

# Distributed Personalized Empirical Risk Minimization

Anonymous Author(s)\*

## ABSTRACT

This paper advocates a new paradigm Personalized Empirical Risk Minimization (PERM) to facilitate learning from heterogeneous data sources without imposing stringent constraints on computational resources shared by participating devices. In PERM, we aim at learning a distinct model for each client by personalizing the aggregation of local empirical losses by effectively estimating the statistical discrepancy among data distributions, which entails optimal statistical accuracy for all local distributions and overcomes the data heterogeneity issue. To learn personalized models at scale, we propose a distributed algorithm that replaces the standard model averaging with model shuffling to simultaneously optimize PERM objectives for all devices. This also allows to learn distinct model architectures (e.g., neural networks with different number of parameters) for different clients, thus confining to underlying memory and compute resources of individual clients. We rigorously analyze the convergence of proposed algorithm and conduct experiments that corroborate the effectiveness of proposed paradigm.

## KEYWORDS

federated learning

### ACM Reference Format:

Anonymous Author(s). 2018. Distributed Personalized Empirical Risk Minimization. In *Proceedings of Make sure to enter the correct conference title from your rights confirmation email (Conference acronym 'XX)*. ACM, New York, NY, USA, 24 pages. <https://doi.org/XXXXXXX.XXXXXXX>

## 1 INTRODUCTION

Recently *federated learning* (FL) has emerged as an alternative paradigm to centralized learning to encourage federated model sharing and create a framework to support edge intelligence by shifting model training and inference from data centers to potentially scattered—and perhaps self-interested—systems where data is generated. While undoubtedly being a better paradigm than centralized learning, FL usually suffers from *heterogeneity* of data and compute resources among participants. To mitigate the negative effect of data heterogeneity (non-IIDness), two common approaches are clustering and personalization. The key idea behind the clustering based methods [16] is to partition the devices into clusters (coalitions) of similar data distributions and then learn a single shared model for all clients within each cluster. While appealing, the partitioning methods are limited to heuristic ideas such as clustering based on geographical distribution of devices without taking

the actual data distributions into account. In personalization based methods [6, 8, 10, 17, 22], the idea is to learn a distinct *personalized model* for each device alongside the global model, which can be unified as minimizing a bi-level optimization problem [9]. Existing personalization techniques also suffer from two issues. First, as the number of clients grow, while the number of training data increases, the number of parameters to be learned increases which limits to increase the number of clients beyond a certain point to balance data and overall model complexity tradeoff— a phenomenon known as incidental parameters problem [13]. Moreover, since the knowledge transfer among data sources happens through the single global model, it might lead to suboptimal results. To see this, consider an extreme example, where half of the users have identical data distributions, say  $\mathcal{D}$ , while other half share a data distribution that is completely different, say  $-\mathcal{D}$ . In this case, the global model obtained by optimizing averaged local losses converges to a solution that suffers from low test accuracy on all local distributions and it would be preferable to learn a model for each client solely based on its local data.

The aforementioned issues lead to a fundamental question: *What is the best strategy to learn from heterogeneous data sources to achieve optimal statistical accuracy w.r.t. each data source, without imposing stringent constraints on computational resources shared by participating devices?* To answer this question, we advocate a new paradigm dubbed as Personalized Empirical Risk Minimization (PERM) to facilitate learning from massively fragmented private data under resource constraints. Motivated by generalization bounds in multiple source domain adaptation [2, 5, 12, 18], in PERM we aim to learn a distinct model for each client by *personalizing the aggregation* of empirical losses of different data sources. That is, for each client, its personalized model is learnt on a distinct weighted-ERM problem, and the mixing weight is determined by the discrepancy between distributions. While PERM overcomes the data heterogeneity issue, the number of optimization problems (i.e., distinct personalized ERMs) to be solved scales linearly with the number of data sources. To simultaneously optimize all objectives in a scalable and computationally efficient manner, we propose a novel idea which replaces the standard *model averaging* in federated learning with *model shuffling* and establish its convergence rate. This also allows us to learn distinct model architectures (e.g., neural networks with different number of parameters) for different clients.

## 2 PERSONALIZED EMPIRICAL RISK MINIMIZATION

In this section we formally state the problem and introduce PERM as an ideal paradigm to learning from heterogeneous data sources. We assume there are  $N$  distributed devices where each holds a distinct data shard  $\mathcal{S}_i = \{(\mathbf{x}_{i,j}, y_{i,j})\}_{j=1}^{n_i}$  with  $n_i$  training samples that are realized by a source distribution  $\mathcal{D}_i$  over instance space  $\Xi = \mathcal{X} \times \mathcal{Y}$ . The data distributions across the devices are not independently and identically distributed (non-IID or *heterogeneous*), i.e.,  $\mathcal{D}_1 \neq \mathcal{D}_2 \neq \dots \neq \mathcal{D}_m$ , and each distribution corresponds to a local

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for components of this work owned by others than ACM must be honored. Abstracting with credit is permitted. To copy otherwise, or republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Request permissions from [permissions@acm.org](mailto:permissions@acm.org).

Conference acronym 'XX, June 03–05, 2018, Woodstock, NY

© 2018 Association for Computing Machinery.  
ACM ISBN 978-1-4503-XXXX-X/18/06...\$15.00  
<https://doi.org/XXXXXXX.XXXXXXX>

generalization error or true risk  $\mathcal{L}_i(h) = \mathbb{E}_{(x,y) \sim \mathcal{D}_i} [\ell(h(x), y)]$  on unseen samples for any model  $h \in \mathcal{H}$ , where  $\mathcal{H}$  is the hypothesis set (e.g., a linear model or a deep neural network) and a given convex or non-convex loss function  $\ell : \mathcal{Y} \times \mathcal{Y} \rightarrow \mathbb{R}^+$ . We use  $\widehat{\mathcal{L}}_i(h) = (1/n_i) \sum_{(x,y) \in \mathcal{S}_i} \ell(h(x), y)$  to denote the local empirical risk or training loss at  $i$ th data shard  $\mathcal{S}_i$  with  $n_i$  samples. We seek to collaboratively learn a model or personalized models that entail a good generalization on all local distributions, minimizing true risk  $\mathcal{L}_i(\cdot)$ ,  $i = 1, \dots, N$  for all data sources (all-for-all [7]).

A simple non-personalized solution for FL aims to minimize a (weighted) empirical risk over all data shards in a communication-efficient manner [11]:

$$\arg \min_{h \in \mathcal{H}} \sum_{i=1}^N \alpha(i) \widehat{\mathcal{L}}_i(h) \text{ with } \alpha \in \Delta_N, \quad (\text{WERM})$$

where  $\Delta_N = \{\alpha \in \mathbb{R}_+^N \mid \sum_{i=1}^N \alpha(i) = 1\}$  denotes the simplex set.

To motivate our proposal, let us consider the empirical loss  $\sum_{i=1}^m \alpha(i) \widehat{\mathcal{L}}_i(h)$  in WERM with fixed mixing weights  $\alpha \in \Delta_N$ , and denote the optimal solution by  $\widehat{h}_\alpha$ . The excess risk of the learned model  $\widehat{h}_\alpha$  on  $i$ th local distribution  $\mathcal{D}_i$  w.r.t. the optimal local model  $h_i^* = \arg \min_{h \in \mathcal{H}} \mathcal{L}_{\mathcal{D}_i}(h)$  (i.e., all-for-one) can be bounded by [12]

$$\begin{aligned} \mathcal{L}_i(\widehat{h}_\alpha) \leq & \mathcal{L}_i(h_i^*) + \sum_{j=1}^N \alpha(j) R_j(\mathcal{H}) + 2 \sum_{j=1}^N \alpha(j) \text{disc}_{\mathcal{H}}(\mathcal{D}_j, \mathcal{D}_i) \\ & + O\left(\sqrt{\sum_{j=1}^N \frac{\alpha_j^2}{n_j}}\right) \end{aligned} \quad (\text{GEN})$$

where  $R_j(\mathcal{H})$  is the empirical Radamacher complexity  $\mathcal{H}$  w.r.t.  $\mathcal{S}_i$ , and  $\text{disc}_{\mathcal{H}}(\mathcal{D}_i, \mathcal{D}_j)$  is a pseudo-distance on the set of probability measures on  $\Xi$  to assess the discrepancy between the distributions  $\mathcal{D}_i$  and  $\mathcal{D}_j$  with respect to the hypothesis class  $\mathcal{H}$  as defined below [2]:

**DEFINITION 1.** For a model space  $\mathcal{H}$  and  $\mathcal{D}, \mathcal{D}'$  two probability distributions on  $\Xi$ ,

$$\text{disc}_{\mathcal{H}}(\mathcal{D}, \mathcal{D}') = \sup_{h \in \mathcal{H}} |\mathbb{E}_{\xi \sim \mathcal{D}}(\ell(h, \xi)) - \mathbb{E}_{\xi' \sim \mathcal{D}'}(\ell(h, \xi'))|$$

Intuitively, the discrepancy between the two distributions is large, if there exists a predictor that performs well on one of them and badly on the other. On the other hand, if all functions in the hypothesis class perform similarly on both, then  $\mathcal{D}$  and  $\mathcal{D}'$  have low discrepancy.

From GEN, it can be observed that mismatch between pairs of distributions limits the benefits of ERM on all distributions. Indeed, the generalization risk will blow up when the distribution divergence term  $\text{disc}_{\mathcal{H}}(\mathcal{D}_i, \mathcal{D}_j)$  is large. When the discrepancies are small or almost zero, it leads to an ideal sample complexity  $1/\sqrt{n}$  where  $n = n_1 + n_2 + \dots + n_N$  is the total number of samples, which could have been obtained in IID setting with  $\alpha(j) = 1/N$ . Also, we note that even if global model achieves small training error over union of all data (e.g., overparametrized setting), the divergence term still remains which illustrates the poor performance of global model on  $\mathcal{D}_i$ . This implies that even personalization of global model as in WERM can not entail a good generalization on all local distributions as no positive knowledge transfer among data sources occurs.

Interestingly the bound suggest that seeking optimal accuracy on each local distribution requires choosing a distinct mixing of local losses for each client  $i$  that minimizes the right-hand side of GEN. This indicates that in an ideal setting, we can achieve the best accuracy for each local distribution  $\mathcal{D}_i$  by personalizing the WERM, i.e., (i) first estimating  $\alpha_i$ ,  $i = 1, 2, \dots, N$  for each client individually, then (ii) solving a variant of WERM for each client with obtained mixing parameters:

$$\arg \min_{h \in \mathcal{H}_i} \sum_{j=1}^N \alpha_i(j) \widehat{\mathcal{L}}_j(h) \quad \text{for } i = 1, 2, \dots, N \quad (\text{PERM})$$

By doing this each device achieves the optimal local generalization error by learning who to learn with based on number of samples at each source and mismatch between its data distribution with other clients. Compared to WERM, in PERM since we solve a different aggregated empirical loss for each client, we can pick a different model space/model architecture  $\mathcal{H}_i$  for each client to meet its available computational resources.

While this two-stage method is guaranteed to entail optimal test accuracy for all local distributions  $\mathcal{D}_i$ , however, making it scalable requires overcoming two issues. First, estimating the statistical discrepancies between each pair of data sources (i.e.,  $\alpha_i$ ,  $i = 1, \dots, N$ ) is a computing burden as it requires solving  $O(N^2)$  difference of (non)-convex functions in a distributed manner and requires enough samples from each source to entail good accuracy on estimating pairwise discrepancies. Second, we need to solve  $N$  variants of the optimization problem in PERM, possibly each with a different model space, which is infeasible when the number of devices is huge (e.g., cross-device federated learning). In the next section, we propose a simple yet effective idea to overcome these issues in a computationally efficient manner.

### 3 PERM AT SCALE VIA MODEL SHUFFLING

In this section, we propose a method to efficiently estimate the empirical discrepancies among data sources, followed by model shuffling idea to simultaneously solve  $N$  versions of PERM to learn a model for each client. We first start by proposing a two-stage algorithm: estimating mixing parameters followed by model shuffling. Then, we propose a single loop unified algorithm which simultaneously optimizes the mixing weights and personalized models. We assume that the model space  $\mathcal{H}$  is a parameterized by a convex set  $\mathcal{W} \subseteq \mathbb{R}^d$  and use  $f_i(w) := \widehat{\mathcal{L}}_i(w) = \sum_{(x,y) \in \mathcal{S}_i} \ell(w; (x, y))$  to denote the empirical loss at  $i$ th data shard.

#### 3.1 Warmup: a two-stage algorithm

**Stage 1: Mixing parameters estimation.** In the first stage we aim to efficiently estimate the pairwise discrepancy among local distributions to construct mixing parameters  $\alpha_i$ ,  $i = 1, 2, \dots, N$ . From generalization bound GEN and Definition 1, a direct solution to estimate  $\alpha_i$  is to solve the following convex-nonconcave minimax problem for each client:

$$\arg \min_{\alpha \in \Delta_N} \sum_{j=1}^N \alpha(j) \max_{w \in \mathcal{W}} |f_i(w) - f_j(w)| + \sum_{j=1}^N \frac{\alpha(j)^2}{n_j}. \quad (1)$$

However, solving the above minimax problem itself is already challenging: the inner maximization loop is a nonconcave (or difference of convex) problem, so most of existing minimax algorithms will

fail on this problem. To our best knowledge, the only provable deterministic algorithm is [26], and but it is hard to generalize it to stochastic and distributed fashion.

To overcome aforementioned issues, we make two reasonable relaxations. First, inspired by notion of average drift at optimum as a right metric to measure the effect of data heterogeneity in federated learning [23], we propose to measure discrepancy at the optimal solution obtained by solving averaged global objective, i.e.,  $\mathbf{w}^* := \arg \min_{\mathbf{w} \in \mathcal{W}} (1/N) \sum_{i=1}^N f_i(\mathbf{w})$ . That is, we use  $|f_i(\mathbf{w}^*) - f_j(\mathbf{w}^*)|$  to replace  $\sup_{\mathbf{w}} |f_i(\mathbf{w}) - f_j(\mathbf{w})|$ . By doing this, solving the supremum problem reduces to a simple minimization problem for each client, given the optimal global solution. Then we relax the difference between function values to the difference of gradients, i.e., gradient dissimilarity. To see this, for the difference of pair of losses  $f_i(\mathbf{w}) - f_j(\mathbf{w})$ , it is  $2L$  smooth, hence we can upper bound  $|f_i(\mathbf{w}) - f_j(\mathbf{w})|$  as:

$$\begin{aligned} & |f_i(\mathbf{w}) - f_j(\mathbf{w})| \\ & \leq |(f_i(\mathbf{w}^*) - f_j(\mathbf{w}^*)) + \langle \nabla f_i(\mathbf{w}^*) - \nabla f_j(\mathbf{w}^*), \mathbf{w} - \mathbf{w}^* \rangle| + L \|\mathbf{w} - \mathbf{w}^*\|^2 \\ & \leq |f_i(\mathbf{w}^*) - f_j(\mathbf{w}^*)| + \|\nabla f_i(\mathbf{w}^*) - \nabla f_j(\mathbf{w}^*)\| D + LD^2 \end{aligned}$$

where  $D$  is the size of the domain. The above relaxation lead to solving the following tractable optimization problem to decide the per-client mixing parameters:

$$\begin{aligned} \alpha_i^* = \arg \min_{\alpha \in \Delta_N} g_i(\mathbf{w}^*, \alpha) := & \sum_{j=1}^N \alpha(j) \|\nabla f_i(\mathbf{w}^*) - \nabla f_j(\mathbf{w}^*)\|^2 \\ & + \lambda \sum_{j=1}^N \alpha(j)^2 / n_j \end{aligned} \quad (2)$$

where we added a regularization parameter  $\lambda$  and used the squared of gradient dissimilarity for computational convenience. Thus, obtaining all  $N$  mixing parameters requires solving a single ERM to obtain optimal global solution and  $N$  variants of (2). To get the optimal solution in a communication reduced manner, we adapt the Local SGD algorithm [21] (or FedAvg [19]) and find the optimal solution in intermittent communication setting [24] where the clients work in parallel and are allowed to make  $K$  stochastic updates between two communication rounds for  $R$  consecutive rounds. The detailed steps are given in Algorithm A1 in Appendix A for completeness. After obtaining the global model  $\mathbf{w}^R$  we optimize over  $\alpha$  in  $g_i(\mathbf{w}^R, \alpha)$  using  $T_\alpha$  iterations of GD to get  $\hat{\alpha}_i$ . Actually, we will show that as long as  $\mathbf{w}^R$  converge to  $\mathbf{w}^*$ ,  $\hat{\alpha}_i, i = 1, \dots, N$  converges to solution of (2) very fast.

**DEFINITION 2 (GRADIENT DISSIMILARITY).** We define the following quantities to measure the gradient dissimilarity among local functions:

$$\begin{aligned} \zeta_{i,j}(\mathbf{w}) & := \|\nabla f_i(\mathbf{w}) - \nabla f_j(\mathbf{w})\|^2, \quad \bar{\zeta}_j(\mathbf{w}) := \zeta_{i,j}(\mathbf{w}), \\ \zeta & := \sup_{\mathbf{w} \in \mathcal{W}} \max_{i \in [N]} \|\nabla f_i(\mathbf{w}) - (1/N) \sum_{j=1}^N \nabla f_j(\mathbf{w})\|^2. \end{aligned}$$

The following theorem gives the convergence rate of estimated discrepancies to optimal counterparts.

**THEOREM 1.** Assume each  $f_i$  is  $L$ -smooth and  $\mu$ -strongly-convex. If we run Algorithm A1 on  $F(\mathbf{w}) := \frac{1}{N} \sum_{j=1}^N f_j(\mathbf{w})$  with  $\gamma = \Theta\left(\frac{\log(RK)}{\mu RK}\right)$

for  $R$  rounds with synchronization gap  $K$ , it holds that  $\forall i \in [N]$

$$\mathbb{E} \|\hat{\alpha}_i - \alpha_i^*\|^2 \leq \tilde{O} \left( \exp\left(-\frac{T_\alpha}{\kappa_g}\right) + \kappa_g^2 \bar{\zeta}_i(\mathbf{w}^*) L^2 \left( \frac{D^2}{RK} + \frac{\kappa \zeta^2}{\mu^2 R^2} + \frac{\delta^2}{\mu^2 NRK} \right) \right)$$

The implication of Theorem 1 is that even we solve (2) at  $\mathbf{w}^R$ , the algorithm will eventually converge to optimal solution of (2) at  $\mathbf{w}^*$ . **Stage 2: Scalable personalized optimization with model shuffling.** After obtaining the per client mixing parameters, in the second stage to learn personalized models, we aim at solving  $N$  different personalized variants of PERM i.e.,  $\Phi(\hat{\alpha}_1, \nu), \dots, \Phi(\hat{\alpha}_N, \nu)$  where

$$\min_{\nu \in \mathcal{W}} \Phi(\hat{\alpha}_i, \nu) := \frac{1}{N} \sum_{j=1}^N \hat{\alpha}_i(j) f_j(\nu). \quad (3)$$

Here we devise an iterative algorithm based on Local SGD to solve these  $N$  optimization problems in parallel with *no extra overhead*. The idea is to replace the model averaging in vanilla distributed (Local) SGD with *model shuffling*. Specifically, as shown in Algorithm 1 the algorithm proceeds for  $R$  epochs where each epoch runs for  $N$  communication rounds. At the beginning of each epoch  $r$  the server generates a random permutation  $\sigma_r$  over  $N$  clients. At each communication round  $j$  within the epoch, the server sends the model of client  $i$  to client  $i_j = (i+j)\%N$  in the permutation  $\sigma_r$  along with  $\alpha_i(i_j)$ . After receiving a model from server, the client updates the received model for  $K$  local steps and returns it back to the server. As it can be seen, the updates of each loss  $\Phi(\hat{\alpha}_i, \nu), i = 1, 2, \dots, N$  during an epoch is equivalent to sequentially processing individual losses in (3) which can be considered as permutation-based SGD but with the difference that each component now is updated for  $K$  steps. By *interleaving the permutations*, we are able to simultaneously optimize all  $N$  objectives. The following theorem establishes the convergence rate

**THEOREM 2.** Assume each  $f_i$  is  $L$ -smooth,  $\mu$ -strongly-convex and with gradient bounded by  $G$ , i.e.,  $\sup_{\mathbf{w} \in \mathcal{W}} \|\nabla f_i(\mathbf{w})\| \leq G$ . Let  $\alpha_i^*$  be the solution of (2). Then if we run Algorithm 1 on the  $\hat{\alpha}_i$  obtained from Algorithm A1, then Algorithm 1 with  $\eta = \Theta\left(\frac{\log(NKR^3)}{\mu R}\right)$  will output the solution  $\hat{\nu}_i, \forall i \in [N]$ , such that with probability at least  $1 - p$ , the following statement holds:

$$\begin{aligned} & \mathbb{E}[\Phi(\alpha_i^*, \hat{\nu}_i) - \Phi(\alpha_i^*, \nu^*(\alpha_i^*))] \\ & \leq \tilde{O} \left( \frac{D^2 L}{NKR^2} + \frac{L\delta^2}{\mu^2 R} + \left( \frac{L^4 + N}{\mu^4 R^2} \right) LG^2 N \log(1/p) \right) \\ & \quad + \kappa^2 L \tilde{O} \left( \exp\left(-\frac{T_\alpha}{\kappa_g}\right) + \kappa_g^2 \bar{\zeta}_i(\mathbf{w}^*) L^2 \left( \frac{D^2}{RK} + \frac{\kappa \zeta^2}{\mu^2 R^2} + \frac{\delta^2}{\mu^2 NRK} \right) \right), \end{aligned}$$

where the expectation is taken over randomness of Algorithm A1.

