

Layerwise and Dimensionwise Local Adaptive Method for Federated Learning

Anonymous Authors¹

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. Under the orchestration of a central server, the data needs to remain decentralized, as it can not be shared among clients or with the central server. Then, 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 the particular case of training Deep, as in multilayered Neural Networks, under such settings, we propose in this paper, FED-LAMB, a novel Federated Learning method based on a Layerwise and Dimensionwise updates of the local models. A periodic averaging is added to obtain estimates of the desired global model parameters. We provide a thorough finite time convergence analysis for our algorithm, substantiated by numerical runs on benchmark datasets.

1. Introduction

A growing and important task while learning models on observed data, is the ability to train the latter over a large number of clients which could either be devices or distinct entities. In the paradigm of Federated Learning (FL) (Konečný et al., 2016; McMahan et al., 2017), the focus of our paper, a central server orchestrates the optimization over those clients under the constraint that the data can neither be centralized nor shared among the clients. Most modern machine learning tasks can be casted as a 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 θ the global model parameter taking value in Θ , a subset of \mathbb{R}^d . While this formulation recalls that of distributed optimization, the core principle of FL is different than standard distributed paradigm.

FL currently suffers from two bottlenecks: communication efficiency and privacy. We focus on the former in this paper. While local updates, updates during which each client learn their local models, can reduce drastically the number of communication rounds between the central server and devices, new techniques are still necessary to tackle the challenge of communication due to, e.g., wireless bandwidth. Some quantization (Alistarh et al., 2017; Wangni et al., 2018) or compression (Lin et al., 2017) methods allow to decrease the number of bits communicated at each round and are efficient methods in a distributed setting. The other approach one can take is to accelerate the local training on each device and thus sending a better local model to the server at each round, thus reducing the number of communication rounds needed to get a well-trained global model.

Under the important setting of heterogenous data, i.e. the data among each device can be distributed according to different distributions, current local optimization algorithms are perfectible. One of the most popular frameworks for FL is using multiple local Stochastic Gradient Descent (SGD) steps in each device, sending those local models to the server that computes the average over those received local model parameters and broadcasts it back to the devices. This method is called FEDAVG (McMahan et al., 2017).

In Chen et al. (2020), the authors motivate the usage of adaptive gradient optimization methods as a better alternative to the standard SGD inner loop in FEDAVG. They propose an adaptive gradient method, namely LOCAL AMSGRAD, with communication cost sublinear in T that is guaranteed to converge to stationary points in $\mathcal{O}(\sqrt{d/Tn})$, where T is the number of communication rounds, d is the overall dimension of the problem and n corresponds to the number of clients available.

¹Anonymous Institution, Anonymous City, Anonymous Region, Anonymous Country. Correspondence to: Anonymous Author <anon.email@domain.com>.

Preliminary work. Under review by the International Conference on Machine Learning (ICML). Do not distribute.

Based on recent progress in adaptive methods for accelerating the training procedure, see [You et al. \(2019\)](#), we propose a variant of LOCAL AMSGRAD integrating dimensionwise and layerwise adaptive learning rate in each device’s local update. Our contributions are as follows:

- We develop a novel optimization algorithm for federated learning, namely FED-LAMB, following a principled layerwise adaptation strategy to accelerate training of deep neural networks. Our method is provably and empirically communication-efficient for compositional structural models.
- We provide a rigorous theoretical understanding of the non asymptotic convergence rate of FED-LAMB. Based on the recent progress on nonconvex stochastic optimization, we derive for a any finite number of rounds performed by our method, a characterization of the rate at which the classical suboptimality condition, *i.e.*, the second order moment of the gradient of the objective function, decreases. Our bound in $\mathcal{O}\left(\sqrt{\frac{L}{n}} \frac{1}{L\sqrt{R}}\right)$ matches state of the art methods in Federated Learning reaching a sublinear convergence in R , the total number of rounds.
- We exhibit the advantages of our method to reach similar, or better, test accuracy than baseline methods with less number of communication rounds, on several benchmarks supervised learning methods on both homogeneous and heterogeneous settings.

After having established a literature review of both realms of federated and adaptive learning in subsection 1.1, we develop in Section 2, our method, namely FED-LAMB, based on the computation per layer and per dimension, of a scaling factor in the traditional stepsize of AMSGrad. Theoretical understanding of our method’s behaviour with respect to convergence towards a stationary point is developed in Section 3. We present numerical illustrations showing the advantages of our method in Section 4.

1.1. Related Work

Adaptive gradient methods. In recent study on stochastic nonconvex optimization, adaptive methods have proven to be the spearhead in many applications. Those 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. Most used examples include AMSGRAD ([Reddi et al., 2018](#)), ADAM ([Kingma & Ba, 2015](#)), RMSPROP ([Tieleman & Hinton, 2012](#)), ADADELTA ([Zeiler, 2012](#)), and NADAM ([Dozat, 2016](#)).

Their popularity and efficiency are due to their great performance at training deep neural networks. They

generally combine the idea of adaptivity from ADA-GRAD ([Duchi et al., 2011](#); [McMahan & 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. Its update has a notable feature: it leverages an anisotropic learning rate depending on the magnitude of the gradient for each dimension which helps in exploiting the geometry of the data.

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 & Ba \(2015\)](#) for more details. Consequently, ([Kingma & 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 and is presented in Algorithm 1. The difference between ADAM and AMSGRAD lies in Line 1.

Algorithm 1 AMSGRAD ([Reddi et al., 2018](#))

```

1: Required: parameter  $\beta_1, \beta_2$ , and  $\eta_t$ .
2: Init:  $w_1 \in \Theta \subseteq \mathbb{R}^d$  and  $v_0 = \epsilon \mathbf{1} \in \mathbb{R}^d$ .
3: for  $t = 1$  to  $T$  do
4:   Get mini-batch stochastic gradient  $g_t$  at  $w_t$ .
5:    $\theta_t = \beta_1 \theta_{t-1} + (1 - \beta_1) g_t$ .
6:    $v_t = \beta_2 v_{t-1} + (1 - \beta_2) g_t^2$ .
7:    $\hat{v}_t = \max(\hat{v}_{t-1}, v_t)$ .
8:    $w_{t+1} = w_t - \eta_t \frac{\theta_t}{\sqrt{\hat{v}_t}}$ . (element-wise division)
9: end for
```

A natural extension of Algorithm 1 has been developed in [You et al. \(2019\)](#) specifically for multi layered neural network. A principled layerwise adaptation strategy to accelerate training of deep neural networks using large mini-batches is proposed using either a standard stochastic gradient update or a generalized adaptive method under the setting of a classical single server empirical risk minimization problem. 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, layerwise 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-AVG – extending parallel SGD with local updates performed on each device. In FED-AVG, each worker updates their own model parameters locally using SGD, and the local models are synchronized by periodic averaging on the central parameter server.

2. Layerwise and Dimensionwise Adaptive Methods

Beforehand, it is important to provide useful and important notations used throughout our paper.

Notations: We denote by θ the vector of parameters taking values in \mathbb{R}^d . For each layer $\ell \in \llbracket L \rrbracket$, where L is the total number of layers of the neural networks, and each coordinate $j \in \llbracket p_\ell \rrbracket$ where p_ℓ is the dimension per layer ℓ , we note $\theta^{\ell,j}$ its j th coordinate. The gradient of f with respect to θ^ℓ is denoted by $\nabla_\ell f(\theta)$. The index $i \in \llbracket n \rrbracket$ denotes the index of the worker i in our federated framework. r and t are used as the round and local iteration numbers respectively. The smoothness per layer is denoted by L_ℓ for each layer $\ell \in \llbracket L \rrbracket$. We note for each communication $r > 0$, the set of randomly drawn devices D^r performing local updates.

2.1. AMSGrad, Local AMSGrad and Periodic Averaging

Under our Federated setting, we stress on the important of reducing the communication cost at each round 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 shift the computation burden from the many clients to the central server as much as possible. This 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 & Cong (2017). The heuristic efficiency of local updates using SGD and periodic averaging has been studied in Stich (2018); 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, such as 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.

2.2. Layerwise and Dimensionwise Learning with Periodic Averaging

Recall that our original problem is the following optimization task:

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

where $f_i(\theta)$ is the loss function associated to the client $i \in \llbracket n \rrbracket$ and is parameterized, in our paper, by a deep neural network. The multilayer and nonconvex nature of the loss function implies having recourse to particular optimization methods in order to efficiently train our model. Besides, the distributed and clients low bandwidth constraints are strong motivations for improving existing methods performing (1).

Based on the periodic averaging and local AMSGrad algorithms, presented prior, we propose a layerwise and dimensionwise local AMS algorithm which is depicted in Figure 1 and detailed in Algorithm 2, which is a natural adaptation of the vanilla AMSGrad method, for *multilayer* neural networks under the *federated* setting. In particular, while Line 2 and Line 2 corresponds to the standard approximation of the first and second moments, via the smooth updates allowed by the tuning parameters β_1 and β_2 respectively and that both Lines 2-2 are correct the biases of those estimates, the final local update in Line 2 is novel and corresponds to the specialization per layer of our federated method. Note that a scaling factor is applied to the learning rate α_r at each round $r > 0$ via the quantity $\phi(\|\theta_{r,i}^{\ell,t-1}\|)$ depending on the dimensionwise and layerwise quantity computed in Line 2. This function is user designed and can be, for instance, set to the identity function. In other words, we normalize the gradient in each layer according to the magnitude of the layer’s weight.

The adaptivity of our federated learning method is thus manifold. There occurs a per dimension normalization with respect to the square root of the second moment used in

adaptive gradient methods and a layerwise normalization obtained via the final local update (Line 2).

