Byzantine-Resilient High-Dimensional SGD with Local Iterations on Heterogeneous Data

Deepesh Data 1 Suhas Diggavi 1

Abstract

We study stochastic gradient descent (SGD) with local iterations in the presence of Byzantine clients, motivated by the federated learning. The clients, instead of communicating with the server in every iteration, maintain their local models, which they update by taking several SGD iterations based on their own datasets and then communicate the net update with the server, thereby achieving communication-efficiency. Furthermore, only a subset of clients communicates with the server at synchronization times. The Byzantine clients may collude and send arbitrary vectors to the server to disrupt the learning process. To combat the adversary, we employ an efficient highdimensional robust mean estimation algorithm at the server to filter-out corrupt vectors; and to analyze the outlier-filtering procedure, we develop a novel matrix concentration result that may be of independent interest. We provide convergence analyses for both strongly-convex and non-convex smooth objectives in the heterogeneous data setting. We believe that ours is the first Byzantineresilient local SGD algorithm and analysis with non-trivial guarantees. We corroborate our theoretical results with preliminary experiments for neural network training.

1. Introduction

In the *federated learning* (FL) paradigm (Konecný, 2017; Konecný et al., 2016; McMahan et al., 2017; Mohri et al., 2019), several clients (e.g., mobiles devices, organizations, etc.) collaboratively learn a machine learning model, where the training process is facilitated by the data held by the participating clients (without data centralization) and is coordinated by a central server (e.g., the service provider). Due to its many advantages over the traditional centralized learning

Proceedings of the 38^{th} International Conference on Machine Learning, PMLR 139, 2021. Copyright 2021 by the author(s).

(Dean et al., 2012) (e.g., training a machine learning model without collecting the clients' data, which, in addition to reducing the communication load on the network, provides a basic level of privacy to clients' data), FL has emerged as an active area of research recently; see (Kairouz et al., 2019) for a detailed survey. Stochastic gradient descent (SGD) has become a de facto standard in optimization for training machine learning models at such a large scale (Bottou, 2010; Kairouz et al., 2019; McMahan et al., 2017), where clients iteratively communicate the gradient updates with the central server, which aggregates the gradients, updates the learning model, and sends the aggregated gradient back to the clients. The promise of FL comes with its own set of challenges (Kairouz et al., 2019): (i) optimizing with heterogeneous data at different clients – the local datasets at clients may be "non-i.i.d.", i.e., can be thought of as being generated from different underlying distributions; (ii) slow and unreliable network connections between server and clients. so communication in every iteration may not be feasible; (iii) availability of only a subset of clients for training at a given time (maybe due to low connectivity, as clients may be in different geographic locations); and (iv) robustness against malicious/Byzantine clients who may send incorrect gradient updates to the server to disrupt the training process. In this paper, we propose and analyze an SGD algorithm that simultaneously addresses all these challenges. First we setup the problem, put our work in context with the related work, and then summarize our contributions.

We consider an empirical risk minimization problem, where data is stored at R clients, each having a different dataset (with no probabilistic assumption on data generation); client $r \in [R]$ has dataset \mathcal{D}_r . Let $F_r : \mathbb{R}^d \to \mathbb{R}$ denote the local loss function associated with the dataset \mathcal{D}_r , which is defined as $F_r(\boldsymbol{x}) \triangleq \mathbb{E}_{i \in U[n_r]}[F_{r,i}(\boldsymbol{x})]$, where $n_r = |\mathcal{D}_r|$, i is uniformly distributed over $[n_r] \triangleq \{1, 2, \dots, n_r\}$, and $F_{r,i}(\boldsymbol{x})$ is the loss associated with the i'th data point at client r with respect to (w.r.t.) \boldsymbol{x} . Our goal is to solve the following minimization problem:

$$\arg\min_{\boldsymbol{x}\in\mathcal{C}} \left(F(\boldsymbol{x}) \triangleq \frac{1}{R} \sum_{r=1}^{R} \mathbb{E}_{i\in_{U}[n_r]}[F_{r,i}(\boldsymbol{x})] \right), \quad (1)$$

where $\mathcal{C}\subseteq\mathbb{R}^d$ denotes the parameter space that is either equal to \mathbb{R}^d or a compact and convex set.

¹University of California, Los Angeles, USA. Correspondence to: Deepesh Data <deepesh.data@gmail.com>.

In the absence of the above-mentioned FL challenges, we can minimize (1) using distributed *vanilla* SGD, where in any iteration, server broadcasts the current model parameters to all clients, each of them then samples a stochastic gradient from its local dataset and sends it back to the server, who aggregates the received gradients and updates the global model. However, this simple solution does not satisfy the FL challenges, as *every* client communicates with the server (i.e., no sampling of clients) in *every* SGD iteration (i.e., no local iterations), and furthermore, this solution breaks down even with a single malicious client (Blanchard et al., 2017).

Related work. Recent work have proposed variants of the above-described vanilla SGD that address *some* of the FL challenges. The algorithms in (Basu et al., 2019; Haddadpour & Mahdavi, 2019; Haddadpour et al., 2019; Karimireddy et al., 2020; Khaled et al., 2020; Li et al., 2020; Sahu et al., 2020; Yu et al., 2019b) work under different heterogeneity assumptions but do not provide any robustness to malicious clients. On the other hand, (Alistarh et al., 2018; Blanchard et al., 2017; Chen et al., 2017; Data & Diggavi, 2020b; Su & Xu, 2019; Xie et al., 2019b; Yin et al., 2018; 2019) provide robustness, but with no local iterations or sampling of clients; furthermore, they assume homogeneous (either same or i.i.d.) data across all clients. A different line of work (Chen et al., 2018; Data & Diggavi, 2019; 2020a; Data et al., 2019; 2021; Ghosh et al., 2019; Li et al., 2019a; Rajput et al., 2019) provide robustness with heterogeneous data, but without local iterations or sampling of clients: Chen et al. (2018), Rajput et al. (2019), Data et al. (2019; 2021) use coding across datasets, which is hard to implement in FL; Li et al. (2019a) change the objective function and adds a regularizer term to combat the adversary; Ghosh et al. (2019) effectively reduce the heterogeneous problem to a homogeneous problem by clustering, and then learning happens within each cluster having homogeneous data; and Data & Diggavi (2020a) studied SGD with heterogeneous data under the same assumptions as ours, but without local iterations or client sampling. Incorporating local iterations with Byzantine adversaries makes it significantly more challenging as it requires deriving a new matrix concentration bound (see Theorem 2) and different convergence analyses.

