

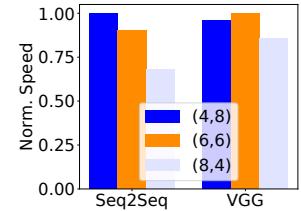


Figure 2: Training speed under diff. PS-to-worker ratios

cantly outperforms other representative schedulers in various scenarios, *e.g.*, 44.1% improvement in average job completion time as compared to the widely adopted DRF scheduler. We also demonstrate DL^2 ’s scaling overhead and generality (§6).

2 Background and Motivation

2.1 The Parameter Server Framework

We focus on the parameter server (PS) architecture [37], which is widely adopted in distributed ML learning frameworks for parallel training, such as in MXNet [20], TensorFlow [15], PaddlePaddle [11] and Angel [32]. Note that DL^2 can also be extended to all-reduce, as discussed in §7. In the PS architecture, the model, *e.g.*, a deep neural network (DNN), is partitioned among multiple parameter servers (PSs) and training dataset is split among workers (*i.e.*, in the representative data parallel training model). The data partition at each worker is divided into mini-batches: each worker processes a mini-batch locally and computes model parameter changes, typically expressed as gradients. The gradients are pushed to PSs which maintain global model parameters. We focus on synchronous training, where the workers’ training progresses are synchronized and PSs update the global model after receiving gradients from all workers in each iteration. Updated parameters are sent back to the workers. A worker starts the next training iteration/step by processing the next mini-batch with the updated parameters. After all mini-batches in the entire dataset have been processed once, one *training epoch* is done. The input dataset is usually trained for multiple epochs until the model converges.

2.2 Motivations

The typical workflow for a user to train a model in a DL cluster is as follows: The user specifies how many PSs and workers she/he wishes to use and the amount of resources (*e.g.*, GPU, CPU) each PS/worker needs, and then submits the job to the scheduler (*e.g.*, Borg [59], YARN [58], Mesos [31]). The scheduler allocates PSs and workers to

the job according to both user demand and its scheduling strategy, and the allocated resources then remain fixed over the entire training course of the job. This workflow has two limitations, as illustrated below.

Difficulty in setting the right worker/PS numbers.

How does a job’s training speed improve when more PSs and workers are added to the job? We train 3 classical models, *i.e.*, ResNet-50 [30], VGG-16 [53] and Seq2Seq [23], in our testbed of 6 machines (see §6 for hardware details), and measure their training speeds (in terms of the number of samples trained per unit time), when increasing the number of workers and keeping the number of PSs equal to the worker number. Each worker uses 1 GPU, 4 CPU cores, 10GB memory and each PS has 4 CPU cores, 10GB memory. In Fig. 1, the speedup is calculated by dividing the training speed achieved using multiple workers/PSs (they are deployed onto physical machines in a load-balanced fashion) by the training speed obtained using one worker and one PS colocated on a single machine. We observe a trend of decreasing return, *i.e.*, adding PSs/workers does not improve the training speed linearly. This is because more communication overhead is incurred when there are more PSs or workers.

On the other hand, is an equal number of PSs and workers (as a general rule of thumb) always the best? We fix the total number of PSs and workers to be 12 and measure the training speed of two models under different combinations of PS/worker numbers (*i.e.*, 4:8, 6:6, 8:4) [49]. Fig. 2 shows that Seq2Seq achieves highest training speed when there are 4 PSs and 8 workers, while VGG-16 is trained fastest with 6 PSs and 6 workers.

From the above, we see that it is challenging to reason about which job will have the largest marginal gain from extra resources and what the best PS-to-worker ratio is, as they are affected by many factors, *e.g.*, allocated resources, models. Existing schedulers largely side-step this problem and leave it to the user to decide how many PSs/workers to use.

Static resource allocation. The GPU cluster resources are often not fully utilized: when a training job is completed, the resources it releases (*e.g.*, expensive GPUs) may become idle, rather than being exploited by remaining jobs that are still running. Fig. 3 shows the GPU utilization during a 24-hour interval in a production DL cluster with about 1000 P100 GPU cards (company name removed due to anonymity requirement), whose job traces will be used in our evaluation (§6). We see that the GPU utilization level varies significantly over time, providing opportunity for dynamic resource scaling out/in in training jobs when cluster load is low/high.

We advocate dynamic adjustment of worker/PS numbers in training jobs over time, to maximally utilize available resources in the DL cluster to expedite job completion. With this, we further do not require users to submit

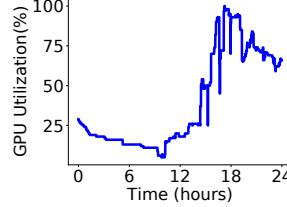


Figure 3: GPU utilization in a production DL cluster

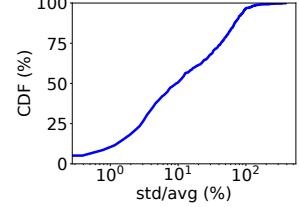


Figure 4: Variation of training completion time

the number of workers/PSs they want to use in their jobs (who nonetheless may not be at the best position to decide that), but will decide the best worker/PS numbers for each user at each time based on both global resource availability and individual jobs’ performance.

White-box heuristics. There have been existing studies which explicitly model detailed relationship between the training speed and resources within jobs, and design scheduling heuristics based on the resource-speed model, *e.g.*, SLAQ [67], Optimus [49] and OASiS [18]. They have two limitations. *First*, in order to derive an accurate performance model, the modeling process is coupled tightly with ML framework implementation, and re-modeling is often needed when the framework changes (*e.g.*, adding new features or adopting optimization). For example, Optimus models computation and communication as two separate procedures during one training step; its model needs to be rebuilt when new features are incorporated into ML frameworks, *e.g.*, overlapping backward computation with communication, gradient compression [20].

Second, explicit performance models are built without considering interference in multi-tenant GPU clusters. For example, SLAQ [67] and Optimus [49] assume no network congestion on PSs, and OASiS [18] and Optimus [49] assume that the available bandwidth is a constant. However, we observe that the speed for training the same model may vary significantly. Fig. 4 shows the performance variation (*i.e.*, the standard deviation of completion time of a training job divided by average completion time of the job over its multiple runs) of 898 DL jobs from the production ML cluster trace. The average variation is 27.3% and the variation for some jobs (3.5% of all jobs) even exceeds 100%. Besides, explicitly modeling interference among ML jobs is also very difficult [17], as each additional dimension (neural network structure, parallelism architecture, runtime isolation, etc.) increases complexity.

In contrast to white-box model-based schedulers, we resort to a black-box approach and design an RL-based resource scheduler: it automatically learns end-to-end resource allocation policy without requiring expert heuristics and without explicitly modeling the ML framework, the workload, and the interference.

2.3 Deep Reinforcement Learning

DRL has been widely used for sequential decision making in an unknown environment, where the agent learns a policy to optimize a cumulative reward by trial-and-error interactions with the environment [55]. In each iteration, the agent observes the current state of the environment and then chooses an action based on the current policy. The environment moves to a new state and reveals the reward, and the policy is updated based on the received reward.

Existing DRL-based schedulers for resource allocation [39][21][41][42] generate a large amount of traces for offline DRL model training, typically by building an explicit resource-performance model for jobs and using it to estimate job progress based on the allocated resources, in the offline simulation environment. The need for model rebuilding (due to ML system changes) and inaccuracy (due to interference) of the performance model degrade the quality of the DRL policy learned (see Fig. 9). Another possibility is to use available historical traces for offline DRL training. However, due to the large decision space of resource allocation (exponential with the amount of resources), historical traces usually do not include feedback for all possible decisions produced by the DRL policy [39][41][17].

