

Dynamic Stale Synchronous Parallel Distributed Training for Deep Learning

Xing Zhao*, Aijun An*, Junfeng Liu[†], Bao Xin Chen*

*Department of Electrical Engineering and Computer Science

York University, Toronto, Canada

{xingzhao, aan, baoxchen}@cse.yorku.ca

[†]Platform Computing

IBM Canada, Markham, Canada

jfliu@ca.ibm.com

Abstract—Deep learning is a popular machine learning technique and has been applied to many real-world problems, ranging from computer vision to natural language processing. However, training a deep neural network is very time-consuming, especially on big data. It has become difficult for a single machine to train a large model over large datasets. A popular solution is to distribute and parallelize the training process across multiple machines using the parameter server framework. In this paper, we present a distributed paradigm on the parameter server framework called *Dynamic Stale Synchronous Parallel* (DSSP) which improves the state-of-the-art *Stale Synchronous Parallel* (SSP) paradigm by dynamically determining the staleness threshold at the run time. Conventionally to run distributed training in SSP, the user needs to specify a particular staleness threshold as a hyper-parameter. However, a user does not usually know how to set the threshold and thus often finds a threshold value through trial and error, which is time-consuming. Based on workers' recent processing time, our approach DSSP adaptively adjusts the threshold per iteration at running time to reduce the waiting time of faster workers for synchronization of the globally shared parameters (the weights of the model), and consequently increases the frequency of parameters updates (increases iteration throughput), which speeds up the convergence rate. We compare DSSP with other paradigms such as *Bulk Synchronous Parallel* (BSP), *Asynchronous Parallel* (ASP), and SSP by running deep neural networks (DNN) models over GPU clusters in both homogeneous and heterogeneous environments. The results show that in a heterogeneous environment where the cluster consists of mixed models of GPUs, DSSP converges to a higher accuracy much earlier than SSP and BSP and performs similarly to ASP. In a homogeneous distributed cluster, DSSP has more stable and slightly better performance than SSP and ASP, and converges much faster than BSP.

Keywords-distributed deep learning, parameter server, BSP, ASP, SSP, GPU cluster.

I. INTRODUCTION

The parameter server framework [1] [2] has been developed to support distributed training of large-scale machine learning (ML) models (such as deep neural networks [3] [4] [5]) over very large data sets, such as Microsoft COCO [6], ImageNet 1K [3] and ImageNet 22K [7]. Training a deep model using a large-scale cluster with an efficient distributed paradigm reduces the training time from weeks on a single server to days or hours.

Since the DistBelief [1] framework was developed in 2012, distributed machine learning has attracted the attention of many ML researchers and system engineers. In 2014, the Parameter Server architecture [8] was launched. Its coarse-grained parallelism shows a significant speedup in convergence over 6000 servers. In 2015, a fine-grained parallel and distributed system Petuum [9] was developed to support customized distributed training for particular machine learning algorithms instead of providing a general distributed framework to many machine learning algorithms based on, e.g., Hadoop [10] and Spark [11]. By then a global competition has begun on developing efficient distributed machine learning training platforms. Baidu published the distributed training system PaddlePaddle [12] in 2015 for deep learning, which inherits the parameter server framework. Alibaba released KunPeng [13] in 2017, an variation of the parameter server, which was claimed as an universal distributed platform. Due to its efficient scalable network communication design, the parameter server framework can be found in most distributed platforms in practice regardless whether they are implemented in fine-grained or coarse-grained parallelism.

In a nutshell, the **parameter server framework** consists of a logic *server* and many *workers*. Workers are all connected to the server. The server usually maintains the *globally shared weights* by aggregating weight updates from the workers and updating the global weights. It provides a central storage for the workers to upload their computed updates (by the `push` operation) and fetch the up-to-date global weights (by the `pull` operation). The parameter server framework supports *model parallelism* and *data parallelism* [14] [15]. In model parallelism, a ML model can be partitioned and its components are assigned to a server group (i.e., a distributed logic server) and workers. A worker computes the gradients for a server based on its assigned model partition and data. However, it is difficult to decouple a model due to dependencies between components of the model (e.g., layers of DNN) and the nature of the optimization method (e.g., stochastic gradient descent) [16]. Thus, model parallelism is not commonly seen in practice. In this work, we focus on data parallelism, in which the

training data is partitioned based on the number of workers and a partition is assigned to each worker. A worker machine usually contains a replica of the ML model and is assigned an equal-sized partition of the entire training data. Each worker iterates the following steps: ① computing the gradients based on a sample or a mini-batch and its local parameters (e.g., weights), ② sending the gradients as an update to the server, ③ retrieving the latest global weights from the server and ④ assigning the retrieved weights as its local weights. We call the time span accumulated from ① to ④ the *iteration interval*. From the server’s perspective, an iteration interval of a worker is the time period between two consecutive updates it receives from the worker.

A. Distributed paradigms for updating the parameters

There are three paradigms for updating the model parameters during distributed deep learning with data parallelism. They are Bulk Synchronous Parallel (BSP) [17], Stale Synchronous Parallel (SSP) [18] [19] and Asynchronous Parallel (ASP) [1] [20] [7].

1) *Bulk Synchronous Parallel (BSP)*: In BSP, all workers send their computed gradients to the parameter server for the global weight update and wait for each other at the end of every iteration for synchronization. Once the parameter server receives gradients from all workers and updates the global weights, it sends the latest global weights (parameters) to the workers before each worker starts a new iteration. In this paradigm, every worker starts a new iteration based on the same version of the global weights from the server, that is, the weights are *consistent* among all workers. BSP generally achieves the best accuracy among the three paradigms but takes the most training time due to its frequent synchronization among workers. Synchronization incurs waiting time for faster workers. Since ML has the fault tolerant property [16] (that is, it is robust against minor errors in intermediate calculations) when it uses the iterative convergent optimization algorithm such as stochastic gradient descent (SGD) [21], a more flexible paradigm that uses less synchronization can be applied.

2) *Asynchronous Parallel (ASP)*: In ASP, all workers send their computed gradients to the parameter server at each iteration but no synchronization is required. Workers do not wait for each other and simply run independently during the entire training. In this case, some slower workers will bring the delayed or staled gradient updates to the globally shared weights on the parameter server. The delayed gradients introduce errors to the iterative convergent method. Consequently, it prolongs the convergent rate of a training model and even diverges the learning of the model when the staled updates are from very old iterations. Without any synchronization, each worker may obtain a different version of the global weight from the parameter server at the end of its iteration. From the system’s point of view, the global weights are *inconsistent* to all workers at the beginning of

their iterations. In contrast to BSP, ASP has the least training time for a given number of epochs, but usually yields a much lower accuracy due to the missing synchronization step among workers. Nonetheless, zero or less synchronization for workers usually diffuses the convergence of a DNN model [22]. Hence, ASP is not stable in terms of the model convergence.

