

# Personalized Federated Learning with Moreau Envelopes

Canh T. Dinh<sup>1</sup>, Nguyen H. Tran<sup>1</sup>, Tuan Dung Nguyen<sup>1,2</sup>

<sup>1</sup>The University of Sydney, Australia  
tdin6081@uni.sydney.edu.au, nguyen.tran@sydney.edu.au

<sup>2</sup>The University of Melbourne, Australia  
tuandungn@unimelb.edu.au

## Abstract

Federated learning (FL) is a decentralized and privacy-preserving machine learning technique in which a group of clients collaborate with a server to learn a global model without sharing clients' data. One challenge associated with FL is statistical diversity among clients, which restricts the global model from delivering good performance on each client's task. To address this, we propose an algorithm for personalized FL (pFedMe) using Moreau envelopes as clients' regularized loss functions, which help decouple personalized model optimization from the global model learning in a bi-level problem stylized for personalized FL. Theoretically, we show that pFedMe's convergence rate is state-of-the-art: achieving quadratic speedup for strongly convex and sublinear speedup of order  $2/3$  for smooth non-convex objectives. Experimentally, we verify that pFedMe excels at empirical performance compared with the vanilla FedAvg and Per-FedAvg, a meta-learning based personalized FL algorithm.

## 1 Introduction

The abundance of data generated in a massive number of hand-held devices these days has stimulated the development of Federated learning (FL) [1]. The setting of FL is a network of clients connected to a server, and its goal is to build a global model from clients' data in a privacy-preserving and communication-efficient way. The current techniques that attempt to fulfill this goal mostly follow three steps: (i) at each communication iteration, the server sends the current global model to clients; (ii) the clients update their local models using their local data; (iii) the server collects the latest local models from a subset of sampled clients in order to update a new global model, repeated until convergence [1–4].

Despite its advantages of data privacy and communication reduction, FL faces a main challenge that affects its performance and convergence rate: statistical diversity, which means that data distributions among clients are distinct (i.e., non-i.i.d.). Thus, the global model, which is trained using these non-i.i.d. data, is hardly well-generalized on each client's data. This particular behaviour has been reported in [5, 6], which showed that when the statistical diversity increases, generalization errors of the global model on clients' local data also increase significantly. On the other hand, individual learning without FL (i.e., no client collaboration) will also have large generalization error due to insufficient data. These raise the question: *How can we leverage the global model in FL to find a "personalized model" that is stylized for each client's data?*

Motivated by critical roles of personalized models in several business applications of healthcare, finance, and AI services [5], we address this question by proposing a new FL scheme for personalization, which minimizes the Moreau envelopes [7] of clients' loss functions. With this scheme,

clients not only contribute to building the “reference” global model as in the standard FL, but also leverage the reference model to optimize their personalized models w.r.t. local data. Geometrically, the global model in this scheme can be considered as a “central point” where all clients agree to meet, and personalized models are the points in different directions that clients follow according to their heterogeneous data distributions.

**Our key contributions** in this work are summarized as follows. First, we formulate a new bi-level optimization problem designed for personalized FL (pFedMe) by using the Moreau envelope as a regularized loss function. The bi-level structure of pFedMe has a key advantage: decoupling the process of optimizing personalized models from learning the global model. Thus, pFedMe updates the global model similarly to the standard FL algorithm such as FedAvg [1], yet parallelly optimizes the personalized models with low complexity.

Second, we exploit the convexity-preserving and smoothness-enabled properties of the Moreau envelopes to facilitate the convergence analysis of pFedMe, which characterizes both client-sampling and client-drift errors: two notorious issues in FL [3]. With carefully tuned hyperparameters, pFedMe can obtain the state-of-the-art quadratic speedup (resp. sublinear speedup of order  $2/3$ ), compared with the existing works with linear speedup (resp. sublinear speedup of order  $1/2$ ), for strongly convex (resp. smooth nonconvex) objective.

Finally, we empirically evaluate the performance of pFedMe using both real and synthetic datasets that capture the statistical diversity of clients’ data. We show that pFedMe outperforms the vanilla FedAvg and a meta-learning based personalized FL algorithm Per-FedAvg [8] in terms of convergence rate and local accuracy.

## 2 Related Work

**FL and challenges.** One of the first FL algorithms is FedAvg [1], which uses local SGD updates and builds a global model from a subset of clients with non-i.i.d. data. Subsequently, one-shot FL [9] allows the global model to learn in one single round of communication. To address the limitations on communications in a FL network, [10, 11] introduced quantization methods, while [12–14] proposed performing multiple local optimization rounds before sending the local models to the server. In addition, the problem of statistical diversity has been addressed in [15–20]. Preserving privacy in FL has been studied in [21–25].

**Personalized FL: mixing models, contextualization, meta-learning, and multi-task learning.** Multiple approaches have been proposed to achieve personalization in FL. One such approach is *mixing* the global and local models. [26] combined the optimization of the local and global models in its L2GD algorithm. [27] introduced three personalization approaches to: user clustering, data interpolation, and model interpolation. While the first two approaches need meta-features from all clients that make them not feasible in FL due to privacy concern, the last approach was used in [6] to create an adaptive personalized federated learning (APFL) algorithm, which attempted to mix a user’s local model with the global model. One personalization method used in neural networks is FedPer [28], in which a network is divided into base and personalized layers, and while the base layers are trained by the server, both types of layers will be trained by users to create a personalized model. Regarding using a model in different *contexts*, in the next-character prediction task in [29], the requirement to predict differently among devices raises a need to inspect more features about the context of client devices during training, which was studied in [30]. [31] achieves personalization on each user in a fully decentralized network using asynchronous gossip algorithms with assumptions on network topology and similarity between users. The concept of personalization can also be linked to *meta-learning*. Per-FedAvg [8], influenced by Model-Agnostic Meta-Learning (MAML) [32], built an initial meta-model that can be updated effectively after one more gradient descent step. During meta-optimization, however, MAML theoretically requires computing the Hessian term, which is computationally prohibitive; therefore, several works including [32–34] attempted to approximate the Hessian matrix. [35] based its framework, ARUBA, on online convex optimization and meta-learning, which can be integrated into FL to improve personalization. [36] discovered that FedAvg can be interpreted as meta-learning and proposed combining FedAvg with Reptile [33] for FL personalization. The application of federated meta-learning in recommender systems was studied in [37]. Finally, *multi-task learning* can be used for personalization: [20] introduced a federated multi-task framework

called MOCHA, addressing both systems and statistical heterogeneity. For more details about FL, its challenges, and personalization approaches, we refer the readers to comprehensive surveys in [38, 39].

### 3 Personalized Federated Learning with Moreau Envelopes (pFedMe)

#### 3.1 pFedMe: Problem Formulation

In conventional FL, there are  $N$  clients communicating with a server to solve the following problem:

$$\min_{w \in \mathbb{R}^d} \left\{ f(w) := \frac{1}{N} \sum_{i=1}^N f_i(w) \right\} \quad (1)$$

to find a *global model*  $w$ . The function  $f_i : \mathbb{R}^d \rightarrow \mathbb{R}$ ,  $i = 1, \dots, N$ , denotes the expected loss over the data distribution of the client  $i$ :

$$f_i(w) = \mathbb{E}_{\xi_i} [\tilde{f}_i(w; \xi_i)],$$

where  $\xi_i$  is a random data sample drawn according to the distribution of client  $i$  and  $\tilde{f}_i(w; \xi_i)$  is a loss function corresponding to this sample and  $w$ . In FL, since clients' data possibly come from different environments, contexts, and applications, clients can have non-i.i.d. data distributions, i.e., the distributions of  $\xi_i$  and  $\xi_j$ ,  $i \neq j$ , are distinct.

Instead of solving the traditional FL problem (1), we take a different approach by using a regularized loss function with  $l_2$ -norm for each client as follows

$$f_i(\theta_i) + \frac{\lambda}{2} \|\theta_i - w\|^2, \quad (2)$$

where  $\theta_i$  denotes the *personalized model* of client  $i$  and  $\lambda$  is a regularization parameter that controls the strength of  $w$  to the personalized model. While large  $\lambda$  can benefit clients with unreliable data from the abundant data aggregation, small  $\lambda$  helps clients with sufficient useful data prioritize personalization. Note that  $\lambda \in (0, \infty)$  to avoid extreme cases of  $\lambda = 0$ , i.e., no FL, or  $\lambda = \infty$ , i.e., no personalized FL. Overall, the idea is allowing clients to pursue their own models with different directions, but not to stay far away from the “reference point”  $w$ , to which every client contributes. Based on this, the personalized FL can be formulated as a bi-level problem:

$$\text{pFedMe} : \min_{w \in \mathbb{R}^d} \left\{ F(w) := \frac{1}{N} \sum_{i=1}^N F_i(w) \right\}, \text{ where } F_i(w) = \min_{\theta_i \in \mathbb{R}^d} \left\{ f_i(\theta_i) + \frac{\lambda}{2} \|\theta_i - w\|^2 \right\}.$$

In pFedMe, while  $w$  is found by exploiting the data aggregation from multiple clients at the outer level,  $\theta_i$  is optimized with respect to (w.r.t) client  $i$ 's data distribution and is maintained a bounded distance from  $w$  at the inner level. The definition of  $F_i(w)$  is the well-known Moreau envelope, which facilitates several learning algorithm designs [40, 41]. The optimal personalized model, which is the unique solution to the inner problem of pFedMe and also known as the proximal operator in the literature, is defined as follows:

$$\hat{\theta}_i(w) := \text{prox}_{f_i/\lambda}(w) = \arg \min_{\theta_i \in \mathbb{R}^d} \left\{ f_i(\theta_i) + \frac{\lambda}{2} \|\theta_i - w\|^2 \right\}. \quad (3)$$

For comparison, we consider Per-FedAvg [8], which arguably has the closest formulation to pFedMe:

$$\min_{w \in \mathbb{R}^d} \left\{ F(w) := \frac{1}{N} \sum_{i=1}^N f_i(\theta_i(w)) \right\}, \text{ where } \theta_i(w) = w - \alpha \nabla f_i(w). \quad (4)$$

Based on the MAML framework [32], Per-FedAvg aims to find a global model  $w$  which client  $i$  can use as an *initialization* to perform one more step of gradient update (with step size  $\alpha$ ) w.r.t its own loss function to obtain its personalized model  $\theta_i(w)$ .

Compared to Per-FedAvg, our problem has a similar meaning of  $w$  as a “meta-model”, but instead of using  $w$  as the initialization, we parallelly pursue both the personalized and global models by solving a bi-level problem, which has several benefits. First, while Per-FedAvg is optimized for one-step gradient update for its personalized model, pFedMe is agnostic to the inner optimizer, which means (3) can be solved using any iterative approach with multi-step updates. Second, by re-writing the personalized model update of Per-FedAvg as

$$\theta_i(w) = w - \alpha \nabla f_i(w) = \arg \min_{\theta_i \in \mathbb{R}^d} \left\{ \langle \nabla f_i(w), \theta_i - w \rangle + \frac{1}{2\alpha} \|\theta_i - w\|^2 \right\}, \quad (5)$$

where we use  $\langle x, y \rangle$  for the inner product of two vectors  $x$  and  $y$ , we can see that apart from the similar regularization term, Per-FedAvg only optimizes the first-order approximation of  $f_i$ , whereas pFedMe directly minimizes  $f_i$  in (3). Third, Per-FedAvg (or generally several MAML-based methods) requires computing or estimating Hessian matrix, whereas pFedMe only needs gradient calculation using first-order approach, as will be shown in the next section.

**Assumption 1** (Strong convexity and smoothness).  $f_i$  is either (a)  $\mu$ -strongly convex or (b) nonconvex and  $L$ -smooth (i.e.,  $L$ -Lipschitz gradient), respectively, as follows when  $\forall w, w'$ :

$$\begin{aligned} (a) \quad & f_i(w) \geq f_i(w') + \langle \nabla f_i(w'), w - w' \rangle + \frac{\mu}{2} \|w - w'\|^2, \\ (b) \quad & \|\nabla f_i(w) - \nabla f_i(w')\| \leq L \|w - w'\|. \end{aligned}$$

**Assumption 2** (Bounded variance). The variance of stochastic gradients in each client is bounded

$$\mathbb{E}_{\xi_i} [\|\nabla \tilde{f}_i(w; \xi_i) - \nabla f_i(w)\|^2] \leq \gamma_f^2.$$

**Assumption 3** (Bounded diversity). The variance of local gradients to global gradient is bounded

$$\frac{1}{N} \sum_{i=1}^N \|\nabla f_i(w) - \nabla f(w)\|^2 \leq \sigma_f^2.$$

While Assumption 1 is standard for convergence analysis, Assumptions 2 and 3 are widely used in FL context in which  $\gamma_f^2$  and  $\sigma_f^2$  quantify the sampling noise and the diversity of client's data distribution, respectively [3, 8, 42, 43]. Note that we avoid using the uniformly bounded gradient assumption, i.e.,  $\|\nabla f_i(w)\| \leq G, \forall i$ , which was used in several related works [6, 8]. It was shown that this assumption is not satisfied in the unconstrained strongly convex minimization [44, 45].

Finally, we review several useful properties of the Moreau envelope such as smoothing and preserving convexity as follows (see the review and proof for the convex case in [40, 46, 47] and for nonconvex smooth case in [48], respectively):

**Proposition 1.** If  $f_i$  is convex or nonconvex with  $L$ -Lipschitz  $\nabla f_i$ , then  $\nabla F_i$  is  $L_F$ -smooth with  $L_F = \lambda$  (with the condition that  $\lambda > 2L$  for nonconvex  $L$ -smooth  $f_i$ ), and

$$\nabla F_i(w) = \lambda(w - \hat{\theta}_i(w)). \quad (6)$$

Furthermore, if  $f_i$  is  $\mu$ -strongly convex, then  $F_i$  is  $\mu_F$ -strongly convex with  $\mu_F = \frac{\lambda\mu}{\lambda+\mu}$ .

### 3.2 pFedMe: Algorithm

In this section, we propose an algorithm, presented in Alg. 1, to solve pFedMe. Similar to conventional FL algorithms such as FedAvg [1], at each communication round  $t$ , the server broadcasts the latest global model  $w_t$  to all clients. Then, after all clients perform  $R$  local updates, the server will receive the latest local models from a uniformly sampled subset  $\mathcal{S}^t$  of clients to perform the model averaging. Note that we use an additional parameter  $\beta$  for global model update, which includes FedAvg's model averaging when  $\beta = 1$ . Though a similar parameter at the server side was also used in [3, 49], it will be shown that pFedMe can obtain better speedup convergence rates.

