

# FED-LAMB: Layer-wise and Dimension-wise Locally Adaptive Optimization

Belhal Karimi, Xiaoyun Li, Ping Li

Cognitive Computing Lab  
Baidu Research  
10900 NE 8th St. Bellevue, WA 98004, USA  
{belhalkarimi, xiaoyunli, liping11}@baidu.com

## Abstract

In the emerging paradigm of federated learning (FL), large amount of clients, such as mobile devices, are used to train possibly high-dimensional models on their respective data. Due to the low bandwidth of mobile devices, decentralized optimization methods need to shift the computation burden from those clients to the computation server while preserving *privacy* and reasonable *communication cost*. In this paper, we focus on the training of deep, as in multi-layered, neural networks, under the FL settings. We present FED-LAMB, a novel federated learning method based on a *layer-wise* and *dimension-wise* updates of the local models, alleviating the nonconvexity and the multi-layered nature of the optimization task at hand. We provide a thorough finite-time convergence analysis for FED-LAMB characterizing how fast its gradient decreases, which improves the communication efficiency compared with the baseline method. We provide experimental results on various datasets and models, under both iid and non-iid settings, to show that the proposed FED-LAMB achieves faster convergence speed and better generalization performance, compared to the state-of-the-art.

# 1 Introduction

A growing and important task while learning models on observed data, is the ability to train over a large number of clients which could either be personal devices or distinct entities. In the paradigm of Federated Learning (FL) (Konečný et al., 2016; McMahan et al., 2017), a central server orchestrates the optimization over those clients under the constraint that the data can neither be centralized nor shared among the clients. This is more computationally efficient, since more computing resources are used; also, this is a very practical scenario which allows individual data holders (e.g., mobile devices) to train a model jointly without leaking the private data. In this paper, as most modern machine learning tasks, we consider the large finite-sum optimization problem written as

$$\min_{\theta \in \Theta} \frac{1}{n} \sum_{i=1}^n f_i(\theta), \quad (1)$$

where  $n$  denotes the number of workers,  $f_i$  represents the average loss for worker  $i$  and  $\theta$  the global model parameter taking value in  $\Theta$ , a subset of  $\mathbb{R}^p$ . While this formulation recalls that of standard distributed optimization, the core principle of FL is different from the traditional distributed paradigm, in that FL allows local models to perform multiple local updates on the local models before the global aggregation.

In this paper, we mainly address two important aspects of FL: communication efficiency and model performance, which, in some sense, are also coherent to each other. While local updates can effectively reduce the number of communication rounds between the central server and devices, new techniques are still necessary to tackle the challenge of communication between devices and server, due to, e.g., wireless bandwidth. Some quantization (Alistarh et al., 2017; Wangni et al., 2018), compression (Lin et al., 2018) and sketching (Rothchild et al.) methods allow to decrease the number of bits communicated at each round. Another direction to improve the communication cost of FL methods is to design better algorithms to accelerate local training, such that better local models are sent to the server at each round. This could lead to reduced communication (to reach a certain accuracy) as well as possibly improved overall learning performance. In this work, we will propose an accelerated local model training framework which achieves both communication reduction and improved empirical learning performance.

One of the most popular frameworks for FL is called Fed-SGD (McMahan et al., 2017): we adopt multiple local Stochastic Gradient Descent (SGD) steps in each device, send those local models to the server that computes the average over those received local model parameters, and broadcasts it back to the devices. As mentioned above, momentum can be added to the local SGD training for faster convergence and better learning performance (Yu et al., 2019). In this work, we focus on an alternative framework as proposed by Chen et al. (2020), where AMSGrad, a popular adaptive optimization method, is deployed locally instead of SGD. Adaptive methods have shown success in many deep learning tasks for fast convergence and good accuracy. Chen et al. (2020) showed that when adapted to the FL setting, the so-called “Local AMS” algorithm has communication cost sublinear in  $R$ , that is guaranteed to converge to a stationary point with rate  $\mathcal{O}(\sqrt{p/Rn})$ , where  $R$  is the number of communication rounds,  $p$  is the model dimensionality and  $n$  corresponds to the number of federated clients. Specifically, in Local AMS, each round the global server not only aggregates the local models, but also averages and broadcasts the local second moment estimations, which is a crucial ingredient in AMSGrad controlling the dimension-wise learning rates. Thus, this step can be regarded as a natural remedy to data heterogeneity (i.e., the data on local workers follows different probability distributions), which is a common scenario in practice that affects the performance of FL algorithms (Li et al., 2020; Liang et al., 2019; Karimireddy et al., 2019).

**Summary of Contributions.** Based on the recent progress in accelerating adaptive methods for efficient training (You et al., 2020), we propose an improved algorithm of Local AMSGrad (Chen et al., 2020), integrating both dimension-wise and layer-wise adaptive learning rates in each device’s local update. More specifically:

- We develop FED-LAMB, a novel optimization framework for federated learning, following a principled layer-wise adaptation strategy to accelerate training of deep neural networks. The algorithm enjoys fast convergence and good performance from the adaptivity from both AMSGrad and layer-wise learning rate adjustment.
- We provide theoretical analysis on the non-asymptotic convergence rate of FED-LAMB. Our rate,  $\mathcal{O}\left(\sqrt{\frac{p}{n}} \frac{1}{\sqrt{hR}}\right)$ , where  $h$  is the total number of layers and  $p$  denotes the dimension, matches the state of the art methods in federated learning and reaches a sublinear convergence in the total number of communication rounds. It also improves the theoretical communication efficiency compared with the baseline Local AMSGrad approach.
- We conducted numerical experiments under both homogeneous and heterogeneous settings on various benchmark datasets such as FMNIST, CIFAR-10 and TinyImagenet. Our results confirms the fast convergence and communication efficiency of the proposed method. In particular, FED-LAMB reaches similar, or better, test accuracy than the baseline local SGD and local AMS methods, with less number of communication rounds.

**Roadmap.** After establishing a literature review of both realms of federated and adaptive learning in Section 2, we develop in Section 3, our method, namely FED-LAMB, based on the computation per layer and per dimension, of the adaptive learning rate in the traditional AMSGrad. Theoretical understanding of our method’s behaviour with respect to convergence towards a stationary point is developed in Section 4 in nonconvex optimization. We present numerical experiments showing the effectiveness and advantages of our method in Section 5.

## 2 Related Work

Below, we summarize some relevant works on federated learning and adaptive optimization.

**Adaptive gradient methods.** Adaptive methods have proven to be the spearhead in many nonconvex optimization tasks. Gradient based optimization algorithms alleviate the possibly high nonconvexity of the objective function by adaptively updating each coordinate of their learning rate using past gradients. Common used examples include RMSprop (Tieleman and Hinton, 2012), Adadelta (Zeiler, 2012), Adam (Kingma and Ba, 2015), Nadam (Dozat, 2016) and AMSGrad (Reddi et al., 2018). Their popularity and efficiency are due to their great performance at training deep neural networks. They generally combine the idea of adaptivity from AdaGrad (Duchi et al., 2011; McMahan and Streeter, 2010), as explained above, and the idea of momentum from Nesterov’s Method (Nesterov, 2004) or Heavy ball method (Polyak, 1964) using past gradients. AdaGrad displays a great edge when the gradient is sparse compared to other classical methods. The anisotropic nature of this update represented a real breakthrough in the training of high dimensional and non-convex loss functions. This adaptive learning rate helps accelerate the convergence when the gradient vector is sparse (Duchi et al., 2011). Yet, when applying AdaGrad to train deep neural networks, it is observed that the learning rate might decay too fast, see Kingma and Ba (2015) for more details. Consequently, Kingma and Ba (2015) develops Adam leveraging a moving average of the gradients divided by the square root of the second moment of this moving average (element-wise

multiplication). A variant, called AMSGrad described in Reddi et al. (2018), ought to fix Adam failures using a max operator. Beyond improving the convergence speed of optimization methods, several studies including Zhou et al. (2020) also focus on improving their generalization properties. A natural extension of Adam has been developed in You et al. (2020) specifically for multi layered neural network where a principled layer-wise adaptation strategy is employed to accelerate training of deep neural networks using large mini-batches using either SGD or adaptive method under the setting of a classical single server. In simple terms, the idea is based on the observation that in a large deep neural network, the magnitude of the gradient might be too small in comparison with the magnitude of the weight for some layers of the model, hence slowing down the overall convergence. As a consequence, layer-wise adaptive learning rate is applied, such that in each iteration the model can move sufficiently far. This method empirically speeds up the convergence significantly in classical sequential models and can be provably faster than baseline methods.

**Federated learning.** An extension of the well known parameter server framework, where a model is being trained on several servers in a distributed manner, is called federated learning (FL), see Konečný et al. (2016). Here, the central server only plays the role of computing power for aggregation and global update of the model. Compared with the distributed learning paradigm, in federated learning, the data stored in each worker must not be seen by the central server – preserving privacy is key – and the nature of those workers (e.g., mobile devices), combined with their usually large amount, makes communication between the devices and the central server less appealing – communication cost needs to be controlled. Thus, while traditional distributed gradient methods (Recht et al., 2011; Li et al., 2014; Zhao et al., 2020) do not respect those constraints, it has been proposed in McMahan et al. (2017), an algorithm called Federated Averaging – Fed-SGD – extending parallel SGD with local updates performed on each device. In Fed-SGD, each worker updates their own model parameters locally using SGD, and the local models are synchronized by periodic averaging on the central parameter server.

### 3 Layer-wise and Dimension-wise Adaptive Optimization

**Notations.** We denote by  $\theta$  the vector of parameters taking values in  $\mathbb{R}^d$ . For each layer  $\ell \in [\![h]\!]$ , where  $h$  is the total number of layers of the neural networks, and each coordinate  $j \in [\![p_\ell]\!]$  where  $p_\ell$  is the dimension per layer  $\ell$  ( $p := \sum_{\ell=1}^h p_\ell$  denotes the total dimension). We also note  $\theta_{r,i}^{\ell,t}$  its value for layer  $\ell$  at round  $r$ , local iteration  $t$  and for worker  $i$ . The gradient of  $f$  with respect to  $\theta^\ell$  is denoted by  $\nabla_\ell f(\theta)$ . The smoothness per layer is denoted by  $L_\ell$  for each layer  $\ell \in [\![h]\!]$ .

#### 3.1 AMSGrad, Local AMSGrad and Periodic Averaging

Under the federated learning setting, we stress on the importance of reducing, at each round, the communication cost between the central server, used mainly for aggregation purposes, and the many clients used for gradient computation and local updates. Using Periodic Averaging after few local epochs, updating local models on each device, as developed in McMahan et al. (2017) is the gold standard for achieving such communication cost reduction. Intuitively, one rather shifts the computation burden from the many clients to the central server as much as possible. This technique allows for fewer local epochs and a better global model, from a loss minimization (or model fitting) perspective. The premises of that new paradigm are SGD updates performed locally on each device then averaged periodically, see Konečný et al. (2016); Zhou and Cong (2018). The heuristic efficiency of local updates using SGD and periodic averaging has been studied in Stich (2019); Yu et al. (2019) and shown to reach a similar sublinear convergence rate as in the standard

distributed optimization settings. Then, with the growing need of training far more complex models, e.g., deep neural networks, several efficient methods, built upon adaptive gradient algorithms, such as Local AMSGrad in [Chen et al. \(2020\)](#), extended both empirically and theoretically the benefits of performing local updates coupled with periodic averaging, especially when dealing with non-iid local data distribution. Thus, it can be a better alternative than local SGD on many learning tasks.