Xie et al. (2019a) also analyzed SGD in the FL setting, but the approximation error (even in the Byzantine-free setting) of their solution could be as large as $\mathcal{O}(D^2+G^2)$, where G is the gradient bound and D is the diameter of the parameter space that contains the optimal parameters \boldsymbol{x}^* and all the local parameters \boldsymbol{x}^t ever emerged at any client $r \in [R]$ in any iteration $t \in [T]$; this, in our opinion, makes their bound vacuous. In optimization, one would ideally like to have convergence rates depend on D with a factor that decays with the number of iterations, e.g., with $\frac{1}{T}$ or $\frac{1}{\sqrt{T}}$, as also in Theorem 1. In Section 4, we also empirically demonstrate the poor learning performance of their algorithm.

Our contributions. In this paper, we tackle heterogeneity assuming that the gradient dissimilarity among local datasets is bounded (see (6)), and propose and analyze a Byzantine-resilient SGD algorithm (Algorithm 1) with local iterations and client sampling under the bounded variance assumption for SGD (see (2)). We provide convergence analyses for strongly-convex and non-convex smooth objectives.

For strongly-convex objectives, our algorithm can find approximate optimal parameters exponentially (in $\frac{T}{H}$) fast, and for non-convex objectives, it can reach to an approximate stationary point with a speed of $\frac{1}{T/H}$. See Theorem 1 for convergence results. The approximation error in the optimization solution comprises of two terms, one is because to the stochasticity in gradients (due to SGD) and is equal to zero if we work with full-batch gradients, and the other term arises because of heterogeneity in local datasets. See a detailed discussion in Section 2.2 on the approximation error analysis and the convergence rates, and also for the reason behind obtaining rates that are off by a factor of H when compared to vanilla SGD – looking ahead, the reason is working with weak assumptions.

To tackle the malicious behavior of Byzantine clients, we borrow tools from recent advances in high-dimensional robust statistics (Diakonikolas & Kane, 2019; Diakonikolas et al., 2019; Lai et al., 2016; Steinhardt et al., 2018); in particular, we use the polynomial-time outlier-filtering procedure from (Diakonikolas et al., 2019), which was developed for robust mean estimation in high dimensions. In order to use their algorithm (described in Algorithm 2) in our setting that combines Byzantine resilience with local iterations, we develop a novel matrix concentration result (see Theorem 2), which may be of independent interest. As far as we know, this is the first concentration result for stochastic gradients with local iterations on heterogeneous data.

We believe that ours is the first work that combines *local iterations* with *Byzantine-resilience* for SGD and achieves non-trivial results. Not only that, we also analyze our algorithm on *heterogeneous* data and allow *sampling of clients*. Note that the earlier work that provide robustness (without local iterations or sampling of clients) either assume homogeneous data across clients (Alistarh et al., 2018; Blanchard et al., 2017; Chen et al., 2017; Data & Diggavi, 2020b; Su & Xu, 2019; Yin et al., 2018; 2019) or require strong assumptions, such as the bounded gradient assumption on local functions (Xie et al., 2019b); more on this on page 3.

Paper organization. We describe our algorithm and state the convergence results in Section 2. In Section 3, we describe our main technical tool, a new matrix concentration result for analyzing the robust accumulated gradient estimation procedure. We provide empirical evaluation of our method in Section 4. Omitted details/proofs are given in appendices, provided as part of the supplementary material.

2. Problem Setup and Our Results

In this section, we state our assumptions, describe the adversary model and our algorithm, and state our convergence results followed by important remarks about them.

Assumption 1 (Bounded local variances). *The stochastic* gradients sampled from any local dataset have uniformly bounded variance over C for all clients, i.e., there exists a finite σ , such that for all $x \in C$, $r \in [R]$, we have

$$\mathbb{E}_{i \in U[n_r]} \|\nabla F_{r,i}(\boldsymbol{x}) - \nabla F_r(\boldsymbol{x})\|^2 \le \sigma^2.$$
 (2)

It will be helpful to formally define mini-batch stochastic gradients, where instead of computing stochastic gradients based on just one data point, each client samples $b \geq 1$ data points (without replacement) from its local dataset and computes the average of b gradients. For any $x \in \mathbb{R}^d$, $r \in [R]$, $b \in [n_r]$, consider the following set

$$\mathcal{F}_r^{\otimes b}(\boldsymbol{x}) := \left\{ \frac{1}{b} \sum_{i \in \mathcal{H}_b} \nabla F_{r,i}(\boldsymbol{x}) : \mathcal{H}_b \in {\binom{[n_r]}{b}} \right\}. \quad (3)$$

Note that $g_r(x) \in_U \mathcal{F}_r^{\otimes b}(x)$ is a mini-batch stochastic gradient with batch size b at client r. It is not hard to see the following, which hold for all $x \in \mathcal{C}, r \in [R]$:

$$\mathbb{E}\left[\boldsymbol{g}_r(\boldsymbol{x})\right] = \nabla F_r(\boldsymbol{x}),\tag{4}$$

$$\mathbb{E} \|\boldsymbol{q}_r(\boldsymbol{x}) - \nabla F_r(\boldsymbol{x})\|^2 < \sigma^2/b. \tag{5}$$

Assumption 2 (Bounded gradient dissimilarity). The difference of the local gradients $\nabla F_r(\boldsymbol{x}), r \in [R]$ and the global gradient $\nabla F(\boldsymbol{x}) = \frac{1}{R} \sum_{r=1}^R \nabla F_r(\boldsymbol{x})$ is uniformly bounded over \mathbb{R}^d for all clients, i.e., there exists a finite κ , such that

$$\|\nabla F_r(\boldsymbol{x}) - \nabla F(\boldsymbol{x})\|^2 \le \kappa^2, \quad \forall \boldsymbol{x} \in \mathcal{C}, r \in [R].$$
 (6)

Assumption 1 has been standard in SGD literature. Assumption 2 has also been used earlier to bound heterogeneity in datasets; see, for example, (Li et al., 2019b; Yu et al., 2019a), which study decentralized SGD with momentum (without adversaries). Note that when clients compute full-batch gradients, we have $\sigma=0$ in Assumption 1; similarly, when all clients have access to the same dataset as in (Alistarh et al., 2018; Blanchard et al., 2017), we have $\kappa=0$ in Assumption 2. Note that (6) can be seen as a *deterministic* condition on local datasets, under which we derive our results.

A note on Assumption 2. In the presence of Byzantine adversaries, since we do not know which ϵR clients are corrupt, we have to make some structural assumption on the data that can provide relationships among gradients sampled at different nodes for reliable decoding, and Assumption 2 is a natural way to achieve that. There are many alternatives to establish this relationship, e.g., by assuming homogeneous (same or i.i.d.) data across clients (Alistarh et al.,

2018; Blanchard et al., 2017; Chen et al., 2017; Data & Diggavi, 2020b; Su & Xu, 2019; Yin et al., 2018; 2019) or by explicitly introducing redundancy in the system via coding-theoretic solutions (Chen et al., 2018; Data et al., 2021; Rajput et al., 2019); however, these approaches fall short of in the FL setting.

