# ADAPTIVE COMMUNICATION STRATEGIES TO ACHIEVE THE BEST ERROR-RUNTIME TRADE-OFF IN LOCAL-UPDATE SGD

Jianyu Wang <sup>1</sup> Gauri Joshi <sup>1</sup>

#### **ABSTRACT**

Large-scale machine learning training, in particular, distributed stochastic gradient descent, needs to be robust to inherent system variability such as node straggling and random communication delays. This work considers a distributed training framework where each worker node is allowed to perform local model updates and the resulting models are averaged periodically. We analyze the true speed of error convergence with respect to wall-clock time (instead of the number of iterations), and analyze how it is affected by the frequency of averaging. The main contribution is the design of ADACOMM, an adaptive communication strategy that starts with infrequent averaging to save communication delay and improve convergence speed, and then increases the communication frequency in order to achieve a low error floor. Rigorous experiments on training deep neural networks show that ADACOMM can take  $3\times$  less time than fully synchronous SGD and still reach the same final training loss.

#### 1 Introduction

Stochastic gradient descent (SGD) is the backbone of stateof-the-art supervised learning, which is revolutionizing inference and decision-making in many diverse applications. Classical SGD was designed to be run on a single computing node, and its error-convergence with respect to the number of iterations has been extensively analyzed and improved via accelerated SGD methods. Due to the massive training data-sets and neural network architectures used today, it has became imperative to design distributed SGD implementations, where gradient computation and aggregation is parallelized across multiple worker nodes. Although parallelism boosts the amount of data processed per iteration, it exposes SGD to unpredictable node slowdown and communication delays stemming from variability in the computing infrastructure. Thus, there is a critical need to make distributed SGD fast, yet robust to system variability.

Need to Optimize Convergence in terms of Error versus Wall-clock Time. The convergence speed of distributed SGD is a product of two factors: 1) the error in the trained model versus the number of iterations, and 2) the number of iterations completed per second. Traditional single-node SGD analysis focuses on optimizing the first factor, because the second factor is generally a constant when SGD

Proceedings of the  $2^{nd}$  SysML Conference, Palo Alto, CA, USA, 2019. Copyright 2019 by the author(s).

is run on a single dedicated server. In distributed SGD, which is often run on shared cloud infrastructure, the second factor depends on several aspects such as the number of worker nodes, their local computation and communication delays, and the protocol (synchronous, asynchronous or periodic) used to aggregate their gradients. Hence, in order to achieve the fastest convergence speed we need: 1) optimization techniques (eg. variable learning rate) to maximize the error-convergence rate with respect to iterations, and 2) scheduling techniques (eg. straggler mitigation, infrequent communication) to maximize the number of iterations completed per second. These directions are inter-dependent and need to be explored together rather than in isolation. While many works have advanced the first direction, the second is less explored from a theoretical point of view, and the juxtaposition of both is an unexplored problem.

#### Local-Update SGD to Reduce Communication Delays.

A popular distributed SGD implementation is the parameter server framework (Dean et al., 2012; Cui et al., 2014; Li et al., 2014; Gupta et al., 2016; Mitliagkas et al., 2016) where in each iteration, worker nodes compute gradients on one mini-batch of data and a central parameter server aggregates these gradients (synchronously or asynchronously) and updates the parameter vector x. The constant communication between the parameter server and worker nodes in each iteration can be expensive and slow in bandwidth-limited computed environments. Recently proposed distributed SGD frameworks such as Elastic-averaging (Zhang et al., 2015; Chaudhari et al., 2017), Federated Learning (McMahan et al., 2016; Smith et al., 2017b) and decentralized SGD (Lian et al., 2017; Jiang et al., 2017) save this

<sup>&</sup>lt;sup>1</sup>Department of Electrical & Computer Engineering, Carnegie Mellon University, Pittsburgh, PA, USA. Correspondence to: Jianyu Wang <jianyuw1@andrew.cmu.edu>, Gauri Joshi <gaurij@andrew.cmu.edu>.

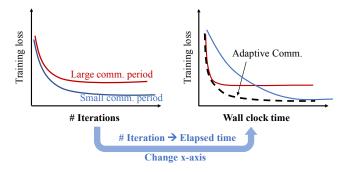


Figure 1. This work departs from the traditional view of considering error-convergence with respect to the number of iterations, and instead considers the true convergence in terms of error versus wall-clock time. Adaptive strategies that start with infrequent model-averaging and increase the communication frequency can achieve the best error-runtime trade-off.

communication cost by allowing worker nodes to perform local updates to the parameter x instead of just computing gradients. The resulting locally trained models (which are different due to variability in training data across nodes) are periodically averaged through a central server, or via direct inter-worker communication. This local-update strategy has been shown to offer significant speedup in deep neural network training (Lian et al., 2017; McMahan et al., 2016).

Error-Runtime Trade-offs in Local-Update SGD. While local updates reduce the communication-delay incurred per iteration, discrepancies between local models can result in an inferior error-convergence. For example, consider the case of periodic-averaging SGD (PASGD) where each of m worker nodes makes  $\tau$  local updates, and the resulting models are averaged after every  $\tau$  iterations (Moritz et al., 2015; Su & Chen, 2015; Chen & Huo, 2016; Seide & Agarwal, 2016; Zhang et al., 2016; Zhou & Cong, 2017; Lin et al., 2018). A larger value of  $\tau$  leads to slower convergence with respect to the number of iterations as illustrated in Figure 1. However, if we look at the true convergence with respect to the wall-clock time, then a larger  $\tau$ , that is, less frequent averaging, saves communication delay and reduces the runtime per iteration. While some recent theoretical works (Zhou & Cong, 2017; Yu et al., 2018; Wang & Joshi, 2018; Stich, 2018) study this dependence of the error-convergence with respect to the number of iterations as  $\tau$  varies, achieving a provably-optimal speed-up in the true convergence with respect to wall-clock time is an open problem that we aim to address in this work.

Need for Adaptive Communication Strategies. In the error-runtime in Figure 1, we observe a trade-off between the convergence speed and the error floor when the number of local updates  $\tau$  is varied. A larger  $\tau$  gives a faster initial drop in the training loss but results in a higher error floor. This calls for *adaptive communication strategies* that start

with a larger  $\tau$  and gradually decrease it as the model reaches closer to convergence. Such an adaptive strategy will offer a win-win in the error-runtime trade-off by achieving fast convergence as well as low error floor. To the best of our knowledge, this is the first work to propose an adaptive communication frequency strategy.

**Main Contributions.** This paper focuses on periodic-averaging local-update SGD (PASGD) and makes the following main contributions:

- We first analyze the runtime per iteration of periodic averaging SGD (PASGD) by modeling local computing time and communication delays as random variables, and quantify its runtime speed-up over fully synchronous SGD. A novel insight from this analysis is that periodic-averaging strategy not only reduces the communication delay but also mitigates synchronization delays in waiting for slow or straggling nodes.
- 2. By combining the runtime analysis error-convergence analysis of PASGD (Wang & Joshi, 2018), we can obtain the error-runtime trade-off for different values of  $\tau$ . Using this combined error-runtime trade-off, we derive an expression of the optimal communication period, which can serve as a useful guideline in practice.
- 3. Based on the observations in runtime and convergence analysis, we develop an adaptive communication scheme: ADACOMM. Experiments on training VGG-16 and ResNet-50 deep neural networks and different settings (with/without momentum, fixed/decaying learning rate) show that ADACOMM can give a 3× runtime speed-up and still reach the same low training loss as fully synchronous SGD.
- 4. We present a convergence analysis for PASGD with variable communication period  $\tau$  and variable learning rate  $\eta$ , generalizing previous work (Wang & Joshi, 2018). This analysis shows that decaying  $\tau$  provides similar convergence benefits as decaying learning rate, the difference being that varying  $\tau$  improves the true convergence with respect to the wall-clock time. Adaptive communication can also be used in conjunction with existing learning rate schedules.

Although we focus on periodic simple-averaging of local models, the insights on error-runtime trade-offs and adaptive communication strategies are directly extendable to other communication-efficient SGD algorithms including Federated Learning (McMahan et al., 2016), Elastic-Averaging (Zhang et al., 2015) and Decentralized averaging (Jiang et al., 2017; Lian et al., 2017), as well as synchronous/asynchronous distributed SGD with a central parameter server (Dean et al., 2012; Cui et al., 2014; Dutta et al., 2018).

#### 2 PROBLEM FRAMEWORK

Empirical Risk Minimization via Mini-batch SGD. Our objective is to minimize an objective function  $F(\mathbf{x})$ , the empirical risk function, with respect to model parameters denoted by  $\mathbf{x} \in \mathbb{R}^d$ . The training dataset is denoted by  $S = \{s_1, \dots, s_N\}$ , where  $s_i$  represents the i-th labeled data point. The objective function can be expressed as the empirical risk calculated using the training data and is given by

$$\min_{\mathbf{x} \in \mathbb{R}^d} \left[ F(\mathbf{x}) := \frac{1}{N} \sum_{i=1}^N f(\mathbf{x}; s_i) \right]$$
 (1)

where  $f(\mathbf{x}; s_i)$  is the composite loss function at the  $i^{th}$  data point. In classic mini-batch stochastic gradient descent (SGD) (Dekel et al., 2012), updates to the parameter vector  $\mathbf{x}$  are performed as follows. If  $\xi_k \subset S$  represents a randomly sampled mini-batch, then the update rule is

$$\mathbf{x}_{k+1} = \mathbf{x}_k - \eta g(\mathbf{x}_k; \xi_k) \tag{2}$$

where  $\eta$  denotes the learning rate and the stochastic gradient is defined as:  $g(\mathbf{x};\xi) = \frac{1}{|\xi|} \sum_{s_i \in \xi} \nabla f(\mathbf{x};s_i)$ . For simplicity, we will use  $g(\mathbf{x}_k)$  instead of  $g(\mathbf{x}_k;\xi_k)$  in the rest of the paper. A complete review of convergence properties of serial SGD can be found in (Bottou et al., 2018).