### 3.2 Layer-wise and Dimension-wise Adaptive Local Update

For training large deep neural networks, [You et al. \(2020\)](#) proposed LAMB, an acceleration framework for both SGD and adaptive algorithm that allows large-batch training of BERT in hours. The idea is based on the observation that, during the training of large neural networks, different layers may have very different (absolute) scales, while the magnitude of the gradients in these layers may be similar. Intuitively, this means that in some iterations, the parameter may move a too small step even when the direction (negative gradient) is correct. Therefore, [You et al. \(2020\)](#) proposed to normalize the gradients in each layer according to the scale of the network layer. Effectively, the algorithm assigns different learning rates to different layers. Empirically, this strategy significantly accelerates the convergence for training large deep models.

Inspired by the merits of local adaptive methods and the success of layer-wise adaptive learning rates, we propose a layer-wise and dimension-wise local AMS algorithm which is detailed in Algorithm 1 and depicted in Figure 1. The proposed algorithm is a natural adaptation of the vanilla AMSGrad method, for multi-layer neural networks under the federated learning setting. Here, “multi-layer” includes a broad class of network architectures that can be parameterized layer-by-

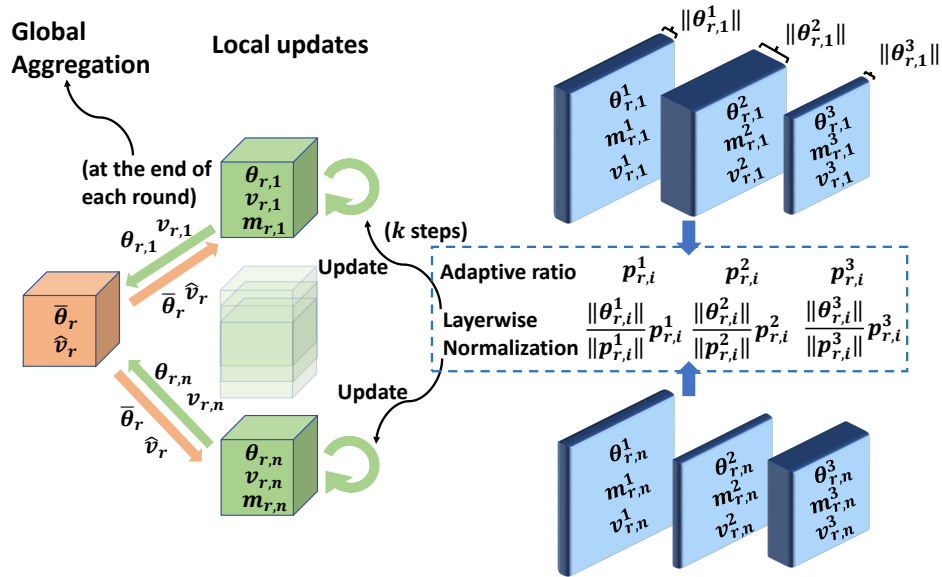


Figure 1: Illustration of FED-LAMB (Algorithm 1), with a three-layer network and  $\phi(x) = x$  as an example. For device  $i$  and each local iteration in round  $r$ , the adaptive ratio of  $j$ -th layer  $p_{r,i}^j$  is normalized according to  $\|\theta_{r,i}^j\|$ , and then used for updating the local model. At the end of each round  $r$ , local worker  $i$  sends  $\theta_{r,i} = [\theta_{r,i}^\ell]_{\ell=1}^h$  and  $v_{r,i}$  to the central server, which transmits back aggregated  $\theta$  and  $\hat{v}$  to local devices to complete a round of training.

---

**Algorithm 1** FED-LAMB for federated learning

---

1: **Input:** parameter  $0 < \beta_1, \beta_2 < 1$ , and learning rate  $\alpha_t$ , weight decaying parameter  $\lambda \in [0, 1]$ .  
2: **Initialize:**  $\theta_{0,i} \in \Theta \subseteq \mathbb{R}^d$ ,  $m_{0,i}^0 = \hat{v}_{0,i}^0 = v_{0,i}^0 = 0$ ,  $\forall i \in \llbracket n \rrbracket$ , and  $\bar{\theta}_0 = \frac{1}{n} \sum_{i=1}^n \theta_{0,i}$ .  
3: **for**  $r = 1, \dots, R$  **do**  
4:   **parallel for device**  $i$  **do**  
5:     Set  $\theta_{r,i}^0 = \bar{\theta}_{r-1}$ .  
6:     Set  $m_{r,i}^0 = m_{r-1,i}^T$ ,  $v_{r,i}^0 = \hat{v}_{r-1}$ .  
7:     **for**  $t = 1, \dots, T$  **do**  
8:       Compute stochastic gradient  $g_{r,i}^t$  at  $\theta_{r,i}^0$ .  
9:        $m_{r,i}^t = \beta_1 m_{r,i}^{t-1} + (1 - \beta_1) g_{r,i}^t$  and  $m_{r,i}^t = m_{r,i}^t / (1 - \beta_1^t)$ .  
10:        $v_{r,i}^t = \beta_2 v_{r-1,i}^t + (1 - \beta_2) (g_{r,i}^t)^2$  and  $v_{r,i}^t = v_{r,i}^t / (1 - \beta_2^t)$ .  
11:       Compute the ratio  $p_{r,i}^t = m_{r,i}^t / (\sqrt{\hat{v}_r^t} + \epsilon)$ .  
12:       Update local model for each layer  $\ell \in \llbracket h \rrbracket$ :

$$\theta_{r,i}^{\ell,t} = \theta_{r,i}^{\ell,t-1} - \alpha_r \phi(\|\theta_{r,i}^{\ell,t-1}\|) (p_{r,i}^{\ell,t} + \lambda \theta_{r,i}^{\ell,t-1}) / \|p_{r,i}^{\ell,t} + \lambda \theta_{r,i}^{\ell,t-1}\|. \quad (2)$$

13:     **end for**  
14:     Devices send  $\theta_{r,i}^T = [\theta_{r,i}^{\ell,T}]_{\ell=1}^h$  and  $v_{r,i}^T$  to server.  
15:     **end for**  
16:     Server computes averages of the local models  $\bar{\theta}_r = [\bar{\theta}_r^{\ell,T}]_{\ell=1}^h = [\frac{1}{n} \sum_{i=1}^n \theta_{r,i}^{\ell,T}]_{\ell=1}^h$  and  $\hat{v}_{r+1} = \max(\hat{v}_r, \frac{1}{n} \sum_{i=1}^n v_{r,i}^T)$  and send them back to the devices.  
17:     **end for**  
18: **Output:** Global model parameter  $\bar{\theta}_R = [\bar{\theta}_R^{\ell,T}]_{\ell=1}^h$ .

---

layer (e.g., CNN, ResNet, Transformers). Lines 8-11 in Algorithm 1 correspond to the standard first and second moment estimation and correction performed by AMSGrad locally. The local model update rule (2) in Line 12 incorporates the layer-wise normalization according to the magnitude of each layer's weights. For example, taking  $\phi$  as the identity function and weight decay  $\lambda = 0$ , then the model  $\theta$  is updated by  $-\frac{\alpha \|\theta\|}{\|p\|} \cdot p$ , instead of  $-\alpha p$  as in the standard AMSGrad update. The norm of each update is precisely  $\alpha \|\theta\|$ . That is, we force the change in each layer's parameter to have norm proportional to the scale of the layer's weight. This is a novel extension of the layer-wise adaptivity to the federated learning framework. Therefore, the proposed FED-LAMB exhibits two-fold adaptivity—a dimension-wise adaptive learning rate with respect to the square root of the second moment used in AMSGrad, and a layer-wise adaptive normalization brought by LAMB. Such two-fold adaptivity has not been considered in FL literature, and we will demonstrate the advantage of this scheme theoretically and empirically in the remainder of the paper.

## 4 Convergence of FED-LAMB

In this section, we develop the theoretical analysis of Algorithm 1. Based on classical result for stochastic nonconvex optimization, we present a collection of results that aims to providing a better understanding of the convergence behavior of our distributed optimization method under the federated learning framework. The main challenges we ought to overcome are manifold: (i) The large amount of decentralized workers working solely on their own data stored locally. (ii) A periodic averaging occurs on the central server pushing each of those clients to send local models after some local iterations. (iii) Each client computes a backpropagation of the main model, i.e., the deep

neural network, and then updates its local version of the model via an adaptive gradient method: the distinctiveness being that those updates are done *dimension-wise* and *layer-wise*. Our analysis encompasses the consideration of those challenges and leads to an informative convergence rates depending on the quantities of interest in our problem: the number of layers of the DNN, the number of communications rounds and the number of clients used under our federated settings.

#### 4.1 Finite Time Analysis of FED-LAMB

In the sequel, the analysis of our scheme we provide is *global*, in the sense that it does not depend on the initialization of our algorithm, and *finite-time*, meaning that it is true for any arbitrary number of communication rounds, in particular small ones. In the particular context of nonconvex stochastic optimization for distributed clients, we assume the following:

**H1.** (*Smoothness per layer*) For  $i \in \llbracket n \rrbracket$  and  $\ell \in \llbracket L \rrbracket$ :  $\|\nabla f_i(\theta^\ell) - \nabla f_i(\vartheta^\ell)\| \leq L_\ell \|\theta^\ell - \vartheta^\ell\|$ .

We add some classical assumption on the gradient of the objective function:

**H2.** (*Unbiased and Bounded gradient*) The stochastic gradient is unbiased for any iteration  $r > 0$ :  $\mathbb{E}[g_r] = \nabla f(\theta_r)$  and is bounded from above, i.e.,  $\|g_t\| \leq M$ .

**H3.** (*Bounded variance*) The variance of the stochastic gradient is bounded for any iteration  $r > 0$  and any dimension  $j \in \llbracket d \rrbracket$ :  $\mathbb{E}[|g_r^j - \nabla f(\theta_r)^j|^2] < \sigma^2$ .

**H4.** (*Bounded Scale*) For any  $a \in \mathbb{R}_+^*$ , there exists strictly positive constants such that  $\phi_m \leq \phi(a) \leq \phi_M$ .

Two important Lemmas are required in the proof of the Theorem above. We also report the complete proof of our bound in the Appendix of this paper.

The first result gives a characterization of the gap between the averaged model, that is computed by the central server in a periodic manner, and each of the local models stored in each client  $i \in \llbracket n \rrbracket$ .

**Lemma 1.** Consider  $\{\bar{\theta}_r\}_{r>0}$ , the sequence of parameters obtained running Algorithm 1. Then for  $i \in \llbracket n \rrbracket$  and  $r > 0$ , the gap  $\|\bar{\theta}_r - \theta_{r,i}\|^2$  satisfies:

$$\|\bar{\theta}_r - \theta_{r,i}\|^2 \leq \alpha_r^2 M^2 \phi_M^2 \frac{(1 - \beta_2)p}{v_0},$$

where  $\phi_M$  is defined in H4 and  $p$  is the total number of dimensions  $p = \sum_{\ell=1}^h p_\ell$ .

The gap is provably bounded by some quantities of interest such as the total dimension of the multi-layered model  $p$ , the learning rate and the assumed upper bound of the gradient, see H2.