The above theorem shows that, even though we run the optimization on  $\Phi(\hat{\alpha}_i, \nu)$ , our obtained model  $\hat{\nu}_i$  will still converge to the optimal solution of  $\Phi(\alpha_i^*, \nu)$ . The convergence rate is contributed from two parts: convergence of  $\hat{\alpha}_i$  (Algorithm A1) and convergence of personalized model  $\hat{\nu}_i$  (Algorithm 1). Notice that, for the convergence rate of  $\hat{\nu}_i$ , we roughly recover the optimal rate of shuffling SGD [1], which is  $O\left(\frac{1}{R^2}\right)$ . However, we suffer from a  $O\left(\frac{\delta^2}{R}\right)$  term since each client runs vanilla SGD on their local data (the SGD-Update procedure in Algorithm 1). A recent work [4] also



**Algorithm 1:** Shuffling Local SGD

**Input:** Clients  $1, \dots, N$ , Number of Local Steps  $K$ , Number of Epoch  $R$ , mixing parameter  $\hat{\alpha}_1, \dots, \hat{\alpha}_N$

**Epoch for**  $r = 0, \dots, R - 1$  **do**

Server generates permutation  $\sigma_r : [N] \mapsto [N]$ .

**parallel for**  $i = 1, \dots, N$  **do**

Client  $i$  sets initial model  $v_i^{r,0} = v_i^r$ .

**for**  $j = 1, \dots, N$  **do**

Set indices  $i_j = \sigma_r((i + j) \bmod N)$ .

Server sends  $v_i^{r,j}$  to Client  $i_j$ .

$v_i^{r,j+1} = \text{SGD-Update}(v_i^{r,j}, \eta, i_j, K, \hat{\alpha}_i)$ .

Client  $i$  does projection:  $v_i^{r+1} = \mathcal{P}_{\mathcal{W}}(v_i^{r,N})$ .

**Output:**  $\hat{v}_i = v_i^R, \forall i \in [N]$ .

**SGD-Update**( $v, \eta, j, K, \alpha$ )

**Initialize**  $v^0 = v$

**for**  $t = 0, \dots, K - 1$  **do**

$v^t = v^{t-1} - \eta \alpha(j) N \nabla f_j(v^{t-1}; \xi_j^{t-1})$

**Output:**  $v^K$

considering client-level shuffling idea. The differences are two-fold: 1) they only consider training one model while we employ shuffling idea to simultaneously learn  $N$  models 2) they work with fixed  $\alpha$  while we have to show that the algorithm can converge to the true optimal solution of  $\Phi(\alpha_i^*, v)$  given that we only optimize on a surrogate function  $\Phi(\hat{\alpha}_i, v)$ .

### 3.2 A unified single loop algorithm

In practice, a single-loop algorithm is preferred since it is easy to implement. We hence turn to introducing a single stage algorithm that jointly optimizes  $\alpha_i$ s and  $v_i$ s as depicted in Algorithm 2 by intertwining the two stages in Algorithm A1 and Algorithm 1 in a single unified method. At each communication round, the clients compute gradients on global model, on their data, after that server collects these gradients and does one step mini-batch SGD update on global model, and then updates mixing parameters. Next we proceed to update the personalized models similar to Algorithm 1. We note that, unlike the two stages method where the mixing parameters are computed at the final global model, here the mixing parameters are updated adaptively based on intermediate global models.

**THEOREM 3.** Assume each  $f_i$  is  $L$ -smooth,  $\mu$ -strongly-convex and with gradient bounded by  $G$ , i.e.,  $\sup_{\mathbf{w} \in \mathcal{W}} \|\nabla f_i(\mathbf{w})\| \leq G$ . Let  $\alpha_i^*$  be the solution of (2). Then if we run Algorithm 2 with  $\eta = \Theta\left(\frac{\log(NKR^3)}{\mu R}\right)$  and  $\gamma = \Theta\left(\frac{\log(NKR^3)}{\mu R}\right)$ , it will output the solution  $\hat{v}_i, \forall i \in [N]$ , such that with probability at least  $1 - p$ , the following statement holds:

$$\begin{aligned} \mathbb{E}[\Phi(\alpha_i^*, \hat{v}_i) - \Phi(\alpha_i^*, v_i^*)] &\leq O\left(\frac{LD^2}{NKR^3}\right) \\ &+ \tilde{O}\left(\left(\frac{\kappa^4 L}{R^2} + \frac{NL}{\mu^2 R^2}\right) G^2 N \log(1/p) + \frac{L\delta^2}{\mu^2 R}\right) \\ &+ \tilde{O}\left(\frac{\kappa^2 \kappa_g^2 L^3 \bar{\xi}_i(\mathbf{w}^*) DG}{R} + LR^2 \left(1 - \frac{1}{\kappa_g}\right)^{T_\alpha} + \frac{L\kappa_g^2 \kappa^2 \bar{\xi}_i(\mathbf{w}^*) \delta^2}{\mu^2 M}\right), \end{aligned}$$

**Algorithm 2:** Single Loop PERM

**Input:** Clients  $1, \dots, N$ , Number of Local Steps  $K$ , Number of Epoch  $R$ , Initial mixing parameter

$\alpha_1^0, \dots, \alpha_N^0 = \bar{\alpha} = [1/N, \dots, 1/N]$ .

**Epoch for**  $r = 0, \dots, R - 1$  **do**

Server generates permutation  $\sigma_r : [N] \mapsto [N]$ .

**parallel for** Client  $i = 1, \dots, N$  **do**

Client  $i$  sets initial model  $v_i^{r,0} = v_i^r$ .

**for**  $j = 1, \dots, N$  **do**

Set indices  $i_j = \sigma_r((i + j) \bmod N)$ .

Server sends  $v_i^{r,j}$  to client  $i_j$ .

$v_i^{r,j+1} = \text{SGD-Update}(v_i^{r,j}, \eta, i_j, K, \alpha_i^r)$ .

// Personalized model update

Client  $i$  does projection:  $v_i^{r+1} = \mathcal{P}_{\mathcal{W}}(v_i^{r,N})$ .

$\mathbf{w}^{r+1} = \mathcal{P}_{\mathcal{W}}(\mathbf{w}^r - \gamma \frac{1}{N} \sum_{i=1}^N \frac{1}{M} \sum_{j=1}^M \nabla f_i(\mathbf{w}^r; \xi_{i,j}^r))$

// Global model update

Compute  $\alpha_i^{r+1}$  by running  $T_\alpha$  steps GD on  $g_i(\mathbf{w}^{r+1}, \alpha)$

//  $\alpha$  update

**Output:**  $\hat{v}_i = v_i^R, \hat{\alpha}_i = \alpha_i^R, \forall i \in [N]$ .

**SGD-Update**( $v, \eta, j, K, \alpha$ )

**Initialize**  $v_j^0 = v$

**for**  $t = 0, \dots, K - 1$  **do**

$v^t = v^{t-1} - \eta \alpha(j) N \nabla f_j(v^{t-1}; \xi_j^{t-1})$

**Output:**  $v^K$

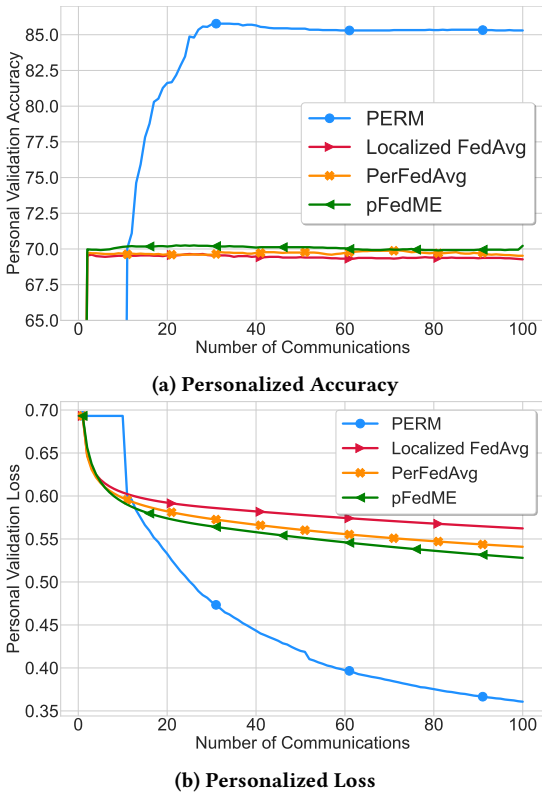
where the expectation is taken over the randomness of stochastic samples in Algorithm 2.

## 4 EXPERIMENT

To demonstrate the superior effectiveness of our proposed single-loop PERM algorithm compared to other existing personalization methods, we conducted an experiment using synthetic data generated according to the following specifications. We consider a scenario with a total of  $N$  clients, where we draw samples from the distribution  $\mathcal{N}(\mu_1, \Sigma_i)$  for half of the clients, denoted by  $i \in [1, \frac{N}{2}]$ , and from  $\mathcal{N}(\mu_2, \Sigma_i)$  for the remaining clients, denoted by  $i \in [\frac{N}{2}, N]$ . Following the approach outlined in [14], we adopt a uniform variate for all samples, with  $\Sigma_{k,k} = k^{-1.2}$ . Subsequently, we generate a labeling model using the distribution  $\mathcal{N}(\mu_w, \Sigma_w)$ .

Given a data sample  $\mathbf{x} \in \mathbb{R}^d$ , the labels are generated as follows: clients  $1, \dots, \frac{N}{2}$  assign labels based on  $y = \text{sign}(\mathbf{w}^\top \mathbf{x})$ , while clients  $\frac{N}{2} + 1, \dots, N$  assign labels based on  $y = \text{sign}(-\mathbf{w}^\top \mathbf{x})$ . For this specific experiment, we set  $\mu_1 = 0.2$ ,  $\mu_2 = -0.2$ , and  $\mu_w = 0.1$ . The data dimension is  $d = 60$ , and there are 2 classes in the output. We have a total of 50 clients, each generating 500 samples following the aforementioned guidelines. We train a logistic regression model on each client's data.

To demonstrate the superiority of our PERM algorithm, we conducted a performance comparison against other prominent personalized approaches, including the fined-tuned model of FedAvg [19] (referred to as localized FedAvg), perFedAg [8], and pFedME [22]. The results in Figure 1 highlight PERM's efficient learning of personalized models for individual clients. In contrast, competing methods



**Figure 1: Comparative Analysis of Personalization methods, including our single-loop PERM algorithm, localized FedAvg, perFedAg, and pFedME, with synthetic data. The disparity in personalized accuracy and loss highlights PERM’s capability in leveraging relevant client correlations.**

relying on globally trained models struggle to match PERM’s effectiveness in highly heterogeneous scenarios, as seen in personalized accuracy and loss. This showcases PERM’s exceptional ability to leverage relevant client learning.

## 5 CONCLUSION

This paper introduces a new *data&system-aware* paradigm for learning from multiple heterogeneous data sources to achieve optimal statistical accuracy across all data distributions without imposing stringent constraints on computational resources shared by participating devices. The proposed PERM schema, though simple, provides an efficient solution to enable each client to learn a personalized model by *learning who to learn with* via personalizing the aggregation of data sources through an efficient empirical statistical discrepancy estimation module. PERM can also be employed in other learning settings with multiple sources of data such as domain adaptation and multi-task learning to entail optimal statistical accuracy. To efficiently solve all aggregated personalized losses, we propose a model shuffling idea to optimizes all losses in parallel. This also enables us to learn models with varying complexity for different devices to meet their available resources.

## REFERENCES

- [1] Kwangjun Ahn, Chulhee Yun, and Suvrit Sra. Sgd with shuffling: optimal rates without component convexity and large epoch requirements. *Advances in Neural Information Processing Systems*, 33:17526–17535, 2020.
- [2] Shai Ben-David, John Blitzer, Koby Crammer, Alex Kulesza, Fernando Pereira, and Jennifer Wortman Vaughan. A theory of learning from different domains. *Machine learning*, 79:151–175, 2010.
- [3] Sebastian Caldas, Peter Wu, Tian Li, Jakub Konečný, H Brendan McMahan, Virginia Smith, and Ameet Talwalkar. Leaf: A benchmark for federated settings. *arXiv preprint arXiv:1812.01097*, 2018.
- [4] Yae Jee Cho, Pranay Sharma, Gauri Joshi, Zheng Xu, Satyen Kale, and Tong Zhang. On the convergence of federated averaging with cyclic client participation. *arXiv preprint arXiv:2302.03109*, 2023.
- [5] Koby Crammer, Michael Kearns, and Jennifer Wortman. Learning from multiple sources. *Journal of Machine Learning Research*, 9(8), 2008.
- [6] Yuyang Deng, Mohammad Mahdi Kamani, and Mehrdad Mahdavi. Adaptive personalized federated learning. *arXiv preprint arXiv:2003.13461*, 2020.
- [7] Mathieu Even, Laurent Massoulié, and Kevin Scaman. On sample optimality in personalized collaborative and federated learning. In *NeurIPS 2022-36th Conference on Neural Information Processing System*, 2022.
- [8] Alireza Fallah, Aryan Mokhtari, and Asuman Ozdaglar. Personalized federated learning: A meta-learning approach. *Advances in Neural Information Processing Systems*, 2020.
- [9] Filip Hanzely, Boxin Zhao, and Mladen Kolar. Personalized federated learning: A unified framework and universal optimization techniques. *arXiv preprint arXiv:2102.09743*, 2021.
- [10] Yutao Huang, Lingyang Chu, Zirui Zhou, Lanjun Wang, Jiangchuan Liu, Jian Pei, and Yong Zhang. Personalized federated learning: An attentive collaboration approach. *arXiv preprint arXiv:2007.03797*, 2020.
- [11] Jakub Konečný, H Brendan McMahan, Felix X Yu, Peter Richtárik, Ananda Theertha Suresh, and Dave Bacon. Federated learning: Strategies for improving communication efficiency. *arXiv preprint arXiv:1610.05492*, 2016.
- [12] Nikola Konstantinov and Christoph Lampert. Robust learning from untrusted sources. In *International conference on machine learning*, pages 3488–3498. PMLR, 2019.
- [13] Tony Lancaster. The incidental parameter problem since 1948. *Journal of econometrics*, 95(2):391–413, 2000.
- [14] Tian Li, Anit Kumar Sahu, Manzil Zaheer, Maziar Sanjabi, Ameet Talwalkar, and Virginia Smith. Federated optimization in heterogeneous networks. *arXiv preprint arXiv:1812.06127*, 2018.
- [15] Tianyi Lin, Chi Jin, and Michael I Jordan. On gradient descent ascent for nonconvex-concave minimax problems. *arXiv preprint arXiv:1906.00331*, 2019.
- [16] Jie Ma, Guodong Long, Tianyi Zhou, Jing Jiang, and Chengqi Zhang. On the convergence of clustered federated learning. *arXiv preprint arXiv:2202.06187*, 2022.
- [17] Yishay Mansour, Mehryar Mohri, Jae Ro, and Ananda Theertha Suresh. Three approaches for personalization with applications to federated learning. *arXiv preprint arXiv:2002.10619*, 2020.
- [18] Yishay Mansour and Mariano Schain. Robust domain adaptation. *Annals of Mathematics and Artificial Intelligence*, 71(4):365–380, 2014.
- [19] Brendan McMahan, Eider Moore, Daniel Ramage, Seth Hampson, and Blaise Agueria y Arcas. Communication-efficient learning of deep networks from decentralized data. In *Artificial intelligence and statistics*, pages 1273–1282. PMLR, 2017.
- [20] Markus Schneider. Probability inequalities for kernel embeddings in sampling without replacement. In *Artificial Intelligence and Statistics*, pages 66–74. PMLR, 2016.
- [21] Sebastian U Stich. Local sgd converges fast and communicates little. In *International Conference on Learning Representations*.
- [22] Canh T Dinh, Nguyen Tran, and Josh Nguyen. Personalized federated learning with moreau envelopes. *Advances in Neural Information Processing Systems*, 33:21394–21405, 2020.
- [23] Jianyu Wang, Rudrajit Das, Gauri Joshi, Satyen Kale, Zheng Xu, and Tong Zhang. On the unreasonable effectiveness of federated averaging with heterogeneous data. *arXiv preprint arXiv:2206.04723*, 2022.
- [24] Blake Woodworth, Jialei Wang, Adam Smith, Brendan McMahan, and Nathan Srebro. Graph oracle models, lower bounds, and gaps for parallel stochastic optimization. *arXiv preprint arXiv:1805.10222*, 2018.
- [25] Blake E Woodworth, Kumar Kshitij Patel, and Nati Srebro. Minibatch vs local sgd for heterogeneous distributed learning. *Advances in Neural Information Processing Systems*, 33:6281–6292, 2020.
- [26] Zi Xu, Huiling Zhang, Yang Xu, and Guanghui Lan. A unified single-loop alternating gradient projection algorithm for nonconvex-concave and convex-nonconcave minimax problems. *Mathematical Programming*, pages 1–72, 2023.

## Supplementary Material

The supplementary material is organized as follows:

- In Section A, we provide the proof of convergence of two-stage implementation of PERM (computing mixing parameters followed by learning personalized models via model shuffling using permutation-based variant of distributed SGD with periodic communication). The technical tools we develop to establish the convergence could be interesting by their own right.
- In Section B, we provide the proof of convergence of single loop implementation of PERM, where estimation of mixing parameters are coupled with learning of personalized models.
- In Section C, we report additional experiments on performance of PERM on EMNIST dataset along with effectiveness of proposed method to estimate mixing parameters.

### A PROOF OF TWO STAGES ALGORITHM

#### A.1 Technical Lemmas

LEMMA 1. Define  $v^*(\alpha) := \arg \min_{v \in \mathcal{W}} \Phi(\alpha, v)$ , and assume  $\Phi(\alpha, \cdot)$  is  $\mu$ -strongly convex and  $L$  smooth. Then,  $v^*(\cdot)$  is  $\kappa$ -Lipschitz.

