

# SCAFFOLD: Stochastic Controlled Averaging for Federated Learning

Sai Praneeth Karimireddy<sup>1,2</sup> Satyen Kale<sup>3</sup> Mehryar Mohri<sup>3,4</sup> Sashank J. Reddi<sup>3</sup> Sebastian U. Stich<sup>1</sup>  
Ananda Theertha Suresh<sup>3</sup>

## Abstract

Federated Averaging (FEDAVG) has emerged as the algorithm of choice for federated learning due to its simplicity and low communication cost. However, in spite of recent research efforts, its performance is not fully understood. We obtain tight convergence rates for FEDAVG and prove that it suffers from ‘client-drift’ when the data is heterogeneous (non-iid), resulting in unstable and slow convergence.

As a solution, we propose a new algorithm (SCAFFOLD) which uses control variates (variance reduction) to correct for the ‘client-drift’ in its local updates. We prove that SCAFFOLD requires significantly fewer communication rounds and is not affected by data heterogeneity or client sampling. Further, we show that (for quadratics) SCAFFOLD can take advantage of similarity in the client’s data yielding even faster convergence. The latter is the first result to quantify the usefulness of local-steps in distributed optimization.

## 1. Introduction

Federated learning has emerged as an important paradigm in modern large-scale machine learning. Unlike in traditional centralized learning where models are trained using large datasets stored in a central server (Dean et al., 2012; Iandola et al., 2016; Goyal et al., 2017), in federated learning, the training data remains distributed over a large number of clients, which may be phones, network sensors, hospitals, or alternative local information sources (Konečný et al., 2016b;a; McMahan et al., 2017; Mohri et al., 2019; Kairouz et al., 2019). A centralized model (referred to as server model) is then trained without ever transmitting

<sup>1</sup>EPFL, Lausanne <sup>2</sup>Based on work performed at Google Research, New York. <sup>3</sup>Google Research, New York <sup>4</sup>Courant Institute, New York. Correspondence to: Sai Praneeth Karimireddy <sai.karimireddy@epfl.ch>.

*Proceedings of the 37<sup>th</sup> International Conference on Machine Learning*, Online, PMLR 119, 2020. Copyright 2020 by the author(s).

client data over the network, thereby ensuring a basic level of privacy. In this work, we investigate stochastic optimization algorithms for federated learning.

The key challenges for federated optimization are 1) dealing with unreliable and relatively slow network connections between the server and the clients, 2) only a small subset of clients being available for training at a given time, and 3) large heterogeneity (non-iid-ness) in the data present on the different clients (Konečný et al., 2016a). The most popular algorithm for this setting is FEDAVG (McMahan et al., 2017). FEDAVG tackles the communication bottleneck by performing multiple local updates on the available clients before communicating to the server. While it has shown success in certain applications, its performance on heterogeneous data is still an active area of research (Li et al., 2018; Yu et al., 2019; Li et al., 2019b; Haddadpour & Mahdavi, 2019; Khaled et al., 2020). We prove that indeed such heterogeneity has a large effect on FEDAVG—it introduces a *drift* in the updates of each client resulting in slow and unstable convergence. Further, we show that this client-drift persists even if full batch gradients are used and all clients participate throughout the training.

As a solution, we propose a new Stochastic Controlled Averaging algorithm (SCAFFOLD) which tries to correct for this client-drift. Intuitively, SCAFFOLD estimates the update direction for the server model ( $c$ ) and the update direction for each client  $c_i$ .<sup>1</sup> The difference  $(c - c_i)$  is then an estimate of the client-drift which is used to correct the local update. This strategy successfully overcomes heterogeneity and converges in significantly fewer rounds of communication. Alternatively, one can see heterogeneity as introducing ‘client-variance’ in the updates across the different clients and SCAFFOLD then performs ‘client-variance reduction’ (Schmidt et al., 2017; Johnson & Zhang, 2013; Defazio et al., 2014). We use this viewpoint to show that SCAFFOLD is relatively unaffected by client sampling.

Finally, while accommodating heterogeneity is important, it is equally important that a method can take advantage of similarities in the client data. We prove that SCAFFOLD indeed has such a property, requiring fewer rounds of com-

<sup>1</sup>We refer to these estimates as *control variates* and the resulting correction technique as stochastic controlled averaging.

munication when the clients are more similar.

**Contributions.** We summarize our main results below.

- We derive tighter convergence rates for FEDAVG than previously known for convex and non-convex functions with client sampling and heterogeneous data.
- We give matching lower bounds to prove that even with no client sampling and full batch gradients, FEDAVG can be slower than SGD due to client-drift.
- We propose a new Stochastic Controlled Averaging algorithm (SCAFFOLD) which corrects for this client-drift. We prove that SCAFFOLD is at least as fast as SGD and converges for arbitrarily heterogeneous data.
- We show SCAFFOLD can additionally take advantage of similarity between the clients to further reduce the communication required, proving the advantage of taking local steps over large-batch SGD for the first time.
- We prove that SCAFFOLD is relatively unaffected by the client sampling obtaining variance reduced rates, making it especially suitable for federated learning.

Finally, we confirm our theoretical results on simulated and real datasets (extended MNIST by Cohen et al. (2017)).

**Related work.** For identical clients, FEDAVG coincides with parallel SGD analyzed by (Zinkevich et al., 2010) who proved asymptotic convergence. Stich (2018) and, more recently Stich & Karimireddy (2019); Patel & Dieuleveut (2019); Khaled et al. (2020), gave a sharper analysis of the same method, under the name of local SGD, also for identical functions. However, there still remains a gap between their upper bounds and the lower bound of Woodworth et al. (2018). The analysis of FEDAVG for heterogeneous clients is more delicate due to the afore-mentioned client-drift, first empirically observed by Zhao et al. (2018). Several analyses bound this drift by assuming bounded gradients (Wang et al., 2019; Yu et al., 2019), or view it as additional noise (Khaled et al., 2020), or assume that the client optima are  $\epsilon$ -close (Li et al., 2018; Haddadpour & Mahdavi, 2019). In a concurrent work, (Liang et al., 2019) propose to use variance reduction to deal with client heterogeneity but still show rates slower than SGD and do not support client sampling. Our method SCAFFOLD can also be seen as an improved version of the distributed optimization algorithm DANE by (Shamir et al., 2014), where a fixed number of (stochastic) gradient steps are used in place of a proximal point update. A more in-depth discussion of related work is given in Appendix A. We summarize the complexities of different methods for heterogeneous clients in Table 2.

## 2. Setup

We formalize the problem as minimizing a sum of stochastic functions, with only access to stochastic samples:

Table 1. Summary of notation used in the paper

$N, S$ , and $i$	total num., sampled num., and index of clients
$R, r$	number, index of communication rounds
$K, k$	number, index of local update steps
$\mathbf{x}^r$	aggregated server model after round $r$
$\mathbf{y}_{i,k}^r$	$i$ th client's model in round $r$ and step $k$
$c^r, c_i^r$	control variate of server, $i$ th client after round $r$

$$\min_{\mathbf{x} \in \mathbb{R}^d} \left\{ f(\mathbf{x}) := \frac{1}{N} \sum_{i=1}^N (f_i(\mathbf{x}) := \mathbb{E}_{\zeta_i}[f_i(\mathbf{x}; \zeta_i)]) \right\}.$$

The functions  $f_i$  represents the loss function on client  $i$ . All our results can be easily extended to the weighted case.

We assume that  $f$  is bounded from below by  $f^*$  and  $f_i$  is  $\beta$ -smooth. Further, we assume  $g_i(\mathbf{x}) := \nabla f_i(\mathbf{x}; \zeta_i)$  is an unbiased stochastic gradient of  $f_i$  with variance bounded by  $\sigma^2$ . For some results, we assume  $\mu \geq 0$  (strong) convexity. Note that  $\sigma$  only bounds the variance *within* clients. We also define two non-standard terminology below.

**(A1)  $(G, B)$ -BHD** or bounded gradient dissimilarity: there exist constants  $G \geq 0$  and  $B \geq 1$  such that

$$\frac{1}{N} \sum_{i=1}^N \|\nabla f_i(\mathbf{x})\|^2 \leq G^2 + B^2 \|\nabla f(\mathbf{x})\|^2, \forall \mathbf{x}.$$

If  $\{f_i\}$  are convex, we can relax the assumption to

$$\frac{1}{N} \sum_{i=1}^N \|\nabla f_i(\mathbf{x})\|^2 \leq G^2 + 2\beta B^2 (f(\mathbf{x}) - f^*), \forall \mathbf{x}.$$

**(A2)  $\delta$ -BHD** or bounded Hessian dissimilarity:

$$\|\nabla^2 f_i(\mathbf{x}) - \nabla^2 f(\mathbf{x})\| \leq \delta, \forall \mathbf{x}.$$

Further,  $f_i$  is  $\delta$ -weakly convex i.e.  $\nabla^2 f_i(\mathbf{x}) \succeq -\delta I$ .

The assumptions A1 and A2 are orthogonal—it is possible to have  $G = 0$  and  $\delta = 2\beta$ , or  $\delta = 0$  but  $G \gg 1$ .

## 3. Convergence of FedAvg

In this section we review FEDAVG and improve its convergence analysis by deriving tighter rates than known before. The scheme consists of two main parts: local updates to the model (1), and aggregating the client updates to update the server model (2). In each round, a subset of clients  $\mathcal{S} \subseteq [N]$  are sampled uniformly. Each of these clients  $i \in \mathcal{S}$  copies the current sever model  $\mathbf{y}_i = \mathbf{x}$  and performs  $K$  local updates of the form:

$$\mathbf{y}_i \leftarrow \mathbf{y}_i - \eta_l g_i(\mathbf{y}_i). \quad (1)$$

Here  $\eta_l$  is the local step-size. Then the clients' updates  $(\mathbf{y}_i - \mathbf{x})$  are aggregated to form the new server model using a global step-size  $\eta_g$  as:

$$\mathbf{x} \leftarrow \mathbf{x} + \frac{\eta_g}{|\mathcal{S}|} \sum_{i \in \mathcal{S}} (\mathbf{y}_i - \mathbf{x}). \quad (2)$$

*Table 2.* Number of communication rounds required to reach  $\epsilon$  accuracy for  $\mu$  strongly convex and non-convex functions (log factors are ignored). Set  $\mu = \epsilon$  for general convex rates.  $(G, B)$  bounds gradient dissimilarity (A1), and  $\delta$  bounds Hessian dissimilarity (A2). Our rates for FEDAVG are more general and tighter than others, even matching the lower bound. However, SGD is still faster ( $B \geq 1$ ). SCAFFOLD does not require any assumptions, is faster than SGD, and is robust to client sampling. Further, when clients become more similar (small  $\delta$ ), SCAFFOLD converges even faster.

Method	Strongly convex	Non-convex	Sampling	Assumptions
SGD (large batch)	$\frac{\sigma^2}{\mu NK\epsilon} + \frac{1}{\mu}$	$\frac{\sigma^2}{NK\epsilon^2} + \frac{1}{\epsilon}$	×	-
FedAvg				
(Li et al., 2019b)	$\frac{\sigma^2}{\mu^2 NK\epsilon} + \frac{G^2 K}{\mu^2 \epsilon}$	-	×	$(G, 0)$ -BGD
(Yu et al., 2019)	-	$\frac{\sigma^2}{NK\epsilon^2} + \frac{G^2 NK}{\epsilon}$	×	$(G, 0)$ -BGD
(Khaled et al., 2020)	$\frac{\sigma^2 + G^2}{\mu NK\epsilon} + \frac{\sigma + G}{\mu\sqrt{\epsilon}} + \frac{NB^2}{\mu}$	-	×	$(G, B)$ -BGD
Ours (Thm. I) <sup>1</sup>	$\frac{M^2}{\mu SK\epsilon} + \frac{G}{\mu\sqrt{\epsilon}} + \frac{B^2}{\mu}$	$\frac{M^2}{SK\epsilon^2} + \frac{G}{\epsilon^{3/2}} + \frac{B^2}{\epsilon}$	✓	$(G, B)$ -BGD
Lower-bound (Thm. II)	$\Omega\left(\frac{\sigma^2}{\mu NK\epsilon} + \frac{G}{\sqrt{\mu\epsilon}}\right)$	?	×	$(G, 1)$ -BGD
FedProx (Li et al., 2018) <sup>2</sup>	$\frac{B^2}{\mu}$	$\frac{B^2}{\epsilon}$ (weakly convex)	✓	$\sigma = 0, (0, B)$ -BGD
DANE (Shamir et al., 2014) <sup>2,3</sup>	$\frac{\delta^2}{\mu^2}$	-	×	$\sigma = 0, \delta$ -BHD
VRL-SGD (Liang et al., 2019)	-	$\frac{N\sigma^2}{K\epsilon^2} + \frac{N}{\epsilon}$	×	-
SCAFFOLD				
Theorem III	$\frac{\sigma^2}{\mu SK\epsilon} + \frac{1}{\mu} + \frac{N}{S}$	$\frac{\sigma^2}{SK\epsilon^2} + \frac{1}{\epsilon}\left(\frac{N}{S}\right)^{\frac{2}{3}}$	✓	-
Theorem IV <sup>3</sup>	$\frac{\sigma^2}{\mu NK\epsilon} + \frac{1}{\mu K} + \frac{\delta}{\mu}$	$\frac{\sigma^2}{NK\epsilon^2} + \frac{1}{K\epsilon} + \frac{\delta}{\epsilon}$	×	$\delta$ -BHD

<sup>1</sup>  $M^2 := \sigma^2 + K(1 - \frac{S}{N})G^2$ . Note that  $\frac{M^2}{S} = \frac{\sigma^2}{N}$  when no sampling ( $S = N$ ).

<sup>2</sup> proximal point method i.e.  $K \gg 1$ .

<sup>3</sup> proved only for quadratic functions.

### 3.1. Rate of convergence

We now state our novel convergence results for functions with bounded dissimilarity (proofs in Appendix D.2).

**Theorem I.** For  $\beta$ -smooth functions  $\{f_i\}$  which satisfy (A1), the output of FEDAVG has expected error smaller than  $\epsilon$  for some values of  $\eta_l, \eta_g, R$  satisfying

- **Strongly convex:**  $\eta_g \geq \sqrt{S}$ ,  $\eta_l \leq \frac{1}{(1+B^2)6\beta K\eta_g}$ , and

$$R = \tilde{\mathcal{O}}\left(\frac{\sigma^2}{\mu KS\epsilon} + (1 - \frac{S}{N})\frac{G^2}{\mu S\epsilon} + \frac{\sqrt{\beta}G}{\mu\sqrt{\epsilon}} + \frac{B^2\beta}{\mu}\right),$$

- **General convex:**  $\eta_g \geq \sqrt{S}$ ,  $\eta_l \leq \frac{1}{(1+B^2)6\beta K\eta_g}$ , and

$$R = \mathcal{O}\left(\frac{\sigma^2 D^2}{KS\epsilon^2} + (1 - \frac{S}{N})\frac{G^2 D^2}{S\epsilon^2} + \frac{\sqrt{\beta}G}{\epsilon^{\frac{3}{2}}} + \frac{B^2\beta D^2}{\epsilon}\right),$$

- **Non-convex:**  $\eta_g \geq \sqrt{S}$ ,  $\eta_l \leq \frac{1}{(1+B^2)6\beta K\eta_g}$ , and

$$R = \mathcal{O}\left(\frac{\beta\sigma^2 F}{KS\epsilon^2} + (1 - \frac{S}{N})\frac{G^2 F}{S\epsilon^2} + \frac{\sqrt{\beta}G}{\epsilon^{\frac{3}{2}}} + \frac{B^2\beta F}{\epsilon}\right),$$

where  $D := \|\mathbf{x}^0 - \mathbf{x}^*\|^2$  and  $F := f(\mathbf{x}^0) - f^*$ .

It is illuminating to compare our rates with those of the simpler iid. case i.e. with  $G = 0$  and  $B = 1$ . Our strongly-convex rates become  $\frac{\sigma^2}{\mu SK\epsilon} + \frac{1}{\mu}$ . In comparison, the best previously known rate for this case was by Stich & Karimireddy (2019) who show a rate of  $\frac{\sigma^2}{\mu SK\epsilon} + \frac{S}{\mu}$ . The main source of improvement in the rates came from the use of two separate step-sizes ( $\eta_l$  and  $\eta_g$ ). By having a larger global step-size  $\eta_g$ , we can use a smaller local step-size  $\eta_l$  thereby reducing the client-drift while still ensuring progress. However, even our improved rates do not match the lower-bound for the identical case of  $\frac{\sigma^2}{\mu SK\epsilon} + \frac{1}{K\mu}$  (Woodworth et al., 2018). We bridge this gap for quadratic functions in Section 6.

We now compare FEDAVG to two other algorithms FedProx by (Li et al., 2018) (aka EASGD by (Zhang et al., 2015)) and to SGD. Suppose that  $G = 0$  and  $\sigma = 0$  i.e. we use full batch gradients and all clients have very similar optima. In such a case, FEDAVG has a complexity of  $\frac{B^2}{\mu}$  which is identical to that of FedProx (Li et al., 2018). Thus, FedProx does not have any theoretical advantage.

Next, suppose that all clients participate (no sampling) with  $S = N$  and there is no variance  $\sigma = 0$ . Then, the

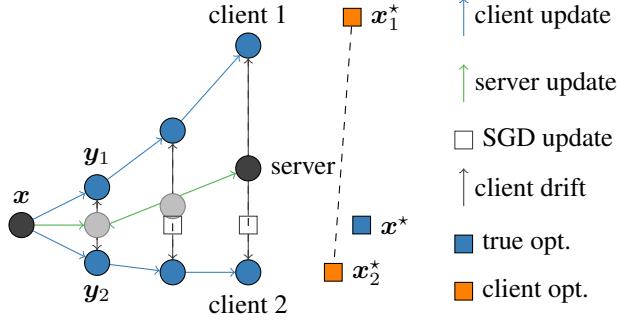


Figure 1. Client-drift in FEDAVG is illustrated for 2 clients with 3 local steps ( $N = 2, K = 3$ ). The local updates  $y_i$  (in blue) move towards the individual client optima  $x_i^*$  (orange square). The server updates (in red) move towards  $\frac{1}{N} \sum_i x_i^*$  instead of to the true optimum  $x^*$  (black square).

above for strongly-convex case simplifies to  $\frac{G}{\mu\sqrt{\epsilon}} + \frac{B^2}{\mu}$ . In comparison, extending the proof of (Khaled et al., 2020) using our techniques gives a worse dependence on  $G$  of  $\frac{G^2}{\mu KN\epsilon} + \frac{G}{\mu\sqrt{\epsilon}}$ . Similarly, for the non-convex case, our rates are tighter and have better dependence on  $G$  than (Yu et al., 2019). However, simply running SGD in this setting would give a communication complexity of  $\frac{\beta}{\mu}$  which is faster, and independent of similarity assumptions. In the next section we examine the necessity of such similarity assumptions.

### 3.2. Lower bounding the effect of heterogeneity

We now show that when the functions  $\{f_i\}$  are distinct, the local updates of FEDAVG on each client experiences *drift* thereby slowing down convergence. We show that the amount of this client drift, and hence the slowdown in the rate of convergence, is exactly determined by the gradient dissimilarity parameter  $G$  in (A1).

We now examine the mechanism by which the client-drift arises (see Fig. 1). Let  $x^*$  be the global optimum of  $f(\mathbf{x})$  and  $x_i^*$  be the optimum of each client's loss function  $f_i(\mathbf{x})$ . In the case of heterogeneous data, it is quite likely that each of these  $x_i^*$  is far away from the other, and from the global optimum  $x^*$ . Even if all the clients start from the same point  $x$ , each of the  $y_i$  will move towards their client optimum  $x_i^*$ . This means that the average of the client updates (which is the server update) moves towards  $\frac{1}{N} \sum_{i=1}^N x_i^*$ . This difference between  $\frac{1}{N} \sum_{i=1}^N x_i^*$  and the true optimum  $x^*$  is exactly the cause of client-drift. To counter this drift, FEDAVG is forced to use much smaller step-sizes which in turn hurts convergence. We can formalize this argument to prove a lower-bound (see Appendix D.4 for proof).

**Theorem II.** For any positive constants  $G$  and  $\mu$ , there exist  $\mu$ -strongly convex functions satisfying A1 for which FEDAVG with  $K \geq 2, \sigma = 0$  and  $N = S$  has an error

**Algorithm 1** SCAFFOLD: Stochastic Controlled Averaging for federated learning

```

1: server input: initial  $\mathbf{x}$  and  $\mathbf{c}$ , and global step-size  $\eta_g$ 
2: client  $i$ 's input:  $\mathbf{c}_i$ , and local step-size  $\eta_l$ 
3: for each round  $r = 1, \dots, R$  do
4:   sample clients  $\mathcal{S} \subseteq \{1, \dots, N\}$ 
5:   communicate  $(\mathbf{x}, \mathbf{c})$  to all clients  $i \in \mathcal{S}$ 
6:   on client  $i \in \mathcal{S}$  in parallel do
7:     initialize local model  $\mathbf{y}_i \leftarrow \mathbf{x}$ 
8:     for  $k = 1, \dots, K$  do
9:       compute mini-batch gradient  $g_i(\mathbf{y}_i)$ 
10:       $\mathbf{y}_i \leftarrow \mathbf{y}_i - \eta_l (g_i(\mathbf{y}_i) - \mathbf{c}_i + \mathbf{c})$ 
11:    end for
12:     $\mathbf{c}_i^+ \leftarrow$  (i)  $g_i(\mathbf{x})$ , or (ii)  $\mathbf{c}_i - \mathbf{c} + \frac{1}{K\eta_l} (\mathbf{x} - \mathbf{y}_i)$ 
13:    communicate  $(\Delta\mathbf{y}_i, \Delta\mathbf{c}_i) \leftarrow (\mathbf{y}_i - \mathbf{x}, \mathbf{c}_i^+ - \mathbf{c}_i)$ 
14:     $\mathbf{c}_i \leftarrow \mathbf{c}_i^+$ 
15:  end on client
16:   $(\Delta\mathbf{x}, \Delta\mathbf{c}) \leftarrow \frac{1}{|\mathcal{S}|} \sum_{i \in \mathcal{S}} (\Delta\mathbf{y}_i, \Delta\mathbf{c}_i)$ 
17:   $\mathbf{x} \leftarrow \mathbf{x} + \eta_g \Delta\mathbf{x}$  and  $\mathbf{c} \leftarrow \mathbf{c} + \frac{|\mathcal{S}|}{N} \Delta\mathbf{c}$ 
18: end for

```

$$f(\mathbf{x}^r) - f(\mathbf{x}^*) \geq \Omega\left(\frac{G^2}{\mu R^2}\right).$$

This implies that the  $\frac{G}{\sqrt{\epsilon}}$  term is unavoidable even if there is no stochasticity. Further, because FEDAVG uses  $RKN$  stochastic gradients, we also have the statistical lower-bound of  $\frac{\sigma^2}{\mu KN\epsilon}$ . Together, these lower bounds prove that the rate derived in Theorem I is nearly optimal (up to dependence on  $\mu$ ). In the next section, we introduce a new method SCAFFOLD to mitigate this client-drift.

### 4. SCAFFOLD algorithm

In this section we first describe SCAFFOLD and then discuss how it solves the problem of client-drift.

**Method.** SCAFFOLD has three main steps: local updates to the client model (3), local updates to the client control variate (4), and aggregating the updates (5). We describe each in more detail.

Along with the server model  $\mathbf{x}$ , SCAFFOLD maintains a state for each client (client control variate  $\mathbf{c}_i$ ) and for the server (server control variate  $\mathbf{c}$ ). These are initialized to ensure that  $\mathbf{c} = \frac{1}{N} \sum \mathbf{c}_i$  and can safely all be initialized to 0. In each round of communication, the server parameters  $(\mathbf{x}, \mathbf{c})$  are communicated to the participating clients  $\mathcal{S} \subset [N]$ . Each participating client  $i \in \mathcal{S}$  initializes its local model with the server model  $\mathbf{y}_i \leftarrow \mathbf{x}$ . Then it makes a pass over its local data performing  $K$  updates of the form:

$$\mathbf{y}_i \leftarrow \mathbf{y}_i - \eta_l (g_i(\mathbf{y}_i) + \mathbf{c} - \mathbf{c}_i). \quad (3)$$

Then, the local control variate  $\mathbf{c}_i$  is also updated. For this,

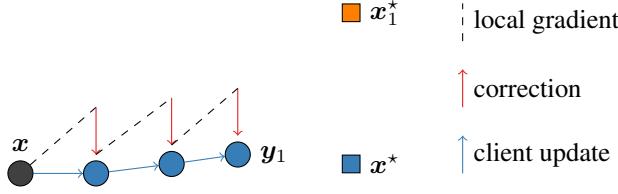


Figure 2. Update steps of SCAFFOLD on a single client. The local gradient (dashed black) points to  $x_1^*$  (orange square), but the correction term ( $c - c_i$ ) (in red) ensures the update moves towards the true optimum  $x^*$  (black square).

we provide two options:

$$c_i^+ \leftarrow \begin{cases} \text{Option I. } g_i(\mathbf{x}), \text{ or} \\ \text{Option II. } c_i - c + \frac{1}{K\eta_l}(\mathbf{x} - \mathbf{y}_i). \end{cases} \quad (4)$$

Option I involves making an additional pass over the local data to compute the gradient at the server model  $\mathbf{x}$ . Option II instead re-uses the previously computed gradients to update the control variate. Option I can be more stable than II depending on the application, but II is cheaper to compute and usually suffices (all our experiments use Option II). The client updates are then aggregated and used to update the server parameters:

$$\begin{aligned} \mathbf{x} &\leftarrow \mathbf{x} + \frac{\eta_g}{|\mathcal{S}|} \sum_{i \in \mathcal{S}} (\mathbf{y}_i - \mathbf{x}), \\ \mathbf{c} &\leftarrow \mathbf{c} + \frac{1}{N} \sum_{i \in \mathcal{S}} (c_i^+ - c_i). \end{aligned} \quad (5)$$

This finishes one round of communication. Note that the clients in SCAFFOLD are *stateful* and retain the value of  $c_i$  across multiple rounds. Further, if  $c_i$  is always set to 0, then SCAFFOLD becomes equivalent to FEDAVG. The full details are summarized in Algorithm 1.

**Usefulness of control variates.** If communication cost was not a concern, the ideal update on client  $i$  would be

$$\mathbf{y}_i \leftarrow \mathbf{y}_i + \frac{1}{N} \sum_j g_j(\mathbf{y}_i). \quad (6)$$

Such an update essentially computes an unbiased gradient of  $f$  and hence becomes equivalent to running FEDAVG in the iid case (which has excellent performance). Unfortunately such an update requires communicating with all clients for every update step. SCAFFOLD instead uses control variates such that

$$\mathbf{c}_j \approx g_j(\mathbf{y}_i) \text{ and } \mathbf{c} \approx \frac{1}{N} \sum_j g_j(\mathbf{y}_i).$$

Then, SCAFFOLD (3) mimics the ideal update (6) with

$$(g_i(\mathbf{y}_i) - \mathbf{c}_i + \mathbf{c}) \approx \frac{1}{N} \sum_j g_j(\mathbf{y}_i).$$

Thus, the local updates of SCAFFOLD remain synchronized and converge for arbitrarily heterogeneous clients.

## 5. Convergence of SCAFFOLD

We state the rate of SCAFFOLD without making any assumption on the similarity between the functions. See Appendix E for the full proof.