Therefore, instead of offline training in a simulated environment, we advocate online RL in the live cluster, exploiting true feedback for resource allocation decisions produced by the DRL agent, to learn a good policy over time. Pure online learning of the policy network model from scratch can result in poor policies at the beginning of learning (see Fig. 10). To avoid poor initial decisions and for the smooth transition from an existing scheduler, we adopt offline supervised learning to bootstrap the DRL policy with the existing scheduling strategy.

3 DL² Overview

The ultimate goal of DL² is to find the best resource allocation policy in a live DL cluster and minimize the average job completion time among all concurrent jobs.

3.1 DL Cluster

In the DL cluster with multiple GPU servers, DL training jobs are submitted over time. Each job runs a distributed ML framework (*e.g.*, MXNet, as in our experiments) to learn a specific DL model by repeatedly training its dataset.

Upon submission of a job, the user, *i.e.*, job owner, provides her/his resource demand to run each worker and each PS, respectively, as well as the total number of training epochs to run. For example, a worker often requires

at least 1 GPU and a PS needs many CPU cores. The total training epoch number to achieve model convergence (*e.g.*, the convergence of loss or accuracy of the model) can be estimated based on expert knowledge or job history.

Depending on resource availability and training speeds, each job may run over a different number of workers and PSs from one time slot to the other (as decided by the scheduler). For synchronous training, to guarantee the same training result (model) while varying the number of workers, we adjust the mini-batch size of each worker, so that the total batch size in a job, as specified by the user, still remains unchanged [49][26]. For asynchronous training, the mini-batch size of each worker remains the same while the number of workers varies (as the global batch size equals each worker’s batch size).

3.2 DL² Scheduler

Our DL-based scheduler, DL², adopts joint offline and online learning of a policy neural network (NN) for making resource allocation decisions to concurrent jobs in the cluster. An overview of DL² is given in Fig. 5.

Offline supervised learning. For warm-up, we use supervised learning to train the policy NN, to initialize a policy whose performance is as good as the existing scheduler in the DL cluster. A small set of historical job runtime traces collected from the cluster are used for supervised learning, to allow the NN to produce similar decisions as made by the existing scheduler. This step is a must due to the poor performance of applying online RL directly (see Fig. 10).

Online reinforcement learning. Online RL works in a time-slotted fashion; each time slot is a scheduling interval, *e.g.*, 1 hour. At the beginning of a scheduling interval, the policy NN takes the information of all the concurrent jobs as input state, and produces the numbers of workers and PSs for each job. The concurrent jobs include new jobs arrived in the previous time slot (after previous scheduling) and jobs which were submitted earlier and whose training has not been completed yet. Workers and PSs are placed on physical machines following the placement policy in the cluster, such as load balancing [51]. Jobs’ training progress is observed at the end of each time slot, and used as the reward to improve the policy network.

4 Detailed Design

4.1 Policy Neural Network

State. The input state to the policy NN is a matrix $s = (\vec{x}, \vec{d}, \vec{e}, \vec{r}, \vec{w}, \vec{u})$, including the following (Fig. 6):

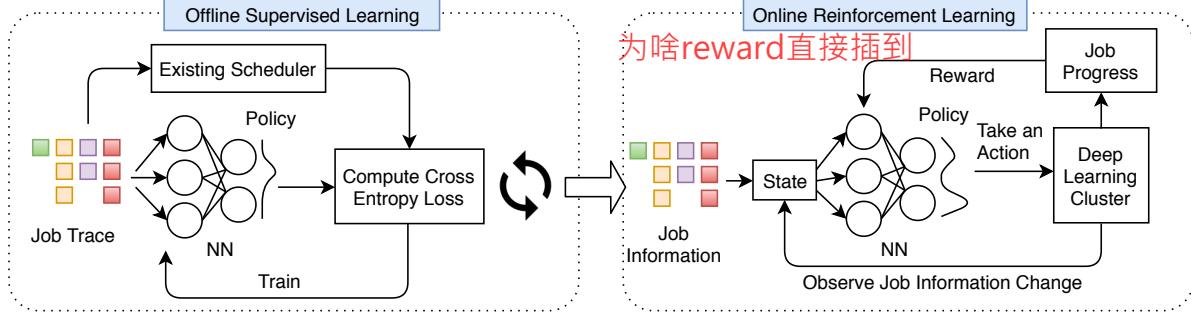


Figure 5: An overview of DL^2

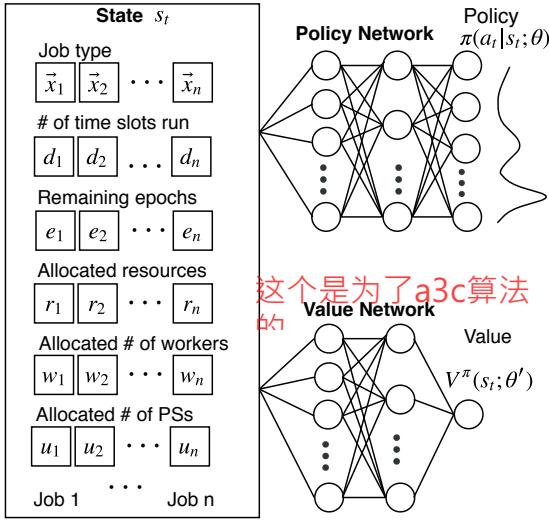


Figure 6: Actor-critic reinforcement learning

- x , a $J \times L$ matrix representing the DL models trained in the jobs, where J is an upper bound of the maximal number of concurrent jobs in a time slot that we are scheduling, and L is the maximal number of training job types in the cluster at all times. We consider DL jobs training similar DNN architecture as the same type in our input. For example, fine-tuning jobs based on the same pre-trained model is common¹ and they can be treated as the same type. Each vector \vec{x}_i in x , $\forall i = 1, \dots, J$, is a one-hot encoding of job i 's type. For example, if there are 3 job types in total and 3 concurrent jobs in each type respectively, then $x = \{[100]; [010]; [001]\}$.

- d , a J -dimensional vector encoding the number of time slots that each job has run in the cluster, for all jobs. For example, d_i is the number of time slots that job i has run.

- \vec{e} , a J -dimensional vector encoding the remaining number of epochs to train for each job. e_i is the difference between user-specified total training epoch number

¹Many computer vision jobs use pre-trained ResNet [30] model as initialization for training on a target dataset. Similarly, natural language understanding jobs use BERT [22] model to initialize training.

for job i and the number of epochs trained till current time slot.

- \vec{r} , a J -dimensional vector representing the amount of dominant resource already allocated to each job in the current time slot. For example, r_i is the amount of dominant resource (the type of resource the job occupies most as compared to the overall capacity of the resource in the cluster) allocated to job i by resource allocation decisions already made through inferences in this time slot.

- \vec{w} and \vec{u} , each of them is a J -dimensional vector where the i -th item is the number of workers (PSs) allocated to job i in the current time slot.

Information of concurrent jobs in different components of the state are ordered according to the jobs' arrival times. The input state does not directly include available resource capacities in the cluster; our scheduler can handle time-varying overall resource capacities in the cluster.

Action. The NN produces a policy $\pi : \pi(a | s; \vec{\theta}) \rightarrow [0, 1]$, which is a probability distribution over the action space. a represents an action, and $\vec{\theta}$ is the current set of parameters in the NN. A straightforward design is to allow each action to specify the numbers of workers/PSs to allocate to all concurrent jobs; this leads to an exponentially large action space, containing all possible worker/PS number combinations. A large action space incurs significant training cost and slow convergence [60].

