# *MindTheStep-AsyncPSGD*:
## Adaptive Asynchronous Parallel Stochastic Gradient Descent

Karl Bäckström, Marina Papatriantafilou, Philippas Tsigas
Dept. of Computer Science and Engineering, Chalmers University of Technology, Gothenburg, Sweden
{bakarl, ptrianta, tsigas}@chalmers.se

*Abstract*—**Stochastic Gradient Descent (SGD) is very useful in optimization problems with high-dimensional non-convex target functions, and hence constitutes an important component of several Machine Learning and Data Analytics methods. Recently there have been significant works on understanding the parallelism inherent to SGD, and its convergence properties. Asynchronous, parallel SGD (*AsyncPSGD*) has received particular attention, due to observed performance benefits. On the other hand, asynchrony implies inherent challenges in understanding the execution of the algorithm and its convergence, stemming from the fact that the contribution of a thread might be based on an old (stale) view of the state. In this work we aim to deepen the understanding of *AsyncPSGD* in order to increase the statistical efficiency in the presence of stale gradients. We propose new models for capturing the nature of the staleness distribution in a practical setting. Using the proposed models, we derive a staleness-adaptive SGD framework, *MindTheStep-AsyncPSGD*, for adapting the step size in an online-fashion, which provably reduces the negative impact of asynchrony. Moreover, we provide general convergence time bounds for a wide class of staleness-adaptive step size strategies for convex target functions. We also provide a detailed empirical study, showing how our approach implies faster convergence for deep learning applications.**

## I. INTRODUCTION

The explosion of data volumes available for Machine Learning (ML) has posed tremendous scalability challenges for machine intelligence systems. Understanding the ability to parallelise, scale and guarantee convergence of basic ML methods under different synchronization and consistency scenarios have recently attracted a significant interest in the literature. The classic Stochastic Gradient Descent (SGD) algorithm is a significant target of research studying its convergence properties under parallelism.

In SGD, the goal is to minimize a function $f : \mathbb{R}^d \to \mathbb{R}$ of a $d$-dimensional vector $x$ using a first-order light-weight iterative optimization approach; i.e., given a randomly chosen starting point $x_0$, SGD repeatedly changes $x$ in the negative direction of a stochastic gradient sample, which provably is the direction in which

the target function is expected to decrease the most. The step size $\alpha_t$ defines how coarse the updates are:

$$x_{t+1} \leftarrow x_t - \alpha_t \nabla F(x_t) \tag{1}$$

SGD is very useful in nonconvex optimization with high-dimensional target functions, and hence constitutes a major part in several ML and Data Analytics methods, such as regression, classification and clustering. In many applications, the target function is differentiable and the gradient can be efficiently computed, e.g. Artificial Neural Networks (ANNs) using Back Propagation [32].

To better utilize modern computing architectures, recent efforts propose *parallel* SGD methods, complemented with different approaches for analyzing the convergence. However, asynchrony poses challenges in understanding the algorithm due to *stale* views of the state of $x$, which leads to reduced *statistical efficiency* in the SGD steps, requiring a larger number of iterations for achieving similar performance. In this work, we focus on increasing the statistical efficiency of the SGD steps, and propose a staleness-adaptive framework *MindTheStep-AsyncPSGD* that adapts parameters to significantly reduce the number of SGD steps required to reach sufficient performance. Our framework is compatible with recent orthogonal works focusing on computational efficiency, such as efficient parameter server architectures [14][9] and efficient gradient communication [5][31].

**Motivation and summary of state-of-the-art**
Many established ML methods, such as ANN training and Regression, constitute of minimizing a function $f(x)$ that takes the form of a finite sum of error terms $L(d; x)$ parameterized by $x$, evaluated at different data points $d$ from a given set $D$ of measurements:

$$f_D(x) = \frac{1}{|D|} \sum_{d \in D} L(d; x) \tag{2}$$

where the parameter vector $x$, encodes previously gathered features from $D$. In this context, SGD typically selects mini-batches $B \subseteq D$ over which $f_B$ is minimized, and is known as *Mini-Batch* Gradient Descent (MBGD). This type of SGD reduces the computational load in each step and hence enables processing of large datasets more efficiently. Moreover, randomly selecting mini-

batches induces stochastic variation in the algorithm, which makes it effective in non-convex problems as well.

A natural approach to distribute work for objective functions of the form (2) is to utilize *data parallelism* [34], where different workers (threads in a multicore system or nodes in a distributed one) run SGD over different subsets of $D$. This will result in differently learned parameter vectors $x$, which are aggregated, commonly in a *shared parameter server* (thread or node). The aggregation typically computes the average of the workers contributions; this approach is referred to as *Synchronous* Parallel SGD (*SyncPSGD*) due to its barrier-based nature. In its simple form, SyncPSGD has scalability issues due to the waiting time that is inherent in the aggregation when different workers compute with different speed. As more workers are introduced to the system, the waiting time will increase unbounded. Requiring only a fixed number of workers in the aggregation, known as $\lambda$-*softsync*, bounds this waiting time [17]. The barrier-based nature of the synchronous approaches to parallel SGD enables a straightforward (yet expensive) linearization making the vast analysis of classical SGD applicable also to the parallel version. As a result, its convergence is well-understood also in the parallel case, which however suffers from the performance-degradation of the barrier mechanisms.

An alternative type of parallelization is *Asynchronous* Parallel SGD (*AsyncPSGD*), in which workers *get* and *update* the shared variable $x$ independently of each other. There are inherent benefits in performance due to that *AsyncPSGD* eliminates waiting time, however the lack of coordination implies that gradients can be computed based on *stale* (old) views of $x$, which are *statistically inefficient*. However, gains in *computational efficiency* due to parallelism and asynchrony can compensate for this, reducing the overall wall-clock computation time.

**Challenges** *AsyncPSGD* shows performance benefits due to allowing workers to continue doing work independently of the progress of other workers. However, asynchrony comes with inherent challenges in understanding the execution of the algorithm and its convergence. In this work we address mainly (i) understanding the impact on the convergence and statistical efficiency of *stale* gradients computed based on old views of $x$ and (ii) how to adapt the step size in SGD to accommodate for the presence of asynchrony and delays in the system.

**Contributions** With the above challenges in mind, in this work we aim to increase the understanding of *AsyncPSGD* and the effect of stale gradients in order to increase the statistical efficiency of the SGD iterations. To achieve this, we find models suitable for capturing the nature of the staleness distribution in a practical setting. Under the proposed models, we derive a staleness-adaptive framework *MindTheStep-AsyncPSGD*

for adapting the step size in the precense of stale gradients. We prove analytically that our framework reduces the negative impact of asynchrony. In addition, we provide an empirical study which shows that our proposed method exhibits faster convergence by reducing the number of required SGD iterations compared to *AsyncPSGD* with constant step size. In some more detail:

- We prove analytically scalability limitations of the standard *SyncPSGD* approach that have been observed empirically in other works.
- We propose a new distribution model for capturing the staleness in *AsyncPSGD*, and show analytically how the optimal parameters can be chosen efficiently. We evaluate our proposed models by measuring the distance to the real staleness distribution observed empirically in a deep learning application, and compare the performance to models proposed in other works.
- Under the proposed distribution models, we derive efficiently computable staleness-adaptive step size functions which we show analytically can control the impact of asynchrony. We show how this enables *tuning* the implicit momentum to any desired value.
- We provide an empirical evaluation of *MindTheStep-AsyncPSGD* using the staleness-adaptive step size function derived from our proposed model, where we observe a significant reduction in the number of SGD iterations required to reach sufficient performance.

Before the presentation of the results in Sections III-VI, we outline preliminaries and background. Following the results-sections, we provide an extensive discussion on related work, conclusions and future work.

## II. PRELIMINARIES

### A. Stochastic Gradient Descent

We consider the optimization problem

$$\underset{x}{\text{minimize}} \quad f(x) \tag{3}$$

for a function $f : \mathbb{R}^d \to \mathbb{R}$. In this context, we focus on methods to address this minimization problem (3) using SGD, defined by (1) for some randomly chosen starting position $x_0$. We assume that the stochastic gradient $\nabla F$ is an unbiased estimator of $\nabla f$, i.e. $\mathbf{E}[\nabla F(x) \mid x] = \nabla f(x)$ for all $x$. This assumption holds for several relevant applications, in particular for problems of the form (2), including regression and ANN training. We assume that the stochastic gradient samples are i.i.d, which is reasonable since the sampling occurs independently by different threads. For the analysis in section V we adopt some additional standard assumptions on smoothness and convexity which we will introduce in that section.

### B. System Model and Asynchronous SGD

We consider a system with $m$ workers (that can be threads in a multicore system or nodes in a distributed

one), which repeatedly compute gradient contributions based on independently drawn data mini-batches from some given data set $D$. We also consider a *shared parameter server* (that can be a thread or a node respectively), which communicates with each of the workers independently, to give state information and get updates that it applies according to the algorithm it follows.

The $m$ asynchronous workers aim at performing SGD updates according to (1). Since each worker $W$ must get a state $x_t$ prior to computing a gradient, there can be intermediate updates from other workers before gradient from $W$ is applied. The number of such updates defines the *staleness* $\tau_t$ corresponding to the gradient $\nabla F(x_t)$.