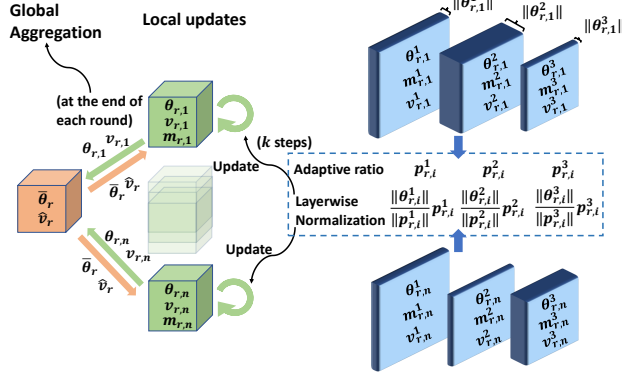


Figure 1. Illustration of Fed-LAMB (Algorithm 2), with a three-layer network and $\phi(x) = x$ as an example. The depth of each network layer represents the norm of its weights. For device i and each local iteration in round r , the adaptive ratio of j -th layer $p_{r,i}^j$ is normalized according to $\|\theta_{r,i}^j\|$, and then used for updating the local model. At the end of each round r , local worker i sends $\theta_{r,i}$ and $v_{r,i}$ to the central server, which transmits back aggregated $\bar{\theta}$ and \bar{v} to local devices to complete a round of training.

Algorithm 2 FED-LAMB for Federated Learning

- 1: **Input:** parameter β_1, β_2 , and learning rate α_t .
- 2: Init: $\theta_0 \in \Theta \subseteq \mathbb{R}^d$, as the global model and $\hat{v}_0 = v_0 = \epsilon \mathbf{1} \in \mathbb{R}^d$ and $\bar{\theta}_0 = \frac{1}{n} \sum_{i=1}^n \theta_0$.
- 3: **for** $r = 1$ to R **do**
- 4: Set $\theta_{r,i}^0 = \bar{\theta}_{r-1}$
- 5: **for parallel for** device $d \in D^r$ **do**
- 6: Compute stochastic gradient $g_{r,i}$ at $\theta_{r,i}$.
- 7: **for** $t = 1$ to T **do**
- 8: $m_{r,i}^t = \beta_1 m_{r-1,i}^{t-1} + (1 - \beta_1) g_{r,i}$.
- 9: $m_{r,i}^t = m_{r,i}^t / (1 - \beta_1^t)$.
- 10: $v_{r,i}^t = \beta_2 v_{r-1,i}^t + (1 - \beta_2) g_{r,i}^2$.
- 11: $v_{r,i}^t = v_{r,i}^t / (1 - \beta_2^t)$.
- 12: Compute ratio $p_{r,i}^t = \frac{m_{r,i}^t}{\sqrt{\hat{v}_{r,i}^t + \epsilon}}$.
- 13: Update local model for each layer $\ell \in \llbracket L \rrbracket$:

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

- 14: **end for**
- 15: Devices send $\theta_{r,i}^T = [\theta_{r,i}^{\ell,T}]_{\ell=1}^L$ and $v_{r,i}^T$ to server.
- 16: **end for**
- 17: Server computes the averages of the local models $\bar{\theta}_r = \frac{1}{n} \sum_{i=1}^n \theta_{r,i}^T$ 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.
- 18: **end for**

3. On The Convergence of FED-LAMB

We develop in this section, the theoretical analysis of Algorithm 2. In Table 3, we recall some important notations that will be used in our following analysis.

R, T	\triangleq	Number of communications rounds and local iterations (resp.)
n, D, i	\triangleq	Total number of clients, portion sampled uniformly and client index
L, ℓ	\triangleq	Total number of layers in the DNN and its index
$\phi(\cdot)$	\triangleq	Scaling factor in FED-LAMBupdate
$\bar{\theta}$	\triangleq	Global model (after periodic averaging)

Table 1. Summary of notations used in the paper.

Based on classical result for stochastic nonconvex optimization, we provide 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:

- The large amount of decentralized workers working solely on their own data stored locally.
- A periodic averaging occurs on the central server pushing each of those clients to send local models after some local iterations.
- Each clients computes a backpropagation of the main model, *i.e.*, the deep neural network, and then update its local version of the model via an adaptive gradient method: the distinctiveness being that those updates are done *dimensionwise* and *layerwise*.

Our analysis encompasses the consideration of those challenges and leads to a 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.

3.1. 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) For $i \in \llbracket n \rrbracket$ and $\ell \in \llbracket L \rrbracket$, f_i is L -smooth: $\|\nabla f_i(\theta) - \nabla f_i(\vartheta)\| \leq L_\ell \|\theta^\ell - \vartheta^\ell\|$.

We add some classical assumption in the unbiased stochastic optimization realm, 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 value $a \in \mathbb{R}_+^*$, there exists strictly positive constants such that $\phi_m \leq \phi(a) \leq \phi_M$.

We now state our main result regarding the non asymptotic convergence analysis of our Algorithm 2:

Theorem 1. Assume **H1-H4**. Consider $\{\bar{\theta}_r\}_{r>0}$, the sequence of parameters obtained running Algorithm 2 with a decreasing learning rate α_r . Then, if the number of local epochs is set to $T = 1$ and $\epsilon = \lambda = 0$, we have:

$$\begin{aligned} & \frac{1}{R} \sum_{r=1}^R \mathbb{E} \left[\left\| \frac{\nabla f(\bar{\theta}_r)}{\sqrt{v_r^t}} \right\|^2 \right] \\ & \leq \sqrt{\frac{M^2 p}{n}} \frac{f(\bar{\theta}_1) - \mathbb{E}[f(\bar{\theta}_{R+1})]}{LR} + \frac{\phi_M \sigma^2}{Rn} \sqrt{\frac{1 - \beta_2}{M^2 p}} \\ & \quad + \alpha \phi_M \sigma L p \sqrt{n} + \frac{\bar{L} \beta_1^2 L (1 - \beta_2) M^2 \phi_M^2 n}{2(1 - \beta_1)^2 v_0} \\ & \quad + \frac{\alpha \beta_1}{1 - \beta_1} \sqrt{(1 - \beta_2) p} \frac{LM^2}{\sqrt{v_0}} + \bar{L} \alpha^2 M^2 \phi_M^2 \frac{(1 - \beta_2) p}{R v_0} \end{aligned} \quad (2)$$

3.2. Important Lemmas for the Proof

Two important Lemmas are required in the proof the above Theorem. We state those two Lemmas above and 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 2. Then for $i \in \llbracket n \rrbracket$ and $r > 0$:

$$\|\bar{\theta}_r - \theta_{r,i}\|^2 \leq \alpha_r^2 M^2 \phi_M^2 \frac{(1 - \beta_2) p}{v_0} \quad (3)$$

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

The gap is provably bounded by some quantities of interest such as the total dimension of the multilayered model p , the

learning rate and the assumed upper bound of the gradient, see H2.

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

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

$$\left\| \frac{\nabla f(\bar{\theta}_r)}{\sqrt{v_r^t}} \right\| \geq \frac{1}{2} \left\| \frac{\nabla f(\bar{\theta}_r)}{\sqrt{v_r^t}} \right\| - \bar{L} \alpha^2 M^2 \phi_M^2 \frac{(1 - \beta_2) p}{v_0} \quad (4)$$

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

We focus in the next subsection on two particular papers of utmost interest for our contribution.

3.3. Comparison with LAMB and Local-AMS

We dedicate the following paragraph to a discussion on the bound derived above in comparison with known results in the literature.

LAMB bound in You et al. (2019): We first start our discussion with the comparison of convergence rate of FED-LAMB with that of LAMB, Theorem 3 in You et al. (2019). The convergence rates of FED-LAMB and LAMB differ in two ways: (i) the suboptimality, or convergence criterion,

used here is $\frac{1}{R} \sum_{r=1}^R \mathbb{E} \left[\left\| \frac{\nabla f(\bar{\theta}_r)}{\sqrt{v_r^t}} \right\|^2 \right]$ for a total number of

rounds R as opposed to $\mathbb{E} [\|\nabla f(\theta_{\mathcal{R}})\|^2]$ for some random termination round \mathcal{R} uniformly drawn from $\llbracket R \rrbracket$. First, note that the characterization is given at the averaged parameters noted $\bar{\theta}_{\mathcal{R}}$ due to our distributed settings. It is thus natural to consider the evolution of our objective function, its gradient to be precise, 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 (2) 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 & Lan, 2013). (ii) Assuming that the convergence criterion in both Theorems is of similar order (which happens for a large enough number of rounds), convergence rate of FED-LAMB displays a similar $\mathcal{O}(\frac{1}{R})$ behaviour for the initialization term, meaning that, despite the distributed (federated) settings, our dimen-

sionwise and layerwise method benefits from the double adaptivity phenomenon explained above and exhibited in the LAMB method of You et al. (2019), performed under a central server setting.

Local-AMS bound in Chen et al. (2020): We now discuss the similarities and differences between the distributed adaptive method developed in Chen et al. (2020) and named local-AMS, and our *deep federated* method, namely FED-LAMB. We first recall their main result:

Theorem 2 (Theorem 5.1 in Chen et al. (2020)). *Under some regularity conditions on the local losses and similar arguments as ours, LOCAL-AMS reaches a stationary point with the following rate:*

$$\begin{aligned} & \frac{1}{R} \sum_{r=1}^R \mathbb{E} \left[\left\| \frac{\nabla f(\bar{\theta}_r)}{\sqrt{v_r^t}} \right\|^2 \right] \\ & \leq 8\sqrt{\frac{p}{Rn}} [f(\bar{\vartheta}_1) - \mathbb{E}[f(\bar{\vartheta}_{R+1})]] \\ & \quad + 8L_s \sqrt{\frac{p}{Rn}} \frac{\sigma^2}{\epsilon} + cst. \end{aligned} \quad (5)$$

where ϵ corresponds to their initialization of the vector \hat{v}_0 and L_s is the sum of the local smoothness constants.