To expedite learning of the NN, we simplify the action definition, and allow the NN to output an action out of the following $3 \times J + 1$ actions through each inference: (i) $(i, 0)$, meaning allocating one worker to job i , (ii) $(i, 1)$, allocating one PS to job i , (iii) $(i, 2)$, allocating one worker and one PS to job i , (iv) a void action which indicates stopping allocating resources in the current time slot (as allocating more resources does not necessarily lead to higher training speed [49]). Since each inference only outputs an incremental amount of resources to be allocated to one of J jobs, we allow multiple inferences over the NN for producing the complete set of resource allocation decisions in each time slot: after producing one action, we update state s , and then use the NN to produce another action; the inference repeats until the resources are used

up or a void action is produced. The void action indicates that further allocating resources to jobs no longer improves training speeds.

Though we produce the worker/PS numbers for each job anew in each time slot, for a job that has been running in the previous time slot, we compare the new and previous numbers and perform dynamic scaling to adjust the deployment numbers only (§5).

NN architecture. The input state matrix s is connected to a fully connected layer with the ReLU [48] function for activation. The number of neurons in this layer is proportional to the size of the state matrix. Output from this layer is aggregated in a hidden fully connected layer, which is then connected to the final output layer. The final output layer uses the softmax function [25] as the activation function. The NN architecture is designed based on empirical training trials.

4.2 Offline Supervised Learning

In offline supervised learning, we use stochastic gradient descent (SGD) [56] to update parameters $\vec{\theta}$ of the policy NN to minimize a loss function, which is the cross entropy of the resource allocation decisions made by the NN and decisions of the existing scheduler in the traces [38]. The NN is repeatedly trained using the trace data, *e.g.*, hundreds of times as in our experiments, such that the policy produced by the NN converges to the policy of the existing scheduler.

4.3 Online Reinforcement Learning

Reward. DL² targets average job completion time minimization in the entire cluster. Job completion time would be a natural reward to observe, but it is only known when a job is finished, which may well be hundreds of time slots later. **The significant feedback delay of the reward is unacceptable for online RL, since the delayed reward provides little guidance to improve the early decisions.** We design a per-timeslot reward to collect more reward samples through the job processes, for more frequent RL model updates to expedite convergence. The per-timeslot reward is the sum of normalized number of epochs that the concurrent jobs have trained in this time slot, where the number of epochs trained in job i (t_i) is normalized over the overall number of epochs to train for the job (E_i):

$$r_t = \sum_{i \in [J]} \frac{t_i}{E_i}, \forall t = 1, \dots, \quad (1)$$

The rationale is that the more epochs a job runs in a time slot, the fewer time slots it takes to complete, and hence maximizing cumulative reward amounts to minimizing average job completion time. The normalization is to prevent bias towards large jobs.

Policy gradient-based learning. In online RL, the policy NN obtained through offline supervised learning is further trained using the REINFORCE algorithm [62], to maximize the expected cumulative discounted reward $\mathbb{E}[\sum_{t=0}^{\infty} \gamma^t r_t]$, where $\gamma \in (0, 1)$ is the discount factor. We model the problem as a non-linear one with long-term impact instead of a traditional linear model with one-round independent feedback, *e.g.*, contextual bandit [36], because actions in different time slots are correlated. The REINFORCE algorithm updates the policy network’s parameters, $\vec{\theta}$, by performing SGD on $\mathbb{E}[\sum_{t=0}^{\infty} -\gamma^t r_t]$. The gradient is:

$$\nabla_{\vec{\theta}} \mathbb{E}_{\pi} \left[\sum_{t=0}^{\infty} -\gamma^t r_t \right] = \mathbb{E}_{\pi} \left[-\nabla_{\vec{\theta}} \log(\pi(a | s; \vec{\theta})) Q(a, s; \vec{\theta}) \right] \quad (2)$$

where the Q value, $Q(a, s; \vec{\theta})$ represents the “quality” of the action a taken in given state s following the policy $\pi(\cdot; \vec{\theta})$, calculated as the expected cumulative discounted reward to obtain after selecting action a at state s following $\pi(\cdot; \vec{\theta})$. Each Q value can be computed (empirically) using a mini-batch of samples [56]. Each sample is a four-tuple, (s, a, s', r) , where s' is the new state after action a is taken in state s .

Note that our system runs differently from standard RL: we do multiple inferences (*i.e.*, produce multiple actions) using the NN in each time slot t ; the input state changes after each inference; we only observe the reward and update the NN once after all inferences in the time slot are done. We can obtain multiple samples in a time slot t , and set the reward in each sample to be the reward (1) observed after all inferences are done in t .

We further adopt a number of techniques to stabilize online RL, expedite policy convergence, and improve the quality of the obtained policy.

Actor-critic. We improve the basic policy gradient-based RL with the actor-critic algorithm [46] (illustrated in Fig. 6), for faster convergence of the policy network. The basic idea is to replace Q value in Eqn. 2 with an advantage, $Q(a, s; \vec{\theta}) - V^{\pi}(s; \vec{\theta})$, where $V^{\pi}(s; \vec{\theta})$ is a value function, representing the expected reward over the actions drawn using policy $\pi(a | s; \vec{\theta})$ at all times starting from time slot t . The advantage shows how much better a specific action is, as compared to the expected reward of taking actions according to $\pi(a | s; \vec{\theta})$ in the current state. Using the advantage in computing the policy gradients ensures a much lower variance in the gradients, such that policy learning is more stable.

The value function is evaluated by a value network, which has the same NN structure as the policy network except that its final output layer is a linear neuron without any activation function [46], and it produces the estimate of value function $V^{\pi}(s; \vec{\theta})$. The input state to the

value network is the same as that to the policy network. We train the value network using the temporal difference method [46].

Job-aware exploration. To obtain a good policy through RL, we need to ensure that the action space is adequately explored (*i.e.*, actions leading to good rewards can be sufficiently produced); as otherwise, the RL may well converge to poor local optimal policy [60][46]. We first adopt a commonly used entropy exploration method, by adding an entropy regularization term $\beta \nabla_{\vec{\theta}} H(\pi(\cdot | s; \vec{\theta}))$ in gradient calculation to update the policy network [46]. In this way, parameters of the policy network, $\vec{\theta}$, is updated towards the direction of higher entropy (implying exploring more of the action space).

During training, we observe a large number of unnecessary or poor explorations (*e.g.*, allocating multiple workers but 0 PS to a job) due to unawareness of job semantics. To improve exploration efficiency, we adopt another technique based on the ϵ -greedy method [55]. At each inference using the policy network, we check the input state: if the input state belongs to one of the poor states that we have identified, with probability $1 - \epsilon$, we apply the resource allocation decisions produced by the policy network, and with probability ϵ , we discard the output from the policy network, but adopt a specified action and observe the reward of this action.

The set of poor input states includes three cases: (i) there exists one job to be scheduled which has been allocated with multiple workers but no PS; (ii) there exists one job which has been allocated multiple PSs but no workers; (iii) there exists one job whose allocated numbers of workers (w) and PSs (u) differ too much, *i.e.*, $w/u > threshold$ or $u/w > threshold$ (the threshold is 10 in our experiments). Our manually specified action upon each of these input states is: (i) allocate one PS to that job; (ii) allocate one more worker to the job; (iii) allocate one more PS or one more worker to that job, to make its worker/PS numbers more even.

Experience replay. It is known that correlation among the samples prevents convergence of an actor-critic model to a good policy [55]. In our online RL, the current policy network determines the following training samples, *e.g.*, if the policy network finds that allocating more workers improves reward, then the following sample sequence will be dominated by those produced from this strategy; this may lead to a bad feedback loop, preventing the exploration of samples with higher rewards.

To alleviate correlation in the observed sample sequence, we adopt experience replay [47] in the actor-critic framework. Specifically, we maintain a replay buffer to store the samples collected in the latest time slots. At the end of each time slot, instead of using all samples collected during this time slot, we select a mini-batch of

samples from the replay buffer to compute the gradient updates, where the samples could be from multiple previous time slots.