Assuming that the read and update operations can be performed atomically (see details in Section IV), under the system model above, the SGD update (1) becomes

$$x_{t+1} \leftarrow x_t - \alpha_t \nabla F(v_t) \qquad (4)$$

where $v_t = x_{t-\tau_t}$ is the thread's *view* of $x$.

We assume that the staleness values $\tau_t$ constitute a *stochastic process* which is influenced by the computation speed of individual threads as well as the scheduler. Unless explicitly specified, we make no particular assumptions on the scheduler or computational speed among threads, except that all delays follow the same distribution with the same expected delay, i.e. $\mathbf{E}[\tau_t] = \bar{\tau}$ for all $t$. We do not require the staleness to be globally upper bounded, only that updates are eventually applied, making our system model *fully asynchronous*.

While we assume above that gradient samples are pairwise independent, it is not reasonable to make the same assumption for the staleness. In fact, a staleness $\tau_t$ is by definition dependent on writing time of concurrent updates, which in turn are dependent on their respective staleness values. For the analysis in Section V, we assume that *stochastic gradients* and *staleness* are uncorrelated, i.e. that the stochastic variation of the gradients does not influence the delays and vice versa. This is also a realistic assumption, since delays are due to computation time and scheduling and the gradient's stochastic variation is due to random draws from a dataset.

### C. Momentum

SGD is typically inefficient in *narrow valleys* when the target function in some neighbourhood increases more rapidly in one direction relative to another. Such neighbourhoods are frequent in target functions that arise in ML applications due to their inherent highly irregular and non-convex nature. Adding *momentum* (5) to SGD has been seen to significantly improve the convergence speed for such functions. SGD with momentum, defined in (5), takes all previous gradient samples into account with exponentially decaying magnitude in its parameter $\mu$. As pointed out in [23], $\mu$ is often left out in parameter tuning, and in some instances even failed to be reported [1]. However, the optimal value of algorithmic parameters such as $\mu$, just like $\alpha$, depends the problem, underlying hardware, as well as the choice of other parameters. Tuning $\mu$ has been shown to significantly improve performance [30], especially under asynchrony [23].

For $\mu \in [0,1]$, SGD with momentum is defined by

$$x_{t+1} \leftarrow x_t + \mu(x_t - x_{t-1}) - \alpha_t \nabla F(x_t) \qquad (5)$$

### III. ON THE SCALABILITY OF SYNC-PSGD

Optimal convergence with *SyncPSGD* requires, as observed empirically in [13], that the mini-batch size is reduced as the number of worker nodes increase. We prove analytically this empirical observation. We show that, from an optimization perspective, the effect of more workers on the convergence is equivalent to using a larger mini-batch size, which we refer to as the *effective* mini-batch size. For maintaining a desired effective mini-batch size, which is the case in many applications[15][22], workers must hence use smaller batches prior to the aggregation. Since the mini-batch size clearly is lower bounded, there is an implied strict upper bound on the number of worker nodes that can leverage the parallelization, which provides a bound on the scalability of the synchronous approach.

In mini-batch GD for target functions $f(x)$ of the form (2) the stochasticity is due to randomly drawing mini-batches $B$ of size $b$ from a dataset $D$ without replacement. For any positive mini-batch size $b$, we have that $F(x) = f_B(x)$ is an unbiased estimator of $f(x)$ since

$$\mathbf{E}[F(x)] = \mathbf{E}[f_B(x)] = \frac{b}{|D|} \sum_i f_{B_i}(x)$$

$$= \frac{b}{|D|} \sum_i \frac{1}{b} \sum_{d \in B_i} L(d;x) = \frac{1}{|D|} \sum_{d \in D} L(d;x) = f(x)$$

Hence, the SGD updates are in expectation representing the entire dataset $D$. Note that we assume $\bigcup B_i = D$. We have, however, that as the batch size $b$ increases, the variation of $F(x)$ diminishes. One can realize this by considering the extreme case $b = |D|$ for which the data sampling is deterministic. Hence, decreasing $b$ induces larger variance for the expectation $\mathbf{E}[F(x)] = f(x)$. This enables SGD to avoid local minima and hence be effective also in non-convex optimization problems.

The optimal value of $b$ is dependent on the problem and requires tuning. In particular, it has been seen that the convergence can suffer if $b$ is too large [15][22].

In the following theorem we show that by increasing the number of worker nodes in SyncPSGD, from an optimization perspective, we get a behavior equivalent to a sequential execution of SGD with a larger mini-batch size, which we refer to as the *effective* mini-batch size.

**Theorem 1.** SyncPSGD *with $m$ workers, all using batch size $b$, is equivalent to a sequential execution of SGD with batch size $m \cdot b$, reffered to as effective batch size.*

The proof appears in the appendix due to space limitations. The main idea is to compute the average of two workers, using batch size $b$, from which it is clear that the result is equivalent to an execution of sequential SGD with batch size $2b$. The result follows inductively.

Since the mini-batch size is clearly lower bounded, Theorem 1 implies that for a sufficiently large number of worker nodes, the effective mini-batch size scales linearly in the number of workers nodes. In order to maintain reasonable mini-batch size with sufficient variation in the updates, this implies a strict upper bound on the number of workers nodes. Moreover, under the assumption that there is an optimal mini-batch size $b^*$ for a given problem, which has been seen to be a common assumption, we have that the maximum number of workers possible in order to achieve optimal convergence is exactly $m = b^*$, each using mini-batch size $b = 1$.

## IV. THE PROPOSED FRAMEWORK

We outline *MindTheStep-AsyncPSGD* for staleness-adaptive steps and analyze how to choose a suitable adaptive step size function under different staleness models. Due to space limitations, proofs appear in the appendix, while brief arguments are presented here instead.

### A. The MindTheStep-AsyncPSGD Framework

We consider a standard parameter-server type of algorithm [14][18], with atomic read and write operations, ensuring that workers acquire consistent views of the state $x$. In a distributed system, the consistency can be realized through the communication protocol. In a multi-core system, where worker nodes are threads and $x$ can be stored on shared memory, consistency can be realized with appropriate synchronization and producer-consumer data structures, with the extra benefit that they can pass pointers to the data (parameter arrays) instead of moving it. In Algorithm 1 we show the pseudocode for *MindTheStep-AsyncPSGD*, describing how standard *AsyncPSGD* using a parameter server (thread or node) is extended with a staleness-adaptive step.

Note that *MindTheStep-AsyncPSGD* as a framework essentially "modularizes" the role of $\alpha$ as a parameter that can configure and tune performance, with criteria and benefits that are analysed in the next subsection.

### B. Tuning the impact of asynchrony

As pointed out in [23], asynchrony and delays introduce *memory* in the behaviour of the algorithms. In particular, in [23, Theorem 2], they quantify this and show its resemblance to momentum, however for

---

**Algorithm 1:** *MindTheStep-AsyncPSGD*

---
1 GLOBAL start point $x_0$, functions $F(x)$ and $\alpha(\tau)$

2 Worker $W$;      8 Parameter server $S$;
3 $(t, x) \leftarrow (0, x_0)$      9 $(t', x) \leftarrow (0, x_0)$
4 **repeat**      10 **repeat**
5     compute $g \leftarrow \nabla F(x)$      11     *receive* $(t, g)$ from a
6     *send* $(t, g)$ to $S$            ready worker $W$
7     *receive* $(t, x)$ from $S$      12     $\tau \leftarrow t' - t$
                                   13     $x \leftarrow x - \alpha(\tau)g$
                                   14     $t' \leftarrow t' + 1$
                                   15     *send* $(t', x)$ to $W$

---

a constant step size. The corresponding result for a stochastic staleness-adaptive step size is formulated here:

**Lemma 1.** *Let $\tau$ be distributed according to some PDF $p$ such that $P[\tau = i] = p(i)$. Then, for an adaptive step size function $\alpha(\tau)$, we have*

$$\mathbf{E}[x_{t+1} - x_t] = \mathbf{E}[x_t - x_{t-1}] + \sum_{i=0}^{\infty} \big(p(i)\alpha(i) - \quad (6)$$
$$p(i+1)\alpha(i+1)\big)\nabla f(x_{t-i-1}) - p(0)\alpha(0)\nabla f(x_t)$$

The proof of Lemma 1 follows the structure the one in [23], now taking into account the adaptive step size. The main takeaways from Lemma 1 are that, under asynchrony, (i) the gradient contribution diminishes as the number of workers increases[1]; (ii) there is a momentum-like term introduced with parameter $\mu = 1$ and (iii) the update depends on the series term:

$$\Sigma_{p,\alpha}^{\nabla} = \sum_{i=0}^{\infty} \big(p(i)\alpha(i) - p(i+1)\alpha(i+1)\big)\nabla f(x_{t-i-1})$$
$$(7)$$

which quantifies the potential impact of stale gradients depending on the distribution of $\tau$.