Note that the end goal is to characterize how fast the gradient of the averaged/global parameter  $\bar{\theta}_r$  goes to zero, but not the averaged gradient. The following Lemma allows us to convert the suboptimality condition  $\left\| \frac{\bar{\nabla} f(\theta_r)}{\sqrt{v_r^t}} \right\|$  to the desired one which is  $\left\| \frac{\nabla f(\bar{\theta}_r)}{\sqrt{v_r^t}} \right\|$ .

**Lemma 2.** Consider  $\{\bar{\theta}_r\}_{r>0}$ , the sequence of parameters obtained running Algorithm 1. Then for  $r > 0$ :

$$\left\| \frac{\bar{\nabla} f(\theta_r)}{\sqrt{v_r^t}} \right\|^2 \geq \frac{1}{2} \left\| \frac{\nabla f(\bar{\theta}_r)}{\sqrt{v_r^t}} \right\|^2 - \bar{L} \alpha^2 M^2 \phi_M^2 \frac{(1 - \beta_2)p}{v_0},$$

where  $M$  is defined in H2,  $p = \sum_{\ell=1}^h p_\ell$  and  $\phi_M$  is defined in H4.



We now state our main result regarding the non asymptotic convergence analysis of our Algorithm 1 for multiple local updates and true for any communication rounds number  $R$ .

**Theorem 1.** Assume **H1-H4**. Consider  $\{\bar{\theta}_r\}_{r>0}$ , the sequence of parameters obtained running Algorithm 1 with a constant learning rate  $\alpha$ . Let the number of local epochs be  $T \geq 1$  and  $\lambda = 0$ . Then, for any round  $R > 0$ , we have

$$\begin{aligned} \frac{1}{R} \sum_{r=1}^R \mathbb{E} \left[ \left\| \frac{\nabla f(\bar{\theta}_r)}{\hat{v}_r^{1/4}} \right\|^2 \right] &\leq \sqrt{\frac{M^2 p}{n}} \frac{\Delta}{h \alpha R} + 4\alpha \left[ \frac{\alpha^2 L}{\sqrt{v_0}} M^2 (T-1)^2 \phi_M^2 (1-\beta_2) p + \phi_M^2 \sqrt{M^2 + p \sigma^2} \right] \\ &\quad + 4\alpha \frac{M^2}{\sqrt{v_0}} + \frac{\phi_M \sigma^2}{R n} \sqrt{\frac{1-\beta_2}{M^2 p}} + 4\alpha \left[ \phi_M \frac{h \sigma^2}{\sqrt{n}} \right] + cst, \end{aligned} \quad (3)$$

where  $\Delta = \mathbb{E}[f(\bar{\theta}_1)] - \min_{\theta \in \Theta} f(\theta)$ .

By choosing a suitable decreasing learning rate, we have the following simplified result.

**Corollary 1.** Under the same setting as Theorem 1, with  $\alpha = \mathcal{O}(\frac{1}{\sqrt{hR}})$ , it holds that

$$\frac{1}{R} \sum_{r=1}^R \mathbb{E} \left[ \left\| \frac{\nabla f(\bar{\theta}_r)}{\hat{v}_t^{1/4}} \right\|^2 \right] \leq \mathcal{O} \left( \sqrt{\frac{p}{n}} \frac{1}{\sqrt{hR}} + \frac{\sigma^2}{R n \sqrt{p}} + \frac{(T-1)^2 p}{h^3 R^{3/2}} \right). \quad (4)$$

The leading two terms display a dependence of the convergence rate of FED-LAMB on the initialization and the variance of the stochastic gradient (see H3), which are common in distributed optimization. The last term involves the number of local updates which relates to the communication efficiency. More discussion will be provided next.

## 4.2 Comparisons

We dedicate the following paragraph to a discussion on the bound (and implications) derived above in comparison with known results most relevant to our interest in literature.

**LAMB bound in You et al. (2020):** We first start our discussion with the comparison of convergence rate of FED-LAMB with that of LAMB, Theorem 3 in You et al. (2020). The convergence rates of FED-LAMB and LAMB differ in two ways: (i) First, note that the characterization, on the suboptimality, or convergence criterion, is given at the averaged parameters noted  $\bar{\theta}_r$  due to our distributed settings. It is thus natural to consider the evolution of our objective function, precisely its gradient, evaluated at some global model values –as opposed to the outcome of a single step drift in the central server paradigm. Besides, for ease of interpretation, the LHS of (3) is summed over all rounds instead of a fictive random termination point. A simple calculation would lead to such characterization found in several nonconvex stochastic optimization paper such as Ghadimi and Lan (2013). (ii) Assuming that the convergence criterion in both Theorems is of similar order (which happens for a large enough number of rounds), the convergence rate of FED-LAMB displays a similar  $\mathcal{O}(1/R)$  behaviour for the initialization term. That said, despite the distributed (federated) setting, our dimension-wise and layer-wise method benefits from the double adaptivity phenomenon explained above and exhibited in LAMB (You et al., 2020), under a central server setting.

**Fed-AMS bound in Chen et al. (2020):** We now discuss the similarities and differences between Fed-AMS, the baseline distributed adaptive method developed in Chen et al. (2020), and our FED-LAMB. For clarity, we restate their main result under our notations:



**Theorem 2** (Corollary 5.2 in [Chen et al. \(2020\)](#)). *Under some regularity conditions on the local losses and similar assumption as ours, with some properly chosen learning rate, when  $R \geq \mathcal{O}(T^3\sqrt{n})$ , Fed-AMS has convergence rate*

$$\frac{1}{R} \sum_{r=1}^R \mathbb{E} \left[ \left\| \frac{\nabla f(\bar{\theta}_r)}{\sqrt{v_r}} \right\|^2 \right] \leq \mathcal{O}\left(\frac{\sqrt{p}}{\sqrt{nR}}\right). \quad (5)$$

Firstly, when the number of rounds  $R$  is sufficiently large, both rates (4) and (5) are dominated by  $\mathcal{O}(\frac{\sqrt{p}}{\sqrt{nR}})$ , matching the convergence rate of the standard AMSGrad (e.g. [Zhou et al. \(2018\)](#)). The acceleration of our layer-wise scheme is exhibited in the  $\mathcal{O}(1/(nR))$  term and the dependence on the number of layers  $\mathcal{O}(1/(h^3 R^{3/2}))$  term. Note that the boundedness assumption is on each dimension in H3 and leads to a manifestation of the term  $\sqrt{p}$  in both rates. This dependency on  $p$  can be removed when H3 is assumed globally, which is also common in optimization literature.

Secondly, in (4), the last term containing the number of local updates  $T$  is small as long as  $T \leq \mathcal{O}(\frac{R^{1/2}h^{5/4}}{(np)^{1/4}})$ . Treating  $p^{1/4}/h = \mathcal{O}(1)$ , the result implies that for a given  $T$ , we can get the same rate of convergence as vanilla AMSGrad with  $R \geq \mathcal{O}(T^2\sqrt{n})$  rounds of communication. From Theorem 2, we know that Fed-AMS requires  $R \geq \mathcal{O}(T^3\sqrt{n})$  rounds to achieve the same rate. This implies that, our FED-LAMB reduces the number of communication rounds needed to reach an  $\epsilon$ -stationary point compared with [Chen et al. \(2020\)](#). Therefore, by leveraging the layer-wise acceleration on local models, our FED-LAMB improve the communication cost of Fed-AMS.

### 4.3 More Discussion

We provide more discussion on the algorithmic and theoretical properties of FED-LAMB from the following aspects:

**Communication and Client Efficiency:** The (sublinear) dependence on the number of communication rounds of our bound matches that of most recent methods in federated learning, e.g. SCAFFOLD ([Karimireddy et al., 2019](#)), a solution to the problem posed by heterogeneity of the data in each client. Yet, contrary to SCAFFOLD, our method only sends bits once per communication round while SCAFFOLD needs to send two vectors, including an additional control variate term from the clients to the central server. Our result also matches the communication bound of [Reddi et al. \(2021\)](#) which adapts Adam ([Kingma and Ba, 2015](#)) to the federated setting. However, the algorithm of [Reddi et al. \(2021\)](#) performs adaptive updates only at the central server, while SGD is still used for local updates. In addition, the  $1/\sqrt{n}$  term in convergence rate of FED-LAMB implies a linear speedup effect in the number of local workers, which also matches the dependency on  $n$  of most federated learning methods.

**Data Heterogeneity:** To demonstrate the effect of non-iid data distribution theoretically, some related works pose assumptions on the global variance and the local variance of the stochastic gradients of the objective function (1) separately, such that both variances appears in the convergence rate. Our analysis can also be easily extended to incorporate the global variance term. While some works including [Karimireddy et al. \(2019\)](#) target on designing specific strategies to alleviate the negative influence of data heterogeneity, we note that our FED-LAMB are, in some sense, naturally capable of balancing the heterogeneity in different local data distributions. As mentioned before, this is largely due to the “moment averaging” step (line 16 in Algorithm 1), where the adaptive

learning rates guided by the second moment estimation are aggregated among local workers periodically. In our experiments, we will see that the advantage of FED-LAMB is consistent under both homogeneous and heterogeneous data setting.

## 5 Numerical Experiments

In this section, we conduct numerical experiments on various datasets and network architectures to justify the effectiveness of our proposed method in practice. Our main objective is to validate the benefit of dimension-wise adaptive learning when integrated with the locally adaptive FL method. We observe that our method empirically confirms its edge in terms of convergence speed. Basically, FED-LAMB reduces the number of rounds and thus the communication cost required to achieve a similar stationary point (or test accuracy) than the baseline methods. In many cases, FED-LAMB could also bring notable improvement in the generalization performance.

**Methods and tuning.** In our experiments, we evaluate four federated learning algorithms:

1. Fed-SGD (McMahan et al., 2017), standard federated averaging with local SGD updates.
2. Fed-AMS (Chen et al., 2020), locally adaptive AMSGrad.
3. Adp-Fed (“Adaptive Federated Optimization”), the federated adaptive algorithm proposed by Reddi et al. (2021). The algorithm performs local SGD updates. At the end of each round  $r$ , the changes in local models,  $\Delta_i = w_{r,i}^T - w_{r,i}^0$ ,  $i = 1, \dots, n$ , are transformed to the central server for an aggregated Adam update. See Appendix A for a precise algorithm description.
4. Our proposed FED-LAMB (Algorithm 1), layer-wise accelerated local AMSGrad.

For all the adaptive methods, we set hyper-parameters  $\beta_1 = 0.9$ ,  $\beta_2 = 0.999$  as the recommended default (Reddi et al., 2018). Regarding the federated learning environment, we use 50 local workers with 0.5 participation rate. That is, we randomly pick half of the workers to be active for training in each round. We set the local mini-batch size as 128. In each round, the training samples are allocated to the active devices, and one local epoch is finished after all the local devices run one epoch over their received samples by mini-batch training.

We tune the initial learning rate  $\alpha$  for each algorithm over a fine grid. For Adp-Fed, there are two learning rates involved. We tune the combination of local learning rate (for local SGD) and global learning rate (for global Adam) over a fine grid. For FED-LAMB the parameter  $\lambda$  in Algorithm 1 controlling the overall scale of the layer-wise gradients is tuned from  $\{0, 0.01, 0.1\}$ , and use the identity function for  $\phi(\cdot)$ . The detailed parameter tuning diagram can be found in Appendix A. For each run, we report the test accuracy with the best  $\alpha$  and  $\lambda$ . The results are averaged over three independent runs, each with same initialization for every method.