**Theorem III.** For any  $\beta$ -smooth functions  $\{f_i\}$ , the output of SCAFFOLD has expected error smaller than  $\epsilon$  for  $\eta_g = \sqrt{S}$  and some values of  $\eta_l, R$  satisfying

- **Strongly convex:**  $\eta_l \leq \min\left(\frac{1}{81\beta K\eta_g}, \frac{S}{15\mu N K\eta_g}\right)$  and

$$R = \tilde{\mathcal{O}}\left(\frac{\sigma^2}{\mu K S \epsilon} + \frac{\beta}{\mu} + \frac{N}{S}\right),$$

- **General convex:**  $\eta_l \leq \frac{1}{81\beta K\eta_g}$  and

$$R = \tilde{\mathcal{O}}\left(\frac{\sigma^2 D^2}{K S \epsilon^2} + \frac{\beta D^2}{\epsilon} + \frac{N F}{S}\right),$$

- **Non-convex:**  $\eta_l \leq \frac{1}{24K\eta_g\beta} \left(\frac{S}{N}\right)^{\frac{2}{3}}$  and

$$R = \mathcal{O}\left(\frac{\beta\sigma^2 F}{K S \epsilon^2} + \left(\frac{N}{S}\right)^{\frac{2}{3}} \frac{\beta F}{\epsilon}\right),$$

where  $D := \|\mathbf{x}^0 - \mathbf{x}^*\|^2$  and  $F := f(\mathbf{x}^0) - f^*$ .

Let us first examine the rates without client sampling ( $S = N$ ). For the strongly convex case, the number of rounds becomes  $\frac{\sigma^2}{\mu N K \epsilon} + \frac{1}{\mu}$ . This rate holds for arbitrarily heterogeneous clients unlike Theorem I and further matches that of SGD with  $K$  times larger batch-size, proving that SCAFFOLD is at least as fast as SGD. These rates also match known lower-bounds for distributed optimization (Arjevani & Shamir, 2015) (up to acceleration) and are unimprovable in general. However in certain cases SCAFFOLD is provably faster than SGD. We show this fact in Section 6.

Now let  $\sigma = 0$ . Then our rates in the strongly-convex case are  $\frac{1}{\mu} + \frac{N}{S}$  and  $\left(\frac{N}{S}\right)^{\frac{2}{3}} \frac{1}{\epsilon}$  in the non-convex case. These exactly match the rates of SAGA (Defazio et al., 2014; Reddi et al., 2016c). In fact, when  $\sigma = 0$ ,  $K = 1$  and  $S = 1$ , the update of SCAFFOLD with option I reduces to SAGA where in each round consists of sampling one client  $f_i$ . Thus SCAFFOLD can be seen as an extension of variance reduction techniques for federated learning, and one could similarly extend SARAH (Nguyen et al., 2017), SPIDER (Fang et al., 2018), etc. Note that standard SGD with client sampling is provably slower and converges at a sub-linear rate even with  $\sigma = 0$ .

**Proof sketch.** For simplicity, assume that  $\sigma = 0$  and consider the ideal update of (6) which uses the full gradient  $\nabla f(\mathbf{y})$  every step. Clearly, this would converge at a linear rate even with  $S = 1$ . FEDAVG would instead use an

update  $\nabla f_i(\mathbf{y})$ . The difference between the ideal update (6) and the FEDAVG update (1) is  $\|\nabla f_i(\mathbf{y}) - \nabla f(\mathbf{y})\|$ . We need a bound on the gradient-dissimilarity as in (A1) to bound this error. SCAFFOLD instead uses the update  $\nabla f_i(\mathbf{y}) - \mathbf{c}_i + \mathbf{c}$ , and the difference from ideal update becomes

$$\sum_i \|(\nabla f_i(\mathbf{y}) - \mathbf{c}_i + \mathbf{c}) - \nabla f(\mathbf{y})\|^2 \leq \sum_i \|\mathbf{c}_i - \nabla f_i(\mathbf{y})\|^2.$$

Thus, the error is independent of how similar or dissimilar the functions  $f_i$  are, and instead only depends on the quality of our approximation  $\mathbf{c}_i \approx \nabla f_i(\mathbf{y})$ . Since  $f_i$  is smooth, we can expect that the gradient  $\nabla f_i(\mathbf{y})$  does not change too fast and hence is easy to approximate. Appendix E translates this intuition into a formal proof.

## 6. Usefulness of local steps

In this section we investigate when and why taking local steps might be useful over simply computing a large-batch gradient in distributed optimization. We will show that when the functions across the clients share some similarity, local steps can take advantage of this and converge faster. For this we consider quadratic functions and express their similarity with the  $\delta$  parameter introduced in (A2).

**Theorem IV.** *For any  $\beta$ -smooth quadratic functions  $\{f_i\}$  with  $\delta$  bounded Hessian dissimilarity (A2), the output of SCAFFOLD with  $S = N$  (no sampling) has error smaller than  $\epsilon$  for  $\eta_g = 1$  and some values of  $\eta_l, R$  satisfying*

- **Strongly convex:**  $\eta_l \leq \frac{1}{15K\delta+8\beta}$  and

$$R = \tilde{\mathcal{O}}\left(\frac{\beta\sigma^2}{\mu KN\epsilon} + \frac{\beta + \delta K}{\mu K}\right),$$

- **Weakly convex:**  $\eta_l \leq \frac{1}{15K\delta+8\beta}$  and

$$R = \mathcal{O}\left(\frac{\beta\sigma^2 F}{KN\epsilon^2} + \frac{(\beta + \delta K)F}{K\epsilon}\right),$$

where we define  $F := (f(\mathbf{x}^0) - f^*)$ .

When  $\sigma = 0$  and  $K$  is large, the complexity of SCAFFOLD becomes  $\frac{\delta}{\mu}$ . In contrast DANE, which being a proximal point method also uses large  $K$ , requires  $(\frac{\delta}{\mu})^2$  rounds (Shamir et al., 2014) which is significantly slower, or needs an additional backtracking-line search to match the rates of SCAFFOLD (Yuan & Li, 2019). Further, Theorem IV is the first result to demonstrate improvement due to similarity for non-convex functions as far as we are aware.

Suppose that  $\{f_i\}$  are identical. Recall that  $\delta$  in (A2) measures the Hessian dissimilarity between functions and so  $\delta = 0$  for this case. Then Theorem IV shows that the complexity of SCAFFOLD is  $\frac{\sigma^2}{\mu KN\epsilon} + \frac{1}{\mu K}$  which (up to acceleration) matches the i.i.d. lower bound of (Woodworth

et al., 2018). In contrast, SGD with  $K$  times larger batch-size would require  $\frac{\sigma^2}{\mu KN\epsilon} + \frac{1}{\mu}$  (note the absence of  $K$  in the second term). Thus, for identical functions, SCAFFOLD (and in fact even FEDAVG) improves linearly with increasing number of local steps. In the other extreme, if the functions are arbitrarily different, we may have  $\delta = 2\beta$ . In this case, the complexity of SCAFFOLD and large-batch SGD match the lower bound of Arjevani & Shamir (2015) for the heterogeneous case.

The above insights can be generalized to when the functions are only somewhat similar. If the Hessians are  $\delta$ -close and  $\sigma = 0$ , then the complexity is  $\frac{\beta + \delta K}{\mu K}$ . This bound implies that the optimum number of local steps one should use is  $K = \frac{\beta}{\delta}$ . Picking a smaller  $K$  increases the communication required whereas increasing it further would only waste computational resources. While this result is intuitive—if the functions are more ‘similar’, local steps are more useful—Theorem IV shows that it is the similarity of the *Hessians* which matters. This is surprising since the Hessians of  $\{f_i\}$  may be identical even if their individual optima  $\mathbf{x}_i^*$  are arbitrarily far away from each other and the gradient-dissimilarity (A1) is unbounded.

**Proof sketch.** Consider a simplified SCAFFOLD update with  $\sigma = 0$  and no sampling ( $S = N$ ):

$$\mathbf{y}_i = \mathbf{y}_i - \eta(\nabla f_i(\mathbf{y}_i) + \nabla f(\mathbf{x}) - \nabla f_i(\mathbf{x})).$$

We would ideally want to perform the update  $\mathbf{y}_i = \mathbf{y}_i - \eta \nabla f(\mathbf{y}_i)$  using the full gradient  $\nabla f(\mathbf{y}_i)$ . We reinterpret the correction term of SCAFFOLD ( $\mathbf{c} - \mathbf{c}_i$ ) as performing the following first order correction to the local gradient  $\nabla f_i(\mathbf{y}_i)$  to make it closer to the full gradient  $\nabla f(\mathbf{y}_i)$ :

$$\begin{aligned} & \underbrace{\nabla f_i(\mathbf{y}_i) - \nabla f_i(\mathbf{x})}_{\approx \nabla^2 f_i(\mathbf{x})(\mathbf{y}_i - \mathbf{x})} + \underbrace{\nabla f(\mathbf{x})}_{\approx \nabla f(\mathbf{y}_i) + \nabla^2 f(\mathbf{x})(\mathbf{x} - \mathbf{y}_i)} \\ & \approx \nabla f(\mathbf{y}_i) + (\nabla^2 f_i(\mathbf{x}) - \nabla^2 f(\mathbf{x}))(\mathbf{y}_i - \mathbf{x}) \\ & \approx \nabla f(\mathbf{y}_i) + \delta(\mathbf{y}_i - \mathbf{x}) \end{aligned}$$

Thus the SCAFFOLD update approximates the ideal update up to an error  $\delta$ . This intuition is proved formally for quadratic functions in Appendix F. Generalizing these results to other functions is a challenging open problem.

## 7. Experiments

We run experiments on both simulated and real datasets to confirm our theory. Our main findings are i) SCAFFOLD consistently outperforms SGD and FEDAVG across all parameter regimes, and ii) the benefit (or harm) of local steps depends on both the algorithm and the similarity of the clients data.

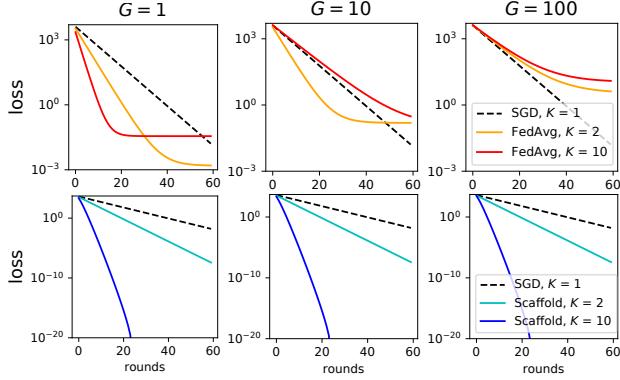


Figure 3. SGD (dashed black), FedAvg (above), and SCAFFOLD (below) on simulated data. FedAvg gets worse as local steps increases with  $K = 10$  (red) worse than  $K = 2$  (orange). It also gets slower as gradient-dissimilarity ( $G$ ) increases (to the right). SCAFFOLD significantly improves with more local steps, with  $K = 10$  (blue) faster than  $K = 2$  (light blue) and SGD. Its performance is identical as we vary heterogeneity ( $G$ ).

## 7.1. Setup

Our simulated experiments uses  $N = 2$  quadratic functions based on our lower-bounds in Theorem II. We use full-batch gradients ( $\sigma = 0$ ) and no client sampling. Our real world experiments run logistic regression (convex) and 2 layer fully connected network (non-convex) on the EMNIST (Cohen et al., 2017). We divide this dataset among  $N = 100$  clients as follows: for  $s\%$  similar data we allocate to each client  $s\%$  i.i.d. data and the remaining  $(100 - s)\%$  by sorting according to label (cf. Hsu et al. (2019)).

We consider four algorithms: SGD, FEDAVG, SCAFFOLD and FEDPROX with SGD as the local solver (Li et al., 2018). On each client SGD uses the full local data to compute a single update, whereas the other algorithms take 5 steps per epoch (batch size is 0.2 of local data). We always use global step-size  $\eta_g = 1$  and tune the local step-size  $\eta_l$  individually for each algorithm. SCAFFOLD uses option II (no extra gradient computations) and FEDPROX has fixed regularization = 1 to keep comparison fair. Additional tuning of the regularization parameter may sometimes yield improved empirical performance.

## 7.2. Simulated results

The results are summarized in Fig. 3. Our simulated data has Hessian difference  $\delta = 1$  (A2) and  $\beta = 1$ . We vary the gradient heterogeneity (A1) as  $G \in [1, 10, 100]$ . For all values of  $G$ , FEDAVG gets slower as we increase the number of local steps. This is explained by the fact that client-drift increases as we increase the number of local steps, hindering progress. Further, as we increase  $G$ , FEDAVG continues to slow down exactly as dictated by Thms. I and II. Note that when heterogeneity is small ( $G = \beta = 1$ ), FEDAVG can be competitive with SGD.

SCAFFOLD is consistently faster than SGD, with  $K = 2$  being twice as fast and  $K = 10$  about 5 times faster. Further, its convergence is completely unaffected by  $G$ , confirming our theory in Thm. III. The former observation that we do not see linear improvement with  $K$  is explained by Thm. IV since we have  $\delta > 0$ . This sub linear improvement is still significantly faster than both SGD and FEDAVG.

## 7.3. EMNIST results

We run extensive experiments on the EMNIST dataset to measure the interplay between the algorithm, number of epochs (local updates), number of participating clients, and the client similarity. Table 3 measures the benefit (or harm) of using more local steps, Table 4 studies the resilience to client sampling, and Table 5 reports preliminary results on neural networks. We are mainly concerned with minimizing the number of *communication rounds*. We observe that

**SCAFFOLD is consistently the best.** Across all range of values tried, we observe that SCAFFOLD outperforms SGD, FEDAVG, and FEDPROX. The latter FEDPROX is always slower than the other local update methods, though in some cases it outperforms SGD. Note that it is possible to improve FEDPROX by carefully tuning the regularization parameter (Li et al., 2018). FEDAVG is always slower than SCAFFOLD and faster than FEDPROX.

**SCAFFOLD > SGD > FedAvg for heterogeneous clients.** When similarity is 0%, FEDAVG gets slower with increasing local steps. If we take more than 5 epochs, its performance is worse than SGD's. SCAFFOLD initially worsens as we increase the number of epochs but then flattens. However, its performance is always better than that of SGD, confirming that it can handle heterogeneous data.

**SCAFFOLD and FedAvg get faster with more similarity, but not SGD.** As similarity of the clients increases, the performance of SGD remains relatively constant. On the other hand, SCAFFOLD and FEDAVG get significantly faster as similarity increases. Further, local steps become much more useful, showing monotonic improvement with the increase in number of epochs. This is because with increasing the i.i.d.ness of the data, both the gradient and Hessian dissimilarity decrease.

**SCAFFOLD is resilient to client sampling.** As we decrease the fraction of clients sampled, SCAFFOLD ,and FEDAVG only show a sub-linear slow-down. They are more resilient to sampling in the case of higher similarity.

**SCAFFOLD outperforms FedAvg on non-convex experiments.** We see that SCAFFOLD is better than FEDAVG in terms of final test accuracy reached, though interestingly FEDAVG seems better than SGD even when similarity is 0. However, much more extensive experiments (beyond current scope) are needed before drawing conclusions.

*Table 3.* Communication rounds to reach 0.5 test accuracy for logistic regression on EMNIST as we vary number of epochs. 1k+ indicates 0.5 accuracy was not reached even after 1k rounds, and similarly an arrowhead indicates that the barplot extends beyond the table. 1 epoch for local update methods corresponds to 5 local steps (0.2 batch size), and 20% of clients are sampled each round. We fix  $\mu = 1$  for FEDPROX and use variant (ii) for SCAFFOLD to ensure all methods are comparable. Across all parameters (epochs and similarity), SCAFFOLD is the fastest method. When similarity is 0 (sorted data), FEDAVG consistently gets worse as we increase the number of epochs, quickly becoming slower than SGD. SCAFFOLD initially gets worse and later stabilizes, but is always at least as fast as SGD. As similarity increases (i.e. data is more shuffled), both FEDAVG and SCAFFOLD significantly outperform SGD though SCAFFOLD is still better than FEDAVG. Further, with higher similarity, both methods benefit from increasing number of epochs.

	Epochs	0% similarity (sorted)		10% similarity		100% similarity (i.i.d.)	
		Num. of rounds	Speedup	Num. of rounds	Speedup	Num. of rounds	Speedup
SGD	1	317	(1×)	365	(1×)	416	(1×)
SCAFFOLD1	77	(4.1×)		62	(5.9×)	60	(6.9×)
	5	(2.1×)		20	(18.2×)	10	(41.6×)
	10	(1.1×)		16	(22.8×)	7	(59.4×)
	20	(1.2×)	11	(33.2×)	4	(104×)	
FEDAVG	1	(1.2×)		74	(4.9×)	83	(5×)
	5	(0.7×)		34	(10.7×)	10	(41.6×)
	10	(0.4×)		25	(14.6×)	6	(69.3×)
	20	(< 0.3×)		18	(20.3×)	4	(104×)
FEDPROX	1	(< 0.3×)		979	(0.4×)	459	(0.9×)
	5	(< 0.3×)		794	(0.5×)	351	(1.2×)
	10	(< 0.3×)		894	(0.4×)	308	(1.4×)
	20	(< 0.3×)		916	(0.4×)	351	(1.2×)

*Table 4.* Communication rounds to reach 0.45 test accuracy for logistic regression on EMNIST as we vary the number of sampled clients. Number of epochs is kept fixed to 5. SCAFFOLD is consistently faster than FEDAVG. As we decrease the number of clients sampled in each round, the increase in number of rounds is sub-linear. This slow-down is better for more similar clients.

	Clients	0% similarity	10% similarity
SCAFFOLD	20%	143 (1.0×)	9 (1.0×)
	5%	290 (2.0×)	13 (1.4×)
	1%	790 (5.5×)	28 (3.1×)
FEDAVG	20%	179 (1.0×)	12 (1.0×)
	5%	334 (1.9×)	17 (1.4×)
	1%	1k+ (5.6+×)	35 (2.9×)

## 8. Conclusion

Our work studied the impact of heterogeneity on the performance of optimization methods for federated learning. Our careful theoretical analysis showed that FEDAVG can be severely hampered by *gradient dissimilarity*, and can be even slower than SGD. We then proposed a new stochastic algorithm (SCAFFOLD) which overcomes gradient dissimilarity using control variates. We demonstrated the effectiveness of SCAFFOLD via strong convergence guar-

*Table 5.* Best test accuracy after 1k rounds with 2-layer fully connected neural network (non-convex) on EMNIST trained with 5 epochs per round (25 steps) for the local methods, and 20% of clients sampled each round. SCAFFOLD has the best accuracy and SGD has the least. SCAFFOLD again outperforms other methods. SGD is unaffected by similarity, whereas the local methods improve with client similarity.

	0% similarity	10% similarity
SGD	0.766	0.764
FEDAVG	0.787	0.828
SCAFFOLD	<b>0.801</b>	<b>0.842</b>

antees and empirical evaluations. Further, we showed that while SCAFFOLD is always at least as fast as SGD, it can be much faster depending on the *Hessian dissimilarity* in our data. Thus, different algorithms can take advantage of (and are limited by) different notions of dissimilarity. We believe that characterizing and isolating various dissimilarities present in real world data can lead to further new algorithms and significant impact on distributed, federated, and decentralized learning.

**Acknowledgments.** We thank Filip Hanzely and Jakub Konečný for discussions regarding variance reduction techniques and Blake Woodworth, Virginia Smith and Kumar Kshitij Patel for suggestions which improved the writing.

## References

- Acar, D. A. E., Zhao, Y., Navarro, R. M., Mattina, M., Whatmough, P. N., and Saligrama, V. Federated learning based on dynamic regularization. In *International Conference on Learning Representations*, 2021.
- Agarwal, N., Suresh, A. T., Yu, F. X., Kumar, S., and McMahan, B. cpSGD: Communication-efficient and differentially-private distributed SGD. In *Proceedings of NeurIPS*, pp. 7575–7586, 2018.
- Arjevani, Y. and Shamir, O. Communication complexity of distributed convex learning and optimization. In *Advances in neural information processing systems*, pp. 1756–1764, 2015.
- Bassily, R., Smith, A., and Thakurta, A. Private empirical risk minimization: Efficient algorithms and tight error bounds. In *2014 IEEE 55th Annual Symposium on Foundations of Computer Science*, pp. 464–473. IEEE, 2014.
- Basu, D., Data, D., Karakus, C., and Diggavi, S. Qsparse-local-SGD: Distributed SGD with quantization, sparsification, and local computations. *arXiv preprint arXiv:1906.02367*, 2019.
- Bonawitz, K., Ivanov, V., Kreuter, B., Marcedone, A., McMahan, H. B., Patel, S., Ramage, D., Segal, A., and Seth, K. Practical secure aggregation for privacy-preserving machine learning. In *Proceedings of the 2017 ACM SIGSAC Conference on Computer and Communications Security*, pp. 1175–1191. ACM, 2017.
- Brisimi, T. S., Chen, R., Mela, T., Olshevsky, A., Paschalidis, I. C., and Shi, W. Federated learning of predictive models from federated electronic health records. *International journal of medical informatics*, 112:59–67, 2018.
- Cen, S., Zhang, H., Chi, Y., Chen, W., and Liu, T.-Y. Convergence of distributed stochastic variance reduced methods without sampling extra data. *arXiv preprint arXiv:1905.12648*, 2019.
- Chaudhuri, K., Monteleoni, C., and Sarwate, A. D. Differentially private empirical risk minimization. *Journal of Machine Learning Research*, 12(Mar):1069–1109, 2011.
- Chen, M., Mathews, R., Ouyang, T., and Beaufays, F. Federated learning of out-of-vocabulary words. *arXiv preprint arXiv:1903.10635*, 2019a.
- Chen, M., Suresh, A. T., Mathews, R., Wong, A., Beaufays, F., Allauzen, C., and Riley, M. Federated learning of N-gram language models. In *Proceedings of the 23rd Conference on Computational Natural Language Learning (CoNLL)*, 2019b.
- Cohen, G., Afshar, S., Tapson, J., and Van Schaik, A. Emnist: Extending mnist to handwritten letters. In *2017 International Joint Conference on Neural Networks (IJCNN)*, pp. 2921–2926. IEEE, 2017.
- Dean, J., Corrado, G., Monga, R., Chen, K., Devin, M., Mao, M., Ranzato, M., Senior, A., Tucker, P., Yang, K., Le, Q. V., and Ng, A. Y. Large scale distributed deep networks. In *Advances in neural information processing systems*, pp. 1223–1231, 2012.
- Defazio, A. and Bottou, L. On the ineffectiveness of variance reduced optimization for deep learning. In *Advances in Neural Information Processing Systems*, pp. 1753–1763, 2019.
- Defazio, A., Bach, F., and Lacoste-Julien, S. SAGA: A fast incremental gradient method with support for non-strongly convex composite objectives. In *Advances in neural information processing systems*, pp. 1646–1654, 2014.
- Fang, C., Li, C. J., Lin, Z., and Zhang, T. SPIDER: Near-optimal non-convex optimization via stochastic path-integrated differential estimator. In *Advances in Neural Information Processing Systems*, pp. 689–699, 2018.
- Glasserman, P. *Monte Carlo methods in financial engineering*, volume 53. Springer Science & Business Media, 2013.
- Goyal, P., Dollár, P., Girshick, R., Noordhuis, P., Wesolowski, L., Kyrola, A., Tulloch, A., Jia, Y., and He, K. Accurate, large minibatch SGD: Training imagenet in 1 hour. *arXiv preprint arXiv:1706.02677*, 2017.
- Haddadpour, F. and Mahdavi, M. On the convergence of local descent methods in federated learning. *arXiv preprint arXiv:1910.14425*, 2019.
- Hanzely, F. and Richtárik, P. One method to rule them all: Variance reduction for data, parameters and many new methods. *arXiv preprint arXiv:1905.11266*, 2019.
- Hard, A., Rao, K., Mathews, R., Beaufays, F., Augenstein, S., Eichner, H., Kiddon, C., and Ramage, D. Federated learning for mobile keyboard prediction. *arXiv preprint arXiv:1811.03604*, 2018.
- Hsu, T.-M. H., Qi, H., and Brown, M. Measuring the effects of non-identical data distribution for federated visual classification. *arXiv preprint arXiv:1909.06335*, 2019.

- Iandola, F. N., Moskewicz, M. W., Ashraf, K., and Keutzer, K. Firecaffe: near-linear acceleration of deep neural network training on compute clusters. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pp. 2592–2600, 2016.
- Johnson, R. and Zhang, T. Accelerating stochastic gradient descent using predictive variance reduction. In *Advances in neural information processing systems*, pp. 315–323, 2013.
- Kairouz, P., McMahan, H. B., Avent, B., Bellet, A., Bennis, M., Bhagoji, A. N., Bonawitz, K., Charles, Z., Cormode, G., Cummings, R., et al. Advances and open problems in federated learning. *arXiv preprint arXiv:1912.04977*, 2019.
- Karimi, H., Nutini, J., and Schmidt, M. Linear convergence of gradient and proximal-gradient methods under the polyak-Łojasiewicz condition. In *Joint European Conference on Machine Learning and Knowledge Discovery in Databases*, pp. 795–811. Springer, 2016.
- Karimireddy, S. P., Rebjock, Q., Stich, S. U., and Jaggi, M. Error feedback fixes SignSGD and other gradient compression schemes. *arXiv preprint arXiv:1901.09847*, 2019.
- Khaled, A., Mishchenko, K., and Richtárik, P. Tighter theory for local SGD on identical and heterogeneous data. In *Proceedings of AISTATS*, 2020.
- Kifer, D., Smith, A., and Thakurta, A. Private convex empirical risk minimization and high-dimensional regression. In *Conference on Learning Theory*, pp. 25–1, 2012.
- Konečný, J., McMahan, H. B., Ramage, D., and Richtárik, P. Federated optimization: Distributed machine learning for on-device intelligence. *arXiv preprint arXiv:1610.02527*, 2016a.
- Konečný, J., McMahan, H. B., Yu, F. X., Richtárik, P., Suresh, A. T., and Bacon, D. Federated learning: Strategies for improving communication efficiency. *arXiv preprint arXiv:1610.05492*, 2016b.
- Kulunchakov, A. and Mairal, J. Estimate sequences for stochastic composite optimization: Variance reduction, acceleration, and robustness to noise. *arXiv preprint arXiv:1901.08788*, 2019.
- Lee, J. D., Lin, Q., Ma, T., and Yang, T. Distributed stochastic variance reduced gradient methods and a lower bound for communication complexity. *arXiv preprint arXiv:1507.07595*, 2015.
- Lei, L. and Jordan, M. Less than a single pass: Stochastically controlled stochastic gradient. In *AISTATS*, pp. 148–156, 2017.
- Li, T., Sahu, A. K., Sanjabi, M., Zaheer, M., Talwalkar, A., and Smith, V. On the convergence of federated optimization in heterogeneous networks. *arXiv preprint arXiv:1812.06127*, 2018.
- Li, T., Sanjabi, M., and Smith, V. Fair resource allocation in federated learning. *arXiv preprint arXiv:1905.10497*, 2019a.
- Li, T., Sahu, A. K., Zaheer, M., Sanjabi, M., Talwalkar, A., and Smith, V. Feddane: A federated newton-type method. *arXiv preprint arXiv:2001.01920*, 2020.
- Li, X., Huang, K., Yang, W., Wang, S., and Zhang, Z. On the convergence of FedAvg on non-iid data. *arXiv preprint arXiv:1907.02189*, 2019b.
- Liang, X., Shen, S., Liu, J., Pan, Z., Chen, E., and Cheng, Y. Variance reduced local sgd with lower communication complexity. *arXiv preprint arXiv:1912.12844*, 2019.
- McMahan, B., Moore, E., Ramage, D., Hampson, S., and y Arcas, B. A. Communication-efficient learning of deep networks from decentralized data. In *Proceedings of AISTATS*, pp. 1273–1282, 2017.
- Mishchenko, K., Gorbunov, E., Takáč, M., and Richtárik, P. Distributed learning with compressed gradient differences. *arXiv preprint arXiv:1901.09269*, 2019.
- Mohri, M., Sivek, G., and Suresh, A. T. Agnostic federated learning. *arXiv preprint arXiv:1902.00146*, 2019.
- Nedich, A., Olshevsky, A., and Shi, W. A geometrically convergent method for distributed optimization over time-varying graphs. In *2016 IEEE 55th Conference on Decision and Control (CDC)*, pp. 1023–1029. IEEE, 2016.
- Nesterov, Y. *Lectures on convex optimization*, volume 137. Springer, 2018.
- Nguyen, L. M., Liu, J., Scheinberg, K., and Takáč, M. Sarah: A novel method for machine learning problems using stochastic recursive gradient. In *Proceedings of the 34th International Conference on Machine Learning-Volume 70*, pp. 2613–2621. JMLR. org, 2017.
- Nguyen, L. M., Scheinberg, K., and Takáč, M. Inexact SARAH algorithm for stochastic optimization. *arXiv preprint arXiv:1811.10105*, 2018.
- Patel, K. K. and Dieuleveut, A. Communication trade-offs for synchronized distributed SGD with large step size. *arXiv preprint arXiv:1904.11325*, 2019.
- Ramaswamy, S., Mathews, R., Rao, K., and Beaufays, F. Federated learning for emoji prediction in a mobile keyboard. *arXiv preprint arXiv:1906.04329*, 2019.