3) *Stale Synchronous Parallel (SSP)*: SSP is a combination of BSP and ASP. It relies on a policy to switch between ASP and BSP dynamically during the training. The policy is to restrict the number of iterations between the fastest worker(s) and the slowest worker(s) to not exceed more than a user-specified staleness threshold s where s is a natural number. This policy ensures that the difference in the number of iterations among workers is no larger than s . Hence, as long as the policy is not violated, there is no waiting time among workers. When the threshold is exceeded, the synchronization is forced on fastest workers which are s iterations ahead of the slowest worker. One effective implementation of SSP from [18] only forces the fastest worker(s) to wait for the slowest worker(s) and allows the rest continue their iterations. In this distribution paradigm, the parameters (the global weights) of an ML model are considered *inconsistent* [23] among the workers at the beginning of their iterations when the policy is not violated. However, the inconsistency is limited to a certain extent by the *small* threshold s so that the ML model can still converge (close to an optimum) [24].

B. Contributions

SSP is an intermediate solution between BSP and ASP. It is faster than BSP, and guarantees convergence, leading to a more accurate model than ASP. However, in SSP the user-specified staleness threshold is fixed, which leads to two problems. First, it is usually hard for the user to specify a good single threshold since user has no knowledge which value is the best. Choosing a good threshold may involve manually searching in an integer range via numerous trials. Also, a DNN model involves many other hyperparameters (such as the number of layers and the number of nodes in each layer). When these parameters change, the same searching trials have to be repeated again to fine-tuning the staleness threshold. Second, a single fixed value may not be suitable for the whole training process. An ill-specified value may cause the fastest workers to wait for longer time than necessary. For example, Figure 2 shows that if the threshold is exceeded at the red solid line, the waiting time for the fastest worker if it starts waiting right away is more than the waiting time if it continues but starts waiting at the green solid line. In fact, the waiting time for it to start waiting at the yellow solid line is the minimum of the three. However, we have to make sure that the difference in iterations between the fastest and slowest workers is not too large. Otherwise, too many staled updates may delay the

convergence of the ML model [25].

To solve these problems, we propose an adaptive SSP scheme named *Dynamic Stale Synchronous Parallel (DSSP)*. Our approach dynamically selects a threshold value from a given range in the training process based on the statistics of the real-time processing speed of distributed computing resources. It allows the threshold to change over time and also allows different workers to have different threshold values, adapting to the run-time environment. To achieve this purpose, we design a *synchronization controller* at the server side to determine how many iterations the current fastest worker should continue running at the end of its iteration upon the excess of the lower bound of a user-specified staleness threshold range. The decision is made at run time by estimating the future waiting times of workers based on the timestamps of their previous *push requests* and selecting a time point in the range that leads to the least estimated waiting time. In this way, we enable the parameter server to dynamically adjust or relax the synchronization of workers during the training based on the run-time environment of the distributed system.

In addition, although experiments have been reported on parameter servers with a variety of ML models, experiments of comparing DNN models under different distributed paradigms are rarely seen in the literature. In this paper, we look into four distributed paradigms (i.e., BSP, ASP, SSP and DSSP) and compare their performance by training three DNN models on two image datasets using MXNet [2] which provides BSP and ASP. We implemented both SSP and DSSP in MXNet, report and analyze our findings from the experiments.

II. RELATED WORK

There are variety of approaches to optimizing the distributed paradigms under the parameter server framework. Generally, they can be categorized into three basis streams: BSP, ASP and SSP. Chen et al. [26] try to optimize the BSP by adding few backup workers. That is to train N workers in BSP, they add c backup workers so that there are $N + c$ workers during the training. By the end of each iteration, the server only takes the first N arrived updates and drops the c slower arrived updates from the stragglers for weight synchronization. In this case, the training data allocated to the c random slower workers are partially wasted in each iteration.

Hadjis et al. [27] optimize ASP from machine learning’s perspective. It adjusts the momentum based on the degree of asynchrony (staleness of the gradients). Then, it uses the tuned momentum to mitigate the divergent direction that staled gradients introduce. Meanwhile, model parallel computing is applied here for better performance where a DNN model is split into two parts: convolutional layers and fully connected layers. Both parts are computed parallelly and concurrently.

Zhang and Kwok [28] propose to use asynchronous distribution to optimize synchronous ADMM algorithm. However, it uses *partial barrier* to make the fastest workers wait for the slowest workers and *bounded delay* to guarantee that the iterations among workers do not exceed a user specified hyperparameter τ which is equivalent to the staleness threshold s of SSP in [18]. Then, all these make the approach rather close to SSP than ASP optimization. Bounded delay also appears in [25] and is elaborated in the rest of this section.

Li et al. [25] introduce *bounded delay* which is similar to SSP except that it takes all workers’ iterations into account instead of letting each worker count its own iteration. In order to keep the ML model (global weights) consistent, iterations in sequence are allowed to run concurrently in parallel under the dependency restriction. Iteration t depends on iteration $t - k$ if iteration t requires the result of iteration $t - k$ in order to proceed. In the bounded delay approach, the number of bounded iterations k is specified by the user, similar to the staleness threshold s in SSP. k means that for a continuous k iterations, every iteration is independent of each other, and they can run concurrently in parallel without waiting for each other. When k is exceeded, the fastest iteration t has to wait for the slower iterations behind $t - k$. Imagine iterations are pre-assigned to workers as tasks, then the bounded delay is equivalent to SSP. For example, we have iterations $\{I_1, I_2, I_3, I_4, I_5, I_6\}$ and two workers P_1 and P_2 . Each iteration I_i completes a min-batch of samples. P_1 receives $\{I_1, I_3, I_5\}$, P_2 receives $\{I_2, I_4, I_6\}$. If $k = 3$, then I_4 depends on I_1 , I_5 depends on I_2 , I_6 depends on I_3 . So P_2 at I_4 has to wait for P_1 finishes I_1 . P_1 at I_5 has to wait for P_2 completes I_2 . Bounded delay is rather an inflexible scheme by pre-scheduling tasks to workers. Although the authors briefly claim that more consistent paradigms can be developed based on the dependency constraint, no further exploration is provided in their paper. Our work extends this direction and presents a flexible scheduling approach in which k is dynamically assigned at the training time. In our dynamic bounded delay paradigm, every optimal bound k yields the least waiting time for workers by optimizing the synchronization frequency. For the continuous k iterations in which every iteration is independent is adjusted dynamically in the running time to reduce waiting time of coming iterations which depend on the earlier ones.

III. DYNAMIC STALE SYNCHRONOUS PARALLEL

In this section, we propose the *Dynamic Stale Synchronous Parallel (DSSP)* synchronization method. Instead of using a single and definite staleness threshold as in SSP, DSSP takes a range for the staleness threshold as input, and dynamically determines an optimal value for the staleness threshold during the run time. The value for the threshold can change over time and adapt to the run-time environment.