PROOF. The proof is similar to Lin et al's result on minimax objective [15]. First, according to optimality conditions we have:

$$\begin{aligned} \langle v - v^*(\alpha), \nabla_v \Phi(\alpha, v^*(\alpha)) \rangle &\geq 0, \\ \langle v - v^*(\alpha'), \nabla_v \Phi(\alpha', v^*(\alpha')) \rangle &\geq 0 \end{aligned}$$

Substituting  $v$  with  $v^*(\alpha')$  and  $v^*(\alpha)$  in the above first and second inequalities respectively yields:

$$\begin{aligned} \langle v^*(\alpha') - v^*(\alpha), \nabla_v \Phi(\alpha, v^*(\alpha)) \rangle &\geq 0, \\ \langle v^*(\alpha) - v^*(\alpha'), \nabla_v \Phi(\alpha', v^*(\alpha')) \rangle &\geq 0 \end{aligned}$$

Adding up the above two inequalities yields:

$$\langle v^*(\alpha') - v^*(\alpha), \nabla_v \Phi(\alpha, v^*(\alpha)) - \nabla_v \Phi(\alpha', v^*(\alpha')) \rangle \geq 0, \quad (4)$$

Since  $\Phi(\alpha, \cdot)$  is  $\mu$  strongly convex, we have:

$$\langle v^*(\alpha') - v^*(\alpha), \nabla_y \Phi(\alpha, v^*(\alpha')) - \nabla_y \Phi(\alpha, v^*(\alpha)) \rangle \geq \mu \|v^*(\alpha') - v^*(\alpha)\|^2. \quad (5)$$

Adding up (4) and (5) yields:

$$\langle v^*(\alpha') - v^*(\alpha), \nabla_y \Phi(\alpha, v^*(\alpha')) - \nabla_y \Phi(\alpha', v^*(\alpha')) \rangle \geq \mu \|v^*(\alpha') - v^*(\alpha)\|^2$$

Finally, using  $L$  smoothness of  $\Phi$  will conclude the proof:

$$\begin{aligned} L \|v^*(\alpha') - v^*(\alpha)\| \|\alpha - \alpha'\| &\geq \mu \|v^*(\alpha') - v^*(\alpha)\|^2 \\ \iff \kappa \|\alpha - \alpha'\| &\geq \|v^*(\alpha') - v^*(\alpha)\| \end{aligned}$$

□

LEMMA 2 (OPTIMALITY GAP). Let  $\Phi(\alpha, v)$  be defined in (3). If we assume each  $f_i$  is  $L$ -smooth and  $\mu$ -strongly convex, then the following statement holds true:

$$\Phi(\alpha_i^*, \hat{v}_i) - \Phi(\alpha_i^*, v_i^*) \leq L \|\hat{v}_i - v^*(\hat{\alpha}_i)\|^2 + \kappa^2 L \|\hat{\alpha}_i - \alpha_i^*\|^2,$$

where  $v_i^* = \arg \min_{v \in \mathcal{W}} \Phi(\alpha_i^*, v)$ .

PROOF. We notice the following fact:

$$\begin{aligned} \Phi(\alpha_i^*, \hat{v}_i) - \Phi(\alpha_i^*, v_i^*) &\leq \underbrace{\langle \nabla_v \Phi(\alpha_i^*, v_i^*), \hat{v}_i - v_i^* \rangle}_{\leq 0} + \frac{L}{2} \|\hat{v}_i - v_i^*\|^2 \\ &\leq \frac{L}{2} \|\hat{v}_i - v_i^*\|^2. \end{aligned}$$

**Algorithm A1:** Discrepancy Estimation at Optimum**Input:** Number of clients  $N$ , number of local steps  $K$ , number of communications rounds  $R$ **for**  $r = 0, \dots, R - 1$  **do**    **parallel for** client  $i = 0, \dots, N - 1$  **do**        Client  $i$  initializes model  $\mathbf{w}_i^{r,0} = \mathbf{w}_i^r$ .        **for**  $t = 0, \dots, K - 1$  **do**             $\mathbf{w}_i^{r,t+1} = \mathbf{w}_i^{r,t} - \gamma \nabla f_i(\mathbf{w}_i^{r,t}; \xi_i^{r,t})$  where  $\xi_i^{r,t}$  is a mini-batch sampled from  $S_i$ .        Client  $i$  sends  $\mathbf{w}_i^{r,K}$  to Server.    Server computes  $\mathbf{w}^{r+1} = \mathcal{P}_{\mathcal{W}} \left( \frac{1}{N} \sum_{i=1}^N \mathbf{w}_i^{r,K} \right)$     Server broadcasts  $\mathbf{w}^{r+1}$  to all clients.Server computes  $\hat{\alpha}_i, i = 1, 2, \dots, N$  by running  $T_\alpha$  steps of GD on  $g_i(\mathbf{w}^R, \alpha)$ .**Output:**  $\hat{\alpha}_1, \dots, \hat{\alpha}_N$ .

where we use the smoothness of  $\Phi$  and optimality condition of  $\mathbf{v}_i^*$  for fixed  $\hat{\alpha}_i$ . Next, due to the  $\kappa$ -Lipschitzness property of  $\mathbf{v}^*(\cdot)$  as shown in Lemma 1, it follows that:

$$\begin{aligned} \frac{L}{2} \|\hat{\mathbf{v}}_i - \mathbf{v}_i^*\|^2 &\leq L \|\hat{\mathbf{v}}_i - \mathbf{v}^*(\hat{\alpha}_i)\|^2 + L \|\mathbf{v}^*(\hat{\alpha}_i) - \mathbf{v}_i^*\|^2 \\ &\leq L \|\hat{\mathbf{v}}_i - \mathbf{v}^*(\hat{\alpha}_i)\|^2 + \kappa^2 L \|\hat{\alpha}_i - \alpha_i^*\|^2. \end{aligned}$$

as desired.  $\square$ **A.2 Proof of Convergence of Theorem 1**

In this section we are going to prove the result in Theorem 1. To this end, we need to show that mixing parameters we compute by first learning the global model and then solving the optimization problem in objective (2) (as depicted in Algorithm A1) converges to optimal values. Notice that in Algorithm A1 we do not solve  $g_i(\mathbf{w}^*, \alpha)$  directly, but optimize  $g_i(\mathbf{w}^R, \alpha)$  on  $\alpha$  for  $T_\alpha$  iterations of GD. Hence, firstly we need to show that optimizing the surrogate function will also guarantee the convergence of output of algorithm  $\hat{\alpha}$  to  $\alpha^*$  by deriving a property of the objective in (2). Formally the property is captured by the following lemma.

**LEMMA 3.** Let  $g(\mathbf{w}, \alpha) := \sum_{j=1}^N \alpha_j \|\nabla f_j(\mathbf{w}) - \nabla f_j(\mathbf{w}^*)\|^2 + \lambda \sum_{j=1}^N \alpha_j^2 / n_j$  and  $\alpha_g^*(\mathbf{w}) = \arg \min_{\alpha \in \Delta^N} g(\mathbf{w}, \alpha)$ . Let  $\mathbf{w}^R$  be the output of Algorithm A1. Then the following statement holds:

$$\left\| \alpha_g^*(\mathbf{w}^R) - \alpha_g^*(\mathbf{w}^*) \right\| \leq \kappa_g^2 \sum_{j=1}^N \left( 2 \|\nabla f_j(\mathbf{w}^*) - \nabla f_j(\mathbf{w}^*)\|^2 + 4L^2 \|\mathbf{w}^R - \mathbf{w}^*\|^2 \right) 4L \|\mathbf{w}^R - \mathbf{w}^*\|^2$$

where  $\kappa_g$  is the condition number of  $g$ .**PROOF.** Define function

$$W(\mathbf{z}, \alpha) = \sum_{j=1}^N \alpha_j z_j + \lambda \sum_{j=1}^N \alpha_j^2 / n_j \quad (6)$$

Apparently,  $W(\mathbf{z}, \alpha)$  is linear in  $\mathbf{z}$  and strongly convex in  $\alpha$ . Then, according to Proposition 1,  $\alpha_W^*(\mathbf{z}) := \arg \min_{\alpha \in \Delta^N} W(\mathbf{z}, \alpha)$  is  $\kappa_g$  lipschitz in  $\mathbf{z}$  where  $\kappa_g = \frac{n_{\max}}{n_{\min}}$ , i.e.,  $\|\alpha_W^*(\mathbf{z}) - \alpha_W^*(\mathbf{z}')\| \leq \kappa_g \|\mathbf{z} - \mathbf{z}'\|$ . Now, let us consider the objective (2):

$$g(\mathbf{w}, \alpha) := \sum_{j=1}^N \alpha_j \|\nabla f_j(\mathbf{w}) - \nabla f_j(\mathbf{w}^*)\|^2 + \lambda \sum_{j=1}^N \alpha_j^2 / n_j$$

We define  $\alpha_g^*(\mathbf{w}) = \arg \min_{\alpha \in \Delta^N} g(\mathbf{w}, \alpha)$ .

We set

$$\begin{aligned} \mathbf{z}^R &= \left[ \|\nabla f_1(\mathbf{w}^R) - \nabla f_1(\mathbf{w}^*)\|^2, \dots, \|\nabla f_N(\mathbf{w}^R) - \nabla f_N(\mathbf{w}^*)\|^2 \right], \\ \mathbf{z}^* &= \left[ \|\nabla f_1(\mathbf{w}^*) - \nabla f_1(\mathbf{w}^*)\|^2, \dots, \|\nabla f_N(\mathbf{w}^*) - \nabla f_N(\mathbf{w}^*)\|^2 \right]. \end{aligned}$$

Then we know that

$$\left\| \alpha_g^*(\mathbf{w}^R) - \alpha_g^*(\mathbf{w}^*) \right\|^2 = \left\| \alpha_w^*(z^R) - \alpha_w^*(z^*) \right\|^2 \leq \kappa_g^2 \left\| z^R - z^* \right\|^2 \quad (7)$$

$$\leq \kappa_g^2 \sum_{j=1}^N \left| \left\| \nabla f_j(\mathbf{w}^R) - \nabla f_j(\mathbf{w}^*) \right\|^2 - \left\| \nabla f_j(\mathbf{w}^*) - \nabla f_j(\mathbf{w}^*) \right\|^2 \right|^2 \quad (8)$$

$$\begin{aligned} &\leq \kappa_g^2 \sum_{j=1}^N \left| \left( \nabla f_j(\mathbf{w}^R) - \nabla f_j(\mathbf{w}^*) + \nabla f_j(\mathbf{w}^*) - \nabla f_j(\mathbf{w}^*) \right) \right. \\ &\quad \left. \times \left( \nabla f_j(\mathbf{w}^R) - \nabla f_j(\mathbf{w}^*) - \nabla f_j(\mathbf{w}^*) + \nabla f_j(\mathbf{w}^*) \right) \right|^2 \\ &\leq \kappa_g^2 \sum_{j=1}^N \left\| \nabla f_j(\mathbf{w}^R) - \nabla f_j(\mathbf{w}^*) + \nabla f_j(\mathbf{w}^*) - \nabla f_j(\mathbf{w}^*) \right\|^2 4L^2 \left\| \mathbf{w}^R - \mathbf{w}^* \right\|^2 \end{aligned}$$

Since  $\left\| \nabla f_j(\mathbf{w}^R) - \nabla f_j(\mathbf{w}^*) \right\| \leq \left\| \nabla f_j(\mathbf{w}^*) - \nabla f_j(\mathbf{w}^*) \right\| + 2L \left\| \mathbf{w}^R - \mathbf{w}^* \right\|$ , we can conclude that

$$\left\| \alpha_g^*(\mathbf{w}^R) - \alpha_g^*(\mathbf{w}^*) \right\| \leq \kappa_g^2 \sum_{j=1}^N \left( 2 \left\| \nabla f_j(\mathbf{w}^*) - \nabla f_j(\mathbf{w}^*) \right\|^2 + 4L^2 \left\| \mathbf{w}^R - \mathbf{w}^* \right\|^2 \right) 4L \left\| \mathbf{w}^R - \mathbf{w}^* \right\|^2.$$

□

With above lemma, to show the convergence of  $\hat{\alpha}$  to  $\alpha^*$ , we do the following decomposition

$$\begin{aligned} \left\| \hat{\alpha} - \alpha^* \right\|^2 &\leq 2 \left\| \hat{\alpha} - \alpha_g^*(\mathbf{w}^R) \right\|^2 + 2 \left\| \alpha_g^*(\mathbf{w}^R) - \alpha_g^*(\mathbf{w}^*) \right\|^2 \\ &\leq 2(1 - \mu\eta\alpha)^K + 2\kappa_g^2 \sum_{j=1}^N \left( 2 \left\| \nabla f_j(\mathbf{w}^*) - \nabla f_j(\mathbf{w}^*) \right\|^2 + 4L^2 \left\| \mathbf{w}^R - \mathbf{w}^* \right\|^2 \right) 4L \left\| \mathbf{w}^R - \mathbf{w}^* \right\|^2. \end{aligned}$$

Now it remains to show the convergence of Local SGD *last iterate*  $\mathbf{w}^R$  to optimal solution  $\mathbf{w}^*$ . By convention, we use  $\mathbf{w}^t = \frac{1}{N} \sum_{i=1}^N \mathbf{w}_i^t$  to denote the virtual average iterates.

LEMMA 4 (ONE ITERATION ANALYSIS OF LOCAL SGD). *Under the condition of Theorem 1, the following statement holds true for any  $t \in [T]$ :*

$$\begin{aligned} \mathbb{E} \left\| \mathbf{w}^{r,t+1} - \mathbf{w}^* \right\|^2 &\leq (1 - \mu\gamma) \mathbb{E} \left\| \mathbf{w}^{r,t} - \mathbf{w}^* \right\|^2 - (2\gamma - 4\gamma^2 L) \mathbb{E} (F(\mathbf{w}^*) - F(\mathbf{w}^{r,t})) \\ &\quad + (\gamma L + 2\gamma^2 L^2) \frac{1}{N} \sum_{j=1}^N \left\| \mathbf{w}_j^{r,t} - \mathbf{w}^{r,t} \right\|^2 + \gamma^2 \frac{\delta^2}{N}. \end{aligned}$$

PROOF. According to updating rule in Algorithm A1, we have the following identity:

$$\mathbb{E} \left\| \mathbf{w}^{r,t+1} - \mathbf{w}^* \right\|^2 = \mathbb{E} \left\| \mathbf{w}^{r,t} - \mathbf{w}^* \right\|^2 - 2\gamma \mathbb{E} \left\langle \frac{1}{N} \sum_{j=1}^N \nabla f_j(\mathbf{w}_j^{r,t}; z_j^{r,t}), \mathbf{w}^{r,t} - \mathbf{w}^* \right\rangle + \gamma^2 \mathbb{E} \left\| \frac{1}{N} \sum_{j=1}^N \nabla f_j(\mathbf{w}_j^{r,t}) \right\|^2 \quad (9)$$

$$= \underbrace{\mathbb{E} \left\| \mathbf{w}^t - \mathbf{w}^* \right\|^2 - 2\gamma \left\langle \frac{1}{N} \sum_{j=1}^N \nabla f_j(\mathbf{w}_j^{r,t}), \mathbf{w}^{r,t} - \mathbf{w}^* \right\rangle}_{T_1} + \underbrace{\gamma^2 \mathbb{E} \left\| \frac{1}{N} \sum_{j=1}^N \nabla f_j(\mathbf{w}_j^{r,t}; z_j^{r,t}) \right\|^2}_{T_2} + \gamma^2 \frac{\delta^2}{N}. \quad (10)$$

For  $T_1$ , since each  $f_j$  is  $L$  smooth and  $\mu$  strongly convex, we have:

$$\begin{aligned} -2\gamma \left\langle \frac{1}{N} \sum_{j=1}^N \nabla f_j(\mathbf{w}_j^{r,t}), \mathbf{w}^t - \mathbf{w}^* \right\rangle &= -2\gamma \left\langle \frac{1}{N} \sum_{j=1}^N \nabla f_j(\mathbf{w}_j^t), \mathbf{w}^t - \mathbf{w}_j^t + \mathbf{w}_j^t - \mathbf{w}^* \right\rangle \\ &\leq -2\gamma \left\langle \frac{1}{N} \sum_{j=1}^N \nabla f_j(\mathbf{w}_j^t), \mathbf{w}^{r,t} - \mathbf{w}_j^{r,t} + \mathbf{w}_j^{r,t} - \mathbf{w}^* \right\rangle \\ &\leq 2\gamma \frac{1}{N} \sum_{j=1}^N \left( f_j(\mathbf{w}^*) - f_j(\mathbf{w}^{r,t}) - \frac{\mu}{2} \left\| \mathbf{w}_j^{r,t} - \mathbf{w}^* \right\|^2 + \frac{L}{2} \left\| \mathbf{w}_j^{r,t} - \mathbf{w}^{r,t} \right\|^2 \right). \end{aligned}$$



Due to Jensen's inequality we know:  $-\frac{1}{N} \sum_{j=1}^N \frac{\mu}{2} \|\mathbf{w}_i^t - \mathbf{w}^*\|^2 \leq -\frac{\mu}{2} \|\mathbf{w}^t - \mathbf{w}^*\|^2$ . Hence we know:

$$-2\gamma \left\langle \frac{1}{N} \sum_{j=1}^N \nabla f_j(\mathbf{w}_i^{r,t}), \mathbf{w}^{r,t} - \mathbf{w}^* \right\rangle \leq 2\gamma \left( F(\mathbf{w}^*) - F(\mathbf{w}^{r,t}) - \frac{\mu}{2} \|\mathbf{w}^{r,t} - \mathbf{w}^*\|^2 + \frac{L}{2} \frac{1}{N} \sum_{j=1}^N \|\mathbf{w}_i^{r,t} - \mathbf{w}^{r,t}\|^2 \right).$$

For  $T_2$ , we have:

$$\begin{aligned} \mathbb{E} \left\| \frac{1}{N} \sum_{j=1}^N \nabla f_j(\mathbf{w}_i^{r,t}) \right\|^2 &= 2\mathbb{E} \left\| \frac{1}{N} \sum_{j=1}^N \nabla f_j(\mathbf{w}_i^{r,t}) - \nabla F(\mathbf{w}^{r,t}) \right\|^2 + 2\mathbb{E} \|\nabla F(\mathbf{w}^{r,t})\|^2 \\ &\leq 2L^2 \frac{1}{N} \sum_{j=1}^N \mathbb{E} \|\mathbf{w}_i^{r,t} - \mathbf{w}^{r,t}\|^2 + 4L (F(\mathbf{w}^{r,t}) - F(\mathbf{w}^*)). \end{aligned}$$

Now, plugging  $T_1$  and  $T_2$  back to (10) yields:

$$\begin{aligned} \mathbb{E} \|\mathbf{w}^{r,t+1} - \mathbf{w}^*\|^2 &\leq (1 - \mu\gamma) \mathbb{E} \|\mathbf{w}^{r,t} - \mathbf{w}^*\|^2 - (2\gamma - 4\gamma^2 L) \mathbb{E} (F(\mathbf{w}^*) - F(\mathbf{w}^{r,t})) \\ &\quad + (\gamma L + 2\gamma^2 L^2) \frac{1}{N} \sum_{j=1}^N \mathbb{E} \|\mathbf{w}_i^{r,t} - \mathbf{w}^{r,t}\|^2 + \gamma^2 \frac{\delta^2}{N}. \end{aligned}$$

□

LEMMA 5. [25, Lemma 8] For the iterates  $\{\mathbf{w}_i^{r,t}\}$  generated in Algorithm A1, the following statement holds true:

$$\frac{1}{N} \sum_{j=1}^N \mathbb{E} \|\mathbf{w}_i^{r,t} - \mathbf{w}^{r,t}\|^2 \leq 3K\gamma^2 \delta^2 + 6K^2 \gamma^2 \zeta^2.$$