The issue of diminishing gradient contributions as the number of workers increase can in theory be resolved by choosing a larger $\alpha$. However, this would require step sizes proportional to $p(0)^{-1}$, which rapidly grows out of bounds as the number of workers increase. Since large $\alpha$ can significantly impact the statistical efficiency of the SGD steps in practice and in fact needs to be carefully tuned, this poses a scalability limitation.

This is where *MindTheStep-AsyncPSGD* can help tune the impact of asynchrony, as we show in the following.

**Momentum from geometric $\tau$.** Assuming a geometrically distributed $\tau$, the series $\Sigma_{p,\alpha}^{\nabla}$ is manifested in the convergence behaviour in the form of asynchrony-induced memory with a *momentum* effect; see Theorem 3 of [23], repeated here for self-containment:

---

[1]Here it is assumed that $p(0)$ tends to zero as the number of workers increases. This is easily realized for our proposed $CMP$ $\tau$ model (12). For the geometric staleness model we confirm empirically in section VI that this assumption holds in practice, recall that $p(0) = p$.

**Theorem 2** ([23]). *Let all $\tau_t$ be geometrically distributed with parameter $p$, i.e. $\mathbf{P}[\tau = k] = p(1-p)^k$. Then, for a constant $\alpha$, the expected update (4) becomes*

$$\mathbf{E}[x_{t+1} - x_t] = (1-p)\mathbf{E}[x_t - x_{t-1}] - p\alpha\nabla f(x_t) \quad (8)$$

The statement of Theorem 2 is easily confirmed by substituting $p(i)$ in (7) with constant $\alpha$ with the geometric PDF, which yields $\Sigma_{p,\alpha}^{\nabla} = -p\mathbf{E}[x_t - x_{t-1}]$.

Eq. (8) resembles the definition of momentum, with expected implicit asynchrony-induced momentum of magnitude $\mu = 1 - p$. As the number of workers grow and $p$ tends to 0, Theorem 2 suggests an implicit momentum that approaches 1. This would imply a scalability limitation since the parameter $\mu$ requires careful tuning.

Assuming a geometric staleness model, we show in the following theorem how *MindTheStep-AsyncPSGD* with a particular step size function resolves this issue.

**Theorem 3.** *Let staleness $\tau \in Geom(p)$ and*

$$\alpha_t = C^{-\tau_t} p^{-1} \alpha \quad (9)$$

*where $\alpha$ is a parameter to be chosen suitably. Then*

$$\mathbf{E}[x_{t+1} - x_t] = \mu_{C,p}\mathbf{E}[x_t - x_{t-1}] - \alpha\nabla f(x_t)$$

*and the implicit asynchrony-induced momentum is*

$$\mu_{C,p} = 2 - (1-p)/C \quad (10)$$

This is confirmed by substituting $\alpha(\tau)$ in (7) with the adaptive step(9). Note that the expected implicit momentum vanishes for $C = (1-p)/2$. More generally:

**Corollary 1.** *Any desired momentum $\mu^*$ is, in expectation, implicitly induced by asynchrony by using the staleness-adaptive step size in (9) with*

$$C = (1-p)/(2-\mu^*) \quad (11)$$

**Applicability of geometric $\tau$.** Each gradient staleness is comprised by two parts, one of which is the staleness $\tau_C$ which counts the number of gradients applied from other workers concurrent with the gradient computation. The second part of the staleness, which we denote $\tau_S$, counts, after the gradient computation of a worker finishes, the number of gradients from other workers which are applied first, which is decided by the order with which the workers are scheduled to apply their updates. The complete staleness of a gradient is $\tau = \tau_C + \tau_S$. Note that, if we assume a uniform fair stochastic scheduler, then $\tau_S$ is decided exactly by the number of Bernoulli trials until a specific gradient is chosen, hence $\tau_S \in Geom(\cdot)$. The geometric $\tau$ model is therefore applicable for problems where $\tau_C << \tau_S$, i.e. when the gradient computation time typically is smaller than the time it takes to apply a computed gradient (eq. 4).

Now consider also relevant applications of SGD where the gradient computation time $\tau_C$ is far from negligible, e.g the increasingly popular Deep Learning, which typically includes ANN training with BackProp [32] for gradient computation. The BackProp algorithm requires in the best case multiple multiplications of matrices of dimension $d$, which by far dominates the SGD update step (4) which consists of exactly $d$ floating point multiplications and additions. For such applications the geometric $\tau$ model is hence not sufficient; we confirm this empirically in Section VI. In the following, we propose a class of $\tau$ distributions which is more suitable.

**Conway-Maxwell-Poisson (CMP) $\tau$.** Considering applications with time-consuming gradient computation such as ANN training, we aim to find a suitable staleness model. Since now we consider (i) that $\tau_C >> \tau_S$ and (ii) that applying a computed gradient is relatively fast, we can consider the completion of gradient computations as rare arrival events. This opts for a variant of the Poisson distribution, such as the CMP distribution which in addition to Poisson has a parameter $\nu$ which controls the rate of decay. We have that $\tau \in CMP(\lambda, \nu)$ if

$$P[\tau = i] = \frac{1}{Z(\lambda,\nu)}\frac{\lambda^i}{(i!)^\nu} \ , \ Z(\lambda,\nu) = \sum_{j=0}^{\infty}\frac{\lambda^i}{(j!)^\nu} \quad (12)$$

which reduces to the Poisson distribution in the special case $\nu = 1$, i.e if $\tau \in CMP(\lambda, 1)$ then $\tau \in Poi(\lambda)$. For the remainder of this section we aim to further investigate the behaviour of parallelism in SGD under the CMP and Poisson models, and propose an adaptive step size strategy to reduce the negative impact and improve the statistical efficiency under asynchrony.

In a homogeneous system with $m$ equally powerful worker nodes/threads, we expect that the most frequent staleness observation (the distribution mode) should relate to the number of workers. More precisely, since a sequential execution would always have $\tau = 0$, an appropriate choice of $\tau$ distribution should have the mode $m - 1$. For the CMP distribution, we have that if $\tau \in CMP(\lambda, \nu)$ then the mode of $\tau$ is $\lfloor \lambda^{1/\nu} \rfloor$, and we therefore hypothesize the following relation:

$$\lambda^{1/\nu} = m \quad (13)$$

For the special case $\nu = 1$, i.e. a Poisson $\tau$ model, (13) enables us to immediately choose an appropriate value for $\lambda$ given the number of workers $m$. In general, (13) simplifies the parameter search when fitting a CMP distribution model to a one-dimensional line search, which is in practice a significant complexity reduction.

**$\tau$-adaptive $\alpha$.** In the following, we argue analytically about how to choose an adaptive step size function $\alpha(\tau)$ for reducing the negative impact of stale gradients. We will see how a certain $\tau$-adaptive step size can bound

the magnitude of $\Sigma_{p,\alpha}^{\nabla}$ (7), and even tune the implicit asynchrony-induced momentum to any desired value.

**Theorem 4.** *Assume* $\tau \in CMP(\lambda,\nu)$*, and let the adaptive step size function be defined as follows:*

$$\alpha(\tau) = C\lambda^{-\tau}(\tau!)^{\nu}\alpha \tag{14}$$

*for any constant C. Then we have* $\Sigma_{p,\alpha}^{\nabla} = 0$.

Theorem 4 shows how a simple and tunable $\tau$-adaptive step size mitigates the $\Sigma_{p,\alpha}^{\nabla}$ quantity. The proof consists of confirming that each contribution of the sum $\Sigma_{p,\alpha}^{\nabla}$ (7) vanishes when applying the definition of the CMP distribution (12) and the adaptive step size (14).

However, from Lemma 1, we see that even though $\Sigma_{p,\alpha}^{\nabla}$ is mitigated by the adaptive step size (14), the SGD steps still have a fixed implicit momentum term of magnitude $\mu = 1$. We show in Theorem 5 how the implicit momentum can be tuned to any desired value through a particular choice of $\alpha(\tau)$.

**Theorem 5.** *Assume* $\tau \in CMP(\lambda,\nu)$*. Then, we have that* $\Sigma_{p,\alpha}^{\nabla}$ *in expectation takes the form of asynchrony-induced momentum of magnitude exactly K, i.e.*

$$\Sigma_{p,\alpha}^{\nabla} = K\mathbf{E}[x_t - x_{t-1}]$$

*when using the adaptive step size function:*

$$\alpha(\tau) = c(\tau)\lambda^{-\tau}(\tau!)^{\nu}\alpha \tag{15}$$

*where*

$$c(\tau) = 1 - \frac{K}{\alpha e^{\lambda}}\sum_{j=0}^{\tau-1}\frac{\lambda^j}{(j!)^{\nu}} \tag{16}$$

Theorem 5 shows how the series term $\Sigma_{p,\alpha}^{\nabla}$ can take the form of momentum of desired magnitude by using a particular $\tau$-adaptive step size. The main idea of the proof is the observation that the contributions of $\Sigma_{p,\alpha}^{\nabla}$ are simplified by the particular choice of the adaptive factor $c(\tau)$ (15), and the result follows from the definition of expectation. The $c(\tau)$ contains a sum that is $\mathcal{O}(\tau)$ in computation time. This indicates that such an adaptive step size function might not scale well, since $\tau$ is expected to be in the magnitude of $m$. In the following Corollary we show how this is resolved by the corresponding $\alpha(\tau)$ under the Poisson $\tau$-model.