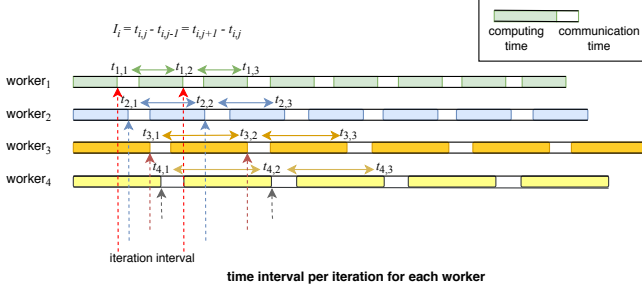


Figure 1: Iteration *intervals* measured by timestamps of push requests from workers. A dotted line represents the time for a push request from a worker to the server. An *interval* consists of *communication period* (blank block) and *gradient computation period* (solid block).

A. Problem statement

Given a *lower bound* and an *upper bound* of staleness thresholds s_L and s_U , DSSP finds an *optimal* threshold $s^* \in [s_L, s_U]$ for a worker dynamically, which yields the *minimum waiting time* for the worker to *synchronize* with others, based on the *iteration time* collected from each worker at the run time.

In other words, DSSP finds an integer $r^* \in R$, where $R = [0, r_{max}]$, $r_{max} = s_U - s_L$ and $r^* = s^* - s_L$. That is, DSSP finds an *optimal spot* on the time line R to bound the workers for synchronization. Empirically we can find a $r = s - s_L$ that is close to r^* , which is the same as finding $s \in [s_L, s_U]$ close to s^* .

B. Assumption

An *iteration interval* of a worker is the time period between two *consecutive* updates (i.e., push requests) the server receives from the worker. We can measure the length of an iteration interval of a worker by using the timestamps of the push requests sent by the worker (see Figure 1).

We assume that the iteration intervals of a worker in continuous iterations in a short time period are very similar. That is, for contiguous iterations of a worker in a short time period, each iteration has the similar processing time which includes computing gradients over a mini-batch, sending gradients to and receiving updated weights (parameters) from the parameter server. Therefore, if we use the *most recent intervals* to estimate the length of next intervals, the error is small under this assumption. Even if the network experiences some instabilities in a short period and we may make some wrong predictions which may lead to some latent updates, DSSP can still converge due to the error tolerance of an iterative-convergent method such as Parallelized SGD [21].

C. Method

The proposed DSSP method is described in Algorithm 1. The algorithm contains two parts: one for workers and

Algorithm 1 Dynamic Staled Gradient Method

Worker p at iteration t_p

- 1: Wait until receiving OK from Server
- 2: pull weights w^s from Sever
- 3: Replace local weights w^{t_p} with w^s
- 4: Gradient $g^{t_p} \leftarrow \frac{1}{m} \sum_{i=1}^m \partial l_{loss}((x_i, y_i), w^{t_p})$
 $\{m: \text{size of mini-batch } M \text{ and } (x_i, y_i) \in M\}$
- 5: push g^{t_p} to Server

Server at iteration t_s

- Upon receiving push request with g^{t_p} from worker p ;
 - r_p stores the number of extra iterations worker p is allowed beyond s_L , initialized to zero at the very beginning;
 - t_i stores the number of push requests received from worker i so far
- 1: $t_p = t_p + 1$
 - 2: Update the server weights w^{t_s} with g^{t_p} . If some other workers send their updates at the same time, their gradients are *aggregated* before updating w^{t_s}
 - 3: **if** ($r_p > 0$) **then**
 - 4: $r_p = r_p - 1$
 - 5: Send OK to worker p
 - 6: **else**
 - 7: Find the *slowest* and *fastest* workers based on array t
 - 8: **if** ($t_p - t_{slowest} \leq s_L$) **then**
 - 9: Send OK to worker p
 - 10: **else**
 - 11: **if** t_p is the *fastest* worker **then**
 - 12: $r_p \leftarrow \text{synchronization_controller}(clock_p^{push}, r_p)$
 $\{clock_p^{push}: \text{timestamp of push request from worker } p\}$
 - 13: **if** ($r_p > 0$) **then**
 - 14: Send OK to worker p
 - 15: **end if**
 - 16: **end if**
 - 17: Wait until the *slowest* worker sends the next push request(s) so that $t_p - t_{slowest} \leq s_L$. After updating the server weights w^{t_s} with (*aggregated*) gradients, send OK to worker p
 - 18: **end if**
 - 19: **end if**
-

the other for the server. Each worker is assigned a partition of the training data, and computes parameter updates (i.e., gradients) iteratively with the partition using the stochastic gradient descent (SGD) method. In each iteration, a mini-batch of the partition is used to compute the gradients based on the current local weights. The worker then sends the gradients to the server through a *push request* and waits for the server to send back the OK signal. After receiving OK , the worker pulls the weights from the server and replaces its

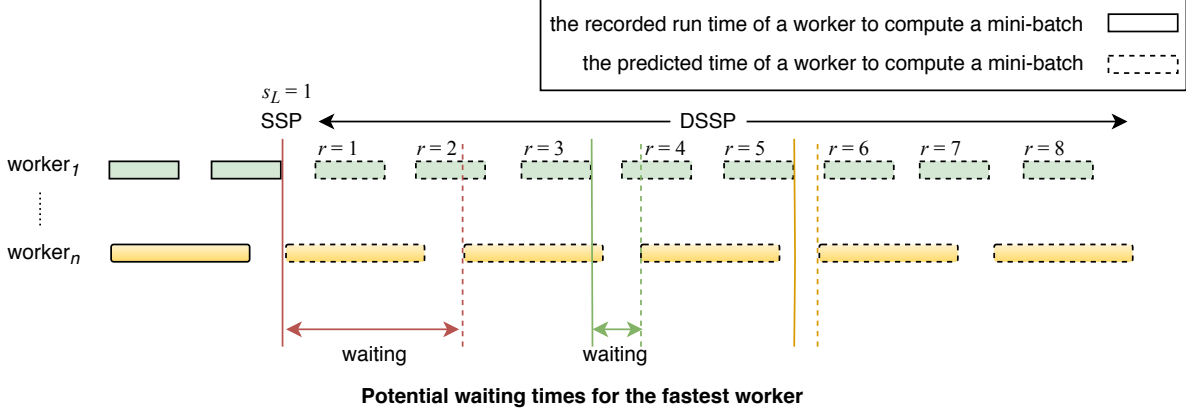


Figure 2: Prediction module finding the *least waiting time* for the *fastest worker* via iteration time intervals of workers. A solid line represents a boundary to stop the fastest workers continuing new iterations for *synchronization* and a dash line represents the end of waiting when the slowest worker completes its running iteration. The solid line is drawn upon a fastest worker sends a `push` request to the server and waits for the `OK` signal from the server. Once `OK` is received, it pulls the new updated weight from the server and starts a new iteration where the dash line is drawn. The dash line also indicates the time that the *slowest worker* receives a new updated weight via `pull` request and starts a new iteration. Worker₁ is the *fastest worker* and the worker_n is the *slowest worker*. The colored block represents one iteration time. Following SSP, worker₁ has to stop at the red solid line. DSSP compares each r value and finds the r^* which gives the least waiting time. Here, $r^* = 3$ if $r \in R = [0, 4]$. DSSP allows worker₁ to run 3 more iterations and stop at the green solid line.

local weights with the global weights from the server. The training at the worker continues with the next mini-batch of the data partition based on the new weights.