Specifically, our algorithm, which aims to solve the bi-level problem pFedMe, has two key differences compared with FedAvg, which aims to solve (1). First, at the inner level, each client  $i$  solves (3) to obtain its personalized model  $\hat{\theta}_i(w_{i,r}^t)$  where  $w_{i,r}^t$  denotes the *local model* of the client  $i$  at the global round  $t$  and local round  $r$ . Similar to FedAvg, the purpose of local models is to contribute to building global model with reduced communication rounds between clients and server. Second, at the outer level, the local update of client  $i$  using gradient descent is with respect to  $F_i$  (instead of  $f_i$ ) as the following

$$w_{i,r+1}^t = w_{i,r}^t - \eta \nabla F_i(w_{i,r}^t),$$

where  $\eta$  is the learning rate and  $\nabla F_i(w_{i,r}^t)$  is calculated according to (6) using the current personalized model  $\hat{\theta}_i(w_{i,r}^t)$ .

For the practical algorithm, we use a  $\delta$ -approximation of  $\hat{\theta}_i(w_{i,r}^t)$ , denoted by  $\tilde{\theta}_i(w_{i,r}^t)$  satisfying  $\mathbb{E}[\|\tilde{\theta}_i(w_{i,r}^t) - \hat{\theta}_i(w_{i,r}^t)\|] \leq \delta$ , and correspondingly use  $\lambda(w_{i,r}^t - \tilde{\theta}_i(w_{i,r}^t))$  to approximate  $\nabla F_i(w_{i,r}^t)$

---

**Algorithm 1** pFedMe: Personalized Federated Learning using Moreau Envelope Algorithm

---

```

1: input:  $T, R, S, \lambda, \eta, \beta, w^0$ 
2: for  $t = 0$  to  $T - 1$  do                                     ▷ Global communication rounds
3:   Server sends  $w_t$  to all clients
4:   for all  $i = 1$  to  $N$  do
5:      $w_{i,0}^t = w_t$ 
6:     for  $r = 0$  to  $R - 1$  do                                     ▷ Local update rounds
7:       Sample a fresh mini-batch  $\mathcal{D}_i$  with size  $|\mathcal{D}|$  and minimize  $\tilde{h}_i(\theta_i; w_{i,r}^t, \mathcal{D}_i)$ , defined in
          (7), up to an accuracy level according to (8) to find a  $\delta$ -approximate  $\tilde{\theta}_i(w_{i,r}^t)$ 
8:        $w_{i,r+1}^t = w_{i,r}^t - \eta \lambda (w_{i,r}^t - \tilde{\theta}_i(w_{i,r}^t))$ 
9:       Server uniformly samples a subset of clients  $\mathcal{S}^t$  with size  $S$ , and each of the sampled client
          sends the local model  $w_{i,R}^t, \forall i \in \mathcal{S}^t$ , to the server
10:    Server updates the global model:  $w_{t+1} = (1 - \beta)w_t + \beta \sum_{i \in \mathcal{S}^t} \frac{w_{i,R}^t}{S}$ 

```

---

(c.f. line 8). The reason of using the  $\delta$ -approximate  $\tilde{\theta}_i(w_{i,r}^t)$  is two-fold. First, obtaining  $\hat{\theta}_i(w_{i,r}^t)$  according to (3) usually needs the gradient  $\nabla f_i(\theta_i)$ , which, however, requires the distribution of  $\xi_i$ . In practice, we use the following unbiased estimate of  $\nabla f_i(\theta_i)$  by sampling a mini-batch of data  $\mathcal{D}_i$

$$\nabla \tilde{f}_i(\theta_i, \mathcal{D}_i) := \frac{1}{|\mathcal{D}_i|} \sum_{\xi_i \in \mathcal{D}_i} \nabla \tilde{f}_i(\theta_i; \xi_i)$$

such that  $\mathbb{E}[\nabla \tilde{f}_i(\theta_i, \mathcal{D}_i)] = \nabla f_i(\theta_i)$ . Second, in general, it is not straightforward to obtain  $\hat{\theta}_i(w_{i,r}^t)$  in closed-form. Instead we usually use iterative first-order approach to obtain an approximate  $\tilde{\theta}_i(w_{i,r}^t)$  with high accuracy. Defining

$$\tilde{h}_i(\theta_i; w_{i,r}^t, \mathcal{D}_i) := \tilde{f}_i(\theta_i; \mathcal{D}_i) + \frac{\lambda}{2} \|\theta_i - w_{i,r}^t\|^2, \quad (7)$$

suppose we choose  $\lambda$  such that  $\tilde{h}_i(\theta_i; w_{i,r}^t, \mathcal{D}_i)$  is strongly convex with a condition number  $\kappa$  (which quantifies how hard to optimize (7)), then we can apply gradient descent (resp. Nesterov's accelerated gradient descent) to obtain  $\tilde{\theta}_i(w_{i,r}^t)$  such that

$$\|\nabla \tilde{h}_i(\tilde{\theta}_i; w_{i,r}^t, \mathcal{D}_i)\| \leq \nu, \quad (8)$$

with the number of  $\nabla \tilde{h}_i$  computations  $K := \mathcal{O}(\kappa \log(\frac{d}{\nu}))$  (resp.  $\mathcal{O}(\sqrt{\kappa} \log(\frac{d}{\nu}))$ ) [50], where  $d$  is the diameter of the search space,  $\nu$  is an accuracy level, and  $\mathcal{O}(\cdot)$  hides constants. The computation complexity of each client in pFedMe is  $K$  times that in FedAvg. In the following lemma, we show how  $\delta$  can be adjusted by controlling the (i) sampling noise using mini-batch size  $|\mathcal{D}|$  and (ii) accuracy level  $\nu$ .

**Lemma 1.** *Let  $\tilde{\theta}_i(w_{i,r}^t)$  be a solution to (8), we have*

$$\mathbb{E} \left[ \|\tilde{\theta}_i(w_{i,r}^t) - \hat{\theta}_i(w_{i,r}^t)\|^2 \right] \leq \delta^2 := \begin{cases} \frac{2}{(\lambda + \mu)^2} \left( \frac{\gamma_f^2}{|\mathcal{D}|} + \nu \right), & \text{if Assumption 1(a) holds;} \\ \frac{2}{(\lambda - L)^2} \left( \frac{\gamma_f^2}{|\mathcal{D}|} + \nu \right), & \text{if Assumption 1(b) holds, and } \lambda > L. \end{cases}$$

## 4 pFedMe: Convergence Analysis

In this section, we present the convergence of pFedMe. We first prove an intermediate result.

**Lemma 2.** *Recall the definition of the Moreau envelope  $F_i$  in pFedMe.*

(a) *Let Assumption 1(a) hold, then we have*

$$\frac{1}{N} \sum_{i=1}^N \|\nabla F_i(w) - \nabla F(w)\|^2 \leq 4L_F(F(w) - F(w^*)) + 2 \underbrace{\frac{1}{N} \sum_{i=1}^N \|\nabla F_i(w^*)\|^2}_{=:\sigma_{F,1}^2}.$$

(b) If Assumption 1(b) holds and additionally  $\lambda > 2\sqrt{2}L$ , we have

$$\frac{1}{N} \sum_{i=1}^N \|\nabla F_i(w) - \nabla F(w)\|^2 \leq \frac{8L^2}{\lambda^2 - 8L^2} \|\nabla F(w)\|^2 + 2 \underbrace{\frac{\lambda^2}{\lambda^2 - 8L^2} \sigma_f^2}_{=:\sigma_{F,2}^2}.$$

This lemma provides the bounded diversity of  $F_i$ , characterized by the variances  $\sigma_{F,1}^2$  and  $\sigma_{F,2}^2$ , for strongly convex and nonconvex smooth  $f_i$ , respectively. While  $\sigma_{F,2}^2$  is related to  $\sigma_f^2$  that needs to be bounded in Assumption 3,  $\sigma_{F,1}^2$  is measured only at the unique solution  $w^*$  to pFedMe (for strongly convex  $F_i$ ,  $w^*$  always exists), and thus  $\sigma_{F,1}^2$  is finite. These bounds are tight in the sense that  $\sigma_{F,1}^2 = \sigma_{F,2}^2 = 0$  when data distribution of clients are i.i.d.

**Theorem 1** (Strongly convex pFedMe's convergence). *Let Assumptions 1(a) and 2 hold. If  $T \geq \frac{2}{\eta_1 \mu_F}$ , there exists an  $\eta \leq \frac{\hat{\eta}_1}{\beta R}$ , where  $\hat{\eta}_1 := \frac{1}{6L_F(3+128\kappa_F/\beta)}$  with  $\beta \geq 1$ , such that*

$$\begin{aligned} (a) \mathbb{E} [F(\bar{w}^T) - F(w^*)] &\leq \mathcal{O}(\mathbb{E} [F(\bar{w}^T) - F(w^*)]) := \\ &\mathcal{O}(\Delta_0 \mu_F e^{-\hat{\eta}_1 \mu_F T/2}) + \tilde{\mathcal{O}}\left(\frac{(N/S-1)\sigma_{F,1}^2}{\mu_F T N}\right) + \tilde{\mathcal{O}}\left(\frac{(R\sigma_{F,1}^2 + \delta^2 \lambda^2) \kappa_F}{R(T\beta \mu_F)^2}\right) + \mathcal{O}\left(\frac{\lambda^2 \delta^2}{\mu_F}\right) \\ (b) \frac{1}{N} \sum_{i=1}^N \mathbb{E} [\|\tilde{\theta}_i^T(w_T) - w^*\|^2] &\leq \frac{1}{\mu_F} \mathcal{O}(\mathbb{E} [F(\bar{w}_T) - F^*]) + \mathcal{O}\left(\frac{\sigma_{F,1}^2}{\lambda^2} + \delta^2\right), \end{aligned}$$

where  $\Delta_0 := \|w_0 - w^*\|^2$ ,  $\kappa_F := \frac{L_F}{\mu_F}$ ,  $\bar{w}_T := \sum_{t=0}^{T-1} \alpha_t w_t / A_T$  with  $\alpha_t := (1 - \eta \mu_F / 2)^{-(t+1)}$  and  $A_T := \sum_{t=0}^{T-1} \alpha_t$ , and  $\tilde{\mathcal{O}}(\cdot)$  hides both constants and polylogarithmic factors.

**Corollary 1.** *When there is no client sampling (i.e.,  $S = N$ ), we can choose either (i)  $\beta = \Theta(\sqrt{N}/T)$  if  $\sqrt{N} \geq T$  (i.e., massive clients) or (ii)  $\beta = \Theta(N\sqrt{R})$  otherwise, to obtain either linear speedup  $\mathcal{O}(1/(TRN))$  or quadratic speedup  $\mathcal{O}(1/(TRN^2))$  w.r.t computation rounds, respectively.*

**Remark 1.** Theorem 1 (a) shows the convergence of the global model w.r.t four error terms, where the expectation is w.r.t the randomness of mini-batch and client samplings. While the first term shows that a carefully chosen constant step size can reduce the initial error  $\|w_0 - w^*\|^2$  linearly, the last term means that pFedMe converges towards a  $\frac{\lambda^2 \delta^2}{\mu_F}$ -neighbourhood of  $w^*$ , due to the approximation error  $\delta$  at each local round. The second error term is due to the client sampling, which obviously is 0 when  $S = N$ . If we choose  $S$  such that  $S/N$  corresponds to a fixed ratio, e.g., 0.5, then we can obtain a linear speedup  $\mathcal{O}(1/(TN))$  w.r.t communication rounds for client sampling error. The third error term is due to client drift with multiple local updates. According to Corollary 1, we are able to obtain the quadratic speedup, while most of existing FL convergence analysis of strongly convex loss functions can only achieve linear speedup [3, 6, 19]. Theorem 1 (b) shows the convergence of personalized models in average to a ball of center  $w^*$  and radius  $\mathcal{O}(\frac{\lambda^2 \delta^2}{\mu_F} + \frac{\sigma_{F,1}^2}{\lambda^2} + \delta^2)$ , which shows that  $\lambda$  can be controlled to trade off reducing the errors between  $\delta^2$  and  $\sigma_{F,1}^2$ .

**Theorem 2** (Nonconvex and smooth pFedMe's convergence). *Let Assumptions 1(b), 2, and 3 hold. If  $\eta \leq \frac{\hat{\eta}_2}{\beta R}$ , where  $\hat{\eta}_2 := \frac{1}{75L_F \lambda^2}$  with  $\lambda \geq \sqrt{8L^2 + 1}$  and  $\beta \geq 1$ , then we have*

$$\begin{aligned} (a) \mathbb{E} [\|\nabla F(w_{t^*})\|^2] &\leq \mathcal{O}(\mathbb{E} [\|\nabla F(w_{t^*})\|^2]) := \\ &\mathcal{O}\left(\frac{\Delta_F}{\hat{\eta}_2 T} + \frac{(\Delta_F L_F \sigma_{F,2}^2 (N/S-1))^{\frac{1}{2}}}{\sqrt{TN}} + \frac{(\Delta_F)^{\frac{2}{3}} (R\sigma_{F,2}^2 + \lambda^2 \delta^2)^{\frac{1}{3}}}{\beta^{\frac{4}{3}} R^{\frac{1}{3}} T^{\frac{2}{3}}} + \lambda^2 \delta^2\right) \\ (b) \frac{1}{N} \sum_{i=1}^N \mathbb{E} [\|\tilde{\theta}_i^{t^*}(w_{t^*}) - w_{t^*}^*\|^2] &\leq \mathcal{O}(\mathbb{E} [\|\nabla F(w_{t^*})\|^2]) + \mathcal{O}\left(\frac{\sigma_{F,2}^2}{\lambda^2} + \delta^2\right), \end{aligned}$$

where  $\Delta_F := F(w_0) - F^*$ , and  $t^* \in \{0, \dots, T-1\}$  is sampled uniformly.

**Corollary 2.** *When there is no client sampling, we can choose  $\beta = \Theta(N^{1/2} R^{1/4})$  and  $\Theta(T^{1/3}) = \Theta((NR)^{2/3})$  to obtain a sublinear speed-up of  $\mathcal{O}(1/(TRN)^{2/3})$ .*

**Remark 2.** Theorem 2 shows a similar convergence structure to that of Theorem 1, but with a sublinear rate for nonconvex case. According to Corollary 2, we are able to obtain the sublinear speedup  $\mathcal{O}(1/(TRN)^{2/3})$ , while most of existing convergence analysis for nonconvex FL can only achieve a sublinear speed-up of  $\mathcal{O}(1/\sqrt{TRN})$  [3, 6, 49].

## 5 Experimental Results and Discussion

In this section, we validate the performance of pFedMe when the data distributions are heterogeneous and non-i.i.d. We first observe the effect of hyperparameters  $R$ ,  $K$ ,  $|\mathcal{D}|$ ,  $\lambda$ , and  $\beta$  on the convergence of pFedMe. We then compare pFedMe with FedAvg and Per-FedAvg in both  $\mu$ -strongly convex and nonconvex settings.