Assuming bounded gradients of local functions (i.e., $\|\nabla F_r(\boldsymbol{x})\| \leq G$ for some finite G) is a common assumption in literature with heterogeneous data; see, for example, (Li et al., 2020; Yu et al., 2019b, without adversaries) and (Xie et al., 2019b, with adversaries). Note that under this assumption, we can trivially bound the heterogeneity among local datasets by $\|\nabla F_r(\boldsymbol{x}) - \nabla F_s(\boldsymbol{x})\| \leq 2G$. So, assuming bounded gradients not only simplifies the analysis but also obscures the effect of heterogeneity on the convergence bounds, which Assumption 2 clearly brings out. ¹

Bounds on σ^2 and κ^2 in the statistical heterogeneous model. Since all our results (matrix concentration and convergence) are given in terms of σ and κ , to show the clear dependence of our results on the dimensionality of the problem, we bound these quantities in the statistical heterogeneous data model under different distributional assumptions on local gradients; see Appendix E for more details, where we prove the following: For the SGD variance bound, we show that if local gradients have sub-Gaussian distribution, then $\sigma = \mathcal{O}(\sqrt{d \log(d)})$. For the gradient dissimilarity bound, we show that if either the local gradients have sub-exponential distribution and each worker has at least $n = \Omega(d \log(nd))$ data points or local gradients have sub-Gaussian distribution and $n \in \mathbb{N}$ is arbitrary, then $\kappa \leq \kappa_{\text{mean}} + \mathcal{O}(\sqrt{d \log(nd)/n})$, where κ_{mean} denotes the distance of the expected local gradients from the global gradient. Note that we make distributional assumptions on data generation *only* to derive bounds on σ , κ ; otherwise, all our results hold for arbitrary datasets satisfying (5), (6).

Adversary model. Throughout the paper, we assume that ϵ denotes the fraction of the K communicating clients that are corrupt, i.e., at most ϵK (out of K) clients that communicate with the server at synchronization indices may be corrupt, where $K \leq R$ is the number of clients chosen at synchronization indices. This translates to, in the worst case, having $\frac{\epsilon K}{R}$ fraction (i.e., a total of ϵK) of corrupt nodes in the entire system, as in the worst-case, all the corrupt nodes can be selected in a communication round; however, in practice, due to several constraints, such as the unreliable network connection (for which the adversary has no control over), we cannot expect that the server will select all corrupt nodes in all iterations. The corrupt clients may collude and arbitrarily

¹See (Khaled et al., 2020) for a detailed discussion on the inappropriateness of making bounded gradient assumption in heterogeneous data settings and how it obscures the effect of heterogeneity on convergence rates (even without robustness).

Algorithm 1 Byzantine-Resilient SGD with Local Iterations

- 1: Initialize. Set t:=0, $\boldsymbol{x}_r^0:=\boldsymbol{0}, \forall r\in[R]$, and $\boldsymbol{x}:=\boldsymbol{0}$. Here, x denotes the global model and x_r^0 denotes the local model at client r at time 0. Fix a constant step-size η and a mini-batch size b.
- while $(t \leq T)$ do
- Server selects an arbitrary subset $\mathcal{K} \subseteq [R]$ of $|\mathcal{K}| =$ K clients and sends x to all clients in K.
- All clients $r \in \mathcal{K}$ do in parallel: 4:
- Set $\boldsymbol{x}_r^t = \boldsymbol{x}$. 5:
- 6: while (true) do
- 7: Take a mini-batch stochastic gradient $q_r(x_r^t) \in U$ $\mathcal{F}^{\otimes b}(\boldsymbol{x}_r^t)$ and update the local model:

$$\boldsymbol{x}_r^{t+1} \leftarrow \boldsymbol{x}_r^t - \eta \boldsymbol{g}_r(\boldsymbol{x}_r^t)); \quad t \leftarrow (t+1).$$

- if $(t \in \mathcal{I}_T)$ then 8:
- Let $\widetilde{\boldsymbol{x}}_r^t = \boldsymbol{x}_r^t$, if client r is honest, otherwise can 9: be an arbitrary vector in \mathbb{R}^d .
- Send $\widetilde{\boldsymbol{x}}_r^t$ to the server and break the inner while 10:
- 11: end if
- 12: end while
- At Server: 13:
- Receive $\{\tilde{\boldsymbol{x}}_r, r \in \mathcal{K}\}$ from the clients in \mathcal{K} . 14:
- 15:
- For every $r \in \mathcal{K}$, let $\widetilde{\boldsymbol{g}}_{r,\mathrm{accu}} := (\widetilde{\boldsymbol{x}}_r \boldsymbol{x})/\eta$. Apply the decoding algorithm RAGE (see Algo-16: rithm 2) on $\{\widetilde{\boldsymbol{g}}_{r,\text{accu}}, r \in \mathcal{K}\}$. Let

$$\widehat{\boldsymbol{g}}_{\text{accu}} := \text{RAGE}(\widetilde{\boldsymbol{g}}_{r \text{ accu}}, r \in \mathcal{K}).$$

- Update the global model $\boldsymbol{x} \leftarrow \Pi_{\mathcal{C}}(\boldsymbol{x} \eta \widehat{\boldsymbol{g}}_{\text{accu}})$, where 17: $\Pi_{\mathcal{C}}$ denotes the projection operator onto the set \mathcal{C} .
- 18: end while

deviate from their pre-specified programs: at synchronization indices, instead of sending the true stochastic gradients (or local models), corrupt clients may send adversarially chosen vectors to the server.

2.1. Main Results

Let $\mathcal{I}_T = \{t_1, t_2, \dots, t_k, \dots\}$, with $t_1 = 0$, denote the set of synchronization indices (where $\max_{i>1} |t_{i+1} - t_i| = H$) when the server arbitrarily selects a subset of $K \leq R$ clients (denoted by $\mathcal{K} \subseteq [R]$) and sends the global model (denoted by x) to them; each client $r \in \mathcal{K}$ updates its local model x_r by taking SGD steps based on its local dataset until the next synchronization time, when all clients in K send their local models to the server. Note that some of these clients may be corrupt and may send arbitrary vectors.² Server employs

a decoding RAGE and update the global model x based on that. We present our Byzantine-resilient SGD algorithm with local iterations in Algorithm 1.

Our convergence results are for both strongly-convex and non-convex smooth objectives, and we state them in the following theorem. Since our main focus in this paper is on combining Byzantine resilience with local iterations, to avoid the technical complications arising due to the projection operator (in line 17), we prove our results assuming that the parameter space \mathcal{C} is equal to \mathbb{R}^d . The analysis involving the projection can be done using the techniques in (Yin et al., 2018).