On the server side, once the server receives a *push request* from a worker p , it updates its weights with the gradients from worker p . It then determines whether to allow worker p to continue. If yes, it will send worker p an `OK` signal; otherwise, it postpones sending the `OK` signal until the slowest worker catches up.

To determine whether to allow a worker p to continue, the server stores *the number of push requests* received from each worker and finds the *slowest worker*. If the number of push requests of worker p is no more than s_L iterations away from the slowest worker, the server allows worker p to continue by sending `OK` to worker p (Lines 7-9 in Algorithm 1). Otherwise, if worker p is currently the *fastest worker*¹, the server calls the *synchronization_controller* procedure to determine whether it allows worker p to continue with extra iterations.

Algorithm 2 describes the *synchronization_controller* procedure. It stores in table \mathcal{A} the *timestamps* of the *two latest push requests* from all workers, and uses the information in \mathcal{A} to simulate the next r_{max} iterations of worker p and the slowest worker, where r_{max} is the maximum number of extra iterations allowed for a worker to be ahead of the slowest worker beyond the lower bound of the staleness threshold. With the simulated timestamps, it finds a time point r^* in the

¹The reason we call the procedure only for the current fastest worker is to save the server's computation time.

range of $[0, r_{max}]$ that minimizes the simulated waiting time of worker p (Line 8 in Algorithm 2). The value r^* is returned to the caller (the Server part of Algorithm 1) and stored as r_p for worker p . For example, in Figure 2, suppose worker n is the slowest worker and we are currently processing worker 1 (i.e., $p = 1$). The green boundary yields the least waiting time for worker 1. Then worker 1 should continue running 3 more iterations once s_L is exceeded. In this case, 3 is the r^* returned to the server procedure of Algorithm 1. If 0 is returned, it indicates that the current iteration yields the least waiting time, and the worker should wait for the slowest worker at the current iteration.

In future iterations when worker p sends a push request, if $r_p > 0$, the server sends the `OK` signal right after updating the global weights with the gradients sent by the worker and decreases r_p by 1. In this way, even if worker p is not the fastest worker in that iteration of the server, as long as its $r_p > 0$ (due to being the fastest worker in a previous iteration), it can still perform extra iterations beyond s_L . Thus, our method is flexible in that different workers may have different thresholds, and also the threshold for a worker can change over time, depending on the run-time environment.

IV. THEORETICAL ANALYSIS

In this section, we prove the convergence of SGD under DSSP by showing that DSSP shares the same *regret bound* $O(\sqrt{T})$ as SSP from [18]. That is, SGD converges in expectation when the number of iterations T is large under DSSP. We first present the theorem of SSP. Based on

Algorithm 2 synchronization_controller

Input: $push_p^t$: timestamp of push request of worker p for sending its iteration t 's update to the server

Output: r^* , number of extra iterations that worker p is allowed to run

{Table \mathcal{A} stores the timestamps of two latest push requests by all workers, where $\mathcal{A}[i][0]$ stores the timestamp of the latest push request by worker i and $\mathcal{A}[i][1]$ stores the timestamp of the second latest push request by worker i }

- 1: $\mathcal{A}[p][1] \leftarrow \mathcal{A}[p][0]$
 - 2: $\mathcal{A}[p][0] \leftarrow push_p^t$
 - 3: Find the *slowest worker* from table \mathcal{A}
 - 4: Compute the length of the latest iteration interval of worker p : $\mathcal{I}_p \leftarrow \mathcal{A}[p][0] - \mathcal{A}[p][1]$
 - 5: Compute the length of the latest iteration interval of the *slowest worker*: $\mathcal{I}_{slowest} \leftarrow \mathcal{A}[slowest][0] - \mathcal{A}[slowest][1]$
 - 6: Simulate next r_{max} iterations for worker p based on \mathcal{I}_p and $\mathcal{A}[p][0]$ by storing the r_{max} simulated timestamps in array Sim_p so that:
 $Sim_p[0] \leftarrow \mathcal{A}[p][0]$
 $Sim_p[i] \leftarrow Sim_p[0] + i \times \mathcal{I}_p$ where $0 < i \leq r_{max}$
{ r_{max} : the maximum extra iterations allowed}
 - 7: Repeat the above step for the *slowest worker* and store the r_{max} simulated timestamps in array $Sim_{slowest}$ with $Sim_{slowest}[0] \leftarrow \mathcal{A}[slowest][0] + \mathcal{I}_{slowest}$
 - 8: Find the simulated time point r^* for the index of $Sim_p[r]$ that minimizes $|Sim_{slowest}[k] - Sim_p[r]|$ for all $k \in [0, r_{max}]$ and $r \in [0, r_{max}]$
 - 9: **return** r^*
-

the theorem, we show that DSSP has a bound on regret. Therefore, DSSP supports SGD convergence following the same conditions as SSP.

Theorem 1 (adapted from [18]. SGD under SSP):

Suppose function $f(w) := \sum_{t=1}^T f_t(w)$ is a convex function and $\forall f_t(w)$ is also convex. We use iterative convergent optimization algorithm (gradient descent) on one component ∇f_t at a time to search for the minimizer w^* under SSP with the staleness threshold s and P workers. Let $v_t := -\eta_t \nabla f_t(\tilde{w}_t)$ where $\eta_t = \frac{\sigma}{\sqrt{t}}$ and $\sigma = \frac{F}{L\sqrt{2(s+1)P}}$. Here \tilde{w}_t represents the noisy state of the globally shared weight. F and L are constants. Assume f_t are L -Lipschitz with constant L and the distance $D(w||w')$ between two multidimensional points w and w' is bounded such that $D(w||w') := \frac{1}{2}\|w - w'\|_2^2 \leq F^2$ where F is constant. We have a bound on the regret

$$R[X] := \sum_{t=1}^T f_t(\tilde{w}_t) - f(w^*) \leq 4FL\sqrt{2(s+1)PT} \quad (1)$$

Thus, $R[X] = O(\sqrt{T})$ since $\lim_{T \rightarrow \infty} \frac{R[X]}{T} = 0$