### 5.1 Experimental Settings

We consider a classification problem using both real (MNIST) and synthetic datasets. MNIST [51] is a handwritten digit dataset containing 10 labels and 70,000 instances. Due to the limitation on MNIST’s data size, we distribute the complete dataset to  $N = 20$  clients. To model a heterogeneous setting in terms of local data sizes and classes, each client is allocated a different local data size in the range of [1165, 3834], and only has 2 of the 10 labels. For synthetic data, we adopt the data generation and distribution procedure from [15], using two parameters  $\bar{\alpha}$  and  $\bar{\beta}$  to control how much the local model and the dataset of each client differ, respectively. Specifically, the dataset serves a 10-class classifier using 60-dimensional real-valued data. We generate a synthetic dataset with  $\bar{\alpha} = 0.5$  and  $\bar{\beta} = 0.5$ . Each client’s data size is in the range of [250, 25810]. Finally, we distribute the data to  $N = 100$  clients according to the power law in [15].

We fix the subset of clients  $S = 5$  for MNIST, and  $S = 10$  for Synthetic. We compare the algorithms using both cases of the same and fine-tuned learning rates, batch sizes, and number of local and global iterations. For  $\mu$ -strongly convex setting, we consider a  $l_2$ -regularized multinomial logistic regression model (MLR) with the softmax activation and cross-entropy loss functions. For nonconvex case, a two-layer deep neural network (DNN) is implemented with hidden layer of size 100 for MNIST and 20 for Synthetic using ReLU activation and a softmax layer at the end. For pFedMe, we use gradient descent to obtain  $\delta$ -approximate  $\tilde{\theta}_i(w_{i,r}^t)$  and the personalized model is evaluated on the personalized parameter  $\tilde{\theta}_i$  while the global model is evaluated on  $w$ . For the comparison with Per-FedAvg, we use its personalized model which is the local model after taking an SGD step from the global model.

All datasets are split randomly with 75% and 25% for training and testing, respectively. All experiments were conducted using PyTorch [52] version 1.4.0. The code and datasets are available online<sup>1</sup>.

### 5.2 Effect of hyperparameters

To understand how different hyperparameters such as  $R$ ,  $K$ ,  $|\mathcal{D}|$ ,  $\lambda$ , and  $\beta$  affect the convergence of pFedMe in both  $\mu$ -strongly convex and nonconvex settings, we conduct various experiments on MNIST dataset with  $\eta = 0.005$  and  $S = 5$ .

**Effects of local computation rounds  $R$ :** When the communication is relatively costly, the server tends to allow users to have more local computations, which can lead to less global model updates and thus faster convergence. Therefore, we monitor the behavior of pFedMe using a number of values of  $R$ , which results in Fig. 1. The results show that larger values of  $R$  have a benefit on the convergence of both the personalized and the global models. There is, nevertheless, a trade-off between the computations and communications: while larger  $R$  requires more computations at local users, smaller  $R$  needs more global communication rounds to converge. To balance this trade-off, we fix  $R = 20$  and evaluate the effect of other hyperparameters accordingly.

**Effects of computation complexity  $K$ :** As  $K$  allows for approximately finding the personalized model  $\theta$ ,  $K$  is also considered as a hyper-parameter of pFedMe. In Fig. 2, only the value of  $K$  is changed during the experiments. We observe that pFedMe requires a small value of  $K$  (around 3 to 5 steps) to approximately compute the personalized model. Larger values of  $K$ , such as 7, do not

<sup>1</sup><https://github.com/CharlieDinh/pFedMe>

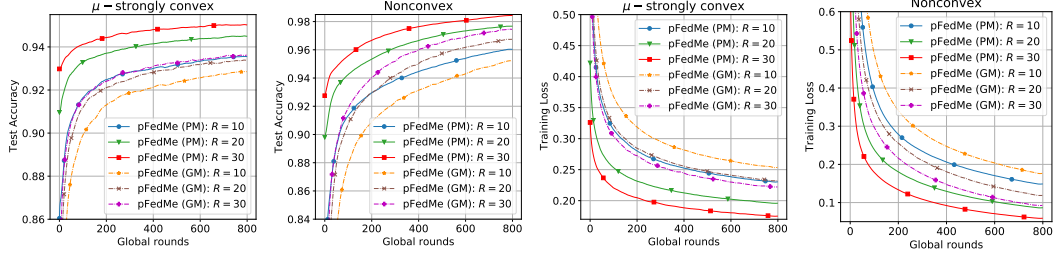


Figure 1: Effect of  $R$  on the convergence of pFedMe in  $\mu$ -strongly convex and nonconvex settings on MNIST ( $|\mathcal{D}| = 20$ ,  $\lambda = 15$ ,  $K = 5$ ,  $\beta = 1$ ).

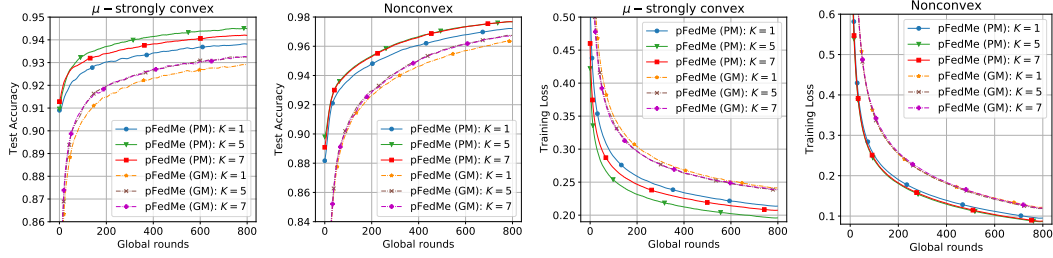


Figure 2: Effect of  $K$  on the convergence of pFedMe in  $\mu$ -strongly convex and nonconvex settings on MNIST ( $|\mathcal{D}| = 20$ ,  $\lambda = 15$ ,  $R = 20$ ,  $\beta = 1$ ).

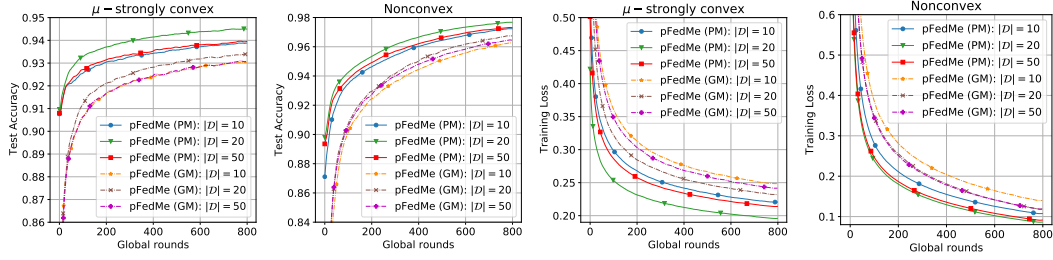


Figure 3: Effect of  $|\mathcal{D}|$  on the convergence of pFedMe in  $\mu$ -strongly convex and nonconvex settings on MNIST ( $\lambda = 15$ ,  $R = 20$ ,  $K = 5$ ,  $\beta = 1$ ).

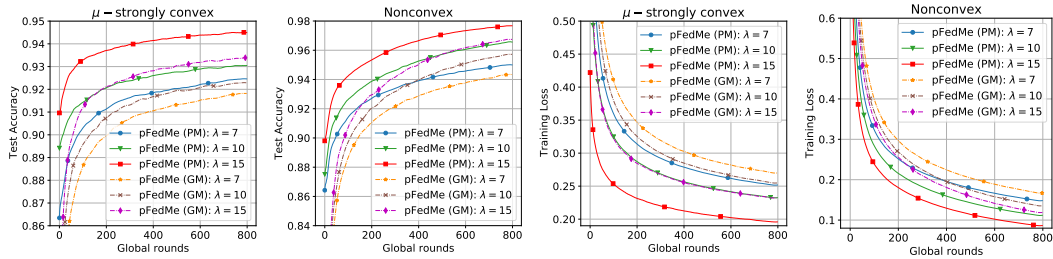


Figure 4: Effect of  $\lambda$  on the convergence of pFedMe in  $\mu$ -strongly convex and nonconvex settings on MNIST ( $|\mathcal{D}| = 20$ ,  $R = 20$ ,  $K = 5$ ,  $\beta = 1$ ).

show the improvement on the convergence of the personalized model nor the global model. Similar to  $R$ , larger  $K$  also requires more user's computation, which has negative effects on user energy consumption. Therefore, the value of  $K = 5$  is chosen for the remaining experiments.



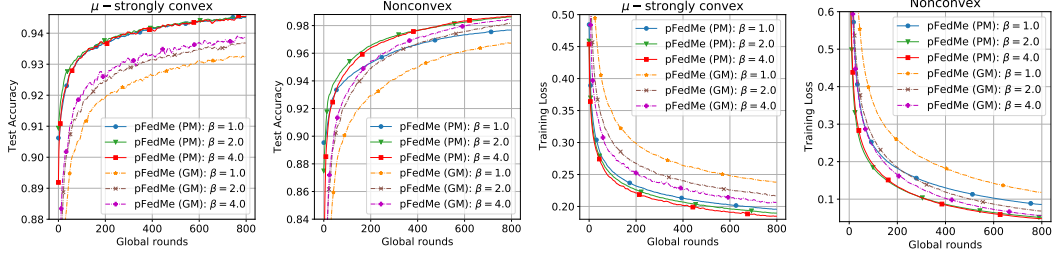


Figure 5: Effect of  $\beta$  on the convergence of pFedMe in  $\mu$ -strongly convex and nonconvex settings on MNIST ( $|\mathcal{D}| = 20$ ,  $\lambda = 15$ ,  $R = 20$ ,  $K = 5$ ).

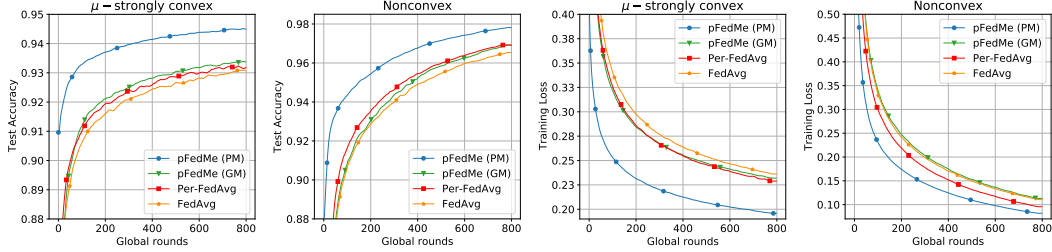


Figure 6: Performance comparison of pFedMe, FedAvg, and Per-FedAvg in  $\mu$ -strongly convex and nonconvex settings using MNIST ( $\eta = 0.005$ ,  $|\mathcal{D}| = 20$ ,  $S = 5$ ,  $\beta = 1$  for all experiments).

**Effects of Mini-Batch size  $|\mathcal{D}|$ :** As mentioned in the Lemma 1,  $|\mathcal{D}|$  is one of the parameters which can be controlled to adjust the value of  $\delta$ . In Fig. 3, when the size of the mini-batch is increased, pFedMe has the higher convergence rate. However, very large  $|\mathcal{D}|$  will not only slow the convergence of pFedMe but also requires higher computations at the local users. During the experiments, the value of  $|\mathcal{D}|$  is configured as a constant value equal to 20.

**Effects of regularization  $\lambda$ :** Fig. 4 shows the convergence rate of pFedMe with different values of  $\lambda$ . In all settings, larger  $\lambda$  allows for faster convergence; however, we also observe that the significantly large  $\lambda$  will hurt the performance of pFedMe by making pFedMe diverge. Therefore,  $\lambda$  should be tuned carefully depending on the dataset. We fix  $\lambda = 15$  for all scenarios with MNIST.

**Effects of  $\beta$ :**

Fig. 5 illustrates how  $\beta$  ( $\beta \geq 1$ ) affects both the personalized and global models. It is noted that when  $\beta = 1$ , it is similar to model averaging of FedAvg. According to the figure, it is beneficial to shift the value of  $\beta$  to be larger as it allows pFedMe to converge faster, especially the global model. However, turning  $\beta$  carefully is also significant to prevent the divergence and instability of pFedMe. For example, when  $\beta$  moves to the large value, to stabilize the global model as well as the personalized model, the smaller value of  $\eta$  needs to be considered. Alternatively,  $\beta$  and  $\eta$  should be adjusted in inverse proportion to reach the stability of pFedMe.

### 5.3 Performance Comparison

In order to highlight the empirical performance of pFedMe, we perform several comparisons between pFedMe, FedAvg, and Per-FedAvg. We first use the same parameters for all algorithms as an initial comparison. As algorithms behave differently when hyperparameters are changed, we conduct a grid search on a wide range of hyperparameters to figure out the combination of fine-tuned parameters that achieves the highest test accuracy w.r.t. each algorithm. We use both personalized model (PM) and the global model (GM) of pFedMe for comparisons.

The comparisons for MNIST dataset are shown in Fig. 6 (the same hyperparameters) and Table. 1 (fine-tuned hyperparameters). Fig. 6 shows that the pFedMe's personalized models in strongly convex setting are 1.1%, 1.3%, and 1.5% more accurate than its global model, Per-FedAvg, and FedAvg,

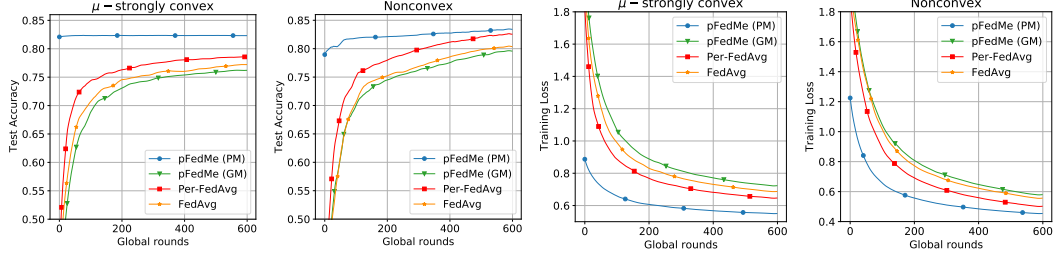


Figure 7: Performance comparison of pFedMe, FedAvg, and Per-FedAvg in  $\mu$ -strongly convex and nonconvex settings using Synthetic ( $\eta = 0.005$ ,  $|\mathcal{D}| = 20$ ,  $S = 10$ ,  $\beta = 1$  for all experiments).