**Periodic-Averaging SGD (PASGD).** We consider a distributed SGD framework with m worker nodes where all workers can communicate with others via a central server or via direct inter-worker communication. In periodic-averaging SGD, all workers start at the same initial point  $\mathbf{x}_1$ . Each worker performs  $\tau$  local mini-batch SGD updates according to (2), and the local models are averaged by a fusion node or by performing an all-node broadcast. The workers then update their local models with the averaged model, as illustrated in Figure 2. Thus, the overall update rule at the  $i^{th}$  worker is given by

$$\mathbf{x}_{k+1}^{(i)} = \begin{cases} \frac{1}{m} \sum_{j=1}^{m} [\mathbf{x}_k^{(j)} - \eta g(\mathbf{x}_k^{(j)})], & k \mod \tau = 0\\ \mathbf{x}_k^{(i)} - \eta g(\mathbf{x}_k^{(i)}), & \text{otherwise} \end{cases}$$
(3)

where  $\mathbf{x}_k^{(i)}$  denote the model parameters in the *i*-th worker after k iterations and  $\tau$  is defined as the communication period. Note that the iteration index k corresponds to the local iterations, and not the number of averaging steps.

**Special Case** ( $\tau = 1$ ): **Fully Synchronous SGD.** When  $\tau = 1$ , that is, the local models are synchronized after every iteration, periodic-averaging SGD is equivalent to fully synchronous SGD which has the update rule

$$\mathbf{x}_{k+1} = \mathbf{x}_k - \eta \left[ \frac{1}{m} \sum_{i=1}^m g(\mathbf{x}_k; \xi_k^{(i)}) \right]. \tag{4}$$

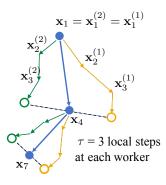


Figure 2. Illustration of PASGD in the model parameter space for m=2 workers. The discrepancy between the local models increases with the number of local updates,  $\tau=3$ .



Figure 3. Illustration of PASGD in the time space for m=2 and  $\tau=3$ . Lengths of the colored arrows at the  $i^{th}$  worker are  $Y_{i,k}$ , the local-update times, which are i.i.d. across workers and updates. The blue block represents the communication delay for each model-averaging step.

The analysis of fully synchronous SGD is identical to serial SGD with m-fold large mini-batch size.

Local Computation Times and Communication Delay. In order to analyze the effect of  $\tau$  on the expected runtime per iteration, we consider the following delay model. The time taken by the  $i^{th}$  worker to compute a mini-batch gradient at the  $k^{th}$  local-step is modeled a random variable  $Y_{i,k} \sim F_Y$ , assumed to be i.i.d. across workers and minibatches. The communication delay is a random variable D for each all-node broadcast, as illustrated in Figure 3. The value of random variable D can depend on the number of workers as follows.

$$D = D_0 \cdot s(m) \tag{5}$$

where  $D_0$  represents the time taken for each inter-node communication, and s(m) describes how the delay scales with the number of workers, which depends on the implementation and system characteristics. For example, in the parameter server framework, the communication delay can be proportional to  $2\log_2(m)$  by exploiting a reduction tree structure (Iandola et al., 2016). We assume that s(m) is known beforehand for the communication-efficient distributed SGD framework under consideration.

Convergence Criteria. In the error-convergence analysis, since the objective function is non-convex, we use the expected gradient norm as a an indicator of convergence following (Ghadimi & Lan, 2013; Bottou et al., 2018). We say

the algorithm achieves an  $\epsilon$ -suboptimal solution if:

$$\mathbb{E}\left[\min_{k\in[1,K]}\left\|\nabla F(\mathbf{x}_{k})\right\|^{2}\right] \leq \epsilon. \tag{6}$$

When  $\epsilon$  is arbitrarily small, this condition can guarantee the algorithm converges to a stationary point.

## 3 JOINTLY ANALYZING RUNTIME AND ERROR-CONVERGENCE

#### 3.1 Runtime Analysis

We now present a comparison of the runtime per iteration of periodic-averaging SGD with fully synchronous SGD to illustrate how increasing  $\tau$  can lead to a large runtime speedup. Another interesting effect of performing more local update  $\tau$  is that it mitigates the slowdown due to straggling worker nodes.

Runtime Per Iteration of Fully Synchronous SGD. Fully synchronous SGD is equivalent to periodic-averaging SGD with  $\tau=1$ . Each of the m workers computes the gradient of one mini-batch and updates the parameter vector  $\mathbf{x}$ , which takes time  $Y_{i,1}$  at the  $i^{th}$  worker<sup>1</sup>. After all workers finish their local updates, an all-node broadcast is performed to synchronize and average the models. Thus, the total time to complete each iteration is given by

$$T_{\text{sync}} = \max(Y_{1,1}, Y_{2,1}, \dots, Y_{m,1}) + D$$
 (7)

$$\mathbb{E}\left[T_{\text{sync}}\right] = \mathbb{E}[Y_{m:m}] + \mathbb{E}[D] \tag{8}$$

where  $Y_{i,1}$  are i.i.d. random variables with probability distribution  $F_Y$  and D is the communication delay. The term  $Y_{m:m}$  denotes the highest order statistic of m i.i.d. random variables (David & Nagaraja, 2003).

Runtime Per Iteration of Periodic-Averaging SGD (PASGD). In periodic-averaging SGD, each worker performs  $\tau$  local updates before communicating with other workers. Let us denote the average local computation time at the  $i^{th}$  worker by

$$\overline{Y}_i = \frac{Y_{i,1} + Y_{i,2} + \dots Y_{i,\tau}}{\tau} \tag{9}$$

Since the communication delay D is amortized over  $\tau$  iterations, the average computation time per iteration is

$$T_{\text{P-Avg}} = \max(\overline{Y}_1, \overline{Y}_2, \dots, \overline{Y}_m) + \frac{D}{\tau}$$
 (10)

$$\mathbb{E}[T_{\text{P-Avg}}] = \mathbb{E}[\overline{Y}_{m:m}] + \frac{\mathbb{E}[D]}{\tau}$$
(11)

The value of the first term  $\overline{Y}_{m:m}$  and how it compares with  $Y_{m:m}$  depends on the probability distribution  $F_Y$  of Y.

#### 3.2 Runtime Benefits of Periodic Averaging Strategy

**Speed-up over fully synchronous SGD.** We evaluate the speed-up of periodic-averaging SGD over fully synchronous SGD for different Y and D to demonstrate how the relative value of computation versus communication delays affects the speed-up. Consider the simplest case where Y and D are constants and define  $\alpha = D/Y$ , the communication/computation ratio. Besides systems aspects such as network bandwidth and computing capacity, for deep neural network training, this ratio  $\alpha$  also depends on the size of the neural network model and the mini-batch size. See Figure 8 for a comparison of the communication/computation delays of common deep neural network architectures. Then  $\overline{Y}$ ,  $Y_{m:m}$ ,  $\overline{Y}_{m:m}$  are all equal to Y, and the ratio of  $\mathbb{E}[T_{\text{Sync}}]$  and  $\mathbb{E}[T_{\text{P-Ayg}}]$  is given by

$$\frac{\mathbb{E}\left[T_{\text{sync}}\right]}{\mathbb{E}\left[T_{\text{P-Avg}}\right]} = \frac{Y+D}{Y+D/\tau} = \frac{1+\alpha}{1+\alpha/\tau}$$
 (12)

Figure 4 shows the speed-up for different values of  $\alpha$  and  $\tau$ . When D is comparable with Y ( $\alpha = 0.9$ ), periodicaveraging SGD (PASGD) can be almost twice as fast as fully synchronous SGD.

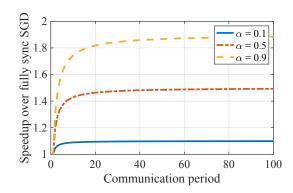


Figure 4. The speed-up offered by using periodic-averaging SGD increases with  $\tau$  (the communication period) and with the communication/computation delay ratio  $\alpha=D/Y$ , where D is the all-node broadcast delay and Y is the time taken for each local update at a worker.

Straggler Mitigation due to Local Updates. Suppose that Y is exponentially distributed with mean y and variance  $y^2$ . For fully synchronous SGD, the term  $\mathbb{E}[Y_{m:m}]$  in (8) is equal to  $y \sum_{i=1}^m 1/i$ , which is approximately equal to  $y \log m$ . Thus, the expected runtime per iteration of fully synchronous SGD (8) increases logarithmically with the number of workers m. Let us compare this with the scaling of the runtime of periodic-averaging SGD (11). Here,  $\overline{Y}$  (9) is an Erlang random variable with mean y and variable  $y^2/\tau$ . Since the variance is  $\tau$  times smaller than that of Y, the maximum order statistic  $\mathbb{E}[\overline{Y}_{m:m}]$  is smaller than

<sup>&</sup>lt;sup>1</sup>Instead of local updates, typical implementations of fully synchronous SGD have a central server that performs the update. Here we compare PASGD with fully synchronous SGD without a central parameter server.

 $\mathbb{E}[Y_{m:m}]$ . Figure 5 shows the probability distribution of  $T_{\text{sync}}$  and  $T_{\text{P-Avg}}$  for exponentially distributed Y. Observe that  $T_{\text{P-Avg}}$  has a much lighter tail. This is because the effect of the variability in Y on  $T_{\text{P-Avg}}$  is reduced due to the Y in (8) being replaced by  $\overline{Y}$  (which has lower variance) in (11).

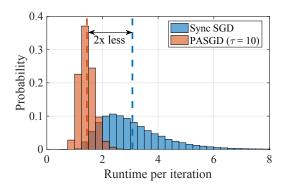


Figure 5. Probability distribution of runtime per iteration, where communication delay D=1, mean computation time y=1, and number of workers m=16. Dash lines represent the mean values.

#### 3.3 Joint Analysis with Error-convergence

In this subsection, we combine the runtime analysis with previous error-convergence analysis for PASGD (Wang & Joshi, 2018). Due to space limitations, we state the necessary theoretical assumptions in the Appendix; the assumptions are similar to previous works (Zhou & Cong, 2017; Wang & Joshi, 2018) on the convergence of local-update SGD algorithms.