LEMMA 6 (LAST ITERATE CONVERGENCE OF LOCAL SGD). Under the conditions of Theorem 1, the following statement holds true for the iterates in Algorithm A1:

$$\mathbb{E} \|\mathbf{w}^R - \mathbf{w}^*\|^2 \leq (1 - \mu\gamma)^{RK} \mathbb{E} \|\mathbf{w}^0 - \mathbf{w}^*\|^2 + \frac{1}{\mu\gamma} (\gamma L + 2\gamma^2 L^2) (3K\gamma^2 \delta^2 + 6K^2 \gamma^2 \zeta^2) + \frac{\gamma \delta^2}{\mu N}$$

PROOF. We first unroll the recursion in Lemma 4 from  $t = K$  to 0, within one communication round:

$$\begin{aligned} \mathbb{E} \|\mathbf{w}^{r,K} - \mathbf{w}^*\|^2 &= (1 - \mu\gamma)^K \mathbb{E} \|\mathbf{w}^{r,0} - \mathbf{w}^*\|^2 - \sum_{t=0}^{K-1} (1 - \mu\gamma)^{K-t} (2\gamma - 4\gamma^2 L) \mathbb{E} (F(\mathbf{w}^*) - F(\mathbf{w}^{r,t})) \\ &\quad + \sum_{t=0}^{K-1} (1 - \mu\gamma)^{K-t} (\gamma L + 2\gamma^2 L^2) \frac{1}{N} \sum_{j=1}^N \mathbb{E} \|\mathbf{w}_i^{r,t} - \mathbf{w}^{r,t}\|^2 + \sum_{t=0}^{K-1} (1 - \mu\gamma)^{K-t} \gamma^2 \frac{\delta^2}{N} \end{aligned}$$

Since we choose  $\gamma \leq \frac{1}{2L}$ , we know  $\sum_{t=0}^{K-1} (1 - \mu\gamma)^{K-t} (2\gamma - 4\gamma^2 L) \mathbb{E} (F(\mathbf{w}^*) - F(\mathbf{w}^{r,t})) \geq 0$ . Plugging in Lemma 5 yields:

$$\mathbb{E} \|\mathbf{w}^R - \mathbf{w}^*\|^2 = (1 - \mu\gamma)^{RK} \mathbb{E} \|\mathbf{w}^0 - \mathbf{w}^*\|^2 + \frac{1}{\mu\gamma} (\gamma L + 2\gamma^2 L^2) (3K\gamma^2 \delta^2 + 6K^2 \gamma^2 \zeta^2) + \frac{\gamma \delta^2}{\mu N},$$

Plugging in  $\gamma = \frac{\log(RK)}{\mu RK}$  gives the convergence rate:

$$\mathbb{E} \|\mathbf{w}^R - \mathbf{w}^*\|^2 \leq \tilde{O} \left( \frac{\mathbb{E} \|\mathbf{w}^0 - \mathbf{w}^*\|^2}{RK} + \kappa \left( \frac{\delta^2}{\mu^2 R^2 K} + \frac{\zeta^2}{\mu^2 R^2} \right) + \frac{\delta^2}{\mu^2 NRK} \right),$$

which concludes the proof. □

Equipped with above results, we are now ready to provide the convergence of main theorem.

PROOF OF THEOREM 1. The proof simply follows from Lemma 3:

$$\begin{aligned} \|\hat{\boldsymbol{\alpha}} - \boldsymbol{\alpha}^*\|^2 &\leq 2 \|\hat{\boldsymbol{\alpha}} - \boldsymbol{\alpha}_g^*(\mathbf{w}^R)\|^2 + 2 \|\boldsymbol{\alpha}_g^*(\mathbf{w}^R) - \boldsymbol{\alpha}_g^*(\mathbf{w}^*)\|^2 \\ &\leq 2(1 - \mu\eta\alpha)^{T\alpha} + 8L\kappa_g^2 \sum_{j=1}^N \left( 2 \|\nabla f_i(\mathbf{w}^*) - \nabla f_j(\mathbf{w}^*)\|^2 + 4L^2 \|\mathbf{w}^R - \mathbf{w}^*\|^2 \right) \|\mathbf{w}^R - \mathbf{w}^*\|^2 \\ &\leq 2(1 - \mu\eta\alpha)^{T\alpha} + 8L\kappa_g^2 \left( 2\bar{\zeta}_i(\mathbf{w}^*) + 4NL^2 \|\mathbf{w}^R - \mathbf{w}^*\|^2 \right) \|\mathbf{w}^R - \mathbf{w}^*\|^2 \end{aligned}$$

**Algorithm A2:** Shuffling Local SGD (One Client)

**Input:** Clients  $0, \dots, N-1$ , Number of Local Steps  $K$ , Number of Epoch  $R$ , Mixing parameter  $\hat{\alpha}$

**Epoch for**  $r = 0, \dots, R-1$  **do**

Server generates permutation  $\sigma_r : [N] \mapsto [N]$ .

Client sets initial model  $\mathbf{v}^{r,0} = \mathbf{v}^r$ .

**for**  $j = 0, \dots, N-1$  **do**

Server sends  $\mathbf{v}^{r,j}$  to Client  $\sigma_r(j)$ .

$\mathbf{v}^{r,j+1} = \text{SGD-Update}(\mathbf{v}^{r,j}, \eta, \sigma_r(j), K, \hat{\alpha})$ .

Client  $i$  does projection:  $\mathbf{v}^{r+1} = \mathcal{P}_{\mathcal{W}}(\mathbf{v}^{r,N})$ .

**Output:**  $\hat{\mathbf{v}} = \mathbf{v}^R$ .

**SGD-Update**( $\mathbf{v}, \eta, j, K, \alpha$ )

**Initialize**  $\mathbf{v}^0 = \mathbf{v}$

**for**  $t = 0, \dots, K-1$  **do**

$\mathbf{v}^t = \mathbf{v}^{t-1} - \eta\alpha(j)N\nabla f_j(\mathbf{v}^{t-1}; \xi^{t-1})$

**Output**  $\mathbf{v}^K$

Plugging in the convergence of  $\|\mathbf{w}^R - \mathbf{w}^*\|^2$  from Lemma 6, and the stepsize  $\eta\alpha = \frac{1}{L_g}$  for  $\alpha$  yields:

$$\mathbb{E}\|\alpha_i^R - \alpha_i^*\|^2 \leq \bar{O}\left(\exp\left(-\frac{T\alpha}{\kappa_g}\right) + \kappa_g^2 \bar{\zeta}_i(\mathbf{w}^*)L^2\left(\frac{D^2}{RK} + \kappa\left(\frac{\delta^2}{\mu^2 R^2 K} + \frac{\zeta^2}{\mu^2 R^2}\right) + \frac{\delta^2}{\mu^2 NRK}\right)\right).$$

□

### A.3 Proof of Convergence of Shuffling Local SGD

In this section, we are going to prove the convergence of proposed shuffled variant of Local SGD (Theorem 2). The whole proof framework follows the analysis of vanilla shuffling SGD, but notice that there are two differences. First, in vanilla shuffling SGD, in each epoch, algorithm only updates on each component function  $f_j$  once, while here we have to take  $K$  steps of SGD update on each component function. Second, we are considering a weighted sum objective in contrary to averaged objective in [4], which means we need to rescale the objective when we apply without-replacement concentration inequality. Even though our algorithm solves models for  $N$  clients, for the sake of simplicity, throughout the proof we only show the convergence of one client's model. The algorithm from one client point of view is described in Algorithm A2, where we drop the client index for notational convenience.

**PROPOSITION 1.** Assume a sequence  $\{\mathbf{w}^t\}_{t=1}^K$  is obtained by

$$\mathbf{w}^t = \mathbf{w}^{t-1} - \eta\alpha N\nabla f(\mathbf{w}^{t-1}; \xi^{t-1}), \quad t = 1, \dots, K,$$

then we have

$$\mathbf{w}^{t+1} = \mathbf{w}^0 - \left(\sum_{\tau=0}^t \prod_{t'=t}^{\tau+1} (\mathbf{I} - \alpha N \eta \mathbf{H}_{t'})\right) \eta\alpha N \nabla f(\mathbf{w}^0) - \sum_{\tau=0}^t \prod_{t'=t}^{\tau+1} (\mathbf{I} - \alpha N \eta \mathbf{H}_{t'}) \eta\alpha N \delta^t, \quad \forall 0 \leq t \leq K-1,$$

where  $\delta^t := \nabla f(\mathbf{w}^t; \xi^t) - \nabla f(\mathbf{w}^t)$ , and by convention, we define  $\prod_{j=a}^b \mathbf{A}_j = \mathbf{I}$  if  $a < b$ .

**PROOF.** According to updating rule, we have:

$$\begin{aligned} \mathbf{w}^{t+1} - \mathbf{w}^0 &= \mathbf{w}^t - \mathbf{w}^0 - \eta\alpha N \nabla f(\mathbf{w}^t; \xi^t) \\ &= \mathbf{w}^t - \mathbf{w}^0 - \eta\alpha N \nabla f(\mathbf{w}^t) - \eta\alpha N \delta^t \\ &= \mathbf{w}^t - \mathbf{w}^0 - \eta\alpha N \nabla f(\mathbf{w}^0) - \eta\alpha N (\nabla f(\mathbf{w}^t) - \nabla f(\mathbf{w}^0)) - \eta\alpha N \delta^t. \end{aligned}$$

Since  $f$  is  $L$  smooth, and according to Mean Value Theorem, there is a matrix  $\mathbf{H}_t$  satisfying  $\mu \mathbf{I} \leq \mathbf{H}_t \leq L \mathbf{I}$ , such that  $\nabla f(\mathbf{w}^t) - \nabla f(\mathbf{w}^0) = \mathbf{H}_t(\mathbf{w}^t - \mathbf{w}^0)$ . Hence we have:

$$\mathbf{w}^{t+1} - \mathbf{w}^0 = (\mathbf{I} - \eta\alpha N \mathbf{H}_t)(\mathbf{w}^t - \mathbf{w}^0) - \eta\alpha N \nabla f(\mathbf{w}^0) - \eta\alpha N \delta^t.$$

Unrolling the recursion from  $t$  to 0 will conclude the proof. □

The following lemma establishes the updating rule of models between epochs  $r$  and  $r+1$ . For notational convenience, whenever there is no confusion, we drop the superscript  $r$  in  $\sigma^r$ .

LEMMA 7 (ONE EPOCH UPDATING RULE). Let  $\mathbf{v}^r$  and  $\mathbf{v}^{r+1}$  be two iterates generated by Shuffling Local SGD (Algorithm A2), then the following updating rule holds:

$$\mathbf{v}^{r+1} = \mathbf{v}^r - \sum_{j=1}^N \prod_{j'=N}^{j+1} (\mathbf{I} - \mathbf{Q}_{j'} \mathbf{H}_{j'}) (\mathbf{Q}_j \nabla f_{\sigma(j)}(\mathbf{v}^r) - \boldsymbol{\delta}_j),$$

where

$$\mathbf{Q}_j := \left( \sum_{\tau=0}^{K-1} \prod_{t'=K-1}^{\tau+1} (\mathbf{I} - \eta \hat{\alpha}(\sigma(j)) N \mathbf{H}_{t'}) \right) \eta \hat{\alpha}(\sigma(j)) N,$$

$$\boldsymbol{\delta}_j := \sum_{\tau=0}^{K-1} \prod_{t'=t}^{\tau+1} (\mathbf{I} - \hat{\alpha}(\sigma(j)) N \eta \mathbf{H}_{t'}) \eta \hat{\alpha}(\sigma(j)) N \boldsymbol{\delta}_{\sigma(j)}^t,$$

by convention, we define  $\prod_{j=a}^b \mathbf{A}_j = \mathbf{I}$  if  $a < b$ .

PROOF. According to Proposition 1, we have

$$\mathbf{v}^{r,j+1} = \mathbf{v}^{r,j} - \left( \sum_{\tau=0}^{K-1} \prod_{t'=t}^{\tau+1} (\mathbf{I} - \hat{\alpha}(\sigma(j)) N \eta \mathbf{H}_{t'}) \right) \eta \hat{\alpha}(\sigma(j)) N \nabla f_{\sigma(j)}(\mathbf{v}^{r,j})$$

$$- \sum_{\tau=0}^{K-1} \prod_{t'=t}^{\tau+1} (\mathbf{I} - \hat{\alpha}(\sigma(j)) N \eta \mathbf{H}_{t'}) \eta \hat{\alpha}(\sigma(j)) N \boldsymbol{\delta}_{\sigma(j)}^t.$$

Plugging our definition of  $\mathbf{Q}_j$  and  $\boldsymbol{\delta}_j$  yields:

$$\mathbf{v}^{r,j+1} - \mathbf{v}^r = \mathbf{v}^{r,j} - \mathbf{v}^r - \mathbf{Q}_j \nabla f_{\sigma(j)}(\mathbf{v}^{r,j}) - \boldsymbol{\delta}_j.$$

Following the same reasoning in the proof of Proposition 1 will conclude the proof.  $\square$

LEMMA 8 (SUMMATION BY PARTS). Let  $\mathbf{A}_j$  and  $\mathbf{B}_j$  be complex valued matrices. Then the following fact holds:

$$\sum_{j=1}^N \mathbf{A}_j \mathbf{B}_j = \mathbf{A}_N \sum_{j=1}^N \mathbf{B}_j - \sum_{n=1}^{N-1} (\mathbf{A}_{n+1} - \mathbf{A}_n) \sum_{j=1}^n \mathbf{B}_j.$$

PROPOSITION 2 (SPECTRAL BOUND OF POLYNOMIAL EXPANSION). Given a collection of matrices  $\{\mathbf{A}_t\}$  and  $\{\mathbf{B}_t\}$ , such that  $\mathbf{A}_t \leq \mathbf{L}\mathbf{I}$  and  $\mathbf{B}_t \leq \mathbf{L}\mathbf{I}$ , the following bound hold:

$$\left\| \prod_{t=l}^h (\mathbf{I} - a \mathbf{A}_t) - \mathbf{I} \right\| \leq \sum_{m=1}^{h-l} \left( \frac{e(h-l)}{m} \right)^m (aL)^m,$$

$$\left\| \prod_{t=l}^h (\mathbf{I} - a \mathbf{A}_t) - \prod_{t=l}^h (\mathbf{I} - b \mathbf{B}_t) \right\| \leq \sum_{m=1}^{h-l} \left( \frac{e(h-l)}{m} \right)^m (aL)^m + \sum_{m=1}^{h-l} \left( \frac{e(h-l)}{m} \right)^m (bL)^m.$$

PROOF. We start with proving the first statement. Expanding the product yields:

$$\prod_{t=l}^h (\mathbf{I} - a \mathbf{A}_t) = \mathbf{I} + \sum_{m=1}^{h-l} (-1)^m a^m \sum_{|S|=m, |S| \subseteq \{l, \dots, h\}} \prod_{m' \in S} \mathbf{A}_{m'}.$$

Hence we have:

$$\left\| \prod_{t=l}^h (\mathbf{I} - a \mathbf{A}_t) - \mathbf{I} \right\| = \left\| \sum_{m=1}^{h-l} (-1)^m a^m \sum_{|S|=m, |S| \subseteq \{l, \dots, h\}} \prod_{m' \in S} \mathbf{A}_{m'} \right\| \leq \sum_{m=1}^{h-l} \binom{h-l}{m} (aL)^m$$

According to the upper bound for binomial coefficients:  $\binom{h-l}{m} \leq \left( \frac{e(h-l)}{m} \right)^m$ , we have:

$$\sum_{m=1}^{h-l} \binom{h-l}{m} (aL)^m \leq \sum_{m=1}^{h-l} \left( \frac{e(h-l)}{m} \right)^m (aL)^m.$$

Then we switch to the second one. Using the same expanding product yields:

$$\begin{aligned}
& \left\| \prod_{t=l}^h (\mathbf{I} - a\mathbf{A}_t) - \prod_{t=l}^h (\mathbf{I} - b\mathbf{B}_t) \right\| \\
&= \left\| \sum_{m=1}^{h-l} (-1)^m a^m \sum_{|S|=m, |S| \subseteq \{l, \dots, h\}} \prod_{m' \in S} \mathbf{A}_{m'} - \sum_{m=1}^{h-l} (-1)^m b^m \sum_{|S|=m, |S| \subseteq \{l, \dots, h\}} \prod_{m' \in S} \mathbf{B}_{m'} \right\| \\
&\leq \sum_{m=1}^{h-l} \binom{h-l}{m} (aL)^m + \sum_{m=1}^{h-l} \binom{h-l}{m} (bL)^m \\
&\leq \sum_{m=1}^{h-l} \left( \frac{e(h-l)}{m} \right)^m (aL)^m + \sum_{m=1}^{h-l} \left( \frac{e(h-l)}{m} \right)^m (bL)^m.
\end{aligned}$$

□

The following concentration result is the key to bound variance during shuffling updating. The original result holds for average of gradients, and we will later on generalize it to an arbitrary weighted sum of gradients.

LEMMA 9 ([20, THEOREM 2]). *Suppose  $n \geq 2$ . Let  $\mathbf{g}_1, \mathbf{g}_2, \dots, \mathbf{g}_n \in \mathbb{R}^d$  satisfy  $\|\mathbf{g}_j\| \leq G$  for all  $j$ . Let  $\bar{\mathbf{g}} = \frac{1}{n} \sum_{j=1}^n \mathbf{g}_j$ . Let  $\sigma \in S_n$  be a uniform random permutation of  $n$  elements. Then, for  $i \leq n$ , with probability at least  $1 - p$ , we have*

$$\left\| \frac{1}{i} \sum_{j=1}^i \mathbf{g}_{\sigma(j)} - \bar{\mathbf{g}} \right\| \leq G \sqrt{\frac{8(1 - \frac{i-1}{n}) \log \frac{2}{p}}{i}}.$$

LEMMA 10 (CONCENTRATION OF PARTIAL SUM OF GRADIENTS). *Given a uniformly randomly generated permutation  $\sigma$ , and simplex vector  $\hat{\alpha}$ , if we assume each  $\sup_{\mathbf{v} \in \mathcal{W}} \|\nabla f_j(\mathbf{v})\| \leq G$ , then the following statement holds true:*

$$\left\| \sum_{j=0}^n \hat{\alpha}(\sigma(j)) \nabla f_{\sigma(j)}(\mathbf{v}^r) \right\| \leq G \sqrt{8n \log(1/p)} + \frac{n}{N} \|\nabla \Phi(\hat{\alpha}, \mathbf{v}^r)\|.$$

PROOF. The proof works by re-writing weighted sum of vectors to average of the these vectors:

$$\begin{aligned}
\left\| \sum_{j=0}^n \hat{\alpha}(\sigma(j)) \nabla f_{\sigma(j)}(\mathbf{v}^r) \right\| &= \frac{1}{N} \left\| \sum_{j=0}^n \hat{\alpha}(\sigma(j)) N \nabla f_{\sigma(j)}(\mathbf{v}^r) \right\| \\
&= \frac{1}{N} \left( \left\| \sum_{j=0}^n \hat{\alpha}(\sigma(j)) N \nabla f_{\sigma(j)}(\mathbf{v}^r) - n \nabla \Phi(\hat{\alpha}, \mathbf{v}^r) \right\| + n \|\nabla \Phi(\hat{\alpha}, \mathbf{v}^r)\| \right) \\
&\leq G \sqrt{8n \log(1/p)} + \frac{n}{N} \|\nabla \Phi(\hat{\alpha}, \mathbf{v}^r)\|.
\end{aligned}$$

□

PROPOSITION 3 (SPECTRAL NORM BOUND OF  $\mathbf{Q}_j$ ). *Let  $\mathbf{Q}_j$  be defined in (11). Then the following bound for the spectral norm of  $\mathbf{Q}_j$  holds true for all  $j \in [N]$ :*

$$\|\mathbf{Q}_j\| \leq \eta \hat{\alpha}(\sigma(j)) N K (1 + \eta N L)^K$$

PROOF. The proof can be completed by writing down the definitin of  $\mathbf{Q}_j$  and applying Cauchy-Schwartz inequality:

$$\begin{aligned}
\|\mathbf{Q}_j\| &= \left\| \left( \sum_{\tau=0}^{K-1} \prod_{t'=K-1}^{\tau+1} (\mathbf{I} - \eta \hat{\alpha}(\sigma(j)) N' \mathbf{H}_{t'}) \right) \eta \hat{\alpha}(\sigma(j)) N \right\| \\
&\leq \eta \hat{\alpha}(\sigma(j)) N \sum_{\tau=0}^{K-1} \prod_{t'=K-1}^{\tau+1} \|(\mathbf{I} - \eta \hat{\alpha}(\sigma(j)) N' \mathbf{H}_{t'})\| \\
&\leq \eta \hat{\alpha}(\sigma(j)) N \sum_{\tau=0}^{K-1} \prod_{t'=K-1}^{\tau+1} (1 + \eta \hat{\alpha}(\sigma(j)) N' L) \\
&\leq \eta \hat{\alpha}(\sigma(j)) N K (1 + \eta N L)^K.
\end{aligned}$$



The last step is due to  $\eta NL \leq \frac{1}{K}$ . □

The following lemma establishes the bound regarding cumulative update between two epochs, namely,  $\mathbf{v}^{r+1} - \mathbf{v}^r$ . In particular, Lemma 11 below shows that: (a) in shuffling Local SGD, our update from  $\mathbf{v}^r$  to  $\mathbf{v}^{r+1}$  approximates performing  $NK$  times of gradient descent with  $\hat{\alpha}(j)N\nabla f_{\sigma(j)}(\mathbf{v}^r)$ , namely, the bias is controlled, and (b) the update itself is bounded, and can be related to the norm of full gradient.