Theorem 1 (Mini-Batch Local Stochastic Gradient Descent). Let K_t denote the set of K clients that are active at any given time $t \in [0:T]$ and ϵ denote the fraction of corrupt clients in K_t . For a global objective function $F: \mathbb{R}^d \to \mathbb{R}$, let Algorithm 1 generate a sequence of iterates $\{x_r^t: t \in [0:T], r \in \mathcal{K}_t\}$ when running with a fixed step-size $\eta = \frac{1}{8HL}$. Fix any constant $\epsilon' > 0$. If $\epsilon \leq \frac{1}{3} - \epsilon'$, then with probability $1 - \frac{T}{H} \exp(-\frac{\epsilon'^2(1-\epsilon)K}{16})$, the sequence of average iterates $\{ \boldsymbol{x}^t = \frac{1}{K} \sum_{r \in \mathcal{K}_t} \boldsymbol{x}_r^t : t \in [0:T] \}$ satisfy the following convergence guarantees:

• Strongly-convex: If F is L-smooth for $L \geq 0$, and μ -strongly convex for $\mu > 0$, we get:

$$\mathbb{E} \left\| \boldsymbol{x}^T - \boldsymbol{x}^* \right\|^2 \leq \left(1 - \frac{\mu}{16HL}\right)^T \left\| \boldsymbol{x}^0 - \boldsymbol{x}^* \right\|^2 + \frac{13}{\mu^2} \Gamma.$$

• Non-convex: If F is L-smooth for $L \ge 0$, we get:

$$\frac{1}{T} \sum_{t=0}^{T} \mathbb{E} \left\| \nabla F(\boldsymbol{x}^t) \right\|^2 \leq \frac{\left[\mathbb{E}[F(\boldsymbol{x}^0)] - \mathbb{E}[F(\boldsymbol{x}^*)] \right]}{\frac{T}{16HL}} + \frac{9}{2} \Gamma.$$

In both the bounds above, $\Gamma=\left(\frac{3\Upsilon^2}{H}+\frac{11H\sigma^2}{b}+36H\kappa^2\right)$ with $\Upsilon^2=\mathcal{O}\left(\sigma_0^2(\epsilon+\epsilon')\right)$, where $\sigma_0^2=$ $\frac{25H^2\sigma^2}{\hbar\epsilon'}\left(1+\frac{3d}{2K}\right)+28H^2\kappa^2$, and expectation is taken over the sampling of mini-batch stochastic gradients.

We prove the strongly-convex part of Theorem 1 in Appendix B and the non-convex part in Appendix C. In addition to other complications arising due to handling Byzantine clients together with local iterations, our proof deviates from the standard proofs for local SGD: We need to show two recurrences, which arise because at synchronization indices, server performs decoding to filter-out the corrupt clients, while at other indices there is no decoding, as there is no communication. The proof of the first recurrence is significantly more involved than that of the other one.

Because of this and for the purpose of analysis, we can assume, without loss of generality, that in between the synchronization indices, the corrupt clients sample stochastic gradients and update their local parameters honestly.

²Note that the only disruption that the corrupt clients can cause in the training process is during the gradient aggregation at synchronization indices by sending adversarially chosen vectors to the server, and we give unlimited power to the adversary for that.

 $egin{aligned} & {}^3F(oldsymbol{y}) \leq F(oldsymbol{x}) + \langle
abla F(oldsymbol{x}), oldsymbol{y} - oldsymbol{x}
angle + rac{L}{2} \|oldsymbol{x} - oldsymbol{y}\|^2, orall oldsymbol{x}, oldsymbol{y} \in \mathbb{R}^d. \ & {}^4F(oldsymbol{y}) \geq F(oldsymbol{x}) + \langle
abla F(oldsymbol{x}), oldsymbol{y} - oldsymbol{x}
angle + rac{L}{2} \|oldsymbol{x} - oldsymbol{y}\|^2, orall oldsymbol{x}, oldsymbol{y} \in \mathbb{R}^d. \end{aligned}$

2.2. Important Remarks About Theorem 1

Failure probability. The failure probability of our algorithm is at most $\frac{T}{H} \exp(-\frac{\epsilon'^2(1-\epsilon)K}{16})$, which though scales linearly with T, also goes down exponentially with K. As a result, in settings such as federated learning, where number of clients could be large (e.g., in tens/hundreds of millions) and server samples tens of thousands of them, we can get a very small probability of error, even if run our algorithm for a long time. Note that the error probability is due to the *stochastic* sampling of gradients, and if we want a "zero" probability of error, we can run full-batch GD (yielding an error of $\Gamma = \mathcal{O}(H\kappa^2)$); we analyze that in Appendix D with a much simplified analysis than that of Theorem 1.

Analysis of the approximation error. In Theorem 1, the approximation error Γ essentially consists of two types of error terms: $\Gamma_1 = \mathcal{O}\left(\frac{H\sigma^2}{b\epsilon'}\left(1 + \frac{3d}{2K}\right)(\epsilon + \epsilon')\right)$ and $\Gamma_2 =$ $\mathcal{O}(H\kappa^2)$, where Γ_1 arises due to stochastic sampling of gradients and Γ_2 arises due to dissimilarity in the local datasets. Observe that Γ_1 decreases as we increase the batch size b of stochastic gradients and becomes zero if we take full-batch gradients (which implies $\sigma = 0$), as is the case in Theorem 4 in Appendix D. Note that even though the variance (and gradient dissimilarity) of accumulation of H gradients blows up by a factor of H^2 , still both Γ_1 and Γ_2 have a *linear* dependence on the number of local iterations H. Observe that since we are working with heterogeneous datasets, the presence of gradient dissimilarity bound κ^2 (which captures the heterogeneity) in the approximation error is inevitable, and will always show up when bounding the deviation of the true "global" gradient from the decoded one in the presence of Byzantine clients, even when H=1.

Convergence rates. In the strongly-convex case, Algorithm 1 approximately finds the optimal parameters \boldsymbol{x}^* (within Γ error) with $\left(1-\frac{\mu}{16HL}\right)^T$ speed. Note that $\left(1-\frac{\mu}{16HL}\right)^T \leq \exp^{-\frac{\mu}{16L}\frac{T}{H}}$, which implies an exponentially fast (in T/H) convergence rate. In the non-convex case, Algorithm 1 reaches to a stationary point (within Γ error) with a speed of $\frac{1}{T/H}$. Note that the convergence rates of vanilla SGD (i.e., without local iterations and in Byzantine-free settings) are exponential (in T) and $\frac{1}{T}$ for strongly-convex and non-convex objectives, respectively; whereas, our convergence rates are affected by the number of local iterations H. The reason for this is precisely because we

need $\eta \leq \frac{1}{8HL}$ to bound the drift in local parameters across clients; see Lemma 2. Instead, if we had assumed a stronger bounded gradient assumption (which trivially bound the heterogeneity, as explained on page 3), then Lemma 2 would hold for a constant step-size (e.g., $\eta = \frac{1}{2L}$ would suffice), which would lead to vanilla SGD like convergence rates.