**Theorem 1** (Error-runtime Convergence of PASGD). For PASGD, under certain assumptions (stated in the Appendix), if the learning rate satisfies  $\eta L + \eta^2 L^2 \tau(\tau - 1) \leq 1$ , Y and D are constants, and all workers are initialized at the same point  $\mathbf{x}_1$ , then after total T wall-clock time, the minimal expected squared gradient norm within T time interval will be bounded by:

$$\frac{2\left[F(\mathbf{x}_1) - F_{inf}\right]}{\eta T} \left(Y + \frac{D}{\tau}\right) + \frac{\eta L \sigma^2}{m} + \eta^2 L^2 \sigma^2 (\tau - 1)$$
(13)

where L is the Lipschitz constant of the objective function and  $\sigma^2$  is the variance bound of mini-batch stochastic gradients.

The proof of Theorem 1 is presented in the Appendix. From the optimization error upper bound (13), one can easily observe the error-runtime trade-off for different communication periods. While a larger  $\tau$  reduces the runtime per iteration and let the first term in (13) become smaller, it also adds additional noise and increases the last term. In Figure 6,

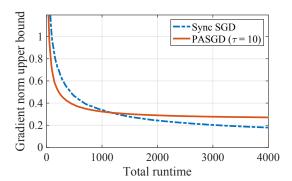


Figure 6. Illustration of theoretical error bound versus runtime in Theorem 1. The runtime per iteration is generated under the same parameters as Figure 5. Other constants in (13) are set as follows:  $F(\mathbf{x}_1) = 1$ ,  $F_{\text{inf}} = 0$ ,  $\eta = 0.08$ , L = 1,  $\sigma^2 = 1$ .

we plot theoretical bounds for both fully synchronous SGD  $(\tau=1)$  and PASGD. It is shown that although PASGD with  $\tau=10$  starts with a rapid drop, it will eventually converge to a high error floor. This theoretical result is also corroborated by experiments in Section 5. Another direct outcome of Theorem 1 is the determination of the best communication period that balances the first and last terms in (13). We will discuss the selection of communication period later in Section 4.1.

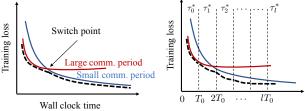
## 4 ADACOMM: PROPOSED ADAPTIVE COMMUNICATION STRATEGY

Inspired by the clear trade-off in the learning curve in Figure 6, it would be better to have an adaptive communication strategy that starts with infrequent communication to improve convergence speed, and then increases the frequency to achieve a low error floor. In this section, we are going to develop the proposed adaptive communication scheme.

The basic idea to adapt the communication is to choose the communication period that minimizes the optimization error at each wall-clock time. One way to achieve the idea is switching between the learning curves at their intersections. However, without prior knowledge of various curves, it would be difficult to determine the switch points.

Instead, we divide the whole training procedure into uniform wall-clock time intervals with the same length  $T_0$ . At the beginning of each time interval, we select the best value of  $\tau$  that has the fastest decay rate in the next  $T_0$  wall-clock time. If the interval length  $T_0$  is small enough and the best choice of communication period for each interval can be precisely estimated, then this adaptive scheme should achieve a winwin in the error-runtime trade-off as illustrated in Figure 7.

After setting the interval length, the next question is how



- (a) Switch between curves.
- (b) Choose the best  $\tau$  for each time interval.

Figure 7. Illustration of communication period adaptation strategies. Dash line denotes the learning curve using adaptive communication.

to estimate the best communication period for each time interval. In Section 4.1 we use the error-runtime analysis in Section 3.3 to find the best  $\tau$  at each time.

#### 4.1 Determining the Best Communication Period for Each Time Interval

From Theorem 1, it can be observed that there is an optimal value  $\tau^*$  that minimizes the optimization error bound at given wall-clock time. In particular, consider the simplest setting where Y and D are constants. Then, by minimizing the upper bound (13) over  $\tau$ , we obtain the following.

**Theorem 2.** For PASGD, under the same assumptions as Theorem 1, the optimization error upper bound in (13) at time T is minimized when the communication period is

$$\tau^* = \sqrt{\frac{2(F(\mathbf{x}_1) - F_{inf})D}{\eta^3 L^2 \sigma^2 T}}.$$
 (14)

The proof is straightforward by setting the derivative of (13) to zero. We present the details in the Appendix. Suppose all workers starts from the same initial point  $\mathbf{x}_1 = \mathbf{x}_{t=0}$  where subscript t denotes the wall-clock time. Directly applying Theorem 2 to the first time interval, then the best choice of communication period is:

$$\tau_0 = \sqrt{\frac{2(F(\mathbf{x}_{t=0}) - F_{\text{inf}})D}{\eta^3 L^2 \sigma^2 T_0}}.$$
 (15)

Similarly, for the l-th time interval, workers can be viewed as restarting training at a new initial point  $\mathbf{x}_{t=lT_0}$ . Applying Theorem 2 again, we have

$$\tau_{l} = \sqrt{\frac{2(F(\mathbf{x}_{t=lT_{0}}) - F_{\text{inf}})D}{\eta^{3}L^{2}\sigma^{2}T_{0}}}.$$
 (16)

Comparing (15) and (16), it is easy to see the generated communication period sequence decreases along with the objective value  $F(\mathbf{x}_t)$  when the learning rate is fixed. This

result is consistent with the intuition that the trade-off between error-convergence and communication-efficiency varies over time. Compared to the initial phase of training, the benefit of using a large communication period diminishes as the model reaches close to convergence. At this later stage, a lower error floor is more preferable to speeding up the runtime.

Remark 1 (Connection to Decaying Learning Rate). Using a fixed learning rate in SGD leads to an error floor at convergence. To further reduce the error, practical SGD implementations generally decay the learning rate or increase the mini-batch size (Smith et al., 2017a; Goyal et al., 2017). As we saw from the convergence analysis Theorem 1, performing local updates adds additional noise in stochastic gradients, resulting in a higher error floor convergence. Decaying the communication period can gradually reduce the variance of gradients and yield a similar improvement in convergence. Thus, adaptive communication strategies are similar in spirit to decaying learning rate or increasing minibatch size. The key difference is that here we are optimizing the true error convergence with respect to wall-clock time rather than the number iterations.

#### 4.2 Practical Considerations

Although (15) and (16) provide useful insights about how to adapt  $\tau$  over time, it is still difficult to directly use them in practice due to the Lipschitz constant L and the gradient variance bound  $\sigma^2$  being unknown. For deep neural networks, estimating these constants can be difficult and unreliable due to the highly non-convex and high-dimensional loss surface. As an alternative, we propose a simpler rule where we approximate  $F_{\rm inf}$  by 0, and divide (16) by (15) to obtain the basic communication period update rule:

Basic update rule 
$$\tau_l = \left[ \sqrt{\frac{F(\mathbf{x}_{t=lT_0})}{F(\mathbf{x}_{t=0})}} \tau_0 \right]$$
 (17)

where  $\lceil a \rceil$  is the ceil function to round a to the nearest integer  $\geq a$ . Since the objective function values (i.e., training loss)  $F(\mathbf{x}_{t=lT_0})$  and  $F(\mathbf{x}_{t=0})$  can be easily obtained in the training, the only remaining thing now is to determine the initial communication period  $\tau_0$ . We obtain a heuristic estimate of  $\tau_0$  by a simple grid search over different  $\tau$  run for one or two epochs each.

#### 4.3 Refinements to the Proposed Adaptive Strategy

#### 4.3.1 Faster Decay When Training Saturates

The communication period update rule (17) tends to give a decreasing sequence  $\{\tau_l\}$ . Nonetheless, it is possible that the best value of  $\tau_l$  for next time interval is larger than the current one due to random noise in the training process. Besides, when the training loss gets stuck on plateaus and

decreases very slowly, (17) will result in  $\tau_l$  saturating at the same value for a long time. To address this issue, we borrow a idea used in classic SGD where the learning rate is decayed by a factor  $\gamma$  when the training loss saturates for several epochs (Goyal et al., 2017). Similarly, in the our scheme, the communication period will be multiplied by  $\gamma < 1$  when the  $\tau_l$  given by (17) is not strictly less than  $\tau_{l-1}$ . To be specific, the communication period for the  $l^{th}$  time interval will be determined as follows:

$$\tau_{l} = \begin{cases} \left\lceil \sqrt{\frac{F(\mathbf{x}_{t=l}\tau_{0})}{F(\mathbf{x}_{t=0})}} \tau_{0} \right\rceil, & \text{if } \left\lceil \sqrt{\frac{F(\mathbf{x}_{t=l}\tau_{0})}{F(\mathbf{x}_{t=0})}} \tau_{0} \right\rceil < \tau_{l-1} \\ \gamma \tau_{l-1}, & \text{otherwise} \end{cases}$$
(18)

In the experiments,  $\gamma=1/2$  turns out to be a good choice. One can obtain a more aggressive decay in  $\tau_l$  by either reducing the value of  $\gamma$  or introducing a slack variable s in the condition, such as  $\lceil \sqrt{\frac{F(\mathbf{x}_{t=l}\tau_0)}{F(\mathbf{x}_{t=0})}} \tau_0 \rceil + s < \tau_{l-1}$ .

#### 4.3.2 Incorporating Adaptive Learning Rate

So far we consider a fixed learning rate  $\eta$  for the local SGD updates at the workers. We now present an adaptive communication strategy that adjusts  $\tau_l$  for a given variable learning rate schedule, in order to obtain the best error-runtime trade-off. Suppose  $\eta_l$  denotes the learning rate for the  $l^{th}$  time interval. Then, combining (15) and (16) again, we have

$$\tau_l = \left[ \sqrt{\frac{\eta_0^3}{\eta_l^3} \frac{F(\overline{\mathbf{x}}_{t=lT_0})}{F(\overline{\mathbf{x}}_{t=0})}} \tau_0 \right]. \tag{19}$$

Observe that when the learning rate becomes smaller, the communication period  $\tau_l$  increases. This result corresponds the intuition that a small learning rate reduces the discrepancy between the local models, and hence is more tolerant to large communication periods.

Equation (19) states that the communication period should be proportional to  $(\eta_0/\eta_l)^{3/2}$ . However, in practice, it is common to decay the learning rate 10 times after some given number of epochs. The dramatic change of learning rate may push the communication period to an unreasonably large value. In the experiments, we observe that when applying (19), the communication period can increase to  $\tau=1000$  which causes the training loss to diverge.