Table 1: Comparison using fine-tuned hyperparameters. We fix  $|\mathcal{D}| = 20$ ,  $R = 20$ ,  $K = 5$ , and  $T = 800$  for MNIST, and  $T = 600$  for Synthetic,  $\beta = 2$  for pFedMe ( $\hat{\alpha}$  and  $\hat{\beta}$  are learning rates of Per-FedAvg).

| Algorithm  | Model | MNIST     |                                   |                                    | Synthetic |                                   |                                    |
|------------|-------|-----------|-----------------------------------|------------------------------------|-----------|-----------------------------------|------------------------------------|
|            |       | $\lambda$ | $\eta(\hat{\alpha}, \hat{\beta})$ | Accuracy (%)                       | $\lambda$ | $\eta(\hat{\alpha}, \hat{\beta})$ | Accuracy (%)                       |
| FedAvg     | MLR   |           | 0.02                              | $93.96 \pm 0.02$                   |           | 0.02                              | $77.62 \pm 0.11$                   |
| Per-FedAvg | MLR   |           | 0.03, 0.003                       | $94.37 \pm 0.04$                   |           | 0.02, 0.002                       | $81.49 \pm 0.09$                   |
| pFedMe-GM  | MLR   | 15        | 0.01                              | $94.18 \pm 0.06$                   | 20        | 0.01                              | $78.65 \pm 0.25$                   |
| pFedMe-PM  | MLR   | 15        | 0.01                              | <b><math>95.62 \pm 0.04</math></b> | 20        | 0.01                              | <b><math>83.20 \pm 0.06</math></b> |
| FedAvg     | DNN   |           | 0.02                              | $98.79 \pm 0.03$                   |           | 0.03                              | $83.64 \pm 0.22$                   |
| Per-FedAvg | DNN   |           | 0.02, 0.001                       | $98.90 \pm 0.02$                   |           | 0.01, 0.001                       | $85.01 \pm 0.10$                   |
| pFedMe-GM  | DNN   | 30        | 0.01                              | $99.16 \pm 0.03$                   | 30        | 0.01                              | $84.17 \pm 0.35$                   |
| pFedMe-PM  | DNN   | 30        | 0.01                              | <b><math>99.46 \pm 0.01</math></b> | 30        | 0.01                              | <b><math>86.36 \pm 0.15</math></b> |

respectively. The corresponding figures for nonconvex setting are 0.9%, 0.9%, and 1.3%. Table. 1 shows that when using fine-tuned hyperparameters, the pFedMe’s personalized model is the best performer in all settings.

For Synthetic dataset, the comparisons for the utilizing the same parameters and the fine-tuned parameter are presented in Fig. 7 and Table. 1, respectively. In Fig. 7, even though the global model of pFedMe is less well-performed than others concerning testing accuracy and training loss, pFedMe’s personalized model still shows its advantages as achieving the highest testing accuracy and smallest training loss. Fig. 7 shows that pFedMe’s personalized model is 6.1%, 3.8%, and 5.2% more accurate than its global model, Per-FedAvg, and FedAvg, respectively. The corresponding figures for the nonconvex setting are 3.9%, 0.7%, and 3.1%. In addition, with fine-tuned hyperparameters in Table. 1, the personalized model of pFedMe beats others in all settings while the global model of pFedMe only performs better than FedAvg.

From the experimental results, when the data among clients are non-i.i.d, both pFedMe and Per-Avg gain higher testing accuracy than FedAvg as they allow the global model to be personalized for a specific client. However, by optimizing the personalized model approximately with multiple gradient updates and avoiding computing the Hessian matrix, the personalized model of pFedMe is more advantageous than Per-FedAvg in terms of the convergence rate and the computation complexity.

## 6 Conclusion

In this paper, we propose pFedMe as a personalized FL algorithm that can adapt to the statistical diversity issue to improve the FL performance. Our approach makes use of the Moreau envelope function which helps decompose the personalized model optimization from global model learning, which allows pFedMe to update the global model similarly to FedAvg, yet in parallel to optimize the personalized model w.r.t each client’s local data distribution. Theoretical results show that pFedMe

can achieve the state-of-the-art convergence speedup rate. Experimental results demonstrate that pFedMe outperforms the vanilla FedAvg and the meta-learning based personalized FL algorithm Per-FedAvg in both convex and non-convex settings, using both real and synthetic datasets. Finally, the degree to which personalization becomes provably useful is a topic of experimental research, as parameters will need to be adapted to each dataset and federated setting.

## A Proof of the Results

In this section, we first provide some existing results useful for following proofs. We then present the proofs of Lemma 1, Lemma 2, Theorem 1, and Theorem 2.

### A.1 Review of useful existing results

**Proposition 2.** [53, Theorems 2.1.5 and 2.1.10] *If a function  $F_i(\cdot)$  is  $L_F$ -smooth and  $\mu_F$ -strongly convex,  $\forall w, w'$ , we have the following useful inequalities, in respective order,*

$$\begin{aligned}\|\nabla F_i(w) - \nabla F_i(w')\|^2 &\leq 2L_F(F_i(w) - F_i(w') - \langle \nabla F_i(w'), w - w' \rangle) \\ \mu_F \|w - w'\| &\leq \|\nabla F_i(w) - \nabla F_i(w')\|.\end{aligned}$$

where  $w^*$  is the solution to problem  $\min_{w \in \mathbb{R}^d} F_i(w)$ , i.e.,  $\nabla F_i(w^*) = 0$ .

**Proposition 3.** *For any vector  $x_i \in \mathbb{R}^d$ ,  $i = 1, \dots, M$ , by Jensen's inequality, we have*

$$\left\| \sum_{i=1}^M x_i \right\|^2 \leq M \sum_{i=1}^M \|x_i\|^2.$$

### A.2 Proof of Lemma 1

*Proof.* We first prove case (a). Let  $h_i(\theta_i; w_{i,r}^t) := f_i(\theta_i) + \frac{\lambda}{2} \|\theta_i - w_{i,r}^t\|^2$ . Then  $h_i(\theta_i; w_{i,r}^t)$  is  $(\lambda + \mu)$ -strongly convex with its unique solution  $\hat{\theta}_i(w_{i,r}^t)$ . Then, by Proposition 2, we have

$$\begin{aligned}\|\tilde{\theta}_i(w_{i,r}^t) - \hat{\theta}_i(w_{i,r}^t)\|^2 &\leq \frac{1}{(\lambda + \mu)^2} \|\nabla h_i(\tilde{\theta}_i; w_{i,r}^t)\|^2 \\ &\leq \frac{2}{(\lambda + \mu)^2} \left( \|\nabla h_i(\tilde{\theta}_i; w_{i,r}^t) - \nabla \tilde{h}_i(\tilde{\theta}_i; w_{i,r}^t, \mathcal{D}_i)\|^2 + \|\nabla \tilde{h}_i(\tilde{\theta}_i; w_{i,r}^t, \mathcal{D}_i)\|^2 \right) \\ &\leq \frac{2}{(\lambda + \mu)^2} \left( \|\nabla \tilde{f}_i(\tilde{\theta}_i; \mathcal{D}_i) - \nabla f_i(\tilde{\theta}_i)\|^2 + \nu \right) \\ &= \frac{2}{(\lambda + \mu)^2} \left( \frac{1}{|\mathcal{D}|^2} \left\| \sum_{\xi_i \in \mathcal{D}_i} \nabla \tilde{f}_i(\tilde{\theta}_i; \xi_i) - \nabla f_i(\tilde{\theta}_i) \right\|^2 + \nu \right),\end{aligned}$$

where the second inequality is by Proposition 3. Taking expectation to both sides, we have

$$\begin{aligned}\mathbb{E} \left[ \|\tilde{\theta}_i(w_{i,r}^t) - \hat{\theta}_i(w_{i,r}^t)\|^2 \right] &= \frac{2}{(\lambda + \mu)^2} \left( \frac{1}{|\mathcal{D}|^2} \sum_{\xi_i \in \mathcal{D}_i} \mathbb{E}_{\xi_i} \left[ \|\nabla \tilde{f}_i(\tilde{\theta}_i; \xi_i) - \nabla f_i(\tilde{\theta}_i)\|^2 \right] + \nu \right) \\ &\leq \frac{2}{(\lambda + \mu)^2} \left( \frac{\gamma_f^2}{|\mathcal{D}|} + \nu \right),\end{aligned}$$

where the first equality is due to  $\mathbb{E} \left[ \left\| \sum_{i=1}^M X_i - \mathbb{E}[X_i] \right\|^2 \right] = \sum_{i=1}^M \mathbb{E} [\|X_i - \mathbb{E}[X_i]\|^2]$  with  $M$  independent random variables  $X_i$  and the unbiased estimate  $\mathbb{E} [\nabla \tilde{f}_i(\tilde{\theta}_i; \xi_i)] = \nabla f_i(\tilde{\theta}_i)$ , and the last inequality is due to Assumption 2.

The proof of case (b) follows similarly, considering that  $h_i(\theta_i; w_{i,r}^t)$  is  $(\lambda - L)$ -strongly convex.  $\square$

### A.3 Proof of Lemma 2

*Proof.* We first prove case (a).

$$\begin{aligned}
\frac{1}{N} \sum_{i=1}^N \|\nabla F_i(w)\|^2 &\leq \frac{1}{N} \sum_{i=1}^N 2 \left( \|\nabla F_i(w) - \nabla F_i(w^*)\|^2 + \|\nabla F_i(w^*)\|^2 \right) \\
&\leq 4L_F(F(w) - F(w^*)) + \frac{2}{N} \sum_{i=1}^N \|\nabla F_i(w^*)\|^2,
\end{aligned}$$

where the first and the second inequalities are due to Propositions 3 and 2, respectively.

We next prove case (b):

$$\begin{aligned}
&\|\nabla F_i(w) - \nabla F(w)\|^2 \\
&= \left\| \lambda(w - \hat{\theta}_i(w)) - \frac{1}{N} \sum_{j=1}^N \lambda(w - \hat{\theta}_j(w)) \right\|^2 \\
&= \left\| \nabla f_i(\hat{\theta}_i(w)) - \frac{1}{N} \sum_{j=1}^N \nabla f_j(\hat{\theta}_j(w)) \right\|^2 \\
&= 2 \left\| \nabla f_i(\hat{\theta}_i(w)) - \frac{1}{N} \sum_{j=1}^N \nabla f_j(\hat{\theta}_i(w)) \right\|^2 + 2 \left\| \frac{1}{N} \sum_{j=1}^N \nabla f_j(\hat{\theta}_i(w)) - \nabla f_j(\hat{\theta}_j(w)) \right\|^2,
\end{aligned}$$

where the second inequality is due to the first-order condition  $\nabla f_i(\hat{\theta}_i(w)) - \lambda(w - \hat{\theta}_i(w)) = 0$ , and the last one is due to Proposition 3. Taking the average over the number of clients, we have

$$\frac{1}{N} \sum_{i=1}^N \|\nabla F_i(w) - \nabla F(w)\|^2 \leq 2\sigma_f^2 + \frac{2}{N^2} \sum_{i=1}^N \sum_{j=1}^N \|\nabla f_j(\hat{\theta}_i(w)) - \nabla f_j(\hat{\theta}_j(w))\|^2 \quad (9)$$

$$\leq 2\sigma_f^2 + \frac{2L^2}{N^2} \sum_{i=1}^N \sum_{j=1}^N \|\hat{\theta}_i(w) - \hat{\theta}_j(w)\|^2 \quad (10)$$

$$\leq 2\sigma_f^2 + \frac{2L^2}{N^2} \sum_{i=1}^N \sum_{j=1}^N 2 \left( \|\hat{\theta}_i(w) - w\|^2 + \|\hat{\theta}_j(w) - w\|^2 \right) \quad (11)$$

$$\leq 2\sigma_f^2 + \frac{2L^2}{N^2} \sum_{i=1}^N \sum_{j=1}^N \frac{2}{\lambda^2} \left( \|\nabla F_i(w)\|^2 + \|\nabla F_j(w)\|^2 \right) \quad (12)$$

$$\begin{aligned}
&= 2\sigma_f^2 + \frac{8L^2}{\lambda^2} \frac{1}{N} \sum_{i=1}^N \|\nabla F_i(w)\|^2 \\
&= 2\sigma_f^2 + \frac{8L^2}{\lambda^2} \left[ \frac{1}{N} \sum_{i=1}^N \|\nabla F_i(w) - \nabla F(w)\|^2 + \|\nabla F(w)\|^2 \right] \quad (13)
\end{aligned}$$

where (9) is due to Assumption 3 and Proposition 3, which is also used for (11), (10) is due to  $L$ -smoothness of  $f_i(\cdot)$ , (12) is due to Proposition 1, (13) is by the fact that  $\mathbb{E}[\|X\|^2] = \mathbb{E}[\|X - \mathbb{E}[X]\|^2] + \mathbb{E}[\|\mathbb{E}[X]\|^2]$  for any vector of random variable  $X$ . Finally, by re-arranging the terms of (13), we obtain

$$\frac{1}{N} \sum_{i=1}^N \|\nabla F_i(w) - \nabla F(w)\|^2 \leq \frac{2\lambda^2}{\lambda^2 - 8L^2} \sigma_f^2 + \frac{8L^2}{\lambda^2 - 8L^2} \|\nabla F(w)\|^2.$$

□

#### A.4 Proof of Theorem 1

We first define additional notations for the ease of analysis. We next provide supporting lemmas, and finally we will combine them to complete the proof of Theorem 1.

#### A.4.1 Additional notations

We re-write the local update as follows

$$w_{i,r+1}^t = w_{i,r}^t - \underbrace{\eta \lambda(w_{i,r}^t - \tilde{\theta}_i(w_{i,r}^t))}_{=: g_{i,r}^t}$$

which implies

$$\eta \sum_{r=0}^{R-1} g_{i,r}^t = \sum_{r=0}^{R-1} (w_{i,r}^t - w_{i,r+1}^t) = w_{i,0}^t - w_{i,R}^t = w_t - w_{i,R}^t,$$

where  $g_{i,r}^t$  can be considered as the biased estimate of  $\nabla F_i(w_{i,r}^t)$  since  $\mathbb{E}[g_{i,r}^t] \neq \nabla F_i(w_{i,r}^t)$ . We also re-write the global update as follows

$$\begin{aligned} w_{t+1} &= (1 - \beta)w_t + \frac{\beta}{S} \sum_{i \in S^t} w_{i,R}^t \\ &= w_t - \frac{\beta}{S} \sum_{i \in S^t} (w_t - w_{i,R}^t) \\ &= w_t - \underbrace{\eta \beta R}_{=: \tilde{\eta}} \underbrace{\frac{1}{SR} \sum_{i \in S^t} \sum_{r=0}^{R-1} g_{i,r}^t}_{=: g_t}, \end{aligned}$$

where  $\tilde{\eta}$  and  $g_t$  can be interpreted as the step size and approximate stochastic gradient, respectively, of the global update.

#### A.4.2 Supporting lemmas

**Lemma 3** (One-step global update). *Let Assumption 1(b) hold. We have*

$$\begin{aligned} \mathbb{E}[\|w_{t+1} - w^*\|^2] &\leq \left(1 - \frac{\tilde{\eta}\mu_F}{2}\right) \mathbb{E}[\|w_t - w^*\|^2] - \tilde{\eta}(2 - 6L_F\tilde{\eta}) \mathbb{E}[F(w_t) - F(w^*)] \\ &\quad + \frac{\tilde{\eta}(3\tilde{\eta} + 2/\mu_F)}{NR} \sum_{i,r}^{N,R} \mathbb{E}[\|g_{i,r} - \nabla F_i(w_t)\|^2] + 3\tilde{\eta}^2 \mathbb{E}\left[\left\|\frac{1}{S} \sum_{i \in S^t} \nabla F_i(w_t) - \nabla F(w_t)\right\|^2\right], \end{aligned}$$

where  $\sum_{i,r}^{N,R}$  is used as an alternative for  $\sum_{i=1}^N \sum_{r=0}^{R-1}$ .