3. Robust Accumulated Gradient Estimation

In this section, first we discuss the inadequacy of traditional methods (such as coordinate-wise median and trimmed-mean) for filtering corrupt gradients in our setting, and then we motivate and describe the robust accumulated gradient estimation (RAGE) procedure that we use in Algorithm 1 as a subroutine at every synchronization index. Then we prove our new matrix concentration result that is required to establish the performance guarantee of RAGE.

Inadequacy of median and trimmed-mean: Coordinatewise median (med) and trimmed-mean (trimmean) are the two widely used robust estimation procedures that are easy to describe and implement, and they have been employed earlier for robust gradient aggregation in distributed optimization; see, for example, (Yin et al., 2018; 2019, i.i.d. data setting) and (Xie et al., 2019a, FL setting). Below we argue that these methods give poor performance in FL settings for learning high-dimensional models; we also validate this claim through experiments in Section 4.

- For the simple task of robust mean estimation with inputs coming a unit covariance distribution, med and trimmean have an error that scales with the dimension as \sqrt{d} (Diakonikolas et al., 2019; Lai et al., 2016); when we apply these methods in each SGD iteration, this error translates to a large sub-optimality gap in the convergence rate.
- The adversary may corrupt samples in a way that they preserve the norm of the original uncorrupted samples, but have different adversarially chosen directions (these are called directional attacks); since the performance of these methods are based on the magnitude of the samples, they cannot distinguish between the corrupt and uncorrupt samples.
- When taking coordinate-wise median, for estimating each coordinate, we use only a *single* sample and discard the rest. This is not a good idea in large-scale settings with non-i.i.d. data, such as FL, where there are potentially millions of clients, and if we somehow are able to use samples from *all* (or most of the) honest clients, we could get a significant reduction in variance of stochastic gradients. In med, we do not take advantage of this variance reduction, which leads to a performance degradation, which may be detrimental for performance due to heterogeneity in data. The same reason also applies to the robust gradient aggregation method (KRUM) adopted in (Blanchard et al., 2017), which also uses only one of the input gradients and discards the rest, giving poor performance.

 $[\]overline{}^5$ As a concrete scenario, say the total number of devices is R=10 million and the server selects K=10,000 of them. Then, even if we want robustness against one million malicious clients, the total probability of failure of our algorithm would still be less than $\frac{T}{H}e^{-30}$, which even if $T=10^6$ and H=1, would still be less than 10^{-7} . Note that the bound on probability of error in Theorem 1 is a worst-case bound, and in practice, our algorithm succeeds with moderate parameter values; see, for example, Section 4 for our experimental setup and the results.

Robust mean estimation: The above limitations of traditional methods motivate us to employ modern tools from high-dimensional robust statistics (Diakonikolas & Kane, 2019; Diakonikolas et al., 2019; Lai et al., 2016). In particular, we use the polynomial-time outlier-filtering procedure for high-dimensional robust mean estimation (RME) from (Diakonikolas et al., 2019) for robust gradient aggregation in Algorithm 1. For clear exposition of the ideas behind their algorithm, we use a version of their algorithm as described in Algorithm 2, which is from (Li, 2019). The crucial observation in these RME algorithms is that if the empirical mean of the samples is far from their true mean, then the empirical covariance matrix has high largest eigenvalue. So, the idea is to iteratively filter out samples that have large projection on the principal eigenvector of the empirical covariance matrix, and keep on doing it until the largest eigenvalue of the empirical covariance matrix becomes sufficiently small (line 7). This is done via a soft-removal method, where we assign weights (confidence score) to the samples and down-weighting those that have large projection (line 10) - in each iteration t, at least one sample (whose projection $\tau_i^{(t)}$ is the maximum) gets 0 weight. In the end, take the weighted average of the surviving samples.⁶

The RME algorithms overcome most of the abovementioned limitations of traditional methods, except for that their guarantees are not directly applicable to our setting. This is because the error guarantee of RME algorithms are given in terms of concentration of the good samples around their sample mean, which is easy to bound if good samples come from the same distribution. Note that our setup significantly deviates from this, where not only the input samples (which are accumulated gradients) come from different distributions (as clients have heterogeneous data), but each of them is also a sum of H stochastic gradients (due to local iterations). Since local iterations cause local parameters to drift from each other, bounding the concentration of good samples requires bounding this drift.

To this end, we develop a novel matrix concentration inequality that first shows an existence of a large subset of uncorrupted accumulated stochastic gradients and then bounds their concentration around the sample mean; see (7) in Theorem 2 below. As far as we know, this is the first matrix concentration result in an FL setting.

First we setup the notation. Let Algorithm 1 generate a sequence of iterates $\{x_r^t: t \in [0:T], r \in \mathcal{K}_t\}$ when Algorithm 2 Robust Accumulated Gradient Estimation (RAGE) (Diakonikolas et al., 2019; Li, 2019)

- 1: **Input:** K vectors $\boldsymbol{g}_1, \boldsymbol{g}_2, \dots, \boldsymbol{g}_K \in \mathbb{R}^d$ such that there is a subset of them $\mathcal{S} \subset [K]$ with $|\mathcal{S}| \geq \frac{2K}{3}$ having bounded covariance $\lambda_{\max}\left(\frac{1}{|\mathcal{S}|}\sum_{i\in\mathcal{S}}\left(\boldsymbol{g}_{i}-\boldsymbol{g}_{\mathcal{S}}\right)\left(\boldsymbol{g}_{i}-\boldsymbol{g}_{\mathcal{S}}\right)^{T}\right)\leq\sigma_{0}^{2},$ where $oldsymbol{g}_{\mathcal{S}} = rac{1}{|\mathcal{S}|} \sum_{i \in \mathcal{S}} oldsymbol{g}_i.$
- 2: For any $\boldsymbol{w} \in [0,1]^K$ with $\|\boldsymbol{w}\|_1 > 0$, define

$$\begin{split} & \boldsymbol{\mu}(\boldsymbol{w}) = \sum_{i=1}^K \frac{w_i}{\|\boldsymbol{w}\|_1} \boldsymbol{g}_i \\ & \boldsymbol{\Sigma}(\boldsymbol{w}) = \sum_{i=1}^K \frac{w_i}{\|\boldsymbol{w}\|_1} (\boldsymbol{g}_i - \boldsymbol{\mu}(\boldsymbol{w})) (\boldsymbol{g}_i - \boldsymbol{\mu}(\boldsymbol{w}))^T \end{split}$$

- 3: Let $\boldsymbol{w}^{(0)} = [\frac{1}{K}, \dots, \frac{1}{K}]$ be a length K vector. 4: Let $C \geq 11$ be a universal constant. 5: Let $\boldsymbol{\Sigma}^{(0)} = \boldsymbol{\Sigma}(\boldsymbol{w}^{(0)})$.