- Reddi, S. J., Hefny, A., Sra, S., Poczos, B., and Smola, A. Stochastic variance reduction for nonconvex optimization. In *International conference on machine learning*, pp. 314–323, 2016a.
- Reddi, S. J., Konečný, J., Richtárik, P., Póczos, B., and Smola, A. Aide: Fast and communication efficient distributed optimization. *arXiv preprint arXiv:1608.06879*, 2016b.
- Reddi, S. J., Sra, S., Póczos, B., and Smola, A. Fast incremental method for smooth nonconvex optimization. In *2016 IEEE 55th Conference on Decision and Control (CDC)*, pp. 1971–1977. IEEE, 2016c.
- Safran, I. and Shamir, O. How good is sgd with random shuffling? *arXiv preprint arXiv:1908.00045*, 2019.
- Samarakoon, S., Bennis, M., Saad, W., and Debbah, M. Federated learning for ultra-reliable low-latency v2v communications. In *2018 IEEE Global Communications Conference (GLOBECOM)*, pp. 1–7. IEEE, 2018.
- Schmidt, M., Le Roux, N., and Bach, F. Minimizing finite sums with the stochastic average gradient. *Mathematical Programming*, 162(1-2):83–112, 2017.
- Shamir, O., Srebro, N., and Zhang, T. Communication-efficient distributed optimization using an approximate newton-type method. In *International conference on machine learning*, pp. 1000–1008, 2014.
- Shi, W., Ling, Q., Wu, G., and Yin, W. EXTRA: An exact first-order algorithm for decentralized consensus optimization. *SIAM Journal on Optimization*, 25(2):944–966, 2015.
- Smith, V., Forte, S., Ma, C., Takáč, M., Jordan, M., and Jaggi, M. CoCoA: A general framework for communication-efficient distributed optimization. *Journal of Machine Learning Research*, 18(230):1–49, 2018.
- Stich, S. U. Local SGD converges fast and communicates little. *arXiv preprint arXiv:1805.09767*, 2018.
- Stich, S. U. Unified optimal analysis of the (stochastic) gradient method. *arXiv preprint arXiv:1907.04232*, 2019.
- Stich, S. U. and Karimireddy, S. P. The error-feedback framework: Better rates for SGD with delayed gradients and compressed communication. *arXiv preprint arXiv:1909.05350*, 2019.
- Stich, S. U., Cordonnier, J.-B., and Jaggi, M. Sparsified SGD with memory. In *Advances in Neural Information Processing Systems*, pp. 4447–4458, 2018.
- Suresh, A. T., Yu, F. X., Kumar, S., and McMahan, H. B. Distributed mean estimation with limited communication. In *Proceedings of the 34th International Conference on Machine Learning-Volume 70*, pp. 3329–3337. JMLR.org, 2017.
- Tran-Dinh, Q., Pham, N. H., Phan, D. T., and Nguyen, L. M. Hybrid stochastic gradient descent algorithms for stochastic nonconvex optimization. *arXiv preprint arXiv:1905.05920*, 2019.
- Vaswani, S., Bach, F., and Schmidt, M. Fast and faster convergence of sgd for over-parameterized models and an accelerated perceptron. In *The 22nd International Conference on Artificial Intelligence and Statistics*, pp. 1195–1204, 2019.
- Wang, S., Tuor, T., Salonidis, T., Leung, K. K., Makaya, C., He, T., and Chan, K. Adaptive federated learning in resource constrained edge computing systems. *IEEE Journal on Selected Areas in Communications*, 37(6):1205–1221, 2019.
- Woodworth, B. E., Wang, J., Smith, A., McMahan, H. B., and Srebro, N. Graph oracle models, lower bounds, and gaps for parallel stochastic optimization. In *Advances in neural information processing systems*, pp. 8496–8506, 2018.
- Yang, T., Andrew, G., Eichner, H., Sun, H., Li, W., Kong, N., Ramage, D., and Beaufays, F. Applied federated learning: Improving google keyboard query suggestions. *arXiv preprint arXiv:1812.02903*, 2018.
- Yu, H., Yang, S., and Zhu, S. Parallel restarted SGD with faster convergence and less communication: Demystifying why model averaging works for deep learning. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 33, pp. 5693–5700, 2019.
- Yuan, X.-T. and Li, P. On convergence of distributed approximate newton methods: Globalization, sharper bounds and beyond. *arXiv preprint arXiv:1908.02246*, 2019.
- Zhang, L., Mahdavi, M., and Jin, R. Linear convergence with condition number independent access of full gradients. In *Advances in Neural Information Processing Systems*, pp. 980–988, 2013a.
- Zhang, S., Choromanska, A. E., and LeCun, Y. Deep learning with elastic averaging sgd. In *Advances in neural information processing systems*, pp. 685–693, 2015.
- Zhang, Y., Duchi, J. C., and Wainwright, M. J. Communication-efficient algorithms for statistical optimization. *The Journal of Machine Learning Research*, 14(1):3321–3363, 2013b.

Zhao, Y., Li, M., Lai, L., Suda, N., Civin, D., and Chandra, V. Federated learning with non-iid data. *arXiv preprint arXiv:1806.00582*, 2018.

Zinkevich, M., Weimer, M., Li, L., and Smola, A. J. Parallelized stochastic gradient descent. In *Advances in neural information processing systems*, pp. 2595–2603, 2010.

# Appendix

## A. Related work and significance

**Federated learning.** As stated earlier, federated learning involves learning a centralized model from distributed client data. This centralized model benefits from all client data and can often result in a beneficial performance e.g. in including next word prediction (Hard et al., 2018; Yang et al., 2018), emoji prediction (Ramaswamy et al., 2019), decoder models (Chen et al., 2019b), vocabulary estimation (Chen et al., 2019a), low latency vehicle-to-vehicle communication (Samarakoon et al., 2018), and predictive models in health (Brisimi et al., 2018). Nevertheless, federated learning raises several types of issues and has been the topic of multiple research efforts studying the issues of generalization and fairness (Mohri et al., 2019; Li et al., 2019a), the design of more efficient communication strategies (Konečný et al., 2016b;a; Suresh et al., 2017; Stich et al., 2018; Karimireddy et al., 2019; Basu et al., 2019), the study of lower bounds (Woodworth et al., 2018), differential privacy guarantees (Agarwal et al., 2018), security (Bonawitz et al., 2017), etc. We refer to Kairouz et al. (2019) for an in-depth survey of this area.

**Convergence of FEDAVG** For identical clients, FEDAVG coincides with parallel SGD analyzed by (Zinkevich et al., 2010) who proved asymptotic convergence. Stich (2018) and, more recently Stich & Karimireddy (2019); Patel & Dieuleveut (2019); Khaled et al. (2020), gave a sharper analysis of the same method, under the name of local SGD, also for identical functions. However, there still remains a gap between their upper bounds and the lower bound of Woodworth et al. (2018). The analysis of FEDAVG for heterogeneous clients is more delicate due to the afore-mentioned client-drift, first empirically observed by Zhao et al. (2018). Several analyses bound this drift by assuming bounded gradients (Wang et al., 2019; Yu et al., 2019), or view it as additional noise (Khaled et al., 2020), or assume that the client optima are  $\epsilon$ -close (Li et al., 2018; Haddadpour & Mahdavi, 2019). In a concurrent work, (Liang et al., 2019) propose to use variance reduction to deal with client heterogeneity but still show rates slower than SGD. We summarize the communication complexities of different methods for heterogeneous clients in Table 2.

**Variance reduction.** The use of *control variates* is a classical technique to reduce variance in Monte Carlo sampling methods (cf. (Glasserman, 2013)). In optimization, they were used for finite-sum minimization by SVRG (Johnson & Zhang, 2013; Zhang et al., 2013a) and then in SAGA (Defazio et al., 2014) to simplify the linearly convergent method SAG (Schmidt et al., 2017). Numerous variations and extensions of the technique are studied in (Hanzely & Richtárik, 2019). Starting from (Reddi et al., 2016a), control variates have also frequently been used to reduce variance in finite-sum non-convex settings (Reddi et al., 2016c; Nguyen et al., 2018; Fang et al., 2018; Tran-Dinh et al., 2019). Further, they are used to obtain linearly converging decentralized algorithms under the guise of ‘gradient-tracking’ in (Shi et al., 2015; Nedich et al., 2016) and for gradient compression as ‘compressed-differences’ in (Mishchenko et al., 2019). Our method can be viewed as seeking to remove the ‘client-variance’ in the gradients across the clients, though there still remains additional stochasticity as in (Kulunchakov & Mairal, 2019), which is important in deep learning (Defazio & Bottou, 2019).

**Distributed optimization.** The problem of client-drift we described is a common phenomenon in distributed optimization. In fact, classic techniques such as ADMM mitigate this drift, though they are not applicable in federated learning. For well structured convex problems, CoCoA uses the dual variable as the control variates, enabling flexible distributed methods (Smith et al., 2018). DANE by (Shamir et al., 2014) obtain a closely related primal only algorithm, which was later accelerated by Reddi et al. (2016b) and recently extended to federated learning (Li et al., 2020). SCAFFOLD can be viewed as an improved version of DANE where a fixed number of (stochastic) gradient steps are used in place of a proximal point update. In a similar spirit, distributed variance reduction techniques have been proposed for the finite-sum case (Lee et al., 2015; Konečný et al., 2016a; Cen et al., 2019). However, these methods are restricted to finite-sums and are not applicable to the stochastic setting studied here.

## B. Technicalities

We examine some additional definitions and introduce some technical lemmas.

### B.1. Additional definitions

We make precise a few definitions and explain some of their implications.

(A3)  $f_i$  is  $\mu$ -convex for  $\mu \geq 0$  and satisfies:

$$\langle \nabla f_i(\mathbf{x}), \mathbf{y} - \mathbf{x} \rangle \leq -\left(f_i(\mathbf{x}) - f_i(\mathbf{y}) + \frac{\mu}{2}\|\mathbf{x} - \mathbf{y}\|^2\right), \text{ for any } i, \mathbf{x}, \mathbf{y}.$$

Here, we allow that  $\mu = 0$  (we refer to this case as the general convex case as opposed to strongly convex). It is also possible to generalize all proofs here to the weaker notion of PL-strong convexity (Karimi et al., 2016).

(A4)  $g_i(\mathbf{x}) := \nabla f_i(\mathbf{x}; \zeta_i)$  is unbiased stochastic gradient of  $f_i$  with **bounded variance**

$$\mathbb{E}_{\zeta_i}[\|g_i(\mathbf{x}) - \nabla f_i(\mathbf{x})\|^2] \leq \sigma^2, \text{ for any } i, \mathbf{x}.$$

Note that (A4) only bounds the variance within the same client, but not the variance across the clients.

(A5)  $\{f_i\}$  are  $\beta$ -smooth and satisfy:

$$\|\nabla f_i(\mathbf{x}) - \nabla f_i(\mathbf{y})\| \leq \beta\|\mathbf{x} - \mathbf{y}\|, \text{ for any } i, \mathbf{x}, \mathbf{y}. \quad (7)$$

The assumption (A5) also implies the following quadratic upper bound on  $f_i$

$$f_i(\mathbf{y}) \leq f_i(\mathbf{x}) + \langle \nabla f_i(\mathbf{x}), \mathbf{y} - \mathbf{x} \rangle + \frac{\beta}{2}\|\mathbf{y} - \mathbf{x}\|^2. \quad (8)$$

If additionally the function  $\{f_i\}$  are convex and  $\mathbf{x}^*$  is an optimum of  $f$ , (A5) implies (via Nesterov (2018), Theorem 2.1.5)

$$\frac{1}{2\beta N} \sum_{i=1}^N \|\nabla f_i(\mathbf{x}) - \nabla f_i(\mathbf{x}^*)\|^2 \leq f(\mathbf{x}) - f^*. \quad (9)$$

Further, if  $f_i$  is twice-differentiable, (A5) implies that  $\|\nabla^2 f_i(\mathbf{x})\| \leq \beta$  for any  $\mathbf{x}$ .

### B.2. Some technical lemmas

Now we cover some technical lemmas which are useful for computations later on. The two lemmas below are useful to unroll recursions and derive convergence rates. The first one is a slightly improved (and simplified) version of (Stich, 2019, Theorem 2). It is straightforward to remove the additional logarithmic terms if we use a varying step-size (Kulunchakov & Mairal, 2019, Lemma 13).

**Lemma 1** (linear convergence rate). *For every non-negative sequence  $\{d_{r-1}\}_{r \geq 1}$  and any parameters  $\mu > 0$ ,  $\eta_{\max} \in (0, 1/\mu]$ ,  $c \geq 0$ ,  $R \geq \frac{1}{2\eta_{\max}\mu}$ , there exists a constant step-size  $\eta \leq \eta_{\max}$  and weights  $w_r := (1 - \mu\eta)^{1-r}$  such that for  $W_R := \sum_{r=1}^{R+1} w_r$ ,*

$$\Psi_R := \frac{1}{W_R} \sum_{r=1}^{R+1} \left( \frac{w_r}{\eta} (1 - \mu\eta) d_{r-1} - \frac{w_r}{\eta} d_r + c\eta w_r \right) = \tilde{\mathcal{O}} \left( \mu d_0 \exp(-\mu\eta_{\max} R) + \frac{c}{\mu R} \right).$$

*Proof.* By substituting the value of  $w_r$ , we observe that we end up with a telescoping sum and estimate

$$\Psi_R = \frac{1}{\eta W_R} \sum_{r=1}^{R+1} (w_{r-1} d_{r-1} - w_r d_r) + \frac{c\eta}{W_R} \sum_{r=1}^{R+1} w_r \leq \frac{d_0}{\eta W_R} + c\eta.$$

When  $R \geq \frac{1}{2\mu\eta}$ ,  $(1 - \mu\eta)^R \leq \exp(-\mu\eta R) \leq \frac{2}{3}$ . For such an  $R$ , we can lower bound  $\eta W_R$  using

$$\eta W_R = \eta(1 - \mu\eta)^{-R} \sum_{r=0}^R (1 - \mu\eta)^r = \eta(1 - \mu\eta)^{-R} \frac{1 - (1 - \mu\eta)^R}{\mu\eta} \geq (1 - \mu\eta)^{-R} \frac{1}{3\mu}.$$

This proves that for all  $R \geq \frac{1}{2\mu\eta}$ ,

$$\Psi_R \leq 3\mu d_0 (1 - \mu\eta)^R + c\eta \leq 3\mu d_0 \exp(-\mu\eta R) + c\eta.$$

The lemma now follows by carefully tuning  $\eta$ . Consider the following two cases depending on the magnitude of  $R$  and  $\eta_{\max}$ :

- Suppose  $\frac{1}{2\mu R} \leq \eta_{\max} \leq \frac{\log(\max(1, \mu^2 R d_0 / c))}{\mu R}$ . Then we can choose  $\eta = \eta_{\max}$ ,

$$\Psi_R \leq 3\mu d_0 \exp[-\mu\eta_{\max} R] + c\eta_{\max} \leq 3\mu d_0 \exp[-\mu\eta_{\max} R] + \tilde{\mathcal{O}}\left(\frac{c}{\mu R}\right).$$

- Instead if  $\eta_{\max} > \frac{\log(\max(1, \mu^2 R d_0 / c))}{\mu R}$ , we pick  $\eta = \frac{\log(\max(1, \mu^2 R d_0 / c))}{\mu R}$  to claim that

$$\Psi_R \leq 3\mu d_0 \exp[-\log(\max(1, \mu^2 R d_0 / c))] + \tilde{\mathcal{O}}\left(\frac{c}{\mu R}\right) \leq \tilde{\mathcal{O}}\left(\frac{c}{\mu R}\right).$$

□

The next lemma is an extension of (Stich & Karimireddy, 2019, Lemma 13), (Kulunchakov & Mairal, 2019, Lemma 13) and is useful to derive convergence rates for general convex functions ( $\mu = 0$ ) and non-convex functions.

**Lemma 2** (sub-linear convergence rate). *For every non-negative sequence  $\{d_{r-1}\}_{r \geq 1}$  and any parameters  $\eta_{\max} \geq 0$ ,  $c \geq 0$ ,  $R \geq 0$ , there exists a constant step-size  $\eta \leq \eta_{\max}$  and weights  $w_r = 1$  such that,*

$$\Psi_R := \frac{1}{R+1} \sum_{r=1}^{R+1} \left( \frac{d_{r-1}}{\eta} - \frac{d_r}{\eta} + c_1\eta + c_2\eta^2 \right) \leq \frac{d_0}{\eta_{\max}(R+1)} + \frac{2\sqrt{c_1 d_0}}{\sqrt{R+1}} + 2 \left( \frac{d_0}{R+1} \right)^{\frac{2}{3}} c_2^{\frac{1}{3}}.$$

*Proof.* Unrolling the sum, we can simplify

$$\Psi_R \leq \frac{d_0}{\eta(R+1)} + c_1\eta + c_2\eta^2.$$

Similar to the strongly convex case (Lemma 1), we distinguish the following cases:

- When  $R+1 \leq \frac{d_0}{c_1\eta_{\max}^2}$ , and  $R+1 \leq \frac{d_0}{c_2\eta_{\max}^3}$  we pick  $\eta = \eta_{\max}$  to claim

$$\Psi_R \leq \frac{d_0}{\eta_{\max}(R+1)} + c_1\eta_{\max} + c_2\eta_{\max}^2 \leq \frac{d_0}{\eta_{\max}(R+1)} + \frac{\sqrt{c_1 d_0}}{\sqrt{R+1}} + \left( \frac{d_0}{R+1} \right)^{\frac{2}{3}} c_2^{\frac{1}{3}}.$$

- In the other case, we have  $\eta_{\max}^2 \geq \frac{d_0}{c_1(R+1)}$  or  $\eta_{\max}^3 \geq \frac{d_0}{c_2(R+1)}$ . We choose  $\eta = \min\left\{\sqrt{\frac{d_0}{c_1(R+1)}}, \sqrt[3]{\frac{d_0}{c_2(R+1)}}\right\}$  to prove

$$\Psi_R \leq \frac{d_0}{\eta(R+1)} + c\eta = \frac{2\sqrt{c_1 d_0}}{\sqrt{R+1}} + 2\sqrt[3]{\frac{d_0^2 c_2}{(R+1)^2}}.$$

□

Next, we state a relaxed triangle inequality true for the squared  $\ell_2$  norm.

**Lemma 3** (relaxed triangle inequality). *Let  $\{\mathbf{v}_1, \dots, \mathbf{v}_\tau\}$  be  $\tau$  vectors in  $\mathbb{R}^d$ . Then the following are true:*

1.  $\|\mathbf{v}_i + \mathbf{v}_j\|^2 \leq (1+a)\|\mathbf{v}_i\|^2 + (1+\frac{1}{a})\|\mathbf{v}_j\|^2$  for any  $a > 0$ , and
2.  $\|\sum_{i=1}^\tau \mathbf{v}_i\|^2 \leq \tau \sum_{i=1}^\tau \|\mathbf{v}_i\|^2$ .

*Proof.* The proof of the first statement for any  $a > 0$  follows from the identity:

$$\|\mathbf{v}_i + \mathbf{v}_j\|^2 = (1+a)\|\mathbf{v}_i\|^2 + (1+\frac{1}{a})\|\mathbf{v}_j\|^2 - \|\sqrt{a}\mathbf{v}_i + \frac{1}{\sqrt{a}}\mathbf{v}_j\|^2.$$

For the second inequality, we use the convexity of  $\mathbf{x} \rightarrow \|\mathbf{x}\|^2$  and Jensen's inequality

$$\left\| \frac{1}{\tau} \sum_{i=1}^\tau \mathbf{v}_i \right\|^2 \leq \frac{1}{\tau} \sum_{i=1}^\tau \|\mathbf{v}_i\|^2.$$

□

Next we state an elementary lemma about expectations of norms of random vectors.

**Lemma 4** (separating mean and variance). *Let  $\{\Xi_1, \dots, \Xi_\tau\}$  be  $\tau$  random variables in  $\mathbb{R}^d$  which are not necessarily independent. First suppose that their mean is  $\mathbb{E}[\Xi_i] = \xi_i$  and variance is bounded as  $\mathbb{E}[\|\Xi_i - \xi_i\|^2] \leq \sigma^2$ . Then, the following holds*

$$\mathbb{E}[\|\sum_{i=1}^\tau \Xi_i\|^2] \leq \|\sum_{i=1}^\tau \xi_i\|^2 + \tau^2 \sigma^2.$$

Now instead suppose that their conditional mean is  $\mathbb{E}[\Xi_i | \Xi_{i-1}, \dots, \Xi_1] = \xi_i$  i.e. the variables  $\{\Xi_i - \xi_i\}$  form a martingale difference sequence, and the variance is bounded by  $\mathbb{E}[\|\Xi_i - \xi_i\|^2] \leq \sigma^2$  as before. Then we can show the tighter bound

$$\mathbb{E}[\|\sum_{i=1}^\tau \Xi_i\|^2] \leq 2\|\sum_{i=1}^\tau \xi_i\|^2 + 2\tau\sigma^2.$$

*Proof.* For any random variable  $X$ ,  $\mathbb{E}[X^2] = (\mathbb{E}[X - \mathbb{E}[X]])^2 + (\mathbb{E}[X])^2$  implying

$$\mathbb{E}[\|\sum_{i=1}^\tau \Xi_i\|^2] = \|\sum_{i=1}^\tau \xi_i\|^2 + \mathbb{E}[\|\sum_{i=1}^\tau \Xi_i - \xi_i\|^2].$$

Expanding the above expression using relaxed triangle inequality (Lemma 3) proves the first claim:

$$\mathbb{E}[\|\sum_{i=1}^\tau \Xi_i - \xi_i\|^2] \leq \tau \sum_{i=1}^\tau \mathbb{E}[\|\Xi_i - \xi_i\|^2] \leq \tau^2 \sigma^2.$$

For the second statement,  $\xi_i$  is not deterministic and depends on  $\Xi_{i-1}, \dots, \Xi_1$ . Hence we have to resort to the cruder relaxed triangle inequality to claim

$$\mathbb{E}[\|\sum_{i=1}^\tau \Xi_i\|^2] \leq 2\|\sum_{i=1}^\tau \xi_i\|^2 + 2\mathbb{E}[\|\sum_{i=1}^\tau \Xi_i - \xi_i\|^2]$$

and then use the tighter expansion of the second term:

$$\mathbb{E}[\|\sum_{i=1}^\tau \Xi_i - \xi_i\|^2] = \sum_{i,j} \mathbb{E}[(\Xi_i - \xi_i)^\top (\Xi_j - \xi_j)] = \sum_i \mathbb{E}[\|\Xi_i - \xi_i\|^2] \leq \tau\sigma^2.$$

The cross terms in the above expression have zero mean since  $\{\Xi_i - \xi_i\}$  form a martingale difference sequence. □

## C. Properties of convex functions

We now study two lemmas which hold for any smooth and strongly-convex functions. The first is a generalization of the standard strong convexity inequality (A3), but can handle gradients computed at slightly perturbed points.

**Lemma 5** (perturbed strong convexity). *The following holds for any  $\beta$ -smooth and  $\mu$ -strongly convex function  $h$ , and any  $\mathbf{x}, \mathbf{y}, \mathbf{z}$  in the domain of  $h$ :*

$$\langle \nabla h(\mathbf{x}), \mathbf{z} - \mathbf{y} \rangle \geq h(\mathbf{z}) - h(\mathbf{y}) + \frac{\mu}{4} \|\mathbf{y} - \mathbf{z}\|^2 - \beta \|\mathbf{z} - \mathbf{x}\|^2.$$

*Proof.* Given any  $\mathbf{x}$ ,  $\mathbf{y}$ , and  $\mathbf{z}$ , we get the following two inequalities using smoothness and strong convexity of  $h$ :

$$\begin{aligned} \langle \nabla h(\mathbf{x}), \mathbf{z} - \mathbf{x} \rangle &\geq h(\mathbf{z}) - h(\mathbf{x}) - \frac{\beta}{2} \|\mathbf{z} - \mathbf{x}\|^2 \\ \langle \nabla h(\mathbf{x}), \mathbf{x} - \mathbf{y} \rangle &\geq h(\mathbf{x}) - h(\mathbf{y}) + \frac{\mu}{2} \|\mathbf{y} - \mathbf{x}\|^2. \end{aligned}$$

Further, applying the relaxed triangle inequality gives

$$\frac{\mu}{2} \|\mathbf{y} - \mathbf{x}\|^2 \geq \frac{\mu}{4} \|\mathbf{y} - \mathbf{z}\|^2 - \frac{\mu}{2} \|\mathbf{x} - \mathbf{z}\|^2.$$

Combining all the inequalities together we have

$$\langle \nabla h(\mathbf{x}), \mathbf{z} - \mathbf{y} \rangle \geq h(\mathbf{z}) - h(\mathbf{y}) + \frac{\mu}{4} \|\mathbf{y} - \mathbf{z}\|^2 - \frac{\beta + \mu}{2} \|\mathbf{z} - \mathbf{x}\|^2.$$

The lemma follows since  $\beta \geq \mu$ . □

Here, we see that a gradient step is a contractive operator.