*Proof.* Denote the expectation conditioning on all randomness prior to round  $t$  by  $\mathbb{E}_t$ . We have

$$\begin{aligned} \mathbb{E}_t[\|w_{t+1} - w^*\|^2] &= \mathbb{E}_t[\|w_t - \tilde{\eta}g_t - w^*\|^2] \\ &= \|w_t - w^*\|^2 - 2\tilde{\eta} \mathbb{E}_t[\langle g_t, w_t - w^* \rangle] + \tilde{\eta}^2 \mathbb{E}_t[\|g_t\|^2]. \end{aligned} \quad (14)$$

We first take expectation of the second term of (14) w.r.t client sampling

$$\begin{aligned} -\mathbb{E}_{S_t}[\langle g_t, w_t - w^* \rangle] &= -\langle \mathbb{E}_{S_t}[g_t], w_t - w^* \rangle \\ &= -\frac{1}{NR} \sum_{i,r}^{N,R} \left( \langle g_{i,r}^t - \nabla F_i(w_t), w_t - w^* \rangle + \langle \nabla F_i(w_t), w_t - w^* \rangle \right), \end{aligned} \quad (15)$$

where the second equality is obtained by having  $\mathbb{E}_{S_t}[g_t] = \mathbb{E}_{S_t}[\frac{1}{SR} \sum_{i,r}^{S^t,R} g_{i,r}^t] = \frac{1}{SR} \sum_{i,r}^{N,R} g_{i,r}^t \mathbb{E}_{S_t}[\mathbb{I}_{i \in S_t}] = \frac{1}{NR} \sum_{i,r}^{N,R} g_{i,r}^t$ , where  $\mathbb{I}_A$  is the indicator function of an event  $A$  and thus  $\mathbb{E}_{S_t}[\mathbb{I}_{i \in S_t}] = S/N$  due to uniform sampling. We then bound two terms of (15) as follows

$$-\frac{1}{N} \sum_{i=1}^N \langle \nabla F_i(w_t), w_t - w^* \rangle \leq F(w^*) - F(w_t) - \frac{\mu_F}{2} \|w_t - w^*\|^2 \quad (16)$$

$$-\frac{2}{NR} \sum_{i,r}^{N,R} \langle g_{i,r}^t - \nabla F_i(w_t), w_t - w^* \rangle \leq \frac{1}{NR} \sum_{i,r}^{N,R} \left( \frac{2}{\mu_F} \|g_{i,r}^t - \nabla F_i(w_t)\|^2 + \frac{\mu_F}{2} \|w_t - w^*\|^2 \right) \quad (17)$$

where the first and second inequalities are due to  $\mu_F$ -strongly convex  $F_i(\cdot)$  and the Peter Paul inequality, respectively.

We next take expectation of the last term of (14) w.r.t client sampling

$$\begin{aligned}
\mathbb{E}_{\mathcal{S}_t} [\|g_t\|^2] &= \mathbb{E}_{\mathcal{S}_t} \left\| \frac{1}{SR} \sum_{i,r}^{S^t R} g_{i,r}^t \right\|^2 \\
&\leq 3\mathbb{E}_{\mathcal{S}_t} \left[ \left\| \frac{1}{SR} \sum_{i,r}^{S^t R} g_{i,r}^t - \nabla F_i(w_t) \right\|^2 + \left\| \frac{1}{S} \sum_{i \in \mathcal{S}_t} \nabla F_i(w_t) - \nabla F(w_t) \right\|^2 + \|\nabla F(w_t)\|^2 \right] \\
&\leq \frac{3}{NR} \sum_{i,r}^{N,R} \|g_{i,r}^t - \nabla F_i(w_t)\|^2 + 3\mathbb{E}_{\mathcal{S}_t} \left\| \frac{1}{S} \sum_{i \in \mathcal{S}_t} \nabla F_i(w_t) - \nabla F(w_t) \right\|^2 + 6L_F(F(w_t) - F(w^*)),
\end{aligned} \tag{18}$$

where the first inequality is by Proposition 3, and the second inequality is by Proposition 2 and

$$\begin{aligned}
\mathbb{E}_{\mathcal{S}_t} \left[ \left\| \frac{1}{SR} \sum_{i,r}^{S^t R} g_{i,r}^t - \nabla F_i(w_t) \right\|^2 \right] &\leq \frac{1}{SR} \mathbb{E}_{\mathcal{S}_t} \left[ \sum_{i,r}^{S^t R} \|g_{i,r}^t - \nabla F_i(w_t)\|^2 \right] \\
&= \frac{1}{SR} \sum_{i,r}^{N,R} \|g_{i,r}^t - \nabla F_i(w_t)\|^2 \mathbb{E}_{\mathcal{S}_t} [\mathbb{I}_{i \in \mathcal{S}_t}] \\
&= \frac{1}{NR} \sum_{i,r}^{N,R} \|g_{i,r}^t - \nabla F_i(w_t)\|^2.
\end{aligned}$$

By substituting (16), (17), and (18) into (14), and take expectation with all history, we finish the proof.  $\square$

**Lemma 4** (Bounded diversity of  $F_i$  w.r.t client sampling).

$$\mathbb{E}_{\mathcal{S}_t} \left\| \frac{1}{S} \sum_{i \in \mathcal{S}_t} \nabla F_i(w_t) - \nabla F(w_t) \right\|^2 \leq \frac{N/S - 1}{N - 1} \sum_{i=1}^N \frac{1}{N} \|\nabla F_i(w_t) - \nabla F(w_t)\|^2.$$

*Proof.* We use similar proof arguments in [18, Lemma 5] as follows

$$\begin{aligned}
\mathbb{E}_{\mathcal{S}_t} \left\| \frac{1}{S} \sum_{i \in \mathcal{S}_t} \nabla F_i(w_t) - \nabla F(w_t) \right\|^2 &= \frac{1}{S^2} \mathbb{E}_{\mathcal{S}_t} \left\| \sum_{i=1}^N \mathbb{I}_{i \in \mathcal{S}_t} (\nabla F_i(w_t) - \nabla F(w_t)) \right\|^2 \\
&= \frac{1}{S^2} \left[ \sum_{i=1}^N \mathbb{E}_{\mathcal{S}_t} [\mathbb{I}_{i \in \mathcal{S}_t}] \|\nabla F_i(w_t) - \nabla F(w_t)\|^2 \right. \\
&\quad \left. + \sum_{i \neq j} \mathbb{E}_{\mathcal{S}_t} [\mathbb{I}_{i \in \mathcal{S}_t} \mathbb{I}_{j \in \mathcal{S}_t}] \langle \nabla F_i(w_t) - \nabla F(w_t), \nabla F_j(w_t) - \nabla F(w_t) \rangle \right] \\
&= \frac{1}{SN} \sum_{i=1}^N \|\nabla F_i(w_t) - \nabla F(w_t)\|^2 + \sum_{i \neq j} \frac{S-1}{SN(N-1)} \langle \nabla F_i(w_t) - \nabla F(w_t), \nabla F_j(w_t) - \nabla F(w_t) \rangle \\
&= \frac{1}{SN} \left( 1 - \frac{S-1}{N-1} \right) \sum_{i=1}^N \|\nabla F_i(w_t) - \nabla F(w_t)\|^2 \\
&= \frac{N/S - 1}{N - 1} \sum_{i=1}^N \frac{1}{N} \|\nabla F_i(w_t) - \nabla F(w_t)\|^2,
\end{aligned}$$

where the third equality is due to  $\mathbb{E}_{\mathcal{S}_t} [\mathbb{I}_{i \in \mathcal{S}_t}] = \mathbb{P}(i \in \mathcal{S}_t) = \frac{S}{N}$  and  $\mathbb{E}_{\mathcal{S}_t} [\mathbb{I}_{i \in \mathcal{S}_t} \mathbb{I}_{j \in \mathcal{S}_t}] = \mathbb{P}(i, j \in \mathcal{S}_t) = \frac{S(S-1)}{N(N-1)}$  for all  $i \neq j$ , and the fourth equality is by  $\sum_{i=1}^N \|\nabla F_i(w_t) - \nabla F(w_t)\|^2 + \sum_{i \neq j} \langle \nabla F_i(w_t) - \nabla F(w_t), \nabla F_j(w_t) - \nabla F(w_t) \rangle = 0$ .  $\square$

**Lemma 5** (Bounded client drift error). *If  $\tilde{\eta} \leq \frac{\beta}{2L_F} \Leftrightarrow \eta \leq \frac{1}{2RL_F}$ , we have*

$$\frac{1}{NR} \sum_{i,r}^{N,R} \mathbb{E} [\|g_{i,r}^t - \nabla F_i(w_t)\|^2] \leq 2\lambda^2\delta^2 + \frac{16L_F^2\tilde{\eta}^2}{\beta^2} \left( 3 \frac{1}{N} \sum_{i=1}^N \mathbb{E} [\|\nabla F_i(w_t)\|^2] + \frac{2\lambda^2\delta^2}{R} \right).$$

*Proof.*

$$\begin{aligned} \mathbb{E} [\|g_{i,r}^t - \nabla F_i(w_t)\|^2] &\leq 2\mathbb{E} [\|g_{i,r}^t - \nabla F_i(w_{i,r}^t)\|^2 + \|\nabla F_i(w_{i,r}^t) - \nabla F_i(w_t)\|^2] \\ &\leq 2 \left( \lambda^2 \mathbb{E} [\|\tilde{\theta}_i(w_{i,r}^t) - \hat{\theta}_i(w_{i,r}^t)\|^2] + L_F^2 \mathbb{E} [\|w_{i,r}^t - w_t\|^2] \right) \\ &\leq 2 \left( \lambda^2 \delta^2 + L_F^2 \mathbb{E} [\|w_{i,r}^t - w_t\|^2] \right), \end{aligned} \quad (19)$$

where the first and second inequalities are due to Propositions 3 and 2, respectively. We next bound the drift of local update of client  $i$  from global model  $\|w_{i,r}^t - w_t\|^2$  as follows

$$\begin{aligned} \mathbb{E} [\|w_{i,r}^t - w_t\|^2] &= \mathbb{E} [\|w_{i,r-1}^t - w_t - \eta g_{i,r-1}^t\|^2] \\ &\leq 2\mathbb{E} [\|w_{i,r-1}^t - w_t - \eta \nabla F_i(w_t)\|^2 + \eta^2 \|g_{i,r-1}^t - \nabla F_i(w_t)\|^2] \\ &\leq 2 \left( 1 + \frac{1}{2R} \right) \mathbb{E} [\|w_{i,r-1}^t - w_t\|^2] + 2(1 + 2R)\eta^2 \mathbb{E} [\|\nabla F_i(w_t)\|^2] \\ &\quad + 4\eta^2 \left( \lambda^2 \delta^2 + L_F^2 \mathbb{E} [\|w_{i,r-1}^t - w_t\|^2] \right) \\ &= 2 \left( 1 + \frac{1}{2R} + 2\eta^2 L_F^2 \right) \mathbb{E} [\|w_{i,r-1}^t - w_t\|^2] + 2(1 + 2R)\eta^2 \mathbb{E} [\|\nabla F_i(w_t)\|^2] + 4\eta^2 \lambda^2 \delta^2 \\ &\leq 2 \left( 1 + \frac{1}{R} \right) \mathbb{E} [\|w_{i,r-1}^t - w_t\|^2] + 2(1 + 2R)\eta^2 \mathbb{E} [\|\nabla F_i(w_t)\|^2] + 4\eta^2 \lambda^2 \delta^2 \end{aligned} \quad (20)$$

$$\leq \left( \frac{6\tilde{\eta}^2}{\beta^2 R} \mathbb{E} [\|\nabla F_i(w_t)\|^2] + \frac{4\tilde{\eta}^2 \lambda^2 \delta^2}{\beta^2 R^2} \right) \sum_{r=0}^{R-1} 2 \left( 1 + \frac{1}{R} \right)^r \quad (21)$$

$$\leq \frac{8\tilde{\eta}^2}{\beta^2} \left( 3\mathbb{E} [\|\nabla F_i(w_t)\|^2] + \frac{2\lambda^2 \delta^2}{R} \right), \quad (22)$$

where (20) is by having  $2\eta^2 L_F^2 = 2L_F^2 \frac{\tilde{\eta}^2}{\beta^2 R^2} \leq \frac{1}{2R^2} \leq \frac{1}{2R}$  when  $\tilde{\eta}^2 \leq \frac{\beta^2}{4L_F^2}$ , for all  $R \geq 1$ . (21) is due to unrolling (20) recursively, and  $2(1 + 2R)\eta^2 = 2(1 + 2R) \frac{\tilde{\eta}^2}{\beta^2 R^2} \leq \frac{6\tilde{\eta}^2}{\beta^2 R}$  because  $\frac{1+2R}{R} \leq 3$  when  $R \geq 1$ . We have (22) because  $\sum_{r=0}^{R-1} (1 + 1/R)^r = \frac{(1+1/R)^R - 1}{1/R} \leq \frac{e-1}{1/R} \leq 2R$ , by using the facts that  $\sum_{i=0}^{n-1} x^i = \frac{x^n - 1}{x - 1}$  and  $(1 + \frac{x}{n})^n \leq e^x$  for any  $x \in \mathbb{R}, n \in \mathbb{N}$ . Substituting (22) to (19), we obtain

$$\mathbb{E} [\|g_{i,r}^t - \nabla F_i(w_t)\|^2] \leq 2\lambda^2\delta^2 + \frac{16\tilde{\eta}^2 L_F^2}{\beta^2} \left( 3\mathbb{E} [\|\nabla F_i(w_t)\|^2] + \frac{2\lambda^2\delta^2}{R} \right). \quad (23)$$

By taking average over  $N$  and  $R$ , we finish the proof.  $\square$

#### A.4.3 Completing the proof of Theorem 1

*Proof.* Before proving the main theorem, we derive the first auxiliary result:

$$\mathbb{E} \left[ \left\| \frac{1}{S} \sum_{i \in S^t} \nabla F_i(w_t) - \nabla F(w_t) \right\|^2 \right] \leq \frac{N/S - 1}{N - 1} \sum_{i=1}^N \frac{1}{N} \mathbb{E} [\|\nabla F_i(w_t) - \nabla F(w_t)\|^2] \quad (24)$$

$$\leq \frac{N/S - 1}{N - 1} \left( 4L_F \mathbb{E} [F(w_t) - F(w^*)] + 2\sigma_{F,1}^2 \right), \quad (25)$$

where (24) is by Lemma 4 and (25) is by Lemma 2 (a).