- 6: Let t = 0.

- 7: while $\lambda_{\max}(\mathbf{\Sigma}(\boldsymbol{w}^{(t)})) > C\sigma_0^2$ do 8: Let $\boldsymbol{v}^{(t)}$ be the principal eigenvector of $\mathbf{\Sigma}(\boldsymbol{w}^{(t)})$. 9: For $i \in [K]$, define $\tau_i^{(t)} = \left\langle \boldsymbol{v}^{(t)}, \boldsymbol{g}_i \boldsymbol{\mu}(\boldsymbol{w}^{(t)}) \right\rangle^2$.
- For $i \in [K]$, compute $w_i^{(t+1)} = \left(1 \frac{\tau_i^{(t)}}{\tau_{\max}^{(t)}}\right) w_i^{(t)}$, 10: where $\tau_{\max}^{(t)} = \max_{i:w^{(t)} > 0} \tau_i^{(t)}$.
- 11:
- 12: end while
- 13: **return** $\hat{g} = \sum_{i=1}^{K} \frac{w_i^{(t)}}{\|\mathbf{w}^{(t)}\|_1} g_i$.

running with a fixed step-size $\eta \leq \frac{1}{8HL}$, where \mathcal{K}_t denotes the set of K clients that are active at time $t \in [0:T]$. Take any two consecutive synchronization indices $t_k, t_{k+1} \in \mathcal{I}_T$. Note that $|t_{k+1} - t_k| \le H$. For an honest client $r \in \mathcal{K}_{t_k}$, let $g_{r, \text{accu}}^{t_k, t_{k+1}} := \sum_{t=t_k}^{t_{k+1}-1} g_r(x_r^t)$ denote the sum of local mini-batch stochastic gradients sampled by client r between time t_k and t_{k+1} , where $\boldsymbol{g}_r(\boldsymbol{x}_r^t) \in_U \mathcal{F}_r^{\otimes b}(\boldsymbol{x}_r^t)$ satisfies (4), (5). At iteration t_{k+1} , every honest client $r \in \mathcal{K}_{t_k}$ reports its local model $x_r^{t_{k+1}}$ to the server, from which server computes $g_{r,\text{accu}}^{t_k,t_{k+1}}$ (see line 15 of Algorithm 1), whereas, the corrupt clients may report arbitrary and adversarially chosen vectors in \mathbb{R}^d . Server does not know the identities of the corrupt clients, and its goal is to produce an estimate $\widehat{m{g}}_{ ext{accu}}^{t_k,t_{k+1}}$ of the average accumulated gradients from honest clients.

Theorem 2 (Matrix concentration). Suppose an ϵ fraction of K clients that communicate with the server are corrupt. In the setting described above, suppose we are given $K \leq$ R accumulated gradients $\widetilde{g}_{r,\text{accu}}^{t_k,t_{k+1}}, r \in \mathcal{K}_{t_k}$ in \mathbb{R}^d , where $\widetilde{m{g}}_{r, ext{accu}}^{t_k,t_{k+1}} = m{g}_{r, ext{accu}}^{t_k,t_{k+1}}$ if r'th client is honest, otherwise can be arbitrary. For any $\epsilon' > 0$, if $(\epsilon + \epsilon') \leq \frac{1}{3}$, then with probability $1 - \exp(-\frac{\epsilon'^2(1-\epsilon)K}{16})$, there exists a subset $S \subseteq$

⁶Note that the outlier-filtering procedure described in Algorithm 2 is intuitive and easy to understand. There are better algorithms that are also more efficient and can achieve better guarantees; see, for example, (Dong et al., 2019). All these algorithms require the same bounded matrix concentration assumption that we show in Theorem 2, thus making them applicable to use as a subroutine in Algorithm 1 without requiring any modification in our analysis.

 \mathcal{K}_{t_k} of uncorrupted gradients of size $(1 - (\epsilon + \epsilon'))K$ s.t.

$$\lambda_{\max} \left(\frac{1}{|\mathcal{S}|} \sum_{i \in \mathcal{S}} (\boldsymbol{g}_i - \boldsymbol{g}_{\mathcal{S}}) (\boldsymbol{g}_i - \boldsymbol{g}_{\mathcal{S}})^T \right)$$

$$\leq \frac{25H^2 \sigma^2}{b\epsilon'} \left(1 + \frac{3d}{2K} \right) + 28H^2 \kappa^2, \quad (7)$$

where, for $i \in \mathcal{S}$, $\boldsymbol{g}_i = \boldsymbol{g}_{i,\text{accu}}^{t_k,t_{k+1}}, \boldsymbol{g}_{\mathcal{S}} = \frac{1}{|\mathcal{S}|} \sum_{i \in \mathcal{S}} \boldsymbol{g}_{i,\text{accu}}^{t_k,t_{k+1}},$ and λ_{\max} denotes the largest eigenvalue.

Theorem 2 establishes the concentration results required for the RME algorithm (described in Algorithm 2) that we employ in Algorithm 1. This RME algorithm takes a collection of vectors as input, out of which an unknown large subset (at least a $\frac{2}{3}$ -fraction) is promised to be well-concentrated around its sample mean, and outputs an estimate of the sample mean. The formal guarantee is given as follows:

Theorem 3 (Outlier-filtering algorithm (Diakonikolas et al., 2019)). Under the same setting and notation of Theorem 2, we can find an estimate $\hat{\mathbf{g}}$ of $\mathbf{g}_{\mathcal{S}}$ in polynomial-time with probability 1, such that $\|\hat{\mathbf{g}} - \mathbf{g}_{\mathcal{S}}\| \leq \mathcal{O}\left(\sigma_0\sqrt{\epsilon + \epsilon'}\right)$, where $\sigma_0^2 = \frac{25H^2\sigma^2}{b\epsilon'}\left(1 + \frac{3d}{2K}\right) + 28H^2\kappa^2$.

Note that, instead of the RME algorithm, if we use med or trimmean, we would get an extra multiplicative factor of \sqrt{d} in the upper-bound on $\|\widehat{g} - g_S\|$ above.