**Corollary 2.** *Assuming* $\tau \in Pois(\lambda)$*, the series term* $\Sigma_{p,\alpha}^{\nabla}$ *takes the form of implicit momentum of magnitude K when using the adaptive step size function:*

$$\alpha(\tau) = \left(1 - \frac{K}{\alpha}\frac{\Gamma(\tau,\lambda)}{\Gamma(\tau)}\right)\lambda^{-\tau}\tau!\alpha \tag{17}$$

*where* $\Gamma(\cdot)$ *and* $\Gamma(\cdot,\cdot)$ *are the Gamma and Upper Incomplete Gamma function, respectively.*

Corollary 2 shows how the series $\Sigma_{p,\alpha}^{\nabla}$ is in expectation replaced by momentum of any desired magnitude. Assuming Poisson $\tau$, the $\mathcal{O}(\tau)$ sum in (16) is replaced in (17) by the Gamma and Upper Incomplete Gamma function, for which there exist efficient ($\mathcal{O}(1)$) and accurate numerical approximation methods [12].

## V. CONVEX CONVERGENCE ANALYSIS

In this section we analyze the convergence time of *MindTheStep-AsyncPSGD*-type algorithms for convex and smooth optimization problems. Here, too, due to space limitations, the proofs appear the appendix, while brief intuitive arguments are presented instead.

Consider the optimization problem (3) where an acceptable solution $x$ satisfies $\epsilon$-convergence, defined as

$$\|x - x^*\|^2 \leq \epsilon \tag{18}$$

We assume that the problem is addressed using *MindTheStep-AsyncPSGD* under the system model described in Section II. Note that we consider a staleness-adaptive step size, hence $\alpha_t = \alpha(\tau_t)$ is stochastic.

For the analysis in this section, we consider strong convexity and smoothness, specified in *Assumption 1*. These analytical requirements are common in convergence analysis for convex problems [10][29][4][6].

**Assumption 1.** *We assume that the objective function* $f$ *is, in expectation with respect to the stochastic gradients, strongly convex with parameter c with L-Lipschitz continuous gradients and that the second momentum of the stochastic gradient is upper bounded.*

$$\mathbf{E}\left[(x - y)^T(\nabla f(x) - \nabla f(y)) \mid x,y\right] \geq c\|x - y\|^2 \tag{19}$$

$$\mathbf{E}\left[\|\nabla F(x) - \nabla F(y)\|\right] \mid x,y] \leq L\|x - y\| \tag{20}$$

$$\mathbf{E}\left[\|\nabla F(x)\|^2 \mid x\right] \leq M^2 \tag{21}$$

The assumption (19) is standard in convex optimization and ensures that gradient-based methods will converge to a global optimum. Lipschitz continuity (20) is a type of strong continuity which bounds the rate with which the gradients can vary. Due to that $\mathbf{E}[\nabla F(x^*)] = 0$, (21) can be interpreted as bounding the variance of the gradient norm around the optimum $x^*$.

In addition to our system model in Section II, we make the following assumption on the staleness process:

**Assumption 2.** *The staleness process* $(\tau_i)$ *is non-anticipative, i.e. mean-independent of the outcome of future states of the algorithm (e.g. future delays and gradients). In particular, we have*

$$\mathbf{E}[\tau_i \mid \tau_t] = \mathbf{E}[\tau_i] \text{ for all } i < t$$

Assumption 2 is reasonable considering that the staleness (i.e. scheduler's decisions) at time $i$ should not

be considered to be influenced by staleness values $\tau_t$ of gradients yet to be computed.

Under Assumptions 1 and 2 above, we give a general bound on the number of iterations sufficient for expected $\epsilon$-convergence in the following theorem:

**Theorem 6.** *Consider the unconstrained convex optimization problem of (3). Under Assumptions 1 and 2, for any $\epsilon > 0$, there is a sufficiently large number $T$ of asynchronous SGD updates of the form (4) such that*

$$T \leq \left( 2\big(c - LM\epsilon^{-1/2}\mathbf{E}\left[\tau\alpha\right]\big)\mathbf{E}\left[\alpha\right] - \epsilon^{-1}M^2\mathbf{E}\left[\alpha^2\right] \right)^{-1} \ln\left(\|x_0 - x^*\|^2\epsilon^{-1}\right) \quad (22)$$

*for which we have* $\mathbf{E}[\|x_T - x^*\|^2] < \epsilon$

The main idea in the proof of Theorem 6 is to bound $\|x_{t+1} - x^*\|/\|x_t - x^*\|$, which quantifies the improvement of each SGD step. The statement then follows from a recursive argument.

**Corollary 3.** *Consider the optimization problem and the conditions as in Theorem 6. There exists a choice of a step size $\alpha$ such that the convergence time $T$ is in the magnitude of $\mathcal{O}\left(\mathbf{E}\left[\tau\right]\right)$ (remember $\mathbf{E}\left[\tau\right]$ is denoted by $\bar{\tau}$). In particular, letting $\alpha$ be*

$$\alpha = \theta\frac{c\epsilon M^{-1}}{M + 2L\sqrt{\epsilon}\bar{\tau}} \quad (23)$$

*for a tunable factor $\theta \in (0,2)$, there exists a $T$ such that*

$$T \leq \frac{M + 2L\sqrt{\epsilon}\bar{\tau}}{\theta(2 - \theta)c^2M^{-1}\epsilon} \ln(\epsilon^{-1}\|x_0 - x^*\|^2) \quad (24)$$

The results in Theorem 6 and Corollary 3 are related to the results presented in [10] and [4]. The main differences are that in our analysis we tighten the bound with a factor $(2-\theta)^{-1}$, expand the allowed step size interval, as well as relax the *maximum staleness* assumption and reduce the magnitude of the bound from linear in the *maximum* staleness $\mathcal{O}(\hat{\tau})$ to the *expected* $\mathcal{O}(\bar{\tau})$.

In the following corollary, we give a general bound assuming *any* non-increasing step size function $\alpha(\tau)$.

**Corollary 4.** *Under the same conditions as Theorem 6, let $\alpha_t = \alpha(\tau_t)$ be a non-increasing function of $\tau_t$. Then we have the following bound on the expected number of iterations until convergence:*

$$T \leq \big(2c\mathbf{E}\left[\alpha\right] - \epsilon^{-1}M\left(M + 2L\sqrt{\epsilon}\bar{\tau}\right)\mathbf{E}\left[\alpha^2\right]\big)^{-1} \cdot \ln(\epsilon^{-1}\|x_0 - x^*\|^2) \quad (25)$$

Corollary 4 describes a general convergence bound for any step size function $\alpha(\tau)$ which decays in $\tau$. We see that such step size functions also achieve the asymptotic $\mathcal{O}(\bar{\tau}^{-1})$ bound, similar to the one for a constant $\alpha$ (24).
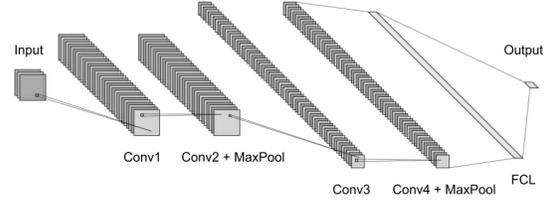


Fig. 1. CNN architecture; Four convolutional layers with $3\times 3$ kernels, with intermediate MaxPool layers. First two convolutions have 32 filters, the last two 64. The architecture has two fully connected layers, one with 256 neurons, and the output layer with 10 neurons.

## VI. Experimental Study

In this section we aim to evaluate the results derived in section IV in a practical setting. This is achieved by (i) measuring the accuracy and scalability of the proposed $\tau$-models (ii) evaluating the convergence properties of *MindTheStep-AsyncPSGD* with an adaptive step size function derived under the CMP/Poisson $\tau$ models.

**Setup.** We apply *MindTheStep-AsyncPSGD* for training a 4-layer Convolutional Neural Network (CNN) architecture (see Fig. 1) on the common image classification benchmark dataset CIFAR10 [16]. The performance of the CNN is measured as the *cross entropy* between the true and the predicted class distribution. The algorithm is evaluated on a setup with a 36-thread Intel Xeon CPU and 64GB memory. The implementation is in Python 2.7 and uses the standard Python multiprocessor library as well as TensorFlow [1] for gradient computation.

**CMP/Poisson $\tau$.** We evaluate the $\tau$ models (Poisson, CMP) proposed in section IV by comparing with the $\tau$ distribution observed in practice for different number of workers. We compare our proposed $\tau$ models with distributions proposed in other works, namely the geometric $\tau$ model [23] and the bounded uniform $\tau$ model [29].