Theorem 2 (SGD under DSSP): Following all conditions and assumptions from Theorem 1, we add a new term $R = [0, s_U - s_L]$, the *range* of the staleness threshold. Let $r \in R$ and $r \geq \forall r' \in R$. We have a bound on the regret

$$R[X] := \sum_{t=1}^T f_t(\tilde{w}_t) - f(w^*) \leq 4FL\sqrt{2(s_L + r + 1)PT} \quad (2)$$

Thus, $R[X] = O(\sqrt{T})$ since $\lim_{T \rightarrow \infty} \frac{R[X]}{T} = 0$

Proof: Since we follow all conditions and assumptions from Theorem 1, we need to show the newly added range R does not change the regret bound of SSP. In DSSP, the threshold is dynamically changing between s_L and s_U where $r \in R$ and $R = [0, s_U - s_L]$. We know that SSP with threshold s_L has a bound on regret according to (1). We only extend the threshold s_L of Theorem 1 to $s_L + r$ where r is the largest number from R . Suppose we set a fixed threshold s' for SSP, then our DSSP can be deducted to SSP when we set $s' = s_L + r$. Thus, we have an upper bound on regret of SSP with threshold s' . ■

V. EXPERIMENTS

We evaluate the performance of DSSP compared to other three distributed paradigms. We aim to find whether DSSP converges faster than SSP on average and whether it can maintain the predictive accuracy of SSP.

A. Experiment setup

1) *Hardware:* We conducted experiments on the SOSCIP GPU cluster [29] with up to four IBM POWER8 servers running Ubuntu 16.04. Each server has four NVIDIA P100 GPUs. Each server has 512 GB ram and 2×10 cores. The servers are connected with Infiniband EDR. Each server connects directly to a switch with dedicated 100 Gbps bandwidth.

We also set up a virtual cluster with a mixed GPU models by creating two Docker containers running Ubuntu 16.04 on a server with NVIDIA GTX1060 and GTX1080 Ti. The server has 64 GB ram and 8 cores. Each container is assigned with a dedicated GPU.

2) *Dataset:* We used CIFAR-10 and CIFAR-100 datasets [30] for image classification tasks. Both datasets have 50,000 training images and 10,000 test images. CIFAR-10 has 10 classes, while CIFAR-100 has 100 classes.

3) *Models:* We used a downsized AlexNet [3], ResNet-50 and ResNet-110 [31] as our deep neural network structure to evaluate the four distributed paradigms. We reduced the original AlexNet structure to a network with 3 convolutional layers and 2 fully connected layers to achieve faster convergence within 24 hours (which is the time limit we are allowed to run for each job on the SOSCIP cluster). We set the staleness threshold $s_L = 3$ and the range $R = [0, 12]$ for DSSP which is equivalent to the corresponding threshold range [3, 15] for SSP.

4) *Implementation:* We ran each paradigm on 4 servers. Each server represents a worker which has 4 GPUs. Each GPU loads a copy of the DNN model. Thus, there are 16 replica models for 4 workers. Each worker collects the computed gradients from 4 GPUs by the end of every iteration and sends the sum of the gradients to the parameter server. One of the 4 servers is also elected to run the parameter server when the training starts from the very beginning as defined in MXNet. We ran each experiment three times and chose the medium result based on the test accuracy.

B. Results and discussion

We used batch size 128, learning rate 0.001 in 300 epochs to train the downsized AlexNet on CIFAR-10. Figure 3a shows that DSSP, SSP and ASP converge much faster than BSP, and that DSSP and SSP converge to a higher accuracy than ASP. BSP is the slowest to complete the 300 epochs. The performance of DSSP and averaged SSP are similar, with DSSP converging a little bit faster to a bit higher accuracy. DSSP and averaged SSP complete 300 epochs almost at the same time. Note that this result is expected because the result of averaged SSP is the average over the results from 13 different threshold values from 3 to 15, and when its threshold is large, it is very fast, much faster than DSSP with a threshold range of [3,15]. However, a larger threshold of SSP incurs more staler gradients, which implies more noises and decreases the quality of iterations [18]. Theoretically, as the threshold s of SSP increases, the rate of convergence decreases per iteration update [18]. Figure 3b compares DSSP with individual SSPs with different threshold values. It shows that DSSP converges a bit faster to a higher accuracy than almost all of the SSPs except for one.

For ResNet-50 and ResNet-110 training on CIFAR-100, we used batch size 128, learning rate 0.05 and decay 0.1 twice at epoch 200 and 250 in 300 epochs for both. In Figure 3c, DSSP has the same convergence rate as ASP and SSP, and they converges much faster than BSP although BSP completes 300 epochs faster than others on both ResNets. Again, DSSP converges a little faster and archives a bit higher accuracy than averaged SSP in Figure 3e.

Four distributed paradigms on both ResNets behave in an opposite way compared to the downsized AlexNet which has fully connected layers in terms of the time taken to complete 300 epochs. The order from fastest to slowest is BSP, SSP, DSSP and ASP.

Based on the empirical results, we observe two opposite trends of ASP, DSSP, SSP and BSP with respect to their performance. The trends can be classified by the architecture of DNNs: ones that contain fully connected layers and ones that do not. Note that we do not count the final fully connected softmax layer as the fully connected layer over the discussion.

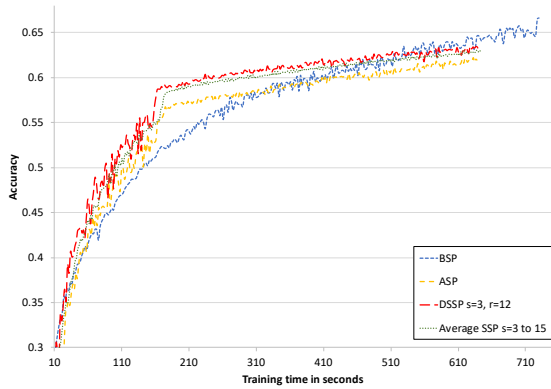
1) *DNNs with fully connected layers (AlexNet):* DSSP converges to a higher test accuracy faster than ASP, BSP and average SSPs in its corresponding staleness threshold range. ASP has the largest iteration throughput and its convergence rate is close to DSSP but it usually converges to a very low accuracy (the lowest of four paradigms) and diverges sometimes (see Figure 3a). DSSP performs between SSP and BSP in terms of final test accuracy. In this category, our DSSP converges faster than SSP and ASP to a higher accuracy. We know BSP guarantees the convergence and its accuracy is the same as using a single machine. Thus, it has no consistency errors caused by delayed updates. Given abundant time of training, BSP can reach the highest accuracy among all distributed paradigms. We do not discuss BSP in detail here since our focus is to show the benefits that our DSSP brings compared to SSP and ASP, both cost less training time than BSP.