The first two terms of their results and ours, standard in convergence analysis, displays a strong dependence of the convergence rate on the initialization, which tends to be forgotten, and the bounded variance assumption of the stochastic gradient, see H3. The acceleration of our layerwise scheme is exhibited in the dependence on $\frac{1}{R}$ in those two terms, while results in Chen et al. (2020) are of order $\frac{1}{\sqrt{R}}$. Note that the boundedness assumption is done on each dimension in H3 and leads to a manifestation of the term \sqrt{p} in both rates. This can be handled for simplicity and clarity of the results when H3 is assumed globally. It is important to note that their rate include an analysis with respect to the number of local updates, noted T in this paper. While our result is derived for a single local update $T = 1$, we acknowledge that exhibiting a dependency on T can be interesting, particularly in order to see the impact of several local updates which are, in practice, in different directions –as a byproduct of heterogeneous local training samples. We leave this investigation to future work.

In light of the prior remarks, we now give a simple variant of Theorem 1 in the following Corollary where we assume an upperbound on the norm of the second moment approximate $\sqrt{v_r^t}$:

H5. *For $t > 0$ and $r > 0$, there exists some constant such that $\|\sqrt{v_r^t}\| \leq V^2$.*

Corollary 1. *Assume H1-H4. Consider $\{\bar{\theta}_r\}_{r>0}$, the sequence of parameters obtained running Algorithm 2 and set $\alpha_r = \mathcal{O}(\frac{1}{L\sqrt{R}})$. Then, if the number of local epochs is set to $T = 1$, $\epsilon = \lambda = 0$, and $R \geq$ we have, under H5:*

$$\frac{1}{R} \mathbb{E} [\|\nabla f(\bar{\theta}_R)\|^2] \leq \mathcal{O} \left(\sqrt{\frac{p}{n}} \frac{1}{L\sqrt{R}} \right) \quad (6)$$

We now discuss our bound regarding several aspects in order to gain understanding of the advantages of our method:

- **Communication Complexity:** The (sublinear) dependence on the number of communication rounds of our bound matches that of most recent methods in Federated Learning, see Karimireddy et al. (2019) developing SCAFFOLD, a solution to the problem posed by heterogeneity of the data in each client, and of Reddi et al. (2020), adapting state of the art method in optimization, here ADAM, under the federated setting. 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. Novelty in our bound occurs considering the sublinear dependence on the number of layer L and the dependence on the total number of dimension of our problem p . For the latter, the dependency can be simplified with stronger assumption on the control of the variance of the stochastic gradient. For the former,
- **Homogeneous and Heterogeneous Data:** A common assumption regarding the stochastic gradient, and its variance, considers an upper-bound of the global variance, in the sense that it applies to both the aggregated objective function (1) and each of its local component. Another theoretical approach is to set apart a local variance, for each local loss, and global variance, for their sum, see Chen et al. (2020) for instance. Heterogeneity is of utmost importance in FL since client may store radically different data points. Existing methods can lead to poor convergence as detailed in (Li et al., 2020; Liang et al., 2019). Improvements are proposed through the use of gradient tracking techniques performed locally as seen in (Haddadpour et al., 2020a; Horváth et al., 2019; Karimireddy et al., 2019).
- **Dependence on the dimension p :** The dependence on the overall size of the vectors being transmitted back and forth from the central to the devices can be drastically improved. Indeed, recent efficient techniques aim at reducing the number of bits communicated at each round through sketches or compression techniques, see

for instance (Haddadpour et al., 2020b; Ivkin et al., 2019; Li et al., 2019) to name a few. This can be a great addition to FED-LAMB, and is naturally compatible, but is not the main focus of our contribution.

4. Experiments

In this section, we conduct numerical experiments on various datasets and network architectures to testify the effectiveness of our proposed method in practice. Our main objective is to validate the benefit of dimensionwise adaptive learning when integrated with adaptive Fed-AMS. We observe in the sequel that our method empirically confirms its edge in terms of convergence speed. Basically, our proposed method reduces the number of rounds and thus the communication cost required to achieve similar stationary points than baseline methods.

Settings. In our experiment, we will evaluate three federated learning algorithms: 1) Fed-SGD, 2) Fed-AMS and 3) our proposed Fed-LAMB (Algorithm 2), where the first two serve as the baseline methods. For adaptive methods 2) and 3), we set $\beta_1 = 0.9$, $\beta_2 = 0.999$ as default and recommended (Reddi et al., 2018). Regarding federated learning environment, we use 50 local workers with 0.5 participation rate. That means, we randomly pick half of the workers to be active for training in each round. To best accord with real scenarios where the local training batch size is usually limited, we set a relatively small local update batch size as 32. 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 batch training. We test different number of local epochs in our experiments. For each dataset and number of local epochs, we tune the constant learning rate α for each algorithm over a fine grid in logarithm scale. For Fed-LAMB, the parameter λ in Algorithm 2 controlling the overall scale of the layerwise gradients is tuned from $\{0, 0.01, 0.1\}$. For each run, we take the model performance with the best α and λ . The reported results are averaged over three independent runs each with same initialization.

Models. We test the performance of different federated learning algorithms on MNIST (LeCun, 1998) and CIFAR10 (Krizhevsky & Hinton, 2009) image classification datasets. For MNIST, we apply 1) a simple multilayer 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. For CIFAR10, we implement: 1) a CNN with three convolutional layers followed by two fully-connected layers, and 2) a ResNet-9 model proposed by (He et al., 2016). All the networks use ReLU

as the activation function.

4.1. MNIST with Multilayer Perceptron and Convolutional Neural Network

In Figure 2, we start by presenting the test accuracy on MNIST dataset trained by MLP. We compare each method under the heterogeneous (non-iid) data distribution settings, where each device only receives samples of one digit (out of ten). This is known to be the scenario where federated learning is harder to generalize well, see (McMahan et al., 2017). First of all, we observe that our proposed method FED-LAMB outperforms Fed-AMS and Fed-SGD with both 1 and 5 local epochs, illustrating its improvement over baseline federated methods. In addition, though it is not the main comparison of this paper, Fed-SGD slightly generalizes better than Fed-AMS, which means that SGD might be sufficient for this rather simple model. Yet, Fed-LAMB overachieves both methods on this task in terms of test accuracies. Similar comparison are drawn with respect to the training and testing losses.

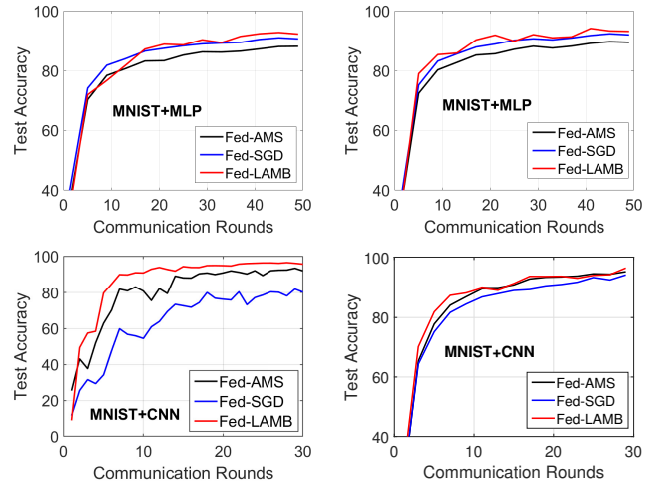


Figure 2. **Top Row:** Test accuracy on MNIST+MLP, with non-iid data distribution. **Bottom Row:** Test accuracy on MNIST+CNN, with non-iid data distribution. **Left panel:** 1 local epoch. **Right panel:** 5 local epochs.

We also evaluate the various algorithms on larger models. Since Fed-LAMB is specifically designed for multi-layer deep learning networks, we expect it to show more substantial advantage on larger network architectures. In Figure 2, we present the results on MNIST with CNN, under *non-iid data distribution*. In general, we again see that Fed-LAMB outperforms Fed-AMS and Fed-SGD. The advantage is in particular significant with 1 local epoch, where Fed-SGD generalizes poorly. Importantly, we would like to address the acceleration effect of Fed-LAMB, in the early stage of

training. We observe that Fed-LAMB converges faster than the vanilla Fed-AMS at first few communication rounds, in both cases. As a side note, the poor performance of Fed-SGD is, to some extent, consistent with prior literature and practical numerical experiments showing that adaptive gradient methods usually perform better than simple SGD in training large deep learning models. We refer the readers to a collection of prior studies such as (Chen et al., 2020; Reddi et al., 2020).

4.2. CIFAR-10 with Convolutional Neural Network and Residual Neural Network

In Figure 3, we report the test accuracies of a Convolutional Neural Network trained on CIFAR-10 dataset, where the data is iid allocated among clients. When we run only 1 local epoch per device, we observe a clear advantage of FED-LAMBover Fed-AMS on both test accuracy and convergence speed. Note that Fed-SGD again fails to achieve close performance as those two adaptive gradient methods. Increasing the number of local iteration in each device per communication round leads to similar observations: our newly introduced method, namely FED-LAMB, converges the fastest, with similar generalization error as Fed-AMS.