**Lemma 6** (contractive mapping). *For any  $\beta$ -smooth and  $\mu$ -strongly convex function  $h$ , points  $\mathbf{x}, \mathbf{y}$  in the domain of  $h$ , and step-size  $\eta \leq \frac{1}{\beta}$ , the following is true*

$$\|\mathbf{x} - \eta \nabla h(\mathbf{x}) - \mathbf{y} + \eta \nabla h(\mathbf{y})\|^2 \leq (1 - \mu\eta) \|\mathbf{x} - \mathbf{y}\|^2.$$

*Proof.*

$$\begin{aligned} \|\mathbf{x} - \eta \nabla h(\mathbf{x}) - \mathbf{y} + \eta \nabla h(\mathbf{y})\|^2 &= \|\mathbf{x} - \mathbf{y}\|^2 + \eta^2 \|\nabla h(\mathbf{x}) - \nabla h(\mathbf{y})\|^2 - 2\eta \langle \nabla h(\mathbf{x}) - \nabla h(\mathbf{y}), \mathbf{x} - \mathbf{y} \rangle \\ &\stackrel{(A5)}{\leq} \|\mathbf{x} - \mathbf{y}\|^2 + (\eta^2 \beta - 2\eta) \langle \nabla h(\mathbf{x}) - \nabla h(\mathbf{y}), \mathbf{x} - \mathbf{y} \rangle. \end{aligned}$$

Recall our bound on the step-size  $\eta \leq \frac{1}{\beta}$  which implies that  $(\eta^2 \beta - 2\eta) \leq -\eta$ . Finally, apply the  $\mu$ -strong convexity of  $h$  to get

$$-\eta \langle \nabla h(\mathbf{x}) - \nabla h(\mathbf{y}), \mathbf{x} - \mathbf{y} \rangle \leq -\eta \mu \|\mathbf{x} - \mathbf{y}\|^2.$$

□

## D. Convergence of FEDAVG

**Algorithm 2** FEDAVG: Federated Averaging

---

```

1: server input: initial  $\mathbf{x}$ , and global step-size  $\eta_g$ 
2: client  $i$ 's input: local step-size  $\eta_l$ 
3: for each round  $r = 1, \dots, R$  do
4:   sample clients  $\mathcal{S} \subseteq \{1, \dots, N\}$ 
5:   communicate  $\mathbf{x}$  to all clients  $i \in \mathcal{S}$ 
6:   on client  $i \in \mathcal{S}$  in parallel do
7:     initialize local model  $\mathbf{y}_i \leftarrow \mathbf{x}$ 
8:     for  $k = 1, \dots, K$  do
9:       compute mini-batch gradient  $g_i(\mathbf{y}_i)$ 
10:       $\mathbf{y}_i \leftarrow \mathbf{y}_i - \eta_l g_i(\mathbf{y}_i)$ 
11:    end for
12:    communicate  $\Delta\mathbf{y}_i \leftarrow \mathbf{y}_i - \mathbf{x}$ 
13:  end on client
14:   $\Delta\mathbf{x} \leftarrow \frac{1}{|\mathcal{S}|} \sum_{i \in \mathcal{S}} \Delta\mathbf{y}_i$ 
15:   $\mathbf{x} \leftarrow \mathbf{x} + \eta_g \Delta\mathbf{x}$ 
16: end for

```

---

We outline the FEDAVG method in Algorithm 2. In round  $r$  we sample  $\mathcal{S}^r \subseteq [N]$  clients with  $|\mathcal{S}^r| = S$  and then perform the following updates:

- Starting from the shared global parameters  $\mathbf{y}_{i,0}^r = \mathbf{x}^{r-1}$ , we update the local parameters for  $k \in [K]$

$$\mathbf{y}_{i,k}^r = \mathbf{y}_{i,k-1}^r - \eta_l g_i(\mathbf{y}_{i,k-1}^r). \quad (10)$$

- Compute the new global parameters using only updates from the clients  $i \in \mathcal{S}^r$  and a global step-size  $\eta_g$ :

$$\mathbf{x}^r = \mathbf{x}^{r-1} + \frac{\eta_g}{S} \sum_{i \in \mathcal{S}^r} (\mathbf{y}_{i,K}^r - \mathbf{x}^{r-1}). \quad (11)$$

Finally, for some weights  $\{w_r\}$ , we output

$$\bar{\mathbf{x}}^R = \mathbf{x}^{r-1} \text{ with probability } \frac{w_r}{\sum_\tau w_\tau} \text{ for } r \in \{1, \dots, R+1\}. \quad (12)$$

### D.1. Bounding heterogeneity

Recall our bound on the gradient dissimilarity:

$$\frac{1}{N} \sum_{i=1}^N \|\nabla f_i(\mathbf{x})\|^2 \leq G^2 + B^2 \|\nabla f(\mathbf{x})\|^2. \quad (13)$$

If  $\{f_i\}$  are convex, we can relax the assumption to

$$\frac{1}{N} \sum_{i=1}^N \|\nabla f_i(\mathbf{x})\|^2 \leq G^2 + 2\beta B^2 (f(\mathbf{x}) - f^*). \quad (14)$$

We defined two variants of the bounds on the heterogeneity depending of whether the functions are convex or not. Suppose that the functions  $f$  is indeed convex as in (A3) and  $\beta$ -smooth as in (A5), then it is straightforward to see that (13) implies (14). Thus for convex functions, (A1) is mildly weaker. Suppose that the functions  $\{f_1, \dots, f_N\}$  are convex and  $\beta$ -smooth.

Then (14) is satisfied with  $B^2 = 2$  since

$$\begin{aligned} \frac{1}{N} \sum_{i=1}^N \|\nabla f_i(\mathbf{x})\|^2 &\leq \frac{2}{N} \sum_{i=1}^N \|\nabla f_i(\mathbf{x}^*)\|^2 + \frac{2}{N} \sum_{i=1}^N \|\nabla f_i(\mathbf{x}) - \nabla f_i(\mathbf{x}^*)\|^2 \\ &\stackrel{(9)}{\leq} \underbrace{\frac{2}{N} \sum_{i=1}^N \|\nabla f_i(\mathbf{x}^*)\|^2}_{=: \sigma_f^2} + 4\beta(f(\mathbf{x}) - f^*). \end{aligned}$$

Thus,  $(G, B)$ -BGD (14) is equivalent to the heterogeneity assumption of (Mishchenko et al., 2019) with  $G^2 = \sigma_f^2$ . Instead, if we have the stronger assumption (A1) but the functions are possibly non-convex, then  $G = \epsilon$  corresponds to the **local dissimilarity** defined in (Li et al., 2018). Note that assuming  $G$  is negligible is quite strong and corresponds to the strong-growth condition (Vaswani et al., 2019).

## D.2. Rates of convergence (Theorem I)

We first restate Theorem I with some additional details and then see its proof.

**Theorem V.** Suppose that the functions  $\{f_i\}$  satisfies assumptions A4, A5, and A1. Then, in each of the following cases, there exist weights  $\{w_r\}$  and local step-sizes  $\eta_l$  such that for any  $\eta_g \geq 1$  the output of FEDAVG (12)  $\bar{\mathbf{x}}^R$  satisfies

- **Strongly convex:**  $f_i$  satisfies (A3) for  $\mu > 0$ ,  $\eta_l \leq \frac{1}{8(1+B^2)\beta K \eta_g}$ ,  $R \geq \frac{8(1+B^2)\beta}{\mu}$  then

$$\mathbb{E}[f(\bar{\mathbf{x}}^R)] - f(\mathbf{x}^*) \leq \tilde{\mathcal{O}}\left(\frac{M^2}{\mu RKS} + \frac{\beta G^2}{\mu^2 R^2} + \mu D^2 \exp(-\frac{\mu}{16(1+B^2)\beta} R)\right),$$

- **General convex:**  $f_i$  satisfies (A3) for  $\mu = 0$ ,  $\eta_l \leq \frac{1}{(1+B^2)8\beta K \eta_g}$ ,  $R \geq 1$  then

$$\mathbb{E}[f(\bar{\mathbf{x}}^R)] - f(\mathbf{x}^*) \leq \mathcal{O}\left(\frac{MD}{\sqrt{RKS}} + \frac{D^{4/3}(\beta G^2)^{1/3}}{(R+1)^{2/3}} + \frac{B^2 \beta D^2}{R}\right),$$

- **Non-convex:**  $f_i$  satisfies (A1) and  $\eta_l \leq \frac{1}{(1+B^2)8\beta K \eta_g}$ , then

$$\mathbb{E}[\|\nabla f(\bar{\mathbf{x}}^R)\|^2] \leq \mathcal{O}\left(\frac{\beta M \sqrt{F}}{\sqrt{RKS}} + \frac{F^{2/3}(\beta G^2)^{1/3}}{(R+1)^{2/3}} + \frac{B^2 \beta F}{R}\right),$$

where  $M^2 := \sigma^2(1 + \frac{S}{\eta_g^2}) + K(1 - \frac{S}{N})G^2$ ,  $D := \|\mathbf{x}^0 - \mathbf{x}^*\|$ , and  $F := f(\mathbf{x}^0) - f^*$ .

## D.3. Proof of convergence

We will only prove the rate of convergence for convex functions here. The corresponding rates for non-convex functions are easy to derive following the techniques in the rest of the paper.

**Lemma 7. (one round progress)** Suppose our functions satisfies assumptions (A1) and (A3)–(A5). For any step-size satisfying  $\eta_l \leq \frac{1}{(1+B^2)8\beta K \eta_g}$  and effective step-size  $\tilde{\eta} := K \eta_g \eta_l$ , the updates of FEDAVG satisfy

$$\mathbb{E}\|\mathbf{x}^r - \mathbf{x}^*\|^2 \leq (1 - \frac{\mu \tilde{\eta}}{2}) \mathbb{E}\|\mathbf{x}^{r-1} - \mathbf{x}^*\|^2 + (\frac{1}{KS}) \tilde{\eta}^2 \sigma^2 + (1 - \frac{S}{N}) \frac{4\tilde{\eta}^2}{S} G^2 - \tilde{\eta}(\mathbb{E}[f(\mathbf{x}^{r-1})] - f(\mathbf{x}^*)) + 3\beta \tilde{\eta} \mathcal{E}_r,$$

where  $\mathcal{E}_r$  is the drift caused by the local updates on the clients defined to be

$$\mathcal{E}_r := \frac{1}{KN} \sum_{k=1}^K \sum_{i=1}^N \mathbb{E}_r[\|\mathbf{y}_{i,k}^r - \mathbf{x}^{r-1}\|^2].$$

*Proof.* We start with the observation that the updates (10) and (11) imply that the server update in round  $r$  can be written as below (dropping the superscripts everywhere)

$$\Delta \mathbf{x} = -\frac{\tilde{\eta}}{KS} \sum_{k,i \in \mathcal{S}} g_i(\mathbf{y}_{i,k-1}) \text{ and } \mathbb{E}[\Delta \mathbf{x}] = -\frac{\tilde{\eta}}{KN} \sum_{k,i} \mathbb{E}[\nabla f_i(\mathbf{y}_{i,k-1})].$$

We adopt the convention that summations are always over  $k \in [K]$  or  $i \in [N]$  unless otherwise stated. Expanding using above observing, we proceed as<sup>2</sup>

$$\begin{aligned} \mathbb{E}_{r-1} \|\mathbf{x} + \Delta \mathbf{x} - \mathbf{x}^*\|^2 &= \|\mathbf{x} - \mathbf{x}^*\|^2 - \frac{2\tilde{\eta}}{KN} \sum_{k,i} \langle \nabla f_i(\mathbf{y}_{i,k-1}), \mathbf{x} - \mathbf{x}^* \rangle + \tilde{\eta}^2 \mathbb{E}_{r-1} \left\| \frac{1}{KS} \sum_{k,i \in \mathcal{S}} g_i(\mathbf{y}_{i,k-1}) \right\|^2 \\ &\stackrel{\text{Lem. 4}}{\leq} \underbrace{\|\mathbf{x} - \mathbf{x}^*\|^2 - \frac{2\tilde{\eta}}{KN} \sum_{k,i} \langle \nabla f_i(\mathbf{y}_{i,k-1}), \mathbf{x} - \mathbf{x}^* \rangle}_{\mathcal{A}_1} \\ &\quad + \underbrace{\tilde{\eta}^2 \mathbb{E}_{r-1} \left\| \frac{1}{KS} \sum_{k,i \in \mathcal{S}} \nabla f_i(\mathbf{y}_{i,k-1}) \right\|^2 + \frac{\tilde{\eta}^2 \sigma^2}{KS}}_{\mathcal{A}_2}. \end{aligned}$$

We can directly apply Lemma 5 with  $h = f_i$ ,  $\mathbf{x} = \mathbf{y}_{i,k-1}$ ,  $\mathbf{y} = \mathbf{x}^*$ , and  $\mathbf{z} = \mathbf{x}$  to the first term  $\mathcal{A}_1$

$$\begin{aligned} \mathcal{A}_1 &= \frac{2\tilde{\eta}}{KN} \sum_{k,i} \langle \nabla f_i(\mathbf{y}_{i,k-1}), \mathbf{x}^* - \mathbf{x} \rangle \\ &\leq \frac{2\tilde{\eta}}{KN} \sum_{k,i} \left( f_i(\mathbf{x}^*) - f_i(\mathbf{x}) + \beta \|\mathbf{y}_{i,k-1} - \mathbf{x}\|^2 - \frac{\mu}{4} \|\mathbf{x} - \mathbf{x}^*\|^2 \right) \\ &= -2\tilde{\eta} \left( f(\mathbf{x}) - f(\mathbf{x}^*) + \frac{\mu}{4} \|\mathbf{x} - \mathbf{x}^*\|^2 \right) + 2\beta\tilde{\eta}\mathcal{E}. \end{aligned}$$

For the second term  $\mathcal{A}_2$ , we repeatedly apply the relaxed triangle inequality (Lemma 4)

$$\begin{aligned} \mathcal{A}_2 &= \tilde{\eta}^2 \mathbb{E}_{r-1} \left\| \frac{1}{KS} \sum_{k,i \in \mathcal{S}} \nabla f_i(\mathbf{y}_{i,k-1}) - \nabla f_i(\mathbf{x}) + \nabla f_i(\mathbf{x}) \right\|^2 \\ &\leq 2\tilde{\eta}^2 \mathbb{E}_{r-1} \left\| \frac{1}{KS} \sum_{k,i \in \mathcal{S}} \nabla f_i(\mathbf{y}_{i,k-1}) - \nabla f_i(\mathbf{x}) \right\|^2 + 2\tilde{\eta}^2 \mathbb{E}_{r-1} \left\| \frac{1}{S} \sum_{i \in \mathcal{S}} \nabla f_i(\mathbf{x}) \right\|^2 \\ &\leq \frac{2\tilde{\eta}^2}{KN} \sum_{i,k} \mathbb{E}_{r-1} \|\nabla f_i(\mathbf{y}_{i,k-1}) - \nabla f_i(\mathbf{x})\|^2 + 2\tilde{\eta}^2 \mathbb{E}_{r-1} \left\| \frac{1}{S} \sum_{i \in \mathcal{S}} \nabla f_i(\mathbf{x}) - \nabla f(\mathbf{x}) + \nabla f(\mathbf{x}) \right\|^2 \\ &\leq \frac{2\tilde{\eta}^2 \beta^2}{KN} \sum_{i,k} \mathbb{E}_{r-1} \|\mathbf{y}_{i,k-1} - \mathbf{x}\|^2 + 2\tilde{\eta}^2 \|\nabla f(\mathbf{x})\|^2 + (1 - \frac{S}{N}) 4\tilde{\eta}^2 \frac{1}{SN} \sum_i \|\nabla f_i(\mathbf{x})\|^2 \\ &\leq 2\tilde{\eta}^2 \beta^2 \mathcal{E} + 8\tilde{\eta}^2 \beta (B^2 + 1) (f(\mathbf{x}) - f(\mathbf{x}^*)) + (1 - \frac{S}{N}) \frac{4\tilde{\eta}^2}{S} G^2 \end{aligned}$$

The last step used Assumption  $(G, B)$ -BGD assumption (14) that  $\frac{1}{N} \sum_{i=1}^N \|\nabla f_i(\mathbf{x})\|^2 \leq G^2 + 2\beta B^2 (f(\mathbf{x}) - f^*)$ . The extra  $(1 - \frac{S}{N})$  improvement we get is due to sampling the functions  $\{f_i\}$  without replacement. Plugging back the bounds on  $\mathcal{A}_1$  and  $\mathcal{A}_2$ ,

$$\begin{aligned} \mathbb{E}_{r-1} \|\mathbf{x} + \Delta \mathbf{x} - \mathbf{x}^*\|^2 &\leq (1 - \frac{\mu\tilde{\eta}}{2}) \|\mathbf{x} - \mathbf{x}^*\|^2 - (2\tilde{\eta} - 8\beta\tilde{\eta}^2(B^2 + 1)) (f(\mathbf{x}) - f(\mathbf{x}^*)) \\ &\quad + (1 + \tilde{\eta}\beta) 2\beta\tilde{\eta}\mathcal{E} + \frac{1}{KS} \tilde{\eta}^2 \sigma^2 + (1 - \frac{S}{N}) \frac{4\tilde{\eta}^2}{S} G^2. \end{aligned}$$

The lemma now follows by observing that  $8\beta\tilde{\eta}(B^2 + 1) \leq 1$  and that  $B \geq 0$ .  $\square$

<sup>2</sup>We use the notation  $\mathbb{E}_{r-1}[\cdot]$  to mean conditioned on filtration  $r$  i.e. on all the randomness generated prior to round  $r$ .

**Lemma 8. (bounded drift)** Suppose our functions satisfies assumptions (A1) and (A3)–(A5). Then the updates of FEDAVG for any step-size satisfying  $\eta_l \leq \frac{1}{(1+B^2)8\beta K\eta_g}$  have bounded drift:

$$3\beta\tilde{\eta}\mathcal{E}_r \leq \frac{2\tilde{\eta}}{3}(\mathbb{E}[f(\mathbf{x}^{r-1})]) - f(\mathbf{x}^*) + \frac{\tilde{\eta}^2\sigma^2}{2K\eta_g^2} + 18\beta\tilde{\eta}^3G^2.$$

*Proof.* If  $K = 1$ , the lemma trivially holds since  $\mathbf{y}_{i,0} = \mathbf{x}$  for all  $i \in [N]$  and  $\mathcal{E}_r = 0$ . Assume  $K \geq 2$  here on. Recall that the local update made on client  $i$  is  $\mathbf{y}_{i,k} = \mathbf{y}_{i,k-1} - \eta_l g_i(\mathbf{y}_{i,k-1})$ . Then,

$$\begin{aligned} \mathbb{E}\|\mathbf{y}_{i,k} - \mathbf{x}\|^2 &= \mathbb{E}\|\mathbf{y}_{i,k-1} - \mathbf{x} - \eta_l g_i(\mathbf{y}_{i,k-1})\|^2 \\ &\leq \mathbb{E}\|\mathbf{y}_{i,k-1} - \mathbf{x} - \eta_l \nabla f_i(\mathbf{y}_{i,k-1})\|^2 + \eta_l^2\sigma^2 \\ &\leq (1 + \frac{1}{K-1}) \mathbb{E}\|\mathbf{y}_{i,k-1} - \mathbf{x}\|^2 + K\eta_l^2\|\nabla f_i(\mathbf{y}_{i,k-1})\|^2 + \eta_l^2\sigma^2 \\ &= (1 + \frac{1}{K-1}) \mathbb{E}\|\mathbf{y}_{i,k-1} - \mathbf{x}\|^2 + \frac{\tilde{\eta}^2}{\eta_g K}\|\nabla f_i(\mathbf{y}_{i,k-1})\|^2 + \frac{\tilde{\eta}^2\sigma^2}{K^2\eta_g^2} \\ &\leq (1 + \frac{1}{K-1}) \mathbb{E}\|\mathbf{y}_{i,k-1} - \mathbf{x}\|^2 + \frac{2\tilde{\eta}^2}{\eta_g K}\|\nabla f_i(\mathbf{y}_{i,k-1}) - \nabla f_i(\mathbf{x})\|^2 + \frac{2\tilde{\eta}^2}{\eta_g K}\|\nabla f_i(\mathbf{x})\|^2 + \frac{\tilde{\eta}^2\sigma^2}{K^2\eta_g^2} \\ &\leq (1 + \frac{1}{K-1} + \frac{2\tilde{\eta}^2\beta^2}{\eta_g K}) \mathbb{E}\|\mathbf{y}_{i,k-1} - \mathbf{x}\|^2 + \frac{2\tilde{\eta}^2}{\eta_g K}\|\nabla f_i(\mathbf{x})\|^2 + \frac{\tilde{\eta}^2\sigma^2}{K^2\eta_g^2} \\ &\leq (1 + \frac{2}{(K-1)}) \mathbb{E}\|\mathbf{y}_{i,k-1} - \mathbf{x}\|^2 + \frac{2\tilde{\eta}^2}{\eta_g K}\|\nabla f_i(\mathbf{x})\|^2 + \frac{\tilde{\eta}^2\sigma^2}{K^2\eta_g^2}. \end{aligned}$$

In the above proof we separated the mean and the variance in the first inequality, then used the relaxed triangle inequality with  $a = \frac{1}{K-1}$  in the next inequality. Next equality uses the definition of  $\tilde{\eta}$ , and the rest follow from the Lipschitzness of the gradient. Unrolling the recursion above,

$$\begin{aligned} \mathbb{E}\|\mathbf{y}_{i,k} - \mathbf{x}\|^2 &\leq \sum_{\tau=1}^{k-1} (\frac{2\tilde{\eta}^2}{\eta_g K}\|\nabla f_i(\mathbf{x})\|^2 + \frac{\tilde{\eta}^2\sigma^2}{K^2\eta_g^2})(1 + \frac{2}{(K-1)})^\tau \\ &\leq (\frac{2\tilde{\eta}^2}{\eta_g K}\|\nabla f_i(\mathbf{x})\|^2 + \frac{\tilde{\eta}^2\sigma^2}{K^2\eta_g^2})3K. \end{aligned}$$

Averaging over  $i$  and  $k$ , multiplying by  $3\beta\tilde{\eta}$  and then using Assumption A1,

$$\begin{aligned} 3\beta\tilde{\eta}\mathcal{E}_r &\leq \frac{1}{N} \sum_i 18\beta\tilde{\eta}^3\|\nabla f_i(\mathbf{x})\|^2 + \frac{3\beta\tilde{\eta}^3\sigma^2}{K\eta_g^2} \\ &\leq 18\beta\tilde{\eta}^3G^2 + \frac{3\beta\tilde{\eta}^3\sigma^2}{K\eta_g^2} + 36\beta^2\tilde{\eta}^3B^2(f(\mathbf{x}) - f(\mathbf{x}^*)) \end{aligned}$$

The lemma now follows from our assumption that  $8(B^2 + 1)\beta\tilde{\eta} \leq 1$ .  $\square$

**Proof of Theorems I, V** Adding the statements of Lemmas 7 and 8, we get

$$\begin{aligned} \mathbb{E}\|\mathbf{x} + \Delta\mathbf{x} - \mathbf{x}^*\|^2 &\leq (1 - \frac{\mu\tilde{\eta}}{2}) \mathbb{E}\|\mathbf{x} - \mathbf{x}^*\|^2 + (\frac{1}{KS})\tilde{\eta}^2\sigma^2 + (1 - \frac{S}{N})\frac{4\tilde{\eta}^2}{S}G^2 - \tilde{\eta}(\mathbb{E}[f(\mathbf{x})] - f(\mathbf{x}^*)) \\ &\quad + \frac{2\tilde{\eta}}{3}(\mathbb{E}[f(\mathbf{x})]) - f(\mathbf{x}^*) + \frac{\tilde{\eta}^2\sigma^2}{2K\eta_g^2} + 18\beta\tilde{\eta}^3G^2 \\ &= (1 - \frac{\mu\tilde{\eta}}{2}) \mathbb{E}\|\mathbf{x} - \mathbf{x}^*\|^2 - \frac{\tilde{\eta}}{3}(\mathbb{E}[f(\mathbf{x})] - f(\mathbf{x}^*)) + \tilde{\eta}^2 \left( \frac{\sigma^2}{KS}(1 + \frac{S}{\eta_g^2}) + \frac{4G^2}{S}(1 - \frac{S}{N}) + 18\beta\tilde{\eta}G^2 \right). \end{aligned}$$

Moving the  $(f(\mathbf{x}) - f(\mathbf{x}^*))$  term and dividing throughout by  $\frac{\tilde{\eta}}{3}$ , we get the following bound for any  $\tilde{\eta} \leq \frac{1}{8(1+B^2)\beta}$

$$\mathbb{E}[f(\mathbf{x}^{r-1})] - f(\mathbf{x}^*) \leq \frac{3}{\tilde{\eta}}(1 - \frac{\mu\tilde{\eta}}{2})\|\mathbf{x}^{r-1} - \mathbf{x}^*\|^2 - \frac{3}{\tilde{\eta}}\|\mathbf{x}^r - \mathbf{x}^*\|^2 + 3\tilde{\eta} \left( \frac{\sigma^2}{KS}(1 + \frac{S}{\eta_g^2}) + \frac{4G^2}{S}(1 - \frac{S}{N}) + 18\beta\tilde{\eta}G^2 \right).$$

If  $\mu = 0$  (general convex), we can directly apply Lemma 2. Otherwise, by averaging using weights  $w_r = (1 - \frac{\mu\tilde{\eta}}{2})^{1-r}$  and using the same weights to pick output  $\bar{x}^R$ , we can simplify the above recursive bound (see proof of Lem. 1) to prove that for any  $\tilde{\eta}$  satisfying  $\frac{1}{\mu R} \leq \tilde{\eta} \leq \frac{1}{8(1+B^2)\beta}$

$$\mathbb{E}[f(\bar{x}^R)] - f(\mathbf{x}^*) \leq \underbrace{3\|\mathbf{x}^0 - \mathbf{x}^*\|^2}_{=:d} \mu \exp(-\frac{\tilde{\eta}}{2}\mu R) + \underbrace{\tilde{\eta} \left( \frac{2\sigma^2}{KS} \left(1 + \frac{S}{\eta_g^2}\right) + \frac{8G^2}{S} \left(1 - \frac{S}{N}\right) \right)}_{=:c_1} + \underbrace{\tilde{\eta}^2 (36\beta G^2)}_{=:c_2}$$

Now, the choice of  $\tilde{\eta} = \min\left\{\frac{\log(\max(1, \mu^2 R d / c_1))}{\mu R}, \frac{1}{(1+B^2)8\beta}\right\}$  yields the desired rate. The proof of the non-convex case is very similar and also relies on Lemma 2.

#### D.4. Lower bound for FEDAVG (Theorem II)

We first formalize the class of algorithms we look at before proving our lower bound.

**(A6)** We assume that FEDAVG is run with  $\eta_g = 1$ ,  $K > 1$ , and arbitrary possibly adaptive positive step-sizes  $\{\eta_1, \dots, \eta_R\}$  are used with  $\eta_r \leq \frac{1}{\mu}$  and fixed within a round for all clients. Further, the server update is a convex combination of the client updates with non-adaptive weights.

Note that we only prove the lower bound here for  $\eta_g = 1$ . In fact, by taking  $\eta_g$  infinitely large and scaling  $\eta_l \propto \frac{1}{K\eta_g}$  such that the effective step size  $\tilde{\eta} = \eta_l\eta_g K$  remains constant, FEDAVG reduces to the simple large batch SGD method. Hence, proving a lower bound for arbitrary  $\eta_g$  is not possible, but also is of questionable relevance. Further, note that when  $\sigma^2 = 0$ , the upper bound in Theorem V uses  $\eta_g = 1$  and hence the lower bound serves to show that our analysis is tight.

Below we state a more formal version of Theorem II.

**Theorem VI.** *For any positive constants  $G$ ,  $\mu$ , there exist  $\mu$ -strongly convex functions satisfying A1 for which that the output of FEDAVG satisfying A6 has the error for any  $r \geq 1$ :*

$$f(\mathbf{x}^r) - f(\mathbf{x}^*) \geq \Omega\left(\min(f(\mathbf{x}^0) - f(\mathbf{x}^*), \frac{G^2}{\mu R^2})\right).$$

*Proof.* Consider the following simple one-dimensional functions for any given  $\mu$  and  $G$ :

$$f_1(x) := \mu x^2 + Gx, \text{ and } f_2(x) := -Gx,$$

with  $f(x) = \frac{1}{2}(f_1(x) + f_2(x)) = \frac{\mu}{2}x^2$  and optimum at  $x = 0$ . Clearly  $f$  is  $\mu$ -strongly convex and further  $f_1$  and  $f_2$  satisfy A1 with  $B = 3$ . Note that we chose  $f_2$  to be a linear function (not strongly convex) to simplify computations. The calculations made here can be extended with slightly more work for  $(\tilde{f}_2 = \frac{\mu}{2}x^2 - Gx)$  (e.g. see Theorem 1 of (Safran & Shamir, 2019)).