To avoid this issue, we propose the adaptive strategy given by (20) below. This strategy can also be justified by theoretical analysis. Suppose that in  $l^{th}$  time interval, the objective function has a local Lipschitz smoothness  $L_l$ . Then, by using the approximation  $\eta_l L_l \approx 1$ , which is common in SGD literature (Balles et al., 2016), we derive the following

adaptive strategy:

$$\tau_{l} = \left[ \sqrt{\frac{\eta_{0}^{3} L_{0}^{2}}{\eta_{l}^{3} L_{l}^{2}} \frac{F(\mathbf{x}_{t=lT_{0}})}{F(\mathbf{x}_{t=0})}} \tau_{0} \right] \approx \left[ \sqrt{\frac{\eta_{0}}{\eta_{l}} \frac{F(\mathbf{x}_{t=lT_{0}})}{F(\mathbf{x}_{t=0})}} \tau_{0} \right]. \tag{20}$$

Apart from coupling the communication period with learning rate, when to decay the learning rate is another key design factor. In order to eliminate the noise introduced by local updates, we choose to first gradually decay the communication period to 1 and then decay the learning rate as usual. For example, if the learning rate is scheduled to be decayed at the  $80^{th}$  epoch but at that time the communication period  $\tau$  is still larger than 1, then we will continue use the current learning rate until  $\tau=1$ .

### 4.4 Theoretical Guarantees for the Convergence of ADACOMM

In this subsection, we are going to provide a convergence guarantee for the proposed adaptive communication scheme by extending the error analysis for PASGD. Without loss of generality, we will analyze an arbitrary communication period sequence  $\{\tau_0,\ldots,\tau_R\}$ , where R represents the total communication rounds<sup>2</sup>. It will be shown that a decreasing sequence of  $\tau$  is beneficial to the error-convergence rate.

**Theorem 3 (Convergence of adaptive communication scheme).** For PASGD with adaptive communication period and adaptive learning rate, suppose the learning rate remains same in each local update period. If the following conditions are satisfied as  $R \to \infty$ ,

$$\sum_{r=0}^{R} \eta_r \tau_r \to \infty, \ \sum_{r=0}^{R} \eta_r^2 \tau_r < \infty, \ \sum_{r=0}^{R} \eta_r^3 \tau_r^2 < \infty, \quad (21)$$

then the averaged model  $\overline{\mathbf{x}}$  is guaranteed to converge to a stationary point:

$$\mathbb{E}\left[\frac{\sum_{r=0}^{R-1} \eta_r \sum_{k=1}^{\tau_r} \|\nabla F(\overline{\mathbf{x}}_{s_r+k})\|^2}{\sum_{r=0}^{R-1} \eta_r \tau_r}\right] \to 0$$
 (22)

where 
$$s_r = \sum_{j=0}^{r-1} \tau_j$$
.

The proof details and a non-asymptotic result (similar to Theorem 1 but with variable  $\tau$ ) are provided in Appendix. In order to understand the meaning of condition (21), let us first consider the case when  $\tau_0 = \cdots = \tau_R$  is a constant. In this case, the convergence condition is identical to mini-batch

 $<sup>^2</sup>$ Note that in the error analysis, the subscripts of communication period and learning rate represent the index of local update periods rather than the index of the  $T_0$ -length wall-clock time intervals as considered in Sections 4.1-4.3.

SGD (Bottou et al., 2018):

$$\sum_{r=0}^{R} \eta_r \to \infty, \ \sum_{r=0}^{R} \eta_r^2 < \infty. \tag{23}$$

As long as the communication period sequence is bounded, it is trivial to adapt the learning rate scheme in mini-batch SGD (23) to satisfy (21). In particular, when the communication period sequence is decreasing, the last two terms in (21) will become easier to be satisfied and put less constraints on the learning rate sequence.

#### 5 EXPERIMENTAL RESULTS

#### 5.1 Experimental Setting

**Platform.** The proposed adaptive communication scheme was implemented in Pytorch (Paszke et al., 2017) with Mpi4Py (Dalcín et al., 2005). All experiments were conducted on a local cluster with 4 worker nodes, each of which has an NVIDIA TitanX GPU and a 16-core Intel Xeon CPU. Worker nodes are connected via a 40 Gbps (5000 Mb/s) Ethernet interface. Due to space limitations, additional results with 8 nodes are listed in Appendix A.

**Dataset.** We evaluate our method for image classification tasks on CIFAR10 and CIFAR100 dataset (Krizhevsky, 2009), which consists of 50,000 training images and 10,000 validation images in 10 and 100 classes respectively. Each worker machine is assigned with a partition which will be randomly shuffled after every epoch.

**Model.** We choose to train deep neural networks VGG-16 (Simonyan & Zisserman, 2014) and ResNet-50 (He et al., 2016) from scratch  $^3$ . These two neural networks have different architectures and parameter sizes, thus resulting in different performance of periodic-averaging. As shown in Figure 8, for VGG-16, the communication time is about 4 times higher than the computation time. Thus, compared to ResNet-50, it requires a larger  $\tau$  in order to reduce the runtime-per-iteration and achieve fast convergence. Similar high communication/computation ratio is common in literature, see (Lin et al., 2018; Harlap et al., 2018).

**Hyperparameter Choice.** Mini-batch size on each worker is 128. Therefore, the total mini-batch size per iteration is 512. The initial learning rates for VGG-16 and ResNet-50 are 0.2 and 0.4 respectively. The weight decay for both networks is 0.0005. In the variable learning rate setting, we decay the learning rate by 10 after  $80^{\text{th}}/120^{\text{th}}/160^{\text{th}}/200^{\text{th}}$  epochs. We set the time interval length  $T_0$  as 60 seconds (about 10 epochs for the initial communication period).

Metrics. We compare the performance of proposed adaptive

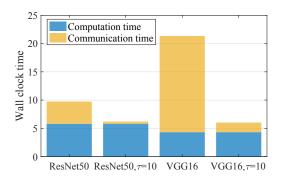


Figure 8. Wall-clock time to finish 100 iterations in a cluster with 4 worker nodes. To achieve the same level communication/computation ratio, VGG-16 requires larger communication period than ResNet-50.

communication scheme with following methods with a fixed communication period: (1) Baseline: fully synchronous SGD ( $\tau=1$ ); (2) Extreme high throughput case where  $\tau=100$ ; (3) Manually tuned case where a moderate value of  $\tau$  is selected after trial runs with different communication periods. Instead of training for a fixed number of epochs, we train all methods for sufficiently long time to convergence and compare the training loss and test accuracy, both of which are recorded after every 100 iterations.

#### 5.2 Adaptive Communication in PASGD

We first validate the effectiveness of ADACOMM which uses the communication period update rule (18) combined with (20) on original PASGD without momentum.

Figure 9 presents the results for VGG-16 for both fixed and variable learning rates. A large communication period  $\tau$  initially results in a rapid drop in the error, but the error finally converges to higher floor. By adapting  $\tau$ , the proposed ADACOMM scheme strikes the best error-runtime trade-off in all settings. In Figure 9a, while fully synchronous SGD takes 33.5 minutes to reach  $4\times 10^{-3}$  training loss, ADACOMM costs 15.5 minutes achieving more than  $2\times$  speedup. Similarly, in Figure 9b, ADACOMM takes 11.5 minutes to reach  $4.5\times 10^{-2}$  training loss achieving  $3.3\times$  speedup over fully synchronous SGD (38.0 minutes).

However, for ResNet-50, the communication overhead is no longer the bottleneck. For fixed communication period, the negative effect of performing local updates becomes more obvious and cancels the benefit of low communication delay (see Figures 10b and 10c). It is not surprising to see fully synchronous SGD is nearly the best one in the error-runtime plot among all fixed- $\tau$  methods. Even in this extreme case, adaptive communication can still have a competitive performance. When combined with learning rate decay, the adaptive scheme is about 1.3 times faster

<sup>&</sup>lt;sup>3</sup>The implementations of VGG-16 and ResNet-50 follow this GitHub repository: https://github.com/meliketoy/wide-resnet.pytorch

than fully synchronous SGD (see Figure 10a, 15.0 versus 21.5 minutes to achieve  $3 \times 10^{-2}$  training loss).

Table 1 lists the test accuracies in different settings; we report the best accuracy within a time budget for each setting. The results show that adaptive communication method have better generalization than fully synchronous SGD. In the variable learning rate case, the adaptive method even gives the better test accuracy than PASGD with the best fixed  $\tau$ .

#### 5.3 Adaptive Communication in Momentum SGD

The adaptive communication scheme is proposed based on the joint error-runtime analysis for PASGD without momentum. However, it can also be extended to other SGD variants, and in this subsection, we show that the proposed method works well for SGD with momentum.

#### 5.3.1 Block Momentum in periodic-averaging

Before presenting the empirical results, we describe how to introduce momentum in PASGD. A naive way is to apply the momentum independently to each local model, where each worker maintains an independent momentum buffer, which is the latest change in the parameter vector  $\mathbf{x}$ . However, this does not account for the potential dramatic change in  $\mathbf{x}$  at each averaging step. When local models are synchronized, the local momentum buffer will contain the update steps before averaging, resulting in a large momentum term in the first SGD step of the each local update period. When  $\tau$  is large, this large momentum term can side-track the SGD descent direction resulting in slower convergence.

To address this issue, a *block momentum scheme* was proposed in (Chen & Huo, 2016) and applied to speech recognition tasks. The basic idea is to treat the local updates in each communication period as one big gradient step between two synchronized models, and to introduce a global momentum for this big accumulated step. The update rule can be written as follows in terms of the momentum  $\mathbf{u}_i$ :

$$\mathbf{u}_{i} = \beta_{\text{glob}} \mathbf{u}_{i-1} + \mathcal{G}_{i} \tag{24}$$

$$\mathbf{x}_{(j+1)\tau+1} = \mathbf{x}_{j\tau+1} - \eta_j \mathbf{u}_j \tag{25}$$

where  $\mathcal{G}_j = \frac{1}{m} \sum_{i=1}^m \sum_{k=1}^m g(\mathbf{x}_{j\tau+k}^{(i)})$  represents the accumulated gradients in the  $j^{th}$  local update period and  $\beta_{\text{glob}}$  denotes the global momentum factor. Moreover, workers can also conduct momentum SGD on local models, but their local momentum buffer will be cleared at the beginning of each local update period. That is, we restart momentum SGD on local models after every averaging step. The same strategy was also suggested in Microsoft's CNTK framework (Seide & Agarwal, 2016). In our experiments, we set the global momentum factor as 0.3 and local momentum factor as 0.9 following (Lin et al., 2018). In the fully synchronous case, there is no need to introduce the block

Table 1. Best test accuracies on CIFAR10 in different settings (SGD without momentum).