**Datasets and models.** We run our experiments on four popular benchmark image classification datasets: MNIST (LeCun, 1998), Fashion MNIST (FMNIST) (Xiao et al., 2017), CIFAR-10 (Krizhevsky, 2009) and TinyImageNet (Deng et al., 2009). The MNIST dataset contains 60000 training samples and 10000 test samples, from 10 classes of handwritten digits from 0 to 9. The FMNIST dataset has the same sample size and train/test split as MNIST, but the samples are fashion products (e.g., dress, bags) which makes it harder to train than MNIST. The CIFAR-10 dataset has 50000 training images and 10000 test images, from 10 classes. The TinyImageNet is a subset of the ImageNet dataset. It includes 100000  $64 \times 64$  images for from 200 classes for training and 10000 for testing. For MNIST, we apply 1) a simple multi-layer perceptron (MLP), which has

one hidden layer containing 200 cells with dropout; 2) Convolutional Neural Network (CNN), which has two max-pooled convolutional layers followed by a dropout layer and two fully-connected layers with 320 and 50 cells respectively. This CNN is also implemented for FMNIST. For CIFAR-10 and TinyImageNet, we implement ResNet-18 network [He et al. \(2016\)](#).

### 5.1 Comparison under iid setting

In Figure 2, we report the test accuracy of MLP trained on MNIST, as well as CNN trained on MNIST and FMNIST, where the data is i.i.d. allocated among the clients. We test 1 local epoch and 3 local epochs (more local iterations). In all the figures, we observe a clear advantage of FED-LAMB over the competing methods in terms of the convergence speed. In particular, we can see that FED-LAMB is able to achieve the same accuracy with fewest number of communication rounds, thus improving the communication efficiency. For instance, this can be observed as follows: on MNIST + CNN (1 local epoch), Fed-AMS requires 20 rounds to achieves 90% accuracy, while FED-LAMB only takes 5 rounds. This implies a 75% reduction in the communication cost. Moreover, on MNIST, FED-LAMB also leads to improved generalisation performance, i.e., test accuracy. Note that, the result on MLP to a good extent provides a straightforward illustration on the benefit of layer-wise adaptivity in FED-LAMB since compared with Fed-AMS, the only difference is that the learning rate becomes adaptive to the scale of the single hidden layer in FED-LAMB. Lastly, we see a faster convergence with 3 local epochs. For iid data distribution, this is expected and consistent with previous results in literature (e.g., [Reddi et al. \(2021\)](#)).

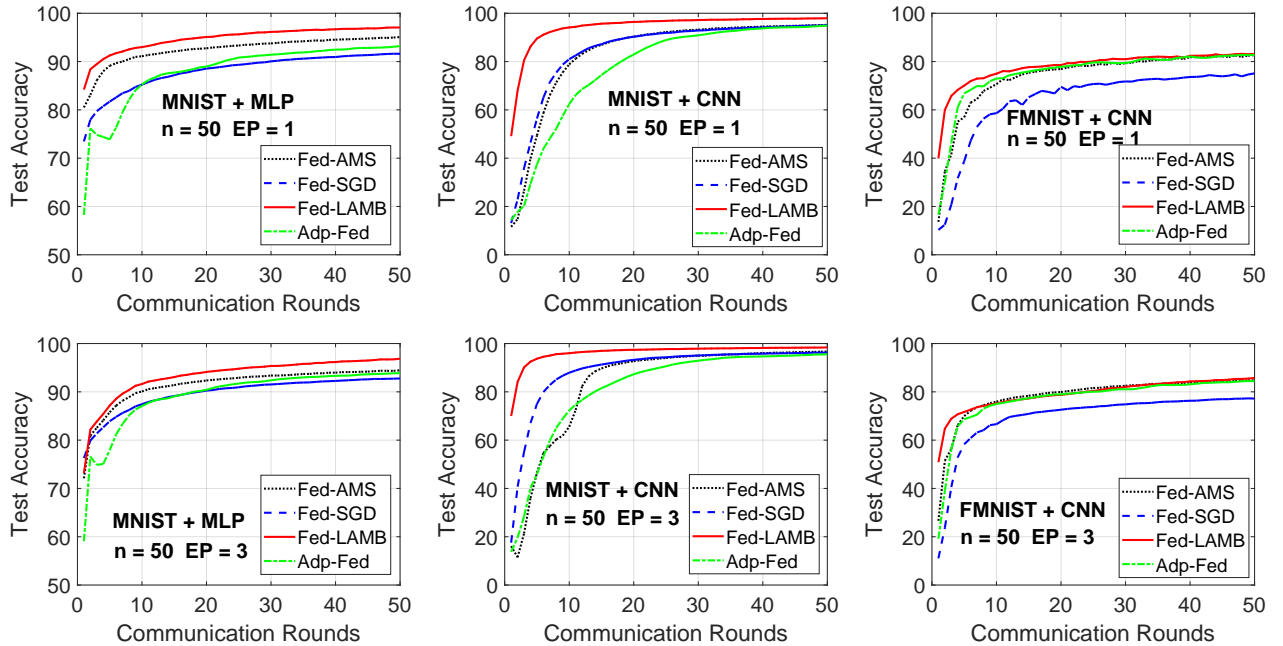


Figure 2: **i.i.d. data setting.** Test accuracy on MNIST and FMNIST against the number of communication rounds with 50 clients. **1st row:** 1 local epoch. **2nd row:** 3 local epochs. We see that FED-LAMB converges faster to better solution (higher test accuracy) in all cases.

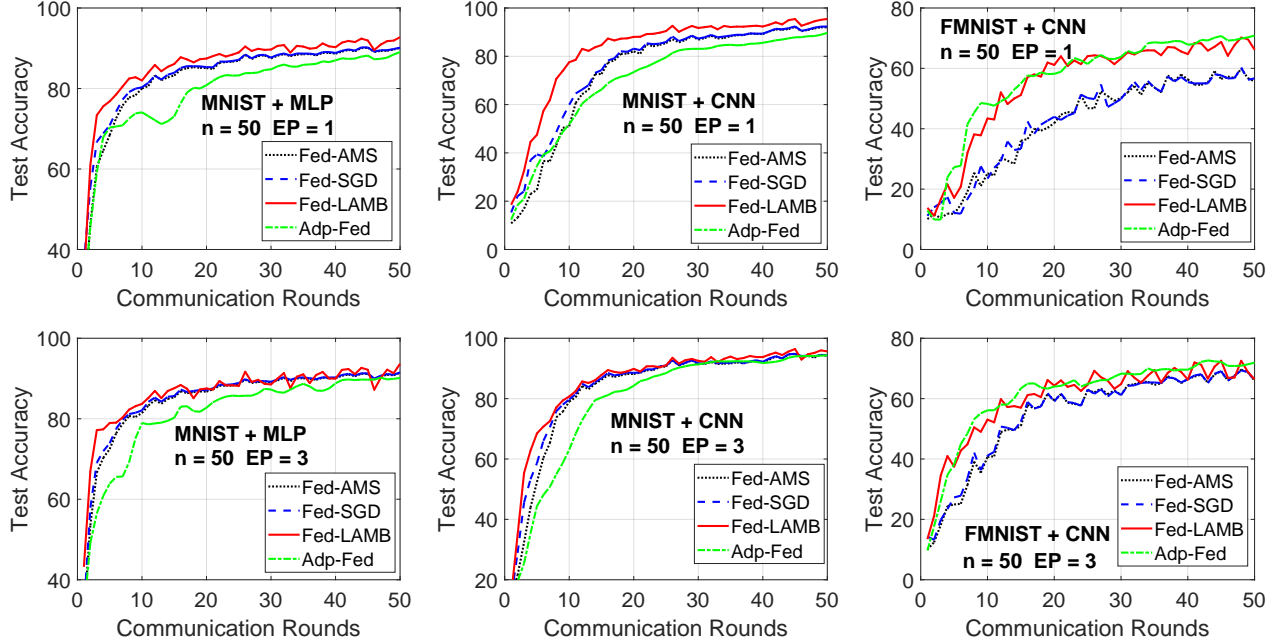


Figure 3: **non-i.i.d. data setting.** Test accuracy on MNIST and FMNIST against the number of communication rounds, with 50 clients. **1st row:** 1 local epoch. **2nd row:** 3 local epochs. We see that FED-LAMB converges faster to better solution (higher test accuracy) in all cases.

## 5.2 Comparison under non-iid setting

In Figure 3, we provide the results on MNIST and FMNIST dataset, when the local data has non-iid distribution (i.e., under data heterogeneity). In particular, in each round of federated training, every local device only receives samples from one class (out of ten). This is known to be the scenario where federated learning is harder to generalize well (McMahan et al., 2017), thus an important case for the empirical evaluation of FL methods.

First of all, from Figure 3, we see that for experiments with 1 local epoch, in all cases our proposed FED-LAMB outperforms the three baseline methods. Similar to the plots for iid data setting, FED-LAMB provides faster convergence speed and achieves higher test accuracy than Fed-SGD and Fed-AMS. The advantage is especially significant for the CNN model, e.g., it improves the accuracy of Fed-SGD and Fed-AMS by more than 10% on FMNIST. On both two datasets, FED-LAMB saves around 50% communication rounds to reach a same accuracy level as Fed-AMS. The other baseline method, Adp-Fed, performs as good as our FED-LAMB on FMNIST, but worse than other methods on MNIST.

The relative comparison is basically the same when we conduct 3 local epochs. Yet, the advantage of FED-LAMB becomes less significant than what we observed in Figure 2 with iid local data distribution. One plausible reason is that when the local data is highly non-iid as in our case, the fast convergence of the local models as in FED-LAMB might not be as beneficial when we allow too many local updates per round. Intuitively speaking, learning the local models “too fast” might not be a good thing to the globally aggregated model, since each local model is trained only with a few classes of data points, i.e., local models target at different loss functions.

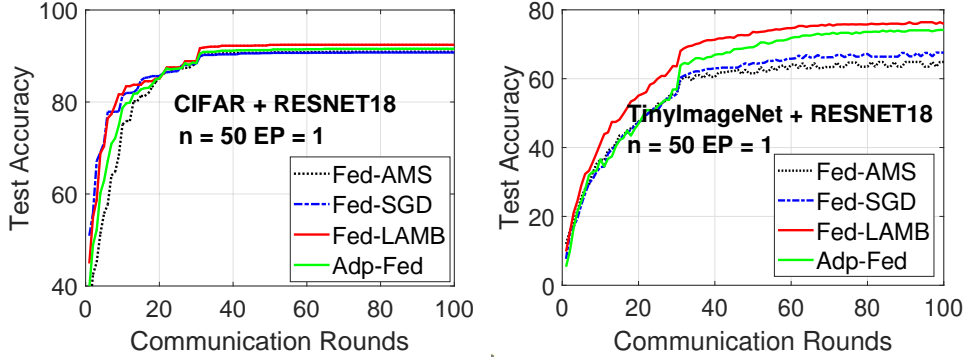


Figure 4: **non-i.i.d. data setting.** Test accuracy on CIFAR-10 + ResNet-18 and TinyImageNet + ResNet-18 with 50 clients.