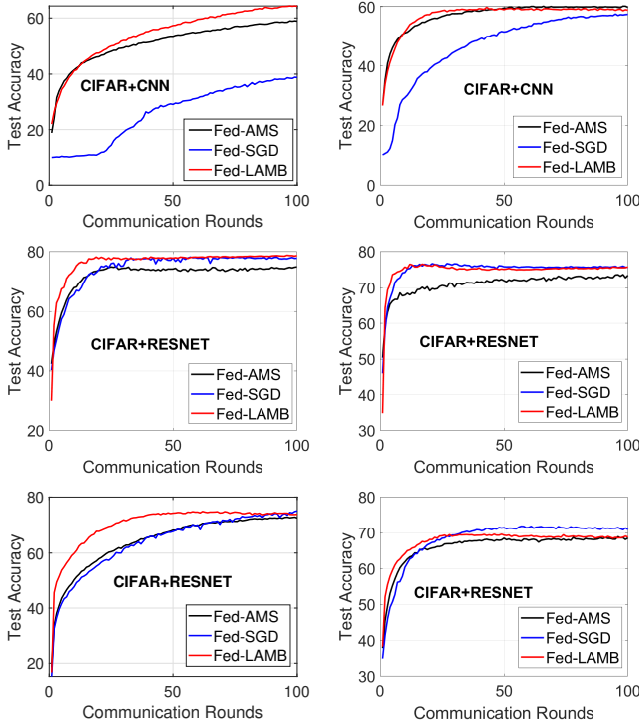


Figure 3. Test accuracy on CIFAR+CNN, with iid data distribution. Left panel: 1 local epoch. Right panel: 3 local epochs.

Lastly, we test the algorithms on CIFAR-10 using a Residual

Neural Network, the ResNet-9 model. We observe Figure 3 that while our method FED-LAMB improves the performance of vanilla local AMS method under various settings regarding the number clients and of local iterations, we also observe that in this particular example, Fed-SGD is reaching similar, and sometimes better, test performances.

5. Conclusion

We study in this contribution a doubly adaptive method in the particular framework of federated learning. Built upon the success of periodic averaging, and of state-of-the-arts adaptive gradient methods for single server non-convex stochastic optimization, we derive FED-LAMB, a distributed AMSGrad method that performs local updates on each worker and periodically averages local models. Besides, a core component of FED-LAMB, when the trained model is a deep neural network, is a layerwise update of each local model. Proved convergence guarantees of our scheme are provided in our contribution and exhibits a sub-linear dependence on the total number of communications rounds, the number of clients and the number of layers of the model. The multiple benefits of periodic averaging, adaptive optimization methods and layerwise updates are displayed in our bounds. We empirically confirm the advantage of our algorithm over baselines methods on a panel of numerical experiments.

References

- Alistarh, D., Grubic, D., Li, J., Tomioka, R., and Vojnovic, M. Qsgd: Communication-efficient sgd via gradient quantization and encoding. In *Advances in Neural Information Processing Systems*, pp. 1709–1720, 2017.
- Chen, X., Li, X., and Li, P. Toward communication efficient adaptive gradient method. In *ACM-IMS Foundations of Data Science Conference (FODS)*, Seattle, WA, 2020.
- Dozat, T. Incorporating nesterov momentum into adam. *ICLR (Workshop Track)*, 2016.
- Duchi, J., Hazan, E., and Singer, Y. Adaptive subgradient methods for online learning and stochastic optimization. *Journal of Machine Learning Research (JMLR)*, 2011.
- Ghadimi, S. and Lan, G. Stochastic first- and zeroth-order methods for nonconvex stochastic programming. *SIAM Journal on Optimization*, 23(4):2341–2368, 2013.
- Haddadpour, F., Kamani, M. M., Mokhtari, A., and Mahdavi, M. Federated learning with compression: Unified analysis and sharp guarantees. *arXiv preprint arXiv:2007.01154*, 2020a.

Haddadpour, F., Karimi, B., Li, P., and Li, X. FedSketch:

- Communication-efficient and private federated learning via sketching. *arXiv preprint arXiv:2008.04975*, 2020b.
- He, K., Zhang, X., Ren, S., and Sun, J. Deep residual learning for image recognition. In *2016 IEEE Conference on Computer Vision and Pattern Recognition, CVPR 2016, Las Vegas, NV, USA, June 27-30, 2016*, pp. 770–778. IEEE Computer Society, 2016.
- Horváth, S., Kovalev, D., Mishchenko, K., Stich, S., and Richtárik, P. Stochastic distributed learning with gradient quantization and variance reduction. *arXiv preprint arXiv:1904.05115*, 2019.
- Ivkin, N., Rothchild, D., Ullah, E., Braverman, V., Stoica, I., and Arora, R. Communication-efficient distributed SGD with sketching. In *Advances in Neural Information Processing Systems (NeurIPS)*, pp. 13144–13154, Vancouver, Canada, 2019.
- Karimireddy, S. P., Kale, S., Mohri, M., Reddi, S. J., Stich, S. U., and Suresh, A. T. Scaffold: Stochastic controlled averaging for on-device federated learning. *arXiv preprint arXiv:1910.06378*, 2019.
- Kingma, D. P. and Ba, J. Adam: A method for stochastic optimization. *ICLR*, 2015.
- Konečný, J., McMahan, H. B., Yu, F. X., Richtárik, P., Suresh, A. T., and Bacon, D. Federated learning: Strategies for improving communication efficiency. *arXiv preprint arXiv:1610.05492*, 2016.
- Krizhevsky, A. and Hinton, G. Learning multiple layers of features from tiny images. *Master's thesis, Department of Computer Science, University of Toronto*, 2009.
- LeCun, Y. The mnist database of handwritten digits. <http://yann.lecun.com/exdb/mnist/>, 1998.
- 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 *11th {USENIX} Symposium on Operating Systems Design and Implementation ({OSDI} 14)*, pp. 583–598, 2014.
- Li, T., Liu, Z., Sekar, V., and Smith, V. Privacy for free: Communication-efficient learning with differential privacy using sketches. *arXiv preprint arXiv:1911.00972*, 2019.
- Li, T., Sahu, A. K., Talwalkar, A., and Smith, V. Federated learning: Challenges, methods, and future directions. *IEEE Signal Process. Mag.*, 37(3):50–60, 2020.
- Liang, X., Shen, S., Liu, J., Pan, Z., Chen, E., and Cheng, Y. Variance reduced local sgd with lower communication complexity. *arXiv preprint arXiv:1912.12844*, 2019.
- Lin, Y., Han, S., Mao, H., Wang, Y., and Dally, W. J. Deep gradient compression: Reducing the communication bandwidth for distributed training. *arXiv preprint arXiv:1712.01887*, 2017.
- McMahan, B., Moore, E., Ramage, D., Hampson, S., and y Arcas, B. A. Communication-efficient learning of deep networks from decentralized data. In *Artificial Intelligence and Statistics*, pp. 1273–1282. PMLR, 2017.
- McMahan, H. B. and Streeter, M. J. Adaptive bound optimization for online convex optimization. *COLT*, 2010.
- Nesterov, Y. Introductory lectures on convex optimization: A basic course. *Springer*, 2004.
- Polyak, B. T. Some methods of speeding up the convergence of iteration methods. *Mathematics and Mathematical Physics*, 1964.
- Recht, B., Re, C., Wright, S., and Niu, F. Hogwild!: A lock-free approach to parallelizing stochastic gradient descent. *Advances in neural information processing systems*, 24: 693–701, 2011.
- Reddi, S., Charles, Z., Zaheer, M., Garrett, Z., Rush, K., Konečný, J., Kumar, S., and McMahan, H. B. Adaptive federated optimization. *arXiv preprint arXiv:2003.00295*, 2020.
- Reddi, S. J., Kale, S., and Kumar, S. On the convergence of adam and beyond. *ICLR*, 2018.
- Stich, S. U. Local sgd converges fast and communicates little. *arXiv preprint arXiv:1805.09767*, 2018.
- Tieleman, T. and Hinton, G. Rmsprop: Divide the gradient by a running average of its recent magnitude. *COURS-ERA: Neural Networks for Machine Learning*, 2012.
- Wangni, J., Wang, J., Liu, J., and Zhang, T. Gradient sparsification for communication-efficient distributed optimization. In *Advances in Neural Information Processing Systems*, pp. 1299–1309, 2018.
- You, Y., Li, J., Reddi, S., Hseu, J., Kumar, S., Bhojanapalli, S., Song, X., Demmel, J., Keutzer, K., and Hsieh, C.-J. Large batch optimization for deep learning: Training bert in 76 minutes. *arXiv preprint arXiv:1904.00962*, 2019.
- Yu, H., Jin, R., and Yang, S. On the linear speedup analysis of communication efficient momentum sgd for distributed non-convex optimization. *arXiv preprint arXiv:1905.03817*, 2019.
- Zeiler, M. D. Adadelta: An adaptive learning rate method. *arXiv:1212.5701*, 2012.

Zhao, W., Xie, D., Jia, R., Qian, Y., Ding, R., Sun, M., and
Li, P. Distributed hierarchical gpu parameter server for
massive scale deep learning ads systems. *arXiv preprint*
arXiv:2003.05622, 2020.

Zhou, F. and Cong, G. On the convergence properties
of a k -step averaging stochastic gradient descent al-
gorithm for nonconvex optimization. *arXiv preprint*
arXiv:1708.01012, 2017.