MODEL	METHODS	FIXED LR	VARIABLE LR
VGG-16	$\begin{array}{c c} \tau = 1 \\ \tau = 20 \\ \tau = 100 \\ \text{ADACOMM} \end{array}$	90.5 <b>92.25</b> 92.0 91.1	92.75 92.5 92.4 <b>92.85</b>
RESNET- 50	$\begin{array}{c c} \tau = 1 \\ \tau = 5 \\ \tau = 100 \\ \text{ADACOMM} \end{array}$	88.76 <b>90.42</b> 88.66 89.57	92.26 92.26 91.8 <b>92.42</b>

momentum and we simply follow the common practice setting the momentum factor as 0.9.

#### 5.3.2 ADACOMM plus Block Momentum

In Figure 11, we apply our adaptive communication strategy in PASGD with block momentum and observe significant performance gain on CIFAR10/100. In particular, the adaptive communication scheme has the fastest convergence rate with respect to wall-clock time in the whole training process. While fully synchronous SGD gets stuck with a plateau before the first learning rate decay, the training loss of adaptive method continuously decreases until converging. For VGG-16 in Figure 11b, ADACOMM is  $3.5\times$  faster (in terms of wall-clock time) than fully synchronous SGD in reaching a  $3\times10^{-3}$  training loss. For ResNet-50 in Figure 11a, ADACOMM takes 15.8 minutes to get  $2\times10^{-2}$  training loss which is 2 times faster than fully synchronous SGD (32.6 minutes).

#### 6 CONCLUDING REMARKS

The design of communication-efficient SGD algorithms that are robust to system variability is vital to scaling machine learning training to resource-limited computing nodes. This paper is one of the first to analyze the convergence of error with respect to wall-clock time instead of number of iterations by accounting for the effect of computation and communication delays on the runtime per iteration. We present a theoretical analysis of the error-runtime trade-off for periodic-averaging SGD (PASGD), where each node performs local updates and their models are averaged after every  $\tau$  iterations. Based on the joint error-runtime analysis, we design the first (to the best of our knowledge) adaptive communication strategy called ADACOMM for distributed deep learning. Experimental results using VGGNet and ResNet show that the proposed method can achieve up to a 3× improvement in runtime, while achieving the same error floor as fully synchronous SGD.

Going beyond periodic-averaging SGD, our idea of adapting

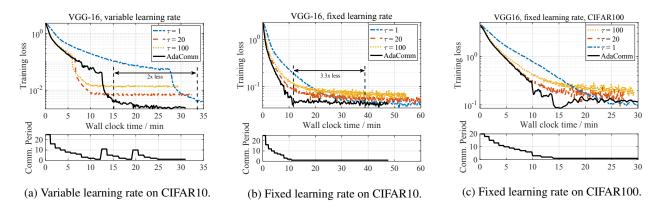


Figure 9. ADACOMM on VGG-16: Achieves  $3.3 \times$  speedup over fully synchronous SGD (in (b), 11.5 versus 38.0 minutes to achieve  $4.5 \times 10^{-2}$  training loss).

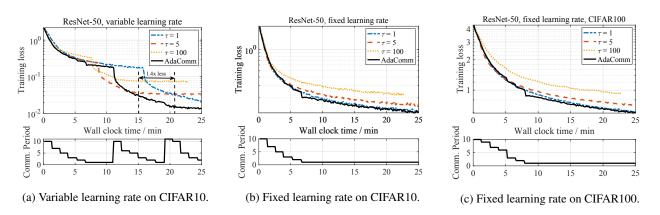


Figure 10. ADACOMM on ResNet-50: Achieves  $1.4 \times$  speedup over Sync SGD (in (a), 15.0 versus 21.5 minutes to achieve  $3 \times 10^{-2}$  training loss).

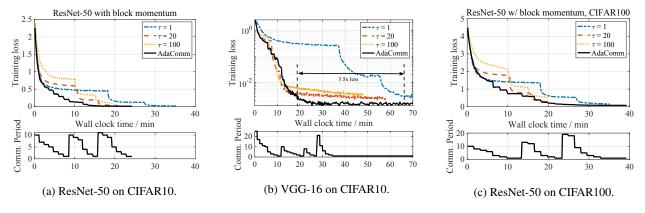


Figure 11. ADACOMM with block momentum achieves  $3.5 \times$  speedup over Sync SGD (in (b), 19.0 versus 66.7 minutes to achieve  $3 \times 10^{-3}$  training loss).

frequency of averaging distributed SGD updates can be easily extended to other SGD frameworks including elastic-averaging (Zhang et al., 2015), decentralized SGD (e.g., adapting network sparsity) (Lian et al., 2017) and parameter server-based training (e.g., adapting asynchrony).

#### ACKNOWLEDGMENTS

The authors thank Prof. Greg Ganger for helpful discussions. This work was partially supported by NSF CCF-1850029 and an IBM Faculty Award. Experiments were conducted on clusters provided by the Parallel Data Lab at CMU.

#### REFERENCES

- Balles, L., Romero, J., and Hennig, P. Coupling adaptive batch sizes with learning rates. *arXiv preprint* arXiv:1612.05086, 2016.
- Bottou, L., Curtis, F. E., and Nocedal, J. Optimization methods for large-scale machine learning. *SIAM Review*, 60(2):223–311, 2018.
- Chaudhari, P., Baldassi, C., Zecchina, R., Soatto, S., Talwalkar, A., and Oberman, A. Parle: parallelizing stochastic gradient descent. *arXiv preprint arXiv:1707.00424*, 2017.
- Chen, K. and Huo, Q. Scalable training of deep learning machines by incremental block training with intra-block parallel optimization and blockwise model-update filtering. In *Acoustics, Speech and Signal Processing (ICASSP)*, 2016 IEEE International Conference on, pp. 5880–5884. IEEE, 2016.
- Cui, H., Cipar, J., Ho, Q., Kim, J. K., Lee, S., Kumar, A., Wei, J., Dai, W., Ganger, G. R., Gibbons, P. B., et al. Exploiting bounded staleness to speed up big data analytics. In 2014 USENIX Annual Technical Conference (USENIX ATC 14), pp. 37–48, 2014.
- Dalcín, L., Paz, R., and Storti, M. MPI for python. *Journal of Parallel and Distributed Computing*, 65(9):1108–1115, 2005.
- David, H. A. and Nagaraja, H. N. *Order statistics*. John Wiley, Hoboken, N.J., 2003.
- Dean, J., Corrado, G., Monga, R., Chen, K., Devin, M., Mao, M., Senior, A., Tucker, P., Yang, K., Le, Q. V., et al. Large scale distributed deep networks. In *Advances in neural information processing systems*, pp. 1223–1231, 2012.
- Dekel, O., Gilad-Bachrach, R., Shamir, O., and Xiao, L. Optimal distributed online prediction using mini-batches. *Journal of Machine Learning Research*, 13(Jan):165–202, 2012.
- Dutta, S., Joshi, G., Ghosh, S., Dube, P., and Nagpurkar, P. Slow and stale gradients can win the race: Errorruntime trade-offs in distributed SGD. arXiv preprint arXiv:1803.01113, 2018.
- Ghadimi, S. and Lan, G. Stochastic first-and zeroth-order methods for nonconvex stochastic programming. *SIAM Journal on Optimization*, 23(4):2341–2368, 2013.
- Goyal, P., Dollár, P., Girshick, R., Noordhuis, P., Wesolowski, L., Kyrola, A., Tulloch, A., Jia, Y., and He, K. Accurate, large minibatch SGD: training ImageNet in 1 hour. arXiv preprint arXiv:1706.02677, 2017.

- Gupta, S., Zhang, W., and Wang, F. Model accuracy and runtime tradeoff in distributed deep learning: A systematic study. In *IEEE 16th International Conference on Data Mining (ICDM)*, pp. 171–180. IEEE, 2016.
- Harlap, A., Narayanan, D., Phanishayee, A., Seshadri, V.,
  Devanur, N., Ganger, G., and Gibbons, P. Pipedream:
  Fast and efficient pipeline parallel dnn training. arXiv preprint arXiv:1806.03377, 2018.
- He, K., Zhang, X., Ren, S., and Sun, J. Deep residual learning for image recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pp. 770–778, 2016.
- Iandola, F. N., Moskewicz, M. W., Ashraf, K., and Keutzer, K. Firecaffe: near-linear acceleration of deep neural network training on compute clusters. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pp. 2592–2600, 2016.
- Jiang, Z., Balu, A., Hegde, C., and Sarkar, S. Collaborative deep learning in fixed topology networks. In *Advances in Neural Information Processing Systems*, pp. 5906–5916, 2017.
- Krizhevsky, A. Learning multiple layers of features from tiny images. Technical report, Citeseer, 2009.
- Li, M., Andersen, D. G., Park, J. W., Smola, A. J., Ahmed, A., Josifovski, V., Long, J., Shekita, E. J., and Su, B.-Y. Scaling distributed machine learning with the parameter server. In *OSDI*, volume 14, pp. 583–598, 2014.
- Lian, X., Zhang, C., Zhang, H., Hsieh, C.-J., Zhang, W., and Liu, J. Can decentralized algorithms outperform centralized algorithms? a case study for decentralized parallel stochastic gradient descent. In *Advances in Neural Information Processing Systems*, pp. 5336–5346, 2017.
- Lin, T., Stich, S. U., and Jaggi, M. Don't use large minibatches, use local SGD. *arXiv preprint arXiv:1808.07217*, 2018.
- McMahan, H. B., Moore, E., Ramage, D., Hampson, S., et al. Communication-efficient learning of deep networks from decentralized data. *arXiv preprint arXiv:1602.05629*, 2016.
- Mitliagkas, I., Zhang, C., Hadjis, S., and Ré, C. Asynchrony begets momentum, with an application to deep learning. In 54th Annual Allerton Conference on Communication, Control, and Computing (Allerton), pp. 997–1004. IEEE, 2016.
- Moritz, P., Nishihara, R., Stoica, I., and Jordan, M. I. SparkNet: Training deep networks in spark. *arXiv* preprint arXiv:1511.06051, 2015.