Let us start FEDAVG from  $x^0 > 0$ . A single local update for  $f_1$  and  $f_2$  in round  $r \geq 1$  is respectively

$$y_1 = y_1 - \eta_r(2\mu x + G) \text{ and } y_2 = y_2 + \eta_r G.$$

Then, straightforward computations show that the update at the end of round  $r$  is of the following form for some averaging weight  $\alpha \in [0, 1]$

$$x^r = x^{r-1}((1 - \alpha)(1 - 2\mu\eta_r)^K + \alpha) + \eta_r G \sum_{\tau=0}^{K-1} (\alpha - (1 - \alpha)(1 - 2\mu\eta_r)^\tau).$$

Since  $\alpha$  was picked obliviously, we can assume that  $\alpha \leq 0.5$ . If indeed  $\alpha > 0.5$ , we can swap the definitions of  $f_1$  and  $f_2$  and the sign of  $x^0$ . With this, we can simplify as

$$\begin{aligned} x^r &\geq x^{r-1} \frac{(1 - 2\mu\eta_r)^K + 1}{2} + \frac{\eta_r G}{2} \sum_{\tau=0}^{K-1} (1 - (1 - 2\mu\eta_r)^\tau) \\ &\geq x^{r-1}(1 - 2\mu\eta_r)^K + \frac{\eta_r G}{2} \sum_{\tau=0}^{K-1} (1 - (1 - 2\mu\eta_r)^\tau). \end{aligned}$$

Observe that in the above expression, the right hand side is increasing with  $\eta_r$ —this represents the effect of the client drift and increases the error as the step-size increases. The left hand side decreases with  $\eta_r$ —this is the usual convergence observed due to taking gradient steps. The rest of the proof is to show that even with a careful balancing of the two terms, the effect of  $G$  cannot be removed. Lemma 9 performs exactly such a computation to prove that for any  $r \geq 1$ ,

$$x^r \geq c \min(x_0, \frac{G}{\mu R}).$$

We finish the proof by noting that  $f(x^r) = \frac{\mu}{2}(x^r)^2$ .  $\square$

**Lemma 9.** Suppose that for all  $r \geq 1$ ,  $\eta_r \leq \frac{1}{\mu}$  and the following is true:

$$x^r \geq x^{r-1}(1 - 2\mu\eta_r)^K + \frac{\eta_r G}{2} \sum_{\tau=0}^{K-1} (1 - (1 - 2\mu\eta_r)^\tau).$$

Then, there exists a constant  $c > 0$  such that for any sequence of step-sizes  $\{\eta^r\}$ :

$$x^r \geq c \min(x_0, \frac{G}{\mu R})$$

*Proof.* Define  $\gamma_r = \mu\eta_r R(K-1)$ . Such a  $\gamma_r$  exists and is positive since  $K \geq 2$ . Then,  $\gamma_r$  satisfies

$$(1 - 2\mu\eta_r)^{\frac{K-1}{2}} = (1 - \frac{2\gamma_r}{R(K-1)})^{\frac{K-1}{2}} \leq \exp(-\gamma_r/R).$$

We then have

$$\begin{aligned} x^r &\geq x^{r-1}(1 - 2\mu\eta_r)^K + \frac{\eta_r G}{2} \sum_{\tau=0}^{K-1} (1 - (1 - 2\mu\eta_r)^\tau) \\ &\geq x^{r-1}(1 - 2\mu\eta_r)^K + \frac{\eta_r G}{2} \sum_{\tau=(K-1)/2}^{K-1} (1 - (1 - 2\mu\eta_r)^\tau) \\ &\geq x^{r-1}(1 - 2\mu\eta_r)^K + \frac{\gamma_r G}{4\mu} (1 - \exp(-\gamma_r/R)). \end{aligned}$$

The second inequality follows because  $\eta_r \leq \frac{1}{\mu}$  implies that  $(1 - (1 - 2\mu\eta_r)^\tau)$  is always positive. If  $\gamma_r \geq R/8$ , then we have a constant  $c_1 \in (0, 1/32)$  which satisfies

$$x^r \geq \frac{c_1 G}{\mu}. \tag{15}$$

On the other hand, if  $\gamma_r < R/8$ , we have a tighter inequality

$$(1 - 2\mu\eta_r)^{\frac{K-1}{2}} = (1 - \frac{2\gamma_r}{R(K-1)})^{\frac{K-1}{2}} \leq 1 - \frac{\gamma_r}{R},$$

implying that

$$\begin{aligned} x^r &\geq x^{r-1} \left(1 - \frac{2\gamma_r}{R(K-1)}\right)^K + \frac{\gamma_r^2 G}{4R\mu} \\ &\geq x^{r-1} \left(1 - \frac{4\gamma_r}{R}\right) + \frac{\gamma_r^2 G}{4\mu R}. \end{aligned} \tag{16}$$

The last step used Bernoulli's inequality and the fact that  $K-1 \leq K/2$  for  $K \geq 2$ . Observe that in the above expression, the right hand side is increasing with  $\gamma_r$ —this represents the effect of the client drift and increases the error as the step-size increases. The left hand side decreases with  $\gamma_r$ —this is the usual convergence observed due to taking gradient steps. The rest of the proof is to show that even with a careful balancing of the two terms, the effect of  $G$  cannot be removed.

Suppose that all rounds after  $r_0 \geq 0$  have a small step-size i.e.  $\gamma_r \leq R/8$  for all  $r > r_0$  and hence satisfies (16). Then we will prove via induction that

$$x^r \geq \min(c_r x^{r_0}, \frac{G}{256\mu R}), \text{ for constants } c_r := (1 - \frac{1}{2R})^{r-r_0}. \quad (17)$$

For  $r = r_0$ , (17) is trivially satisfied. Now for  $r > r_0$ ,

$$\begin{aligned} x^r &\geq x^{r-1}(1 - \frac{4\gamma_r}{R}) + \frac{\gamma_r^2 G}{4\mu R} \\ &\geq \min\left(x^{r-1}(1 - \frac{1}{2R}), \frac{G}{256\mu R}\right) \\ &= \min\left(c_r x^{r_0}, \frac{G}{256\mu R}\right). \end{aligned}$$

The first step is because of (16) and the last step uses the induction hypothesis. The second step considers two cases for  $\gamma_r$ : either  $\gamma_r \leq \frac{1}{8}$  and  $(1 - \frac{1}{2R}) \geq (1 - \frac{1}{2R})$ , or  $\gamma_r^2 \geq \frac{1}{64}$ . Finally note that  $c^r \geq \frac{1}{2}$  using Bernoulli's inequality. We have hence proved

$$x^R \geq \min\left(\frac{1}{2}x^{r_0}, \frac{G}{256\mu R}\right)$$

Now suppose  $\gamma_{r_0} > R/8$ . Then (15) implies that  $x^R \geq \frac{cG}{\mu R}$  for some constant  $c > 0$ . If instead no such  $r_0 \geq 1$  exists, then we can set  $r_0 = 0$ . Now finally observe that the previous proof did not make any assumption on  $R$ , and in fact the inequality stated above holds for all  $r \geq 1$ .  $\square$

## E. Convergence of SCAFFOLD

We first restate the convergence theorem more formally, then prove the result for the convex case, and then for non-convex case. Throughout the proof, we will focus on the harder option II. The proofs for SCAFFOLD with option I are nearly identical and so we skip them.

**Theorem VII.** Suppose that the functions  $\{f_i\}$  satisfies assumptions A4 and A5. Then, in each of the following cases, there exist weights  $\{w_r\}$  and local step-sizes  $\eta_l$  such that for any  $\eta_g \geq 1$  the output (22) of SCAFFOLD satisfies:

- **Strongly convex:**  $f_i$  satisfies (A3) for  $\mu > 0$ ,  $\eta_l \leq \min\left(\frac{1}{81\beta K \eta_g}, \frac{S}{15\mu N K \eta_g}\right)$ ,  $R \geq \max\left(\frac{162\beta}{\mu}, \frac{30N}{S}\right)$  then

$$\mathbb{E}[f(\bar{\mathbf{x}}^R)] - f(\mathbf{x}^*) \leq \tilde{\mathcal{O}}\left(\frac{\sigma^2}{\mu R K S}(1 + \frac{S}{\eta_g^2}) + \frac{N\mu}{S}\tilde{D}^2 \exp\left(-\min\left\{\frac{S}{30N}, \frac{\mu}{162\beta}\right\}R\right)\right).$$

- **General convex:**  $f_i$  satisfies (A3) for  $\mu = 0$ ,  $\eta_l \leq \frac{1}{81\beta K \eta_g}$ ,  $R \geq 1$  then

$$\mathbb{E}[f(\bar{\mathbf{x}}^R)] - f(\mathbf{x}^*) \leq \mathcal{O}\left(\frac{\sigma\tilde{D}}{\sqrt{R K S}}\left(\sqrt{1 + \frac{S}{\eta_g^2}}\right) + \sqrt{\frac{N}{S}}\frac{\beta\tilde{D}^2}{R}\right),$$

- **Non-convex:**  $\eta_l \leq \frac{1}{24K\eta_g\beta}\left(\frac{S}{N}\right)^{\frac{2}{3}}$ , and  $R \geq 1$ , then

$$\mathbb{E}[\|\nabla f(\bar{\mathbf{x}}^R)\|^2] \leq \mathcal{O}\left(\frac{\sigma\sqrt{F}}{\sqrt{R K S}}\left(\sqrt{1 + \frac{N}{\eta_g^2}}\right) + \frac{\beta F}{R}\left(\frac{N}{S}\right)^{\frac{2}{3}}\right).$$

Here  $\tilde{D}^2 := (\|\mathbf{x}^0 - \mathbf{x}^*\|^2 + \frac{1}{2N\beta^2}\sum_{i=1}^N\|\mathbf{c}_i^0 - \nabla f_i(\mathbf{x}^*)\|^2)$  and  $F := (f(\mathbf{x}_0) - f(\mathbf{x}^*))$ .

**Remark 10.** Note that the  $\tilde{D}^2$  defined above involves an additional term  $\frac{1}{2N\beta^2}\sum_{i=1}^N\|\mathbf{c}_i^0 - \nabla f_i(\mathbf{x}^*)\|^2$ . This is standard in variance reduction methods (Johnson & Zhang, 2013; Defazio et al., 2014; Hanzely & Richtárik, 2019). Theoretically, we will use a warm-start strategy to set  $\mathbf{c}_i^0$  and in the first  $N/S$  rounds, we compute  $\mathbf{c}_i^0 = \mathbf{g}_i(\mathbf{x}^0)$  over a batch size of size  $K$ . Then, using smoothness of  $f_i$ , we can bound this additional term as

$$\frac{1}{2N\beta^2}\sum_{i=1}^N\|\mathbf{c}_i^0 - \nabla f_i(\mathbf{x}^*)\|^2 \leq \frac{1}{\beta}(f(\mathbf{x}^0) - f^*) + \frac{\sigma^2}{K\beta^2} \leq D^2 + \frac{\sigma^2}{K\beta^2}.$$

Thus, the asymptotic rates of SCAFFOLD for general convex functions only incurs an additive term of the order of  $O(\sqrt{\frac{N}{S}}\frac{1}{R})$ . For strongly convex functions, we only see the affects in the logarithmic terms.

**Remark 11.** When  $\sigma = 0$  i.e. when clients compute full gradients, the communication complexity of SCAFFOLD is: i) for strongly convex case it is  $\tilde{\mathcal{O}}\left(\frac{N}{S} + \frac{\beta}{\mu}\right)$ , ii) for general convex functions it is  $O\left(\sqrt{\frac{N}{S}}\frac{\beta}{R}\right)$ , <sup>3</sup> and iii) for non-convex functions it is  $O\left(\frac{N^{2/3}}{S}\frac{\beta}{R}\right)$ . In comparison, the follow up work of FedDyn (Acar et al., 2021) proves communication complexity matching ours in the convex and strongly convex settings, but a worse  $O\left(\frac{N}{S}\frac{\beta}{R}\right)$  complexity in the non-convex settings (all when  $\sigma = 0$ ).

We will rewrite SCAFFOLD using notation which is convenient for the proofs:  $\{\mathbf{y}_i\}$  represent the client models,  $\mathbf{x}$  is the aggregate server model, and  $\mathbf{c}_i$  and  $\mathbf{c}$  are the client and server control variates. For an equivalent description which is easier to implement, we refer to Algorithm 1. The server maintains a global control variate  $\mathbf{c}$  as before and each client maintains its own control variate  $\mathbf{c}_i$ . In round  $r$ , a subset of clients  $\mathcal{S}^r$  of size  $S$  are sampled uniformly from  $\{1, \dots, N\}$ . Suppose that every client performs the following updates

- Starting from the shared global parameters  $\mathbf{y}_{i,r}^0 = \mathbf{x}^{r-1}$ , we update the local parameters for  $k \in [K]$

$$\mathbf{y}_{i,k}^r = \mathbf{y}_{i,k-1}^r - \eta_l \mathbf{v}_{i,k}^r, \quad \text{where } \mathbf{v}_{i,k}^r := g_i(\mathbf{y}_{i,k-1}^r) - \mathbf{c}_i^{r-1} + \mathbf{c}^{r-1} \quad (18)$$

<sup>3</sup>A previous version of the paper showed a worse dependence of  $O\left(\frac{N}{S}\frac{\beta}{R}\right)$  due to sub-optimal choice of step-size  $\eta$ .

- Update the control iterates using (option II):

$$\tilde{\mathbf{c}}_i^r = \mathbf{c}^{r-1} - \mathbf{c}_i^{r-1} + \frac{1}{K\eta_l}(\mathbf{x}^{r-1} - \mathbf{x}_{i,K}^r) = \frac{1}{K} \sum_{k=1}^K g_i(\mathbf{y}_{i,k-1}^r). \quad (19)$$

We update the local control variates only for clients  $i \in \mathcal{S}^r$

$$\mathbf{c}_i^r = \begin{cases} \tilde{\mathbf{c}}_i^r & \text{if } i \in \mathcal{S}^r \\ \mathbf{c}_i^{r-1} & \text{otherwise.} \end{cases} \quad (20)$$

- Compute the new global parameters and global control variate using only updates from the clients  $i \in \mathcal{S}^r$ :

$$\mathbf{x}^r = \mathbf{x}^{r-1} + \frac{\eta_g}{S} \sum_{i \in \mathcal{S}^r} (\mathbf{y}_{i,K}^r - \mathbf{x}^{r-1}) \quad \text{and} \quad \mathbf{c}^r = \frac{1}{N} \sum_{i=1}^N \mathbf{c}_i^r = \frac{1}{N} \left( \sum_{i \in \mathcal{S}^r} \mathbf{c}_i^r + \sum_{j \notin \mathcal{S}^r} \mathbf{c}_j^{r-1} \right). \quad (21)$$

Finally, for some weights  $\{w_r\}$ , we output

$$\bar{\mathbf{x}}^R = \mathbf{x}^{r-1} \text{ with probability } \frac{w_r}{\sum_\tau w_\tau} \text{ for } r \in \{1, \dots, R+1\}. \quad (22)$$

Note that the clients are agnostic to the sampling and their updates are identical to when all clients are participating. Also note that the control variate choice (19) corresponds to (option II) of Algorithm 1. Further, the updates of the clients  $i \notin \mathcal{S}^r$  is forgotten and is defined only to make the proofs easier. While actually implementing the method, only clients  $i \in \mathcal{S}^r$  participate and the rest remain inactive (see Algorithm 1).

### E.1. Convergence of SCAFFOLD for convex functions (Theorem III)

We will first bound the variance of SCAFFOLD update in Lemma 12, then see how sampling of clients effects our control variates in Lemma 13, and finally bound the amount of client-drift in Lemma 14. We will then use these three lemmas to prove the progress in a single round in Lemma 15. Combining this progress with Lemmas 1 and 2 gives us the desired rates.

**Additional definitions.** Before proceeding with the proof of our lemmas, we need some additional definitions of the various errors we track. As before, we define the effective step-size to be

$$\tilde{\eta} := K\eta_l\eta_g.$$

We define client-drift to be how much the clients move from their starting point:

$$\mathcal{E}_r := \frac{1}{KN} \sum_{k=1}^K \sum_{i=1}^N \mathbb{E}[\|\mathbf{y}_{i,k}^r - \mathbf{x}^{r-1}\|^2]. \quad (23)$$

Because we are sampling the clients, not all the client control-variates get updated every round. This leads to some ‘lag’ which we call control-lag:

$$\mathcal{C}_r := \frac{1}{N} \sum_{j=1}^N \mathbb{E}[\mathbb{E}[\mathbf{c}_i^r] - \nabla f_i(\mathbf{x}^*)]^2. \quad (24)$$

**Variance of server update.** We study how the variance of the server update can be bounded.

**Lemma 12.** *For updates (18)–(21), we can bound the variance of the server update in any round  $r$  and any  $\tilde{\eta} := \eta_l\eta_g K \geq 0$  as follows*

$$\mathbb{E}[\|\mathbf{x}^r - \mathbf{x}^{r-1}\|^2] \leq 8\beta\tilde{\eta}^2(\mathbb{E}[f(\mathbf{x}^{r-1})] - f(\mathbf{x}^*)) + 8\tilde{\eta}^2\mathcal{C}_{r-1} + 4\tilde{\eta}^2\beta^2\mathcal{E}_r + \frac{12\tilde{\eta}^2\sigma^2}{KS}.$$

*Proof.* The server update in round  $r$  can be written as follows (dropping the superscript  $r$  everywhere)

$$\mathbb{E}\|\Delta\mathbf{x}\|^2 = \mathbb{E}\left\|-\frac{\tilde{\eta}}{KS} \sum_{k,i \in \mathcal{S}} \mathbf{v}_{i,k}\right\|^2 = \mathbb{E}\left\|\frac{\tilde{\eta}}{KS} \sum_{k,i \in \mathcal{S}} (g_i(\mathbf{y}_{i,k-1}) + \mathbf{c} - \mathbf{c}_i)\right\|^2,$$

which can then be expanded as

$$\begin{aligned} \mathbb{E}\|\Delta\mathbf{x}\|^2 &\leq \mathbb{E}\left\|\frac{\tilde{\eta}}{KS} \sum_{k,i \in \mathcal{S}} (g_i(\mathbf{y}_{i,k-1}) + \mathbf{c} - \mathbf{c}_i)\right\|^2 \\ &\leq 4\mathbb{E}\left\|\frac{\tilde{\eta}}{KS} \sum_{k,i \in \mathcal{S}} g_i(\mathbf{y}_{i,k-1}) - \nabla f_i(\mathbf{x})\right\|^2 + 4\tilde{\eta}^2 \mathbb{E}\|\mathbf{c}\|^2 + 4\mathbb{E}\left\|\frac{\tilde{\eta}}{KS} \sum_{k,i \in \mathcal{S}} \nabla f_i(\mathbf{x}^*) - \mathbf{c}_i\right\|^2 \\ &\quad + 4\mathbb{E}\left\|\frac{\tilde{\eta}}{KS} \sum_{k,i \in \mathcal{S}} \nabla f_i(\mathbf{x}) - \nabla f_i(\mathbf{x}^*)\right\|^2 \\ &\stackrel{(9)}{\leq} 4\mathbb{E}\left\|\frac{\tilde{\eta}}{KS} \sum_{k,i \in \mathcal{S}} g_i(\mathbf{y}_{i,k-1}) - \nabla f_i(\mathbf{x})\right\|^2 + 4\tilde{\eta}^2 \mathbb{E}\|\mathbf{c}\|^2 + 4\mathbb{E}\left\|\frac{\tilde{\eta}}{S} \sum_{i \in \mathcal{S}} \nabla f_i(\mathbf{x}^*) - \mathbf{c}_i\right\|^2 \\ &\quad + 8\beta\tilde{\eta}^2(\mathbb{E}[f(\mathbf{x})] - f(\mathbf{x}^*)) \\ &\leq 4\mathbb{E}\left\|\frac{\tilde{\eta}}{KS} \sum_{k,i \in \mathcal{S}} \nabla f_i(\mathbf{y}_{i,k-1}) - \nabla f_i(\mathbf{x})\right\|^2 + 4\tilde{\eta}^2 \mathbb{E}\|\mathbf{c}\|^2 + 4\left\|\frac{\tilde{\eta}}{S} \sum_{i \in \mathcal{S}} \nabla f_i(\mathbf{x}^*) - \mathbb{E}[\mathbf{c}_i]\right\|^2 \\ &\quad + 8\beta\tilde{\eta}^2(\mathbb{E}[f(\mathbf{x})] - f(\mathbf{x}^*)) + \frac{12\tilde{\eta}^2\sigma^2}{KS}. \end{aligned}$$

The inequality before the last used the smoothness of  $\{f_i\}$ . The last inequality which separates the mean and the variance is an application of Lemma 4: the variance of  $(\frac{1}{KS} \sum_{k,i \in \mathcal{S}} g_i(\mathbf{y}_{i,k-1}))$  is bounded by  $\sigma^2/KS$ . Similarly,  $\mathbf{c}_j$  as defined in (19) for any  $j \in [N]$  has variance smaller than  $\sigma^2/K$  and hence the variance of  $(\frac{1}{S} \sum_{i \in \mathcal{S}} \mathbf{c}_i)$  is smaller than  $\sigma^2/KS$ .

Using Lemma 3.2 twice to simplify:

$$\begin{aligned} \mathbb{E}\|\Delta\mathbf{x}\|^2 &\leq \frac{4\tilde{\eta}^2}{KN} \sum_{k,i} \mathbb{E}\|\nabla f_i(\mathbf{y}_{i,k-1}) - \nabla f_i(\mathbf{x})\|^2 + 4\tilde{\eta}^2 \mathbb{E}\|\mathbf{c}\|^2 + \frac{4\tilde{\eta}^2}{N} \sum_i \|\nabla f_i(\mathbf{x}^*) - \mathbb{E}[\mathbf{c}_i]\|^2 \\ &\quad + 8\beta\tilde{\eta}^2(\mathbb{E}[f(\mathbf{x})] - f(\mathbf{x}^*)) + \frac{12\tilde{\eta}^2\sigma^2}{KS} \\ &\leq \underbrace{\frac{4\tilde{\eta}^2}{KN} \sum_{k,i} \mathbb{E}\|\nabla f_i(\mathbf{y}_{i,k-1}) - \nabla f_i(\mathbf{x})\|^2}_{\mathcal{T}_1} + \frac{8\tilde{\eta}^2}{N} \sum_i \|\nabla f_i(\mathbf{x}^*) - \mathbb{E}[\mathbf{c}_i]\|^2 \\ &\quad + 8\beta\tilde{\eta}^2(\mathbb{E}[f(\mathbf{x})] - f(\mathbf{x}^*)) + \frac{12\tilde{\eta}^2\sigma^2}{KS}. \end{aligned}$$

The second step follows because  $\mathbf{c} = \frac{1}{N} \sum_i \mathbf{c}_i$ . Since the gradient of  $f_i$  is  $\beta$ -Lipschitz,  $\mathcal{T}_1 \leq \frac{\beta^2 4\tilde{\eta}^2}{KN} \sum_{k,i} \mathbb{E}\|\mathbf{y}_{i,k-1} - \mathbf{x}\|^2 = 4\tilde{\eta}^2 \beta^2 \mathcal{E}$ . The definition of the error in the control variate  $\mathcal{C}_{r-1} := \frac{1}{N} \sum_{j=1}^N \mathbb{E}\|\mathbb{E}[\mathbf{c}_i] - \nabla f_i(\mathbf{x}^*)\|^2$  completes the proof.  $\square$

**Change in control lag.** We have previously related the variance of the server update to the control lag. We now examine how the control-lag grows each round.

**Lemma 13.** *For updates (18)–(21) with the control update (19) and assumptions A3–A5, the following holds true for any  $\tilde{\eta} := \eta\eta_g K \in [0, 1/\beta]$ :*

$$\mathcal{C}_r \leq (1 - \frac{S}{N})\mathcal{C}_{r-1} + \frac{S}{N}(4\beta(\mathbb{E}[f(\mathbf{x}^{r-1})] - f(\mathbf{x}^*)) + 2\beta^2\mathcal{E}_r).$$

*Proof.* Recall that after round  $r$ , the control update rule (19) implies that  $\mathbf{c}_i^r$  is set as per

$$\mathbf{c}_i^r = \begin{cases} \mathbf{c}_i^{r-1} & \text{if } i \notin \mathcal{S}^r \text{ i.e. with probability } (1 - \frac{S}{N}), \\ \frac{1}{K} \sum_{k=1}^K g_i(\mathbf{y}_{i,k-1}^r) & \text{with probability } \frac{S}{N}. \end{cases}$$

Taking expectations on both sides yields

$$\mathbb{E}[\mathbf{c}_i^r] = (1 - \frac{S}{N}) \mathbb{E}[\mathbf{c}_i^{r-1}] + \frac{S}{KN} \sum_{k=1}^K \mathbb{E}[\nabla f_i(\mathbf{y}_{i,k-1}^r)], \quad \forall i \in [N].$$

Plugging the above expression in the definition of  $\mathcal{C}_r$  we get

$$\begin{aligned} \mathcal{C}_r &= \frac{1}{N} \sum_{i=1}^N \|\mathbb{E}[\mathbf{c}_i^r] - \nabla f_i(\mathbf{x}^*)\|^2 \\ &= \frac{1}{N} \sum_{i=1}^N \|(1 - \frac{S}{N})(\mathbb{E}[\mathbf{c}_i^{r-1}] - \nabla f_i(\mathbf{x}^*)) + \frac{S}{N} (\frac{1}{K} \sum_{k=1}^K \mathbb{E}[\nabla f_i(\mathbf{y}_{i,k-1}^r)] - \nabla f_i(\mathbf{x}^*))\|^2 \\ &\leq (1 - \frac{S}{N}) \mathcal{C}_{r-1} + \frac{S}{N^2 K} \sum_{k=1}^K \mathbb{E}\|\nabla f_i(\mathbf{y}_{i,k-1}^r) - \nabla f_i(\mathbf{x}^*)\|^2. \end{aligned}$$

The final step applied Jensen's inequality twice. We can then further simplify using the relaxed triangle inequality as

$$\begin{aligned} \mathbb{E}_{r-1}[\mathcal{C}_r] &\leq \left(1 - \frac{S}{N}\right) \mathcal{C}_{r-1} + \frac{S}{N^2 K} \sum_{i,k} \mathbb{E}\|\nabla f_i(\mathbf{y}_{i,k-1}^r) - \nabla f_i(\mathbf{x}^*)\|^2 \\ &\leq \left(1 - \frac{S}{N}\right) \mathcal{C}_{r-1} + \frac{2S}{N^2} \sum_i \mathbb{E}\|\nabla f_i(\mathbf{x}^{r-1}) - \nabla f_i(\mathbf{x}^*)\|^2 + \frac{2S}{N^2 K} \sum_{i,k} \mathbb{E}\|\nabla f_i(\mathbf{y}_{i,k-1}^r) - \nabla f_i(\mathbf{x}^{r-1})\|^2 \\ &\stackrel{(7)}{\leq} \left(1 - \frac{S}{N}\right) \mathcal{C}_{r-1} + \frac{2S}{N^2} \sum_i \mathbb{E}\|\nabla f_i(\mathbf{x}^{r-1}) - \nabla f_i(\mathbf{x}^*)\|^2 + \frac{2S}{N^2 K} \beta^2 \sum_{i,k} \mathbb{E}\|\mathbf{y}_{i,k-1}^r - \mathbf{x}^{r-1}\|^2 \\ &\stackrel{(9)}{\leq} \left(1 - \frac{S}{N}\right) \mathcal{C}_{r-1} + \frac{S}{N} (4\beta(\mathbb{E}[f(\mathbf{x}^{r-1})] - f(\mathbf{x}^*)) + \beta^2 \mathcal{E}_r). \end{aligned}$$