2) *DNNs without fully connected layers (ResNet-50, ResNet-110):* DSSP converges faster than average SSP in its corresponding range on very deep neural networks. ASP appears to be a strong rival but it has no guarantee to converge as addressed in section I-A2. BSP delivers the highest iteration throughput. However, it converges slower and to a lower accuracy than other three paradigms mostly (see Figure 3c, 3e). The iteration throughputs of ASP, DSSP, SSP and BSP are in ascending order. In this category, the convergence rate of DSSP, ASP and SSP are very close. DSSP performs slightly above the average SSPs where the threshold s starts from 3 to 15 (see Figure 3d, 3f).

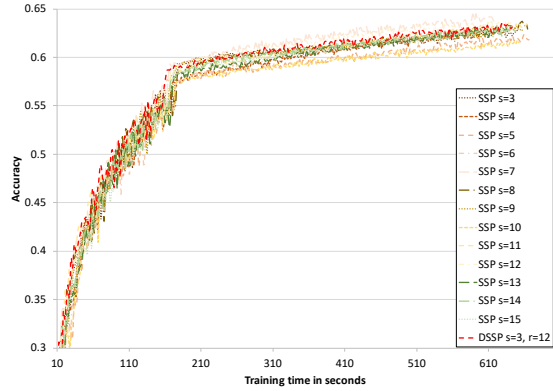
C. Demystify the difference

Below we answer two questions: (1) Why does the iteration throughput have the opposite trends for ASP, DSSP, SSP and BSP on DNNs with and without fully connected layers? (2) Why do pure convolutional neural networks (CNNs) receive a higher accuracy from DSSP, SSP and ASP than BSP? We temporarily name DNNs with fully connected layers as DNNs and pure CNNs as CNNs for the convenience of discussion.

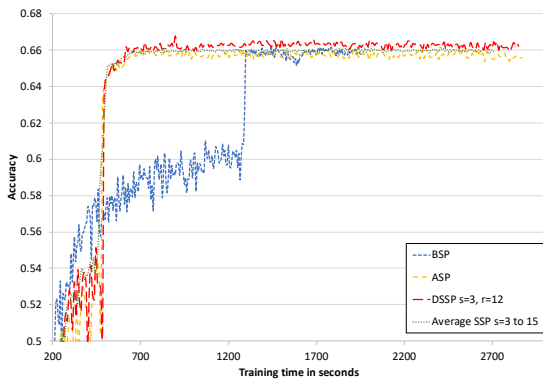
To answer the first question, we observe the difference between the two types of DNNs (with or without fully connected layers): ① A fully connected layer requires more parameters than a convolutional layer which uses shared parameters [32]. DNNs with fully connected layers have a large number of model parameters that need to be transmitted between workers and the server for updates. ② Convolutional layers require intensive computing time for matrix dot product operations while computing for fully connected layers involves simple linear algebra operations [33]. CNNs that only use convolutional layers take a lot of computing time, while their relatively smaller-size model parameters cost less data transmission time between workers and the server than DNNs. Moreover, when the ratio of computing time and communication time per iteration is



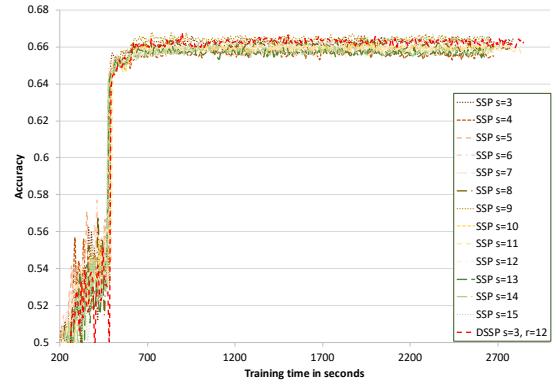
(a) All paradigms run on downsized AlexNet



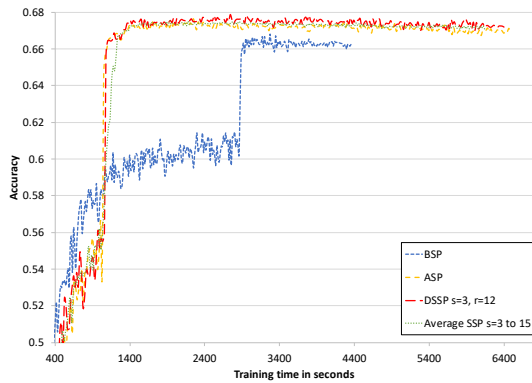
(b) DSSP and SSPs run on downsized AlexNet



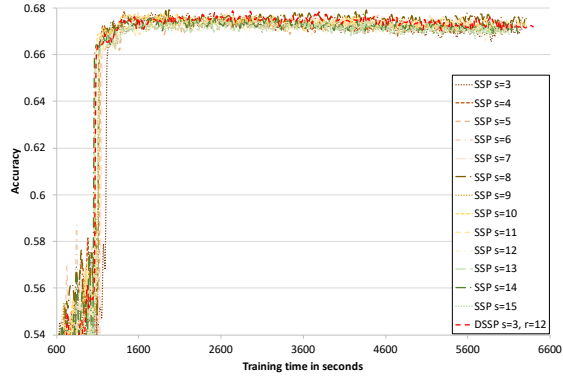
(c) All paradigms run on ResNet-50



(d) DSSP and SSPs run on ResNet-50



(e) All paradigms run on ResNet-110



(f) DSSP and SSPs run on ResNet-110

Figure 3: Distributed paradigms comparison on downsized AlexNet, ResNet-50 and ResNet-110 training for 300 epochs. Downsized AlexNet is trained on CIFAR-10 and both ResNets are trained on CIFAR-100. Average SSP on the right column is derived by averaging SSPs with threshold from 3 to 15 on the left column. Faster convergence to a targeted high accuracy indicates less training time is required for the paradigm.

small, less time can be saved per iteration for workers since computing time per iteration (or one mini-batch) is fixed for a model per worker. To the contrary, when the ratio is large, the communication time per iteration for each worker can be shifted by asynchronous-like parallel schemes and more time can be saved. Therefore, DSSP, SSP and ASP take less

training time on DNNs whereas BSP costs the least training time on CNNs.

To the second question, the answer lies in the difference between fully connected layers and convolutional layers. Fully connected layers are easy to overfit the data set due to its large number of parameters [34]. Thus, any error

introduced by staled updates can cause many parameters diverge in non-uniform convergence [16]. Informally, fully connected layers overfit to the errors injected by delayed updates or noise. Convolutional layers have less parameters due to the use of filters (shared parameters). For image classification tasks, a commonly used trick to train CNNs on a small data set is to increase the data by distorting the existing images and saving them [35] since CNNs are able to tolerate certain scale variations [36]. Distortion can be done by rotating the image, setting one or two of RGB pixels to zero or adding Gaussian noise to the image [37]. It is basically to introduce noise to images so that CNN models receive enhanced training and the predictions are improved. The errors caused by (not too) staled updates can give the same effect to the training model as the distortion. Figures 3c, 3e are good evidence to support that. Furthermore, [38] empirically shows that adding gradient noises improves the accuracy for training very deep neural networks which also happened in our ResNet-110 experiments (in Figure 3e).