In Figure 4, we present the experiment results on CIFAR-10 and TinyImageNet datasets trained by ResNet-18. When training these two models, we decrease the learning rate by 10 at 30 and 70 communication rounds. From Figure 4, we can draw similar conclusion as before: the proposed FED-LAMB is the best method in terms of both convergence speed and generalization accuracy. In particular, on TinyImageNet, we see that FED-LAMB has a significant advantage over all three baselines. Although Adp-Fed performs better than Fed-SGD and Fed-AMS, it is considerably worse than FED-LAMB. For reference, we also report the overall test accuracy at the end of training in Table 1. FED-LAMB achieves the highest accuracy on both datasets. On TinyImageNet, FED-LAMB reaches 76% after 100 communication rounds, against around 65% for Fed-AMS, 68% for Fed-SGD and 74% for Adp-Fed.

Table 1: Test Accuracy on ResNet-18 Network.

	Fed-SGD	Fed-AMS	Adp-Fed	<b>FED-LAMB</b>
CIFAR-10	90.75 $\pm$ 0.48	90.93 $\pm$ 0.22	91.57 $\pm$ 0.38	92.44 $\pm$ 0.53
TinyImageNet	67.58 $\pm$ 0.21	64.86 $\pm$ 0.83	74.17 $\pm$ 0.43	76.00 $\pm$ 0.26

## 6 Conclusion

We study a doubly adaptive method in the particular framework of federated learning. Built upon the success of periodic averaging, and of state-of-the-art adaptive gradient methods for single server nonconvex stochastic optimization, we derive FED-LAMB, a distributed AMSGrad method that performs local updates on each worker and periodically averages local models. When the trained model is a deep neural network, a core component of our method, FED-LAMB, is a *layer-wise* update of each local model. The main contribution of our paper is thus a federated learning optimization algorithm that leverages a double level of adaptivity: the first one stemming from a *dimension-wise* adaptivity of adaptive gradient methods, extended to their distributed (and local) counterpart, the second one is due to a *layer-wise* adaptivity making use of the particular compositionality of the considered model. Proved convergence guarantees of our scheme are provided in our contribution, and exhibit a sublinear dependence on the total number of communications rounds, and a linear speedup against the number of clients. Extensive experiments on various datasets and models, under both iid and non-iid data settings, validates that FED-LAMB is able to provide faster convergence which in turn could lead to reduced communication cost. Moreover, in many cases, FED-LAMB can also improve the overall performance of federated learning over prior methods.

# Appendix

## A Hyper-parameter Tuning and Algorithms

### A.1 The Adp-Fed Algorithm (Reddi et al., 2021)

The Adp-Fed (Adaptive Federated Optimization) is one of the baseline methods compared with FED-LAMB in our paper. The algorithm is given in Algorithm 2. The key difference between Adp-Fed and Fed-AMS (Chen et al., 2020) is that, in Adp-Fed, each local worker runs local SGD (Line 8), and an Adam optimizer is maintained for the global adaptive optimization (Line 15). In the Fed-AMS framework (as well as our FED-LAMB), each clients runs local (adaptive) AMSGrad method, and the global model is simply obtained by averaging the local models.

---

#### Algorithm 2 Adaptive Federated Optimization (Reddi et al., 2021)

---

```

1: Input: parameter  $0 < \beta_1, \beta_2 < 1$ , and learning rate  $\alpha_t$ , weight decaying parameter  $\lambda \in [0, 1]$ .
2: Initialize:  $\theta_{0,i} \in \Theta \subseteq \mathbb{R}^d$ ,  $m_0 = 0$ ,  $v_0 = 0$ ,  $\forall i \in \llbracket n \rrbracket$ , and  $\theta_0 = \frac{1}{n} \sum_{i=1}^n \theta_{0,i}$ .
3: for  $r = 1, \dots, R$  do
4:   parallel for device  $i$  do:
5:     Set  $\theta_{r,i}^0 = \theta_{r-1}$ .
6:     for  $t = 1, \dots, T$  do
7:       Compute stochastic gradient  $g_{r,i}^t$  at  $\theta_{r,i}^0$ .
8:        $\theta_{r,i}^t = \theta_{r,i}^{t-1} - \eta g_{r,i}^t$ 
9:     end for
10:    Devices send  $\Delta_{r,i} = \theta_{r,i}^T - \theta_{r,i}^0$  to server.
11:  end for
12:  Server computes  $\bar{\Delta}_r = \frac{1}{n} \sum_{i=1}^n \Delta_{r,i}$ 
13:   $m_r = \beta_1 m_{r-1} + (1 - \beta_1) \bar{\Delta}_r$ 
14:   $v_r = \beta_2 v_{r-1} + (1 - \beta_2) \bar{\Delta}_r^2$ 
15:   $\theta_r = \theta_{r-1} + \eta g \frac{m_r}{\sqrt{v_r} + \epsilon}$ 
16: end for
17: Output: Global model parameter  $\theta_R$ .

```

---

### A.2 Hyper-parameter Tuning

In our empirical study, we tune the learning rate of each algorithm carefully such that the best performance is achieved. The search grids in all our experiments are provided in Table 2.

Table 2: Search grids of the learning rate.

	Learning rate range
Fed-SGD	[0.001, 0.003, 0.005, 0.01, 0.03, 0.05, 0.1, 0.3, 0.5]
Fed-AMS	[0.0001, 0.0003, 0.0005, 0.001, 0.003, 0.005, 0.01, 0.03, 0.05, 0.1]
FED-LAMB	[0.001, 0.003, 0.005, 0.01, 0.03, 0.05, 0.1, 0.3, 0.5]
Adp-Fed	Local $\eta_l$ : [0.0001, 0.0003, 0.0005, 0.001, 0.003, 0.005, 0.01, 0.03, 0.05, 0.1, 0.3, 0.5]
	Global $\eta_g$ : [0.0001, 0.0003, 0.0005, 0.001, 0.003, 0.005, 0.01, 0.03, 0.05, 0.1]

## B Theoretical Analysis

We first recall in Table 3 some important notations that will be used in our following analysis.

---

$R, T$	$:=$	Number of communications rounds and local iterations (resp.)
$n, D, i$	$:=$	Total number of clients, portion sampled uniformly and client index
$\mathbf{h}, \ell$	$:=$	Total number of layers in the DNN and its index
$\phi(\cdot)$	$:=$	Scaling factor in FED-LAMBupdate
$\bar{\theta}$	$:=$	Global model (after periodic averaging)
$p_{r,i}^t$	$:=$	ratio computed at round $r$ , local iteration $t$ and for device $i$ . $p_{r,i}^{\ell,t}$ denotes its component at layer $\ell$

---

Table 3: Summary of notations used in the paper.

We now provide the proofs for the theoretical results of the main paper, including the intermediary Lemmas and the main convergence result, Theorem 1.

### B.1 Intermediary Lemmas

**Lemma.** Consider  $\{\bar{\theta}_r\}_{r>0}$ , the sequence of parameters obtained running Algorithm 1. Then for  $i \in \llbracket n \rrbracket$ :

$$\|\bar{\theta}_r - \theta_{r,i}\|^2 \leq \alpha^2 M^2 \phi_M^2 \frac{(1 - \beta_2)p}{v_0},$$

where  $\phi_M$  is defined in H4 and  $p$  is the total number of dimensions  $p = \sum_{\ell=1}^{\mathbf{h}} p_\ell$ .

*Proof.* Assuming the simplest case when  $T = 1$ , i.e., one local iteration, then by construction of Algorithm 1, we have for all  $\ell \in \llbracket \mathbf{h} \rrbracket$ ,  $i \in \llbracket n \rrbracket$  and  $r > 0$ :

$$\theta_{r,i}^\ell = \bar{\theta}_r^\ell - \alpha \phi(\|\theta_{r,i}^{\ell,t-1}\|) p_{r,i}^j / \|p_{r,i}^\ell\| = \bar{\theta}_r^\ell - \alpha \phi(\|\theta_{r,i}^{\ell,t-1}\|) \frac{m_{r,i}^t}{\sqrt{v_r^t}} \frac{1}{\|p_{r,i}^\ell\|}$$

leading to

$$\begin{aligned} \|\bar{\theta}_r - \theta_{r,i}\|^2 &= \sum_{\ell=1}^{\mathbf{h}} \left\langle \bar{\theta}_r^\ell - \theta_{r,i}^\ell \mid \bar{\theta}_r^\ell - \theta_{r,i}^\ell \right\rangle \\ &\leq \alpha^2 M^2 \phi_M^2 \frac{(1 - \beta_2)p}{v_0}, \end{aligned}$$

which concludes the proof.  $\square$

**Lemma.** Consider  $\{\bar{\theta}_r\}_{r>0}$ , the sequence of parameters obtained running Algorithm 1. Then for  $r > 0$ :

$$\left\| \frac{\nabla f(\bar{\theta}_r)}{\sqrt{v_r^t}} \right\|^2 \geq \frac{1}{2} \left\| \frac{\nabla f(\bar{\theta}_r)}{\sqrt{v_r^t}} \right\|^2 - \bar{L} \alpha^2 M^2 \phi_M^2 \frac{(1 - \beta_2)p}{v_0},$$

where  $M$  is defined in H2,  $p$  is the total number of dimensions  $p = \sum_{\ell=1}^{\mathbf{h}} p_\ell$  and  $\phi_M$  is defined in H4.



*Proof.* Consider the following sequence:

$$\left\| \frac{\bar{\nabla} f(\theta_r)}{\sqrt{v_r^t}} \right\|^2 \geq \frac{1}{2} \left\| \frac{\nabla f(\bar{\theta}_r)}{\sqrt{v_r^t}} \right\|^2 - \left\| \frac{\bar{\nabla} f(\theta_r) - \nabla f(\bar{\theta}_r)}{\sqrt{v_r^t}} \right\|^2,$$

where the inequality is due to the Cauchy-Schwartz inequality.

Under the smoothness assumption H1 and using Lemma 1, we have

$$\begin{aligned} \left\| \frac{\bar{\nabla} f(\theta_r)}{\sqrt{v_r^t}} \right\|^2 &\geq \frac{1}{2} \left\| \frac{\nabla f(\bar{\theta}_r)}{\sqrt{v_r^t}} \right\|^2 - \left\| \frac{\bar{\nabla} f(\theta_r) - \nabla f(\bar{\theta}_r)}{\sqrt{v_r^t}} \right\|^2 \\ &\geq \frac{1}{2} \left\| \frac{\nabla f(\bar{\theta}_r)}{\sqrt{v_r^t}} \right\|^2 - \bar{L} \alpha^2 M^2 \phi_M^2 \frac{(1 - \beta_2)p}{v_0}, \end{aligned}$$

which concludes the proof.  $\square$

## B.2 Proof of Theorem 1

We now develop a proof for the two intermediary lemmas, Lemma 1 and Lemma 2, in the case when each local model is obtained after more than one local update. Then the two quantities, either the gap between the periodically averaged parameter and each local update, i.e.,  $\|\bar{\theta}_r - \theta_{r,i}\|^2$ , and the ratio of the average gradient, more particularly its relation to the gradient of the average global model (i.e.,  $\left\| \frac{\bar{\nabla} f(\theta_r)}{\sqrt{v_r^t}} \right\|$  and  $\left\| \frac{\nabla f(\bar{\theta}_r)}{\sqrt{v_r^t}} \right\|$ ), are impacted.