The second auxiliary result is as follows

$$\frac{\tilde{\eta}(3\tilde{\eta} + 2/\mu_F)}{NR} \sum_{i,r}^{N,R} \mathbb{E} [\|g_{i,r}^t - \nabla F_i(w_t)\|^2]$$

$$\leq \tilde{\eta} \frac{16\delta^2\lambda^2}{\mu_F} + \frac{\tilde{\eta}^3}{\beta^2} \frac{128L_F^2}{\mu_F} \sum_{i=1}^N \frac{1}{N} \left( 3\mathbb{E} [\|\nabla F_i(w_t)\|^2] + \frac{2\delta^2\lambda^2}{R} \right) \quad (26)$$

$$\leq \tilde{\eta} \frac{16\delta^2\lambda^2}{\mu_F} + \frac{\tilde{\eta}^3}{\beta^2} \frac{128L_F^2}{\mu_F} \sum_{i=1}^N \frac{1}{N} \left( 6\mathbb{E} [\|\nabla F_i(w_t) - \nabla F_i(w^*)\|^2] + 6\mathbb{E} [\|\nabla F_i(w^*)\|^2] + \frac{2\delta^2\lambda^2}{R} \right) \quad (27)$$

$$\leq \tilde{\eta} \frac{16\delta^2\lambda^2}{\mu_F} + \frac{\tilde{\eta}^3}{\beta^2} \frac{128L_F^2}{\mu_F} \left( 12L_F \mathbb{E} [F(w_t) - F(w^*)] + \frac{2(3R\sigma_{F,1}^2 + \delta^2\lambda^2)}{R} \right) \quad (28)$$

$$\leq \tilde{\eta} \frac{16\delta^2\lambda^2}{\mu_F} + \frac{\tilde{\eta}^2}{\beta} 768\kappa_F L_F \mathbb{E} [F(w_t) - F(w^*)] + \frac{\tilde{\eta}^3}{\beta^2} \frac{256(3R\sigma_{F,1}^2 + \delta^2\lambda^2)\kappa_F}{R}, \quad (29)$$

where we have (26) by using Lemma 5 and  $3\tilde{\eta} + 2/\mu_F \leq 8/\mu_F$  when  $\tilde{\eta} \leq 2/\mu_F$ . (27) is by the fact that  $\mathbb{E} [\|X\|^2] = \mathbb{E} [\|X - \mathbb{E}[X]\|^2] + \mathbb{E} [\|\mathbb{E}[X]\|^2]$  for any vector of random variable  $X$ . (28) is due to Lemma 2 and  $\|\nabla F(w_t)\|^2 \leq 2L_F(F(w_t) - F(w^*))$  by  $L_F$ -smoothness of  $F(\cdot)$ . (29) is due to  $\tilde{\eta} \leq \frac{\beta}{2L_F}$  and  $\kappa_F := \frac{L_F}{\mu_F}$ .

By substituting (25) and (28) into Lemma 3, we have

$$\begin{aligned} \mathbb{E} [\|w_{t+1} - w^*\|^2] &\leq \\ &\left(1 - \frac{\tilde{\eta}\mu_F}{2}\right) \mathbb{E} [\|w_t - w^*\|^2] - \tilde{\eta} \overbrace{\left[2 - \tilde{\eta} L_F \left(6 + 12 \frac{N/S - 1}{N - 1} + \frac{768\kappa_F}{\beta}\right)\right]}^{\geq 1 \text{ when } \tilde{\eta} \text{ satisfied (31)}} \mathbb{E} [F(w_t) - F(w^*)] \\ &\quad + \underbrace{\tilde{\eta} \frac{16\delta^2\lambda^2}{\mu_F}}_{=:C_1} + \underbrace{\tilde{\eta}^2 \frac{6\sigma_{F,1}^2(N/S - 1)}{N - 1}}_{=:C_2} + \underbrace{\frac{\tilde{\eta}^3}{\beta^2} \frac{256(3R\sigma_{F,1}^2 + \delta^2\lambda^2)\kappa_F}{R}}_{=:C_3} \\ &\leq \left(1 - \frac{\tilde{\eta}\mu_F}{2}\right) \mathbb{E} [\|w_t - w^*\|^2] - \tilde{\eta} \mathbb{E} [F(w_t) - F(w^*)] + \tilde{\eta} C_1 + \tilde{\eta}^2 C_2 + \frac{\tilde{\eta}^3}{\beta^2} C_3, \end{aligned} \quad (30)$$

where we have (30) by using the fact that  $\frac{N/S-1}{N-1} \leq 1$  for the following inequality

$$2 - \tilde{\eta} L_F \left(6 + 12 \frac{N/S - 1}{N - 1} + \frac{768\kappa_F}{\beta}\right) \geq 2 - 6\tilde{\eta} L_F \left(3 + \frac{128\kappa_F}{\beta}\right) \geq 1$$

with the condition

$$\tilde{\eta} \leq \frac{1}{6L_F(3 + 128\kappa_F/\beta)} =: \hat{\eta}_1. \quad (31)$$

We note that  $\hat{\eta}_1 \leq \min\left\{\frac{\beta}{2L_F}, \frac{2}{\mu_F}\right\}$  with  $\beta \geq 1$  and  $L_F \geq \mu_F$ .

Let  $\Delta_t := \|w_t - w^*\|^2$ . By re-arranging the terms and multiplying both sides of (30) with  $\frac{\alpha_t}{\tilde{\eta}A_T}$ , where  $A_T := \sum_{t=0}^{T-1} \alpha_t$ , then we have

$$\begin{aligned} \sum_{t=0}^{T-1} \frac{\alpha_t \mathbb{E} [F(w_t)]}{A_T} - F(w^*) &\leq \sum_{t=0}^{T-1} \mathbb{E} \left[ \left(1 - \frac{\tilde{\eta}\mu_F}{2}\right) \frac{\alpha_t \Delta_t}{\tilde{\eta}A_T} - \frac{\alpha_t \Delta_{t+1}}{\tilde{\eta}A_T} \right] + \frac{\tilde{\eta}^2}{\beta^2} C_3 + \tilde{\eta} C_2 + C_1 \\ &\leq \sum_{t=0}^{T-1} \mathbb{E} \left[ \frac{\alpha_{t-1} \Delta_t - \alpha_t \Delta_{t+1}}{\tilde{\eta}A_T} \right] + \frac{\tilde{\eta}^2}{\beta^2} C_3 + \tilde{\eta} C_2 + C_1 \end{aligned} \quad (32)$$

$$\begin{aligned} &= \frac{1}{\tilde{\eta}A_T} \Delta_0 - \frac{\alpha_{T-1}}{\tilde{\eta}A_T} \mathbb{E} [\Delta_T] + \frac{\tilde{\eta}^2}{\beta^2} C_3 + \tilde{\eta} C_2 + C_1 \\ &\leq \mu_F e^{-\tilde{\eta}\mu_F T/2} \Delta_0 - \frac{\mu_F}{2} \mathbb{E} [\Delta_T] + \frac{\tilde{\eta}^2}{\beta^2} C_3 + \tilde{\eta} C_2 + C_1, \end{aligned} \quad (33)$$



where we have (32) because in order for telescoping, we choose  $\left(1 - \frac{\tilde{\eta}\mu_F}{2}\right)\alpha_t = \alpha_{t-1}$ , and thus  $\alpha_t = \left(1 - \frac{\tilde{\eta}\mu_F}{2}\right)^{-(t+1)}$  by recursive update. Regarding to (33), we have

$$\begin{aligned} A_T &= \sum_{t=0}^{T-1} \left(1 - \frac{\tilde{\eta}\mu_F}{2}\right)^{-(t+1)} \\ &= \left(1 - \frac{\tilde{\eta}\mu_F}{2}\right)^{-T} \sum_{t=0}^{T-1} \left(1 - \frac{\tilde{\eta}\mu_F}{2}\right)^t \\ &= a_{T-1} \frac{1 - \left(1 - \frac{\tilde{\eta}\mu_F}{2}\right)^T}{\tilde{\eta}\mu_F/2} \end{aligned}$$

which implies

$$\frac{a_{T-1}}{\tilde{\eta}\mu_F} \leq A_T \leq \frac{2a_{T-1}}{\tilde{\eta}\mu_F},$$

where the first inequality is due to the fact that  $\left(1 - \frac{\tilde{\eta}\mu_F}{2}\right)^T \leq \exp(-\tilde{\eta}\mu_F T/2) \leq \exp(-1) \leq 1/2$  by setting  $\tilde{\eta}T \geq \frac{2}{\mu_F}$  and the second inequality is due to  $1 - \left(1 - \frac{\tilde{\eta}\mu_F}{2}\right)^T \leq 1$ ; thus we have  $\frac{\alpha_{T-1}}{\tilde{\eta}A_T} \geq \frac{\mu_F}{2}$  and  $\frac{1}{\tilde{\eta}A_T} \leq \mu_F \left(1 - \frac{\tilde{\eta}\mu_F}{2}\right)^T \leq \mu_F e^{-\tilde{\eta}\mu_F T/2}$ .

Due to the convexity of  $F(\cdot)$ , (33) implies

$$\mathbb{E} \left[ F \left( \sum_{t=0}^{T-1} \frac{\alpha_t}{A_T} w_t \right) \right] - F(w^*) + \frac{\mu_F}{2} \mathbb{E} [\Delta_T] \leq \mu_F \Delta_0 e^{-\tilde{\eta}\mu_F T/2} + \frac{\tilde{\eta}^2}{\beta^2} C_3 + \tilde{\eta} C_2 + C_1 \quad (34)$$

which implies

$$\mathbb{E} [F(\bar{w}_T) - F(w^*)] \leq \mu_F \Delta_0 e^{-\tilde{\eta}\mu_F T/2} + \frac{\tilde{\eta}^2}{\beta^2} C_3 + \tilde{\eta} C_2 + C_1. \quad (35)$$

Next, using the techniques in [3, 54, 55], we consider following cases:

- If  $\hat{\eta}_1 \geq \max \left\{ \frac{2 \ln(\mu_F^2 \Delta_0 T/2C_2)}{\mu_F T}, \frac{2}{\mu_F T} \right\} =: \eta'$ , then we choose  $\tilde{\eta} = \eta'$ ; thus, having

$$\begin{aligned} \mathbb{E} [F(\bar{w}_T) - F(w^*)] &\leq \mu_F \Delta_0 e^{-\ln(\mu_F^2 \Delta_0 T/2C_2)} + \eta' C_2 + \frac{\eta'^2}{\beta^2} C_3 + C_1 \\ &\leq \tilde{\mathcal{O}} \left( \frac{C_2}{T\mu_F} \right) + \tilde{\mathcal{O}} \left( \frac{C_3}{T^2 \beta^2 \mu_F^2} \right) + C_1. \end{aligned}$$

- If  $\frac{2}{\mu_F T} \leq \hat{\eta}_1 \leq \frac{2 \ln(\mu_F^2 \Delta_0 T/2C_2)}{\mu_F T}$ , then we choose  $\tilde{\eta} = \hat{\eta}_1$ ; thus, having

$$\mathbb{E} [F(\bar{w}_T) - F(w^*)] \leq \mu_F \Delta_0 e^{-\hat{\eta}_1 \mu_F T/2} + \tilde{\mathcal{O}} \left( \frac{C_2}{T\mu_F} \right) + \tilde{\mathcal{O}} \left( \frac{C_3}{T^2 \beta^2 \mu_F^2} \right) + C_1.$$

Combining two cases, we obtain

$$\begin{aligned} \mathbb{E} [F(\bar{w}_T) - F(w^*)] &\leq \mathcal{O}(\mathbb{E} [F(\bar{w}_T) - F(w^*)]) := \\ &\mathcal{O}(\Delta_0 \mu_F e^{-\hat{\eta}_1 \mu_F T/2}) + \tilde{\mathcal{O}} \left( \frac{(N/S-1)\sigma_{F,1}^2}{\mu_F T N} \right) + \tilde{\mathcal{O}} \left( \frac{(R\sigma_{F,1}^2 + \delta^2 \lambda^2) \kappa_F}{R(T\beta\mu_F)^2} \right) + \mathcal{O} \left( \frac{\lambda^2 \delta^2}{\mu_F} \right), \end{aligned}$$

which finishes the proof of part (a). We next prove part (b) as follows

$$\begin{aligned}
& \mathbb{E} \left[ \|\tilde{\theta}_i^T(w_T) - w^*\|^2 \right] \\
& \leq 3 \mathbb{E} \left[ \|\tilde{\theta}_i^T(w_T) - \hat{\theta}_i^T(w_T)\|^2 + \|\hat{\theta}_i^T(w_T) - w_T\|^2 + \|w_T - w^*\|^2 \right] \\
& \leq 3 \left( \delta^2 + \frac{1}{\lambda^2} \mathbb{E} [\|\nabla F_i(w_T)\|^2] + \mathbb{E} [\|w_T - w^*\|^2] \right) \\
& \leq 3 \left( \delta^2 + \frac{2}{\lambda^2} \mathbb{E} [\|\nabla F_i(w_T) - \nabla F_i(w^*)\|^2 + \|\nabla F_i(w^*)\|^2] + \mathbb{E} [\|w_T - w^*\|^2] \right) \\
& \leq 3 \left( \delta^2 + 3 \mathbb{E} [\|w_T - w^*\|^2] + \frac{2}{\lambda^2} \|\nabla F_i(w^*)\|^2 \right),
\end{aligned}$$

where the last inequality is due to smoothness of  $F_i$  with  $L_F = \lambda$  according to Proposition 1. Take the average over  $N$  clients, we have

$$\begin{aligned}
\frac{1}{N} \sum_{i=1}^N \mathbb{E} [\|\tilde{\theta}_i^T(w_T) - w^*\|^2] & \leq 9 \mathbb{E} [\|w_T - w^*\|^2] + \frac{6\sigma_{F,1}^2}{\lambda^2} + 3\delta^2 \\
& \leq \frac{1}{\mu_F} \mathcal{O}(\mathbb{E}[F(\bar{w}_T) - F(w^*)]) + \mathcal{O}\left(\frac{\sigma_{F,1}^2}{\lambda^2} + \delta^2\right),
\end{aligned}$$

where the last inequality is by using (34) and (35), we can easily obtain

$$\begin{aligned}
\mathbb{E} [\|w_T - w^*\|^2] & \leq \frac{2}{\mu_F} \left( \mu_F \Delta_0 e^{-\tilde{\eta}\mu_F T/2} + \frac{\tilde{\eta}^2}{\beta^2} C_3 + \tilde{\eta} C_2 + C_1 \right) \\
& = \frac{1}{\mu_F} \mathcal{O}(\mathbb{E}[F(\bar{w}_T) - F(w^*)]).
\end{aligned}$$