The distribution parameters in Table I are found through an exhaustive search where we aim to minimize the Bhattacharyya distance to the $\tau$ distribution observed in practice. Note that: (i) For the Poisson $\tau$ model, as hypothesized in Section IV, the distribution parameter $\lambda$ indeed corresponds well to the number of worker nodes. From Fig. 2 we see that the proposed CMP and Poisson $\tau$ models by far outperforms the geometrical and uniform $\tau$ models, in particular for larger number of workers. (ii) As mentioned in Footnote 1, we confirm in Table I that $\mathbf{P}[\tau = 0]$, i.e. $p$, decays as $m$ increases. (iii) We see in Fig. 2 that the CMP $\tau$ model outperforms the others in terms of accuracy and scalability. The CMP distribution parameter $\nu$ is found through a 1-d search, and using the assumption (13) the other parameter $\lambda$ is calculated. The result in Fig. 2 therefore validates the assumption (13).

**Convergence with $\tau$-adaptive $\alpha$.** We evaluate *MindTheStep-AsyncPSGD* compared with standard
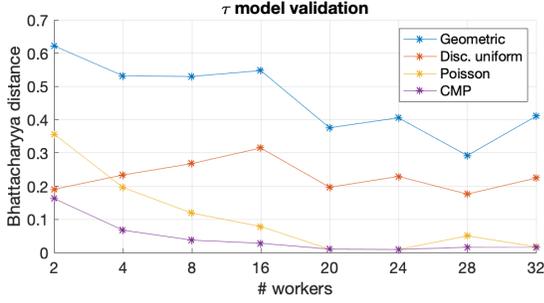
Fig. 2. Bhattacharyya distance of different $\tau$ models compared to the observed distribution. The graph shows that the CMP $\tau$ model is the most accurate in all tests, with the Poisson $\tau$ model as a close second. The uniform and geometric $\tau$ models are persistently less accurate, and show poor scalability in comparison.

| $\tau$ model | 2 | 4 | 8 | 16 | 20 | 24 | 28 | 32 |
|---|---|---|---|---|---|---|---|---|
| $p$ (Geom) | 0.34 | 0.21 | 0.12 | 0.06 | 0.05 | 0.04 | 0.04 | 0.03 |
| $\hat{\tau}$ (Unif) | 2 | 5 | 11 | 22 | 31 | 37 | 48 | 48 |
| $\lambda$ (Pois) | 2.0 | 4.0 | 8.0 | 16.0 | 19.7 | 23.8 | 26.5 | 32 |
| $\nu$ (CMP) | 6.28 | 5.26 | 4.18 | 3.48 | 0.93 | 0.95 | 0.39 | 0.87 |

TABLE I
OPTIMAL DISTRIBUTION PARAMETERS FOR DIFFERENT NUMBER OF WORKERS

*AsyncPSGD* by measuring the number of *epochs* required until a certain error threshold is reached, epochs being the number of passes through the dataset. The number of SGD iterations in one epoch is $\lceil |D|/b \rceil$ where $|D|$ is the size of the dataset and $b$ the batch size. In our experiments we have $\lceil |D|/b \rceil = 469$. We consider performance in terms of *statistical efficiency*, i.e. the statistical benefit of each SGD step. In practice, the approach can be applied to any orthogonal work focusing on computational efficiency, such as efficient parameter server architectures [14][9] and efficient gradient communication and quantization [5][31].

We compare standard *AsyncPSGD* with constant step size $\alpha_c = 0.01$ to *MindTheStep-AsyncPSGD* with an adaptive step size function according to (17) with $\alpha = \alpha_c$, $K = 1$, and $\lambda = m$. In addition, we bound the step size $\alpha(\tau) \leq 5 \cdot \alpha_c$ to mitigate issues with numerical instability in the SGD algorithms, and (very infrequent) gradients with $\tau > 150$ are not applied.

In principle, given a sufficiently small $\alpha_c$, speedup can always be achieved by using an adaptive step size strategy $\alpha(\tau)$ which overall increases the average step size. To ensure a fair comparison, the adaptive step size function $\alpha(\tau)$ is normalized so that

$$\mathbf{E}_\tau[\alpha(\tau)] = \alpha_c \qquad (26)$$

where the expectation is taken over the real $\tau$ distribution observed in the system. Enforcing (26) ensures that any potential speedup is achieved due to how the step size function $\alpha(\tau)$ adaptively changes the impact of gradients
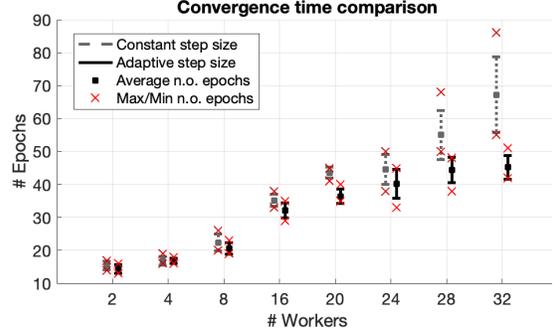


Fig. 3. *AsyncPSGD* vs. *MindTheStep-AsyncPSGD* comparison. The plot shows the n.o. epochs required until sufficient performance (cross-entropy loss $\leq 0.05$). The statistics are computed based on 5 runs, and the bar height corresponds to the standard deviation.

depending on their staleness, and not because of the overall magnitude of the step size.

Fig. 3 shows how *MindTheStep-AsyncPSGD* exhibits persistently faster convergence for different number of workers. For many workers ($m = 28,32$) *MindTheStep-AsyncPSGD* requires significantly fewer epochs compared to standard *AsyncPSGD* to achieve sufficient performance. Observe that for $m = 32$ the average speedup is $\times 1.5$ while the worst-case is $\times 1.7$.

## VII. RELATED WORK

Orthogonal to this work, there are numerous works dedicated to optimizing the effectiveness of SGD by utilizing data sparsity, topology of the search space, and other properties of the problems. One example is introducing momentum to the updates, originally proposed in [25], however not in the context of SGD. Apart from this, there are several variation of SGD in the sequential case introducing adaptivness to aspects of the problem topology, such as Adagrad, Adadelta, RMSprop, Adam, AdaMax, and Nadam (cf. [27] and references therein).

In [23] Mitliagkas et al. show that under certain stochastic delay models, asynchrony has an effect on convergence similar to momentum, referred to as asynchrony-induced or implicit momentum, where more workers imply a larger magnitude of the effect. In [20] these similarities are investigated further, and it is shown that *AsyncPSGD* and momentum shows different convergence rates in general and that *AsyncPSGD* is in fact faster in expectation. Since it has been seen [30] that the magnitude of momentum can have significant impact on convergence, the result by Mitliagkas et al. would imply a harsh scalability limitation of AsyncPSGD. In this paper, we show that under the same $\tau$ model as in [23], *MindTheStep-AsyncPSGD* can in theory mitigate this issue, and even allow the expected asynchrony-induced momentum to be tuned implicitly by the rate of adaptation. In addition, in this

work we propose a different class of $\tau$ distribution models, and show how they better capture the real $\tau$ values observed in a deep learning application. From our proposed models we derive an adaptive step size function $\alpha(\tau)$ which we show significantly reduces the number of SGD steps required for convergence.

Below we give a brief overview of works on synchronous distributed SGD. Under smoothness and convexity assumptions, in [34] and [3], synchronous distributed SGD with data-parallelism was observed and proven to accelerate convergence. This was implemented on a larger scale by Dekel et al. [11] where the convergence rates were improved under stronger analytical assumptions. In [14] and [17] the synchronization is relaxed using a Stale Synchronous Parameter Server with a tunable staleness threshold in order to reduce the waiting-time, which is shown to outperform synchronous SGD. In [13] Gupta et al. give a rigorous empirical investigation of practical trade-offs the number of workers, mini-batch size and staleness; the results provide useful insights in scalability limitations in synchronous methods with averaging. We address this issue in this paper from a theoretical standpoint and explain the results observed in practice. This is discussed in detail in Section III.

The study of numerical methods under parallelism is not new, and sparked due to the works by Bertsekas and Tsitsiklis [7] in 1989. Recent works [8][19] show under various analytical assumptions that the convergence of Async-PSGD is not significantly affected by asynchrony and that the noise introduced by delays is asymptotically negligible compared to the noise from the stochastic gradients. This is confirmed in [8] for convex problems (linear and logistic regression) for a small number of cores. In [19] Lian et al. relax the theoretical assumptions and establish convergence rates for non-convex minimization problems, assuming bounded gradient delays and number of workers. Lock-free Async-PSGD in shared-memory, i.e. HOGWILD!, was proposed by Niu et al. [26] and was shown to achieve near-optimal convergence rates assuming sparse gradients. Properties of Async-PSGD with sparse or component-wise updates have since been rigorously studied in recent literature due to the performance benefits of lock-freedom [28][24][10]. The gradient sparsity assumption was relaxed in the recent work [4] which magnified the convergence time bound in the order of magnitude $\sim \sqrt{d}$, $d$ being the problem dimensionality.

Delayed optimization in completely asynchronous first-order optimization algorithms was analyzed initially in [2], where Agarwal et al. introduce step sizes which diminish over the progression of SGD, depending on the maximum staleness allowed in the system, but not adaptive to the actual delays observed. In comparison, in this work we relax the maximum staleness restriction and

derive a strategy for adapting the step size depending on the actual staleness values observed in the system in an online fashion. Adaptiveness to delayed updates during execution was proposed and analyzed in [21] under assumptions of gradient sparsity and *read* and *write* operations having the same relative ordering. A similar approach was used in [33], however for synchronous SGD with the *softsync* protocol. In [33] statistical speedup is observed in some cases for a limited number of worker nodes, however by using *momentum SGD*, which is not the case in their theoretical analysis, and step size decaying schedules on top of the staleness-adaptive step size.