LEMMA 11. *During the dynamic of Algorithm A2, the following statements hold true with probability at least  $1 - p$ :*

(a)

$$\left\| \sum_{j=1}^N \mathbf{Q}_j \nabla f_{\sigma(j)}(\mathbf{v}^r) - \eta NK \sum_{j=1}^N \hat{\alpha}(j) \nabla f_{\sigma(j)}(\mathbf{v}^r) \right\|^2 \leq 10\eta^2 N^2 K^2 \left( \frac{e}{4R - e} \right)^2 \|\nabla \Phi(\hat{\alpha}, \mathbf{v}^r)\|^2 + 128\eta^2 N^3 K^2 \left( \frac{e}{4R - e} \right)^2 G^2 \log(1/p).$$

(b) for any  $N'$  such that  $0 \leq N' < N$

$$\left\| \sum_{j=1}^{N'-1} \mathbf{Q}_j \nabla f_{\sigma(j)}(\mathbf{v}^r) \right\| \leq 3e\eta NK \left( \|\nabla \Phi(\hat{\alpha}, \mathbf{v}^r)\| + G\sqrt{8N \log(1/p)} \right),$$

where

$$\mathbf{Q}_j := \left( \sum_{\tau=0}^{K-1} \prod_{t'=K-1}^{\tau+1} (\mathbf{I} - \eta \hat{\alpha}(\sigma(j)) N' \mathbf{H}_{t'}) \right) \eta \hat{\alpha}(\sigma(j)) N. \quad (11)$$

PROOF. We start with proving statement (a). Let  $\mathbf{A}_j = \frac{\mathbf{Q}_j}{\hat{\alpha}(\sigma(j))}$  and  $\mathbf{B}_j = \hat{\alpha}(\sigma(j)) \nabla f_{\sigma(j)}(\mathbf{v}^r)$ , applying the identity of summation by parts yields:

$$\sum_{j=1}^N \mathbf{Q}_j \nabla f_{\sigma(j)}(\mathbf{v}^r) = \frac{\mathbf{Q}_{N-1}}{\hat{\alpha}(\sigma(N-1))} \sum_{j=1}^N \hat{\alpha}(\sigma(j)) \nabla f_{\sigma(j)}(\mathbf{v}^r) - \sum_{n=1}^{N-1} \left( \frac{\mathbf{Q}_{n+1}}{\hat{\alpha}(\sigma(n+1))} - \frac{\mathbf{Q}_n}{\hat{\alpha}(\sigma(n))} \right) \sum_{j=1}^n \hat{\alpha}(\sigma(j)) \nabla f_{\sigma(j)}(\mathbf{v}^r)$$

$$\left\| \sum_{j=1}^N \mathbf{Q}_j \nabla f_{\sigma(j)}(\mathbf{v}^r) - \eta NK \sum_{j=1}^N \hat{\alpha}(j) \nabla f_{\sigma(j)}(\mathbf{v}^r) \right\|^2 \leq 2 \underbrace{\left\| \left( \frac{\mathbf{Q}_{N-1}}{\hat{\alpha}(\sigma(N-1))} - \eta NK \mathbf{I} \right) \sum_{j=1}^N \hat{\alpha}(\sigma(j)) \nabla f_{\sigma(j)}(\mathbf{v}^r) \right\|^2}_{T_1} + 2 \underbrace{\left\| \sum_{n=1}^{N-1} \left( \frac{\mathbf{Q}_{n+1}}{\hat{\alpha}(\sigma(n+1))} - \frac{\mathbf{Q}_n}{\hat{\alpha}(\sigma(n))} \right) \sum_{j=1}^n \hat{\alpha}(\sigma(j)) \nabla f_{\sigma(j)}(\mathbf{v}^r) \right\|^2}_{T_2}.$$

According to Proposition 2, we have:

$$\left\| \prod_{t'=K-1}^{\tau+1} (\mathbf{I} - \eta \hat{\alpha}(\sigma(j)) N \mathbf{H}_{t'}) - \mathbf{I} \right\| \leq \sum_{m=1}^{K-2-\tau} \left( \frac{e^{(K-2-\tau)}}{m} \eta \hat{\alpha}(\sigma(j)) NL \right)^m.$$

Since we choose  $\eta \leq \frac{1}{4NKRL}$ , we have:

$$\left\| \prod_{t'=K-1}^{\tau+1} (\mathbf{I} - \eta \hat{\alpha}(\sigma(j)) N \mathbf{H}_{t'}) - \mathbf{I} \right\| \leq \sum_{m=1}^{K-2-\tau} \left( \frac{e}{4Rm} \right)^m \leq \frac{e}{4R - e}, \quad (12)$$

where we use the fact that  $\sum_{m=1}^{K-2-\tau} \left(\frac{e}{4Rm}\right)^m \leq \sum_{m=1}^{K-2-\tau} \left(\frac{e}{4R}\right)^m \leq \frac{e}{4R} \frac{1}{1-e/4R}$ . Hence we know:

$$\begin{aligned}
T_1 &\leq \left\| \left( \frac{\mathcal{Q}_{N-1}}{\hat{\alpha}(\sigma(N-1))} - \eta N K \mathbf{I} \right) \right\|^2 \left\| \sum_{j=1}^N \hat{\alpha}(\sigma(j)) \nabla f_{\sigma(j)}(\mathbf{v}^r) \right\|^2 \\
&\leq \left\| \left( \sum_{\tau=0}^{K-1} \prod_{t'=K-1}^{\tau+1} (\mathbf{I} - \eta \hat{\alpha}(\sigma(N-1)) N \mathbf{H}_{t'}) \right) \eta N - \eta N K \mathbf{I} \right\|^2 \left\| \sum_{j=1}^N \hat{\alpha}(\sigma(j)) \nabla f_{\sigma(j)}(\mathbf{v}^r) \right\|^2 \\
&\leq \eta^2 N^2 K \sum_{\tau=0}^{K-1} \left\| \prod_{t'=K-1}^{\tau+1} (\mathbf{I} - \eta \hat{\alpha}(\sigma(N-1)) N \mathbf{H}_{t'}) - \mathbf{I} \right\|^2 \left\| \sum_{j=1}^N \hat{\alpha}(\sigma(j)) \nabla f_{\sigma(j)}(\mathbf{v}^r) \right\|^2 \\
&\leq \eta^2 N^2 K^2 \left( \frac{e}{4R-e} \right)^2 \left\| \sum_{j=1}^N \hat{\alpha}(\sigma(j)) \nabla f_{\sigma(j)}(\mathbf{v}^r) \right\|^2.
\end{aligned}$$

Thus we have:

$$T_1 \leq \eta^2 N^2 K^2 \left( \frac{e}{4R-e} \right)^2 \|\nabla \Phi(\hat{\alpha}, \mathbf{v}^r)\|^2.$$

For  $T_2$ , we first examine the bound of  $\frac{\mathcal{Q}_{n+1}}{\hat{\alpha}(\sigma(n+1))} - \frac{\mathcal{Q}_n}{\hat{\alpha}(\sigma(n))}$ :

$$\begin{aligned}
\left\| \frac{\mathcal{Q}_{n+1}}{\hat{\alpha}(\sigma(n+1))} - \frac{\mathcal{Q}_n}{\hat{\alpha}(\sigma(n))} \right\| &= \left\| \left( \sum_{\tau=0}^{K-1} \prod_{t'=K-1}^{\tau+1} (\mathbf{I} - \eta \hat{\alpha}(\sigma(n+1)) N \mathbf{H}_{t'}) \right) \eta N - \left( \sum_{\tau=0}^{K-1} \prod_{t'=K-1}^{\tau+1} (\mathbf{I} - \eta \hat{\alpha}(\sigma(n)) N \mathbf{H}_{t'}) \right) \eta N \right\| \\
&= \eta N \left\| \left( \sum_{\tau=0}^{K-1} \prod_{t'=K-1}^{\tau+1} (\mathbf{I} - \eta \hat{\alpha}(\sigma(n+1)) N \mathbf{H}_{t'}) \right) - \left( \sum_{\tau=0}^{K-1} \prod_{t'=K-1}^{\tau+1} (\mathbf{I} - \eta \hat{\alpha}(\sigma(n)) N \mathbf{H}_{t'}) \right) \right\| \\
&= \eta N \left\| \left( \sum_{\tau=0}^{K-1} \prod_{t'=K-1}^{\tau+1} (\mathbf{I} - \eta \hat{\alpha}(\sigma(n+1)) N \mathbf{H}_{t'}) \right) - \left( \sum_{\tau=0}^{K-1} \prod_{t'=K-1}^{\tau+1} (\mathbf{I} - \eta \hat{\alpha}(\sigma(n)) N \mathbf{H}_{t'}) \right) \right\| \\
&\leq \eta N \sum_{\tau=0}^{K-1} \left( \sum_{m=1}^{K-2-\tau} \left( \frac{e(K-2-\tau)}{m} \eta \hat{\alpha}(\sigma(n)) N L \right)^m + \sum_{m=1}^{K-2-\tau} \left( \frac{e(K-2-\tau)}{m} \eta \hat{\alpha}(\sigma(n+1)) N L \right)^m \right).
\end{aligned}$$

where we evoke Proposition 2 at last step. Given that  $\eta \leq \frac{1}{4NKRL}$  we have:

$$\begin{aligned}
\left\| \frac{\mathcal{Q}_{n+1}}{\hat{\alpha}(\sigma(n+1))} - \frac{\mathcal{Q}_n}{\hat{\alpha}(\sigma(n))} \right\| &\leq \eta N \sum_{\tau=0}^{K-1} \left( \sum_{m=1}^{K-2-\tau} \left( \frac{e}{4Rm} \hat{\alpha}(\sigma(n)) \right)^m + \sum_{m=1}^{K-2-\tau} \left( \frac{e}{4Rm} \hat{\alpha}(\sigma(n+1)) \right)^m \right) \\
&\leq \eta N K \left( \frac{\hat{\alpha}(\sigma(n)) e}{4R-e} + \frac{\hat{\alpha}(\sigma(n+1)) e}{4R-e} \right).
\end{aligned}$$

where we use the reasoning in (12). Hence for  $\sqrt{T_2}$ :

$$\begin{aligned}
\sqrt{T_2} &\leq \eta N K \sum_{n=1}^{N-1} \left( \frac{\hat{\alpha}(\sigma(n)) e}{4R-e} + \frac{\hat{\alpha}(\sigma(n+1)) e}{4R-e} \right) \left\| \sum_{j=1}^n \hat{\alpha}(\sigma(j)) \nabla f_{\sigma(j)}(\mathbf{v}^r) \right\| \\
&\leq \eta N K \frac{e}{4R-e} \sum_{n=1}^{N-1} (\hat{\alpha}(\sigma(n)) + \hat{\alpha}(\sigma(n+1))) \left( G \sqrt{8n \log(1/p)} + \frac{n}{N} \|\nabla \Phi(\hat{\alpha}, \mathbf{v}^r)\| \right) \\
&\leq \eta N K \frac{2e}{4R-e} \left( G \sqrt{8N \log(1/p)} + \|\nabla \Phi(\hat{\alpha}, \mathbf{v}^r)\| \right).
\end{aligned}$$

where at last step we evoke Lemma 10. So we can conclude  $T_2 \leq 2\eta^2 N^2 K^2 \left( \frac{2e}{4R-e} \right)^2 \left( G^2 8N \log(1/p) + \|\nabla \Phi(\hat{\alpha}, \mathbf{v}^r)\|^2 \right)$ . Putting the bounds of  $T_1$  and  $T_2$  together will conclude the proof for (a).

Now we switch to proving (b). Once again by the summation of parts identity we have:

$$\sum_{j=1}^{N'} \mathcal{Q}_j \nabla f_{\sigma(j)}(\mathbf{v}^r) = \frac{\mathcal{Q}_{N'}}{\hat{\alpha}(\sigma(N'))} \sum_{j=1}^{N'} \hat{\alpha}(\sigma(j)) \nabla f_{\sigma(j)}(\mathbf{v}^r) - \sum_{n=1}^{N'-1} \left( \frac{\mathcal{Q}_{n+1}}{\hat{\alpha}(\sigma(n+1))} - \frac{\mathcal{Q}_n}{\hat{\alpha}(\sigma(n))} \right) \sum_{j=1}^n \hat{\alpha}(\sigma(j)) \nabla f_{\sigma(j)}(\mathbf{v}^r).$$

Taking the norm of both side yields:

$$\begin{aligned} \left\| \sum_{j=1}^{N'} \mathcal{Q}_j \nabla f_{\sigma(j)}(\mathbf{v}^r) \right\| &= \underbrace{\left\| \frac{\mathcal{Q}_{N'}}{\hat{\alpha}(\sigma(N'))} \sum_{j=1}^{N'} \hat{\alpha}(\sigma(j)) \nabla f_{\sigma(j)}(\mathbf{v}^r) \right\|}_B \\ &+ \underbrace{\left\| \sum_{n=1}^{N'-1} \left( \frac{\mathcal{Q}_{n+1}}{\hat{\alpha}(\sigma(n+1))} - \frac{\mathcal{Q}_n}{\hat{\alpha}(\sigma(n))} \right) \sum_{j=1}^n \hat{\alpha}(\sigma(j)) \nabla f_{\sigma(j)}(\mathbf{v}^r) \right\|}_C. \end{aligned}$$

Plugging our developed bound for  $\|\mathcal{Q}_{N'}\|$  and  $\left\| \sum_{n=1}^{N'-1} \left( \frac{\mathcal{Q}_{n+1}}{\hat{\alpha}(\sigma(n+1))} - \frac{\mathcal{Q}_n}{\hat{\alpha}(\sigma(n))} \right) \right\|$  yields:

$$\begin{aligned} B &\leq \left\| \frac{\mathcal{Q}_{N'}}{\hat{\alpha}(\sigma(N'))} \right\| \left\| \sum_{j=1}^{N'} \hat{\alpha}(\sigma(j)) \nabla f_{\sigma(j)}(\mathbf{v}^r) \right\| \\ &\leq \eta NK(1 + \eta NL)^K \left( G\sqrt{8N' \log(1/p)} + \frac{N'}{N} \|\nabla\Phi(\hat{\boldsymbol{\alpha}}, \mathbf{v}^r)\| \right). \end{aligned}$$

where at last step we evoke Lemma 10. And for C, we use the similar reasoning:

$$\begin{aligned} C &\leq \sum_{n=1}^{N'-1} \left\| \left( \frac{\mathcal{Q}_{n+1}}{\hat{\alpha}(\sigma(n+1))} - \frac{\mathcal{Q}_n}{\hat{\alpha}(\sigma(n))} \right) \right\| \left\| \sum_{j=1}^n \hat{\alpha}(\sigma(j)) \nabla f_{\sigma(j)}(\mathbf{v}^r) \right\| \\ &\leq \sum_{n=1}^{N'-1} \eta NK \frac{e}{4R-e} (\hat{\alpha}(\sigma(n+1)) + \hat{\alpha}(\sigma(n))) \left( G\sqrt{8n \log(1/p)} + \frac{n}{N} \|\nabla\Phi(\hat{\boldsymbol{\alpha}}, \mathbf{v}^r)\| \right) \\ &\leq 2\eta NK \frac{e}{4R-e} \left( G\sqrt{8N \log(1/p)} + \|\nabla\Phi(\hat{\boldsymbol{\alpha}}, \mathbf{v}^r)\| \right). \end{aligned}$$

Putting these pieces together yields:

$$\left\| \sum_{j=1}^{N'} \mathcal{Q}_j \nabla f_{\sigma(j)}(\mathbf{v}^r) \right\| \leq 3e\eta NK \left( \|\nabla\Phi(\hat{\boldsymbol{\alpha}}, \mathbf{v}^r)\| + G\sqrt{8N \log(1/p)} \right).$$

□

LEMMA 12. *During the dynamic of Algorithm A2, the following statements hold true with probability at least  $1 - p$ :*

$$\left\| \sum_{n=1}^{N-1} \left( \prod_{j'=N}^{n+2} (\mathbf{I} - \mathcal{Q}_{j'} \mathbf{H}_{j'}) \right) \mathcal{Q}_{n+1} \mathbf{H}_{n+1} \sum_{j=1}^n \mathcal{Q}_j \nabla f_{\sigma(j)}(\mathbf{v}^r) \right\|^2 \leq 18e^6 \eta^4 N^4 K^4 L^4 \left( \|\nabla\Phi(\hat{\boldsymbol{\alpha}}, \mathbf{v}^r)\|^2 + 8G^2 N \log(1/p) \right)$$

PROOF. We first apply Cauchy-Schwartz inequality:

$$\begin{aligned} &\left\| \sum_{n=1}^{N-1} \left( \prod_{j'=N}^{n+2} (\mathbf{I} - \mathcal{Q}_{j'} \mathbf{H}_{j'}) \right) \mathcal{Q}_{n+1} \mathbf{H}_{n+1} \sum_{j=1}^n \mathcal{Q}_j \nabla f_{\sigma(j)}(\mathbf{v}^r) \right\| \\ &\leq \sum_{n=1}^{N-1} \left\| \left( \prod_{j'=N}^{n+2} (\mathbf{I} - \mathcal{Q}_{j'} \mathbf{H}_{j'}) \right) \right\| \left\| \mathcal{Q}_{n+1} \mathbf{H}_{n+1} \right\| \left\| \sum_{j=1}^n \mathcal{Q}_j \nabla f_{\sigma(j)}(\mathbf{v}^r) \right\| \\ &\leq \left( 1 + \eta NK + \eta^2 N^2 KL \right)^{2N} \eta NLK(1 + \eta NL)^{KL} \sum_{n=1}^{N-1} \hat{\alpha}(\sigma(n+1)) \left\| \sum_{j=1}^n \mathcal{Q}_j \nabla f_{\sigma(j)}(\mathbf{v}^r) \right\| \\ &\leq e^2 \eta NKL^2 \sum_{n=1}^{N-1} \hat{\alpha}(\sigma(n+1)) \left\| \sum_{j=1}^n \mathcal{Q}_j \nabla f_{\sigma(j)}(\mathbf{v}^r) \right\|. \end{aligned}$$