A. Theoretical Analysis

A.1. Intermediary Lemmas

Lemma. Consider $\{\bar{\theta}_r\}_{r>0}$, the sequence of parameters obtained running Algorithm 2. Then for $i \in \llbracket n \rrbracket$:

$$\|\bar{\theta}_r - \theta_{r,i}\|^2 \leq \alpha^2 M^2 \phi_M^2 \frac{(1 - \beta_2)p}{v_0} \quad (7)$$

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

Proof. Assuming the simplest case when $T = 1$, i.e. one local iteration, then by construction of Algorithm 2, we have for all $\ell \in \llbracket L \rrbracket$, $i \in \llbracket n \rrbracket$ and $r > 0$:

$$\theta_{r,i}^\ell = \bar{\theta}_r^\ell - \alpha \phi(\|\theta_{r,i}^{\ell,t-1}\|) p_{r,i}^j / \|p_{r,i}^\ell\| = \bar{\theta}_r^\ell - \alpha \phi(\|\theta_{r,i}^{\ell,t-1}\|) \frac{m_{r,i}^t}{\sqrt{v_r^t}} \frac{1}{\|p_{r,i}^\ell\|} \quad (8)$$

leading to

$$\begin{aligned} \|\bar{\theta}_r - \theta_{r,i}\|^2 &= \sum_{\ell=1}^L \langle \bar{\theta}_r^\ell - \theta_{r,i}^\ell | \bar{\theta}_r^\ell - \theta_{r,i}^\ell \rangle \\ &\leq \alpha^2 M^2 \phi_M^2 \frac{(1 - \beta_2)p}{v_0} \end{aligned} \quad (9)$$

which concludes the proof. \square

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

$$\left\| \frac{\bar{\nabla} f(\bar{\theta}_r)}{\sqrt{v_r^t}} \right\| \geq \frac{1}{2} \left\| \frac{\nabla f(\bar{\theta}_r)}{\sqrt{v_r^t}} \right\| - \bar{L} \alpha^2 M^2 \phi_M^2 \frac{(1 - \beta_2)p}{v_0} \quad (10)$$

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

Proof. Consider the following sequence:

$$\left\| \frac{\bar{\nabla} f(\bar{\theta}_r)}{\sqrt{v_r^t}} \right\| \geq \frac{1}{2} \left\| \frac{\nabla f(\bar{\theta}_r)}{\sqrt{v_r^t}} \right\| - \left\| \frac{\bar{\nabla} f(\bar{\theta}_r) - \nabla f(\bar{\theta}_r)}{\sqrt{v_r^t}} \right\| \quad (11)$$

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

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

$$\begin{aligned} \left\| \frac{\bar{\nabla} f(\bar{\theta}_r)}{\sqrt{v_r^t}} \right\| &\geq \frac{1}{2} \left\| \frac{\nabla f(\bar{\theta}_r)}{\sqrt{v_r^t}} \right\| - \left\| \frac{\bar{\nabla} f(\bar{\theta}_r) - \nabla f(\bar{\theta}_r)}{\sqrt{v_r^t}} \right\| \\ &\geq \frac{1}{2} \left\| \frac{\nabla f(\bar{\theta}_r)}{\sqrt{v_r^t}} \right\| - \bar{L} \alpha^2 M^2 \phi_M^2 \frac{(1 - \beta_2)p}{v_0} \end{aligned} \quad (12)$$

concluding the proof. \square

A.2. Proof of Theorem 1

Theorem. Consider $\{\bar{\theta}_r\}_{r>0}$, the sequence of parameters obtained running Algorithm 2. Then, if the number of local epochs is set to $T = 1$ and $\epsilon = \lambda = 0$, we have:

$$\frac{1}{R} \sum_{r=1}^R \mathbb{E}[\|\nabla f(\bar{\theta}_r)\|^2] \leq dd \quad (13)$$

Case with $T = 1, \epsilon = 0$ and $\lambda = 0$: Using H1, we have:

$$\begin{aligned} f(\bar{\vartheta}_{r+1}) &\leq f(\bar{\vartheta}_r) + \langle \nabla f(\bar{\vartheta}_r) | \bar{\vartheta}_{r+1} - \bar{\vartheta}_r \rangle + \sum_{\ell=1}^L \frac{L_\ell}{2} \|\bar{\vartheta}_{r+1}^\ell - \bar{\vartheta}_r^\ell\|^2 \\ &\leq f(\bar{\vartheta}_r) + \sum_{\ell=1}^L \sum_{j=1}^{p_\ell} \nabla_\ell f(\bar{\vartheta}_r)^j (\bar{\vartheta}_{r+1}^{\ell,j} - \bar{\vartheta}_r^{\ell,j}) + \sum_{\ell=1}^L \frac{L_\ell}{2} \|\bar{\vartheta}_{r+1}^\ell - \bar{\vartheta}_r^\ell\|^2 \end{aligned} \quad (14)$$

Taking expectations on both sides leads to:

$$-\mathbb{E}[\langle \nabla f(\bar{\vartheta}_r) | \bar{\vartheta}_{r+1} - \bar{\vartheta}_r \rangle] \leq \mathbb{E}[f(\bar{\vartheta}_r) - f(\bar{\vartheta}_{r+1})] + \sum_{\ell=1}^L \frac{L_\ell}{2} \mathbb{E}[\|\bar{\vartheta}_{r+1}^\ell - \bar{\vartheta}_r^\ell\|^2] \quad (15)$$

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:

$$\bar{\vartheta}_r = \bar{\theta}_r + \frac{\beta_1}{1 - \beta_1} (\bar{\theta}_r - \bar{\theta}_{r-1}) \quad (16)$$

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

$$\bar{\vartheta}_{r+1}^\ell - \bar{\vartheta}_r^\ell = \frac{1}{1 - \beta_1} (\bar{\theta}_{r+1}^\ell - \bar{\theta}_r^\ell) - \frac{\beta_1}{1 - \beta_1} (\bar{\theta}_r^\ell - \bar{\theta}_{r-1}^\ell) \quad (17)$$

$$= \frac{\alpha_r}{1 - \beta_1} \frac{1}{n} \sum_{i=1}^n \frac{\phi(\|\theta_{r,i}^\ell\|)}{\|p_{r,i}^\ell\|} p_{r,i}^\ell - \frac{\alpha_{r-1}}{1 - \beta_1} \frac{1}{n} \sum_{i=1}^n \frac{\phi(\|\theta_{r-1,i}^\ell\|)}{\|p_{r-1,i}^\ell\|} p_{r-1,i}^\ell \quad (18)$$

$$= \frac{\alpha\beta_1}{1 - \beta_1} \frac{1}{n} \sum_{i=1}^n \left(\frac{\phi(\|\theta_{r,i}^\ell\|)}{\sqrt{v_r^t} \|p_{r,i}^\ell\|} - \frac{\phi(\|\theta_{r-1,i}^\ell\|)}{\sqrt{v_{r-1}^t} \|p_{r-1,i}^\ell\|} \right) m_{r-1}^t + \frac{\alpha}{n} \sum_{i=1}^n \frac{\phi(\|\theta_{r,i}^\ell\|)}{\sqrt{v_r^t} \|p_{r,i}^\ell\|} g_{r,i} \quad (19)$$

where we have assumed a constant learning rate α .

We note for all $\theta \in \Theta$, the majorant $G > 0$ such that $\phi(\|\theta\|) \leq G$. Then, following (15), we obtain:

$$-\mathbb{E}[\langle \nabla f(\bar{\vartheta}_r) | \bar{\vartheta}_{r+1} - \bar{\vartheta}_r \rangle] \leq \mathbb{E}[f(\bar{\vartheta}_r) - f(\bar{\vartheta}_{r+1})] + \sum_{\ell=1}^L \frac{L_\ell}{2} \mathbb{E}[\|\bar{\vartheta}_{r+1}^\ell - \bar{\vartheta}_r^\ell\|^2] \quad (20)$$

Developing the LHS of (20) using (17) leads to

$$\langle \nabla f(\bar{\vartheta}_r) | \bar{\vartheta}_{r+1} - \bar{\vartheta}_r \rangle = \sum_{\ell=1}^L \sum_{j=1}^{p_\ell} \nabla_\ell f(\bar{\vartheta}_r)^j (\bar{\vartheta}_{r+1}^{\ell,j} - \bar{\vartheta}_r^{\ell,j}) \quad (21)$$

$$= \frac{\alpha\beta_1}{1 - \beta_1} \frac{1}{n} \sum_{\ell=1}^L \sum_{j=1}^{p_\ell} \nabla_\ell f(\bar{\vartheta}_r)^j \left[\sum_{i=1}^n \left(\frac{\phi(\|\theta_{r,i}^\ell\|)}{\sqrt{v_r^t} \|p_{r,i}^\ell\|} - \frac{\phi(\|\theta_{r-1,i}^\ell\|)}{\sqrt{v_{r-1}^t} \|p_{r-1,i}^\ell\|} \right) m_{r-1}^t \right] \quad (22)$$

$$- \underbrace{\frac{\alpha}{n} \sum_{\ell=1}^L \sum_{j=1}^{p_\ell} \nabla_\ell f(\bar{\vartheta}_r)^j \sum_{i=1}^n \frac{\phi(\|\theta_{r,i}^\ell\|)}{\sqrt{v_r^t} \|p_{r,i}^\ell\|} g_{r,i}^j}_{=A_1} \quad (23)$$

Term A_1 : Since we have that $\|p_{r,i}^\ell\| \leq \sqrt{\frac{p_\ell}{1-\beta_2}}$ and $1/\sqrt{v_r^t} \leq 1/\sqrt{v_0}$, using H2, we develop the term A_1 as follows:

$$A_1 \leq -\frac{\alpha}{n} \sum_{\ell=1}^L \sum_{j=1}^{p_\ell} \nabla_\ell f(\bar{\vartheta}_r)^j \sum_{i=1}^n \frac{\phi(\|\theta_{r,i}^\ell\|)}{\sqrt{v_r^t} \|p_{r,i}^\ell\|} g_{r,i}^{l,j} \quad (24)$$

$$\leq -\frac{\alpha}{n} \sum_{\ell=1}^L \sqrt{\frac{1-\beta_2}{M^2 p_\ell}} \sum_{i=1}^n \sum_{j=1}^{p_\ell} \phi(\|\theta_{r,i}^\ell\|) \nabla_\ell f(\bar{\vartheta}_r)^j g_{r,i}^{\ell,j} \quad (25)$$

$$- \frac{\alpha}{n} \sum_{\ell=1}^L \sum_{i=1}^n \sum_{j=1}^{p_\ell} \left(\phi(\|\theta_{r,i}^\ell\|) \nabla_\ell f(\bar{\vartheta}_r)^j \frac{p_{r,i}^\ell}{\|p_{r,i}^\ell\|} \right) \mathbf{1}(\text{sign}(\nabla_\ell f(\bar{\vartheta}_r)^j) \neq \text{sign}(p_{r,i}^\ell)) \quad (26)$$