5 Dynamic Scaling

Though node addition and deletion are supported in system design in the literature [37][29][50], existing open-source distributed machine learning frameworks (*e.g.*, TensorFlow [15], MXNet [20], Caffe [5]) do not support dynamic worker/PS adjustment in a running job. To adjust the number of workers/PSs in a job, a simple and general approach is checkpointing (*e.g.*, Optimus [49]): terminate a training job and save global model parameters as a checkpoint image; then restart the job with a new deployment of PSs and workers, and the saved model parameters. Checkpointing and restarting add additional delay to the training process [50]. For example, it takes 1 minute to checkpoint and stop training, and another 5 minutes to completely restore training of a DSSM model [52], due to data re-preprocessing before training starts. The overhead is significant when the frequency of resource scaling is high (*e.g.*, every hour). The other approach is to resize resources without terminating training process. As an example, we improve the MXNet framework [20] to enable dynamic “hot” scaling.

Challenges. In the parameter server architecture, each PS maintains a subset of the parameters in the global model. When the number of PSs changes, the global parameters need to be migrated among the PSs (for load balancing), and workers should be informed in time to send parameter updates to the correct PSs. When the number of workers changes, the new connections between new workers and the PSs should be established. The key challenges are: (1) *correctness*, *i.e.*, a consistent copy of the global model parameters should be maintained while parameters are moved across the PSs, and workers always send gradients to correct PSs; (2) *high performance*, *i.e.*, we should ensure that interruption to training is minimized and the PSs are load balanced.

Scaling Steps. We add a *coordinator* module into the MXNet framework, which works with DL² scheduler to handle joining of new workers or PSs and termination of existing ones. We demonstrate our design using the case of adding a new PS into an existing job. The steps are shown in Fig. 7.

1) *Registration.* When a new PS is launched, it registers itself with the coordinator by sending an ‘INC_SERVER’ request message. The PS will then receive its ID in the job, the global parameters it is responsible to maintain, and the current list of workers and PSs to establish connections with. After that, the PS starts functioning, awaiting workers’ parameter updates and further instructions

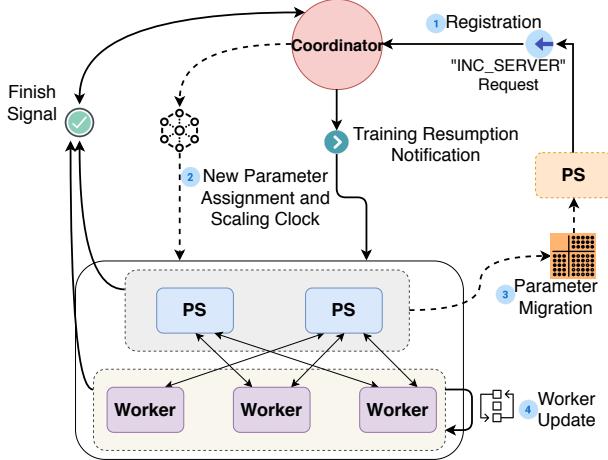


Figure 7: Steps for adding one PS into a running job

from the coordinator (*e.g.*, parameter migration).

2) *Parameter assignment.* Upon receiving a registration request, the coordinator updates its list of workers and PSs, and computes parameter assignment to the new PS. A best-fit algorithm is adopted: move part of the parameters on each existing PS to the new PS, such that all PSs maintain nearly the same number of parameters, while minimizing parameter movement across the PSs.

In order to keep a consistent copy of global model parameters when migrating parameters among PSs, we maintain a version counter for parameters. For PSs, the version counter is the number of parameter updates; for workers, the version counter is received from PSs when pulling updated parameters. To decide when PSs should migrate parameters, we calculate a scaling clock based on current version counter and round trip time between the coordinator and PSs/workers.

The coordinator sends new parameter assignment among PSs and the scaling clock to all PSs and workers.

3) *Parameter migration.* At each PS, when the version counter of parameters reaches the scaling clock received from the coordinator, the PS moves its parameters to the new PS according to the parameter assignment decisions received². Once parameter migration among all PSs is completed, the coordinator notifies all workers to resume training.

4) *Worker update.* At each worker, once its version counter equals the scaling clock received from the coordinator, the worker suspends its push/pull operations and awaits notification for completion of parameter migration. Upon notification from the coordinator, the workers update their parameter-PS mapping, establish connections with the new PS, and resume the training process.

In case of removing a PS, the scheduler chooses the

²For asynchronous training, the PS may need to buffer push or pull requests from workers and redirect them to the new PS.

PS to be removed by keeping the load balanced among the physical machines. The chosen PS sends a removal request to the coordinator. Similar steps as 2)3)4) above are then carried out, except that parameters in the removed PS are moved to other PSs, using the best-fit algorithm.

To add a new worker into an existing job, the coordinator sends the current parameter-PS mapping in the response to the worker’s registration message. It also notifies all PSs the addition of the new worker for building connections. The worker starts operation after training dataset is copied. For worker removal, the scheduler chooses the worker to be removed by keeping the load balanced across physical machines. The coordinator receives a removal request from the worker, and then broadcasts it to all workers and PSs for updating their node lists. The mini-batch size of workers is adjusted so as to keep total batch size the same.

6 Evaluation

6.1 DL² Implementation

We implement DL² as a custom scheduler on Kubernetes [9]. We run each training job using the MXNet framework [20]. Workers and PSs are running on Docker containers. Training data of jobs are stored in HDFS 2.8 [3]. The scheduler constantly queries cluster resources and job states (*e.g.*, training speeds) and instructs deployment of a new job or resource adjustment in an existing job via Kubernetes API server. Mapping the cluster and job states to a scheduling decision takes less than 3ms.

For each new job, DL² launches its coordinator, workers, and PSs on machines decided by the default placement strategy of the cluster (*i.e.*, load balancing). The coordinator is informed of the workers and PSs in the job via Kubernetes API. When a worker/PS container is launched on a machine, an agent in the container starts execution. It queries the readiness of other containers of the same job via Kubernetes API and starts user-provided training scripts after all other containers are ready. The agent also monitors the training status, *e.g.*, the number of trained steps, accuracy, and training speed.

6.2 Methodology

Testbed. Our testbed includes 13 GPU/CPU servers connected by a Dell Networking Z9100-ON 100GbE switch. Each server has one Intel E5-1660 v4 CPU, two GTX 1080Ti GPUs, 48GB RAM, one MCX413A-GCAT 50GbE NIC, one 480GB SSD, and one 4TB HDD. Each server runs Ubuntu 14.04 LTS and Docker 17.09-ce [7].

Trace. We use patterns from a 75-day real-world job trace collected from a large production DL cluster with a few thousands of GPUs and thousands of jobs, to drive

Table 1: DL Jobs in Evaluation

Model	Application domain	Dataset
ResNet-50 [30]	image classification	ImageNet [8]
VGG-16 [53]	image classification	ImageNet [8]
ResNeXt-110 [64]	image classification	CIFAR10 [2]
Inception-BN [57]	image classification	Caltech [1]
Seq2Seq [23]	machine translation	WMT17 [4]
CTC [33]	sentence classification	mr [19]
DSSM [52]	word representation	text8 [43]
WLM [14]	language modeling	PTB [12]

our testbed experiments and simulation studies. Fig. 8 (a) shows the job arrival rate (number of jobs arrived per time slot, *i.e.*, 20 minutes) during a typical week. Fig. 8 (b) shows the distribution of job duration: over a half of jobs run for more than an hour and some for days; the average job duration is 147 minutes.

Due to security and privacy concerns of the company, the job source code is not available, and we do not know job details (*e.g.*, model architecture). So we select 8 categories of ML models for experiments, from official MXNet tutorials [10], with representative application domains, different architectures and parameter sizes [10], as shown in Table 1. Each worker in different jobs uses at most 2 GPUs and 1-4 CPU cores, and each PS uses 1-4 CPU cores.