We proceed by applying the bound from Lemma 11 (b):

$$\left\| \sum_{j=1}^n \mathbf{Q}_j \nabla f_{\sigma(j)}(\mathbf{v}^r) \right\| \leq 3e\eta NK \left( \|\nabla\Phi(\hat{\boldsymbol{\alpha}}, \mathbf{v}^r)\| + G\sqrt{8N \log(1/p)} \right).$$

Therefore, it follows that:

$$\begin{aligned} & \left\| \sum_{n=1}^{N-1} \left( \prod_{j'=N}^{n+2} (\mathbf{I} - \mathbf{Q}_{j'} \mathbf{H}_{j'}) \right) \mathbf{Q}_{n+1} \mathbf{H}_{n+1} \sum_{j=1}^n \mathbf{Q}_j \nabla f_{\sigma(j)}(\mathbf{v}^r) \right\| \\ & \leq e^2 \eta NK L^2 \sum_{n=1}^{N-1} \hat{\alpha}(\sigma(n+1)) \cdot 3e\eta NK \left( \|\nabla\Phi(\hat{\boldsymbol{\alpha}}, \mathbf{v}^r)\| + G\sqrt{8N \log(1/p)} \right) \\ & \leq 3e^3 \eta^2 N^2 K^2 L^2 \left( \|\nabla\Phi(\hat{\boldsymbol{\alpha}}, \mathbf{v}^r)\| + G\sqrt{8N \log(1/p)} \right) \end{aligned}$$

□

LEMMA 13 (NOISE BOUND). *During the dynamic of Algorithm A2, the following statement for gradient noises holds true with probability at least  $1 - p$ :*

$$\mathbb{E} \left\| \sum_{j=1}^N \prod_{j'=N}^{j+1} (\mathbf{I} - \mathbf{Q}_{j'} \mathbf{H}_{j'}) \boldsymbol{\delta}_j \right\| \leq \eta NK e^2 \delta,$$

where

$$\boldsymbol{\delta}_j := \sum_{\tau=0}^{K-1} \prod_{t'=t}^{\tau+1} (\mathbf{I} - \hat{\alpha}(\sigma(j)) N \eta \mathbf{H}_{t'}) \eta \hat{\alpha}(\sigma(j)) N \boldsymbol{\delta}_{\sigma(j)}^t.$$

PROOF. According to triangle and Cauchy-Schwartz inequalities we have:

$$\begin{aligned} \left\| \sum_{j=1}^N \prod_{j'=N}^{j+1} (\mathbf{I} - \mathbf{Q}_{j'} \mathbf{H}_{j'}) \boldsymbol{\delta}_j \right\| & \leq \sum_{j=1}^N \prod_{j'=N}^{j+1} \left\| (\mathbf{I} - \mathbf{Q}_{j'} \mathbf{H}_{j'}) \right\| \|\boldsymbol{\delta}_j\| \\ & \leq \sum_{j=1}^N \left( 1 + (\eta \hat{\alpha}(\sigma(j)) NK (1 + \eta NL^K) L) \right)^N \|\boldsymbol{\delta}_j\| \\ & \leq \sum_{j=1}^N \underbrace{\left( 1 + (\eta \hat{\alpha}(\sigma(j)) NK (1 + \eta NL^K) L) \right)^N}_{\leq e} \cdot \underbrace{\eta \hat{\alpha}(\sigma(j)) NK (1 + \eta NL^K) \delta}_{\leq e} \\ & \leq \eta NK e^2 \delta. \end{aligned}$$

□

## A.4 Proof of Theorem 2

PROOF. For notational convenience, let us define

$$\begin{aligned} \mathbf{g}^r & := \sum_{j=1}^N \prod_{j'=N}^{j+1} (\mathbf{I} - \mathbf{Q}_{j'} \mathbf{H}_{j'}) \mathbf{Q}_j \nabla f_{\sigma(j)}(\mathbf{v}^r), \\ \boldsymbol{\delta}^r & := \sum_{j=1}^N \prod_{j'=N}^{j+1} (\mathbf{I} - \mathbf{Q}_{j'} \mathbf{H}_{j'}) \boldsymbol{\delta}_j. \end{aligned}$$

Then we recall the updating rule of  $\mathbf{v}$  (Lemma 7):

$$\mathbf{v}^{r+1} = \mathcal{P}_{\mathcal{W}}(\mathbf{v}^r - \mathbf{g}^r - \boldsymbol{\delta}^r)$$



Hence we have:

$$\begin{aligned}
\mathbb{E} \|\mathbf{v}^{r+1} - \mathbf{v}^*(\hat{\boldsymbol{\alpha}})\|^2 &= \mathbb{E} \|\mathcal{P}_{\mathcal{W}}(\mathbf{v}^r - \mathbf{g}^r - \boldsymbol{\delta}^r - \mathbf{v}^*(\hat{\boldsymbol{\alpha}}))\|^2 \\
&\leq \mathbb{E} \|\mathbf{v}^r - \mathbf{g}^r - \boldsymbol{\delta}^r - \mathbf{v}^*(\hat{\boldsymbol{\alpha}})\|^2 \\
&\leq \mathbb{E} \|\mathbf{v}^r - \mathbf{v}^*(\hat{\boldsymbol{\alpha}})\|^2 - 2\mathbb{E}\langle \mathbf{g}^r, \mathbf{v}^r - \mathbf{v}^*(\hat{\boldsymbol{\alpha}}) \rangle + \mathbb{E} \|\mathbf{g}^r\|^2 + \mathbb{E} \|\boldsymbol{\delta}^r\|^2 \\
&\leq \mathbb{E} \|\mathbf{v}^r - \mathbf{v}^*(\hat{\boldsymbol{\alpha}})\|^2 - 2\mathbb{E}\langle \eta NK \nabla \Phi(\hat{\boldsymbol{\alpha}}, \mathbf{v}^r), \mathbf{v}^r - \mathbf{v}^*(\hat{\boldsymbol{\alpha}}) \rangle - 2\mathbb{E}\langle \mathbf{g}^r - \eta NK \nabla \Phi(\hat{\boldsymbol{\alpha}}, \mathbf{v}^r), \mathbf{v}^r - \mathbf{v}^*(\hat{\boldsymbol{\alpha}}) \rangle \\
&\quad + \mathbb{E} \|\mathbf{g}^r\|^2 + \mathbb{E} \|\boldsymbol{\delta}^r\|^2.
\end{aligned}$$

Now, applying strongly convexity of  $\Phi(\hat{\boldsymbol{\alpha}}, \cdot)$  and Cauchy-Schwartz inequality yields:

$$\begin{aligned}
\mathbb{E} \|\mathbf{v}^{r+1} - \mathbf{v}^*(\hat{\boldsymbol{\alpha}})\|^2 &\leq (1 - \mu\eta NK) \mathbb{E} \|\mathbf{v}^r - \mathbf{v}^*(\hat{\boldsymbol{\alpha}})\|^2 - \eta NK \mathbb{E} [\Phi(\hat{\boldsymbol{\alpha}}, \mathbf{v}^r) - \Phi(\hat{\boldsymbol{\alpha}}, \mathbf{v}^*(\hat{\boldsymbol{\alpha}}))] \\
&\quad + \frac{1}{2} \left( \frac{1}{\mu\eta NK} \mathbb{E} \|\mathbf{g}^r - \eta NK \nabla \Phi(\hat{\boldsymbol{\alpha}}, \mathbf{v}^r)\|^2 + \mu\eta NK \mathbb{E} \|\mathbf{v}^r - \mathbf{v}^*(\hat{\boldsymbol{\alpha}})\|^2 \right) \\
&\quad + \mathbb{E} \|\mathbf{g}^r\|^2 + \mathbb{E} \|\boldsymbol{\delta}^r\|^2 \\
&\leq (1 - \frac{1}{2}\mu\eta NK) \mathbb{E} \|\mathbf{v}^r - \mathbf{v}^*(\hat{\boldsymbol{\alpha}})\|^2 - \eta NK \mathbb{E} [\Phi(\hat{\boldsymbol{\alpha}}, \mathbf{v}^r) - \Phi(\hat{\boldsymbol{\alpha}}, \mathbf{v}^*(\hat{\boldsymbol{\alpha}}))] \\
&\quad + \frac{1}{2\mu\eta NK} \mathbb{E} \|\mathbf{g}^r - \eta NK \nabla \Phi(\hat{\boldsymbol{\alpha}}, \mathbf{v}^r)\|^2 \\
&\quad + 2\mathbb{E} \|\mathbf{g}^r - \eta NK \nabla \Phi(\hat{\boldsymbol{\alpha}}, \mathbf{v}^r)\|^2 + 2\mathbb{E} \|\eta NK \nabla \Phi(\hat{\boldsymbol{\alpha}}, \mathbf{v}^r)\|^2 + \mathbb{E} \|\boldsymbol{\delta}^r\|^2.
\end{aligned}$$

Since  $\Phi(\hat{\boldsymbol{\alpha}}, \cdot)$  is  $L$  smooth, we have:  $\mathbb{E} \|\nabla \Phi(\hat{\boldsymbol{\alpha}}, \mathbf{v}^r)\|^2 \leq 2L \mathbb{E} [\Phi(\hat{\boldsymbol{\alpha}}, \mathbf{v}^r) - \Phi(\hat{\boldsymbol{\alpha}}, \mathbf{v}^*(\hat{\boldsymbol{\alpha}}))]$ . Therefore, we have:

$$\mathbb{E} \|\mathbf{v}^{r+1} - \mathbf{v}^*(\hat{\boldsymbol{\alpha}})\|^2 \leq (1 - \frac{1}{2}\mu\eta NK) \mathbb{E} \|\mathbf{v}^r - \mathbf{v}^*(\hat{\boldsymbol{\alpha}})\|^2 - (\eta NK - 4\eta^2 N^2 K^2 L) \mathbb{E} [\Phi(\hat{\boldsymbol{\alpha}}, \mathbf{v}^r) - \Phi(\hat{\boldsymbol{\alpha}}, \mathbf{v}^*(\hat{\boldsymbol{\alpha}}))] \quad (13)$$

$$+ \left( \frac{1}{2\mu\eta NK} + 2 \right) \mathbb{E} \|\mathbf{g}^r - \eta NK \nabla \Phi(\hat{\boldsymbol{\alpha}}, \mathbf{v}^r)\|^2 + \mathbb{E} \|\boldsymbol{\delta}^r\|^2 \quad (14)$$

Now, we examine the term  $\|\mathbf{g}^r - \eta NK \nabla \Phi(\hat{\boldsymbol{\alpha}}, \mathbf{v}^r)\|^2$ . First according to summation by part (Lemma 8) by letting  $\mathbf{A}_j := \prod_{j'=N}^{j+1} (\mathbf{I} - \mathbf{Q}_{j'} \mathbf{H}_{j'})$  and  $\mathbf{B}_j = \mathbf{Q}_j \nabla f_{\sigma(j)}(\mathbf{v}^r)$ , we have:

$$\begin{aligned}
\mathbf{g}^r &= \sum_{j=1}^N \prod_{j'=N}^{j+1} (\mathbf{I} - \mathbf{Q}_{j'} \mathbf{H}_{j'}) \mathbf{Q}_j \nabla f_{\sigma(j)}(\mathbf{v}^r) \\
&= \sum_{j=1}^N \mathbf{A}_j \mathbf{B}_j = \sum_{j=1}^N \mathbf{Q}_j \nabla f_{\sigma(j)}(\mathbf{v}^r) - \sum_{n=1}^{N-1} \left( \prod_{j'=N}^{n+2} (\mathbf{I} - \mathbf{Q}_{j'} \mathbf{H}_{j'}) - \prod_{j'=N}^{n+1} (\mathbf{I} - \mathbf{Q}_{j'} \mathbf{H}_{j'}) \right) \sum_{j=1}^n \mathbf{Q}_j \nabla f_{\sigma(j)}(\mathbf{v}^r) \\
&= \sum_{j=1}^N \mathbf{Q}_j \nabla f_{\sigma(j)}(\mathbf{v}^r) - \sum_{n=1}^{N-1} \left( \prod_{j'=N}^{n+2} (\mathbf{I} - \mathbf{Q}_{j'} \mathbf{H}_{j'}) \right) \mathbf{Q}_{n+1} \mathbf{H}_{n+1} \sum_{j=1}^n \mathbf{Q}_j \nabla f_{\sigma(j)}(\mathbf{v}^r).
\end{aligned}$$

Hence we have:

$$\begin{aligned}
&\|\mathbf{g}^r - \eta NK \nabla \Phi(\hat{\boldsymbol{\alpha}}, \mathbf{v}^r)\|^2 \\
&= \left\| \eta NK \sum_{j=1}^N \hat{\boldsymbol{\alpha}}(\sigma(j)) \nabla f_{\sigma(j)}(\mathbf{v}^r) - \left( \sum_{j=1}^N \mathbf{Q}_j \nabla f_{\sigma(j)}(\mathbf{v}^r) - \sum_{n=1}^{N-1} \left( \prod_{j'=N}^{n+2} (\mathbf{I} - \mathbf{Q}_{j'} \mathbf{H}_{j'}) \right) \mathbf{Q}_{n+1} \mathbf{H}_{n+1} \sum_{j=1}^n \mathbf{Q}_j \nabla f_{\sigma(j)}(\mathbf{v}^r) \right) \right\|^2 \\
&\stackrel{(1)}{=} 2 \left\| \left( \eta NK \sum_{j=1}^N \hat{\boldsymbol{\alpha}}(\sigma(j)) \nabla f_{\sigma(j)}(\mathbf{v}^r) - \sum_{j=1}^N \mathbf{Q}_j \nabla f_{\sigma(j)}(\mathbf{v}^r) \right) \right\|^2 + 2 \left\| \sum_{n=1}^{N-1} \left( \prod_{j'=N}^{n+2} (\mathbf{I} - \mathbf{Q}_{j'} \mathbf{H}_{j'}) \right) \mathbf{Q}_{n+1} \mathbf{H}_{n+1} \sum_{j=1}^n \mathbf{Q}_j \nabla f_{\sigma(j)}(\mathbf{v}^r) \right\|^2 \\
&\stackrel{(2)}{\leq} \left( 20\eta^2 N^2 K^2 \left( \frac{e}{4R - e} \right)^2 + 36e^6 \eta^4 N^4 K^4 L^4 \right) \|\nabla \Phi(\hat{\boldsymbol{\alpha}}, \mathbf{v}^r)\|^2 + 256\eta^2 N^3 K^2 \left( \frac{e}{4R - e} \right)^2 G^2 \log(1/p) \\
&\quad + 244e^6 \eta^4 N^4 K^4 L^4 G^2 N \log(1/p) \\
&\stackrel{(3)}{\leq} \left( 20\eta^2 N^2 K^2 \left( \frac{e}{4R - e} \right)^2 + 36e^6 \eta^4 N^4 K^4 L^4 \right) 2L (\Phi(\hat{\boldsymbol{\alpha}}, \mathbf{v}^r) - \Phi(\hat{\boldsymbol{\alpha}}, \mathbf{v}^*(\hat{\boldsymbol{\alpha}}))) \\
&\quad + \left( 244e^6 \eta^4 N^4 K^4 L^4 + 256\eta^2 N^3 K^2 \left( \frac{e}{4R - e} \right)^2 \right) G^2 N \log(1/p),
\end{aligned}$$

where in (1) we apply Jensen's inequality, in (2) we plug in Lemma 11 (a), and Lemma 12, and in (3) we use the  $L$ -smoothness of  $\Phi$ . Plugging above bound back in (19) yields:

$$\begin{aligned} & \mathbb{E} \|\mathbf{v}^{r+1} - \mathbf{v}^*(\hat{\boldsymbol{\alpha}})\|^2 \\ & \leq \left(1 - \frac{1}{2}\mu\eta NK\right) \mathbb{E} \|\mathbf{v}^r - \mathbf{v}^*(\hat{\boldsymbol{\alpha}})\|^2 + \eta^2 N^2 K^2 e^4 \delta^2 \\ & \quad - \underbrace{\left(\eta NK - 4\eta^2 N^2 K^2 L - \left(\frac{1}{2\mu\eta NK} + 2\right) \left(20\eta^2 N^2 K^2 \left(\frac{e}{4R - e}\right)^2 - 36e^6 \eta^4 N^4 K^4 L^4\right)\right)}_{T_1} \mathbb{E}[\Phi(\hat{\boldsymbol{\alpha}}, \mathbf{v}^r) - \Phi(\hat{\boldsymbol{\alpha}}, \mathbf{v}^*(\hat{\boldsymbol{\alpha}}))] \\ & \quad + \left(\frac{1}{2\mu\eta NK} + 2\right) \left(244e^6 \eta^4 N^4 K^4 L^4 + 256\eta^2 N^3 K^2 \left(\frac{e}{4R - e}\right)^2\right) G^2 N \log(1/p). \end{aligned}$$

Since we choose  $\eta = \frac{4\log(\sqrt{NKR})}{\mu NKR}$ , and large enough epoch number:

$$R \geq \max \left\{ \left(\frac{40}{\mu} + 1\right) e, 16 \log(\sqrt{NKR}), 64\kappa \log(\sqrt{NKR}) \right\},$$

we know that  $T_1 \leq 0$ . We thus have:

$$\begin{aligned} & \mathbb{E} \|\mathbf{v}^{r+1} - \mathbf{v}^*(\hat{\boldsymbol{\alpha}})\|^2 \\ & \leq \left(1 - \frac{1}{2}\mu\eta NK\right) \mathbb{E} \|\mathbf{v}^r - \mathbf{v}^*(\hat{\boldsymbol{\alpha}})\|^2 + \eta^2 N^2 K^2 e^4 \delta^2 + \left(\frac{1}{2\mu\eta NK} + 2\right) \left(244e^6 \eta^4 N^4 K^4 L^4 + 256\eta^2 N^3 K^2 \left(\frac{e}{4R - e}\right)^2\right) G^2 N \log(1/p) \end{aligned}$$

Unrolling the recursion from  $r = R$  to 0:

$$\begin{aligned} & \mathbb{E} \|\mathbf{v}^R - \mathbf{v}^*(\hat{\boldsymbol{\alpha}})\|^2 \\ & \leq \left(1 - \frac{1}{2}\mu\eta NK\right)^R \mathbb{E} \|\mathbf{v}^0 - \mathbf{v}^*(\hat{\boldsymbol{\alpha}})\|^2 + \frac{2}{\mu} \eta NK e^4 \delta^2 \\ & \quad + \frac{1}{\mu} \left(\frac{1}{2\mu\eta NK} + 2\right) \left(488e^6 \eta^3 N^3 K^3 L^4 + 512\eta N^2 K \left(\frac{e}{4R - e}\right)^2\right) G^2 N \log(1/p) \end{aligned}$$

Plugging in our choice of  $\eta$  will conclude the proof:

$$\mathbb{E} \|\mathbf{v}^R - \mathbf{v}^*(\hat{\boldsymbol{\alpha}})\|^2 \leq \tilde{O} \left( \frac{\mathbb{E} \|\mathbf{v}^0 - \mathbf{v}^*(\hat{\boldsymbol{\alpha}})\|^2}{NKR^2} + \frac{\delta^2}{\mu^2 R} + \left(\frac{L^4 + N}{\mu^4 R^2}\right) G^2 N \log(1/p) \right)$$

Finally, according to Lemma 2 we can complete the proof:

$$\begin{aligned} \Phi(\boldsymbol{\alpha}_i^*, \hat{\mathbf{v}}_i) - \Phi(\boldsymbol{\alpha}_i^*, \mathbf{v}_i^*) & \leq L \|\hat{\mathbf{v}}_i - \mathbf{v}_i^*\|^2 + \kappa^2 L \|\hat{\boldsymbol{\alpha}}_i - \boldsymbol{\alpha}^*\|^2 \\ & \leq \tilde{O} \left( \frac{LD^2}{NKR^2} + \frac{L\delta^2}{\mu^2 R} + \left(\frac{L^4 + N}{\mu^4 R^2}\right) LG^2 N \log(1/p) \right) \\ & \quad + \tilde{O} \left( \exp\left(-\frac{T\alpha}{\kappa_g}\right) + \kappa_g^2 \bar{\zeta}_i(\mathbf{w}^*) L^2 \left(\frac{D^2}{RK} + \frac{\kappa \zeta^2}{\mu^2 R^2} + \frac{\delta^2}{\mu^2 NRK}\right) \right), \end{aligned}$$

where we plug in the convergence result from Theorem 1 at last step.  $\square$

## B PROOF OF CONVERGENCE OF SINGLE LOOP ALGORITHM

In this section, we turn to presenting the proof of single loop PERM algorithm (Algorithm 2) where the learning of mixing parameters and personalized models are coupled. Compared to Algorithm A2, here during the optimization of model, the mixing parameters are also being updated. As a result, we need to decouple the two updates which makes the analysis more involved. We begin with some technical lemmas that support the proof of main result.