The last two inequalities follow from smoothness of  $\{f_i\}$  and the definition  $\mathcal{E}_r = \frac{1}{NK} \beta^2 \sum_{i,k} \mathbb{E}\|\mathbf{y}_{i,k-1}^r - \mathbf{x}^{r-1}\|^2$ .  $\square$

**Bounding client-drift.** We will now bound the final source of error which is the client-drift.

**Lemma 14.** Suppose our step-sizes satisfy  $\eta_l \leq \frac{1}{81\beta K \eta_g}$  and  $f_i$  satisfies assumptions A3–A5. Then, for any global  $\eta_g \geq 1$  we can bound the drift as

$$3\beta\tilde{\eta}\mathcal{E}_r \leq \frac{2\tilde{\eta}^2}{3} \mathcal{C}_{r-1} + \frac{\tilde{\eta}}{25\eta_g^2} (\mathbb{E}[f(\mathbf{x}^{r-1})] - f(\mathbf{x}^*)) + \frac{\tilde{\eta}^2}{K\eta_g^2} \sigma^2.$$

*Proof.* First, observe that if  $K = 1$ ,  $\mathcal{E}_r = 0$  since  $\mathbf{y}_{i,0} = \mathbf{x}$  for all  $i \in [N]$  and that  $\mathcal{C}_{r-1}$  and the right hand side are both positive. Thus the lemma is trivially true if  $K = 1$ . For  $K > 1$ , we build a recursive bound of the drift. Starting from the definition of the update (18) and then applying the relaxed triangle inequality, we can expand

$$\begin{aligned} \frac{1}{S} \mathbb{E}_{r-1} \left[ \sum_{i \in \mathcal{S}} \|(\mathbf{y}_i - \eta_l \mathbf{v}_i) - \mathbf{x}\|^2 \right] &= \frac{1}{S} \mathbb{E}_{r-1} \left[ \sum_{i \in \mathcal{S}} \|\mathbf{y}_i - \eta_l g_i(\mathbf{y}_i) + \eta_l \mathbf{c} - \eta_l \mathbf{c}_i - \mathbf{x}\|^2 \right] \\ &\leq \frac{1}{S} \mathbb{E}_{r-1} \left[ \sum_{i \in \mathcal{S}} \|\mathbf{y}_i - \eta_l \nabla f_i(\mathbf{y}_i) + \eta_l \mathbf{c} - \eta_l \mathbf{c}_i - \mathbf{x}\|^2 \right] + \eta_l^2 \sigma^2 \\ &\leq \frac{(1+a)}{S} \mathbb{E}_{r-1} \left[ \underbrace{\sum_{i \in \mathcal{S}} \|\mathbf{y}_i - \eta_l \nabla f_i(\mathbf{y}_i) + \eta_l \nabla f_i(\mathbf{x}) - \mathbf{x}\|^2}_{\mathcal{T}_2} \right] \\ &\quad + (1 + \frac{1}{a}) \eta_l^2 \mathbb{E}_{r-1} \underbrace{\left[ \frac{1}{S} \sum_{i \in \mathcal{S}} \|\mathbf{c} - \mathbf{c}_i + \nabla f_i(\mathbf{x})\|^2 \right]}_{\mathcal{T}_3} + \eta_l^2 \sigma^2. \end{aligned}$$

The final step follows from the relaxed triangle inequality (Lemma 3). Applying the contractive mapping Lemma 6 for  $\eta_l \leq 1/\beta$  shows

$$\mathcal{T}_2 = \frac{1}{S} \sum_{i \in \mathcal{S}} \|\mathbf{y}_i - \eta_l \nabla f_i(\mathbf{y}_i) + \eta_l \nabla f_i(\mathbf{x}) - \mathbf{x}\|^2 \leq \|\mathbf{y}_i - \mathbf{x}\|^2.$$

Once again using our relaxed triangle inequality to expand the other term  $\mathcal{T}_3$ , we get

$$\begin{aligned}
 \mathcal{T}_3 &= \mathbb{E}_{r-1} \left[ \frac{1}{S} \sum_{i \in \mathcal{S}} \|\mathbf{c} - \mathbf{c}_i + \nabla f_i(\mathbf{x})\|^2 \right] \\
 &= \frac{1}{N} \sum_{j=1}^N \|\mathbf{c} - \mathbf{c}_i + \nabla f_i(\mathbf{x})\|^2 \\
 &= \frac{1}{N} \sum_{j=1}^N \|\mathbf{c} - \mathbf{c}_i + \nabla f_i(\mathbf{x}^*) + \nabla f_i(\mathbf{x}) - \nabla f_i(\mathbf{x}^*)\|^2 \\
 &\leq 3\|\mathbf{c}\|^2 + \frac{3}{N} \sum_{j=1}^N \|\mathbf{c}_i - \nabla f_i(\mathbf{x}^*)\|^2 + \frac{3}{N} \sum_{j=1}^N \|\nabla f_i(\mathbf{x}) - \nabla f_i(\mathbf{x}^*)\|^2 \\
 &\leq \frac{6}{N} \sum_{j=1}^N \|\mathbf{c}_i - \nabla f_i(\mathbf{x}^*)\|^2 + \frac{3}{N} \sum_{j=1}^N \|\nabla f_i(\mathbf{x}) - \nabla f_i(\mathbf{x}^*)\|^2 \\
 &\leq \frac{6}{N} \sum_{j=1}^N \|\mathbf{c}_i - \nabla f_i(\mathbf{x}^*)\|^2 + 6\beta(f(\mathbf{x}) - f(\mathbf{x}^*)) .
 \end{aligned}$$

The last step used the smoothness of  $f_i$ . Combining the bounds on  $\mathcal{T}_2$  and  $\mathcal{T}_3$  in the original inequality and using  $a = \frac{1}{K-1}$  gives

$$\begin{aligned}
 \frac{1}{N} \sum_i \mathbb{E} \|\mathbf{y}_{i,k} - \mathbf{x}\|^2 &\leq \frac{(1 + \frac{1}{K-1})}{N} \sum_i \mathbb{E} \|\mathbf{y}_{i,k-1} - \mathbf{x}\|^2 + \eta_l^2 \sigma^2 \\
 &\quad + 6\eta_l^2 K \beta (f(\mathbf{x}) - f(\mathbf{x}^*)) + \frac{6K\eta_l^2}{N} \sum_i \mathbb{E} \|\mathbf{c}_i - \nabla f_i(\mathbf{x}^*)\|^2 .
 \end{aligned}$$

Recall that with the choice of  $\mathbf{c}_i$  in (19), the variance of  $c_i$  is less than  $\frac{\sigma^2}{K}$ . Separating its mean and variance gives

$$\begin{aligned}
 \frac{1}{N} \sum_i \mathbb{E} \|\mathbf{y}_{i,k} - \mathbf{x}\|^2 &\leq \left(1 + \frac{1}{K-1}\right) \frac{1}{N} \sum_i \mathbb{E} \|\mathbf{y}_{i,k-1} - \mathbf{x}\|^2 + 7\eta_l^2 \sigma^2 + \\
 &\quad 6\eta_l^2 K \beta (f(\mathbf{x}) - f(\mathbf{x}^*)) + \frac{6K\eta_l^2}{N} \sum_i \mathbb{E} \|\mathbb{E}[\mathbf{c}_i] - \nabla f_i(\mathbf{x}^*)\|^2 \quad (25)
 \end{aligned}$$

Unrolling the recursion (25), we get the following for any  $k \in \{1, \dots, K\}$

$$\begin{aligned}
 \frac{1}{N} \sum_i \mathbb{E} \|\mathbf{y}_{i,k} - \mathbf{x}\|^2 &\leq (6K\beta\eta_l^2(f(\mathbf{x}) - f(\mathbf{x}^*)) + 6K\eta_l^2\mathcal{C}_{r-1} + 7\beta\eta_l^2\sigma^2) \left( \sum_{\tau=0}^{k-1} (1 + \frac{1}{K-1})^\tau \right) \\
 &\leq (6K\beta\eta_l^2(f(\mathbf{x}) - f(\mathbf{x}^*)) + 6K\eta_l^2\mathcal{C}_{r-1} + 7\beta\eta_l^2\sigma^2)(K-1)((1 + \frac{1}{K-1})^K - 1) \\
 &\leq (6K\beta\eta_l^2(f(\mathbf{x}) - f(\mathbf{x}^*)) + 6K\eta_l^2\mathcal{C}_{r-1} + 7\beta\eta_l^2\sigma^2)3K \\
 &\leq 18K^2\beta\eta_l^2(f(\mathbf{x}) - f(\mathbf{x}^*)) + 18K^2\eta_l^2\mathcal{C}_{r-1} + 21K\beta\eta_l^2\sigma^2 .
 \end{aligned}$$

The inequality  $(K-1)((1 + \frac{1}{K-1})^K - 1) \leq 3K$  can be verified for  $K = 2, 3$  manually. For  $K \geq 4$ ,

$$(K-1)((1 + \frac{1}{K-1})^K - 1) < K(\exp(\frac{K}{K-1}) - 1) \leq K(\exp(\frac{4}{3}) - 1) < 3K .$$

Again averaging over  $k$  and multiplying by  $3\beta$  yields

$$\begin{aligned}
 3\beta\mathcal{E}_r &\leq 54K^2\beta^2\eta_l^2(f(\mathbf{x}) - f(\mathbf{x}^*)) + 54K^2\beta\eta_l^2\mathcal{C}_{r-1} + 63\beta K\eta_l^2\sigma^2 \\
 &= \frac{1}{\eta_g^2} \left( 54\beta^2\tilde{\eta}^2(f(\mathbf{x}) - f(\mathbf{x}^*)) + 54\beta\tilde{\eta}^2\mathcal{C}_{r-1} + 63\beta\tilde{\eta}^2\frac{\sigma^2}{K} \right) \\
 &\leq \frac{1}{\eta_g^2} \left( \frac{1}{25}(f(\mathbf{x}) - f(\mathbf{x}^*)) + \frac{2}{3}\tilde{\eta}\mathcal{C}_{r-1} + \tilde{\eta}\frac{\sigma^2}{K} \right) .
 \end{aligned}$$

The equality follows from the definition  $\tilde{\eta} = K\eta_l\eta_g$ , and the final inequality uses the bound that  $\tilde{\eta} \leq \frac{1}{81\beta}$ .  $\square$

**Progress in one round.** Now that we have a bound on all errors, we can describe our progress.

**Lemma 15.** Suppose assumptions A3–A5 are true. Then the following holds for any step-sizes satisfying  $\eta_g \geq 1$ ,  $\eta_l \leq \min\left(\frac{1}{81\beta K\eta_g}, \frac{S}{15\mu N K\eta_g}\right)$ , and effective step-size  $\tilde{\eta} := K\eta_g\eta_l$

$$\mathbb{E}\left[\|\mathbf{x}^r - \mathbf{x}^*\|^2 + \frac{9N\tilde{\eta}^2}{S}\mathcal{C}_r\right] \leq (1 - \frac{\mu\tilde{\eta}}{2})\left(\mathbb{E}\|\mathbf{x}^{r-1} - \mathbf{x}^*\|^2 + \frac{9N\tilde{\eta}^2}{S}\mathcal{C}_{r-1}\right) - \tilde{\eta}(\mathbb{E}[f(\mathbf{x}^{r-1})] - f(\mathbf{x}^*)) + \frac{12\tilde{\eta}^2}{KS}(1 + \frac{S}{\eta_g^2})\sigma^2.$$

*Proof.* Starting from our server update equation,

$$\Delta\mathbf{x} = -\frac{\tilde{\eta}}{KS} \sum_{k,i \in \mathcal{S}} (g_i(\mathbf{y}_{i,k-1}) + \mathbf{c} - \mathbf{c}_i), \text{ and } \mathbb{E}[\Delta\mathbf{x}] = -\frac{\tilde{\eta}}{KN} \sum_{k,i} g_i(\mathbf{y}_{i,k-1}).$$

We can then apply Lemma 12 to bound the second moment of the server update as

$$\begin{aligned} \mathbb{E}_{r-1}\|\mathbf{x} + \Delta\mathbf{x} - \mathbf{x}^*\|^2 &= \mathbb{E}_{r-1}\|\mathbf{x} - \mathbf{x}^*\|^2 - \frac{2\tilde{\eta}}{KS} \mathbb{E}_{r-1} \sum_{k,i \in \mathcal{S}} \langle \nabla f_i(\mathbf{y}_{i,k-1}), \mathbf{x} - \mathbf{x}^* \rangle + \mathbb{E}_{r-1}\|\Delta\mathbf{x}\|^2 \\ &\leq \underbrace{\frac{2\tilde{\eta}}{KS} \mathbb{E}_{r-1} \sum_{k,i \in \mathcal{S}} \langle \nabla f_i(\mathbf{y}_{i,k-1}), \mathbf{x}^* - \mathbf{x} \rangle}_{\mathcal{T}_4} + \mathbb{E}_{r-1}\|\mathbf{x} - \mathbf{x}^*\|^2 \\ &\quad + 8\beta\tilde{\eta}^2(\mathbb{E}[f(\mathbf{x}^{r-1})] - f(\mathbf{x}^*)) + 8\tilde{\eta}^2\mathcal{C}_{r-1} + 4\tilde{\eta}^2\beta^2\mathcal{E} + \frac{12\tilde{\eta}^2\sigma^2}{KS}. \end{aligned}$$

The term  $\mathcal{T}_4$  can be bounded by using perturbed strong-convexity (Lemma 5) with  $h = f_i$ ,  $\mathbf{x} = \mathbf{y}_{i,k-1}$ ,  $\mathbf{y} = \mathbf{x}^*$ , and  $\mathbf{z} = \mathbf{x}$  to get

$$\begin{aligned} \mathbb{E}[\mathcal{T}_4] &= \frac{2\tilde{\eta}}{KS} \mathbb{E} \sum_{k,i \in \mathcal{S}} \langle \nabla f_i(\mathbf{y}_{i,k-1}), \mathbf{x}^* - \mathbf{x} \rangle \\ &\leq \frac{2\tilde{\eta}}{KS} \mathbb{E} \sum_{k,i \in \mathcal{S}} \left( f_i(\mathbf{x}^*) - f_i(\mathbf{x}) + \beta\|\mathbf{y}_{i,k-1} - \mathbf{x}\|^2 - \frac{\mu}{4}\|\mathbf{x} - \mathbf{x}^*\|^2 \right) \\ &= -2\tilde{\eta} \mathbb{E} \left( f(\mathbf{x}) - f(\mathbf{x}^*) + \frac{\mu}{4}\|\mathbf{x} - \mathbf{x}^*\|^2 \right) + 2\beta\tilde{\eta}\mathcal{E}. \end{aligned}$$

Plugging  $\mathcal{T}_4$  back, we can further simplify the expression to get

$$\begin{aligned} \mathbb{E}\|\mathbf{x} + \Delta\mathbf{x} - \mathbf{x}^*\|^2 &\leq \mathbb{E}\|\mathbf{x} - \mathbf{x}^*\|^2 - 2\tilde{\eta} \left( f(\mathbf{x}) - f(\mathbf{x}^*) + \frac{\mu}{4}\|\mathbf{x} - \mathbf{x}^*\|^2 \right) + 2\beta\tilde{\eta}\mathcal{E} \\ &\quad + \frac{12\tilde{\eta}^2\sigma^2}{KS} + 8\beta\tilde{\eta}^2(\mathbb{E}[f(\mathbf{x}^{r-1})] - f(\mathbf{x}^*)) + 8\tilde{\eta}^2\mathcal{C}_{r-1} + 4\tilde{\eta}^2\beta^2\mathcal{E} \\ &= (1 - \frac{\mu\tilde{\eta}}{2})\|\mathbf{x} - \mathbf{x}^*\|^2 + (8\beta\tilde{\eta}^2 - 2\tilde{\eta})(f(\mathbf{x}) - f(\mathbf{x}^*)) \\ &\quad + \frac{12\tilde{\eta}^2\sigma^2}{KS} + (2\beta\tilde{\eta} + 4\beta^2\tilde{\eta}^2)\mathcal{E} + 8\tilde{\eta}^2\mathcal{C}_{r-1}. \end{aligned}$$

We can use Lemma 13 (scaled by  $9\tilde{\eta}^2\frac{N}{S}$ ) to bound the control-lag

$$9\tilde{\eta}^2\frac{N}{S}\mathcal{C}_r \leq (1 - \frac{\mu\tilde{\eta}}{2})9\tilde{\eta}^2\frac{N}{S}\mathcal{C}_{r-1} + 9(\frac{\mu\tilde{\eta}N}{2S} - 1)\tilde{\eta}^2\mathcal{C}_{r-1} + 9\tilde{\eta}^2(4\beta(\mathbb{E}[f(\mathbf{x}^{r-1})] - f(\mathbf{x}^*)) + 2\beta^2\mathcal{E})$$

Now recall that Lemma 14 bounds the client-drift:

$$3\beta\tilde{\eta}\mathcal{E}_r \leq \frac{2\tilde{\eta}^2}{3}\mathcal{C}_{r-1} + \frac{\tilde{\eta}}{25\eta_g^2}(\mathbb{E}[f(\mathbf{x}^{r-1})] - f(\mathbf{x}^*)) + \frac{\tilde{\eta}^2}{K\eta_g^2}\sigma^2.$$

Adding all three inequalities together,

$$\begin{aligned} \mathbb{E}\|\mathbf{x} + \Delta\mathbf{x} - \mathbf{x}^*\|^2 + \frac{9\tilde{\eta}^2N\mathcal{C}_r}{S} &\leq (1 - \frac{\mu\tilde{\eta}}{2})\left(\mathbb{E}\|\mathbf{x} - \mathbf{x}^*\|^2 + \frac{9\tilde{\eta}^2N\mathcal{C}_{r-1}}{S}\right) + (44\beta\tilde{\eta}^2 - \frac{49}{25}\tilde{\eta})(f(\mathbf{x}) - f(\mathbf{x}^*)) \\ &\quad + \frac{12\tilde{\eta}^2\sigma^2}{KS}(1 + \frac{S}{\eta_g^2}) + (22\beta^2\tilde{\eta}^2 - \beta\tilde{\eta})\mathcal{E} + (\frac{9\mu\tilde{\eta}N}{2S} - \frac{1}{3})\tilde{\eta}^2\mathcal{C}_{r-1} \end{aligned}$$

Finally, the lemma follows from noting that  $\tilde{\eta} \leq \frac{1}{81\beta}$  implies  $44\beta^2\tilde{\eta}^2 \leq \frac{24}{25}\tilde{\beta}$  and  $\tilde{\eta} \leq \frac{S}{15\mu N}$  implies  $\frac{9\mu\tilde{\eta}N}{2S} \leq \frac{1}{3}$ .  $\square$

**The final rate for strongly convex** follows simply by unrolling the recursive bound in Lemma 15 using Lemma 1. Also note that if  $c_i^0 = g_i(\mathbf{x}^0)$ , then  $\frac{\tilde{\eta}N}{S}\mathcal{C}_0$  can be bounded in terms of function sub-optimality  $F$ . For the **general convex** setting, averaging over  $r$  in Lemma 15 with  $\mu = 0$  gives

$$\begin{aligned} \frac{1}{R} \sum_{r=1}^R \mathbb{E}[f(\mathbf{x}^{r-1})] - f(\mathbf{x}^*) &\leq \frac{1}{\tilde{\eta}R} \|\mathbf{x}^0 - \mathbf{x}^*\|^2 + \frac{9N\tilde{\eta}}{SR} \mathcal{C}_0 + \frac{12\tilde{\eta}}{KS} (1 + \frac{S}{\eta_g^2}) \sigma^2 \\ &\leq 4\|\mathbf{x}^0 - \mathbf{x}^*\| \sigma \sqrt{\frac{3(1 + S/\eta_g^2)}{RKS}} \\ &\quad + \sqrt{\frac{N}{S} \frac{\|\mathbf{x}^0 - \mathbf{x}^*\|^2 + 9\mathcal{C}_0}{R}} + \frac{81\beta\|\mathbf{x}^0 - \mathbf{x}^*\|^2}{R} .. \end{aligned}$$

The last step follows from using a step size of  $\tilde{\eta} = \min\left(\frac{1}{81\beta}, \sqrt{\frac{S}{N}}, \frac{\|\mathbf{x}^0 - \mathbf{x}^*\|}{\sigma} \sqrt{\frac{KS}{12R(1 + \frac{S}{\eta_g^2})}}\right)$ .

## E.2. Convergence of SCAFFOLD for non-convex functions (Theorem III)

We now analyze the most general case of SCAFFOLD with option II on functions which are potentially non-convex. Just as in the non-convex proof, we will first bound the variance of the server update in Lemma 16, the change in control lag in Lemma 17 and finally we bound the client-drift in Lemma 18. Combining these three together gives us the progress made in one round in Lemma 19. The final rate is derived from the progress made using Lemma 2.

**Additional notation.** Recall that in round  $r$ , we update the control variate as (19)

$$\mathbf{c}_i^r = \begin{cases} \frac{1}{K} \sum_{k=1}^K g_i(\mathbf{y}_{i,k-1}^r) & \text{if } i \in \mathcal{S}^r, \\ \mathbf{c}_i^{r-1} & \text{otherwise.} \end{cases}$$

We introduce the following notation to keep track of the ‘lag’ in the update of the control variate: define a sequence of parameters  $\{\boldsymbol{\alpha}_{i,k-1}^{r-1}\}$  such that for any  $i \in [N]$  and  $k \in [K]$  we have  $\boldsymbol{\alpha}_{i,k-1}^0 := \mathbf{x}^0$  and for  $r \geq 1$ ,

$$\boldsymbol{\alpha}_{i,k-1}^r := \begin{cases} \mathbf{y}_{i,k-1}^r & \text{if } i \in \mathcal{S}^r, \\ \boldsymbol{\alpha}_{i,k-1}^{r-1} & \text{otherwise.} \end{cases} \quad (26)$$

By the update rule for control variates (19) and the definition of  $\{\boldsymbol{\alpha}_{i,k-1}^{r-1}\}$  above, the following property always holds:

$$\mathbf{c}_i^r = \frac{1}{K} \sum_{k=1}^K g_i(\boldsymbol{\alpha}_{i,k-1}^r).$$

We can then define the following  $\Xi_r$  to be the error in control variate for round  $r$ :

$$\Xi_r := \frac{1}{KN} \sum_{k=1}^K \sum_{i=1}^N \mathbb{E}\|\boldsymbol{\alpha}_{i,k-1}^r - \mathbf{x}^r\|^2. \quad (27)$$

Also recall the closely related definition of client drift caused by local updates:

$$\mathcal{E}_r := \frac{1}{KN} \sum_{k=1}^K \sum_{i=1}^N \mathbb{E}\|\mathbf{y}_{i,k}^r - \mathbf{x}^{r-1}\|^2.$$

**Variance of server update.** Let us analyze how the control variates effect the variance of the aggregate server update.

**Lemma 16.** *For updates (18)–(21) and assumptions A4 and A5, the following holds true for any  $\tilde{\eta} := \eta_l \eta_g K \in [0, 1/\beta]$ :*

$$\mathbb{E}\|\mathbb{E}_{r-1}[\mathbf{x}^r] - \mathbf{x}^{r-1}\|^2 \leq 2\tilde{\eta}^2 \beta^2 \mathcal{E}_r + 2\tilde{\eta}^2 \mathbb{E}\|\nabla f(\mathbf{x}^{r-1})\|^2, \text{ and}$$

$$\mathbb{E}\|\mathbf{x}^r - \mathbf{x}^{r-1}\|^2 \leq 4\tilde{\eta}^2 \beta^2 \mathcal{E}_r + 8\tilde{\eta}^2 \beta^2 \Xi_{r-1} + 4\tilde{\eta}^2 \mathbb{E}\|\nabla f(\mathbf{x}^{r-1})\|^2 + \frac{9\tilde{\eta}^2 \sigma^2}{KS}.$$

*Proof.* Recall that the server update satisfies

$$\mathbb{E}[\Delta \mathbf{x}] = -\frac{\tilde{\eta}}{KN} \sum_{k,i} \mathbb{E}[g_i(\mathbf{y}_{i,k-1})].$$

From the definition of  $\alpha_{i,k-1}^{r-1}$  and dropping the superscript everywhere we have

$$\Delta \mathbf{x} = -\frac{\tilde{\eta}}{KS} \sum_{k,i \in \mathcal{S}} (g_i(\mathbf{y}_{i,k-1}) + \mathbf{c} - \mathbf{c}_i) \text{ where } \mathbf{c}_i = \frac{1}{K} \sum_k g_i(\alpha_{i,k-1}).$$

Taking norm on both sides and separating mean and variance, we proceed as

$$\begin{aligned} \mathbb{E}\|\Delta \mathbf{x}\|^2 &= \mathbb{E}\left\|-\frac{\tilde{\eta}}{KS} \sum_{k,i \in \mathcal{S}} (g_i(\mathbf{y}_{i,k-1}) - g_i(\alpha_{i,k-1}) + \mathbf{c} - \mathbf{c}_i)\right\|^2 \\ &\leq \mathbb{E}\left\|-\frac{\tilde{\eta}}{KS} \sum_{k,i \in \mathcal{S}} (\nabla f_i(\mathbf{y}_{i,k-1}) + \mathbb{E}[\mathbf{c}] - \mathbb{E}[\mathbf{c}_i])\right\|^2 + \frac{9\tilde{\eta}^2\sigma^2}{KS} \\ &\leq \mathbb{E}\left[\frac{\tilde{\eta}^2}{KS} \sum_{k,i \in \mathcal{S}} \left\|\nabla f_i(\mathbf{y}_{i,k-1}) + \mathbb{E}[\mathbf{c}] - \mathbb{E}[\mathbf{c}_i]\right\|^2\right] + \frac{9\tilde{\eta}^2\sigma^2}{KS} \\ &= \frac{\tilde{\eta}^2}{KN} \sum_{k,i} \mathbb{E}\left\|(\nabla f_i(\mathbf{y}_{i,k-1}) - \nabla f_i(\mathbf{x})) + (\mathbb{E}[\mathbf{c}] - \nabla f(\mathbf{x})) + \nabla f(\mathbf{x}) - (\mathbb{E}[\mathbf{c}_i] - \nabla f_i(\mathbf{x}))\right\|^2 + \frac{9\tilde{\eta}^2\sigma^2}{KS} \\ &\leq \frac{4\tilde{\eta}^2}{KN} \sum_{k,i} \mathbb{E}\|\nabla f_i(\mathbf{y}_{i,k-1}) - \nabla f_i(\mathbf{x})\|^2 + \frac{8\tilde{\eta}^2}{KN} \sum_{k,i} \mathbb{E}\|\nabla f_i(\alpha_{i,k-1}) - \nabla f_i(\mathbf{x})\|^2 \\ &\quad + 4\tilde{\eta}^2 \mathbb{E}\|\nabla f(\mathbf{x})\|^2 + \frac{9\tilde{\eta}^2\sigma^2}{KS} \\ &\leq 4\tilde{\eta}^2\beta^2\mathcal{E}_r + 8\beta^2\tilde{\eta}^2\Xi_{r-1} + 4\tilde{\eta}^2 \mathbb{E}\|\nabla f(\mathbf{x})\|^2 + \frac{9\tilde{\eta}^2\sigma^2}{KS}. \end{aligned}$$

In the first inequality, note that the three random variables— $\frac{1}{KS} \sum_{k,i \in \mathcal{S}} g_i(\mathbf{y}_{i,k})$ ,  $\frac{1}{S} \sum_{i \in \mathcal{S}} \mathbf{c}_i$ , and  $\mathbf{c}$ —may not be independent but each have variance smaller than  $\frac{\sigma^2}{KS}$  and so we can apply Lemma 4. The rest of the inequalities follow from repeated applications of the relaxed triangle inequality,  $\beta$ -Lipschitzness of  $f_i$ , and the definition of  $\Xi_{r-1}$  (27). This proves the second statement. The first statement follows from our expression of  $\mathbb{E}_{r-1}[\Delta \mathbf{x}]$  and similar computations.  $\square$

**Lag in the control variates.** We now analyze the ‘lag’ in the control variates due to us sampling only a small subset of clients each round. Because we cannot rely on convexity anymore but only on the Lipschitzness of the gradients, the control-lag increases faster in the non-convex case.