In both testbed experiments and simulations, the jobs are submitted to the cluster following the dynamic pattern in Fig. 8 (a) (with arrival rates scaled down). Upon an arrival event, we randomly select a model from Table 1 and vary its required number of training epochs (tens to hundreds) to generate a job variant, following job running time distribution of the real-world trace (scaled down). For models training on large datasets (*e.g.*, ImageNet [8]), we downscale the datasets so that the training can be finished in a reasonable amount of time. In experiments, 30 jobs are submitted to run in our testbed; in simulations, 500 servers are simulated, and 200 jobs are submitted in the simulated cluster.

Training setting. Our DL-based scheduler is implemented using TensorFlow [15]. The neural network is trained using Adam optimizer [34] with a fixed learning rate of 0.005 for offline supervised learning and 0.0001 for online reinforcement learning, mini-batch size of 256 samples, reward discount factor $\gamma = 0.9$, exploration constant $\epsilon = 0.4$, entropy weight $\beta = 0.1$, and an experience replay buffer of 8192 samples. The network has 2 hidden layers with 256 neurons each. These hyper-parameters (neural network structure, learning rate, mini-batch size, etc.) are chosen based on a few empirical training trials. We refer to one update of the neural network at the end of each time slot as one *step* in this section.

Baselines. We compare DL^2 with the following baselines.

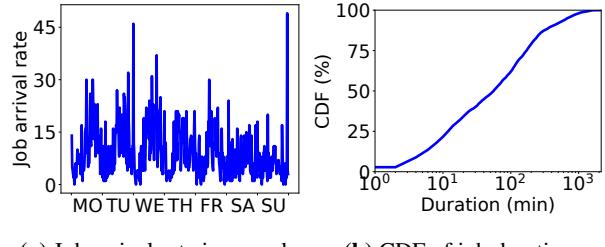


Figure 8: Trace sketch

- Dominant Resource Fairness (DRF) [24]: It allocates resources to jobs based on the fairness of dominant resources. By default, we use DRF as the existing scheduler used to guide supervised learning in DL^2 , since it is widely adopted in existing cluster schedulers, *e.g.*, YARN [58], Mesos [31].

- Tetris [27]: It preferentially allocates resources to jobs with the shortest remaining completion time and highest resource packing efficiency.

- Optimus [49]: It is a customized scheduler for DL workloads, which builds a performance model for deep learning jobs to estimate remaining training time and adopts a greedy heuristic to schedule jobs.

- OfflineRL: The offline reinforcement learning algorithm adopts pure offline training, under the same training settings as online RL in DL^2 , except that the training data are generated by an analytical performance model [49] in a simulation environment (we do not use the trace as it does not contain feedback to all decisions the offline training produces).

Wherever appropriate, we use separate training dataset and validation dataset. Both include job sequences generated using the job arrival and duration distributions from the trace. The random seeds are different when generating the datasets, to ensure that they are different.

6.3 Performance

We first compare the performance of DL^2 with baselines and show the overhead of dynamic scaling using testbed experiments.

Comparison. Fig. 9 shows that DL^2 improves average job completion time by 44.1% when compared to DRF. Tetris performs better than DRF but worse than DL^2 : once it selects a job with the highest score in terms of resource packing and remaining completion time, it always adds tasks to the job until the number of tasks reaches a user-defined threshold. When compared to Optimus, DL^2 achieves 17.5% higher performance, since Optimus' estimation of training speed is inaccurate due to cluster interference and evolved MXNet framework (*e.g.*, communication does not overlap with backward computation in Optimus' model). DL^2 also outperforms OfflineRL by

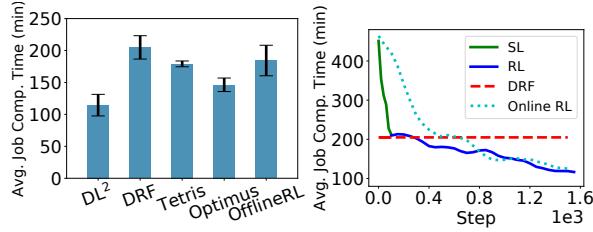


Figure 9: Performance comparison

Figure 10: Training progress comparison

37.9% due to its online training using realistic feedback.

For a better understanding of DL²'s performance gain, Fig. 10 shows how the validated performance keeps improving during the training process, when the policy NN is trained using offline supervised learning only (green curve), online RL only (cyan curve), and offline supervised learning followed by online RL (green+blue). The average job completion time shown at each time slot (*i.e.*, step) is computed over job sequence in the validation dataset, using the policy network trained (on the training dataset) at the current step. We see that with pure online RL, it takes hundreds of steps to achieve the same performance of DRF; with offline supervised learning, the performance quickly converges to a point that is close to DRF's performance within tens of steps (*i.e.*, model updates); as we continue training the NN using online RL, the performance further improves a lot. The performance of DRF is fixed as its strategy does not change over time. Besides smaller job completion time, we also observe that DL² has higher CPU and GPU utilization (similar observation as in [49]).

Scaling overhead. Fig. 11 compares the average training suspension time among all workers when checkpointing and our scaling approach are used, when different numbers of PSs are added to a ResNet-50 [30] job. The training suspension duration at a worker in DL² is from when all the received iteration counts from PSs becomes equal to the scaling clock the worker received from the coordinator, to when the worker resumes training. The checkpoint-based approach takes tens of seconds due to model saving, container relaunching and initialization before restarting training. The overhead in DL² is very small (*i.e.*, tens of milliseconds), even if the time increases linearly with the number of PSs (since we add PSs one by one). We observe similar overhead when removing PSs. On the other hand, adding or removing workers brings little interruption to existing workers, which continue with their training until adjusted training datasets are copied.

We examine detailed time cost for the 4 steps during the scaling process (§5) for adding a PS when training different models. In Fig. 12, the models are listed in increasing order of their model sizes. We observe that the scaling process spends most time in step 3 and step 4, while the

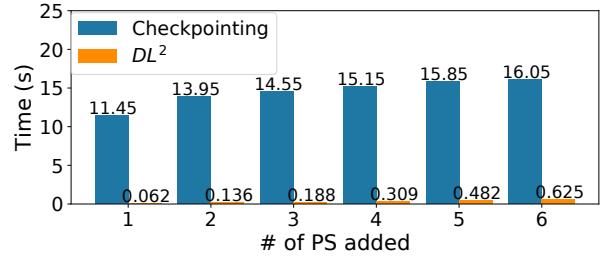


Figure 11: Scaling overhead comparison

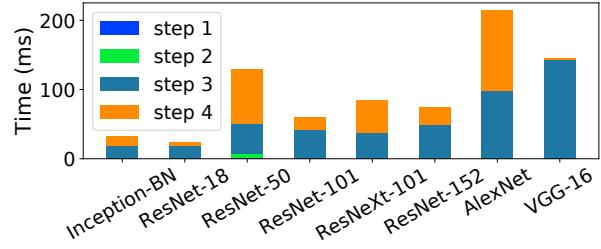


Figure 12: Timing of scaling steps

time for step 1 and step 2 is negligible. The larger a model is, the more time is spent on parameter movement (step 3). Note that only step 4 blocks worker training and is considered as overhead when compared to checkpointing. Step 3 and step 4 may happen concurrently.

In the following, we carry out controlled large-scale simulations to examine various aspects of DL² design.