Taking the expectations on both sides yields and using H4:

$$\mathbb{E}[A_1] \leq -\alpha \sum_{\ell=1}^L \sqrt{\frac{1-\beta_2}{M^2 p_\ell}} \sum_{i=1}^n \sum_{j=1}^{p_\ell} \mathbb{E} \left[\phi(\|\theta_{r,i}^\ell\|) \nabla_\ell f(\bar{\vartheta}_r)^j g_{r,i}^{\ell,j} \right] \quad (27)$$

$$- \frac{\alpha}{n} \sum_{\ell=1}^L \sum_{i=1}^n \sum_{j=1}^{p_\ell} \mathbb{E} \left[\phi(\|\theta_{r,i}^\ell\|) \nabla_\ell f(\bar{\vartheta}_r)^j \frac{p_{r,i}^\ell}{\|p_{r,i}^\ell\|} \mathbf{1}(\text{sign}(\nabla_\ell f(\bar{\vartheta}_r)^j) \neq \text{sign}(p_{r,i}^\ell)) \right] \quad (28)$$

$$\leq -\frac{\alpha}{n} \sum_{\ell=1}^L \phi_m \sqrt{\frac{1-\beta_2}{M^2 p_\ell}} \sum_{i=1}^n \sum_{j=1}^{p_\ell} (\nabla_\ell f(\bar{\vartheta}_r)^j)^2 \quad (29)$$

$$- \frac{\alpha}{n} \sum_{\ell=1}^L \sum_{i=1}^n \sum_{j=1}^{p_\ell} \phi_M \mathbb{E} \left[\left| \nabla_\ell f(\bar{\vartheta}_r)^j \frac{p_{r,i}^\ell}{\|p_{r,i}^\ell\|} \right| \mathbf{1}(\text{sign}(\nabla_\ell f(\bar{\vartheta}_r)^j) \neq \text{sign}(p_{r,i}^\ell)) \right] \quad (30)$$

$$(31)$$

where we have used assumption H4.

Since for any ℓ, i, j , we have

$$\mathbb{E} \left[\left| \nabla_\ell f(\bar{\vartheta}_r)^j \frac{p_{r,i}^\ell}{\|p_{r,i}^\ell\|} \right| \mathbf{1}(\text{sign}(\nabla_\ell f(\bar{\vartheta}_r)^j) \neq \text{sign}(p_{r,i}^\ell)) \right] \leq |\nabla_\ell f(\bar{\vartheta}_r)^j| \mathbb{P}(\text{sign}(\nabla_\ell f(\bar{\vartheta}_r)^j) \neq \text{sign}(p_{r,i}^\ell)) \quad (32)$$

Then, we obtain

$$\mathbb{E}[A_1] \leq -\alpha \phi_m \sqrt{\frac{L(1-\beta_2)}{M^2 p}} \mathbb{E}[\|\nabla f(\bar{\vartheta}_r)\|^2] - \alpha \phi_M \sum_{\ell=1}^L \sum_{i=1}^n \sum_{j=1}^{p_\ell} \frac{\sigma_i^{\ell,j}}{\sqrt{n}} \quad (33)$$

where for any $\theta \in \Theta$ we define $\nabla f(\cdot) = \sum_{i=1}^n \nabla f_i(\cdot)$.

We now need to bound the following terms:

$$A_r^2 := \mathbb{E}[\|\bar{\vartheta}_{r+1} - \bar{\vartheta}_r\|^2] \quad (34)$$

$$A_r^3 := \frac{\alpha \beta_1}{1-\beta_1} \frac{1}{n} \sum_{\ell=1}^L \sum_{j=1}^{p_\ell} \nabla_\ell f(\bar{\vartheta}_r)^j \left[\sum_{i=1}^n \left(\frac{\phi(\|\theta_{r,i}^\ell\|)}{\sqrt{v_r^t} \|p_{r,i}^\ell\|} - \frac{\phi(\|\theta_{r-1,i}^\ell\|)}{\sqrt{v_{r-1}^t} \|p_{r-1,i}^\ell\|} \right) m_{r-1}^t \right] \quad (35)$$

Term A_r^2 : According to definition (16), for each layer $\ell \in \llbracket L \rrbracket$, we have, using the Cauchy-Schwartz inequality, that:

$$\|\bar{\vartheta}_{r+1}^\ell - \bar{\vartheta}_r^\ell\|^2 = \left\| \frac{\alpha\beta_1}{1-\beta_1} \frac{1}{n} \sum_{i=1}^n \left(\frac{\phi(\|\theta_{r,i}^\ell\|)}{\sqrt{v_r^t}\|p_{r,i}^\ell\|} - \frac{\phi(\|\theta_{r-1,i}^\ell\|)}{\sqrt{v_{r-1}^t}\|p_{r-1,i}^\ell\|} \right) m_{r-1}^t + \frac{\alpha}{n} \sum_{i=1}^n \frac{\phi(\|\theta_{r,i}^\ell\|)}{\sqrt{v_r^t}\|p_{r,i}^\ell\|} g_{r,i} \right\|^2 \quad (36)$$

$$\leq 2\frac{\alpha^2}{n^2} \left\| \frac{\beta_1}{1-\beta_1} \sum_{i=1}^n \left(\frac{\phi(\|\theta_{r,i}^\ell\|)}{\sqrt{v_r^t}\|p_{r,i}^\ell\|} - \frac{\phi(\|\theta_{r-1,i}^\ell\|)}{\sqrt{v_{r-1}^t}\|p_{r-1,i}^\ell\|} \right) m_{r-1}^t \right\|^2 + \frac{\alpha^2}{n^2} \left\| \sum_{i=1}^n \frac{\phi(\|\theta_{r,i}^\ell\|)}{\sqrt{v_r^t}\|p_{r,i}^\ell\|} g_{r,i} \right\|^2 \quad (37)$$

Taking the expectation on both sides leads to:

$$\begin{aligned} \mathbb{E}[\|\bar{\vartheta}_{r+1}^\ell - \bar{\vartheta}_r^\ell\|^2] &\leq 2\alpha^2 \mathbb{E} \left[\left\| \frac{\beta_1}{1-\beta_1} \sum_{i=1}^n \left(\frac{\phi(\|\theta_{r,i}^\ell\|)}{\sqrt{v_r^t}\|p_{r,i}^\ell\|} - \frac{\phi(\|\theta_{r-1,i}^\ell\|)}{\sqrt{v_{r-1}^t}\|p_{r-1,i}^\ell\|} \right) m_{r-1}^t \right\|^2 \right] + \frac{\alpha^2}{n^2} \mathbb{E} \left[\left\| \sum_{i=1}^n \frac{\phi(\|\theta_{r,i}^\ell\|)}{\sqrt{v_r^t}\|p_{r,i}^\ell\|} g_{r,i} \right\|^2 \right] \\ &\leq 2\frac{\alpha^2}{n^2} \mathbb{E} \left[\left\| \frac{\beta_1}{1-\beta_1} \sum_{i=1}^n \left(\frac{\phi(\|\theta_{r,i}^\ell\|)}{\sqrt{v_r^t}\|p_{r,i}^\ell\|} - \frac{\phi(\|\theta_{r-1,i}^\ell\|)}{\sqrt{v_{r-1}^t}\|p_{r-1,i}^\ell\|} \right) m_{r-1}^t \right\|^2 \right] \\ &\quad + \frac{\alpha^2}{n^2} \mathbb{E} \left[\left\| \sum_{i=1}^n \sum_{j=1}^p \left\langle \Gamma_{r,i}^j (\nabla f_i(\theta_r)^j + g_{r,i}^j - \nabla f_i(\theta_r)^j) \mid \Gamma_{r,i}^j (\nabla f_i(\theta_r)^j + g_{r,i}^j - \nabla f_i(\theta_r)^j) \right\rangle \right\|^2 \right] \\ &\leq 2\frac{\alpha^2}{n^2} \mathbb{E} \left[\left\| \frac{\beta_1}{1-\beta_1} \sum_{i=1}^n \left(\frac{\phi(\|\theta_{r,i}^\ell\|)}{\sqrt{v_r^t}\|p_{r,i}^\ell\|} - \frac{\phi(\|\theta_{r-1,i}^\ell\|)}{\sqrt{v_{r-1}^t}\|p_{r-1,i}^\ell\|} \right) m_{r-1}^t \right\|^2 \right] \\ &\quad + \frac{\alpha^2}{n^2} \mathbb{E} \left[\left\| \sum_{i=1}^n \frac{\phi(\|\theta_{r,i}^\ell\|)}{\sqrt{v_r^t}\|p_{r,i}^\ell\|} \nabla f_i(\theta_r) \right\|^2 \right] + \frac{\alpha^2}{n^2} \sum_{i=1}^n \sigma_i^2 \mathbb{E} \left[\left\| \frac{\phi(\|\theta_{r,i}^\ell\|)}{\sqrt{v_r^t}\|p_{r,i}^\ell\|} \right\|^2 \right] \end{aligned} \quad (38)$$

where the last line uses assumptions H2 and H3 (unbiased gradient and bounded variance of the stochastic gradient) and