- Paszke, A., Gross, S., Chintala, S., Chanan, G., Yang, E., DeVito, Z., Lin, Z., Desmaison, A., Antiga, L., and Lerer, A. Automatic differentiation in pytorch. In NIPS-W, 2017.
- Seide, F. and Agarwal, A. CNTK: Microsoft's open-source deep-learning toolkit. In *Proceedings of the 22nd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, pp. 2135–2135. ACM, 2016.
- Simonyan, K. and Zisserman, A. Very deep convolutional networks for large-scale image recognition. *arXiv* preprint arXiv:1409.1556, 2014.
- Smith, S. L., Kindermans, P.-J., and Le, Q. V. Don't decay the learning rate, increase the batch size. *arXiv preprint arXiv:1711.00489*, 2017a.
- Smith, V., Chiang, C.-K., Sanjabi, M., and Talwalkar, A. S. Federated multi-task learning. In *Advances in Neural Information Processing Systems*, pp. 4424–4434. 2017b.
- Stich, S. U. Local SGD converges fast and communicates little. *arXiv preprint arXiv:1805.09767*, 2018.
- Su, H. and Chen, H. Experiments on parallel training of deep neural network using model averaging. *arXiv* preprint *arXiv*:1507.01239, 2015.
- Wang, J. and Joshi, G. Cooperative SGD: A unified framework for the design and analysis of communication-efficient SGD algorithms. *arXiv* preprint *arXiv*:1808.07576, 2018.
- Yu, H., Yang, S., and Zhu, S. Parallel restarted SGD for non-convex optimization with faster convergence and less communication. arXiv preprint arXiv:1807.06629, 2018.
- Zhang, J., De Sa, C., Mitliagkas, I., and Ré, C. Parallel SGD: When does averaging help? *arXiv preprint arXiv:1606.07365*, 2016.
- Zhang, S., Choromanska, A. E., and LeCun, Y. Deep learning with elastic averaging SGD. In NIPS'15 Proceedings of the 28th International Conference on Neural Information Processing Systems, pp. 685–693, 2015.
- Zhou, F. and Cong, G. On the convergence properties of a *k*-step averaging stochastic gradient descent algorithm for nonconvex optimization. *arXiv* preprint *arXiv*:1708.01012, 2017.

#### A ADDITIONAL EXPERIMENTAL RESULTS

In the 8 worker case, the communication among nodes is accomplished via Nvidia Collective Communication Library (NCCL). The mini-batch size on each node is 64. The initial learning rate is set as 0.2 for both VGG-16 and ResNet-50. In Figure 12a, while fully synchronous SGD takes 17.5 minutes to reach  $10^{-2}$  training loss, ADACOMM only costs 6.0 minutes achieving about  $2.9 \times$  speedup.

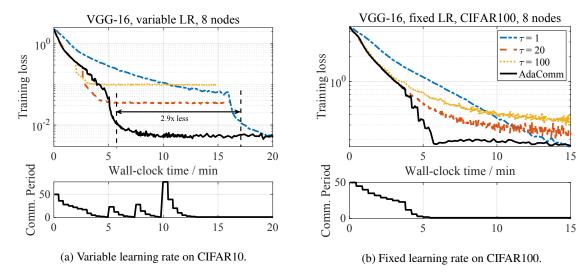


Figure 12. ADACOMM on VGG-16 with 8 workers: Achieves  $2.9\times$  speedup over Sync SGD (in (a), 6.0 versus 17.5 minutes to achieve  $1\times10^{-2}$  training loss). Test accuracies at convergence when using variable learning rate: 92.52% ( $\tau=1$ ), 91.85% ( $\tau=20$ ), 91.15% ( $\tau=100$ ), and 92.72% (AdaComm).

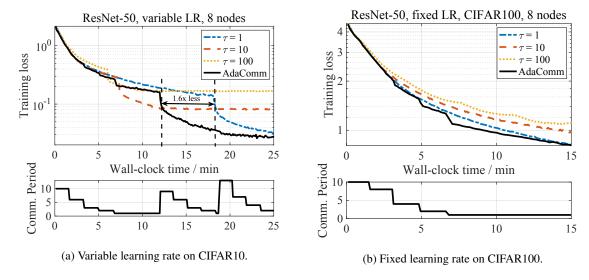


Figure 13. ADACOMM on ResNet-50 with 8 workers: Achieves  $1.6\times$  speedup over Sync SGD (in (a), 11.15 versus 18.25 minutes to achieve  $1\times10^{-1}$  training loss). Test accuracies at convergence when using variable learning rate: 91.93% ( $\tau=1$ ), 91.51% ( $\tau=10$ ), 90.46% ( $\tau=100$ ), and 91.77% (AdaComm).

#### **B** INEFFICIENT LOCAL UPDATES

It is worth noting there is an interesting phenomenon about the convergence of periodic averaging SGD (PASGD). When the learning rate is fixed, PASGD with fine-tuned communication period has better test accuracy than both fully synchronous

SGD and the adaptive method, while its training loss remains higher than the latter two methods (see Figure 9, Figure 10). In particular, on CIFAR100 dataset, we observe about 5% improvement in test accuracy when  $\tau=5$ . To investigate this phenomenon, we evaluate the test accuracy for PASGD ( $\tau=15$ ) in two frequencies: 1) every 135 iterations; 2) every 100 iterations. In the former case, the test accuracy is reported just after the averaging step. However, in the latter case, the test accuracy can come from either the synchronized/averaged model or local models, since 100 cannot be divided by 15.

From Figure 14, it is clear that local model's accuracy is much lower than the synchronized model, even when the algorithm has converged. Thus, we conjecture that the improvement of test accuracy only happens on the synchronized model. That is, after averaging, the test accuracy will undergo a rapid increase but it decreases again in the following local steps due to noise in stochastic gradients. Such behavior may depend on the geometric structure of the loss surface of specific neural networks. The observation also reveals that the local updates are inefficient as they reduces the accuracy and makes no progress. In this sense, it is necessary for PASGD to reduce the gradient variance by either decaying learning rate or decaying communication period.

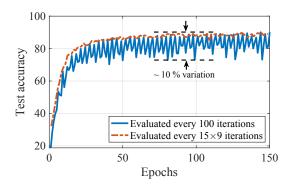


Figure 14. PASGD ( $\tau=15$ ) with ResNet-50 on CIFAR10 (fixed learning rate, no momentum). There exists about 10% accuracy gap between local models and the synchronized model.

#### C ASSUMPTIONS FOR CONVERGENCE ANALYSIS

The convergence analysis is conducted under the following assumptions, which are similar to the assumptions made in previous work on the analysis of PASGD (Zhou & Cong, 2017; Yu et al., 2018; Wang & Joshi, 2018; Stich, 2018). In particular, we make no assumptions on the convexity of the objective function. We also remove the uniform bound assumption for the norm of stochastic gradients.

**Assumption 1** (Lipschitz smooth & lower bound on F). The objective function  $F(\mathbf{x})$  is differentiable and L-Lipschitz smooth, i.e.,  $\|\nabla F(\mathbf{x}) - \nabla F(\mathbf{y})\| \le L \|\mathbf{x} - \mathbf{y}\|$ . The function value is bounded below by a scalar  $F_{inf}$ .

**Assumption 2** (Unbiased estimation). The stochastic gradient evaluated on a mini-batch  $\xi$  is an unbiased estimator of the full batch gradient  $\mathbb{E}_{\xi|\mathbf{x}}[g(\mathbf{x})] = \nabla F(\mathbf{x})$ .

**Assumption 3** (Bounded variance). The variance of stochastic gradient evaluated on a mini-batch  $\xi$  is bounded as

$$\mathbb{E}_{\xi \mid \mathbf{x}} \|g(\mathbf{x}) - \nabla F(\mathbf{x})\|^2 \le \beta \|\nabla F(\mathbf{x})\|^2 + \sigma^2$$

where  $\beta$  and  $\sigma^2$  are non-negative constants and in inverse proportion to the mini-batch size.

#### D PROOF OF THEOREM 2: ERROR-RUNTIME CONVERGENCE OF PASGD

Firstly, let us recall the error-analysis of PASGD. We adapt the theorem from (Wang & Joshi, 2018).

**Lemma 1** (Error-Convergence of PASGD (Wang & Joshi, 2018)). For PASGD, under Assumptions 1 to 3, if the learning rate satisfies  $\eta L + \eta^2 L^2 \tau(\tau - 1) \le 1$  and all workers are initialized at the same point  $\mathbf{x}_1$ , then after K iterations, we have

$$\mathbb{E}\left[\min_{k\in[1,K]}\|\nabla F(\overline{\mathbf{x}}_k)\|^2\right] \le \mathbb{E}\left[\frac{1}{K}\sum_{k=1}^K\|\nabla F(\overline{\mathbf{x}}_k)\|^2\right] \le \frac{2\left[F(\mathbf{x}_1) - F_{inf}\right]}{\eta K} + \frac{\eta L\sigma^2}{m} + \eta^2 L^2\sigma^2(\tau - 1)$$
(26)

where L is the Lipschtiz constant of the objective function,  $\sigma^2$  is the variance bound of mini-batch stochastic gradients and  $\overline{\mathbf{x}}_k$  denotes the averaged model at the  $k^{th}$  iteration.

From the runtime analysis in Section 2, we know that the expected runtime per iteration of PASGD is

$$\mathbb{E}[T_{\text{P-Avg}}] = Y + \frac{D}{\tau}.\tag{27}$$

Accordingly, the total wall-clock time of training K iteration is

$$T = K\left(Y + \frac{D}{\tau}\right). \tag{28}$$

Then, directly substituting  $K = T/\mathbb{E}[T_{P\text{-Avg}}]$  in (26), we complete the proof.