**Theorem.** Assume H1-H4. Consider  $\{\bar{\theta}_r\}_{r>0}$ , the sequence of parameters obtained running Algorithm 1 with a decreasing learning rate  $\alpha$ . Let the number of local epochs be  $T \geq 1$  and  $\lambda = 0$ . Then, at iteration  $\tau$ , we have

$$\begin{aligned} \frac{1}{\tau} \sum_{t=1}^{\tau} \mathbb{E} \left[ \left\| \frac{\nabla f(\bar{\theta}_t)}{\hat{v}_t^{1/4}} \right\|^2 \right] &\leq \sqrt{\frac{M^2 p}{n}} \frac{\mathbb{E}[f(\bar{\theta}_1)] - \min_{\theta \in \Theta} f(\theta)}{\mathfrak{h} \alpha_r \tau} + \frac{\phi_M \sigma^2}{\tau n} \sqrt{\frac{1 - \beta_2}{M^2 p}} \\ &+ 4\alpha \left[ \frac{\alpha^2 L}{\sqrt{v_0}} M^2 (T - 1)^2 \phi_M^2 (1 - \beta_2) p + \frac{M^2}{\sqrt{v_0}} + \phi_M^2 \sqrt{M^2 + p \sigma^2} + \phi_M \frac{\mathfrak{h} \sigma^2}{\sqrt{n}} \right] + cst. \end{aligned}$$

*Proof.* Using H1, we have

$$\begin{aligned} f(\bar{\vartheta}_{r+1}) &\leq f(\bar{\vartheta}_r) + \langle \nabla f(\bar{\vartheta}_r) | \bar{\vartheta}_{r+1} - \bar{\vartheta}_r \rangle + \sum_{\ell=1}^L \frac{L_\ell}{2} \|\bar{\vartheta}_{r+1}^\ell - \bar{\vartheta}_r^\ell\|^2 \\ &\leq f(\bar{\vartheta}_r) + \sum_{\ell=1}^{\mathfrak{h}} \sum_{j=1}^{p_\ell} \nabla_\ell f(\bar{\vartheta}_r)^j (\bar{\vartheta}_{r+1}^{\ell,j} - \bar{\vartheta}_r^{\ell,j}) + \sum_{\ell=1}^L \frac{L_\ell}{2} \|\bar{\vartheta}_{r+1}^\ell - \bar{\vartheta}_r^\ell\|^2. \end{aligned}$$

Taking expectations on both sides leads to

$$-\mathbb{E}[\langle \nabla f(\bar{\vartheta}_r) | \bar{\vartheta}_{r+1} - \bar{\vartheta}_r \rangle] \leq \mathbb{E}[f(\bar{\vartheta}_r) - f(\bar{\vartheta}_{r+1})] + \sum_{\ell=1}^L \frac{L_\ell}{2} \mathbb{E}[\|\bar{\vartheta}_{r+1}^\ell - \bar{\vartheta}_r^\ell\|^2]. \quad (6)$$

Yet, we observe that, using the classical intermediate quantity, used for proving convergence results of adaptive optimization methods, see for instance [Reddi et al. \(2018\)](#), we have

$$\bar{\vartheta}_r = \bar{\theta}_r + \frac{\beta_1}{1 - \beta_1}(\bar{\theta}_r - \bar{\theta}_{r-1}), \quad (7)$$

where  $\bar{\theta}_r$  denotes the average of the local models at round  $r$ . Then for each layer  $\ell$ ,

$$\begin{aligned} \bar{\vartheta}_{r+1}^\ell - \bar{\vartheta}_r^\ell &= \frac{1}{1 - \beta_1}(\bar{\theta}_{r+1}^\ell - \bar{\theta}_r^\ell) - \frac{\beta_1}{1 - \beta_1}(\bar{\theta}_r^\ell - \bar{\theta}_{r-1}^\ell) \\ &= \frac{\alpha_r}{1 - \beta_1} \frac{1}{n} \sum_{i=1}^n \frac{\phi(\|\theta_{r,i}^\ell\|)}{\|p_{r,i}^\ell\|} p_{r,i}^\ell - \frac{\alpha_{r-1}}{1 - \beta_1} \frac{1}{n} \sum_{i=1}^n \frac{\phi(\|\theta_{r-1,i}^\ell\|)}{\|p_{r-1,i}^\ell\|} p_{r-1,i}^\ell \\ &= \frac{\alpha\beta_1}{1 - \beta_1} \frac{1}{n} \sum_{i=1}^n \left( \frac{\phi(\|\theta_{r,i}^\ell\|)}{\sqrt{v_r^t} \|p_{r,i}^\ell\|} - \frac{\phi(\|\theta_{r-1,i}^\ell\|)}{\sqrt{v_{r-1}^t} \|p_{r-1,i}^\ell\|} \right) m_{r-1}^t + \frac{\alpha}{n} \sum_{i=1}^n \frac{\phi(\|\theta_{r,i}^\ell\|)}{\sqrt{v_r^t} \|p_{r,i}^\ell\|} g_{r,i}^t, \end{aligned} \quad (8)$$

where we have assumed a constant learning rate  $\alpha$ .

We note for all  $\theta \in \Theta$ , the majorant  $G > 0$  such that  $\phi(\|\theta\|) \leq G$ . Then, following (6), we obtain

$$-\mathbb{E}[\langle \nabla f(\bar{\vartheta}_r) | \bar{\vartheta}_{r+1} - \bar{\vartheta}_r \rangle] \leq \mathbb{E}[f(\bar{\vartheta}_r) - f(\bar{\vartheta}_{r+1})] + \sum_{\ell=1}^L \frac{L_\ell}{2} \mathbb{E}[\|\bar{\vartheta}_{r+1} - \bar{\vartheta}_r\|^2]. \quad (9)$$

Developing the LHS of (9) using (8) leads to

$$\begin{aligned} \langle \nabla f(\bar{\vartheta}_r) | \bar{\vartheta}_{r+1} - \bar{\vartheta}_r \rangle &= \sum_{\ell=1}^h \sum_{j=1}^{p_\ell} \nabla_\ell f(\bar{\vartheta}_r)^j (\bar{\vartheta}_{r+1}^{\ell,j} - \bar{\vartheta}_r^{\ell,j}) \\ &= \frac{\alpha\beta_1}{1 - \beta_1} \frac{1}{n} \sum_{\ell=1}^h \sum_{j=1}^{p_\ell} \nabla_\ell f(\bar{\vartheta}_r)^j \left[ \sum_{i=1}^n \left( \frac{\phi(\|\theta_{r,i}^\ell\|)}{\sqrt{v_r^t} \|p_{r,i}^\ell\|} - \frac{\phi(\|\theta_{r-1,i}^\ell\|)}{\sqrt{v_{r-1}^t} \|p_{r-1,i}^\ell\|} \right) m_{r-1}^t \right] \\ &\quad - \underbrace{\frac{\alpha}{n} \sum_{\ell=1}^h \sum_{j=1}^{p_\ell} \nabla_\ell f(\bar{\vartheta}_r)^j \sum_{i=1}^n \frac{\phi(\|\theta_{r,i}^\ell\|)}{\sqrt{v_r^t} \|p_{r,i}^\ell\|} g_{r,i}^t}_{=A_1}. \end{aligned} \quad (10)$$

We change all index  $r$  to iteration  $t$ . Suppose  $T$  is the number of local iterations. We can write (10) as

$$A_1 = -\alpha_t \langle \nabla f(\bar{\vartheta}_t), \frac{\bar{g}_t}{\sqrt{\hat{v}_t}} \rangle,$$

where  $\bar{g}_t = \frac{1}{n} \sum_{i=1}^n \bar{g}_{t,i}$ , with  $\bar{g}_{t,i} = \left[ \frac{\phi(\|\theta_{t,i}^1\|)}{\|p_{t,i}^1\|} g_{t,i}^1, \dots, \frac{\phi(\|\theta_{t,i}^L\|)}{\|p_{t,i}^L\|} g_{t,i}^L \right]$  representing the normalized gradient (concatenated by layers) of the  $i$ -th device. It holds that

$$\langle \nabla f(\bar{\vartheta}_t), \frac{\bar{g}_t}{\sqrt{\hat{v}_t}} \rangle = \frac{1}{2} \left\| \frac{\nabla f(\bar{\vartheta}_t)}{\hat{v}_t^{1/4}} \right\|^2 + \frac{1}{2} \left\| \frac{\bar{g}_t}{\hat{v}_t^{1/4}} \right\|^2 - \left\| \frac{\nabla f(\bar{\vartheta}_t) - \bar{g}_t}{\hat{v}_t^{1/4}} \right\|^2. \quad (11)$$

To bound the last term on the RHS, we have

$$\begin{aligned}
\left\| \frac{\nabla f(\bar{\vartheta}_t) - \bar{g}_t}{\hat{v}_t^{1/4}} \right\|^2 &= \left\| \frac{\frac{1}{n} \sum_{i=1}^n (\nabla f(\bar{\vartheta}_t) - \bar{g}_{t,i})}{\hat{v}_t^{1/4}} \right\|^2 \\
&\leq \frac{1}{n} \sum_{i=1}^n \left\| \frac{\nabla f(\bar{\vartheta}_t) - \bar{g}_{t,i}}{\hat{v}_t^{1/4}} \right\|^2 \\
&\leq \frac{2}{n} \sum_{i=1}^n \left( \left\| \frac{\nabla f(\bar{\vartheta}_t) - \nabla f(\bar{\theta}_t)}{\hat{v}_t^{1/4}} \right\|^2 + \left\| \frac{\nabla f(\bar{\theta}_t) - \bar{g}_{t,i}}{\hat{v}_t^{1/4}} \right\|^2 \right).
\end{aligned}$$

By Lipschitz smoothness of the loss function, the first term admits

$$\begin{aligned}
\frac{2}{n} \sum_{i=1}^n \left\| \frac{\nabla f_i(\bar{\vartheta}_t) - \nabla f_i(\bar{\theta}_t)}{\hat{v}_t^{1/4}} \right\|^2 &\leq \frac{2}{n\sqrt{v_0}} \sum_{i=1}^n L_\ell \|\bar{\vartheta}_t - \bar{\theta}_t\|^2 \\
&= \frac{2L_\ell}{n\sqrt{v_0}} \frac{\beta_1^2}{(1-\beta_1)^2} \sum_{i=1}^n \|\bar{\theta}_t - \bar{\theta}_{t-1}\|^2 \\
&\leq \frac{2\alpha_r^2 L_\ell}{n\sqrt{v_0}} \frac{\beta_1^2}{(1-\beta_1)^2} \sum_{l=1}^L \sum_{i=1}^n \left\| \frac{\phi(\|\theta_{t,i}^l\|)}{\|p_{t,i}^l\|} p_{t,i}^l \right\|^2 \\
&\leq \frac{2\alpha_r^2 L_\ell p \phi_M^2}{\sqrt{v_0}} \frac{\beta_1^2}{(1-\beta_1)^2}.
\end{aligned}$$

For the second term,

$$\frac{2}{n} \sum_{i=1}^n \left\| \frac{\nabla f(\bar{\theta}_t) - \bar{g}_{t,i}}{\hat{v}_t^{1/4}} \right\|^2 \leq \frac{4}{n} \underbrace{\left( \sum_{i=1}^n \left\| \frac{\nabla f(\bar{\theta}_t) - \nabla f(\theta_{t,i})}{\hat{v}_t^{1/4}} \right\|^2 \right)}_{B_1} + \underbrace{\sum_{i=1}^n \left\| \frac{\nabla f(\theta_{t,i}) - \bar{g}_{t,i}}{\hat{v}_t^{1/4}} \right\|^2}_{B_2}. \quad (12)$$

Using the smoothness of  $f_i$  we can transform  $B_1$  into consensus error by

$$\begin{aligned}
B_1 &\leq \frac{L}{\sqrt{v_0}} \sum_{i=1}^n \|\bar{\theta}_t - \theta_{t,i}\|^2 \\
&= \frac{\alpha_r^2 L}{\sqrt{v_0}} \sum_{i=1}^n \sum_{l=1}^L \left\| \sum_{j=\lfloor t \rfloor_T + 1}^t \left( \frac{\phi(\|\theta_{j,i}^l\|)}{\|p_{j,i}^l\|} p_{j,i}^l - \frac{1}{n} \sum_{k=1}^n \frac{\phi(\|\theta_{j,k}^l\|)}{\|p_{j,k}^l\|} p_{j,k}^l \right) \right\|^2 \\
&\leq n \frac{\alpha_r^2 L}{\sqrt{v_0}} M^2 (T-1)^2 \phi_M^2 (1-\beta_2) p,
\end{aligned} \quad (13)$$

where the last inequality stems from Lemma 1 in the particular case where  $\theta_{t,i}$  are averaged every  $ct + 1$  local iterations for any integer  $c$ , since  $(t-1) - (\lfloor t \rfloor_T + 1) + 1 \leq T-1$ .