### B.1 Technical Lemmas

**PROPOSITION 4 (BASIC PROPERTIES OF SGD ON SMOOTH STRONGLY CONVEX FUNCTION).** *Let  $\mathbf{w}^t$  to be the  $t$ th iterate of minibatch SGD on smooth and strongly convex function  $F$ , with minibatch size  $M$  and learning rate  $\gamma$ . Also assume the variance is bounded by  $\delta$ . Then the following statements hold true after  $T$  iterations of SGD:*

$$\mathbb{E} \|\nabla F(\mathbf{w}^T)\|^2 \leq 2L(1 - \mu\gamma)^T (F(\mathbf{w}^0) - F(\mathbf{w}^*)) + \frac{2\gamma\kappa\delta^2}{M} \quad (15)$$

$$\mathbb{E}\|\mathbf{w}^{T+1} - \mathbf{w}^T\|^2 \leq 2\gamma^2 L (1 - \mu\gamma)^T (F(\mathbf{w}^0) - F(\mathbf{w}^*)) + \frac{2\gamma^3 \kappa \delta^2}{M} + \frac{\gamma^2 \delta^2}{M} \quad (16)$$

$$\mathbb{E}\|\mathbf{w}^T - \mathbf{w}^*\|^2 \leq \frac{2}{\mu} (1 - \mu\gamma)^T (F(\mathbf{w}^0) - F(\mathbf{w}^*)) + 2\gamma \frac{\delta^2}{\mu^2 M}. \quad (17)$$

LEMMA 14 (BOUNDED ITERATES DIFFERENCE OF  $\alpha$ ). *Let  $\{\alpha_i^r\}$  be iterates generated by Algorithm 2, then under conditions of Theorem 3, the following statement holds:*

$$\|\alpha_i^r - \alpha_i^{r-1}\|^2 \leq 6 \left(1 - \frac{1}{\kappa_g}\right)^{T\alpha} + O\left(\kappa_g^2 L^2 \bar{\zeta}_i(\mathbf{w}^*)\right) \left(\gamma^2 L (1 - \mu\gamma)^r (F(\mathbf{w}^0) - F(\mathbf{w}^*)) + \frac{\gamma^3 \kappa \delta^2}{M} + \frac{\gamma^2 \delta^2}{M}\right)$$

PROOF. Define

$$\mathbf{z}^r = \left[ \|\nabla f_i(\mathbf{w}^r) - \nabla f_1(\mathbf{w}^r)\|^2, \dots, \|\nabla f_i(\mathbf{w}^r) - \nabla f_N(\mathbf{w}^r)\|^2 \right].$$

According to updating rule of  $\alpha$  in Algorithm 2 and Lemma 3 we have:

$$\begin{aligned} \|\alpha_i^r - \alpha_i^{r-1}\|^2 &\leq 3\|\alpha_i^r - \alpha_{g_i}^*(\mathbf{w}^r)\|^2 + 3\|\alpha_{g_i}^*(\mathbf{w}^{r-1}) - \alpha_{g_i}^*(\mathbf{w}^r)\|^2 + 3\|\alpha_{g_i}^*(\mathbf{w}^{r-1}) - \alpha_i^{r-1}\|^2 \\ &\leq 6(1 - \mu_g \eta \alpha)^{T\alpha} + 3\|\alpha_{g_i}^*(\mathbf{w}^{r-1}) - \alpha_{g_i}^*(\mathbf{w}^r)\|^2 \\ &\leq 6(1 - \mu_g \eta \alpha)^{T\alpha} + 3\kappa_g^2 \sum_{j=1}^N \|z_j^{r-1} - z_j^r\|^2 \\ &\leq 6(1 - \mu_g \eta \alpha)^{T\alpha} + 3\kappa_g^2 \sum_{j=1}^N \|\nabla f_i(\mathbf{w}^r) - \nabla f_j(\mathbf{w}^r) + \nabla f_i(\mathbf{w}^{r-1}) - \nabla f_j(\mathbf{w}^{r-1})\|^2 4L^2 \|\mathbf{w}^r - \mathbf{w}^{r-1}\|^2 \end{aligned}$$

where the second inequality follows from (8). Since  $\|\nabla f_i(\mathbf{w}^r) - \nabla f_j(\mathbf{w}^r)\| \leq \|\nabla f_i(\mathbf{w}^*) - \nabla f_j(\mathbf{w}^*)\| + 2L \|\mathbf{w}^r - \mathbf{w}^*\|$ , we can conclude that

$$\begin{aligned} \|\alpha_i^r - \alpha_i^{r-1}\|^2 &\leq 6(1 - \mu_g \eta \alpha)^{T\alpha} \\ &\quad + 12L^2 \kappa_g^2 \sum_{j=1}^N \left(8 \|\nabla f_i(\mathbf{w}^*) - \nabla f_j(\mathbf{w}^*)\|^2 + 8L^2 \|\mathbf{w}^r - \mathbf{w}^*\|^2 + 8L^2 \|\mathbf{w}^{r-1} - \mathbf{w}^*\|^2\right) \|\mathbf{w}^r - \mathbf{w}^{r-1}\|^2 \\ &\leq 6 \left(1 - \frac{1}{\kappa_g}\right)^{T\alpha} + O\left(\kappa_g^2 L^2 \bar{\zeta}_i(\mathbf{w}^*) \|\mathbf{w}^r - \mathbf{w}^{r-1}\|^2\right) \\ &\leq 6 \left(1 - \frac{1}{\kappa_g}\right)^{T\alpha} + O\left(\kappa_g^2 L^2 \bar{\zeta}_i(\mathbf{w}^*)\right) \left(\gamma^2 L (1 - \mu\gamma)^r (F(\mathbf{w}^0) - F(\mathbf{w}^*)) + \frac{\gamma^3 \kappa \delta^2}{M} + \frac{\gamma^2 \delta^2}{M}\right) \end{aligned}$$

where at last step we plug in Proposition 4 (16).  $\square$

LEMMA 15 (CONVERGENCE OF  $\alpha$ ). *Let  $\{\hat{\alpha}_i\}_{i=1}^N$  be the mixing parameters generated by Algorithm 2. Then under the conditions of Theorem 3, the following statement holds:*

$$\|\hat{\alpha}_i - \alpha^*\|^2 \leq 2\left(1 - \frac{1}{\kappa_g}\right)^{T\alpha} + O\left(\kappa_g^2 \bar{\zeta}_i(\mathbf{w}^*) L^2 \frac{2}{\mu} (1 - \mu\gamma)^T + 2\gamma \frac{\delta^2}{\mu^2 M}\right), i \in [N]$$

PROOF. We notice the following decomposition:

$$\begin{aligned} \|\hat{\alpha}_i - \alpha^*\|^2 &= \|\alpha_i^R - \alpha_g^*(\mathbf{w}^*)\|^2 \\ &\leq 2\|\alpha_i^R - \alpha_g^*(\mathbf{w}^R)\|^2 + 2\|\alpha_g^*(\mathbf{w}^R) - \alpha_g^*(\mathbf{w}^*)\|^2 \\ &\leq 2\left(1 - \frac{1}{\kappa_g}\right)^{T\alpha} + O\left(\kappa_g^2 \left(\bar{\zeta}_i(\mathbf{w}^*) + NL^2 \|\mathbf{w}^R - \mathbf{w}^*\|^2\right) 4L \|\mathbf{w}^R - \mathbf{w}^*\|^2\right) \\ &\leq 2\left(1 - \frac{1}{\kappa_g}\right)^{T\alpha} + O\left(\kappa_g^2 \bar{\zeta}_i(\mathbf{w}^*) L^2 \frac{2}{\mu} (1 - \mu\gamma)^T + 2\gamma \frac{\delta^2}{\mu^2 M}\right), \end{aligned}$$

where in the second inequality we apply Lemma 3, and in the third inequality we use Proposition 4 (17).  $\square$

## B.2 Proof of Theorem 3

PROOF. According to Lemma 2, we have:

$$\Phi(\alpha_i^*, \hat{v}_i) - \Phi(\alpha_i^*, v_i^*) \leq L \|\hat{v}_i - v^*(\hat{\alpha}_i)\|^2 + \kappa^2 L \|\hat{\alpha}_i - \alpha_i^*\|^2$$

We first examine the convergence of  $\|\hat{v}_i - v^*(\hat{\alpha}_i)\|^2$ . Applying Cauchy-Schwartz inequality yields:

$$\begin{aligned} \|\mathbf{v}^{r+1} - v^*(\alpha^{r+1})\|^2 &\leq \left(1 + \frac{1}{4a-2}\right) \|\mathbf{v}^{r+1} - v^*(\alpha^r)\|^2 + (1+4a-2) \|v^*(\alpha^{r+1}) - v^*(\alpha^r)\|^2 \\ &\leq \left(1 + \frac{1}{4a-2}\right) \|\mathbf{v}^{r+1} - v^*(\alpha^r)\|^2 + (1+4a-2) \kappa^2 \|\alpha^{r+1} - \alpha^r\|^2 \end{aligned} \quad (18)$$

where  $a = \frac{1}{\mu\eta NK}$ . Similar to the proof of Theorem 2, we first define

$$\begin{aligned} \mathbf{g}^r &:= \sum_{j=1}^N \prod_{j'=N-1}^{j+1} (\mathbf{I} - \mathcal{Q}_{j'} \mathbf{H}_{j'}) \mathcal{Q}_j \nabla f_{\sigma(j)}(\mathbf{v}^r), \\ \delta^r &:= \sum_{j=1}^N \prod_{j'=N-1}^{j+1} (\mathbf{I} - \mathcal{Q}_{j'} \mathbf{H}_{j'}) \delta_j. \end{aligned}$$

Then we recall the updating rule of  $\mathbf{v}$ :

$$\mathbf{v}^{r+1} = \mathcal{P}_{\mathcal{W}}(\mathbf{v}^r - \mathbf{g}^r - \delta^r)$$

Hence we have:

$$\begin{aligned} \mathbb{E} \|\mathbf{v}^{r+1} - v^*(\alpha^r)\|^2 &= \mathbb{E} \|\mathcal{P}_{\mathcal{W}}(\mathbf{v}^r - \mathbf{g}^r - \delta^r - v^*(\alpha^r))\|^2 \\ &\leq \mathbb{E} \|\mathbf{v}^r - \mathbf{g}^r - \delta^r - v^*(\alpha^r)\|^2 \\ &\leq \mathbb{E} \|\mathbf{v}^r - v^*(\alpha^r)\|^2 - 2\mathbb{E} \langle \mathbf{g}^r, \mathbf{v}^r - v^*(\alpha^r) \rangle + \mathbb{E} \|\mathbf{g}^r\|^2 + \mathbb{E} \|\delta^r\|^2 \\ &\leq \mathbb{E} \|\mathbf{v}^r - v^*(\alpha^r)\|^2 - 2\mathbb{E} \langle \eta NK \nabla \Phi(\alpha^r, \mathbf{v}^r), \mathbf{v}^r - v^*(\alpha^r) \rangle \\ &\quad - 2\mathbb{E} \langle \mathbf{g}^r - \eta NK \nabla \Phi(\alpha^r, \mathbf{v}^r), \mathbf{v}^r - v^*(\alpha^r) \rangle + \mathbb{E} \|\mathbf{g}^r\|^2 + \mathbb{E} \|\delta^r\|^2. \end{aligned}$$

Now, applying strongly convexity of  $\Phi(\alpha^r, \cdot)$  and Cauchy-Schwartz inequality yields:

$$\begin{aligned} \mathbb{E} \|\mathbf{v}^{r+1} - v^*(\alpha^r)\|^2 &\leq (1 - \mu\eta NK) \mathbb{E} \|\mathbf{v}^r - v^*(\alpha^r)\|^2 - \eta NK \mathbb{E} [\Phi(\alpha^r, \mathbf{v}^r) - \Phi(\alpha^r, v^*(\hat{\alpha}))] \\ &\quad + \frac{1}{2} \left( \frac{1}{\mu\eta NK} \mathbb{E} \|\mathbf{g}^r - \eta NK \nabla \Phi(\hat{\alpha}, \mathbf{v}^r)\|^2 + \mu\eta NK \mathbb{E} \|\mathbf{v}^r - v^*(\hat{\alpha})\|^2 \right) \\ &\quad + \mathbb{E} \|\mathbf{g}^r\|^2 + \mathbb{E} \|\delta^r\|^2 \\ &\leq \left(1 - \frac{1}{2} \mu\eta NK\right) \mathbb{E} \|\mathbf{v}^r - v^*(\alpha^r)\|^2 - \eta NK \mathbb{E} [\Phi(\alpha^r, \mathbf{v}^r) - \Phi(\alpha^r, v^*(\hat{\alpha}))] \\ &\quad + \frac{1}{2\mu\eta NK} \mathbb{E} \|\mathbf{g}^r - \eta NK \nabla \Phi(\alpha^r, \mathbf{v}^r)\|^2 \\ &\quad + 2\mathbb{E} \|\mathbf{g}^r - \eta NK \nabla \Phi(\alpha^r, \mathbf{v}^r)\|^2 + 2\mathbb{E} \|\eta NK \nabla \Phi(\alpha^r, \mathbf{v}^r)\|^2 + \mathbb{E} \|\delta^r\|^2. \end{aligned}$$

where in the first inequality we applied Cauchy-Schwartz inequality and strongly convexity. Since  $\Phi(\alpha^r, \cdot)$  is  $L$  smooth, we have:  $\mathbb{E} \|\nabla \Phi(\alpha^r, \mathbf{v}^r)\|^2 \leq 2L \mathbb{E} [\Phi(\alpha^r, \mathbf{v}^r) - \Phi(\alpha^r, v^*(\alpha^r))]$ . Therefore, it follows that:

$$\begin{aligned} \mathbb{E} \|\mathbf{v}^{r+1} - v^*(\hat{\alpha})\|^2 &\leq \left(1 - \frac{1}{2} \mu\eta NK\right) \mathbb{E} \|\mathbf{v}^r - v^*(\alpha^r)\|^2 - (\eta NK - 4\eta^2 N^2 K^2 L) \mathbb{E} [\Phi(\alpha^r, \mathbf{v}^r) - \Phi(\alpha^r, v^*(\alpha^r))] \\ &\quad + \left(\frac{1}{2\mu\eta NK} + 2\right) \mathbb{E} \|\mathbf{g}^r - \eta NK \nabla \Phi(\alpha^r, \mathbf{v}^r)\|^2 + \mathbb{E} \|\delta^r\|^2 \end{aligned} \quad (19)$$



Now, we examine the term  $\|g^r - \eta NK \nabla \Phi(\alpha^r, v^r)\|^2$  in the right hand side of above inequality. First according to summation by part (Lemma 8): we let  $A_j := \prod_{j'=N-1}^{j+1} (\mathbf{I} - \mathbf{Q}_{j'} \mathbf{H}_{j'})$  and  $B_j = \mathbf{Q}_j \nabla f_{\sigma(j)}(v^r)$ , then we have:

$$\begin{aligned} g^r &= \sum_{j=1}^N \prod_{j'=N}^{j+1} (\mathbf{I} - \mathbf{Q}_{j'} \mathbf{H}_{j'}) \mathbf{Q}_j \nabla f_{\sigma(j)}(v^r) \\ &= \sum_{j=1}^N A_j B_j = \sum_{j=1}^N \mathbf{Q}_j \nabla f_{\sigma(j)}(v^r) - \sum_{n=1}^{N-1} \left( \prod_{j'=N}^{n+2} (\mathbf{I} - \mathbf{Q}_{j'} \mathbf{H}_{j'}) - \prod_{j'=N}^{n+1} (\mathbf{I} - \mathbf{Q}_{j'} \mathbf{H}_{j'}) \right) \sum_{j=1}^n \mathbf{Q}_j \nabla f_{\sigma(j)}(v^r) \\ &= \sum_{j=1}^N \mathbf{Q}_j \nabla f_{\sigma(j)}(v^r) - \sum_{n=1}^{N-1} \left( \prod_{j'=N}^{n+2} (\mathbf{I} - \mathbf{Q}_{j'} \mathbf{H}_{j'}) \right) \mathbf{Q}_{n+1} \mathbf{H}_{n+1} \sum_{j=1}^n \mathbf{Q}_j \nabla f_{\sigma(j)}(v^r). \end{aligned}$$

Hence we have:

$$\begin{aligned} &\|g^r - \eta NK \nabla \Phi(\hat{\alpha}, v^r)\|^2 \\ &= \left\| \eta NK \nabla \Phi(\hat{\alpha}, v^r) - \sum_{j=1}^N \prod_{j'=N-1}^{j+1} (\mathbf{I} - \mathbf{Q}_{j'} \mathbf{H}_{j'}) \mathbf{Q}_j \nabla f_{\sigma(j)}(v^r) \right\|^2 \\ &= \left\| \eta NK \sum_{j=1}^N \hat{\alpha}(\sigma(j)) \nabla f_{\sigma(j)}(v^r) - \left( \sum_{j=1}^N \mathbf{Q}_j \nabla f_{\sigma(j)}(v^r) - \sum_{n=1}^{N-1} \left( \prod_{j'=N}^{n+2} (\mathbf{I} - \mathbf{Q}_{j'} \mathbf{H}_{j'}) \right) \mathbf{Q}_{n+1} \mathbf{H}_{n+1} \sum_{j=0}^n \mathbf{Q}_j \nabla f_{\sigma(j)}(v^r) \right) \right\|^2 \\ &\stackrel{(1)}{\leq} 2 \left\| \left( \eta NK \sum_{j=1}^N \hat{\alpha}(\sigma(j)) \nabla f_{\sigma(j)}(v^r) - \sum_{j=1}^N \mathbf{Q}_j \nabla f_{\sigma(j)}(v^r) \right) \right\|^2 \\ &\quad + 2 \left\| \sum_{n=1}^{N-1} \left( \prod_{j'=N}^{n+2} (\mathbf{I} - \mathbf{Q}_{j'} \mathbf{H}_{j'}) \right) \mathbf{Q}_{n+1} \mathbf{H}_{n+1} \sum_{j=1}^n \mathbf{Q}_j \nabla f_{\sigma(j)}(v^r) \right\|^2 \\ &\stackrel{(2)}{\leq} \left( 20\eta^2 N^2 K^2 \left( \frac{e}{4R-e} \right)^2 + 36e^6 \eta^4 N^4 K^4 L^4 \right) \|\nabla \Phi(\hat{\alpha}, v^r)\|^2 + 256\eta^2 N^3 K^2 \left( \frac{e}{4R-e} \right)^2 G^2 \log(1/p) \\ &\quad + 244e^6 \eta^4 N^4 K^4 L^4 G^2 N \log(1/p) \\ &\stackrel{(3)}{\leq} \left( 20\eta^2 N^2 K^2 \left( \frac{e}{4R-e} \right)^2 + 36e^6 \eta^4 N^4 K^4 L^4 \right) 2L (\Phi(\hat{\alpha}, v^r) - \Phi(\hat{\alpha}, v^*(\hat{\alpha}))) \\ &\quad + \left( 244e^6 \eta^4 N^4 K^4 L^4 + 256\eta^2 N^3 K^2 \left( \frac{e}{4R-e} \right)^2 \right) G^2 N \log(1/p) \end{aligned}$$

where in (1) we apply Jensen's inequality, in (2) we plug in Lemma 11 (a), and Lemma 12, and in (3) we use the  $L$ -smoothness of  $\Phi$ . Plugging above bound back in (19) yields:

$$\begin{aligned} &\mathbb{E} \|v^{r+1} - v^*(\alpha^r)\|^2 \\ &\leq \left( 1 - \frac{1}{2} \mu \eta NK \right) \mathbb{E} \|v^r - v^*(\alpha^r)\|^2 + \eta^2 N^2 K^2 e^4 \delta^2 \\ &\quad - \left( \eta NK - 4\eta^2 N^2 K^2 L - \left( \frac{1}{2\mu \eta NK} + 2 \right) \left( 20\eta^2 N^2 K^2 \left( \frac{e}{4R-e} \right)^2 - 36e^6 \eta^4 N^4 K^4 L^4 \right) \right) \\ &\quad \times \mathbb{E} [\Phi(\alpha^r, v^r) - \Phi(\alpha^r, v^*(\alpha^r))] \\ &\quad + \left( \frac{1}{2\mu \eta NK} + 2 \right) \left( 244e^6 \eta^4 N^4 K^4 L^4 + 256\eta^2 N^3 K^2 \left( \frac{e}{4R-e} \right)^2 \right) G^2 N \log(1/p). \end{aligned}$$