#### E PROOF OF THEOREM 3: THE BEST COMMUNICATION PERIOD

Taking the derivative of the upper bound (14) with respect to the communication period, we obtain

$$-\frac{2\left[F(\mathbf{x}_1) - F_{\text{inf}}\right]}{\eta T} \frac{\mathbb{E}\left[D\right]}{\tau^2} + \eta^2 L^2 \sigma^2. \tag{29}$$

When the derivative equals to zero, the communication period is

$$\tau^* = \sqrt{\frac{2(F(\mathbf{x}_1) - F_{\text{inf}})\mathbb{E}[D]}{\eta^3 L^2 \sigma^2 T}}.$$
(30)

Since the second derivative of (14) is

$$\frac{4\left[F(\mathbf{x}_1) - F_{\inf}\right]}{nT} \frac{\mathbb{E}\left[D\right]}{\tau^3} > 0,\tag{31}$$

then the optimal value obtained in (30) must be a global minimum.

### F PROOF OF THEOREM 4: ERROR-CONVERGENCE OF ADAPTIVE COMMUNICATION SCHEME

#### F.1 Notations

In order to faciliate the analysis, we would like to first introduce some useful notations. Define matrices  $\mathbf{X}_k, \mathbf{G}_k \in \mathbb{R}^{d \times m}$  that concatenate all local models and gradients:

$$\mathbf{X}_k = [\mathbf{x}_k^{(1)}, \dots, \mathbf{x}_k^{(m)}],\tag{32}$$

$$\mathbf{G}_k = [g(\mathbf{x}_k^{(1)}), \dots, g(\mathbf{x}_k^{(m)})]. \tag{33}$$

Besides, define matrix  $\mathbf{J} = \mathbf{1}\mathbf{1}^{\top}/(\mathbf{1}^{\top}\mathbf{1})$  where  $\mathbf{1}$  denotes the column vector  $[1,1,\ldots,1]^{\top}$ . Unless otherwise stated,  $\mathbf{1}$  is a size m column vector, and the matrix  $\mathbf{J}$  and identity matrix  $\mathbf{I}$  are of size  $m \times m$ , where m is the number of workers.

#### F.2 Proof

Let us first focus on the j-th local update period, where  $j \in \{0, 1, \dots, R\}$ . Without loss of generality, suppose the local index of the  $j^{th}$  local update period starts from 1 and ends with  $\tau_j$ . Then, for the k-th local step in the interested period, we have the following lemma.

**Lemma 2** (Lemma 1 in (Wang & Joshi, 2018)). For PASGD, under Assumptions 1 to 3, at the k-th iteration, we have the following bound for the objective value:

$$\mathbf{E}_{k}\left[F(\overline{\mathbf{x}}_{k+1})\right] - F(\overline{\mathbf{x}}_{k}) \leq -\frac{\eta_{j}}{2} \left\|\nabla F(\overline{\mathbf{x}}_{k})\right\|^{2} - \frac{\eta_{j}}{2} \left[1 - \eta_{j} L\left(\frac{\beta}{m} + 1\right)\right] \cdot \frac{\left\|\nabla F(\mathbf{X}_{k})\right\|_{F}^{2}}{m} + \frac{\eta_{j}^{2} L \sigma^{2}}{2m} + \frac{\eta_{j} L^{2}}{2m} \left\|\mathbf{X}_{k}(\mathbf{I} - \mathbf{J})\right\|_{F}^{2}$$

$$(34)$$

where  $\overline{\mathbf{x}}_k$  denotes the averaged model at the  $k^{th}$  iteration.

Taking the total expectation and summing over all iterates in the j-th local update period, we can obtain

$$\mathbb{E}\left[F(\overline{\mathbf{x}}_{\tau_{j}+1}) - F(\overline{\mathbf{x}}_{1})\right] \leq -\frac{\eta_{j}}{2} \sum_{k=1}^{\tau_{j}} \mathbb{E}\left\|\nabla F(\overline{\mathbf{x}}_{k})\right\|^{2} - \frac{\eta_{j}}{2} \left[1 - \eta_{j} L\left(\frac{\beta}{m} + 1\right)\right] \cdot \sum_{k=1}^{\tau_{j}} \frac{\mathbb{E}\left\|\nabla F(\mathbf{X}_{k})\right\|_{F}^{2}}{m} + \frac{\eta_{j}^{2} L \sigma^{2} \tau_{j}}{2m} + \frac{\eta_{j} L^{2}}{2m} \sum_{k=1}^{\tau_{j}} \mathbb{E}\left\|\mathbf{X}_{k}(\mathbf{I} - \mathbf{J})\right\|_{F}^{2}.$$
(35)

Next, we are going to provide an upper bound for the last term in (35). Note that

$$\mathbf{X}_{k}(\mathbf{I} - \mathbf{J}) = \mathbf{X}_{k-1}(\mathbf{I} - \mathbf{J}) - \eta_{i}\mathbf{G}_{k-1}(\mathbf{I} - \mathbf{J})$$
(36)

$$= \mathbf{X}_{k-2}(\mathbf{I} - \mathbf{J}) - \eta_j \mathbf{G}_{k-2}(\mathbf{I} - \mathbf{J}) - \eta_j \mathbf{G}_{k-1}(\mathbf{I} - \mathbf{J})$$
(37)

$$= \mathbf{X}_1(\mathbf{I} - \mathbf{J}) - \eta_j \sum_{r=1}^{k-1} \mathbf{G}_r(\mathbf{I} - \mathbf{J})$$
(38)

$$= -\eta_j \sum_{r=1}^{k-1} \mathbf{G}_r (\mathbf{I} - \mathbf{J}) \tag{39}$$

where (39) follows the fact that all workers start from the same point at the beginning of each local update period, i.e.,  $\mathbf{X}_1(\mathbf{I} - \mathbf{J}) = 0$ . Accordingly, we have

$$\mathbb{E}\left[\left\|\mathbf{X}_{k}(\mathbf{I}-\mathbf{J})\right\|_{F}^{2}\right] = \eta_{j}^{2}\mathbb{E}\left[\left\|\sum_{r=1}^{k-1}\mathbf{G}_{r}(\mathbf{I}-\mathbf{J})\right\|_{F}^{2}\right]$$
(40)

$$\leq \eta_j^2 \mathbb{E} \left[ \left\| \sum_{r=1}^{k-1} \mathbf{G}_r \right\|_{\mathrm{F}}^2 \right] = \eta_j^2 \sum_{i=1}^m \mathbb{E} \left[ \left\| \sum_{r=1}^{k-1} g(\mathbf{x}_r^{(i)}) \right\|^2 \right]$$
(41)

where the inequality (41) is due to the operator norm of  $(\mathbf{I} - \mathbf{J})$  is less than 1. Furthermore, using the fact  $(a+b)^2 \le 2a^2 + 2b^2$ , one can get

$$\mathbb{E}\left[\left\|\mathbf{X}_{k}(\mathbf{I}-\mathbf{J})\right\|_{F}^{2}\right] \leq \eta_{j}^{2} \sum_{i=1}^{m} \mathbb{E}\left[\left\|\sum_{r=1}^{k-1} \left(g(\mathbf{x}_{r}^{(i)}) - \nabla F(\mathbf{x}_{r}^{(i)})\right) + \sum_{r=1}^{k-1} \nabla F(\mathbf{x}_{r}^{(i)})\right\|^{2}\right]$$
(42)

$$\leq 2\eta_{j}^{2} \sum_{i=1}^{m} \mathbb{E}\left[\left\|\sum_{r=1}^{k-1} \left(g(\mathbf{x}_{r}^{(i)}) - \nabla F(\mathbf{x}_{r}^{(i)})\right)\right\|^{2}\right] + 2\eta_{j}^{2} \sum_{i=1}^{m} \mathbb{E}\left[\left\|\sum_{r=1}^{k-1} \nabla F(\mathbf{x}_{r}^{(i)})\right\|^{2}\right]\right]. \tag{43}$$

For the first term  $T_1$ , since the stochastic gradients are unbiased, all cross terms are zero. Thus, combining with Assumption 3, we have

$$T_1 = 2\eta_j^2 \sum_{i=1}^m \sum_{r=1}^{k-1} \mathbb{E}\left[ \left\| g(\mathbf{x}_r^{(i)}) - \nabla F(\mathbf{x}_r^{(i)}) \right\|^2 \right]$$
 (44)

$$\leq 2\eta_j^2 \sum_{i=1}^m \sum_{r=1}^{k-1} \left[ \beta \mathbb{E}\left[ \left\| \nabla F(\mathbf{x}_r^{(i)}) \right\|^2 \right] + \sigma^2 \right]$$

$$(45)$$

$$=2\eta_j^2 \beta \sum_{r=1}^{k-1} \mathbb{E}\left[\|\nabla F(X_r)\|_F^2\right] + 2\eta_j^2 m(k-1)\sigma^2.$$
(46)

For the second term in (43), directly applying Jensen's inequality, we get

$$T_2 \le 2\eta_j^2(k-1) \sum_{i=1}^m \sum_{r=1}^{k-1} \mathbb{E}\left[ \left\| \nabla F(\mathbf{x}_r^{(i)}) \right\|^2 \right]$$
 (47)

$$=2\eta_{j}^{2}(k-1)\sum_{r=1}^{k-1}\mathbb{E}\left[\left\|\nabla F(\mathbf{X}_{r})\right\|_{F}^{2}\right].$$
(48)

Substituting the bounds of  $T_1$  and  $T_2$  into (43),

$$\mathbb{E}\left[\left\|\mathbf{X}_{k}(\mathbf{I}-\mathbf{J})\right\|_{F}^{2}\right] \leq 2\eta_{j}^{2}\left[\beta+(k-1)\right]\sum_{r=1}^{k-1}\mathbb{E}\left[\left\|\nabla F(\mathbf{X}_{r})\right\|_{F}^{2}\right]+2\eta_{j}^{2}m(k-1)\sigma^{2}.\tag{49}$$