We now develop the expectation of  $B_2$  under the simplification that  $\beta_1 = 0$ :

$$\begin{aligned}
\mathbb{E}[B_2] &= \mathbb{E} \left[ \sum_{i=1}^n \left\| \frac{\nabla f(\theta_{t,i}) - \bar{g}_{t,i}}{\hat{v}_t^{1/4}} \right\|^2 \right] \\
&\leq \frac{nM^2}{\sqrt{v_0}} + n\phi_M^2 \sqrt{M^2 + p\sigma^2} - 2 \sum_{i=1}^n \mathbb{E}[\langle \nabla f(\theta_{t,i}), \bar{g}_{t,i} \rangle / \sqrt{\hat{v}_t}]
\end{aligned}$$

$$\begin{aligned}
&= \frac{nM^2}{\sqrt{v_0}} + n\phi_M^2\sqrt{M^2 + p\sigma^2} - 2 \sum_{i=1}^n \sum_{\ell=1}^L \mathbb{E}[\langle \nabla_\ell f(\theta_{t,i}), \frac{\phi(\|\theta_{t,i}^l\|)}{\|p_{t,i}^l\|} g_{t,i}^l \rangle / \sqrt{\hat{v}_t^l}] \\
&= \frac{nM^2}{\sqrt{v_0}} + n\phi_M^2\sqrt{M^2 + p\sigma^2} - 2 \sum_{i=1}^n \sum_{l=1}^L \sum_{i=1}^{p_l} \mathbb{E}[\nabla_l f(\theta_{t,i})^j \frac{\phi(\|\theta_{t,i}^{l,j}\|)}{\sqrt{\hat{v}_t^{l,j}} \|p_{t,i}^{l,j}\|} g_{t,i}^{l,j}] \\
&\leq \frac{nM^2}{\sqrt{v_0}} + n\phi_M^2\sqrt{M^2 + p\sigma^2} - 2 \sum_{i=1}^n \sum_{l=1}^L \sum_{i=1}^{p_l} \mathbb{E} \left[ \sqrt{\frac{1-\beta_2}{M^2 p_\ell}} \phi(\|\theta_{r,i}^{l,j}\|) \nabla_l f(\theta_{t,i})^j g_{t,i}^{l,j} \right] \\
&\quad - 2 \sum_{i=1}^n \sum_{l=1}^L \sum_{j=1}^{p_l} E \left[ \left( \phi(\|\theta_{r,i}^{l,j}\|) \nabla_l f(\theta_{t,i})^j \frac{g_{r,i}^{t,l,j}}{\|p_{r,i}^{l,j}\|} \right) 1 \left( \text{sign}(\nabla_l f(\theta_{t,i})^j) \neq \text{sign}(g_{r,i}^{t,l,j}) \right) \right],
\end{aligned}$$

where we use assumption H2, H3 and H4. Yet,

$$\begin{aligned}
&- \mathbb{E} \left[ \left( \phi(\|\theta_{r,i}^{l,j}\|) \nabla_l f(\theta_{t,i})^j \frac{g_{r,i}^{t,l,j}}{\|p_{r,i}^{l,j}\|} \right) 1 \left( \text{sign}(\nabla_l f(\theta_{t,i})^j) \neq \text{sign}(g_{r,i}^{t,l,j}) \right) \right] \\
&\leq \phi_M \nabla_l f(\theta_{t,i})^j \mathbb{P} \left[ \text{sign}(\nabla_l f(\theta_{t,i})^j) \neq \text{sign}(g_{r,i}^{t,l,j}) \right].
\end{aligned}$$

Then we have

$$\mathbb{E}[B_2] \leq \frac{nM^2}{\sqrt{v_0}} + n\phi_M^2\sqrt{M^2 + p\sigma^2} - 2\phi_m \sqrt{\frac{1-\beta_2}{M^2 p}} \sum_{i=1}^n \mathbb{E}[\|\nabla f(\theta_{t,i})\|^2] + \phi_M \frac{\mathbf{h}\sigma^2}{\sqrt{n}}$$

Thus, (12) becomes

$$\frac{2}{n} \sum_{i=1}^n \left\| \frac{\nabla f_i(\bar{\theta}_t) - \bar{g}_t}{\hat{v}_t^{1/4}} \right\|^2 \leq 4 \left[ \frac{\alpha_t^2 L l}{\sqrt{v_0}} \alpha_r^2 M^2 (T-1)^2 \phi_M^2 (1-\beta_2) p + \frac{M^2}{\sqrt{v_0}} + \phi_M^2 \sqrt{M^2 + p\sigma^2} + \phi_M \frac{\mathbf{h}\sigma^2}{\sqrt{n}} \right]$$

Substituting all ingredients into (11), we obtain

$$\begin{aligned}
-\alpha_t \mathbb{E}[\langle \nabla f(\bar{\vartheta}_t), \frac{\bar{g}_t}{\sqrt{\hat{v}_t}} \rangle] &\leq -\frac{\alpha_t}{2} \mathbb{E}[\|\frac{\nabla f(\bar{\vartheta}_t)}{\hat{v}_t^{1/4}}\|^2] - \frac{\alpha_t}{2} \mathbb{E}[\|\frac{\bar{g}_t}{\hat{v}_t^{1/4}}\|^2] + \frac{2\alpha_t^3 L_\ell p \phi_M^2}{\sqrt{v_0}} \frac{\beta_1^2}{(1-\beta_1)^2} \\
&\quad + 4 \left[ \frac{\alpha_t^2 L}{\sqrt{v_0}} M^2 (T-1)^2 \phi_M^2 (1-\beta_2) p + \frac{M^2}{\sqrt{v_0}} + \phi_M^2 \sqrt{M^2 + p\sigma^2} + \phi_M \frac{\mathbf{h}\sigma^2}{\sqrt{n}} \right].
\end{aligned}$$

At the same time, we have

$$\begin{aligned}
\mathbb{E}[\|\frac{\bar{g}_t}{\hat{v}_t^{1/4}}\|^2] &= \frac{1}{n^2} \mathbb{E}[\|\sum_{i=1}^n \frac{\bar{g}_{t,i}}{\hat{v}_t^{1/4}}\|^2] \\
&= \frac{1}{n^2} \mathbb{E}[\sum_{l=1}^L \sum_{i=1}^n \|\frac{\phi(\|\theta_{t,i}^l\|)}{\hat{v}_t^{1/4} \|p_{t,i}^l\|} g_{t,i}^l\|^2] \\
&\geq \phi_m^2 (1-\beta_2) \mathbb{E} \left[ \left\| \frac{1}{n} \sum_{i=1}^n \frac{\nabla f(\theta_{t,i})}{\hat{v}_t^{1/4}} \right\|^2 \right] \\
&= \phi_m^2 (1-\beta_2) \mathbb{E} \left[ \left\| \frac{\nabla f(\theta_t)}{\hat{v}_t^{1/4}} \right\|^2 \right].
\end{aligned}$$

Regarding  $\left\| \frac{\nabla f(\theta_t)}{\hat{v}_t^{1/4}} \right\|^2$ , we have

$$\begin{aligned} \left\| \frac{\nabla f(\theta_t)}{\hat{v}_t^{1/4}} \right\|^2 &\geq \frac{1}{2} \left\| \frac{\nabla f(\bar{\theta}_t)}{\hat{v}_t^{1/4}} \right\|^2 - \left\| \frac{\nabla f(\theta_t) - \nabla f(\bar{\theta}_t)}{\hat{v}_t^{1/4}} \right\|^2 \\ &\geq \frac{1}{2} \left\| \frac{\nabla f(\bar{\theta}_t)}{\hat{v}_t^{1/4}} \right\|^2 - \left\| \frac{\frac{1}{n} \sum_{i=1}^n (\nabla f(\theta_{t,i}) - \nabla f(\bar{\theta}_i))}{\hat{v}_t^{1/4}} \right\|^2 \\ &\geq \frac{1}{2} \left\| \frac{\nabla f(\bar{\theta}_t)}{\hat{v}_t^{1/4}} \right\|^2 - \frac{\alpha_t^2 L_\ell}{\sqrt{v_0}} M^2 (T-1)^2 \phi_M^2 (1-\beta_2) p, \end{aligned}$$

where the last line is due to (13). Therefore, we have obtained

$$\begin{aligned} A_1 &\leq -\frac{\phi_m^2 (1-\beta_2)}{2} \left\| \frac{\nabla f(\bar{\theta}_t)}{\hat{v}_t^{1/4}} \right\|^2 + \frac{\alpha_t^2 L_\ell}{\sqrt{v_0}} M^2 (T-1)^2 \phi_m^2 \phi_M^2 (1-\beta_2)^2 p + \frac{2\alpha^3 L_\ell p \phi_M^2}{\sqrt{v_0}} \frac{\beta_1^2}{(1-\beta_1)^2} \\ &\quad + 4\alpha_t \left[ \frac{\alpha_t^2 L}{\sqrt{v_0}} M^2 (T-1)^2 \phi_M^2 (1-\beta_2) p + \frac{M^2}{\sqrt{v_0}} + \phi_M^2 \sqrt{M^2 + p\sigma^2} + \phi_M \frac{\mathbf{h}\sigma^2}{\sqrt{n}} \right]. \end{aligned}$$