6.4 Generality

Training completion time variation. To see how DL² handles practical performance variation (which white-box schedulers may not handle well), we vary the training speeds in each type of jobs to simulate variation in the training completion time of the same type of jobs (the total numbers of epochs to train remain the same). In Fig. 13, the variation indicates how the training speed deviates from the average speed (which can be faster or slower by the respective percentage). We see that Optimus is more sensitive to the variation, as it can be easily stuck in local optimum: its scheduling relies on the convexity of the performance model, but training speed variation often breaks convexity. The average job completion time shown in all simulation figures is in time slots.

Total training epoch estimation. DL² uses the total number of training epochs of jobs as input, estimated by users. The estimated total number of epochs may well be different from the actual numbers of epochs the jobs need to train to achieve model convergence. We examine how DL² performs under different estimation errors: suppose v epochs is fed into DL² as the total epoch number that a job is to train, but $v \cdot (1 + error)$ or $v \cdot (1 - error)$ is the actual number of trained epochs for the job's training

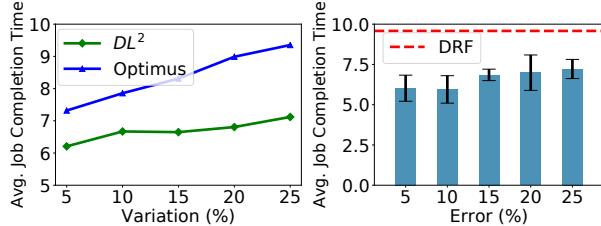


Figure 13: Training completion time variation

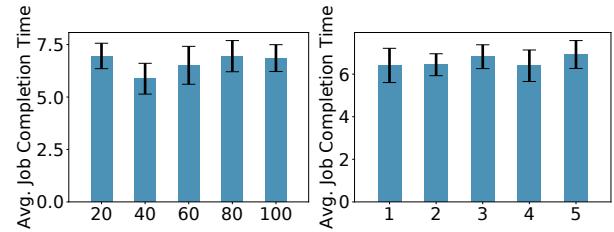


Figure 14: Performance under diff. epoch estimation errors

Figure 17: Concurrent job number

Figure 18: Varying cluster number

Table 2: Effectiveness of Training Techniques

Without	Avg. Job Completion Time	Slowdown (%)
-	5.724 ± 0.844	0
Actor-critic	6.929 ± 0.477	21.1
Exploration	7.372 ± 0.548	28.8
Experience replay	7.988 ± 0.102	39.6

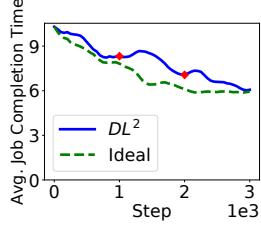


Figure 15: Handling new types of jobs

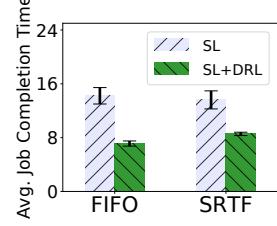


Figure 16: Different existing scheduling strategies

convergence. Fig. 14 shows that the average job completion time increases slightly when the estimation error is larger. It still outperforms DRF (which is oblivious of the estimation errors) by 28% when the error is 20%.

Unseen job types. We investigate whether DL^2 can adapt to jobs training new models. We train the neural network using the first four categories of models (Table 1) in the supervised learning phase and the first 1000 steps of the online RL phase. At step 1000 and step 2000 of the RL phase (*i.e.*, the red dots in Fig. 15), we submit jobs training two new categories of models. In the case of the “ideal” baseline, we train the NN using all categories of jobs in Table 1 from the beginning. Fig. 15 shows the average job completion time achieved using the trained NN at each time respectively, for decision making over the validation dataset. DL^2 gradually achieves the same performance as the “ideal” baseline, showing its capability to handle new types of DL jobs coming on the go.

Other scheduling strategies for supervised learning. We change the default DRF used in supervised learning of DL^2 to two other heuristics, First-In-First-Out (FIFO) and Shortest-Remaining-Time-First (SRTF). Fig. 16 shows average job performance when DL^2 uses each of these strategies in its supervised learning phase, when the NN trained only using supervised learning, or using both supervised learning and online RL, is evaluated on the validation dataset. In both cases, the performance is significantly improved with DL^2 , beyond what the existing scheduling strategy in the cluster can achieve (41.3% speedup in the case of SRTF).

Concurrent job number. We investigate how the maximal number of concurrent jobs to schedule in a time

slot, J specified in the NN input, affects the performance of DL^2 when applying the trained NN (after supervised learning and reinforcement learning) on the validation dataset. The maximal number of uncompleted jobs in all time slots is around 40; when the concurrent job number is larger than J , we schedule the jobs in batches of J jobs, according to their arrival sequence. In Fig. 17, we observe that the performance suffers when J is small, possibly because the NN is not trained on a global view when jobs are fed into the NN in batches in each time slot. Setting J to be large enough to accommodate the maximal number of concurrent jobs gives better results.

6.5 Training Design

SL loss function. We evaluate three common loss functions for supervised learning, *i.e.*, Mean Square, Cross Entropy (the default) and Absolute Difference [13]. We observe similar performance with these loss functions, while adopting Cross Entropy achieves the best performance. This is because Mean Square or Absolute Difference emphasize incorrect or suboptimal output, while only the correct or optimal output contributes to the loss when using Cross Entropy.

Reward function. We evaluate another reward function with DL^2 , which sets the reward of each action (that adds some worker/PS to a job) as the normalized number of epochs trained by the job in the time slot. We find that its performance is 29.1% worse. Our default reward function considers all jobs’ progress, enabling the policy network to learn to schedule from a global perspective.

Actor-critic. To see how the actor-critic algorithm affects training, we remove the value network but only train the policy network. As widely adopted in RL community, we use the exponential moving average of rewards as a baseline in place of the output of the value network in gradient computation of the policy network. As shown in Table 2,

with the value network, the performance is 21.1% better. This is because the average reward is not always an effective baseline; in some cases, even the optimal action leads to a lower reward than the average reward.

Job-aware exploration. We examine how exploration contributes to the performance. From Table 2, we see that without exploration the performance is 28.8% worse, as online RL is stuck in a local optimal policy.

Experience replay. We disable experience replay and see how performance changes. Table 2 shows that the average job completion time is degraded by 39.6%, indicating that experience replay is critical for training.

Federated training. Federated training enables multiple clusters to learn a global DL² model collaboratively. We study how the number of clusters affects the policy training, by implementing the A3C [46] algorithm, which trains a global policy NN using multiple DL² schedulers with different training datasets, each for one cluster. Fig. 18 shows that the global performance remains stable when we increase the number of clusters. We have also observed that with more clusters, the policy NN converges much faster due to the use of more training datasets: if there are x clusters, the NN converges almost x times faster. The preliminary result also suggests the possibility of dividing a single massive cluster into loosely coupled sub-clusters where each runs a DL² scheduler for resource allocation, if scalability issue arises.

7 Discussion and Future Directions

More scheduling features. Besides minimizing average job completion time, DL² can implement other scheduling features by adjusting the learning objective. For example, we can incorporate resource fairness by adding a quantified fairness term in the reward function.

All-reduce architecture. All-reduce architecture [61], where workers train model replicas and exchange updated model parameters directly with each other, is supported in Caffe2 [5], CNTK [6], etc. Though this paper focuses on the PS architecture, DL² can readily handle jobs using all-reduce architecture with minor modification of input state and action space of its NN, e.g., removing the elements related to PSs.

Job placement. While we use the default placement policy in this work, the placement of workers and PSs can potentially be decided by RL too. Using one NN to produce both resource allocation and placement decisions is challenging, mainly because of the significantly larger action space. RL using a hierarchical NN model [44] might be useful in making resource allocation and placement decisions in a hierarchical fashion.

Practical deployment. In practical deployment, the following two issues may need to be considered: (1) adversarial attacks that fool a neural network with malicious

input; (2) neural network monitoring that detects exceptional scheduling. These are interesting directions to explore, with progress in security research and more in-depth understanding of neural networks.