$$\Gamma := \frac{\phi(\|\theta_{r,i}^\ell\|)}{\sqrt{v_r^t}\|p_{r,i}^\ell\|}.$$

On the other hand, using the bound on the gradient H2,

$$\begin{aligned} &\sum_{r=1}^R \mathbb{E} \left[\left\| \frac{\beta_1}{1-\beta_1} \sum_{i=1}^n \left(\frac{\phi(\|\theta_{r,i}^\ell\|)}{\sqrt{v_r^t}\|p_{r,i}^\ell\|} - \frac{\phi(\|\theta_{r-1,i}^\ell\|)}{\sqrt{v_{r-1}^t}\|p_{r-1,i}^\ell\|} \right) m_{r-1}^t \right\|^2 \right] \\ &\leq \frac{\beta_1^2}{(1-\beta_1)^2} M^2 \phi_M^2 \sum_{r=1}^R \mathbb{E} \left[\left\| \sum_{i=1}^n \left(\frac{1}{\sqrt{v_r^t}\|p_{r,i}^\ell\|} - \frac{1}{\sqrt{v_{r-1}^t}\|p_{r-1,i}^\ell\|} \right) \right\|^2 \right] \\ &\leq \frac{\beta_1^2}{(1-\beta_1)^2} \frac{L(1-\beta_2)}{p} M^2 \phi_M^2 \sum_{r=1}^R \mathbb{E} \left[\left\| \sum_{i=1}^n \left(\frac{1}{\sqrt{v_r^t}} - \frac{1}{\sqrt{v_{r-1}^t}} \right) \right\|^2 \right] \\ &\leq \frac{\beta_1^2}{(1-\beta_1)^2} \frac{L(1-\beta_2)}{p} M^2 \phi_M^2 \sum_{r=1}^R \mathbb{E} \left[\left\| \sum_{i=1}^n \sum_{j=1}^p \left(\frac{1}{\sqrt{v_r^{t,j}}} - \frac{1}{\sqrt{v_{r-1}^{t,j}}} \right) \right\|^2 \right] \\ &\leq \frac{\beta_1^2}{(1-\beta_1)^2} \frac{L(1-\beta_2)}{p} M^2 \phi_M^2 \frac{np}{v_0} \end{aligned} \quad (39)$$

where, in the telescopic sum, we have used the initial value v_0 of the non decreasing sequence $\{v_r^t\}_{r>0}$ by construction (max operator).

Combining (39) into (38) and summing over the total number of rounds R yields

$$\begin{aligned} \sum_{r=1}^R A_r^2 &:= \sum_{r=1}^R \mathbb{E}[\|\bar{\vartheta}_{r+1}^\ell - \bar{\vartheta}_r^\ell\|^2] \leq \frac{\beta_1^2}{(1-\beta_1)^2} \frac{\mathbb{L}(1-\beta_2)}{p} M^2 \phi_M^2 \frac{np}{v_0} \\ &\quad + \sum_{r=1}^R \left[\frac{\alpha^2}{n^2} \mathbb{E} \left[\left\| \sum_{i=1}^n \frac{\phi(\|\theta_{r,i}^\ell\|)}{\sqrt{v_r^t} \|p_{r,i}^\ell\|} \nabla f_i(\theta_r) \right\|^2 \right] + \frac{\alpha^2}{n^2} \sum_{i=1}^n \sigma_i^2 \mathbb{E} \left[\left\| \frac{\phi(\|\theta_{r,i}^\ell\|)}{\sqrt{v_r^t} \|p_{r,i}^\ell\|} \right\|^2 \right] \right] \end{aligned} \quad (40)$$

Term A_r^3 : According to similar arguments on the non decreasing sequence involved in the algorithm as in the previous series of calculations, observe that

$$\sum_{r=1}^R A_r^3 \leq \frac{\alpha\beta_1}{1-\beta_1} \sqrt{(1-\beta_2)p} \frac{\mathbb{L}M^2}{\sqrt{v_0}} \quad (41)$$

Plugging (33) into (20) combined with (40) and (41) injected into the original smoothness definition (15) summed over the total number of rounds:

$$-\sum_{r=1}^R \mathbb{E}[\langle \nabla f(\bar{\vartheta}_r) | \bar{\vartheta}_{r+1} - \bar{\vartheta}_r \rangle] \leq \sum_{r=1}^R \mathbb{E}[f(\bar{\vartheta}_r) - f(\bar{\vartheta}_{r+1})] + \sum_{r=1}^R \sum_{\ell=1}^L \frac{L_\ell}{2} \mathbb{E}[\|\bar{\vartheta}_{r+1}^\ell - \bar{\vartheta}_r^\ell\|^2] \quad (42)$$

gives:

$$\begin{aligned} &\sum_{r=1}^R \alpha \phi_m \sqrt{\frac{\mathbb{L}(1-\beta_2)}{M^2 p}} \mathbb{E}[\|\nabla f(\bar{\vartheta}_r)\|^2] - \alpha \phi_M \sum_{\ell=1}^L \sum_{i=1}^n \sum_{j=1}^{p_\ell} \frac{\sigma_i^{\ell,j}}{\sqrt{n}} + \frac{\alpha\beta_1}{1-\beta_1} \sqrt{(1-\beta_2)p} \frac{\mathbb{L}M^2}{\sqrt{v_0}} \\ &\leq \sum_{r=1}^R \mathbb{E}[f(\bar{\vartheta}_r) - f(\bar{\vartheta}_{r+1})] + \sum_{\ell=1}^L \frac{L_\ell}{2} \frac{\beta_1^2}{(1-\beta_1)^2} \frac{\mathbb{L}(1-\beta_2)}{p} M^2 \phi_M^2 \frac{np}{v_0} \\ &\quad - \sum_{r=1}^R \left[\frac{\alpha^2}{n^2} \mathbb{E} \left[\left\| \sum_{i=1}^n \frac{\phi(\|\theta_{r,i}^\ell\|)}{\sqrt{v_r^t} \|p_{r,i}^\ell\|} \nabla f_i(\theta_r) \right\|^2 \right] + \frac{\alpha^2}{n^2} \sum_{i=1}^n \sigma_i^2 \mathbb{E} \left[\left\| \frac{\phi(\|\theta_{r,i}^\ell\|)}{\sqrt{v_r^t} \|p_{r,i}^\ell\|} \right\|^2 \right] \right] \end{aligned} \quad (43)$$

Noting that $\sum_{r=1}^R \mathbb{E}[f(\bar{\vartheta}_r) - f(\bar{\vartheta}_{r+1})] = f(\bar{\vartheta}_1) - \mathbb{E}[f(\bar{\vartheta}_{R+1})]$, we obtain

$$\begin{aligned} &\sum_{r=1}^R \alpha \phi_m \sqrt{\frac{\mathbb{L}(1-\beta_2)}{M^2 p}} \mathbb{E}[\|\nabla f(\bar{\vartheta}_r)\|^2] + \frac{1}{n^2} \sum_{r=1}^R \mathbb{E} \left[\left\| \sum_{i=1}^n \frac{\phi(\|\theta_{r,i}^\ell\|)}{\sqrt{v_r^t} \|p_{r,i}^\ell\|} \nabla f_i(\theta_r) \right\|^2 \right] \\ &\leq f(\bar{\vartheta}_1) - \mathbb{E}[f(\bar{\vartheta}_{R+1})] + \frac{1}{n^2} \sum_{i=1}^n \sigma_i^2 \mathbb{E} \left[\left\| \frac{\phi(\|\theta_{r,i}^\ell\|)}{\sqrt{v_r^t} \|p_{r,i}^\ell\|} \right\|^2 \right] + \alpha \phi_M \sum_{\ell=1}^L \sum_{i=1}^n \sum_{j=1}^{p_\ell} \frac{\sigma_i^{\ell,j}}{\sqrt{n}} + \frac{\alpha\beta_1}{1-\beta_1} \sqrt{(1-\beta_2)p} \frac{\mathbb{L}M^2}{\sqrt{v_0}} \\ &\quad + \sum_{\ell=1}^L \frac{L_\ell}{2} \frac{\beta_1^2}{(1-\beta_1)^2} \frac{\mathbb{L}(1-\beta_2)}{p} M^2 \phi_M^2 \frac{np}{v_0} \end{aligned} \quad (44)$$

leading to

$$\begin{aligned} \sum_{r=1}^R \mathbb{E} \left[\left\| \sum_{i=1}^n \frac{\phi(\|\theta_{r,i}^\ell\|)}{\sqrt{v_r^t} \|p_{r,i}^\ell\|} \nabla f_i(\theta_r) \right\|^2 \right] &\leq \frac{1}{\alpha} [f(\bar{\vartheta}_1) - \mathbb{E}[f(\bar{\vartheta}_{R+1})]] + \sum_{r=1}^R \frac{\alpha}{n^2} \sum_{i=1}^n \sigma_i^2 \mathbb{E} \left[\left\| \frac{\phi(\|\theta_{r,i}^\ell\|)}{\sqrt{v_r^t} \|p_{r,i}^\ell\|} \right\|^2 \right] \\ &\quad + \alpha \phi_M \sigma \mathbb{L} p \sqrt{n} + \frac{\bar{L} \beta_1^2 \mathbb{L} (1 - \beta_2) M^2 \phi_M^2 n}{2(1 - \beta_1)^2 v_0} + \frac{\alpha \beta_1}{1 - \beta_1} \sqrt{(1 - \beta_2) p} \frac{\mathbb{L} M^2}{\sqrt{v_0}} \end{aligned} \quad (45)$$

where $\bar{L} = \sum_{\ell=1}^L L_\ell$ is the sum of all smoothness constants.