3.1. Proof-sketch of Theorem 2 – Matrix Concentration

In order to prove Theorem 2, we use the following result from (Data & Diggavi, 2020a, Lemma 1):

Lemma 1 ((Data & Diggavi, 2020a, Lemma 1)). Suppose there are m independent distributions p_1, p_2, \ldots, p_m in \mathbb{R}^d such that $\mathbb{E}_{\boldsymbol{y} \sim p_i}[\boldsymbol{y}] = \boldsymbol{\mu}_i, i \in [m]$ and each p_i has a bounded variance in all directions, i.e., $\mathbb{E}_{\boldsymbol{y} \sim p_i}[\langle \boldsymbol{y} - \boldsymbol{\mu}_i, \boldsymbol{v} \rangle^2] \leq \sigma_{p_i}^2, \forall \boldsymbol{v} \in \mathbb{R}^d, \|\boldsymbol{v}\| = 1$. Take any $\epsilon' > 0$. Then, given m independent samples $\boldsymbol{y}_1, \boldsymbol{y}_2, \ldots, \boldsymbol{y}_m$, where $\boldsymbol{y}_i \sim p_i$, with probability $1 - \exp(-\epsilon'^2 m/16)$, there is a subset S of $(1 - \epsilon')m$ points such that $\lambda_{\max}(\frac{1}{|S|}\sum_{i \in S}(\boldsymbol{y}_i - \boldsymbol{\mu}_i)(\boldsymbol{y}_i - \boldsymbol{\mu}_i)^T) \leq \frac{4\sigma_{\max}^2}{\epsilon'}(1 + \frac{d}{(1 - \epsilon')m})$, where $\sigma_{p_{\max}}^2 = \max_{i \in [m]} \sigma_{p_i}^2$.

Lemma 1 shows that if we have m independent distributions each having bounded variance, and we take one sample from each of them, then there exists a large subset of these samples that has bounded variance as well. The important thing to note here is that the m samples come from different distributions, which makes it distinct from existing results, such as (Charikar et al., 2017, Proposition B.1), which shows concentration of i.i.d. samples.

Now we give a proof-sketch of Theorem 2 with the help of Lemma 1. A complete proof is provided in Appendix A.

Let $t_k, t_{k+1} \in \mathcal{I}_T$ be any two consecutive synchronization indices. For $i \in \mathcal{K}_{t_k}$ corresponding to an honest client, let

 $Y_i^{t_k}, Y_i^{t_k+1}, \dots, Y_i^{t_{k+1}-1}$ be a sequence of $(t_{k+1}-t_k) \leq H$ (dependent) random variables, where for any $t \in [t_k:t_{k+1}-1]$, the random variable Y_i^t is distributed as

$$Y_i^t \sim \text{Unif}\Big(\mathcal{F}_i^{\otimes b}\big(\boldsymbol{x}_i^t\big(\boldsymbol{x}_i^{t_k}, Y_i^{t_k}, \dots, Y_i^{t-1}\big)\big)\Big).$$
 (8)

Here, Y_i^t corresponds to the mini-batch stochastic gradient sampled from the set $\mathcal{F}_i^{\otimes b} \big(\boldsymbol{x}_i^t \big(\boldsymbol{x}_i^{t_k}, Y_i^{t_k}, \dots, Y_i^{t-1} \big) \big)$, which itself depends on the local parameters $\boldsymbol{x}_i^{t_k}$ (which is a deterministic quantity) at the last synchronization index and the past realizations of $Y_i^{t_k}, \dots, Y_i^{t-1}$. This is because the evolution of local parameters \boldsymbol{x}_i^t depends on $\boldsymbol{x}_i^{t_k}$ and the choice of gradients in between time indices t_k and t-1. Now define $Y_i := \sum_{t=t_k}^{t_{k+1}-1} Y_i^t$. Let p_i be the distribution of Y_i , which we will take when using Lemma 1.

It is not hard to show that for any honest client $i \in \mathcal{K}_{t_k}$, we have $\mathbb{E}\|Y_i - \mathbb{E}[Y_i]\|^2 \leq \frac{H^2\sigma^2}{b}$. It is also easy to see that the hypothesis of Lemma 1 is satisfied with $\mu_i = \mathbb{E}[Y_i], \sigma_{p_i}^2 = \frac{H^2\sigma^2}{b}$ for all honest clients $i \in \mathcal{K}_{t_k}$, i.e., we have $\mathbb{E}_{\boldsymbol{y}_i \sim p_i}[\langle \boldsymbol{y}_i - \mathbb{E}[\boldsymbol{y}_i], \boldsymbol{v} \rangle^2] \leq \frac{H^2\sigma^2}{b}, \forall \boldsymbol{v} \in \mathbb{R}^d, \|\boldsymbol{v}\| = 1$.

We are given K different accumulated gradients (each is a summation of H gradients), out of which at least $(1-\epsilon)K$ are according to the correct distribution. By considering only the uncorrupted gradients (i.e., taking $m=(1-\epsilon)K$), we have from Lemma 1 that there exists a subset $\mathcal{S}\subseteq\mathcal{K}_{t_k}$ of size $(1-\epsilon')(1-\epsilon)K\geq (1-(\epsilon+\epsilon'))K\geq \frac{2K}{3}$ that satisfies (in the following, $\widetilde{\boldsymbol{y}}_i=\boldsymbol{y}_i-\mathbb{E}[\boldsymbol{y}_i]$)

$$\lambda_{\max} \left(\frac{1}{|\mathcal{S}|} \sum_{i \in \mathcal{S}} \widetilde{\boldsymbol{y}}_i \widetilde{\boldsymbol{y}}_i^T \right) \le \widehat{\sigma}_0^2 := \frac{4H^2 \sigma^2}{b\epsilon'} \left(1 + \frac{3d}{2K} \right). \tag{9}$$

Note that (9) bounds the deviation of the points in S from their respective means $\mathbb{E}[y_i]$. However, in (7), we need to bound the deviation of the points in S from their sample mean $\frac{1}{|S|} \sum_{i \in S} y_i$. As it turns out, due to heterogeneity in data and our use of local iterations, this extension is non-trivial and requires some technical work, given next.