**Lemma 17.** *For updates (18)–(21) and assumptions A4, A5, the following holds true for any  $\tilde{\eta} \leq \frac{1}{24\beta}(\frac{S}{N})^\alpha$  for  $\alpha \in [\frac{1}{2}, 1]$  where  $\tilde{\eta} := \eta_l \eta_g K$ :*

$$\Xi_r \leq (1 - \frac{17S}{36N})\Xi_{r-1} + \frac{1}{48\beta^2}(\frac{S}{N})^{2\alpha-1} \|\nabla f(\mathbf{x}^{r-1})\|^2 + \frac{97}{48}(\frac{S}{N})^{2\alpha-1} \mathcal{E}_r + (\frac{S}{N\beta^2}) \frac{\sigma^2}{32KS}.$$

*Proof.* The proof proceeds similar to that of Lemma 13 except that we cannot rely on convexity. Recall that after round  $r$ , the definition of  $\alpha_{i,k-1}^r$  (26) implies that

$$\mathbb{E}_{\mathcal{S}^r}[\alpha_{i,k-1}^r] = (1 - \frac{S}{N})\alpha_{i,k-1}^{r-1} + \frac{S}{N}\mathbf{y}_{i,k-1}^r.$$

Plugging the above expression in the definition of  $\Xi_r$  we get

$$\begin{aligned} \Xi_r &= \frac{1}{KN} \sum_{i,k} \mathbb{E}\|\alpha_{i,k-1}^r - \mathbf{x}^r\|^2 \\ &= \left(1 - \frac{S}{N}\right) \cdot \underbrace{\frac{1}{KN} \sum_i \mathbb{E}\|\alpha_{i,k-1}^{r-1} - \mathbf{x}^r\|^2}_{\mathcal{T}_5} + \frac{S}{N} \cdot \underbrace{\frac{1}{KN} \sum_{k,i} \mathbb{E}\|\mathbf{y}_{i,k-1}^r - \mathbf{x}^r\|^2}_{\mathcal{T}_6}. \end{aligned}$$

We can expand the second term  $\mathcal{T}_6$  with the relaxed triangle inequality to claim

$$\mathcal{T}_6 \leq 2(\mathcal{E}_r + \mathbb{E}\|\Delta\mathbf{x}^r\|^2).$$

We will expand the first term  $\mathcal{T}_5$  to claim for a constant  $b \geq 0$  to be chosen later

$$\begin{aligned} \mathcal{T}_5 &= \frac{1}{KN} \sum_i \mathbb{E}(\|\boldsymbol{\alpha}_{i,k-1}^{r-1} - \mathbf{x}^{r-1}\|^2 + \|\Delta\mathbf{x}^r\|^2 + \mathbb{E}_{r-1} \langle \Delta\mathbf{x}^r, \boldsymbol{\alpha}_{i,k-1}^{r-1} - \mathbf{x}^{r-1} \rangle) \\ &\leq \frac{1}{KN} \sum_i \mathbb{E}(\|\boldsymbol{\alpha}_{i,k-1}^{r-1} - \mathbf{x}^{r-1}\|^2 + \|\Delta\mathbf{x}^r\|^2 + \frac{1}{b} \|\mathbb{E}_{r-1}[\Delta\mathbf{x}^r]\|^2 + b\|\boldsymbol{\alpha}_{i,k-1}^{r-1} - \mathbf{x}^{r-1}\|^2) \end{aligned}$$

where we used Young's inequality which holds for any  $b \geq 0$ . Combining the bounds for  $\mathcal{T}_5$  and  $\mathcal{T}_6$ ,

$$\begin{aligned} \Xi_r &\leq (1 - \frac{S}{N})(1 + b)\Xi_{r-1} + 2\frac{S}{N}\mathcal{E}_r + 2\mathbb{E}\|\Delta\mathbf{x}^r\|^2 + \frac{1}{b}\mathbb{E}\|\mathbb{E}_{r-1}[\Delta\mathbf{x}^r]\|^2 \\ &\leq ((1 - \frac{S}{N})(1 + b) + 16\tilde{\eta}^2\beta^2)\Xi_{r-1} + (\frac{2S}{N} + 8\tilde{\eta}^2\beta^2 + 2\frac{1}{b}\tilde{\eta}^2\beta^2)\mathcal{E}_r + (8 + 2\frac{1}{b})\tilde{\eta}^2\mathbb{E}\|\nabla f(\mathbf{x})\|^2 + \frac{18\tilde{\eta}^2\sigma^2}{KS} \end{aligned}$$

The last inequality applied Lemma 16. Verify that with choice of  $b = \frac{S}{2(N-S)}$ , we have  $(1 - \frac{S}{N})(1 + b) \leq (1 - \frac{S}{2N})$  and  $\frac{1}{b} \leq \frac{2N}{S}$ . Plugging these values along with the bound on the step-size  $16\beta^2\tilde{\eta}^2 \leq \frac{1}{36}(\frac{S}{N})^{2\alpha} \leq \frac{S}{36N}$  completes the lemma.  $\square$

**Bounding the drift.** We will next bound the client drift  $\mathcal{E}_r$ . For this, convexity is not crucial and we will recover a very similar result to Lemma 14 only use the Lipschitzness of the gradient.

**Lemma 18.** Suppose our step-sizes satisfy  $\eta_l \leq \frac{1}{24\beta K\eta_g}$  and  $f_i$  satisfies assumptions A4–A5. Then, for any global  $\eta_g \geq 1$  we can bound the drift as

$$\frac{5}{3}\beta^2\tilde{\eta}\mathcal{E}_r \leq \frac{5}{3}\beta^3\tilde{\eta}^2\Xi_{r-1} + \frac{\tilde{\eta}}{24\eta_g^2}\mathbb{E}\|\nabla f(\mathbf{x}^{r-1})\|^2 + \frac{\tilde{\eta}^2\beta}{4K\eta_g^2}\sigma^2.$$

*Proof.* First, observe that if  $K = 1$ ,  $\mathcal{E}_r = 0$  since  $\mathbf{y}_{i,0} = \mathbf{x}$  for all  $i \in [N]$  and that  $\Xi_{r-1}$  and the right hand side are both positive. Thus the Lemma is trivially true if  $K = 1$  and we will henceforth assume  $K \geq 2$ . Starting from the update rule (18) for  $i \in [N]$  and  $k \in [K]$

$$\begin{aligned} \mathbb{E}\|\mathbf{y}_{i,k} - \mathbf{x}\|^2 &= \mathbb{E}\|\mathbf{y}_{i,k-1} - \eta_l(g_i(\mathbf{y}_{i,k-1}) + \mathbf{c} - \mathbf{c}_i) - \mathbf{x}\|^2 \\ &\leq \mathbb{E}\|\mathbf{y}_{i,k-1} - \eta_l(\nabla f_i(\mathbf{y}_{i,k-1}) + \mathbf{c} - \mathbf{c}_i) - \mathbf{x}\|^2 + \eta_l^2\sigma^2 \\ &\leq (1 + \frac{1}{K-1})\mathbb{E}\|\mathbf{y}_{i,k-1} - \mathbf{x}\|^2 + K\eta_l^2\mathbb{E}\|\nabla f_i(\mathbf{y}_{i,k-1}) + \mathbf{c} - \mathbf{c}_i\|^2 + \eta_l^2\sigma^2 \\ &= (1 + \frac{1}{K-1})\mathbb{E}\|\mathbf{y}_{i,k-1} - \mathbf{x}\|^2 + \eta_l^2\sigma^2 \\ &\quad + K\eta_l^2\mathbb{E}\|\nabla f_i(\mathbf{y}_{i,k-1}) - \nabla f_i(\mathbf{x}) + (\mathbf{c} - \nabla f(\mathbf{x})) + \nabla f(\mathbf{x}) - (\mathbf{c}_i - \nabla f_i(\mathbf{x}))\|^2 \\ &\leq (1 + \frac{1}{K-1})\mathbb{E}\|\mathbf{y}_{i,k-1} - \mathbf{x}\|^2 + 4K\eta_l^2\mathbb{E}\|\nabla f_i(\mathbf{y}_{i,k-1}) - \nabla f_i(\mathbf{x})\|^2 + \eta_l^2\sigma^2 \\ &\quad + 4K\eta_l^2\mathbb{E}\|\mathbf{c} - \nabla f(\mathbf{x})\|^2 + 4K\eta_l^2\mathbb{E}\|\nabla f(\mathbf{x})\|^2 + 4K\eta_l^2\mathbb{E}\|\mathbf{c}_i - \nabla f_i(\mathbf{x})\|^2 \\ &\leq (1 + \frac{1}{K-1} + 4K\beta^2\eta_l^2)\mathbb{E}\|\mathbf{y}_{i,k-1} - \mathbf{x}\|^2 + \eta_l^2\sigma^2 + 4K\eta_l^2\mathbb{E}\|\nabla f(\mathbf{x})\|^2 \\ &\quad + 4K\eta_l^2\mathbb{E}\|\mathbf{c} - \nabla f(\mathbf{x})\|^2 + 4K\eta_l^2\mathbb{E}\|\mathbf{c}_i - \nabla f_i(\mathbf{x})\|^2 \end{aligned}$$

The inequalities above follow from repeated application of the relaxed triangle inequalities and the  $\beta$ -Lipschitzness of  $f_i$ .

Averaging the above over  $i$ , the definition of  $\mathbf{c} = \frac{1}{N} \sum_i \mathbf{c}_i$  and  $\Xi_{r-1}$  (27) gives

$$\begin{aligned} \frac{1}{N} \sum_i \mathbb{E} \|\mathbf{y}_{i,k} - \mathbf{x}\|^2 &\leq (1 + \frac{1}{K-1} + 4K\beta^2\eta_l^2) \frac{1}{N} \sum_i \mathbb{E} \|\mathbf{y}_{i,k-1} - \mathbf{x}\|^2 \\ &\quad + \eta_l^2\sigma^2 + 4K\eta_l^2 \mathbb{E} \|\nabla f(\mathbf{x})\|^2 + 8K\eta_l^2\beta^2\Xi_{r-1} \\ &\leq (\eta_l^2\sigma^2 + 4K\eta_l^2 \mathbb{E} \|\nabla f(\mathbf{x})\|^2 + 8K\eta_l^2\beta^2\Xi_{r-1}) \left( \sum_{\tau=0}^{k-1} (1 + \frac{1}{K-1} + 4K\beta^2\eta_l^2)^{\tau} \right) \\ &= \left( \frac{\tilde{\eta}^2\sigma^2}{K^2\eta_g^2} + \frac{4\tilde{\eta}^2}{K\eta_g^2} \mathbb{E} \|\nabla f(\mathbf{x})\|^2 + \frac{8\tilde{\eta}^2\beta^2}{K\eta_g^2} \Xi_{r-1} \right) \left( \sum_{\tau=0}^{k-1} (1 + \frac{1}{K-1} + \frac{4\beta^2\tilde{\eta}^2}{K\eta_g^2})^{\tau} \right) \\ &\leq \left( \frac{\tilde{\eta}\sigma^2}{24\beta K^2\eta_g^2} + \frac{1}{144\beta^2 K\eta_g^2} \mathbb{E} \|\nabla f(\mathbf{x})\|^2 + \frac{\tilde{\eta}\beta}{3K\eta_g^2} \Xi_{r-1} \right) 3K. \end{aligned}$$

The last inequality used the bound on the step-size  $\beta\tilde{\eta} \leq \frac{1}{24}$ . Averaging over  $k$  and multiplying both sides by  $\frac{5}{3}\beta^2\tilde{\eta}$  yields the lemma statement.  $\square$

**Progress made in each round.** Given that we can bound all sources of error, we can finally prove the progress made in each round.

**Lemma 19.** Suppose the updates (18)–(21) satisfy assumptions A4–A5. For any effective step-size  $\tilde{\eta} := K\eta_g\eta_l$  satisfying  $\tilde{\eta} \leq \frac{1}{24\beta} \left(\frac{S}{N}\right)^{\frac{2}{3}}$ ,

$$\left( \mathbb{E}[f(\mathbf{x}^r)] + 12\beta^3\tilde{\eta}^2\frac{N}{S}\Xi_r \right) \leq \left( \mathbb{E}[f(\mathbf{x}^{r-1})] + 12\beta^3\tilde{\eta}^2\frac{N}{S}\Xi_{r-1} \right) + \frac{5\beta\tilde{\eta}^2\sigma^2}{KS} \left( 1 + \frac{S}{\eta_g^2} \right) - \frac{\tilde{\eta}}{14} \mathbb{E} \|\nabla f(\mathbf{x}^{r-1})\|^2.$$

*Proof.* Starting from the smoothness of  $f$  and taking conditional expectation gives

$$\mathbb{E}_{r-1}[f(\mathbf{x} + \Delta\mathbf{x})] \leq f(\mathbf{x}) + \langle \nabla f(\mathbf{x}), \mathbb{E}_{r-1}[\Delta\mathbf{x}] \rangle + \frac{\beta}{2} \mathbb{E}_{r-1} \|\Delta\mathbf{x}\|^2.$$

We as usual dropped the superscript everywhere. Recall that the server update can be written as

$$\Delta\mathbf{x} = -\frac{\tilde{\eta}}{KS} \sum_{k,i \in \mathcal{S}} (g_i(\mathbf{y}_{i,k-1}) + \mathbf{c} - \mathbf{c}_i), \text{ and } \mathbb{E}_{\mathcal{S}}[\Delta\mathbf{x}] = -\frac{\tilde{\eta}}{KN} \sum_{k,i} g_i(\mathbf{y}_{i,k-1}).$$

Substituting this in the previous inequality and applying Lemma 16 to bound  $\mathbb{E}[\|\Delta\mathbf{x}\|^2]$  gives

$$\begin{aligned} \mathbb{E}[f(\mathbf{x} + \Delta\mathbf{x})] - f(\mathbf{x}) &\leq -\frac{\tilde{\eta}}{KN} \sum_{k,i} \langle \nabla f(\mathbf{x}), \mathbb{E}[\nabla f_i(\mathbf{y}_{i,k-1})] \rangle + \frac{\beta}{2} \mathbb{E} \|\Delta\mathbf{x}\|^2 \\ &\leq -\frac{\tilde{\eta}}{KN} \sum_{k,i} \langle \nabla f(\mathbf{x}), \mathbb{E}[\nabla f_i(\mathbf{y}_{i,k-1})] \rangle + \\ &\quad 2\tilde{\eta}^2\beta^3\mathcal{E}_r + 4\tilde{\eta}^2\beta^3\Xi_{r-1} + 2\beta\tilde{\eta}^2 \mathbb{E} \|\nabla f(\mathbf{x})\|^2 + \frac{9\beta\tilde{\eta}^2\sigma^2}{2KS} \\ &\leq -\frac{\tilde{\eta}}{2} \|\nabla f(\mathbf{x})\|^2 + \frac{\tilde{\eta}}{2} \sum_{i,k} \mathbb{E} \left\| \frac{1}{KN} \sum_{i,k} \nabla f_i(\mathbf{y}_{i,k-1}) - \nabla f(\mathbf{x}) \right\|^2 + \\ &\quad 2\tilde{\eta}^2\beta^3\mathcal{E}_r + 4\tilde{\eta}^2\beta^3\Xi_{r-1} + 2\beta\tilde{\eta}^2 \mathbb{E} \|\nabla f(\mathbf{x})\|^2 + \frac{9\beta\tilde{\eta}^2\sigma^2}{2KS} \\ &\leq -\frac{\tilde{\eta}}{2} \|\nabla f(\mathbf{x})\|^2 + \frac{\tilde{\eta}}{2KN} \sum_{i,k} \mathbb{E} \left\| \nabla f_i(\mathbf{y}_{i,k-1}) - \nabla f_i(\mathbf{x}) \right\|^2 + \\ &\quad 2\tilde{\eta}^2\beta^3\mathcal{E}_r + 4\tilde{\eta}^2\beta^3\Xi_{r-1} + 2\beta\tilde{\eta}^2 \mathbb{E} \|\nabla f(\mathbf{x})\|^2 + \frac{9\beta\tilde{\eta}^2\sigma^2}{2KS} \\ &\leq -\left(\frac{\tilde{\eta}}{2} - 2\beta\tilde{\eta}^2\right) \|\nabla f(\mathbf{x})\|^2 + \left(\frac{\tilde{\eta}}{2} + 2\beta\tilde{\eta}^2\right) \beta^2\mathcal{E}_r + 4\beta^3\tilde{\eta}^2\Xi_{r-1} + \frac{9\beta\tilde{\eta}^2\sigma^2}{2KS}. \end{aligned}$$

The third inequality follows from the observation that  $-ab = \frac{1}{2}((b-a)^2 - a^2) - \frac{1}{2}b^2 \leq \frac{1}{2}((b-a)^2 - a^2)$  for any  $a, b \in \mathbb{R}$ , and the last from the  $\beta$ -Lipschitzness of  $f_i$ . Now we use Lemma 17 to bound  $\Xi_r$  as

$$\begin{aligned} 12\beta^3\tilde{\eta}^2\frac{N}{S}\Xi_r &\leq 12\beta^3\tilde{\eta}^2\frac{N}{S}\left(\left(1 - \frac{17S}{36N}\right)\Xi_{r-1} + \frac{1}{48\beta^2}\left(\frac{S}{N}\right)^{2\alpha-1}\|\nabla f(\mathbf{x}^{r-1})\|^2 + \frac{97}{48}\left(\frac{S}{N}\right)^{2\alpha-1}\mathcal{E}_r + \left(\frac{S}{N\beta^2}\right)\frac{\sigma^2}{32KS}\right) \\ &= 12\beta^3\tilde{\eta}^2\frac{N}{S}\Xi_{r-1} - \frac{17}{3}\beta^3\tilde{\eta}^2\Xi_{r-1} + \frac{1}{4}\beta\tilde{\eta}^2\left(\frac{N}{S}\right)^{2-2\alpha}\|\nabla f(\mathbf{x})\|^2 + \frac{97}{4}\beta^3\tilde{\eta}^2\left(\frac{N}{S}\right)^{2-2\alpha}\mathcal{E}_r + \frac{3\beta\tilde{\eta}^2\sigma^2}{8KS}. \end{aligned}$$

Also recall that Lemma 18 states that

$$\frac{5}{3}\beta^2\tilde{\eta}\mathcal{E}_r \leq \frac{5}{3}\beta^3\tilde{\eta}^2\Xi_{r-1} + \frac{\tilde{\eta}}{24\eta_g^2}\mathbb{E}\|\nabla f(\mathbf{x}^{r-1})\|^2 + \frac{\tilde{\eta}^2\beta}{4K\eta_g^2}\sigma^2.$$

Adding these bounds on  $\Xi_r$  and  $\mathcal{E}_r$  to that of  $\mathbb{E}[f(\mathbf{x} + \Delta\mathbf{x})]$  gives

$$\begin{aligned} (\mathbb{E}[f(\mathbf{x} + \Delta\mathbf{x})] + 12\beta^3\tilde{\eta}^2\frac{N}{S}\Xi_r) &\leq (\mathbb{E}[f(\mathbf{x})] + 12\beta^3\tilde{\eta}^2\frac{N}{S}\Xi_{r-1}) + \left(\frac{5}{3} - \frac{17}{3}\right)\beta^3\tilde{\eta}^2\Xi_{r-1} \\ &\quad - \left(\frac{\tilde{\eta}}{2} - 2\beta\tilde{\eta}^2 - \frac{1}{4}\beta\tilde{\eta}^2\left(\frac{N}{S}\right)^{2-2\alpha}\|\nabla f(\mathbf{x})\|^2 + \left(\frac{\tilde{\eta}}{2} - \frac{5\tilde{\eta}}{3} + 2\beta\tilde{\eta}^2 + \frac{97}{4}\beta\tilde{\eta}^2\left(\frac{N}{S}\right)^{2-2\alpha}\right)\beta^2\mathcal{E}_r + \frac{39\beta\tilde{\eta}^2\sigma^2}{8KS}(1 + \frac{S}{\eta_g^2})\right). \end{aligned}$$

By our choice of  $\alpha = \frac{2}{3}$  and plugging in the bound on step-size  $\beta\tilde{\eta}\left(\frac{N}{S}\right)^{2-2\alpha} \leq \frac{1}{24}$  proves the lemma.  $\square$

The **non-convex rate** of convergence now follows by unrolling the recursion in Lemma 19 and selecting an appropriate step-size  $\tilde{\eta}$  as in Lemma 2. Finally note that if we initialize  $c_i^0 = g_i(\mathbf{x}^0)$  then we have  $\Xi_0 = 0$ .

## F. Usefulness of local steps (Theorem IV)

Let us state our rates of convergence for SCAFFOLD which interpolates between identical and completely heterogeneous clients. In this section, we always set  $\eta_g = 1$  and assume all clients participate ( $S = N$ ).

**Theorem VIII.** Suppose that the functions  $\{f_i\}$  are quadratic and satisfy assumptions A4, A5 and additionally A2. Then, for global step-size  $\eta_g = 1$  in each of the following cases, there exist probabilities  $\{p_k^r\}$  and local step-size  $\eta_l$  such that the output (29) of SCAFFOLD when run with no client sampling ( $S = N$ ) using update (28) satisfies:

- **Strongly convex:**  $f_i$  satisfies (A3) for  $\mu > 0$ ,  $\eta_l \leq \min(\frac{1}{10\beta}, \frac{1}{22\delta K}, \frac{1}{10\mu K})$ ,  $R \geq \max(\frac{20\beta}{\mu}, \frac{44\delta K + 20\mu K}{\mu}, 20K)$  then

$$\mathbb{E}[\|\nabla f(\bar{\mathbf{x}}^R)\|^2] \leq \tilde{\mathcal{O}}\left(\frac{\beta\sigma^2}{\mu RKN} + \mu D^2 \exp\left(-\frac{\mu}{20\beta + 44\delta K + 20\mu K}RK\right)\right).$$

- **General convex:**  $f$  satisfies  $\nabla^2 f \succeq -\delta I$ ,  $\eta_l \leq \min(\frac{1}{10\beta}, \frac{1}{22\delta K})$ , and  $R \geq 1$ , then

$$\mathbb{E}[\|\nabla f(\bar{\mathbf{x}}^R)\|^2] \leq \mathcal{O}\left(\frac{\sigma\sqrt{\beta(f(\mathbf{x}_0) - f^*)}}{\sqrt{RKN}} + \frac{(\beta + \delta K)(f(\mathbf{x}_0) - f^*)}{RK}\right).$$

Note that if  $\delta = 0$ , we match (up to acceleration) the lower bound in (Woodworth et al., 2018). While certainly  $\delta = 0$  when the functions are identical as studied in (Woodworth et al., 2018), our upper-bound is significantly stronger since it is possible that  $\delta = 0$  even for highly heterogeneous functions. For example, objective perturbation (Chaudhuri et al., 2011; Kifer et al., 2012) is an optimal mechanism to achieve differential privacy for smooth convex objectives (Bassily et al., 2014). Intuitively, objective perturbation relies on masking each client's gradients by adding a large random linear term to the objective function. In such a case, we would have high gradient dissimilarity but no Hessian dissimilarity.

Our non-convex convergence rates are the first of their kind as far as we are aware—no previous work shows how one can take advantage of similarity for non-convex functions. However, we should note that non-convex quadratics do not have a global lower-bound on the function value  $f^*$ . We will instead assume that  $f^*$  almost surely lower-bounds the value of  $f(\mathbf{x}^R)$ , implicitly assuming that the iterates remain bounded.

**Outline.** In the rest of this section, we will focus on proving Theorem VIII. We will show how to bound variance in Lemma 23, bound the amount of drift in Lemma 22, and show progress made in one step in Lemma 24. In all of these we do not use convexity, but strongly rely on the functions being quadratics. Then we combine these to derive the progress made by the server in one round—for this we need *weak*-convexity to argue that averaging the parameters does not hurt convergence too much. As before, it is straight-forward to derive rates of convergence from the one-round progress using Lemmas 1 and 2.

### F.1. Additional notation and assumptions

For any matrix  $M$  and vector  $\mathbf{v}$ , let  $\|\mathbf{v}\|_M^2 := \mathbf{v}^\top M \mathbf{v}$ . Since all functions in this section are quadratics, we can assume w.l.o.g they are of the following form:

$$f_i(\mathbf{x}) - f_i(\mathbf{x}_i^*) = \frac{1}{2}\|\mathbf{x} - \mathbf{x}_i^*\|_{A_i}^2 \text{ for } i \in [N], \text{ and } f(\mathbf{x}) = \frac{1}{2}\|\mathbf{x} - \mathbf{x}^*\|_A^2, \text{ for all } \mathbf{x},$$

for some  $\{\mathbf{x}_i^*\}$  and  $\mathbf{x}^*$ ,  $A := \frac{1}{N} \sum_{i=1}^N A_i$ . We also assume that  $A$  is a symmetric matrix though this requirement is easily relaxed. Note that this implies  $f(\mathbf{x}^*) = 0$  and that  $\nabla f_i(\mathbf{x}) = A(\mathbf{x} - \mathbf{x}_i^*)$ . If  $\{f_i\}$  are additionally convex, we have that  $\mathbf{x}_i^*$  is the optimum of  $f_i$  and  $\mathbf{x}^*$  the optimum of  $f$ . However, this is not necessarily true in general.

We will also focus on a simplified version of SCAFFOLD where in each round  $r$ , client  $i$  performs the following update starting from  $\mathbf{y}_{i,0}^r \leftarrow \mathbf{x}^{r-1}$ :

$$\begin{aligned} \mathbf{y}_{i,k}^r &= \mathbf{y}_{i,k-1}^r - \eta(g_i(\mathbf{y}_{i,k-1}^r) + \nabla f(\mathbf{x}^{r-1}) - \nabla f_i(\mathbf{x}^{r-1})), \text{ i.e.} \\ \mathbb{E}_{r-1,k-1}[\mathbf{y}_{i,k}^r] &= \mathbf{y}_{i,k-1}^r - \eta A(\mathbf{y}_{i,k-1}^r - \mathbf{x}^*) - \eta(A_i - A)(\mathbf{y}_{i,k-1}^r - \mathbf{x}^{r-1}), \end{aligned} \tag{28}$$

where the second part is specialized to quadratics and the expectation is conditioned over everything before current step  $k$  of round  $r$ . At the end of each round, as before,  $\mathbf{x}^r = \frac{1}{N} \sum_{i=1}^N \mathbf{y}_{i,K}^r$ . The final output of the algorithm is chosen using probabilities  $\{p_k^r\}$  as

$$\bar{\mathbf{x}}^R = \mathbf{x}_k^r \text{ with probability } p_k^r, \text{ where } \mathbf{x}_k^r := \frac{1}{N} \sum_{i=1}^N \mathbf{y}_{i,k}^r. \quad (29)$$

Note that we are now possibly outputting iterates computed within a single round and that  $\mathbf{x}^r = \mathbf{x}_K^r$ . Beyond this, the update above differs from our usual SCAFFOLD in two key aspects: a) it uses gradients computed at  $\mathbf{x}^{r-1}$  as control variates instead of those at either  $\mathbf{x}^{r-2}$  (as in option I) or  $\mathbf{y}_{i,k}^{r-1}$  (as in option II), and b) it uses full batch gradients to compute its control variates instead of stochastic gradients. The first issue is easy to fix and our proof extends to using both option I or option II using techniques in Section E. The second issue is more technical—using stochastic gradients for control variates couples the randomness across the clients in making the local-updates *biased*. While it may be possible to get around this (cf. (Lei & Jordan, 2017; Nguyen et al., 2017; Tran-Dinh et al., 2019)), we will not attempt to do so in this work. Note that if  $K$  local update steps typically represents running multiple epochs on each client. Hence one additional epoch to compute the control variate  $\nabla f_i(\mathbf{x})$  does not significantly add to the cost.