The work closest to ours is AdaDelay [29] which addresses a particular constrained convex optimization problem, namely training a logistic classifier with projected gradient descent. It utilizes a network of worker nodes computing gradients in parallel which are aggregated at a central parameter server with a step size that is scaled proportionally to $\tau^{-1}$. The staleness model in [29] is a uniform stochastic distribution, which implies a strict upper bound on the delays, making the system partially asynchronous. In comparison, in this work we analyze the convergence of *MindTheStep-AsyncPSGD* for non-convex optimization, relax the bounded gradient staleness assumption, as well as evaluate more delay models both theoretically and empirically. Moreover, we validate our findings experimentally by training a Deep Neural Network (DNN) classifier using real-world dataset, which constitutes a highly non-convex and high-dimensional optimization problem. In addition, we provide convergence analysis in the convex case for *MindTheStep-AsyncPSGD*, where we show explicitly a probabilistic time bound for $\epsilon$-convergence, for any step size function decaying in the staleness $\tau$.

## VIII. CONCLUSIONS AND FUTURE WORK

In this paper, we first analytically confirm scalability limitations of the standard *SyncPSGD*, which were observed empirically in other works; we thus motivate the need to further investigate asynchronous approaches. We propose a new class of $\tau$-distribution models, show analytically how the parameters can be efficiently chosen in a practical setting, and validate the models empirically, as well as compare to models proposed in other works.

We derive and analyse adaptive step size strategies which reduce the impact of asynchrony and stale gradients, using our framework *MindTheStep-AsyncPSGD*. We show that the proposed strategies enable turning asynchrony into implicit asynchrony-induced momentum of desired magnitude. We provide convergence bounds for a wide class of $\tau$-adaptive step size strategies for convex target functions. We validate our findings empirically for a deep learning application and show that *MindTheStep-AsyncPSGD* with our

proposed step size strategy converges significantly faster compared to standard *AsyncPSGD*.

The concept of staleness-adaptive *AsyncPSGD* has been under-explored, despite that, as shown here, it significantly improves scalability and helps maintain statistical efficiency. Continuing to investigate asynchrony-aware SGD, is therefore of interest. Future research directions also include further studying the nature of the staleness, i.e. effect of schedulers and synchronization methods, for understanding the impact of asynchrony and for choosing appropriate adaptive strategies.

## IX. ACKNOWLEDGEMENTS

## REFERENCES

[1] 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 *OSDI*, volume 16, pages 265–283, 2016.

[2] A. Agarwal and J. C. Duchi. Distributed delayed stochastic optimization. In *Advances in Neural Information Processing Systems*, pages 873–881, 2011.

[3] A. Agarwal, M. J. Wainwright, and J. C. Duchi. Distributed dual averaging in networks. In *Advances in Neural Information Processing Systems*, pages 550–558, 2010.

[4] D. Alistarh, C. De Sa, and N. Konstantinov. The convergence of stochastic gradient descent in asynchronous shared memory. In *ACM Symp. on Principles of Distributed Computing*, PODC '18, pages 169–178. ACM, 2018.

[5] D. Alistarh, D. Grubic, J. Li, R. Tomioka, and M. Vojnovic. Qsgd: Communication-efficient sgd via gradient quantization and encoding. In *Advances in Neural Information Processing Systems*, pages 1709–1720, 2017.

[6] D. Alistarh, T. Hoefler, M. Johansson, N. Konstantinov, S. Khirirat, and C. Renggli. The convergence of sparsified gradient methods. In *Advances in Neural Information Processing Systems*, pages 5977–5987, 2018.

[7] D. P. Bertsekas and J. N. Tsitsiklis. *Parallel and distributed computation: numerical methods*, volume 23. Prentice hall Englewood Cliffs, NJ, 1989.

[8] S. Chaturapruek, J. C. Duchi, and C. Ré. Asynchronous stochastic convex optimization: the noise is in the noise and sgd don't care. In *Advances in Neural Inf. Processing Systems*, 2015.

[9] H. Cui, H. Zhang, G. R. Ganger, P. B. Gibbons, and E. P. Xing. Geeps: Scalable deep learning on distributed gpus with a gpu-specialized parameter server. In *Proc. of the 11th European Conf. on Computer Systems*, page 4. ACM, 2016.

[10] C. M. De Sa, C. Zhang, K. Olukotun, C. Ré, and C. Ré. Taming the wild: A unified analysis of hogwild-style algorithms. In *Advances in Neural Information Processing Systems 28*, pages 2674–2682. Curran Associates, Inc., 2015.

[11] O. Dekel, R. Gilad-Bachrach, O. Shamir, and L. Xiao. Optimal distributed online prediction. In *Proc. of the 28th Int'l Conf. on Machine Learning (ICML-11)*, pages 713–720, 2011.

[12] P. Greengard and V. Rokhlin. An algorithm for the evaluation of the incomplete gamma function. *Advances in Computational Mathematics*, 45(1):23–49, 2019.

[13] S. Gupta, W. Zhang, and F. Wang. Model accuracy and runtime tradeoff in distributed deep learning: A systematic study. In *16th Int'l Conf. on Data Mining (ICDM)*,, pages 171–180. IEEE, 2016.

[14] 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*, pages 1223–1231, 2013.

[15] N. S. Keskar, D. Mudigere, J. Nocedal, M. Smelyanskiy, and P. T. P. Tang. On large-batch training for deep learning: Generalization gap and sharp minima. *arXiv:1609.04836*, 2016.

[16] A. Krizhevsky, G. Hinton, et al. Learning multiple layers of features from tiny images. Technical report, Citeseer, 2009.

[17] 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*, pages 2834–2842, 2014.

[18] 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 Symp. on Operating Systems Design and Implementation*, pages 583–598, 2014.

[19] X. Lian, Y. Huang, Y. Li, and J. Liu. Asynchronous parallel stochastic gradient for nonconvex optimization. In *Advances in Neural Information Processing Systems*, pages 2737–2745, 2015.

[20] T. Liu, S. Li, J. Shi, E. Zhou, and T. Zhao. Towards understanding acceleration tradeoff between momentum and asynchrony in nonconvex stochastic optimization. In *Advances in Neural Information Processing Systems 31*, pages 3686–3696. Curran Associates, Inc., 2018.

[21] B. McMahan and M. Streeter. Delay-tolerant algorithms for asynchronous distributed online learning. In *Advances in Neural Information Processing Systems 27*, pages 2915–2923. Curran Associates, Inc., 2014.

[22] D. Mishkin, N. Sergievskiy, and J. Matas. Systematic evaluation of convolution neural network advances on the imagenet. *Computer Vision and Image Understanding*, 161:11–19, 2017.

[23] I. Mitliagkas, C. Zhang, S. Hadjis, and C. Ré. Asynchrony begets momentum, with an application to deep learning. In *Communication, Control, and Computing (Allerton), 2016 54th Annual Allerton Conf. on*, pages 997–1004. IEEE, 2016.

[24] L. M. Nguyen, P. H. Nguyen, M. van Dijk, P. Richtárik, K. Scheinberg, and M. Takáč. Sgd and hogwild! convergence without the bounded gradients assumption. *arXiv:1802.03801*.

[25] B. T. Polyak. Some methods of speeding up the convergence of iteration methods. *USSR Computational Mathematics and Mathematical Physics*, 4(5):1–17, 1964.

[26] 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 (NIPS) 24*, pages 693–701. Curran Associates, Inc., 2011.

[27] S. Ruder. An overview of gradient descent optimization algorithms. *arXiv:1609.04747*, 2016.

[28] S. Sallinen, N. Satish, M. Smelyanskiy, S. S. Sury, and C. Ré. High performance parallel stochastic gradient descent in shared memory. In *Parallel and Distributed Processing Symp., 2016 IEEE Int'l*, pages 873–882. IEEE, 2016.

[29] S. Sra, A. W. Yu, M. Li, and A. Smola. Adadelay: Delay adaptive distributed stochastic optimization. In *Proc. of the 19th Int'l Conf. on Artificial Intelligence and Statistics*, pages 957–965, 2016.

[30] I. Sutskever, J. Martens, G. Dahl, and G. Hinton. On the importance of initialization and momentum in deep learning. In *Int'l Conf. on machine learning*, pages 1139–1147, 2013.

[31] W. Wen, C. Xu, F. Yan, C. Wu, Y. Wang, Y. Chen, and H. Li. Terngrad: Ternary gradients to reduce communication in distributed deep learning. In *Advances in neural information processing systems*, pages 1509–1519, 2017.

[32] P. J. Werbos. Applications of advances in nonlinear sensitivity analysis. In *System modeling and optimization*, pages 762–770. Springer, 1982.

[33] W. Zhang, S. Gupta, X. Lian, and J. Liu. Staleness-aware async-sgd for distributed deep learning. In *Proceedings of the Twenty-Fifth Int'l Joint Conf. on Artificial Intelligence*, IJCAI'16, pages 2350–2356. AAAI Press, 2016.