Consider the following inequality:

$$\frac{1}{n} \sum_{i=1}^n \frac{\phi(\|\theta_{r,i}^\ell\|)}{\sqrt{v_r^t} \|p_{r,i}^\ell\|} \nabla f_i(\theta_r) \leq \phi_M (1 - \beta_2) \frac{\bar{\nabla} f(\theta_r)}{\sqrt{v_r^t}} \quad (46)$$

where $\bar{\nabla} f(\theta_r) := \frac{1}{n} \sum_{i=1}^n \nabla f_i(\theta_r)$. And using the Cauchy-Schwartz inequality we have

$$\left\| \frac{\bar{\nabla} f(\theta_r)}{\sqrt{v_r^t}} \right\| \geq \frac{1}{2} \left\| \frac{\nabla f(\bar{\theta}_r)}{\sqrt{v_r^t}} \right\| - \left\| \frac{\bar{\nabla} f(\theta_r) - \nabla f(\bar{\theta}_r)}{\sqrt{v_r^t}} \right\| \quad (47)$$

Using Lemma 1 and the smoothness assumption H1, we have

$$\begin{aligned} \left\| \frac{\bar{\nabla} f(\theta_r)}{\sqrt{v_r^t}} \right\| &\geq \frac{1}{2} \left\| \frac{\nabla f(\bar{\theta}_r)}{\sqrt{v_r^t}} \right\| - \left\| \frac{\bar{\nabla} f(\theta_r) - \nabla f(\bar{\theta}_r)}{\sqrt{v_r^t}} \right\| \\ &\geq \frac{1}{2} \left\| \frac{\nabla f(\bar{\theta}_r)}{\sqrt{v_r^t}} \right\| - \bar{L} \alpha^2 M^2 \phi_M^2 \frac{(1 - \beta_2) p}{v_0} \end{aligned} \quad (48)$$

which is due to Lemma 2. Plugging the above inequality into (45) and dividing both sides by R yields:

$$\begin{aligned} \frac{1}{R} \sum_{r=1}^R \mathbb{E} \left[\left\| \frac{\nabla f(\bar{\theta}_r)}{\sqrt{v_r^t}} \right\|^2 \right] &\leq \sqrt{\frac{M^2 p}{n}} \frac{f(\bar{\vartheta}_1) - \mathbb{E}[f(\bar{\vartheta}_{R+1})]}{\mathbb{L} \alpha R} + \frac{\alpha}{n^2} \sum_{r=1}^R \sum_{i=1}^n \sigma_i^2 \mathbb{E} \left[\left\| \frac{\phi(\|\theta_{r,i}^\ell\|)}{\sqrt{v_r^t} \|p_{r,i}^\ell\|} \right\|^2 \right] \\ &\quad + \alpha \phi_M \sigma \mathbb{L} p \sqrt{n} + \frac{\bar{L} \beta_1^2 \mathbb{L} (1 - \beta_2) M^2 \phi_M^2 n}{2(1 - \beta_1)^2 v_0} + \frac{\alpha \beta_1}{1 - \beta_1} \sqrt{(1 - \beta_2) p} \frac{\mathbb{L} M^2}{\sqrt{v_0}} + \bar{L} \alpha^2 M^2 \phi_M^2 \frac{(1 - \beta_2) p}{R v_0} \end{aligned} \quad (49)$$

concluding the proof of our main convergence result.

A.3. Proof Corollary 1

Corollary. Assume **H1-H4**. Consider $\{\bar{\theta}_r\}_{r>0}$, the sequence of parameters obtained running Algorithm 2. Then, if the number of local epochs is set to $T = 1$, $\epsilon = \lambda = 0$, and $R \geq$ we have, under **H5**:

$$\frac{1}{R} \mathbb{E} \left[\|\nabla f(\bar{\theta}_R)\|^2 \right] \leq \mathcal{O} \left(\sqrt{\frac{p}{n}} \frac{1}{\mathbb{L} \sqrt{R}} \right) \quad (50)$$

Proof. From the bound in Theorem 1 and with assumption **H5**. □

A.4. Multiple local updates

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.*, $\left\|\frac{\nabla f(\bar{\theta}_r)}{\sqrt{v_r^t}}\right\|$ and $\left\|\frac{\nabla f(\bar{\theta}_r)}{\sqrt{v_r^t}}\right\|$), are impacted.

In previous proof, we change all index r to iteration t . Suppose T is the number of local iterations. We can write (23) as

$$\alpha \langle \nabla f(\bar{\vartheta}_t), \frac{\bar{g}_t}{\sqrt{\hat{v}_t}} \rangle, \quad (51)$$

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

$$\langle \nabla f(\bar{\vartheta}_t), \frac{\bar{g}_t}{\sqrt{\hat{v}_t}} \rangle = \frac{1}{2} \left\| \frac{\nabla f(\bar{\vartheta}_t)}{\hat{v}_t^{1/4}} \right\|^2 + \frac{1}{2} \left\| \frac{\bar{g}_t}{\hat{v}_t^{1/4}} \right\|^2 - \left\| \frac{\nabla f(\bar{\vartheta}_t) - \bar{g}_t}{\hat{v}_t^{1/4}} \right\|^2. \quad (52)$$

To bound the last term on the RHS, we have

$$\left\| \frac{\nabla f(\bar{\vartheta}_t) - \bar{g}_t}{\hat{v}_t^{1/4}} \right\|^2 = \left\| \frac{\frac{1}{n} \sum_{i=1}^n (\nabla f_i(\bar{\vartheta}_t) - \bar{g}_{t,i})}{\hat{v}_t^{1/4}} \right\|^2 \quad (53)$$

$$\leq \frac{1}{n} \sum_{i=1}^n \left\| \frac{\nabla f_i(\bar{\vartheta}_t) - \bar{g}_{t,i}}{\hat{v}_t^{1/4}} \right\|^2 \quad (54)$$

$$\leq \frac{2}{n} \sum_{i=1}^n \left(\left\| \frac{\nabla f_i(\bar{\vartheta}_t) - \nabla f_i(\bar{\theta}_t)}{\hat{v}_t^{1/4}} \right\|^2 + \left\| \frac{\nabla f_i(\bar{\theta}_t) - \bar{g}_{t,i}}{\hat{v}_t^{1/4}} \right\|^2 \right). \quad (55)$$

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

$$\frac{2}{n} \sum_{i=1}^n \left\| \frac{\nabla f_i(\bar{\vartheta}_t) - \nabla f_i(\bar{\theta}_t)}{\hat{v}_t^{1/4}} \right\|^2 \leq \frac{2}{n\sqrt{v_0}} \sum_{i=1}^n L_\ell \|\bar{\vartheta}_t - \bar{\theta}_t\|^2 \quad (56)$$

$$= \frac{2L_\ell}{n\sqrt{v_0}} \frac{\beta_1^2}{(1-\beta_1)^2} \sum_{i=1}^n \|\bar{\theta}_t - \bar{\theta}_{t-1}\|^2 \quad (57)$$

$$\leq \frac{2\alpha^2 L_\ell}{n\sqrt{v_0}} \frac{\beta_1^2}{(1-\beta_1)^2} \sum_{i=1}^n \sum_{l=1}^L \left\| \frac{\phi(\|\theta_{t,i}^l\|)}{\|p_{t,i}^l\|} p_{t,i}^l \right\|^2 \quad (58)$$

$$\leq \frac{2\alpha^2 L_\ell p \phi_M^2}{\sqrt{v_0}} \frac{\beta_1^2}{(1-\beta_1)^2}. \quad (59)$$

For the second term,

$$\frac{2}{n} \sum_{i=1}^n \left\| \frac{\nabla f_i(\bar{\theta}_t) - \bar{g}_{t,i}}{\hat{v}_t^{1/4}} \right\|^2 \leq \frac{4}{n} \left(\underbrace{\sum_{i=1}^n \left\| \frac{\nabla f_i(\bar{\theta}_t) - \nabla f_i(\theta_{t,i})}{\hat{v}_t^{1/4}} \right\|^2}_{B_1} + \underbrace{\sum_{i=1}^n \left\| \frac{\nabla f_i(\theta_{t,i}) - \bar{g}_{t,i}}{\hat{v}_t^{1/4}} \right\|^2}_{B_2} \right). \quad (60)$$

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

$$B_1 \leq \frac{L_\ell}{\sqrt{v_0}} \sum_{i=1}^n \|\bar{\theta}_t - \theta_{t,i}\|^2 \quad (61)$$

$$= \frac{\alpha^2 L_\ell}{\sqrt{v_0}} \sum_{i=1}^n \sum_{l=1}^L \left\| \sum_{j=\lfloor t \rfloor_T + 1}^t \left(\frac{\phi(\|\theta_{j,i}^l\|)}{\|p_{j,i}^l\|} p_{j,i}^l - \frac{1}{n} \sum_{k=1}^n \frac{\phi(\|\theta_{j,k}^l\|)}{\|p_{j,k}^l\|} p_{j,k}^l \right) \right\|^2 \quad (62)$$

$$\leq n \frac{\alpha^2 L_\ell}{\sqrt{v_0}} \alpha^2 M^2 (T-1)^2 \phi_M^2 (1-\beta_2) p \quad (63)$$

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 $(t - 1) - (\lfloor t \rfloor_T + 1) + 1 \leq T - 1$. The expectation of B_2 relates to the variance and bias of the stochastic gradients, which can be bounded by

$$\mathbb{E}[B_2] \leq n \frac{\sigma^2}{\sqrt{\hat{v}_0}}$$

where we use assumption H3.

Then (60) becomes:

$$\frac{2}{n} \sum_{i=1}^n \left\| \frac{\nabla f_i(\bar{\theta}_t) - \bar{g}_{t,i}}{\hat{v}_t^{1/4}} \right\|^2 \leq 4 \left[\frac{\sigma^2}{\sqrt{\hat{v}_0}} + \frac{\alpha^2 L_\ell}{\sqrt{v_0}} \alpha^2 M^2 (T - 1)^2 \phi_M^2 (1 - \beta_2) p \right] \quad (64)$$