Recall the upper bound (35), we further derive the following bound:

$$\sum_{k=1}^{\tau_{j}} \mathbb{E}\left[\|\mathbf{X}_{k}(\mathbf{I} - \mathbf{J})\|_{F}^{2}\right] \leq 2\eta_{j}^{2} \sum_{k=1}^{\tau_{j}} \left[ \left[\beta + (k-1)\right] \sum_{r=1}^{k-1} \mathbb{E}\left[\|\nabla F(\mathbf{X}_{r})\|_{F}^{2}\right] \right] + 2\eta_{j}^{2} m \sigma^{2} \sum_{k=1}^{\tau_{j}} (k-1)$$
(50)

$$=2\eta_{j}^{2} \sum_{k=1}^{\tau_{j}} \left[ \left[ \beta + (k-1) \right] \sum_{r=1}^{k-1} \mathbb{E} \left[ \left\| \nabla F(\mathbf{X}_{r}) \right\|_{F}^{2} \right] + \eta_{j}^{2} m \sigma^{2} \tau_{j} (\tau_{j} - 1)$$
 (51)

$$=2\eta_{j}^{2} \sum_{k=2}^{\tau_{j}} \left[ \left[ \beta + (k-1) \right] \sum_{r=1}^{k-1} \mathbb{E} \left[ \left\| \nabla F(\mathbf{X}_{r}) \right\|_{F}^{2} \right] + \eta_{j}^{2} m \sigma^{2} \tau_{j} (\tau_{j} - 1).$$
 (52)

Then, since  $\sum_{r=1}^{k-1} \mathbb{E}\left[\left\|\nabla F(\mathbf{X}_r)\right\|_{\mathrm{F}}^2\right] \leq \sum_{r=1}^{\tau_j-1} \mathbb{E}\left[\left\|\nabla F(\mathbf{X}_r)\right\|_{\mathrm{F}}^2\right]$ , we have

$$\sum_{k=1}^{\tau_{j}} \mathbb{E}\left[\|\mathbf{X}_{k}(\mathbf{I} - \mathbf{J})\|_{F}^{2}\right] \leq 2\eta_{j}^{2} \sum_{r=1}^{\tau_{j}-1} \mathbb{E}\left[\|\nabla F(\mathbf{X}_{r})\|_{F}^{2}\right] \sum_{k=2}^{\tau_{j}} \left[\beta + (k-1)\right] + \eta_{j}^{2} m \sigma^{2} \tau_{j}(\tau_{j} - 1)$$
(53)

$$= \eta_j^2 \sum_{r=1}^{\tau_j - 1} \mathbb{E}\left[ \|\nabla F(\mathbf{X}_r)\|_F^2 \right] \left[ 2\beta(\tau_j - 1) + \tau_j(\tau_j - 1) \right] + \eta_j^2 m \sigma^2 \tau_j(\tau_j - 1). \tag{54}$$

Plugging (54) into (35),

$$\mathbb{E}\left[F(\overline{\mathbf{x}}_{\tau_{j}+1}) - F(\overline{\mathbf{x}}_{1})\right] \leq -\frac{\eta_{j}}{2} \sum_{k=1}^{\tau_{j}} \mathbb{E} \|\nabla F(\overline{\mathbf{x}}_{k})\|^{2} - \frac{\eta_{j}}{2} \left[1 - \eta_{j} L\left(\frac{\beta}{m} + 1\right)\right] \cdot \sum_{k=1}^{\tau_{j}} \frac{\mathbb{E}\left[\|\nabla F(\mathbf{X}_{k})\|_{F}^{2}\right]}{m} + \frac{\eta_{j}^{2} L \sigma^{2} \tau_{j}}{2m} + \frac{\eta_{j}^{3} L^{2} \sigma^{2} \tau_{j} (\tau_{j} - 1)}{2} \qquad (55)$$

$$\leq -\frac{\eta_{j}}{2} \sum_{k=1}^{\tau_{j}} \mathbb{E} \|\nabla F(\overline{\mathbf{x}}_{k})\|^{2} - \frac{\eta_{j}}{2} \left[1 - \eta_{j} L\left(\frac{\beta}{m} + 1\right)\right] \cdot \sum_{k=1}^{\tau_{j}} \frac{\mathbb{E}\left[\|\nabla F(\mathbf{X}_{k})\|_{F}^{2}\right]}{m} + \frac{\eta_{j}^{2} L \sigma^{2} \tau_{j}}{2m} + \frac{\eta_{j}^{3} L^{2} \sigma^{2} \tau_{j} (\tau_{j} - 1)}{2} \cdot \frac{\eta_{j}^{3} L^$$

Note that when the learning rate satisfies:

$$\eta_j^2 L^2(\tau_j - 1) (2\beta + \tau_j) + \eta_j L\left(\frac{\beta}{m} + 1\right) \le 1,$$
(57)

we have

$$\mathbb{E}\left[F(\overline{\mathbf{x}}_{\tau_j+1}) - F(\overline{\mathbf{x}}_1)\right] \le -\frac{\eta_j}{2} \sum_{k=1}^{\tau_j} \mathbb{E}\left\|\nabla F(\overline{\mathbf{x}}_k)\right\|^2 + \frac{\eta_j^2 L \sigma^2 \tau_j}{2m} + \frac{\eta_j^3 L^2 \sigma^2 \tau_j (\tau_j - 1)}{2}.$$
 (58)

Suppose  $l_j = \sum_{r=0}^{j-1} \tau_r + 1$  is the first index in the j-th local update period. Without loss of generality, we substitute the local index by global index:

$$\mathbb{E}\left[F(\overline{\mathbf{x}}_{l_{j+1}}) - F(\overline{\mathbf{x}}_{l_j})\right] \le -\frac{\eta_j}{2} \sum_{k=1}^{\tau_j} \mathbb{E}\left\|\nabla F(\overline{\mathbf{x}}_{l_j+k-1})\right\|^2 + \frac{\eta_j^2 L \sigma^2 \tau_j}{2m} + \frac{\eta_j^3 L^2 \sigma^2 \tau_j (\tau_j - 1)}{2}.$$
 (59)

Summing over all local periods from j = 0 to j = R, one can obtain

$$\mathbb{E}\left[F(\overline{\mathbf{x}}_{l_R}) - F(\overline{\mathbf{x}}_1)\right] \le -\frac{1}{2} \sum_{j=0}^{R} \eta_j \sum_{k=1}^{\tau_j} \mathbb{E}\left\|\nabla F(\overline{\mathbf{x}}_{l_j+k-1})\right\|^2 + \frac{L\sigma^2}{2m} \sum_{j=0}^{R} \eta_j^2 \tau_j + \frac{L^2\sigma^2}{2} \sum_{j=0}^{R} \eta_j^3 \tau_j (\tau_j - 1). \tag{60}$$

After minor rearranging, it is easy to see

$$\mathbb{E}\left[\sum_{j=0}^{R} \eta_{j} \sum_{k=1}^{\tau_{j}} \left\|\nabla F(\overline{\mathbf{x}}_{l_{j}+k-1})\right\|^{2}\right] \leq 2\left[F(\overline{\mathbf{x}}_{1}) - F^{*}\right] + \frac{L\sigma^{2}}{m} \sum_{j=0}^{R} \eta_{j}^{2} \tau_{j} + L^{2}\sigma^{2} \sum_{j=0}^{R} \eta_{j}^{3} \tau_{j} (\tau_{j} - 1). \tag{61}$$

That is,

$$\mathbb{E}\left[\frac{\sum_{j=0}^{R} \eta_{j} \sum_{k=1}^{\tau_{j}} \left\|\nabla F(\overline{\mathbf{x}}_{l_{j}+k-1})\right\|^{2}}{\sum_{j=0}^{R} \eta_{j} \tau_{j}}\right] \leq \frac{2\left[F(\overline{\mathbf{x}}_{1}) - F^{*}\right]}{\sum_{j=0}^{R} \eta_{j} \tau_{j}} + \frac{L\sigma^{2}}{m} \frac{\sum_{j=0}^{R} \eta_{j}^{2} \tau_{j}}{\sum_{j=0}^{R} \eta_{j} \tau_{j}} + L^{2}\sigma^{2} \frac{\sum_{j=0}^{R} \eta_{j}^{3} \tau_{j}(\tau_{j}-1)}{\sum_{j=0}^{R-1} \eta_{j} \tau_{j}}.$$
 (62)

#### F.3 Asymptotic Result (Theorem 3)

In order to let the upper bound (62) converges to zero as  $R \to \infty$ , a sufficient condition is

$$\lim_{R \to \infty} \sum_{j=0}^{R} \eta_j \tau_j = \infty, \ \lim_{R \to \infty} \sum_{j=0}^{R} \eta_j^2 \tau_j < \infty, \ \lim_{R \to \infty} \sum_{j=0}^{R} \eta_j^3 \tau_j^2 < \infty.$$
 (63)

Here, we complete the proof of Theorem 3.

#### F.4 Simplified Result

We can obtain a simplified result when the learning rate is fixed. To be specific, we have

$$\mathbb{E}\left[\frac{\sum_{j=0}^{R} \sum_{k=1}^{\tau_{j}} \left\|\nabla F(\overline{\mathbf{x}}_{l_{j}+k-1})\right\|^{2}}{\sum_{j=0}^{R} \tau_{j}}\right] \leq \frac{2\left[F(\overline{\mathbf{x}}_{1}) - F^{*}\right]}{\eta \sum_{j=0}^{R} \tau_{j}} + \frac{\eta L \sigma^{2}}{m} \frac{\sum_{j=0}^{R} \tau_{j}}{\sum_{j=0}^{R} \tau_{j}} + \eta^{2} L^{2} \sigma^{2} \frac{\sum_{j=0}^{R-1} \tau_{j}(\tau_{j}-1)}{\sum_{j=0}^{R-1} \tau_{j}}$$
(64)

$$\leq \frac{2\left[F(\overline{\mathbf{x}}_{1}) - F^{*}\right]}{\eta \sum_{j=0}^{R} \tau_{j}} + \frac{\eta L \sigma^{2}}{m} + \eta^{2} L^{2} \sigma^{2} \left(\frac{\sum_{j=0}^{R} \tau_{j}(\tau_{j} - 1)}{\sum_{j=0}^{R-1} \tau_{j}}\right). \tag{65}$$

If we choose the total iterations  $K = \sum_{j=0}^{R} \tau_j$ , then

$$\mathbb{E}\left[\frac{\sum_{k=1}^{K} \|\nabla F(\overline{\mathbf{x}}_{k})\|^{2}}{K}\right] \leq \frac{2\left(F(\overline{\mathbf{x}}_{1}) - F^{*}\right)}{\eta K} + \frac{\eta L \sigma^{2}}{m} + \eta^{2} L^{2} \sigma^{2} \left(\frac{\sum_{j=0}^{R} \tau_{j}^{2}}{\sum_{j=0}^{R} \tau_{j}} - 1\right). \tag{66}$$