[34] M. Zinkevich, M. Weimer, L. Li, and A. J. Smola. Parallelized stochastic gradient descent. In *Advances in neural information processing systems*, pages 2595–2603, 2010.

APPENDIX

*Proof of Theorem 1.* Consider the case with two worker nodes. Assuming that the batches are disjoint, which is likely for large datasets, each SGD step is of the form

$$x_{t+1} = \frac{\left(x_t - \alpha\nabla f_{B_1}(x_t)\right) + \left(x_t - \alpha\nabla f_{B_2}(x_t)\right)}{2}$$

$$= x_t - \frac{\alpha}{2}\left(\nabla f_{B_1}(x_t) + \nabla f_{B_2}(x_t)\right)$$

For mini-batch GD, i.e. a target function of the form (2), and with mini-batch size $b$, the above formula becomes:

$$x_{t+1} = w - \frac{\alpha}{2}\left(\nabla\frac{1}{b}\sum_{d\in B_1}L(d, x_t) + \nabla\frac{1}{b}\sum_{d\in B_2}L(d, x_t)\right)$$

From linearity of the gradient, we have

$$x_{t+1} = w - \alpha\nabla\frac{1}{2b}\sum_{d\in B_1\cup B_2}L(d, x_t)$$

$$= w - \alpha\nabla f_{B_1\cup B_2}(x_t)$$

that corresponds to the SGD step with batch size $2b$. This inductively implies the theorem statement. $\square$

*Proof of Theorem 3.* We have from (4)

$$x_{t+1} - x_t = -\alpha_t\nabla F(v_t)$$

$$= x_t - x_{t-1} - (x_t - x_{t-1}) - \alpha_t\nabla F(v_t)$$

$$= x_t - x_{t-1} + \alpha_t\nabla F(v_{t-1}) - \alpha_t\nabla F(v_t)$$

Since the gradient and staleness processes are independent, we take first expectation conditioned on the staleness

$$\mathbf{E}[x_{t+1} - x_t \mid \tau_t, \tau_{t-1}] = \mathbf{E}[x_t - x_{t-1} \mid \tau_t, \tau_{t-1}]$$
$$+ \alpha_t\nabla f(v_{t-1}) - \alpha_t\nabla f(v_t)$$

Now, take expectation w.r.t. the stochastic staleness $\tau_t, \tau_{t-1}$

$$\mathbf{E}[x_{t+1} - x_t] = \mathbf{E}[x_t - x_{t-1}]$$
$$+ \mathbf{E}[\alpha_t\nabla f(v_{t-1})] - \mathbf{E}[\alpha_t\nabla f(v_t)]$$

$$= \mathbf{E}[x_t - x_{t-1}] + \sum_{i=0}^{\infty}P[\tau = i]\frac{\alpha\nabla f(x_{t-i-1})}{C^i p}$$

$$- \sum_{i=0}^{\infty}P[\tau = i]\frac{\alpha\nabla f(x_{t-i})}{C^i p}$$

$$= \mathbf{E}[x_t - x_{t-1}] + p\sum_{i=0}^{\infty}(1-p)^i\frac{\alpha\nabla f(x_{t-i-1})}{C^i p}$$

$$- p\sum_{i=0}^{\infty}(1-p)^i\frac{\alpha\nabla f(x_{t-i})}{C^i p}$$

$$= \mathbf{E}[x_t - x_{t-1}] - \alpha\nabla f(x_t) + \sum_{i=0}^{\infty}(1-p)^i\frac{\alpha\nabla f(x_{t-i-1})}{C^i}$$

$$- \sum_{i=1}^{\infty}(1-p)^i\frac{\alpha\nabla f(x_{t-i})}{C^i}$$

$$
= \mathbf{E}[x_t - x_{t-1}] - \alpha \nabla f(x_t)
$$
$$
+ \sum_{i=0}^{\infty} \left( \frac{(1-p)^i}{C^i} - \frac{(1-p)^{i+1}}{C^{i+1}} \right) \alpha \nabla f(x_{t-i-1})
$$
$$
= \mathbf{E}[x_t - x_{t-1}] - \alpha \nabla f(x_t)
$$
$$
+ \sum_{i=0}^{\infty} \frac{(1-p)^i}{C^i} \left( 1 - \frac{1-p}{C} \right) \alpha \nabla f(x_{t-i-1})
$$
$$
= \mathbf{E}[x_t - x_{t-1}] - \alpha \nabla f(x_t)
$$
$$
+ \left( 1 - \frac{1-p}{C} \right) \sum_{i=0}^{\infty} \frac{p(1-p)^i}{C^i p^{i+1}} \alpha \nabla f(x_{t-i-1})
$$
$$
= \mathbf{E}[x_t - x_{t-1}] - \alpha \nabla f(x_t)
$$
$$
+ \left( 1 - \frac{1-p}{C} \right) \mathbf{E}[\alpha_t \nabla f(v_{t-1})]
$$
$$
= \mathbf{E}[x_t - x_{t-1}] - \alpha \nabla f(x_t) + \left( 1 - \frac{1-p}{C} \right) \mathbf{E}[x_t - x_{t-1}]
$$
$$
= \left( 2 - \frac{1-p}{C} \right) \mathbf{E}[x_t - x_{t-1}] - \alpha \nabla f(x_t)
$$
$\square$

*Proof of Theorem 4.* We have
$$
\Sigma_{p,\alpha}^{\nabla} = \sum_{i=0}^{\infty} \big( p(i)\alpha(i) - p(i+1)\alpha(i+1) \big) \nabla f(x_{t-i-1})
$$

Substituting $p(i)$ for the $CMP$ PDF (12) gives
$$
\Sigma_{p,\alpha}^{\nabla} = \frac{1}{Z(\lambda,\nu)} \sum_{i=0}^{\infty} \frac{\lambda^i}{(i!)^{\nu}} \left( \alpha(i) - \lambda \frac{\alpha(i+1)}{(i+1)^{\nu}} \right) \nabla f(x_{t-i-1})
$$
(27)

Now, applying the adaptive step size (14) gives
$$
\Sigma_{p,\alpha}^{\nabla} = \frac{C}{Z(\lambda,\nu)} \sum_{i=0}^{\infty} \frac{\lambda^i}{(i!)^{\nu}} \alpha \Big( \lambda^{-i}(i!)^{\nu} -
$$
$$
\frac{\lambda}{(i+1)^{\nu}} \lambda^{-(i+1)}((i+1)!)^{\nu} \Big) \nabla f(x_{t-i-1})
$$
$$
= \frac{C}{Z(\lambda,\nu)} \sum_{i=0}^{\infty} \frac{\lambda^i}{(i!)^{\nu}} \alpha \left( \frac{(i!)^{\nu}}{\lambda^i} - \frac{(i!)^{\nu}}{\lambda^i} \right) \nabla f(x_{t-i-1}) = 0
$$
$\square$

*Proof of Theorem 5.* Let $\Psi(i) = \alpha(i) - \lambda \frac{\alpha(i+1)}{(i+1)^{\nu}}$, and hence
$$
\Sigma_{p,\alpha}^{\nabla} = \frac{1}{Z(\lambda,\nu)} \sum_{i=0}^{\infty} \frac{\lambda^i}{(i!)^{\nu}} \Psi(i) \nabla f(x_{t-i-1})
$$

Applying the adaptive step size (15) gives
$$
\Psi(i) = \frac{i!^{\nu}}{\lambda^i} e^{\lambda} \alpha \left( c(i) - c(i+1) \right)
$$

Now,
$$
\Psi(i) = K \Leftrightarrow c(i) - c(i+1) = \frac{K}{\alpha e^{\lambda}} \frac{\lambda^i}{i!^{\nu}}
$$
$$
\Leftrightarrow c(i) = c(i-1) - \frac{K}{\alpha e^{\lambda}} \frac{\lambda^{i-1}}{(i-1)!^{\nu}}
$$
$$
= c(0) - \frac{K}{\alpha e^{\lambda}} \sum_{j=1}^{i} \frac{\lambda^{i-j}}{(i-j)!^{\nu}} = c(0) - \frac{K}{\alpha e^{\lambda}} \sum_{j=1}^{i} \frac{\lambda^j}{(j)!^{\nu}}
$$

Since $\alpha(0) = \alpha$, we have $c(0) = 1$. Now we have
$$
\Sigma_{p,\alpha}^{\nabla} = K \sum_{i=0}^{\infty} \frac{1}{Z(\lambda,\nu)} \frac{\lambda^i}{(i!)^{\nu}} \nabla f(x_{t-i-1})
$$
$$
= K\mathbf{E}\left[ \nabla f(v_{t-1}) \right] = K\mathbf{E}\left[ x_t - x_{t-1} \right]
$$
$\square$

*Proof of Corollary 2.* Under the Poisson $\tau$ model, which is $CMP$ with $\nu = 1$, (16) rewrites to
$$
c(i) = 1 - \frac{K}{\alpha e^{\lambda}} \sum_{j=0}^{\tau-1} \frac{\lambda^j}{(j!)} = 1 - \frac{K}{\alpha} \frac{\Gamma(i,\lambda)}{(i-1)!}
$$
$$
= 1 - \frac{K}{\alpha} \frac{\Gamma(i,\lambda)}{\Gamma(i)}
$$
$\square$