8 Related Work

Deep reinforcement learning in system research. A number of recent studies use DRL for resource allocation, device placement, and video streaming. Mao *et al.* [39] and Chen *et al.* [21] use DRL for job scheduling in cloud clusters, to minimize average job slowdown. Their NNs select the jobs (single-task jobs) to run with static resource allocation. The NNs are trained offline: multiple job arrival sequences are used as training examples; each example is repeatedly trained for multiple epochs. Mao *et al.* [41][42] learn an NN to schedule graph-based parallel jobs as in Spark, in terms of parallelism level and execution order of tasks in the jobs, using offline training. Adjustment of resources during job execution is not in the scope of the above studies.

Mirhoseini *et al.* [45][44] use DRL to optimize placement of a computation graph, to minimize running time of an individual TensorFlow job. Xu *et al.* [66] use DRL to select routing paths between network nodes for traffic engineering. Mao *et al.* [40] dynamically decide video streaming rates in an adaptive streaming system with DRL. All these studies resort to offline RL training, using data generated by analytical models or simulators. In contrast, we use offline supervised learning to prepare our NN and then online RL to further improve the NN.

ML cluster scheduling. SLAQ [67] adopts online fitting to estimate the training loss of convex algorithms, for scheduling jobs training classical ML models. Dorm [54] uses a utilization-fairness optimizer to schedule ML jobs. These work do not focus on distributed ML jobs using the parameter server architecture. Optimus [49] proposes a dynamic resource scheduler based on online-fitted resource-performance models. Bao *et al.* [18] design an online scheduling algorithm for DL jobs. These studies rely on detailed modeling of DL jobs and simplified assumptions in their design. Gandiva [63] exploits intra-job predictability to time-slice GPUs efficiently across multiple jobs, and dynamically migrate jobs to better-fit GPUs. They do not consider resource allocation adjustment; Resource allocation with GPU sharing will be an intriguing future direction to explore.

9 Conclusions

We present DL², a DL-driven scheduler for DL clusters, which expedites job completion globally with efficient

resource utilization. DL² starts from offline supervised learning, to ensure basic scheduling performance comparable to the existing cluster scheduler, and then runs in the live DL cluster to make online scheduling decisions, while improving its policy through reinforcement learning using live feedback. Our testbed experiments and large-scale trace-driven simulation verify DL²’s low scaling overhead, generality in various scenarios and outperformance over hand-crafted heuristics.

References

- [1] Caltech 256 Dataset. http://www.vision.caltech.edu/Image_Datasets/Caltech256/, 2006.
- [2] The CIFAR-10 Dataset. <https://www.cs.toronto.edu/~kriz/cifar.html>, 2009.
- [3] HDFS. <https://wiki.apache.org/hadoop/HDFS>, 2014.
- [4] WMT 2017. <http://www.statmt.org/wmt17/>, 2017.
- [5] Caffe2. <https://caffe2.ai/>, 2019.
- [6] CNTK. <https://github.com/Microsoft/CNTK>, 2019.
- [7] Docker. <https://www.docker.com/>, 2019.
- [8] ImageNet Dataset. <http://www.image-net.org>, 2019.
- [9] Kubernetes. <https://kubernetes.io>, 2019.
- [10] MXNet Official Examples. <https://github.com/apache/incubator-mxnet/tree/master/example>, 2019.
- [11] PaddlePaddle. <http://www.paddlepaddle.org/>, 2019.
- [12] Penn Tree Bank Dataset. <https://catalog.ldc.upenn.edu/ldc99t42>, 2019.
- [13] Tflearn Objectives. <http://tflearn.org/objectives/>, 2019.
- [14] Word Language Model. https://github.com/apache/incubator-mxnet/tree/master/example/gluon/word_language_model, 2019.
- [15] M. Abadi, P. Barham, J. Chen, Z. Chen, A. Davis, J. Dean, M. Devin, S. Ghemawat, G. Irving, M. Isard, et al. TensorFlow: A System for Large-Scale Machine Learning. In *Proc. of USENIX OSDI*, 2016.
- [16] D. Bahdanau, K. Cho, and Y. Bengio. Neural Machine Translation by Jointly Learning to Align and Translate. In *Proc. of ICLR*, 2015.
- [17] Y. Bao, Y. Peng, and C. Wu. Deep Learning-based Job Placement in Distributed Machine Learning Clusters. In *Proc. of IEEE INFOCOM*, 2019.
- [18] Y. Bao, Y. Peng, C. Wu, and Z. Li. Online Job Scheduling in Distributed Machine Learning Clusters. In *Proc. of IEEE INFOCOM*, 2018.
- [19] P. Bo and L. Lillian. Movie Review Data. <https://www.cs.cornell.edu/people/pabo/movie-review-data/>, 2005.
- [20] T. Chen, M. Li, Y. Li, M. Lin, N. Wang, M. Wang, T. Xiao, B. Xu, C. Zhang, and Z. Zhang. MXNet: A Flexible and Efficient Machine Learning Library for Heterogeneous Distributed Systems. In *NIPS Workshop on Machine Learning Systems (LearningSys)*, 2016.
- [21] W. Chen, Y. Xu, and X. Wu. Deep Reinforcement Learning for Multi-resource Multi-cluster Job Scheduling. In *Proc. of IEEE ICNP*, 2017.
- [22] J. Devlin, M.-W. Chang, K. Lee, and K. Toutanova. Bert: Pre-training of Deep Bidirectional Transformers for Language Understanding. *arXiv preprint arXiv:1810.04805*, 2018.
- [23] J. Gehring, M. Auli, D. Grangier, D. Yarats, and Y. N. Dauphin. Convolutional Sequence to Sequence Learning. In *Proc. of ICML*, 2017.
- [24] A. Ghodsi, M. Zaharia, B. Hindman, A. Konwinski, S. Shenker, and I. Stoica. Dominant Resource Fairness: Fair Allocation of Multiple Resource Types. In *Proc. of USENIX NSDI*, 2011.
- [25] I. Goodfellow, Y. Bengio, A. Courville, and Y. Bengio. *Deep Learning*. MIT press Cambridge, 2016.
- [26] P. Goyal, P. Dollár, R. Girshick, P. Noordhuis, L. Wesolowski, A. Kyrola, A. Tulloch, Y. Jia, and K. He. Accurate, Large Minibatch SGD: Training Imagenet in 1 Hour. *arXiv preprint arXiv:1706.02677*, 2017.
- [27] R. Grandl, G. Ananthanarayanan, S. Kandula, S. Rao, and A. Akella. Multi-Resource Packing for Cluster Schedulers. In *Proc. of ACM SIGCOMM*, 2014.
- [28] A. Graves, A.-r. Mohamed, and G. Hinton. Speech Recognition with Deep Recurrent Neural Networks. In *Proc. of IEEE ICASSP*, 2013.