Finally, we define the following sequence of positive numbers for notation convenience:

$$\begin{aligned} \xi_{i,k}^r &:= (\mathbb{E}_{r-1}[f(\mathbf{y}_{i,k}^r)] - f(\mathbf{x}^*) + \delta(1 + \frac{1}{K})^{K-k} \mathbb{E}_{r-1}\|\mathbf{y}_{i,k}^r - \mathbf{x}^{r-1}\|^2), \text{ and} \\ \tilde{\xi}_{i,k}^r &:= ([f(\mathbb{E}_{r-1}[\mathbf{y}_{i,k}^r])] - f(\mathbf{x}^*) + \delta(1 + \frac{1}{K})^{K-k} \mathbb{E}_{r-1,k-1}\|\mathbb{E}_{r-1}[\mathbf{y}_{i,k}^r] - \mathbf{x}^{r-1}\|^2). \end{aligned}$$

Observe that for  $k = 0$ ,  $\xi_{i,0}^r = \tilde{\xi}_{i,0}^r = f(\mathbf{x}^{r-1}) - f(\mathbf{x}^*)$ .

## F.2. Lemmas tracking errors

**Effect of averaging.** We see how averaging can reduce variance. A similar argument was used in the special case of one-shot averaging in (Zhang et al., 2013b).

**Lemma 20.** *Suppose  $\{f_i\}$  are quadratic functions and assumption A4 is satisfied. Then let  $\mathbf{x}_k^r$  and  $\mathbf{y}_{i,k}^r$  be vectors in step  $k$  and round  $r$  generated using (28)–(29). Then,*

$$\mathbb{E}_{r-1}\|\nabla f(\mathbf{x}_k^r)\|^2 \leq \frac{1}{N} \sum_{i=1}^N \|\nabla f(\mathbb{E}_{r-1}[\mathbf{y}_{i,k}^r])\|^2 + \frac{1}{N^2} \sum_{i=1}^N \mathbb{E}_{r-1}[\|\nabla f(\mathbf{y}_{i,k}^r)\|^2].$$

*Proof.* Observe that the variables  $\{\mathbf{y}_{i,k} - \mathbf{x}\}$  are independent of each other (the only source of randomness is the local gradient computations). The rest of the proof is exactly that of Lemma 4. Dropping superscripts everywhere,

$$\begin{aligned} \mathbb{E}_{r-1}\|A(\mathbf{x}_k^r - \mathbf{x}^*)\|^2 &= \mathbb{E}_{r-1}\left\|\frac{1}{N} \sum_i A(\mathbf{y}_{i,k} - \mathbf{x}^*)\right\|^2 \\ &= \mathbb{E}_{r-1}\left\|\frac{1}{N} \sum_i A(\mathbb{E}_{r-1}[\mathbf{y}_{i,k}] - \mathbf{x}^*)\right\|^2 + \mathbb{E}_{r-1}\left\|\frac{1}{N} \sum_i A(\mathbb{E}_{r-1}[\mathbf{y}_{i,k}] - \mathbf{y}_{i,k})\right\|^2 \\ &= \mathbb{E}_{r-1}\left\|\frac{1}{N} \sum_i A(\mathbb{E}_{r-1}[\mathbf{y}_{i,k}] - \mathbf{x}^*)\right\|^2 + \frac{1}{N^2} \sum_i \mathbb{E}_{r-1}\|A(\mathbb{E}_{r-1}[\mathbf{y}_{i,k}] - \mathbf{y}_{i,k})\|^2 \\ &= \mathbb{E}_{r-1}\left\|\frac{1}{N} \sum_i A(\mathbb{E}_{r-1}[\mathbf{y}_{i,k}] - \mathbf{x}^*)\right\|^2 + \frac{1}{N^2} \sum_i \mathbb{E}_{r-1}\|A(\mathbf{y}_{i,k} - \mathbf{x}^* - \mathbb{E}_{r-1}[\mathbf{y}_{i,k} - \mathbf{x}^*])\|^2 \\ &\leq \mathbb{E}_{r-1}\left\|\frac{1}{N} \sum_i A(\mathbb{E}_{r-1}[\mathbf{y}_{i,k}] - \mathbf{x}^*)\right\|^2 + \frac{1}{N^2} \sum_i \mathbb{E}_{r-1}\|A(\mathbf{y}_{i,k} - \mathbf{x}^*)\|^2. \end{aligned}$$

The third equality was because  $\{\mathbf{y}_{i,k}\}$  are independent of each other conditioned on everything before round  $r$ .  $\square$

We next see the effect of averaging on function values.

**Lemma 21.** *Suppose that  $f$  is  $\delta$  general-convex, then we have:*

$$\frac{1}{N} \sum_{i=1}^n \xi_{i,k}^r \geq \mathbb{E}_{r-1}[f(\mathbf{x}_k^r)] - f(\mathbf{x}^*), \text{ and } \frac{1}{N} \sum_{i=1}^n \tilde{\xi}_{i,k}^r \geq f(\mathbb{E}_{r-1}[\mathbf{x}_k^r]) - f(\mathbf{x}^*).$$

*Proof.* Since  $f$  is  $\delta$ -general convex, it follows that the function  $f(\mathbf{z}) + \delta(1 + \frac{1}{K})^{K-k} \|\mathbf{z} - \mathbf{x}\|_2^2$  is convex in  $\mathbf{z}$  for any  $k \in [K]$ . The lemma now follows directly from using convexity and the definition of  $\mathbf{x}_k^r = \frac{1}{N} \mathbf{y}_{i,k}^r$ .  $\square$

**Bounding drift of one client.** We see how the client drift of SCAFFOLD depends on  $\delta$ .

**Lemma 22.** *For the update (28), assuming (A2) and that  $\{f_i\}$  are quadratics, the following holds for any  $\eta \leq \frac{1}{21\delta K}$*

$$\mathbb{E}_{r-1,k-1} \|\mathbf{y}_{i,k}^r - \mathbf{x}^{r-1}\|^2 \leq (1 + \frac{1}{2K}) \|\mathbf{y}_{i,k-1}^r - \mathbf{x}^{r-1}\|^2 + 7K\eta^2 \|\nabla f(\mathbf{y}_{i,k-1}^r)\|^2 + \eta^2 \sigma^2.$$

*Proof.* Starting from the update step (28)

$$\begin{aligned} \mathbb{E}_{r-1,k-1} \|\mathbf{y}_i^+ - \mathbf{x}\|^2 &\leq \|\mathbf{y}_i - \mathbf{x} - \eta A(\mathbf{y}_i - \mathbf{x}^*) - \eta(A_i - A)(\mathbf{y}_i - \mathbf{x})\|^2 + \eta^2 \sigma^2 \\ &\leq (1 + \frac{1}{7(K-1)}) \|I - \eta(A_i - A)(\mathbf{y}_i - \mathbf{x})\|^2 + 7K\eta^2 \|A(\mathbf{y}_i - \mathbf{x}^*)\|^2 + \eta^2 \sigma^2. \end{aligned}$$

Note that if  $K = 1$ , then the first inequality directly proves the lemma. For the second inequality, we assumed  $K \geq 2$  and then applied our relaxed triangle inequality. By assumption A2, we have the following for  $\eta\delta \leq 1$

$$\|(I - \eta(A_i - A))^2\| = \|I - \eta(A_i - A)\|^2 \leq (1 + \eta\delta)^2 \leq 1 + 3\eta\delta.$$

Using the bound on the step-size  $\eta \leq \frac{1}{21\delta K}$  gives

$$\mathbb{E}_{r-1,k-1} \|\mathbf{y}_i^+ - \mathbf{x}\|^2 \leq (1 + \frac{1}{7K})(1 + \frac{1}{7(K-1)}) \|\mathbf{y}_i - \mathbf{x}\|^2 + 7K\eta^2 \|A(\mathbf{y}_i - \mathbf{x}^*)\|^2 + \eta^2 \sigma^2$$

Simple computations now give the Lemma statement for all  $K \geq 1$ .  $\square$

**Tracking the variance.** We will see how to bound the variance of the output.

**Lemma 23.** *Consider the update (28) for quadratic  $\{f_i\}$  with  $\eta \leq \max(\frac{1}{2\delta K}, \frac{1}{\beta})$ . Then, if further (A2), (A5) and (A4) are satisfied, we have*

$$\mathbb{E}_{r-1} f(\mathbf{x}^r) \leq f(\mathbb{E}_{r-1}[\mathbf{x}^r]) + 3K\beta \frac{\sigma^2}{N}.$$

Further if  $\{f_i\}$  are strongly convex satisfying (A3), we have

$$\mathbb{E}_{r-1} f(\mathbf{x}^r) \leq f(\mathbb{E}_{r-1}[\mathbf{x}^r]) + \beta \frac{\sigma^2}{N} \sum_{k=1}^K (1 - \mu\eta)^{k-1}.$$

*Proof.* We can rewrite the update step (28) as below:

$$\mathbf{y}_{i,k} = \mathbf{y}_{i,k-1} - \eta(A_i(\mathbf{y}_{i,k-1} - \mathbf{x}^*) + (A - A_i)(\mathbf{x} - \mathbf{x}^*)) - \eta\zeta_{i,k},$$

where by the bounded variance assumption A4,  $\zeta_{i,k}$  is a random variable satisfying  $\mathbb{E}_{k-1,r-1}[\zeta_{i,k}] = 0$  and  $\mathbb{E}_{k-1,r-1} \|\zeta_{i,k}\|^2 \leq \sigma^2$ . Subtracting  $\mathbf{x}^*$  from both sides and unrolling the recursion gives

$$\begin{aligned} \mathbf{y}_{i,K} - \mathbf{x}^* &= (I - \eta A_i)(\mathbf{y}_{i,K-1} - \mathbf{x}^*) - \eta((A - A_i)(\mathbf{x} - \mathbf{x}^*) + \zeta_{i,K}) \\ &= (I - \eta A_i)^K (\mathbf{x} - \mathbf{x}^*) - \sum_{k=1}^K \eta(I - \eta A_i)^{k-1} (\zeta_{i,k} + (A - A_i)(\mathbf{x} - \mathbf{x}^*)). \end{aligned}$$

Similarly, the expected iterate satisfies the same equation without the  $\zeta_{i,k}$ ,

$$\mathbb{E}_{r-1} [\mathbf{y}_{i,K}] - \mathbf{x}^* = (I - \eta A_i)^K (\mathbf{x} - \mathbf{x}^*) - \sum_{k=1}^K \eta(I - \eta A_i)^{k-1} (A - A_i)(\mathbf{x} - \mathbf{x}^*).$$

This implies that the difference satisfies

$$\mathbb{E}_{r-1} [\mathbf{y}_{i,K}] - \mathbf{y}_{i,K} = \eta \sum_{k=1}^K (I - \eta A_i)^{k-1} \zeta_{i,k}.$$

We can relate this to the function value as follows:

$$\begin{aligned}
 \mathbb{E}_{r-1} \|\mathbf{x}_K^r - \mathbf{x}^*\|_A^2 &= \|\mathbb{E}_{r-1}[\mathbf{x}_k^r] - \mathbf{x}^*\|_A^2 + \mathbb{E}_{r-1} \|\mathbb{E}_{r-1}[\mathbf{x}_k^r] - \mathbf{x}_K^r\|_A^2 \\
 &= \|\mathbb{E}_{r-1}[\mathbf{x}_k^r] - \mathbf{x}^*\|_A^2 + \mathbb{E}_{r-1} \|\frac{1}{N} \sum_i (\mathbb{E}_{r-1}[\mathbf{y}_{i,K}] - \mathbf{y}_{i,K})\|_A^2 \\
 &= \|\mathbb{E}_{r-1}[\mathbf{x}_k^r] - \mathbf{x}^*\|_A^2 + \eta^2 \mathbb{E}_{r-1} \|\frac{1}{N} \sum_{i,k} (I - \eta A_i)^{k-1} \zeta_{i,k}\|_A^2 \\
 &= \|\mathbb{E}_{r-1}[\mathbf{x}_k^r] - \mathbf{x}^*\|_A^2 + \frac{\eta^2}{N^2} \mathbb{E}_{r-1} \sum_{i,k} \|(I - \eta A_i)^{k-1} \zeta_{i,k}\|_A^2 \\
 &\leq \|\mathbb{E}_{r-1}[\mathbf{x}_k^r] - \mathbf{x}^*\|_A^2 + \frac{\beta \eta^2}{N^2} \mathbb{E}_{r-1} \sum_{i,k} \|(I - \eta A_i)^{k-1} \zeta_{i,k}\|_2^2.
 \end{aligned}$$

The last inequality used smoothness of  $f$  and the one before that relied on the independence of  $\zeta_{i,k}$ . Now, if  $f_i$  is general convex we have for  $\eta \leq \frac{1}{2\delta K}$  that  $I - \eta A_i \preceq (1 + \frac{1}{2K})I$  and hence

$$\|(I - \eta A_i)^{k-1} \zeta_{i,k}\|_2^2 \leq \sigma^2 (1 + \frac{1}{2K})^{2(k-1)} \leq 3\sigma^2.$$

This proves our second statement of the lemma. For strongly convex functions, we have for  $\eta \leq \frac{1}{\beta}$ ,

$$\|(I - \eta A_i)^{k-1} \zeta_{i,k}\|_2^2 \leq \sigma^2 (1 - \eta \mu)^{2(k-1)} \leq \sigma^2 (1 - \eta \mu)^{k-1}.$$

□

### F.3. Lemmas showing progress

**Progress of one client in one step.** Now we focus only on a single client and monitor their progress.

**Lemma 24.** Suppose (A2), (A5) and (A4) hold, and  $\{f_i\}$  are quadratics. Then, the following holds for the update (28) with  $\eta \leq \min(\frac{1}{10\beta}, \frac{1}{22\delta K}, \frac{1}{\mu K})$  with  $\mu = 0$  is  $f$  is non-convex or general-convex

$$\begin{aligned}
 \xi_{i,k}^r &\leq (1 - \frac{\mu\eta}{6}) \xi_{i,k-1}^r - \frac{\eta}{6} \mathbb{E}_{r-1} \|\nabla f(\mathbf{y}_{i,k-1}^r)\|^2 + 7\beta\eta^2\sigma^2, \text{ and} \\
 \tilde{\xi}_{i,k}^r &\leq (1 - \frac{\mu\eta}{6}) \tilde{\xi}_{i,k-1}^r - \frac{\eta}{6} \|\nabla f(\mathbb{E}_{r-1}[\mathbf{y}_{i,k-1}^r])\|^2.
 \end{aligned}$$

*Proof.* Recall that  $\xi_{i,k}^r \geq 0$  is defined to be

$$\xi_{i,k}^r := (\mathbb{E}_{r-1}[f(\mathbf{y}_{i,k}^r)] - f(\mathbf{x}^*) + \delta(1 + \frac{1}{K})^{K-k} \mathbb{E}_{r-1} \|\mathbf{y}_{i,k}^r - \mathbf{x}^{r-1}\|^2).$$

Let us start from the local update step (28) (dropping unnecessary subscripts and superscripts)

$$\begin{aligned}
 \mathbb{E}_{r-1,k-1} \|\mathbf{y}_i^+ - \mathbf{x}^*\|_A^2 &\leq \|\mathbf{y}_i - \mathbf{x}^*\|_A^2 - 2\eta \langle A(\mathbf{y}_i - \mathbf{x}^*), A(\mathbf{y}_i - \mathbf{x}^*) \rangle + 2\eta \langle (A - A_i)(\mathbf{y}_i - \mathbf{x}), A(\mathbf{y}_i - \mathbf{x}^*) \rangle \\
 &\quad + \eta^2 \|A(\mathbf{y}_i - \mathbf{x}^*) + (A_i - A)(\mathbf{y}_i - \mathbf{x})\|_A^2 + \beta\eta^2\sigma^2 \\
 &\leq \|\mathbf{y}_i - \mathbf{x}^*\|_A^2 - \frac{3\eta}{2} \|A(\mathbf{y}_i - \mathbf{x}^*)\|_2^2 + 2\eta \|(A - A_i)(\mathbf{y}_i - \mathbf{x})\|_2^2 \\
 &\quad + 2\eta^2 \|A(\mathbf{y}_i - \mathbf{x}^*)\|_A^2 + 2\eta^2 \|(A_i - A)(\mathbf{y}_i - \mathbf{x})\|_A^2 + \beta\eta^2\sigma^2 \\
 &\leq \|\mathbf{y}_i - \mathbf{x}^*\|_A^2 - (\frac{3\eta}{2} - 2\eta^2\beta) \|A(\mathbf{y}_i - \mathbf{x}^*)\|_2^2 + \beta\eta^2\sigma^2 + \delta^2(2\eta^2\beta + 2\eta) \|\mathbf{y}_i - \mathbf{x}\|_2^2 \\
 &\leq \|\mathbf{y}_i - \mathbf{x}^*\|_A^2 - (\frac{3\eta}{2} - 2\eta^2\beta) \|A(\mathbf{y}_i - \mathbf{x}^*)\|_2^2 + \beta\eta^2\sigma^2 + \frac{\delta}{10K} \|\mathbf{y}_i - \mathbf{x}\|_2^2.
 \end{aligned}$$

The second to last inequality used that  $\|\cdot\|_A^2 \leq \beta \|\cdot\|_2^2$  by (A5) and that  $\|(A - A_i)(\cdot)\|_2^2 \leq \delta^2 \|\cdot\|_2^2$  by (A2). The final inequality used that  $\eta \leq \max(\frac{1}{10\beta}, \frac{1}{22\delta K})$ . Now, multiplying Lemma 22 by  $\delta(1 + \frac{1}{K})^{K-k} \leq \frac{20\delta}{7}$  we have

$$\begin{aligned}
 \delta(1 + \frac{1}{K})^{K-k} \mathbb{E}_{r-1,k-1} \|\mathbf{y}_i^+ - \mathbf{x}\|^2 &\leq \delta(1 + \frac{1}{K})^{K-k} (1 + \frac{1}{2K}) \|\mathbf{y}_i - \mathbf{x}\|^2 + 20\delta K \eta^2 \|A(\mathbf{y}_i - \mathbf{x}^*)\|^2 + 3\delta\eta^2\sigma^2 \\
 &\leq \delta(1 + \frac{1}{K})^{K-k} (1 + \frac{1}{2K} + \frac{1}{10K}) \|\mathbf{y}_i - \mathbf{x}\|^2 - \frac{\delta}{10K} \|\mathbf{y}_i - \mathbf{x}\|^2 \\
 &\quad + 20\delta K \eta^2 \|A(\mathbf{y}_i - \mathbf{x}^*)\|^2 + 3\delta\eta^2\sigma^2 \\
 &\leq (1 - \frac{1}{5K}) \delta(1 + \frac{1}{K})^{K-k+1} (1 + \frac{1}{K}) \|\mathbf{y}_i - \mathbf{x}\|^2 - \frac{\delta}{10K} \|\mathbf{y}_i - \mathbf{x}\|^2 \\
 &\quad + 20\delta K \eta^2 \|A(\mathbf{y}_i - \mathbf{x}^*)\|^2 + 3\delta\eta^2\sigma^2.
 \end{aligned}$$

Adding this to our previous equation gives the following recursive bound:

$$\begin{aligned} & (\mathbb{E}_{r-1,k-1}\|\mathbf{y}_i^+ - \mathbf{x}^*\|_A^2 + \delta(1 + \frac{1}{K})^{K-k} \mathbb{E}_{r-1,k-1}\|\mathbf{y}_i^+ - \mathbf{x}\|^2) \leq \\ & (\|\mathbf{y}_i - \mathbf{x}^*\|_A^2 + (1 - \frac{1}{5K})\delta(1 + \frac{1}{K})^{K-k+1}\|\mathbf{y}_i - \mathbf{x}\|^2) - (\frac{3\eta}{2} - 2\eta^2\beta - 20\delta K\eta^2)\|A(\mathbf{y}_i - \mathbf{x}^*)\|_2^2 + (3\delta + \beta)\eta^2\sigma^2 \end{aligned}$$

The bound on our step-size  $\eta \leq \min(\frac{1}{10\beta}, \frac{1}{22\delta K})$  implies that  $\frac{3\eta}{2} - 2\eta^2\beta - 20\delta K\eta^2 \geq \frac{\eta}{3}$  and recall that  $\delta \leq 2\beta$ . This proves first statement of the lemma for non-strongly convex functions ( $\mu = 0$ ). If additionally  $f$  is strongly-convex with  $\mu > 0$ , we have

$$\eta\|A(\mathbf{y}_i - \mathbf{x}^*)\|_2^2 \geq \frac{\mu\eta}{2}\|\mathbf{y}_i - \mathbf{x}^*\|_A^2 + \frac{\eta}{2}\|A(\mathbf{y}_i - \mathbf{x}^*)\|_2^2.$$

This can be used to tighten the inequality as follows

$$\begin{aligned} & \left(\mathbb{E}_{r-1,k-1}\|\mathbf{y}_i^+ - \mathbf{x}^*\|_A^2 + \delta(1 + \frac{1}{K})^{K-(k-1)} \mathbb{E}_{r-1,k-1}\|\mathbf{y}_i^+ - \mathbf{x}\|^2\right) \leq \\ & ((1 - \frac{\mu\eta}{6})\|\mathbf{y}_i - \mathbf{x}^*\|_A^2 + (1 - \frac{1}{5K})\delta(1 + \frac{1}{K})^{K-k+1}\|\mathbf{y}_i - \mathbf{x}\|^2) - \frac{\eta}{2}\|A(\mathbf{y}_i - \mathbf{x}^*)\|_2^2 + 7\beta\eta^2\sigma^2 \end{aligned}$$

If  $\eta \leq \frac{1}{\mu K}$ , then  $(1 - \frac{1}{5K}) \leq (1 - \frac{\mu\eta}{6})$  and we have the strongly-convex version of the first statement.

Now for the second statement, recall that  $\tilde{\xi}_{i,k}^r \geq 0$  was defined to be

$$\tilde{\xi}_{i,k}^r := ([f(\mathbb{E}_{r-1}[\mathbf{y}_{i,k}^r])] - f(\mathbf{x}^*) + \delta(1 + \frac{1}{K})^{K-k} \mathbb{E}_{r-1}\|\mathbb{E}_{r-1}[\mathbf{y}_{i,k}^r] - \mathbf{x}^{r-1}\|^2).$$

Observe that for quadratics,  $\mathbb{E}_{r-1}[\nabla f(\mathbf{x})] = \nabla f(\mathbb{E}_{r-1}[\mathbf{x}])$ . This implies that the analysis of  $\tilde{\xi}_{i,k}^r$  is essentially of a deterministic process with  $\sigma = 0$ , proving the second statement. It is also straightforward to repeat exactly the above argument to formally verify the second statement.  $\square$

**Server progress in one round.** Now we combine the progress made by each client in one step to calculate the server progress.

**Lemma 25.** Suppose (A2), (A5) and (A4) hold, and  $\{f_i\}$  are quadratics. Then, the following holds for the update (28) with  $\eta \leq \min(\frac{1}{10\beta}, \frac{1}{21\delta K}, \frac{1}{10\mu K})$  and weights  $w_k := (1 - \frac{\mu\eta}{6})^{1-k}$ :

$$\frac{\eta}{6} \sum_{k=1}^K w_k \mathbb{E}_{r-1}\|\nabla f(\mathbf{x}_k^r)\|^2 \leq (f(\mathbb{E}_{r-2}[\mathbf{x}^{r-1}]) - f^*) - w_K(f(\mathbb{E}_{r-1}[\mathbf{x}^r]) - f^*) + \sum_{k=1}^K w_k 8\eta \frac{\sigma^2}{N}.$$

Set  $\mu = 0$  if  $\{f_i\}$ s are not strongly-convex (is only general-convex).

*Proof.* Let us do the non-convex (and general convex) case first. By summing over Lemma 24 we have

$$\frac{\eta}{6} \sum_{k=1}^K \mathbb{E}_{r-1}\|\nabla f(\mathbf{y}_{i,k}^r)\|^2 \leq \xi_{i,0}^r - \xi_{i,K}^r + 7K\beta\eta^2\sigma^2.$$

A similar result holds with  $\sigma = 0$  for  $\mathbb{E}_{r-1}[\mathbf{y}_{i,k}^r]$ . Now, using Lemma 20 we have that

$$\frac{\eta}{6} \sum_{k=1}^K \mathbb{E}_{r-1}\|\nabla f(\mathbf{x}_k^r)\|^2 \leq \underbrace{\frac{1}{N} \sum_{i=1}^N (\tilde{\xi}_{i,0}^r + \frac{1}{N}\xi_{i,0})}_{=: \theta_+^r} - \underbrace{\frac{1}{N} \sum_{i=1}^N (\tilde{\xi}_{i,K}^r + \frac{1}{N}\xi_{i,K})}_{=: \theta_-^r} + 7K\beta\eta^2 \frac{\sigma^2}{N}.$$

Using Lemma 23, we have that

$$\theta_+^r = (1 + \frac{1}{N})(f(\mathbf{x}^{r-1}) - f(\mathbf{x}^*)) \leq f(\mathbb{E}_{r-1}[\mathbf{x}^r]) + \frac{1}{N} \mathbb{E} f(\mathbf{x}^r) - (1 + \frac{1}{N})f(\mathbf{x}^*) + 3K\beta \frac{\sigma^2}{N}.$$

Further, by Lemma 21, we have that

$$\theta_-^r \geq f(\mathbb{E}_{r-1}[\mathbf{x}^r]) + \frac{1}{N} f(\mathbf{x}^r) - (1 + \frac{1}{N})f(\mathbf{x}^*).$$

Combining the above gives:

$$\frac{\eta}{6} \sum_{k=1}^K \mathbb{E}_{r-1} \|\nabla f(\mathbf{x}_k^r)\|^2 \leq f(\mathbb{E}_{r-2}[\mathbf{x}^{r-1}]) - f(\mathbb{E}_{r-1}[\mathbf{x}^r]) + 10\beta K \frac{\sigma^2}{N}.$$

proving the second part of the Lemma for weights  $w_k = 1$ . The proof of strongly convex follows a very similar argument. Unrolling Lemma 24 using weights  $w_k := (1 - \frac{\mu\eta}{6})^{1-k}$  gives

$$\frac{\eta}{6} \sum_{k=1}^K w_k \mathbb{E}_{r-1} \|\nabla f(\mathbf{x}_k^r)\|^2 \leq \theta_+^r - w_K \theta_-^r + \sum_{k=1}^K w_k 7\eta \frac{\sigma^2}{N}.$$

As in the general-convex case, we can use Lemmas 21, 20 and 23 to prove that

$$\frac{\eta}{6} \sum_{k=1}^K w_k \mathbb{E}_{r-1} \|\nabla f(\mathbf{x}_k^r)\|^2 \leq (f(\mathbb{E}_{r-2}[\mathbf{x}^{r-1}]) - f^*) - w_K (f(\mathbb{E}_{r-1}[\mathbf{x}^r]) - f^*) + \sum_{k=1}^K w_k 8\eta \frac{\sigma^2}{N}.$$

□

**Deriving final rates.** The proof of Theorem VIII follows by appropriately unrolling Lemma 25. For general-convex functions, we can simply use Lemma 2 with the probabilities set as  $p_k^r = \frac{1}{KR}$ . For strongly-convex functions, we use  $p_k^r \propto (1 - \frac{\mu\eta}{6})^{1-rk}$  and follow the computations in Lemma 1.