*Proof of Theorem 6.*
$$
\|x_{t+1} - x^*\|^2 = \|x_t - \alpha_t \nabla F(v_t) - x^*\|^2
$$
$$
= \|x_t - x^*\|^2 + \alpha_t^2 \|\nabla F(v_t)\|^2 - 2\alpha_t (x_t - x^*)^T \nabla F(v_t)
$$
$$
= \|x_t - x^*\|^2 + \alpha_t^2 \|\nabla F(v_t)\|^2 - 2\alpha_t (x_t - x^*)^T \nabla F(x_t)
$$
$$
+ 2\alpha_t (x_t - x^*)^T \big( \nabla F(x_t) - \nabla F(v_t) \big)
$$

Under expectation, conditioned on the natural filtration $\mathbb{F}_t^X = \big( (\tau_i)_{i=0}^t, \big( \nabla F(v_i) \big)_{i=0}^t \big)$ of the *past* of the process, we have
$$
\mathbf{E}\big[ \|x_{t+1} - x^*\|^2 \mid \tau_t, \mathbb{F}_{t-1}^X \big] = \|x_t - x^*\|^2
$$
$$
- 2\alpha_t \mathbf{E}\big[ (x_t - x^*)^T \big( \nabla F(x_t) - \nabla F(x^*) \big) |\mathbb{F}_{t-1}^X \big]
$$
$$
+ 2\alpha_t \mathbf{E}\big[ (x_t - x^*)^T \big( \nabla F(x_t) - \nabla F(v_t) \big) |\mathbb{F}_{t-1}^X \big]
$$

Applying the assumptions (19)-(21) gives
$$
\mathbf{E}\big[ \|x_{t+1} - x^*\|^2 | \tau_t, \mathbb{F}_{t-1}^X \big] \leq \|x_t - x^*\|^2 + M^2 \alpha_t^2
$$
$$
- 2\alpha_t c \|x_t - x^*\|^2 + 2\alpha_t L \|x_t - x^*\| \|x_t - v_t\|
$$
$$
= (1 - 2c\alpha_t) \|x_t - x^*\|^2 + M^2 \alpha_t^2
$$
$$
+ 2\alpha_t L \|x_t - x^*\| \|x_t - v_t\|
$$
$$
= (1 - 2c\alpha_t) \|x_t - x^*\|^2 + M^2 \alpha_t^2
$$
$$
+ 2\alpha_t L \|x_t - x^*\| \sum_{i=1}^{\tau_t} x_{t-i+1} - x_{t-i}\|
$$
$$
\leq (1 - 2c\alpha_t) \|x_t - x^*\|^2 + M^2 \alpha_t^2
$$
$$
+ 2\alpha_t L \sum_{i=1}^{\tau_t} \|x_t - x^*\| \alpha_{t-i} \|\nabla F(v_{t-i})\|
$$

The gradient process does not influence the expected delays, so we first consider the expectation conditioned on the gradient process $(\nabla)_0^t := \big( \nabla F(v_i) \big)_{i=0}^t$
$$
\mathbf{E}\big[ \|x_{t+1} - x^*\|^2 | \tau_t, (\nabla)_0^t \big]
$$
$$
\leq (1 - 2c\alpha_t)\mathbf{E}\big[ \|x_t - x^*\|^2 | \tau_t, (\nabla)_0^t \big] + M^2 \alpha_t^2
$$
$$
+ 2L\alpha_t \sum_{i=1}^{\tau_t} \mathbf{E}\big[ \alpha_{t-i} \|x_t - x^*\| | \tau_t, (\nabla)_0^t \big] \|\nabla F(v_{t-i})\|
$$

From the non-anticipativity of the delay process we have
$$
\mathbf{E}\big[ \alpha_{t-i} \|x_t - x^*\| \mid \tau_t, (\nabla)_0^t \big]
$$
$$
= \mathbf{E}\big[ \mathbf{E}\big[ \alpha_{t-i} \|x_t - x^*\| \mid x_t \big] \mid \tau_t, (\nabla)_0^t \big]
$$
$$
= \mathbf{E}\big[ \|x_t - x^*\| \mathbf{E}\big[ \alpha_{t-i} \mid x_t \big] \mid \tau_t, (\nabla)_0^t \big]
$$
$$
= \mathbf{E}\big[ \alpha_{t-i} \big] \mathbf{E}\big[ \|x_t - x^*\| \mid (\nabla)_0^t \big]
$$

Since the delays and gradients are identically distributed we have $\mathbf{E}[\alpha_i] = \mathbf{E}[\alpha_j]$ for all $i,j$. Taking expectation conditioned on the last delay $\tau_t$ and applying Hölder's inequality gives

$$\mathbf{E}\left[\|x_{t+1} - x^*\|^2 \mid \tau_t\right] \leq (1 - 2c\alpha_t)\mathbf{E}\left[\|x_t - x^*\|^2\right] + M^2\alpha_t^2$$
$$+ 2L\tau_t\alpha_t\mathbf{E}[\alpha_t]\sqrt{\mathbf{E}\left[\|x_t - x^*\|^2\right]}\sqrt{\mathbf{E}\left[\|\nabla F(v_t)\|^2\right]}$$

and the full expectation satisfies

$$\mathbf{E}\left[\|x_{t+1} - x^*\|^2\right] \leq (1 - 2c\mathbf{E}[\alpha_t])\mathbf{E}\left[\|x_t - x^*\|^2\right]$$
$$+ M^2\mathbf{E}\left[\alpha_t^2\right] + 2LM\mathbf{E}[\tau_t\alpha_t]\mathbf{E}[\alpha_t]\sqrt{\mathbf{E}\left[\|x_t - x^*\|^2\right]}$$

As long as the process has not yet converged, i.e. $\mathbf{E}\left[\|x_t - x^*\|^2\right] > \epsilon$, we have

$$\mathbf{E}\left[\|x_{t+1} - x^*\|^2\right] \leq \mathbf{E}\left[\|x_t - x^*\|^2\right](1 - 2c\mathbf{E}[\alpha_t]$$
$$+ \epsilon^{-1}M^2\mathbf{E}\left[\alpha_t^2\right] + 2LM\epsilon^{1/2}\mathbf{E}[\tau_t\alpha_t]\mathbf{E}[\alpha_t])$$
$$=: \mathbf{E}\left[\|x_t - x^*\|^2\right](1 - \delta)$$
$$\Rightarrow \mathbf{E}\left[\|x_T - x^*\|^2\right] \leq \mathbf{E}\left[\|x_0 - x^*\|^2\right](1 - \delta)^T$$
$$\Rightarrow T \leq -\ln(1 - \delta)^{-1}\ln\frac{\mathbf{E}\left[\|x_0 - x^*\|^2\right]}{\mathbf{E}\left[\|x_T - x^*\|^2\right]}$$
$$< \delta^{-1}\ln\left(\mathbf{E}\left[\|x_0 - x^*\|^2\right]\epsilon^{-1}\right)$$

for any $T$ such that $\mathbf{E}\left[\|x_T - x^*\|^2\right] > \epsilon$. Equivalently, expected convergence is implied by $T$ exceeding the bound above, which concludes the proof. $\square$

*Proof of Corollary 3.* Let $\rho = \frac{c\epsilon M^{-1}}{M + 2L\sqrt{\epsilon}\bar{\tau}}$. From Theorem 6 we have the improvement factor

$$\delta = 2\left(c - LM\epsilon^{1/2}\mathbf{E}[\tau\alpha]\right)\mathbf{E}[\alpha] - \epsilon^{-1}M^2\mathbf{E}[\alpha^2]$$
$$= 2c\alpha - \epsilon^{-1}M\left(M + 2L\sqrt{\epsilon}\bar{\tau}\right)\alpha^2$$
$$= c\rho^{-1}\alpha(2\rho - \alpha)$$

so $\delta > 0$ when $0 < \alpha < 2\rho$, and the improvement is maximized for $\theta = 1$. Now, using the choice (23) of step size, we have

$$\delta = c\rho^{-1}\theta\rho(2\rho - \theta\rho)$$
$$= \theta(2 - \theta)c\rho$$

Substituting for $\rho$, the convergence bound of Theorem 6 now rewrites to (24) $\square$

*Proof of Corollary 4.* Since $\alpha_t$ is a non-increasing function in $\tau_t$ we have

$$\mathbf{E}[\tau_t\alpha(\tau_t)] = \mathbf{E}[\tau_t\alpha(\tau_t)] - \mathbf{E}[\bar{\tau}\alpha(\tau_t)] + \mathbf{E}[\tau_t]\mathbf{E}[\alpha(\tau_t)]$$
$$= \mathbf{E}[(\tau_t - \bar{\tau})(\alpha(\tau_t) - \alpha(\bar{\tau}))] + \mathbf{E}[\tau]\mathbf{E}[\alpha]$$
$$\leq \mathbf{E}[\tau]\mathbf{E}[\alpha]$$

Using this property, (22) rewrites to (25) $\square$