□

## A.5 Theorem 2

*Proof.* We first prove part (a). Due to the  $L_F$ -smoothness of  $F(\cdot)$ , we have

$$\begin{aligned}
& \mathbb{E} [F(w_{t+1}) - F(w_t)] \\
& \leq \mathbb{E} [\langle \nabla F(w_t), w_{t+1} - w_t \rangle] + \frac{L_F}{2} \mathbb{E} [\|w_{t+1} - w_t\|^2] \\
& = -\tilde{\eta} \mathbb{E} [\langle \nabla F(w_t), g_t \rangle] + \frac{\tilde{\eta}^2 L_F}{2} \mathbb{E} [\|g_t\|^2] \\
& = -\tilde{\eta} \mathbb{E} [\|\nabla F(w_t)\|^2] - \tilde{\eta} \mathbb{E} [\langle \nabla F(w_t), g_t - \nabla F(w_t) \rangle] + \frac{\tilde{\eta}^2 L_F}{2} \mathbb{E} [\|g_t\|^2] \\
& \leq -\tilde{\eta} \mathbb{E} [\|\nabla F(w_t)\|^2] + \frac{\tilde{\eta}}{2} \mathbb{E} [\|\nabla F(w_t)\|^2] + \frac{\tilde{\eta}}{2} \mathbb{E} \left\| \frac{1}{NR} \sum_{i,r}^{N,R} g_{i,r}^t - \nabla F_i(w_t) \right\|^2 + \frac{\tilde{\eta}^2 L_F}{2} \mathbb{E} [\|g_t\|^2]
\end{aligned} \tag{36}$$

$$\begin{aligned}
& \leq -\frac{\tilde{\eta}}{2} \mathbb{E} [\|\nabla F(w_t)\|^2] + \frac{3L_F \tilde{\eta}^2}{2} \mathbb{E} \left\| \frac{1}{S} \sum_{i \in S^t} \nabla F_i(w_t) - \nabla F(w_t) \right\|^2 \\
& \quad + \frac{\tilde{\eta}(1 + 3L_F \tilde{\eta})}{2} \frac{1}{NR} \sum_{i,r}^{N,R} \mathbb{E} [\|g_{i,r}^t - \nabla F_i(w_t)\|^2] + \frac{3\tilde{\eta}^2 L_F}{2} \mathbb{E} [\|\nabla F(w_t)\|^2]
\end{aligned} \tag{37}$$

$$\begin{aligned}
& \leq -\frac{\tilde{\eta}(1 - 3L_F \tilde{\eta})}{2} \mathbb{E} [\|\nabla F(w_t)\|^2] + \frac{3L_F \tilde{\eta}^2}{2} \frac{N/S - 1}{N - 1} \sum_{i=1}^N \frac{1}{N} \mathbb{E} [\|\nabla F_i(w_t) - \nabla F(w_t)\|^2] \\
& \quad + \frac{\tilde{\eta}(1 + 3L_F \tilde{\eta})}{2} \left[ 2\lambda^2 \delta^2 + \frac{16\tilde{\eta}^2 L_F^2}{\beta^2} \left( \frac{2\lambda^2 \delta^2}{R} + 3 \sum_{i=1}^N \frac{1}{N} \mathbb{E} [\|\nabla F_i(w_t) - \nabla F(w_t)\|^2] + 3 \mathbb{E} [\|\nabla F(w_t)\|^2] \right) \right]
\end{aligned} \tag{38}$$

$$\begin{aligned}
&\leq -\frac{\tilde{\eta}(1-3L_F\tilde{\eta})}{2}\mathbb{E}[\|\nabla F(w_t)\|^2] + \frac{3L_F\tilde{\eta}^2}{2}\frac{N/S-1}{N-1}\left(\sigma_{F,2}^2 + \frac{8L^2}{\lambda^2-8L^2}\mathbb{E}[\|\nabla F(w_t)\|^2]\right) \\
&\quad + \frac{\tilde{\eta}(1+3L_F\tilde{\eta})}{2}\left[2\lambda^2\delta^2 + \frac{16\tilde{\eta}^2L_F^2}{\beta^2}\left(\frac{2\lambda^2\delta^2}{R} + 3\sigma_{F,2}^2 + \frac{3\lambda^2}{\lambda^2-8L^2}\mathbb{E}[\|\nabla F(w_t)\|^2]\right)\right] \quad (39) \\
&= -\frac{\tilde{\eta}(1-3L_F\tilde{\eta})}{2}\mathbb{E}[\|\nabla F(w_t)\|^2] + \tilde{\eta}^2L_F\left(\frac{12L^2}{\lambda^2-8L^2}\frac{N/S-1}{N-1} + \frac{24\tilde{\eta}(1+3L_F\tilde{\eta})\lambda^2L_F}{\beta^2(\lambda^2-8L^2)}\right)\mathbb{E}[\|\nabla F(w_t)\|^2] \\
&\quad + \frac{\tilde{\eta}^3}{\beta^2}(1+3L_F\tilde{\eta})\frac{8(3R\sigma_{F,2}^2+2\delta^2\lambda^2)}{R} + \tilde{\eta}^2\sigma_{F,2}^2\left(\frac{3L_F}{2}\frac{N/S-1}{N-1}\right) + \tilde{\eta}(1+3L_F\tilde{\eta})\lambda^2\delta^2 \quad (40)
\end{aligned}$$

$$\begin{aligned}
&\leq -\tilde{\eta}\underbrace{\left[1 - \tilde{\eta}L_F\left(\frac{3}{2} + \frac{12L^2}{\lambda^2-8L^2}\frac{N/S-1}{N-1} + \frac{36\lambda^2}{\lambda^2-8L^2}\right)\right]}_{\geq 1/2 \text{ when } \tilde{\eta} \text{ satisfied (43)}}\mathbb{E}[\|\nabla F(w_t)\|^2] \\
&\quad + \frac{\tilde{\eta}^3}{\beta^2}(1+3L_F\tilde{\eta})\frac{8(3R\sigma_{F,2}^2+2\delta^2\lambda^2)}{R} + \tilde{\eta}^2\frac{3L_F\sigma_{F,2}^2}{2}\frac{N/S-1}{N-1} + \tilde{\eta}(1+3L_F\tilde{\eta})\lambda^2\delta^2 \quad (41)
\end{aligned}$$

$$\begin{aligned}
&\leq -\frac{\tilde{\eta}}{2}\|\nabla F(w_t)\|^2 + \underbrace{\frac{\tilde{\eta}^3}{\beta^2}\frac{16(3R\sigma_{F,2}^2+2\delta^2\lambda^2)}{R}}_{=:C_4} + \underbrace{\tilde{\eta}^2\frac{3L_F\sigma_{F,2}^2}{2}\frac{N/S-1}{N-1}}_{=:C_5} + \underbrace{\tilde{\eta}2\lambda^2\delta^2}_{=:C_6} \quad (42)
\end{aligned}$$

where (36) is due to Cauchy-Swartz and AM-GM inequalities, (37) is by decomposing  $\|g_t\|^2$  into three terms according to (18), and (38) is by using Lemmas 4 and 5, and the fact that  $\mathbb{E}[\|X\|^2] = \mathbb{E}[\|X - \mathbb{E}[X]\|^2] + \mathbb{E}[\|X\|^2]$  for any vector of random variable  $X$ . We have (39) by Lemma 2, and (40) by re-arranging the terms, and (41) by having  $1 + 3L_F\tilde{\eta} \leq 1 + \frac{3\beta}{2} \leq 3\beta$  when  $\tilde{\eta} \leq \frac{\beta}{2L_F}$  according to Lemma 5 and  $\beta \geq 1$ . Finally, we have (42) by using the condition  $\lambda^2 - 8L^2 \geq 1$  and the fact that  $\frac{N/S-1}{N-1} \leq 1$  for the following

$$L_F\left(\frac{3}{2} + \frac{12L^2}{\lambda^2-8L^2}\frac{N/S-1}{N-1} + \frac{36\lambda^2}{\lambda^2-8L^2}\right) \leq \frac{L_F}{2}\left(3 + 24L^2 + 72\lambda^2\right) \leq \frac{L_F}{2}\left(75\lambda^2\right)$$

to get

$$1 - \tilde{\eta}L_F\left(\frac{3}{2} + \frac{12L^2}{\lambda^2-8L^2}\frac{N/S-1}{N-1} + \frac{36\lambda^2}{\lambda^2-8L^2}\right) \geq 1 - \frac{75\tilde{\eta}L_F\lambda^2}{2} \geq \frac{1}{2}$$

with the condition

$$\tilde{\eta} \leq \frac{1}{75L_F\lambda^2} =: \hat{\eta}_2, \quad (43)$$

which also implies  $1 + 3L_F\tilde{\eta} \leq 1 + \frac{1}{25\lambda^2} \leq 2$ .

We note that  $\hat{\eta}_2 \leq \frac{\beta}{2L_F}$  with  $\beta \geq 1$ . By re-arranging the terms of (42) and telescoping, we have

$$\frac{1}{2T}\sum_{t=0}^{T-1}\mathbb{E}[\|\nabla F(w_t)\|^2] \leq \frac{\mathbb{E}[F(w^0) - F(w_T)]}{\tilde{\eta}T} + \frac{\tilde{\eta}^2}{\beta^2}C_4 + \tilde{\eta}C_5 + C_6. \quad (44)$$

Defining  $\Delta_F := F(w^0) - F^*$ , and following the techniques used by [3, 54, 55], we consider two cases:

- If  $\hat{\eta}_2^3 \geq \frac{\beta^2\Delta_F}{TC_4}$  or  $\hat{\eta}_2^2 \geq \frac{\Delta_F}{TC_5}$ , then we choose  $\tilde{\eta} = \min\left\{\left(\frac{\beta^2\Delta_F}{TC_4}\right)^{\frac{1}{3}}, \left(\frac{\Delta_F}{TC_5}\right)^{\frac{1}{2}}\right\}$ ; thus, having

$$\frac{1}{2T}\sum_{t=1}^{T-1}\mathbb{E}[\|\nabla F(w_t)\|^2] \leq \frac{(\Delta_F)^{2/3}C_4^{1/3}}{(\beta^2T)^{2/3}} + \frac{(\Delta_FC_5)^{1/2}}{\sqrt{T}} + C_6.$$

- If  $\hat{\eta}_2^3 \leq \frac{\beta^2\Delta_F}{TC_4}$  and  $\hat{\eta}_2^2 \leq \frac{\Delta_F}{TC_5}$ , then we choose  $\tilde{\eta} = \hat{\eta}_2$ . We have

$$\frac{1}{2T}\sum_{t=0}^{T-1}\mathbb{E}[\|\nabla F(w_t)\|^2] \leq \frac{\Delta_F}{\hat{\eta}_2T} + \frac{(\Delta_F)^{2/3}(C_4)^{1/3}}{(\beta^2T)^{2/3}} + \frac{(\Delta_FC_5)^{1/2}}{\sqrt{T}} + C_6.$$

Combining two cases, and with  $t^*$  uniformly sampled from  $\{0, \dots, T-1, \}$  we have

$$\begin{aligned} \frac{1}{T} \sum_{t=0}^{T-1} \mathbb{E} [\|\nabla F(w_t)\|^2] &= \mathbb{E} [\|\nabla F(w_{t^*})\|^2] \leq \mathcal{O} \left( \mathbb{E} [\|\nabla F(w_{t^*})\|^2] \right) := \\ &\mathcal{O} \left( \frac{\Delta_F}{\hat{\eta}_2 T} + \frac{(\Delta_F)^{\frac{2}{3}} (R\sigma_{F,2}^2 + \lambda^2 \delta^2)^{\frac{1}{3}}}{\beta^{\frac{4}{3}} R^{\frac{1}{3}} T^{\frac{2}{3}}} + \frac{(\Delta_F L_F \sigma_{F,2}^2 (N/S - 1))^{\frac{1}{2}}}{\sqrt{TN}} + \lambda^2 \delta^2 \right) \end{aligned}$$

which proves the first part of Theorem 2.

We next prove part (b) as follows

$$\begin{aligned} \frac{1}{N} \sum_{i=1}^N \mathbb{E} [\|\tilde{\theta}_i^t(w_t) - w_t\|^2] &\leq \frac{1}{N} \sum_{i=1}^N 2\mathbb{E} [\|\tilde{\theta}_i^t(w_t) - \hat{\theta}_i^t\|^2 + \|\hat{\theta}_i^t(w_t) - w_t\|^2] \\ &\leq 2\delta^2 + \frac{2}{N} \sum_{i=1}^N \frac{\mathbb{E} [\|\nabla F_i(w_t)\|^2]}{\lambda^2} \\ &\leq 2\delta^2 + \frac{2}{\lambda^2 - 8L^2} \mathbb{E} [\|\nabla F(w_t)\|^2] + \frac{2\sigma_{F,2}^2}{\lambda^2}, \end{aligned} \quad (45)$$

where the first inequality is due to Proposition (3) and the third inequality is by using the fact that  $\mathbb{E} [\|X\|^2] = \mathbb{E} [\|X - \mathbb{E}[X]\|^2] + \mathbb{E} [\|\mathbb{E}[X]\|^2]$  for any vector of random variable  $X$ , we have

$$\begin{aligned} \frac{1}{N} \sum_{i=1}^N \mathbb{E} [\|\nabla F_i(w_t)\|^2] &= \sum_{i=1}^N \frac{1}{N} \left( \mathbb{E} [\|\nabla F_i(w_t) - \nabla F(w_t)\|^2] + \mathbb{E} [\|\nabla F(w_t)\|^2] \right) \\ &\leq \sigma_{F,2}^2 + \frac{\lambda^2}{\lambda^2 - 8L^2} \mathbb{E} [\|\nabla F(w_t)\|^2]. \end{aligned}$$

Summing (45) from  $t = 0$  to  $T$ , we get

$$\frac{1}{TN} \sum_{i=0}^{T-1} \sum_{i=1}^N \mathbb{E} [\|\tilde{\theta}_i^t - w_t\|^2] \leq \frac{2}{\lambda^2 - 8L^2} \frac{1}{T} \sum_{i=0}^{T-1} \mathbb{E} [\|\nabla F(w_t)\|^2] + 2\delta^2 + \frac{2\sigma_{F,2}^2}{\lambda^2},$$

and with  $t^*$  uniformly sampled from  $\{0, \dots, T-1\}$ , we finish the proof.  $\square$

## Broader Impact

There have been numerous applications of FL in practice. One notable commercial FL usage, which has proved successful in recent years, is in the next-character prediction task on mobile devices. However, we believe this technology promises many more breakthroughs in a number of fields in the near future with the help of personalized FL models. In health care, for example, common causes of a disease can be identified from many patients without the need to have access to their raw data. The development of capable personalized models helps build better predictors on patients' conditions, allowing for faster, more efficient diagnosis and treatment.

As much as FL promises, it also comes with a number of challenges. First, an important societal requirement when deploying such technique is that the server must explain which clients' data will be participated and which will not. The explainability and interpretability of a system are necessary for the sake of public understanding and making informed consent. Second, to successfully preserve privacy, FL has to overcome malicious actors who possibly interfere in the training process during communication. The malicious behaviors include stealing personalized models from the server, perform adversarial attacks such as changing a personalized model on some examples while remaining a good performance on average, and attempt to alter the model. Finally, an effective and unbiased FL system must be aware that data and computational power among clients can be extremely uneven in practice and, therefore, must ensure that the contribution of each client to the global model is adjusted to its level of distribution. These challenges help necessitate future research in decentralized learning in general and personalized FL in particular.