- [29] A. Harlap, A. Tumanov, A. Chung, G. R. Ganger, and P. B. Gibbons. Proteus: Agile ML elasticity through tiered reliability in dynamic resource markets. In *Proc. of ACM EuroSys*, 2017.
- [30] K. He, X. Zhang, S. Ren, and J. Sun. Deep Residual Learning for Image Recognition. In *Proc. of IEEE CVPR*, 2016.
- [31] B. Hindman, A. Konwinski, M. Zaharia, A. Ghodsi, A. D. Joseph, R. H. Katz, S. Shenker, and I. Stoica. Mesos: A Platform for Fine-Grained Resource Sharing in the Data Center. In *Proc. of USENIX NSDI*, 2011.
- [32] J. Jiang, L. Yu, J. Jiang, Y. Liu, and B. Cui. Angel: a New Large-Scale Machine Learning System. *National Science Review*, 2017.
- [33] Y. Kim. Convolutional Neural Networks for Sentence Classification. In *Proc. of SIGDAT EMNLP*, 2014.
- [34] D. P. Kingma and J. Ba. Adam: A Method for Stochastic Optimization. In *Proc. of ICLR*, 2015.
- [35] A. Krizhevsky, I. Sutskever, and G. E. Hinton. ImageNet Classification with Deep Convolutional Neural Networks. In *Proc. of NIPS*, 2012.
- [36] J. Langford and T. Zhang. The Epoch-Greedy Algorithm for Contextual Multi-Armed Bbandits. In *Proc. of NIPS*, 2007.
- [37] M. Li, D. G. Andersen, J. W. Park, A. J. Smola, A. Ahmed, V. Josifovski, J. Long, E. J. Shekita, and B.-Y. Su. Scaling Distributed Machine Learning with the Parameter Server. In *Proc. of USENIX OSDI*, 2014.
- [38] S. Mannor, D. Peleg, and R. Rubinstein. The Cross Entropy Method for Classification. In *Proc. of ICML*, 2005.
- [39] H. Mao, M. Alizadeh, I. Menache, and S. Kanadula. Resource Management with Deep Reinforcement Learning. In *Proc. of ACM HotNets*, 2016.
- [40] H. Mao, R. Netravali, and M. Alizadeh. Neural Adaptive Video Streaming with Pensieve. In *Proc. of ACM SIGCOMM*, 2017.
- [41] H. Mao, M. Schwarzkopf, S. Venkatakrishnan, and M. Alizadeh. Learning Graph-based Cluster Scheduling Algorithms. In *Proc. of SysML*, 2018.
- [42] H. Mao, M. Schwarzkopf, S. B. Venkatakrishnan, Z. Meng, and M. Alizadeh. Learning Scheduling Algorithms for Data Processing Clusters. In *Proc. of ACM SIGCOMM*, 2019.
- [43] M. Matt. text8. <http://mattmahoney.net/dc/>, 2019.
- [44] A. Mirhoseini, A. Goldie, H. Pham, B. Steiner, Q. V. Le, and J. Dean. A Hierarchical Model for Device Placement. In *Proc. of ICLR*, 2018.
- [45] A. Mirhoseini, H. Pham, Q. V. Le, B. Steiner, R. Larsen, Y. Zhou, N. Kumar, M. Norouzi, S. Bengio, and J. Dean. Device Placement Optimization with Reinforcement Learning. In *Proc. of ICML*, 2017.
- [46] V. Mnih, A. P. Badia, M. Mirza, A. Graves, T. Lillicrap, T. Harley, D. Silver, and K. Kavukcuoglu. Asynchronous Methods for Deep Reinforcement Learning. In *Proc. of ICML*, 2016.
- [47] V. Mnih, K. Kavukcuoglu, D. Silver, A. A. Rusu, J. Veness, M. G. Bellemare, A. Graves, M. Riedmiller, A. K. Fidjeland, G. Ostrovski, et al. Human-level Control Through Deep Reinforcement Learning. *Nature*, 2015.
- [48] V. Nair and G. E. Hinton. Rectified Linear Units Improve Restricted Boltzmann Machines. In *Proc. of ICML*, 2010.
- [49] Y. Peng, Y. Bao, Y. Chen, C. Wu, and C. Guo. Optimus: an Efficient Dynamic Resource Scheduler for Deep Learning Clusters. In *Proc. of ACM EuroSys*, 2018.
- [50] A. Qiao, A. Aghayev, W. Yu, H. Chen, Q. Ho, G. A. Gibson, and E. P. Xing. Litz: Elastic Framework for High-Performance Distributed Machine Learning. In *Proc. of USENIX ATC*, 2018.
- [51] M. Randles, D. Lamb, and A. Taleb-Bendiab. A Comparative Study into Distributed Load Balancing Algorithms for Cloud Computing. In *Proc. of IEEE WA/WA*, 2010.
- [52] Y. Shen, X. He, J. Gao, L. Deng, and G. Mesnil. A Latent Semantic Model with Convolutional-Pooling Structure for Information Retrieval. In *Proc. of ACM CIKM*, 2014.
- [53] K. Simonyan and A. Zisserman. Very Deep Convolutional Networks for Large-scale Image Recognition. In *Proc. of ICLR*, 2015.
- [54] P. Sun, Y. Wen, N. B. D. Ta, and S. Yan. Towards Distributed Machine Learning in Shared Clusters: A Dynamically-Partitioned Approach. In *Proc. of IEEE Smart Computing*, 2017.
- [55] R. S. Sutton and A. G. Barto. *Reinforcement Learning: An Introduction*. MIT Press Cambridge, 1998.

- [56] R. S. Sutton, D. A. McAllester, S. P. Singh, and Y. Mansour. Policy Gradient Methods for Reinforcement Learning with Function Approximation. In *Proc. of NIPS*, 2000.
- [57] C. Szegedy, V. Vanhoucke, S. Ioffe, J. Shlens, and Z. Wojna. Rethinking the Inception Architecture for Computer Vision. In *Proc. of IEEE CVPR*, 2016.
- [58] V. K. Vavilapalli, A. C. Murthy, C. Douglas, S. Agarwal, M. Konar, R. Evans, T. Graves, J. Lowe, H. Shah, S. Seth, et al. Apache Hadoop YARN: Yet Another Resource Negotiator. In *Proc. of ACM SoCC*, 2013.
- [59] A. Verma, L. Pedrosa, M. Korupolu, D. Oppenheimer, E. Tune, and J. Wilkes. Large-Scale Cluster Management at Google with Borg. In *Proc. of ACM EuroSys*, 2015.
- [60] O. Vinyals, T. Ewalds, S. Bartunov, P. Georgiev, A. S. Vezhnevets, M. Yeo, A. Makhzani, H. Küttler, J. Agapiou, J. Schrittwieser, et al. StarCraft II: A New Challenge for Reinforcement Learning. *arXiv preprint arXiv:1708.04782*, 2017.
- [61] P. Watcharapichat, V. L. Morales, R. C. Fernandez, and P. Pietzuch. Ako: Decentralised Deep Learning with Partial Gradient Exchange. In *Proc. of ACM SoCC*, 2016.
- [62] R. J. Williams. Simple Statistical Gradient-following Algorithms for Connectionist Reinforcement Learning. *Machine learning*, 1992.
- [63] W. Xiao, R. Bhardwaj, R. Ramjee, M. Sivathanu, N. Kwatra, Z. Han, P. Patel, X. Peng, H. Zhao, Q. Zhang, F. Yang, and L. Zhou. Gandiva: Introspective Cluster Scheduling for Deep Learning. In *Proc. of USENIX OSDI*, 2018.
- [64] S. Xie, R. Girshick, P. Dollár, Z. Tu, and K. He. Aggregated Residual Transformations for Deep Neural Networks. In *Proc. of IEEE CVPR*, 2017.
- [65] E. P. Xing, Q. Ho, W. Dai, J.-K. Kim, J. Wei, S. Lee, X. Zheng, P. Xie, A. Kumar, and Y. Yu. Petuum: A New Platform for Distributed Machine Learning on Big Data. In *Proc. of ACM SIGKDD*, 2015.
- [66] Z. Xu, J. Tang, J. Meng, W. Zhang, Y. Wang, C. H. Liu, and D. Yang. Experience-driven Networking: A Deep Reinforcement Learning based Approach. In *Proc. of IEEE INFOCOM*, 2018.
- [67] H. Zhang, L. Stafman, A. Or, and M. J. Freedman. SLAQ: Quality-Driven Scheduling for Distributed Machine Learning. In *Proc. of SoCC*, 2017.