D. Cluster with mixed GPU models

The results of DSSP and SSPs on ResNet-110 (see Figure 3e) do not show a significant difference on convergence rate on a homogeneous environment where GPUs are identical. Nonetheless, on the heterogeneous environment where we have one GTX1060 and one GTX1080 Ti running on each worker, DSSP converges faster and to a higher accuracy than SSP. We repeated the exact same experiments on ResNet-110 as earlier: use the same hyperparameters setting, run 3 trials on each paradigm and choose the medium one based on the test accuracy. Figure 4 and Table I show that DSSP can reach a higher accuracy significantly faster than SSP. The heterogeneous environment is very common in industry since new GPU models come to the market every year while the old models are still in use. ASP can fully utilize the individual GPU and achieves the largest iteration throughout. However, ASP also introduces the most staled updates among all distributed paradigms. It may converge to a lower accuracy than DSSP when the GPU models' processing capacities are significantly different since the iterations between the fastest worker and the slowest worker are dramatically different. In contrast, DSSP has consistent performance regardless the running environment since it adapts to the environment by adjusting the threshold dynamically for every iteration of workers.

VI. CONCLUSION

We presented dynamic staleness synchronous parallel (DSSP) paradigm for distributed training using the parameter server framework. DSSP improves SSP in the sense that with DSSP a user does not need to provide a specific staleness threshold which is hard to determine in practice, and also that DSSP can dynamically determine the value for the threshold from a range using a lightweight method

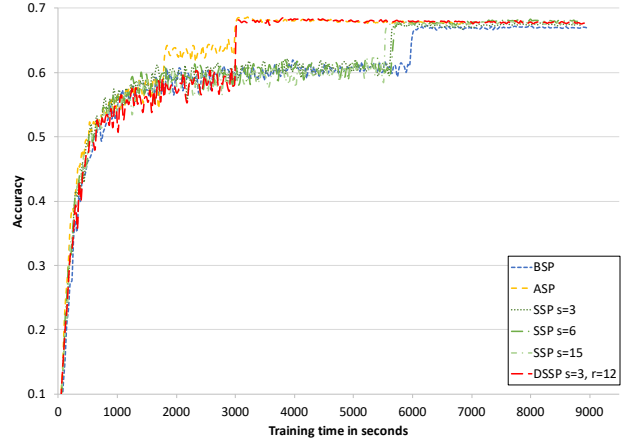


Figure 4: Trained ResNet-110 on CIFAR-100 with two workers on a mixed GPU cluster for 300 epochs. GTX1060 and GTX1080 Ti are assigned to individual worker. Our DSSP converges faster and achieves higher accuracy than SSP.

Distributed Paradigm	Time to reach 0.67 accuracy	Time to reach 0.68 accuracy
BSP	6159.2	—
ASP	2993.1	3017.2
SSP $s=3$	5678.2	—
SSP $s=6$	5703.8	6908.2
SSP $s=15$	5564.9	7255.6
DSSP $s_L=3, r=12$	3016.4	3046.3

Table I: Time in seconds to reach the targeted test accuracy in training. The maximum test accuracy of BSP and SSP with $s=3$ is 0.67. Trained ResNet-110 on CIFAR-100 with two workers for 300 epochs. Each worker has either GTX1080 Ti or GTX1060.

according to the run-time environment. This does not only alleviate the burden of an exact manual staleness threshold selection or multiple trials of hyperparameter selection, but it also provides flexibility of selecting different thresholds for different workers at different times. We provided theoretical analysis on the expected convergence of DSSP which inherits the same regret bound of SSP to show that DSSP converges in theory as long as the range is constant. We evaluated DSSP by training three DNNs on two datasets and compared its results with other distributed paradigms. For DNNs without fully connected layers, DSSP achieves higher accuracy than BSP and slightly better accuracy than averaged SSP. For DNNs with fully connected layers, DSSP generally converges faster than BSP, ASP and averaged SSP to a higher accuracy even though BSP can eventually reach the highest accuracy if it is given more training time. Unlike

ASP, DSSP ensures the convergence of DNNs by limiting the staled delays. DSSP gives significant improvement than SSP and BSP in a heterogeneous environment with mixed models of GPUs, converging much faster to a higher accuracy. DSSP also shows more stable performance on either homogeneous or heterogeneous environment compared to other three distributed paradigms. Furthermore, we discussed the difference in the trends of four distributed paradigms on DNNs with and without fully connected layers and the potential causes. For the further work, we will investigate how DSSP can adapt to an unstable environment where network connections are fluctuating between the servers.

ACKNOWLEDGMENT

This work is funded by the Natural Sciences and Engineering Research Council of Canada (NSERC), IBM Canada and the Big Data Research Analytics and Information Network (BRAIN) Alliance established by Ontario Research Fund - Research Excellence Program (ORF-RE). The experiments reported in this paper were performed on the GPU cluster of SOSCIP. SOSCIP is funded by the Federal Economic Development Agency of Southern Ontario, the Province of Ontario, IBM Canada, Ontario Centres of Excellence, Mitacs and 15 Ontario academic member institutions.