## Acknowledgments and Disclosure of Funding

This research is funded by Vietnam National University HoChiMinh City (VNU-HCM) under grant number DS2020-28-01. Tuan Dung Nguyen's work was supported by the School of Engineering Scholarship at the University of Sydney.

## References

- [1] H. B. McMahan, E. Moore, D. Ramage, S. Hampson, and B. A. y. Arcas, "Communication-Efficient Learning of Deep Networks from Decentralized Data," *arXiv:1602.05629 [cs]*, Feb. 2017. [Online]. Available: <http://arxiv.org/abs/1602.05629>
- [2] M. Mohri, G. Sivek, and A. T. Suresh, "Agnostic Federated Learning," *arXiv:1902.00146 [cs, stat]*, Jan. 2019. [Online]. Available: <http://arxiv.org/abs/1902.00146>
- [3] S. P. Karimireddy *et al.*, "SCAFFOLD: Stochastic Controlled Averaging for Federated Learning," *arXiv:1910.06378 [cs, math, stat]*, Feb. 2020. [Online]. Available: <http://arxiv.org/abs/1910.06378>
- [4] K. Pillutla, S. M. Kakade, and Z. Harchaoui, "Robust Aggregation for Federated Learning," *arXiv:1912.13445 [cs, stat]*, Dec. 2019. [Online]. Available: <http://arxiv.org/abs/1912.13445>
- [5] D. Li and J. Wang, "FedMD: Heterogenous Federated Learning via Model Distillation," *arXiv:1910.03581 [cs, stat]*, Oct. 2019. [Online]. Available: <http://arxiv.org/abs/1910.03581>
- [6] Y. Deng, M. M. Kamani, and M. Mahdavi, "Adaptive Personalized Federated Learning," *arXiv:2003.13461 [cs, stat]*, Mar. 2020. [Online]. Available: <http://arxiv.org/abs/2003.13461>
- [7] J.-J. Moreau, "Propriétés des applications 'prox'," *Compte Rendus Acad. Sci.*, no. 256, pp. 1069–1071, 1963.
- [8] A. Fallah, A. Mokhtari, and A. Ozdaglar, "Personalized Federated Learning: A Meta-Learning Approach," *arXiv:2002.07948 [cs, math, stat]*, Feb. 2020. [Online]. Available: <http://arxiv.org/abs/2002.07948>
- [9] N. Guha, A. Talwalkar, and V. Smith, "One-Shot Federated Learning," *arXiv:1902.11175 [cs, stat]*, Mar. 2019. [Online]. Available: <http://arxiv.org/abs/1902.11175>
- [10] A. Reisizadeh, A. Mokhtari, H. Hassani, A. Jadbabaie, and R. Pedarsani, "FedPAQ: A Communication-Efficient Federated Learning Method with Periodic Averaging and Quantization," *arXiv:1909.13014 [cs, math, stat]*, Mar. 2020. [Online]. Available: <http://arxiv.org/abs/1909.13014>
- [11] X. Dai *et al.*, "Hyper-Sphere Quantization: Communication-Efficient SGD for Federated Learning," *arXiv:1911.04655 [cs, stat]*, Nov. 2019. [Online]. Available: <http://arxiv.org/abs/1911.04655>
- [12] J. Wang and G. Joshi, "Cooperative SGD: A unified Framework for the Design and Analysis of Communication-Efficient SGD Algorithms," *arXiv:1808.07576 [cs, stat]*, Jan. 2019. [Online]. Available: <http://arxiv.org/abs/1808.07576>

- [13] T. Lin, S. U. Stich, K. K. Patel, and M. Jaggi, “Don’t Use Large Mini-Batches, Use Local SGD,” *arXiv:1808.07217 [cs, stat]*, Feb. 2020. [Online]. Available: <http://arxiv.org/abs/1808.07217>
- [14] S. U. Stich, “Local SGD Converges Fast and Communicates Little,” *arXiv:1805.09767 [cs, math]*, May 2019. [Online]. Available: <http://arxiv.org/abs/1805.09767>
- [15] T. Li *et al.*, “Federated Optimization in Heterogeneous Networks,” *arXiv:1812.06127 [cs, stat]*, Sep. 2019. [Online]. Available: <http://arxiv.org/abs/1812.06127>
- [16] Y. Zhao *et al.*, “Federated Learning with Non-IID Data,” *arXiv:1806.00582 [cs, stat]*, Jun. 2018. [Online]. Available: <http://arxiv.org/abs/1806.00582>
- [17] F. Haddadpour and M. Mahdavi, “On the Convergence of Local Descent Methods in Federated Learning,” *arXiv:1910.14425 [cs, stat]*, Dec. 2019. [Online]. Available: <http://arxiv.org/abs/1910.14425>
- [18] X. Li, K. Huang, W. Yang, S. Wang, and Z. Zhang, “On the Convergence of FedAvg on Non-IID Data,” *arXiv:1907.02189 [cs, math, stat]*, Feb. 2020. [Online]. Available: <http://arxiv.org/abs/1907.02189>
- [19] A. Khaled, K. Mishchenko, and P. Richtárik, “Tighter Theory for Local SGD on Identical and Heterogeneous Data,” *arXiv:1909.04746 [cs, math, stat]*, Mar. 2020. [Online]. Available: <http://arxiv.org/abs/1909.04746>
- [20] V. Smith, C.-K. Chiang, M. Sanjabi, and A. Talwalkar, “Federated Multi-Task Learning,” *arXiv:1705.10467 [cs, stat]*, Feb. 2018. [Online]. Available: <http://arxiv.org/abs/1705.10467>
- [21] J. C. Duchi, M. I. Jordan, and M. J. Wainwright, “Privacy Aware Learning,” *J. ACM*, vol. 61, no. 6, pp. 1–57, Dec. 2014. [Online]. Available: <https://dl.acm.org/doi/10.1145/2666468>
- [22] H. B. McMahan, D. Ramage, K. Talwar, and L. Zhang, “Learning Differentially Private Recurrent Language Models,” *arXiv:1710.06963 [cs]*, Feb. 2018. [Online]. Available: <http://arxiv.org/abs/1710.06963>
- [23] W. Zhu, P. Kairouz, B. McMahan, H. Sun, and W. Li, “Federated Heavy Hitters Discovery with Differential Privacy,” *arXiv:1902.08534 [cs]*, Feb. 2020. [Online]. Available: <http://arxiv.org/abs/1902.08534>
- [24] N. Agarwal, A. T. Suresh, F. X. X. Yu, S. Kumar, and B. McMahan, “cpSGD: Communication-efficient and differentially-private distributed SGD,” p. 12.
- [25] Z. Li, V. Sharma, and S. P. Mohanty, “Preserving Data Privacy via Federated Learning: Challenges and Solutions,” *IEEE Consumer Electronics Magazine*, vol. 9, no. 3, pp. 8–16, May 2020. [Online]. Available: <https://ieeexplore.ieee.org/document/9055478/>
- [26] F. Hanzely and P. Richtárik, “Federated Learning of a Mixture of Global and Local Models,” *arXiv:2002.05516 [cs, math, stat]*, Feb. 2020. [Online]. Available: <http://arxiv.org/abs/2002.05516>
- [27] Y. Mansour, M. Mohri, J. Ro, and A. T. Suresh, “Three Approaches for Personalization with Applications to Federated Learning,” *arXiv:2002.10619 [cs, stat]*, Feb. 2020. [Online]. Available: <http://arxiv.org/abs/2002.10619>
- [28] M. G. Arivazhagan, V. Aggarwal, A. K. Singh, and S. Choudhary, “Federated Learning with Personalization Layers,” *arXiv:1912.00818 [cs, stat]*, Dec. 2019. [Online]. Available: <http://arxiv.org/abs/1912.00818>
- [29] A. Hard *et al.*, “Federated Learning for Mobile Keyboard Prediction,” *arXiv:1811.03604 [cs]*, Feb. 2019. [Online]. Available: <http://arxiv.org/abs/1811.03604>
- [30] K. Wang *et al.*, “Federated Evaluation of On-device Personalization,” *arXiv:1910.10252 [cs, stat]*, Oct. 2019. [Online]. Available: <http://arxiv.org/abs/1910.10252>
- [31] P. Vanhaesebrouck, A. Bellet, and M. Tommasi, “Decentralized Collaborative Learning of Personalized Models over Networks,” *arXiv:1610.05202 [cs, stat]*, Feb. 2017. [Online]. Available: <http://arxiv.org/abs/1610.05202>
- [32] C. Finn, P. Abbeel, and S. Levine, “Model-Agnostic Meta-Learning for Fast Adaptation of Deep Networks,” *arXiv:1703.03400 [cs]*, Jul. 2017. [Online]. Available: <http://arxiv.org/abs/1703.03400>
- [33] A. Nichol, J. Achiam, and J. Schulman, “On First-Order Meta-Learning Algorithms,” *arXiv:1803.02999 [cs]*, Oct. 2018. [Online]. Available: <http://arxiv.org/abs/1803.02999>
- [34] A. Fallah, A. Mokhtari, and A. Ozdaglar, “On the Convergence Theory of Gradient-Based Model-Agnostic Meta-Learning Algorithms,” *arXiv:1908.10400 [cs, math, stat]*, Mar. 2020. [Online]. Available: <http://arxiv.org/abs/1908.10400>
- [35] M. Khodak, M.-F. Balcan, and A. Talwalkar, “Adaptive Gradient-Based Meta-Learning Methods,” *arXiv:1906.02717 [cs, stat]*, Dec. 2019. [Online]. Available: <http://arxiv.org/abs/1906.02717>
- [36] Y. Jiang, J. Konečný, K. Rush, and S. Kannan, “Improving Federated Learning Personalization via Model Agnostic Meta Learning,” *arXiv:1909.12488 [cs, stat]*, Sep. 2019. [Online]. Available: <http://arxiv.org/abs/1909.12488>
- [37] F. Chen, Z. Dong, Z. Li, and X. He, “Federated Meta-Learning for Recommendation,” Feb. 2018.

- [38] T. Li, A. K. Sahu, A. Talwalkar, and V. Smith, “Federated Learning: Challenges, Methods, and Future Directions,” *arXiv:1908.07873 [cs, stat]*, Aug. 2019. [Online]. Available: <http://arxiv.org/abs/1908.07873>
- [39] P. Kairouz *et al.*, “Advances and Open Problems in Federated Learning.” *arXiv: 1912.04977*, Dec. 2019.
- [40] H. Lin, J. Mairal, and Z. Harchaoui, “Catalyst Acceleration for First-order Convex Optimization: from Theory to Practice,” *arXiv:1712.05654 [math, stat]*, Jun. 2018. [Online]. Available: <http://arxiv.org/abs/1712.05654>
- [41] P. Zhou, X. Yuan, H. Xu, S. Yan, and J. Feng, “Efficient Meta Learning via Minibatch Proximal Update,” in *Advances in Neural Information Processing Systems 32*, H. Wallach *et al.*, Eds. Curran Associates, Inc., 2019, pp. 1534–1544. [Online]. Available: <http://papers.nips.cc/paper/8432-efficient-meta-learning-via-minibatch-proximal-update.pdf>
- [42] X. Li, W. Yang, S. Wang, and Z. Zhang, “Communication-Efficient Local Decentralized SGD Methods,” *arXiv:1910.09126 [cs, math, stat]*, Feb. 2020. [Online]. Available: <http://arxiv.org/abs/1910.09126>
- [43] H. Yu, R. Jin, and S. Yang, “On the Linear Speedup Analysis of Communication Efficient Momentum SGD for Distributed Non-Convex Optimization,” *arXiv:1905.03817 [cs, math]*, May 2019. [Online]. Available: <http://arxiv.org/abs/1905.03817>
- [44] L. Nguyen *et al.*, “New Convergence Aspects of Stochastic Gradient Algorithms,” *Journal of Machine Learning Research*, vol. 20, Nov. 2019.
- [45] A. Khaled, K. Mishchenko, and P. Richtárik, “First Analysis of Local GD on Heterogeneous Data,” *arXiv:1909.04715 [cs, math, stat]*, Mar. 2020. [Online]. Available: <http://arxiv.org/abs/1909.04715>
- [46] C. Lemaréchal and C. Sagastizábal, “Practical Aspects of the Moreau–Yosida Regularization: Theoretical Preliminaries,” *SIAM J. Optim.*, vol. 7, no. 2, pp. 367–385, May 1997. [Online]. Available: <http://epubs.siam.org/doi/10.1137/S1052623494267127>
- [47] C. Planiden and X. Wang, “Strongly Convex Functions, Moreau Envelopes, and the Generic Nature of Convex Functions with Strong Minimizers,” *SIAM J. Optim.*, vol. 26, no. 2, pp. 1341–1364, Jan. 2016. [Online]. Available: <http://epubs.siam.org/doi/10.1137/15M1035550>
- [48] T. Hoheisel, M. Laborde, A. Oberman, and Department of Mathematics and Statistics, McGill University, Montreal, Canada, “A regularization interpretation of the proximal point method for weakly convex functions,” *Journal of Dynamics & Games*, vol. 7, no. 1, pp. 79–96, 2020. [Online]. Available: <http://aims sciences.org/article/doi/10.3934/jdg.2020005>
- [49] S. Reddi *et al.*, “Adaptive Federated Optimization,” *arXiv:2003.00295 [cs, math, stat]*, Feb. 2020. [Online]. Available: <http://arxiv.org/abs/2003.00295>
- [50] S. Bubeck, “Convex Optimization: Algorithms and Complexity,” *arXiv:1405.4980 [cs, math, stat]*, Nov. 2015. [Online]. Available: <http://arxiv.org/abs/1405.4980>
- [51] Y. Lecun, L. Bottou, Y. Bengio, and P. Haffner, “Gradient-Based Learning Applied to Document Recognition,” *Proceedings of the IEEE*, vol. 86, no. 11, pp. 2278–2324, Nov. 1998.
- [52] A. Paszke *et al.*, “PyTorch: An Imperative Style, High-Performance Deep Learning Library,” in *Advances in Neural Information Processing Systems 32*, H. Wallach *et al.*, Eds. Curran Associates, Inc., 2019, pp. 8026–8037. [Online]. Available: <http://papers.nips.cc/paper/9015-pytorch-an-imperative-style-high-performance-deep-learning-library.pdf>
- [53] Y. Nesterov, *Lectures on convex optimization*. New York, NY: Springer Berlin Heidelberg, 2018. [Online]. Available: <https://www.springer.com/gp/book/9783319915777>
- [54] Y. Arjevani, O. Shamir, and N. Srebro, “A Tight Convergence Analysis for Stochastic Gradient Descent with Delayed Updates,” *arXiv:1806.10188 [cs, math, stat]*, Jun. 2018. [Online]. Available: <http://arxiv.org/abs/1806.10188>
- [55] S. U. Stich, “Unified Optimal Analysis of the (Stochastic) Gradient Method,” Jul. 2019. [Online]. Available: <https://arxiv.org/abs/1907.04232v2>