Substitute back into (10), and leave other derivations unchanged. Assuming  $M \leq 1$ , we have the following

$$\begin{aligned} &\frac{1}{\tau} \sum_{t=1}^{\tau} \mathbb{E} \left[ \left\| \frac{\nabla f(\bar{\theta}_t)}{\hat{v}_t^{1/4}} \right\|^2 \right] \\ &\lesssim \sqrt{\frac{M^2 p}{n}} \frac{f(\bar{\vartheta}_1) - \mathbb{E}[f(\bar{\vartheta}_{\tau+1})]}{\mathbf{h}\alpha_t \tau} + \frac{\alpha_t}{n^2} \sum_{r=1}^{\tau} \sum_{i=1}^n \sigma_i^2 \mathbb{E} \left[ \left\| \frac{\phi(\|\theta_{r,i}^\ell\|)}{\sqrt{v_t} \|p_{r,i}^\ell\|} \right\|^2 \right] + \frac{2\alpha^3 L_\ell p \phi_M^2}{\sqrt{v_0}} \frac{\beta_1^2}{(1-\beta_1)^2} \\ &\quad + 4\alpha_t \left[ \frac{\alpha_t^2 L_\ell}{\sqrt{v_0}} M^2 (T-1)^2 \phi_M^2 (1-\beta_2) p + \frac{M^2}{\sqrt{v_0}} + \phi_M^2 \sqrt{M^2 + p\sigma^2} + \phi_M \frac{\mathbf{h}\sigma^2}{\sqrt{n}} \right] + \frac{\bar{L}\beta_1^2 \mathbf{h}(1-\beta_2) M^2 \phi_M^2 n}{2(1-\beta_1)^2 v_0} \\ &\quad + \frac{\alpha_t \beta_1}{1-\beta_1} \sqrt{(1-\beta_2)p} \frac{\mathbf{h}M^2}{\sqrt{v_0}} + \bar{L}\alpha_t^2 M^2 \phi_M^2 \frac{(1-\beta_2)p}{Tv_0} \\ &\leq \sqrt{\frac{M^2 p}{n}} \frac{\mathbb{E}[f(\bar{\theta}_1)] - \min_{\theta \in \Theta} f(\theta)}{\mathbf{h}\alpha_t \tau} + \frac{\phi_M \sigma^2}{\tau n} \sqrt{\frac{1-\beta_2}{M^2 p}} \\ &\quad + 4\alpha_t \left[ \frac{\alpha_t^2 L_\ell}{\sqrt{v_0}} M^2 (T-1)^2 \phi_M^2 (1-\beta_2) p + \frac{M^2}{\sqrt{v_0}} + \phi_M^2 \sqrt{M^2 + p\sigma^2} + \phi_M \frac{\mathbf{h}\sigma^2}{\sqrt{n}} \right] \\ &\quad + \frac{\alpha_t \beta_1}{1-\beta_1} \sqrt{(1-\beta_2)p} \frac{\mathbf{h}M^2}{\sqrt{v_0}} + \bar{L}\alpha_t^2 M^2 \phi_M^2 \frac{(1-\beta_2)p}{Tv_0} + \frac{\bar{L}\beta_1^2 \mathbf{h}(1-\beta_2) M^2 \phi_M^2 n}{2(1-\beta_1)^2 v_0} + \frac{2\alpha^3 L_\ell p \phi_M^2}{\sqrt{v_0}} \frac{\beta_1^2}{(1-\beta_1)^2}. \end{aligned}$$

And if we set the learning rate to be of order  $\mathcal{O}(\frac{1}{L\sqrt{\tau}})$  then:

$$\frac{1}{\tau} \sum_{t=1}^{\tau} \mathbb{E} \left[ \left\| \frac{\nabla f(\bar{\theta}_t)}{\hat{v}_t^{1/4}} \right\|^2 \right] \leq \mathcal{O} \left( \sqrt{\frac{M^2 p}{n}} \frac{1}{\sqrt{\mathbf{h}\tau}} + \frac{\sigma^2}{\tau n \sqrt{p}} + \frac{(T-1)^2 p}{\tau^3 L^3} \right).$$

This concludes the proof.  $\square$

## References

- Dan Alistarh, Demjan Grubic, Jerry Li, Ryota Tomioka, and Milan Vojnovic. QSGD: communication-efficient SGD via gradient quantization and encoding. In *Advances in Neural Information Processing Systems (NIPS)*, pages 1709–1720, Long Beach, CA, 2017.
- Xiangyi Chen, Xiaoyun Li, and Ping Li. Toward communication efficient adaptive gradient method. In *Proceedings of the ACM-IMS Foundations of Data Science Conference (FODS)*, pages 119–128, Virtual Event, USA, 2020.
- Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. Imagenet: A large-scale hierarchical image database. In *Proceedings of the 2009 IEEE Computer Society Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 248–255, Miami, FL, 2009.
- Timothy Dozat. Incorporating nesterov momentum into Adam. In *Proceedings of the 4th International Conference on Learning Representations (ICLR Workshop)*, San Juan, Puerto Rico, 2016.
- John C. Duchi, Elad Hazan, and Yoram Singer. Adaptive subgradient methods for online learning and stochastic optimization. *J. Mach. Learn. Res.*, 12:2121–2159, 2011.
- Saeed Ghadimi and Guanghui Lan. Stochastic first- and zeroth-order methods for nonconvex stochastic programming. *SIAM J. Optim.*, 23(4):2341–2368, 2013.
- Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *Proceedings of the 2016 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 770–778, Las Vegas, NV, 2016.
- Sai Praneeth Karimireddy, Satyen Kale, Mehryar Mohri, Sashank J Reddi, Sebastian U Stich, and Ananda Theertha Suresh. Scaffold: Stochastic controlled averaging for on-device federated learning. *arXiv preprint arXiv:1910.06378*, 2019.
- Diederik P. Kingma and Jimmy Ba. Adam: A method for stochastic optimization. In *Proceedings of the 3rd International Conference on Learning Representations (ICLR)*, San Diego, CA, 2015.
- Jakub Konečný, H Brendan McMahan, Felix X Yu, Peter Richtárik, Ananda Theertha Suresh, and Dave Bacon. Federated learning: Strategies for improving communication efficiency. *arXiv preprint arXiv:1610.05492*, 2016.
- Alex Krizhevsky. Learning multiple layers of features from tiny images. *Master’s thesis, Department of Computer Science, University of Toronto*, 2009.
- Yann LeCun. The mnist database of handwritten digits. <http://yann.lecun.com/exdb/mnist/>, 1998.
- Mu Li, David G. Andersen, Jun Woo Park, Alexander J. Smola, Amr Ahmed, Vanja Josifovski, James Long, Eugene J. Shekita, and Bor-Yiing Su. Scaling distributed machine learning with the parameter server. In *Proceedings of the 11th USENIX Symposium on Operating Systems Design and Implementation (OSDI)*, pages 583–598, Broomfield, CO, 2014.
- Tian Li, Anit Kumar Sahu, Ameet Talwalkar, and Virginia Smith. Federated learning: Challenges, methods, and future directions. *IEEE Signal Process. Mag.*, 37(3):50–60, 2020.

- Xianfeng Liang, Shuheng Shen, Jingchang Liu, Zhen Pan, Enhong Chen, and Yifei Cheng. Variance reduced local sgd with lower communication complexity. *arXiv preprint arXiv:1912.12844*, 2019.
- Yujun Lin, Song Han, Huizi Mao, Yu Wang, and Bill Dally. Deep gradient compression: Reducing the communication bandwidth for distributed training. In *Proceedings of the 6th International Conference on Learning Representations (ICLR)*, Vancouver, Canada, 2018.
- Brendan McMahan and Matthew J. Streeter. Adaptive bound optimization for online convex optimization. In *Proceedings of the 23rd Conference on Learning Theory (COLT)*, pages 244–256, Haifa, Israel, 2010.
- Brendan McMahan, Eider Moore, Daniel Ramage, Seth Hampson, and Blaise Agüera y Arcas. Communication-efficient learning of deep networks from decentralized data. In *Proceedings of the 20th International Conference on Artificial Intelligence and Statistics (AISTATS)*, pages 1273–1282, Fort Lauderdale, FL, 2017.
- Yurii Nesterov. Introductory lectures on convex optimization: A basic course. *Springer*, 2004.
- B. T. Polyak. Some methods of speeding up the convergence of iteration methods. *Mathematics and Mathematical Physics*, 1964.
- Benjamin Recht, Christopher Ré, Stephen J. Wright, and Feng Niu. Hogwild: A lock-free approach to parallelizing stochastic gradient descent. In *Advances in Neural Information Processing Systems (NIPS)*, pages 693–701, Granada, Spain, 2011.
- Sashank J. Reddi, Satyen Kale, and Sanjiv Kumar. On the convergence of adam and beyond. In *Proceedings of the 6th International Conference on Learning Representations (ICLR)*, Vancouver, Canada, 2018.
- Sashank J. Reddi, Zachary Charles, Manzil Zaheer, Zachary Garrett, Keith Rush, Jakub Konečný, Sanjiv Kumar, and Hugh Brendan McMahan. Adaptive federated optimization. In *Proceedings of the 9th International Conference on Learning Representations (ICLR)*, Virtual Event, Austria, 2021.
- Daniel Rothchild, Ashwinee Panda, Enayat Ullah, Nikita Ivkin, Ion Stoica, Vladimir Braverman, Joseph Gonzalez, and Raman Arora. Fetchsgd: Communication-efficient federated learning with sketching. In *Proceedings of the 37th International Conference on Machine Learning (ICML)*, pages 8253–8265, Virtual Event, 2020.
- Sebastian U. Stich. Local SGD converges fast and communicates little. In *Proceedings of the 7th International Conference on Learning Representations (ICLR)*, New Orleans, LA, 2019.
- T. Tieleman and G. Hinton. Rmsprop: Divide the gradient by a running average of its recent magnitude. *COURSERA: Neural Networks for Machine Learning*, 2012.
- Jianqiao Wangni, Jialei Wang, Ji Liu, and Tong Zhang. Gradient sparsification for communication-efficient distributed optimization. In *Advances in Neural Information Processing Systems (NeurIPS)*, pages 1306–1316, Montréal, Canada, 2018.
- Han Xiao, Kashif Rasul, and Roland Vollgraf. Fashion-MNIST: a novel image dataset for benchmarking machine learning algorithms. *arXiv preprint arXiv:1708.07747*, 2017.



- Yang You, Jing Li, Sashank J. Reddi, Jonathan Hseu, Sanjiv Kumar, Srinadh Bhojanapalli, Xiaodan Song, James Demmel, Kurt Keutzer, and Cho-Jui Hsieh. Large batch optimization for deep learning: Training BERT in 76 minutes. In *Proceedings of the 8th International Conference on Learning Representations (ICLR)*, Addis Ababa, Ethiopia, 2020.
- Hao Yu, Rong Jin, and Sen Yang. On the linear speedup analysis of communication efficient momentum SGD for distributed non-convex optimization. In *Proceedings of the 36th International Conference on Machine Learning (ICML)*, pages 7184–7193, Long Beach, CA, 2019.
- Matthew D Zeiler. Adadelta: an adaptive learning rate method. *arXiv preprint arXiv:1212.5701*, 2012.
- Weijie Zhao, Deping Xie, Ronglai Jia, Yulei Qian, Ruiquan Ding, Mingming Sun, and Ping Li. Distributed hierarchical GPU parameter server for massive scale deep learning ads systems. In *Proceedings of Machine Learning and Systems (MLSys)*, Austin, TX, 2020.
- Dongruo Zhou, Jinghui Chen, Yuan Cao, Yiqi Tang, Ziyang Yang, and Quanquan Gu. On the convergence of adaptive gradient methods for nonconvex optimization. *arXiv preprint arXiv:1808.05671*, 2018.
- Fan Zhou and Guojing Cong. On the convergence properties of a k-step averaging stochastic gradient descent algorithm for nonconvex optimization. In *Proceedings of the Twenty-Seventh International Joint Conference on Artificial Intelligence (IJCAI)*, pages 3219–3227, Stockholm, Sweden, 2018.
- Yingxue Zhou, Belhal Karimi, Jinxing Yu, Zhiqiang Xu, and Ping Li. Towards better generalization of adaptive gradient methods. In *Advances in Neural Information Processing Systems (NeurIPS)*, virtual, 2020.