Since we choose  $\eta = \frac{4 \log(\sqrt{NKR})}{\mu NKR}$ , and

$$R \geq \max \left\{ \frac{3}{8} e, \sqrt[3]{\frac{64\kappa^2 \log(\sqrt{NKR}) e^6}{9\mu}} \right\}$$

Hence we have:

$$\begin{aligned}
& \mathbb{E} \|\mathbf{v}^{r+1} - \mathbf{v}^*(\boldsymbol{\alpha}^r)\|^2 \\
& \leq \left(1 - \frac{1}{2}\mu\eta NK\right) \mathbb{E} \|\mathbf{v}^r - \mathbf{v}^*(\boldsymbol{\alpha}^r)\|^2 + \eta^2 N^2 K^2 e^4 \delta^2 - \frac{1}{2}\eta NK \underbrace{\mathbb{E}[\Phi(\hat{\boldsymbol{\alpha}}, \mathbf{v}^r) - \Phi(\hat{\boldsymbol{\alpha}}, \mathbf{v}^*(\hat{\boldsymbol{\alpha}}))]}_{\geq 0} \\
& \quad + \left(\frac{1}{2\mu\eta NK} + 2\right) \left(244e^6 \eta^4 N^4 K^4 L^4 + 256\eta^2 N^3 K^2 \left(\frac{e}{4R-e}\right)^2\right) G^2 N \log(1/p) \\
& \leq \left(1 - \frac{1}{2}\mu\eta NK\right) \mathbb{E} \|\mathbf{v}^r - \mathbf{v}^*(\hat{\boldsymbol{\alpha}})\|^2 + \eta^2 N^2 K^2 e^4 \delta^2 \\
& \quad + \left(\frac{1}{2\mu\eta NK} + 2\right) \left(244e^6 \eta^4 N^4 K^4 L^4 + 256\eta^2 N^3 K^2 \left(\frac{e}{4R-e}\right)^2\right) G^2 N \log(1/p).
\end{aligned}$$

Putting above inequality back to (18) yields:

$$\begin{aligned}
\|\mathbf{v}^{r+1} - \mathbf{v}^*(\boldsymbol{\alpha}^{r+1})\|^2 & \leq \left(1 - \frac{1}{4a}\right) \|\mathbf{v}^{r+1} - \mathbf{v}^*(\boldsymbol{\alpha}^r)\|^2 + 2\eta^2 N^2 K^2 e^4 \delta^2 + 4a \|\mathbf{v}^*(\boldsymbol{\alpha}^{r+1}) - \mathbf{v}^*(\boldsymbol{\alpha}^r)\|^2 \\
& \quad + 2 \left(\frac{1}{2\mu\eta NK} + 2\right) \left(244e^6 \eta^4 N^4 K^4 L^4 + 256\eta^2 N^3 K^2 \left(\frac{e}{4R-e}\right)^2\right) G^2 N \log(1/p) \\
& \leq \left(1 - \frac{1}{4a}\right) \|\mathbf{v}^{r+1} - \mathbf{v}^*(\boldsymbol{\alpha}^r)\|^2 + 2\eta^2 N^2 K^2 e^4 \delta^2 \\
& \quad + 2 \left(\frac{1}{2\mu\eta NK} + 2\right) \left(244e^6 \eta^4 N^4 K^4 L^4 + 256\eta^2 N^3 K^2 \left(\frac{e}{4R-e}\right)^2\right) G^2 N \log(1/p) \\
& \quad + O\left(\frac{1}{\mu\eta NK} \left( \left(1 - \frac{1}{\kappa_g}\right)^{T_\alpha} + \kappa_g^2 L^2 \bar{\zeta}_i(\mathbf{w}^*) \left( \gamma^2 L (1 - \mu\gamma)^r DG + \frac{\gamma^3 \kappa \delta^2}{M} + \frac{\gamma^2 \delta^2}{M} \right) \right)\right)
\end{aligned}$$

where at second inequality we plug in Lemma 14. Unrolling the recursion from  $r = R$  to 0, and plugging in  $\eta = \frac{4 \log(NKR^3)}{\mu NK R}$  yields:

$$\begin{aligned}
& \|\mathbf{v}^R - \mathbf{v}^*(\boldsymbol{\alpha}^R)\|^2 \\
& \leq \left(1 - \frac{1}{4}\mu\eta NK\right)^R \|\mathbf{v}^0 - \mathbf{v}^*(\boldsymbol{\alpha}^0)\|^2 + \frac{1}{\mu}\eta NK e^4 \delta^2 \\
& \quad + 8 \frac{1}{\mu} \left(\frac{1}{2\mu\eta NK} + 2\right) \left(244e^6 \eta^3 N^3 K^3 L^4 + 256\eta N^2 K \left(\frac{e}{4R-e}\right)^2\right) G^2 N \log(1/p) \\
& \quad + O\left(\frac{1}{\mu\eta NK} \sum_{r=0}^R \left(1 - \frac{1}{4a}\right)^{R-r} \left( \left(1 - \frac{1}{\kappa_g}\right)^{T_\alpha} + \kappa_g^2 L^2 \bar{\zeta}_i(\mathbf{w}^*) \left( \gamma^2 L (1 - \mu\gamma)^r DG + \frac{\gamma^3 \kappa \delta^2}{M} + \frac{\gamma^2 \delta^2}{M} \right) \right)\right) \\
& \leq O\left(\frac{\|\mathbf{v}^0 - \mathbf{v}^*(\boldsymbol{\alpha}^0)\|^2}{NKR^3}\right) + \tilde{O}\left(\left(\frac{\kappa^4}{R^2} + \frac{N}{\mu^2 R^2}\right) G^2 N \log(1/p) + \frac{\delta^2}{\mu R}\right) \\
& \quad + O\left(\frac{1}{\mu\eta NK} \sum_{r=0}^R \left(1 - \frac{\log(NKR^3)}{R}\right)^{R-r} \left( \left(1 - \frac{1}{\kappa_g}\right)^{T_\alpha} + \kappa_g^2 L^2 \bar{\zeta}_i(\mathbf{w}^*) \left( \gamma^2 L (1 - \mu\gamma)^r DG + \frac{\gamma^3 \kappa \delta^2}{M} + \frac{\gamma^2 \delta^2}{M} \right) \right)\right)
\end{aligned}$$

Plugging in  $\gamma = \frac{\log(NKR^3)}{\mu R}$  yields:

$$\begin{aligned}
\|\mathbf{v}^R - \mathbf{v}^*(\boldsymbol{\alpha}^R)\|^2 &\leq O\left(\frac{\|\mathbf{v}^0 - \mathbf{v}^*(\boldsymbol{\alpha}^0)\|^2}{NKR^3}\right) + \tilde{O}\left(\left(\frac{\kappa^4}{R^2} + \frac{N}{\mu^2 R^2}\right)G^2N \log(1/p) + \frac{\delta^2}{\mu R}\right) \\
&\quad + \tilde{O}\left(R\left(\gamma^2 L^2 \sum_{r=0}^R \left(1 - \frac{\log(NKR^3)}{R}\right)^R \kappa_g^2 L^2 \bar{\zeta}_i(\mathbf{w}^*) DG\right.\right. \\
&\quad \left.\left. + \sum_{r=0}^R \left(1 - \frac{\log(NKR^3)}{R}\right)^{R-r} \left(\left(1 - \frac{1}{\kappa_g}\right)^{T_\alpha} + \frac{\kappa_g^2 L^2 \bar{\zeta}_i(\mathbf{w}^*) \gamma^2 \delta^2}{M}\right)\right)\right) \\
&\leq O\left(\frac{D^2}{NKR^3}\right) + \tilde{O}\left(\left(\frac{\kappa^4}{R^2} + \frac{N}{\mu^2 R^2}\right)G^2N \log(1/p) + \frac{\delta^2}{\mu R}\right) \\
&\quad + \tilde{O}\left(\frac{\kappa^2 \kappa_g^2 L^2 \bar{\zeta}_i(\mathbf{w}^*) DG}{R} + R^2 \left(\left(1 - \frac{1}{\kappa_g}\right)^{T_\alpha} + \frac{\kappa_g^2 \kappa^2 \bar{\zeta}_i(\mathbf{w}^*) \delta^2}{\mu^2 MR^2}\right)\right)
\end{aligned}$$

Since  $\hat{\mathbf{v}}_i = \mathbf{v}^R$  and  $\hat{\boldsymbol{\alpha}}_i = \boldsymbol{\alpha}^R$ , we have the convergence of  $\|\hat{\mathbf{v}}_i - \mathbf{v}^*(\hat{\boldsymbol{\alpha}}_i)\|^2$ . Plugging this convergence rate together with the convergence of  $\|\hat{\boldsymbol{\alpha}}_i - \boldsymbol{\alpha}^*\|^2$  from Lemma 15:

$$\begin{aligned}
\|\hat{\boldsymbol{\alpha}}_i - \boldsymbol{\alpha}^*\|^2 &\leq O\left(2\left(1 - \frac{1}{\kappa_g}\right)^{T_\alpha} + O\left(\kappa_g^2 \bar{\zeta}_i(\mathbf{w}^*) L^2 \frac{2}{\mu} (1 - \mu\gamma)^R + 2\gamma \frac{\delta^2}{\mu^2 M}\right)\right) \\
&\leq \tilde{O}\left(\left(1 - \frac{1}{\kappa_g}\right)^{T_\alpha} + O\left(\kappa_g^2 \bar{\zeta}_i(\mathbf{w}^*) L^2 \frac{2}{\mu R} + \frac{\delta^2}{\mu^3 RM}\right)\right)
\end{aligned}$$

together with applying Lemma 2 leads to:

$$\begin{aligned}
\Phi(\boldsymbol{\alpha}_i^*, \hat{\mathbf{v}}_i) - \Phi(\boldsymbol{\alpha}_i^*, \mathbf{v}_i^*) &\leq L \|\hat{\mathbf{v}}_i - \mathbf{v}^*(\hat{\boldsymbol{\alpha}}_i)\|^2 + \kappa^2 L \|\hat{\boldsymbol{\alpha}}_i - \boldsymbol{\alpha}_i^*\|^2 \\
&\leq O\left(L \frac{D^2}{NKR^3}\right) + \tilde{O}\left(\left(\frac{\kappa^4 L}{R^2} + \frac{NL}{\mu^2 R^2}\right)G^2N \log(1/p) + \frac{L\delta^2}{\mu R}\right) \\
&\quad + \tilde{O}\left(\frac{\kappa^2 \kappa_g^2 L^3 \bar{\zeta}_i(\mathbf{w}^*) DG}{R} + LR^2 \left(1 - \frac{1}{\kappa_g}\right)^{T_\alpha} + \frac{L\kappa_g^2 \kappa^2 \bar{\zeta}_i(\mathbf{w}^*) \delta^2}{\mu^2 M}\right) \\
&\quad + \kappa^2 L \tilde{O}\left(\left(1 - \frac{1}{\kappa_g}\right)^{T_\alpha} + \frac{\kappa_g^2 \bar{\zeta}_i(\mathbf{w}^*) \kappa L}{R} + \frac{\delta^2}{\mu^3 RM}\right) \\
&\leq O\left(L \frac{D^2}{NKR^3}\right) + \tilde{O}\left(\left(\frac{\kappa^4 L}{R^2} + \frac{NL}{\mu^2 R^2}\right)G^2N \log(1/p) + \frac{L\delta^2}{\mu R}\right) \\
&\quad + \tilde{O}\left(\frac{\kappa^2 \kappa_g^2 L^3 \bar{\zeta}_i(\mathbf{w}^*) DG}{R} + LR^2 \left(1 - \frac{1}{\kappa_g}\right)^{T_\alpha} + \frac{L\kappa_g^2 \kappa^2 \bar{\zeta}_i(\mathbf{w}^*) \delta^2}{\mu^2 M}\right).
\end{aligned}$$

thus completing the proof as desired.  $\square$

## C ADDITIONAL EXPERIMENTS

In addition to the synthetic dataset discussed in the main body, we run experiments on the EMNIST dataset [3], which is naturally distributed in a federated setting. In this case, we chose 50 clients and use a 2-layer MLP model, each with 200 neurons. We compare the PERM algorithm with the localized model in FedAvg and perFedAvg [8]. As it can be seen in Figure 2, PERM can learn a better personalized model by attending to each client's data according to the similarity of the data distribution between clients. The learned values of  $\boldsymbol{\alpha}$  show that the clients are learning from each others' data, and not focused on their own data only. This signifies that the distribution of data among clients in this dataset is not highly heterogeneous. Note that, since we are using a subset of clients in the EMNIST dataset for the training (only 50 clients for 100 rounds of communication), the results would be sub-optimal. Nonetheless, the experiments are designed to show the effectiveness of different algorithms.

### C.1 The effectiveness of learned mixture weights

To show the effectiveness of two-stage PERM algorithm, as well as the effects of heterogeneity on distribution of data among clients on the learned weights  $\boldsymbol{\alpha}$  in the algorithm, we run this algorithm on MNIST dataset. We use 50 clients, and the model is an MLP, similar to the EMNIST experiment. In this case, once we distribute the data randomly across clients (homogeneous) and once only allocating 1 class per client (highly heterogeneous). As it can be seen from Figure 4, when the data distribution is homogeneous the learned values of  $\boldsymbol{\alpha}$  as

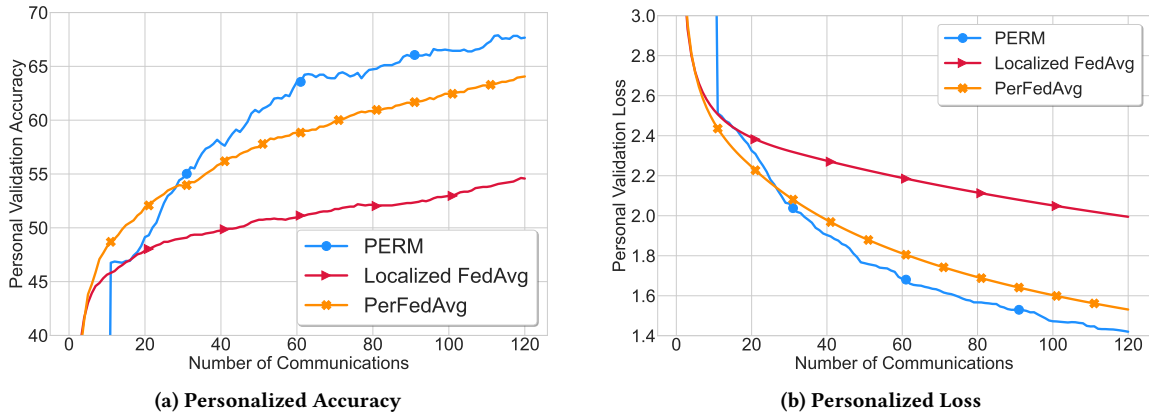


Figure 2: Comparative Analysis of Personalization methods, including our single-loop PERM algorithm, localized FedAvg, and perFedAg, with EMNIST dataset. The disparity in personalized accuracy and loss highlights PERM’s capability in leveraging relevant client correlations.

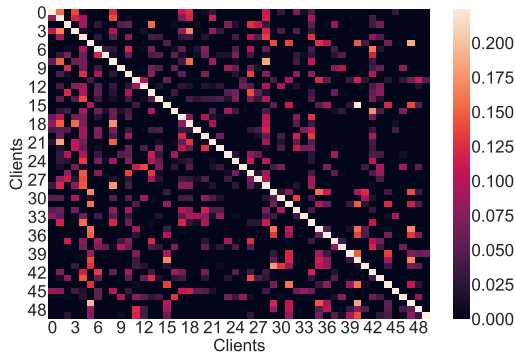


Figure 3: The heat map of the learned  $\alpha$  values for the PERM algorithms. The weights signify that clients mutually benefiting from one another’s data, which also highlight that the distribution of data is not significantly heterogeneous in this dataset.

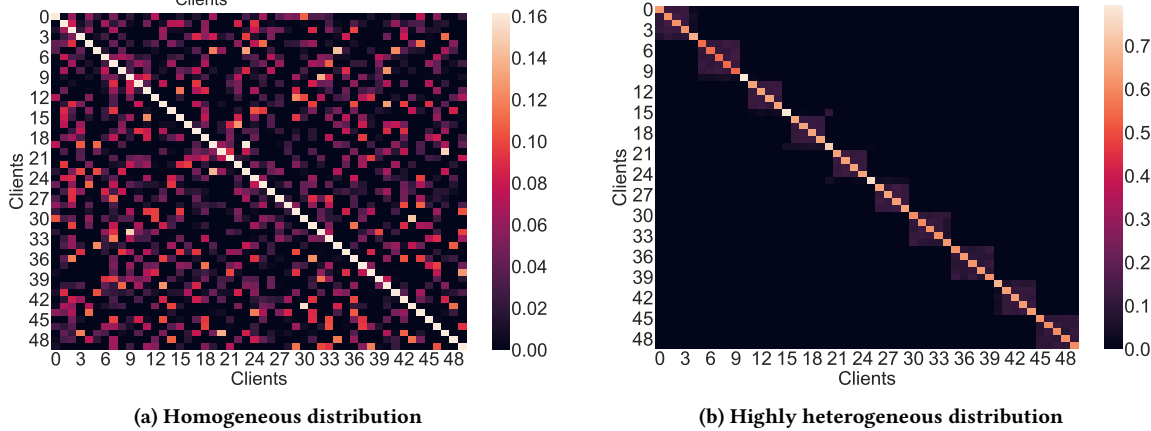


Figure 4: Comparing the performance of two-stage PERM algorithm in learning  $\alpha$  values on heterogeneous and homogeneous data distributions. WE use MNIST dataset across 50 clients with homogeneous and heterogeneous distributions.

diffused across clients. However, when the data is highly heterogeneous, the learned  $\alpha$  values will be highly sparse, indicating that each client is mostly learning from its own data and some other clients with partial distribution similarity.