REFERENCES

- [1] J. Dean, G. Corrado, R. Monga, K. Chen, M. Devin, M. Mao, A. Senior, P. Tucker, K. Yang, Q. V. Le *et al.*, "Large scale distributed deep networks," in *Advances in neural information processing systems*, 2012, pp. 1223–1231.
- [2] 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," *arXiv preprint arXiv:1512.01274*, 2015.
- [3] A. Krizhevsky, I. Sutskever, and G. E. Hinton, "Imagenet classification with deep convolutional neural networks," in *Advances in neural information processing systems*, 2012, pp. 1097–1105.
- [4] C. Szegedy, V. Vanhoucke, S. Ioffe, J. Shlens, and Z. Wojna, "Rethinking the inception architecture for computer vision," in *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2016.
- [5] K. Simonyan and A. Zisserman, "Very deep convolutional networks for large-scale image recognition," *arXiv preprint arXiv:1409.1556*, 2014.
- [6] T.-Y. Lin, M. Maire, S. Belongie, J. Hays, P. Perona, D. Ramanan, P. Dollár, and C. L. Zitnick, "Microsoft coco: Common objects in context," in *European conference on computer vision*. Springer, 2014, pp. 740–755.
- [7] T. M. Chilimbi, Y. Suzue, J. Apacible, and K. Kalyanaraman, "Project adam: Building an efficient and scalable deep learning training system." in *OSDI*, vol. 14, 2014, pp. 571–582.
- [8] 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 *11th {USENIX} Symposium on Operating Systems Design and Implementation ({OSDI} 14)*, 2014, pp. 583–598.
- [9] 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," *IEEE Transactions on Big Data*, vol. 1, no. 2, pp. 49–67, 2015.
- [10] S. Landset, T. M. Khoshgoftaar, A. N. Richter, and T. Hasanin, "A survey of open source tools for machine learning with big data in the hadoop ecosystem," *Journal of Big Data*, vol. 2, no. 1, p. 24, 2015.
- [11] X. Meng, J. Bradley, B. Yavuz, E. Sparks, S. Venkataraman, D. Liu, J. Freeman, D. Tsai, M. Amde, S. Owen *et al.*, "Mllib: Machine learning in apache spark," *The Journal of Machine Learning Research*, vol. 17, no. 1, pp. 1235–1241, 2016.
- [12] J. Mao, W. Xu, Y. Yang, J. Wang, Z. Huang, and A. Yuille, "Deep captioning with multimodal recurrent neural networks (m-rnn)," *arXiv preprint arXiv:1412.6632*, 2014.
- [13] J. Zhou, X. Li, P. Zhao, C. Chen, L. Li, X. Yang, Q. Cui, J. Yu, X. Chen, Y. Ding *et al.*, "Kunpeng: Parameter server based distributed learning systems and its applications in alibaba and ant financial," in *Proceedings of the 23rd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*. ACM, 2017, pp. 1693–1702.
- [14] S. Lee, J. K. Kim, X. Zheng, Q. Ho, G. A. Gibson, and E. P. Xing, "On model parallelization and scheduling strategies for distributed machine learning," in *Advances in neural information processing systems*, 2014, pp. 2834–2842.
- [15] Y. Zhou, Y. Yu, W. Dai, Y. Liang, and E. Xing, "On convergence of model parallel proximal gradient algorithm for stale synchronous parallel system," in *Artificial Intelligence and Statistics*, 2016, pp. 713–722.
- [16] E. P. Xing, Q. Ho, P. Xie, and D. Wei, "Strategies and principles of distributed machine learning on big data," *Engineering*, vol. 2, no. 2, pp. 179–195, 2016.
- [17] A. V. Gerbessiotis and L. G. Valiant, "Direct bulk-synchronous parallel algorithms," *Journal of parallel and distributed computing*, vol. 22, no. 2, pp. 251–267, 1994.
- [18] Q. Ho, J. Cipar, H. Cui, S. Lee, J. K. Kim, P. B. Gibbons, G. A. Gibson, G. Ganger, and E. P. Xing, "More effective distributed ml via a stale synchronous parallel parameter server," in *Advances in neural information processing systems*, 2013, pp. 1223–1231.
- [19] H. Cui, J. Cipar, Q. Ho, J. K. Kim, S. Lee, A. Kumar, J. Wei, W. Dai, G. R. Ganger, P. B. Gibbons *et al.*, "Exploiting bounded staleness to speed up big data analytics." in *USENIX Annual Technical Conference*, 2014, pp. 37–48.
- [20] B. Recht, C. Re, S. Wright, and F. Niu, "Hogwild: A lock-free approach to parallelizing stochastic gradient descent," in *Advances in neural information processing systems*, 2011, pp. 693–701.

- [21] M. Zinkevich, M. Weimer, L. Li, and A. J. Smola, "Parallelized stochastic gradient descent," in *Advances in neural information processing systems*, 2010, pp. 2595–2603.
- [22] L. Wang, Y. Yang, R. Min, and S. Chakradhar, "Accelerating deep neural network training with inconsistent stochastic gradient descent," *Neural Networks*, vol. 93, pp. 219–229, 2017.
- [23] W. Dai, A. Kumar, J. Wei, Q. Ho, G. Gibson, and E. P. Xing, "High-performance distributed ml at scale through parameter server consistency models," in *Twenty-Ninth AAAI Conference on Artificial Intelligence*, 2015.
- [24] J. Wei, W. Dai, A. Qiao, Q. Ho, H. Cui, G. R. Ganger, P. B. Gibbons, G. A. Gibson, and E. P. Xing, "Managed communication and consistency for fast data-parallel iterative analytics," in *Proceedings of the Sixth ACM Symposium on Cloud Computing*. ACM, 2015, pp. 381–394.
- [25] M. Li, D. G. Andersen, A. J. Smola, and K. Yu, "Communication efficient distributed machine learning with the parameter server," in *Advances in Neural Information Processing Systems*, 2014, pp. 19–27.
- [26] J. Chen, X. Pan, R. Monga, S. Bengio, and R. Jozefowicz, "Revisiting distributed synchronous sgd," *arXiv preprint arXiv:1604.00981*, 2016.
- [27] S. Hadjis, C. Zhang, I. Mitliagkas, D. Iter, and C. Ré, "Omnivore: An optimizer for multi-device deep learning on cpus and gpus," *arXiv preprint arXiv:1606.04487*, 2016.
- [28] R. Zhang and J. Kwok, "Asynchronous distributed admm for consensus optimization," in *International Conference on Machine Learning*, 2014, pp. 1701–1709.
- [29] "SOSCIP GPU," accessed 2018-08-01. [Online]. Available: https://docs.scinet.utoronto.ca/index.php/SOSCIP_GPU
- [30] A. Krizhevsky and G. Hinton, "Learning multiple layers of features from tiny images," Citeseer, Tech. Rep., 2009.
- [31] K. He, X. Zhang, S. Ren, and J. Sun, "Deep residual learning for image recognition," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2016, pp. 770–778.
- [32] J. Tompson, R. Goroshin, A. Jain, Y. LeCun, and C. Bregler, "Efficient object localization using convolutional networks," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2015, pp. 648–656.
- [33] A. Krizhevsky, "One weird trick for parallelizing convolutional neural networks," *arXiv preprint arXiv:1404.5997*, 2014.
- [34] N. Srivastava, G. Hinton, A. Krizhevsky, I. Sutskever, and R. Salakhutdinov, "Dropout: a simple way to prevent neural networks from overfitting," *The Journal of Machine Learning Research*, vol. 15, no. 1, pp. 1929–1958, 2014.
- [35] L. Perez and J. Wang, "The effectiveness of data augmentation in image classification using deep learning," *arXiv preprint arXiv:1712.04621*, 2017.
- [36] Y. Xu, T. Xiao, J. Zhang, K. Yang, and Z. Zhang, "Scale-invariant convolutional neural networks," *arXiv preprint arXiv:1411.6369*, 2014.
- [37] L. Kang, P. Ye, Y. Li, and D. Doermann, "Simultaneous estimation of image quality and distortion via multi-task convolutional neural networks," in *Image Processing (ICIP), 2015 IEEE International Conference on*. IEEE, 2015, pp. 2791–2795.
- [38] A. Neelakantan, L. Vilnis, Q. V. Le, I. Sutskever, L. Kaiser, K. Kurach, and J. Martens, "Adding gradient noise improves learning for very deep networks," *arXiv preprint arXiv:1511.06807*, 2015.