From the alternate definition of the largest eigenvalue of symmetric matrices $\mathbf{A} \in \mathbb{R}^{d \times d}$, we have $\lambda_{\max}(\mathbf{A}) = \sup_{\boldsymbol{v} \in \mathbb{R}^{d}, \|\boldsymbol{v}\|=1} \boldsymbol{v}^T \mathbf{A} \boldsymbol{v}$. With this, (9) is equivalent to

$$\sup_{\boldsymbol{v} \in \mathbb{R}^d: ||\boldsymbol{v}|| = 1} \frac{1}{|\mathcal{S}|} \sum_{i \in \mathcal{S}} \langle \boldsymbol{y}_i - \mathbb{E}[\boldsymbol{y}_i], \boldsymbol{v} \rangle^2 \le \widehat{\sigma}_0^2.$$
 (10)

Define $y_S := \frac{1}{|S|} \sum_{i \in S} y_i$ to be the sample mean of points in S. Take an arbitrary unit vector $v \in \mathbb{R}^d$. Using some algebraic manipulations provided in Appendix A, we get

$$\frac{1}{|\mathcal{S}|} \sum_{i \in \mathcal{S}} \langle \boldsymbol{y}_i - \boldsymbol{y}_{\mathcal{S}}, \boldsymbol{v} \rangle^2 \le 6\widehat{\sigma}_0^2 + \frac{4}{|\mathcal{S}|} \sum_{i \in \mathcal{S}} \frac{1}{|\mathcal{S}|} \sum_{i \in \mathcal{S}} \left\| \mathbb{E}[\boldsymbol{y}_j] - \mathbb{E}[\boldsymbol{y}_i] \right\|^2 \tag{11}$$

Using the gradient dissimilarity bound and the L-smoothness of F, we can show that for honest clients $r,s \in \mathcal{K}_{t_k}$, we have $\|\mathbb{E}[\boldsymbol{y}_r] - \mathbb{E}[\boldsymbol{y}_s]\|^2 \leq H \sum_{t=t_k}^{t_{k+1}-1} \left(6\kappa^2 + 3L^2\mathbb{E}\|\boldsymbol{x}_r^t - \boldsymbol{x}_s^t\|^2\right)$. Using this bound in (11) together with some algebraic manipulations, we get

$$\frac{1}{|\mathcal{S}|} \sum_{i \in \mathcal{S}} \langle \boldsymbol{y}_i - \boldsymbol{y}_{\mathcal{S}}, \boldsymbol{v} \rangle^2 \le 6\widehat{\sigma}_0^2 + 24H^2 \kappa^2
+ \frac{12HL^2}{|\mathcal{S}|} \sum_{i \in \mathcal{S}} \frac{1}{|\mathcal{S}|} \sum_{j \in \mathcal{S}} \sum_{t=t_k}^{t_{k+1}-1} \mathbb{E} \|\boldsymbol{x}_r^t - \boldsymbol{x}_s^t\|^2 \tag{12}$$

Now we bound the last term of (12), which is the drift in local parameters at different clients in between any two synchronization indices.

Lemma 2. If
$$\eta \leq \frac{1}{8HL}$$
, we have
$$\sum_{t=t_k}^{t_{k+1}-1} \mathbb{E} \|\boldsymbol{x}_r^t - \boldsymbol{x}_s^t\|^2 \leq 7H^3\eta^2 \left(\frac{\sigma^2}{b} + 3\kappa^2\right).$$

Substituting this in (12) together with some algebraic manipulations provided in Appendix A, we get

$$\frac{1}{|\mathcal{S}|} \sum_{i \in \mathcal{S}} \langle \boldsymbol{y}_i - \boldsymbol{y}_{\mathcal{S}}, \boldsymbol{v} \rangle^2 \leq \frac{25H^2\sigma^2}{b\epsilon'} \Big(1 + \frac{3d}{2K} \Big) + 28H^2\kappa^2.$$

Note that this bound holds for all unit vectors $\boldsymbol{v} \in \mathbb{R}^d$. Now substituting $\boldsymbol{g}_{i,\mathrm{accu}}^{t_k,t_{k+1}} = \boldsymbol{y}_i, \boldsymbol{g}_{\mathcal{S},\mathrm{accu}}^{t_k,t_{k+1}} = \boldsymbol{y}_{\mathcal{S}}$ and using the alternate definition of largest eigenvalue proves Theorem 2.

4. Experiments

In this section, we present preliminary numerical results on a non-convex objective. Additional implementation details can be found in Appendix F in the supplementary material.

Setup: We train a single layer neural network for image classification on the MNIST handwritten digit (from 0-9) dataset. The hidden layer has 25 nodes with ReLU activation function and the output has softmax function. The dimension of the model parameter vector is $19,885.^7$ All clients compute stochastic gradients on a batch-size of 128 in each iteration and communicate the local parameter vectors with the server after taking H=7 local iterations. For all the defense mechanisms, we start with a step-size $\eta=0.08$ and decrease its learning rate by a factor of 0.96 when the difference in the corresponding test accuracies in the last 2 consecutive epochs is less than 0.001.

Heterogeneous datasets: The MNIST dataset has 60,000 training images (with 6000 images of each label) and 10,000 test images (each having $28 \times 28 = 784$ pixels),

and is distributed among the 200 clients in the following *heterogeneous* manner: Each client takes a random permutation of the probability vector [0.8, 0.1, 0.1, 0, 0, 0, 0, 0, 0, 0, 0]. Suppose it obtains a vector \boldsymbol{p} such that $p_i = 0.8, p_j = 0.1, p_k = 0.1$ for some distinct $i, j, k \in [0:9]$ and $p_l = 0$ for the rest of the indices, then it selects *uniformly at random* 800, 100, 100 training images with label i, j, k, respectively.

Adversarial attacks: We have 12.5\% adversarial clients, i.e., 25 out of 200 clients are corrupt, and the corrupt set of clients may change in every iteration. We implement six adversarial attacks: (i) the 'random gradient attack', where local gradients at clients are replaced by independent Gaussian random vectors having the same norm⁸ as the corresponding gradients; (ii) the 'reverse average gradient attack', where corrupt clients send -ve of their average local gradients; (iii) the 'gradient shift attack', where local gradients of corrupt clients are shifted by a scaled (by factor of 50) Gaussian random vector (same for all); (iv) the 'all ones attack', where gradients of the corrupt clients are replaced by the all ones vector; (v) the 'Baruch attack', which was designed in (Baruch et al., 2019) specifically for coordinate-wise trimmed mean (trimmean) (Yin et al., 2018), Krum (Blanchard et al., 2017), and Bulyan (Mhamdi et al., 2018) defenses; and (vi) the 'reverse scaled average gradient attack', where corrupt clients compute the -ve of their average local gradients, scale it by the factor of 50, and then send it.

Performance: We train our neural network under all the above-described adversarial attacks, and demonstrate in Figure 1 the performance of our method (red color) against four other methods for robust gradient aggregation, namely, *coordinate-wise trimmed-mean* (black color) and *coordinate-wise median* (green color), which were used in (Xie et al., 2019a; Yin et al., 2018; 2019), Krum (magenta color), which was proposed in (Blanchard et al., 2017), and Bulyan (cyan color), which was proposed in (Mhamdi et al., 2018). For reference, we also plot (in blue color) the performance of Algorithm 1 with the same setup as above but without adversaries and with no decoding. For each attack, we plot two curves, one for training loss vs. number of epochs and the other for test accuracy vs. number of epochs.

It can be seen from the comparison in Figure 1 that our method consistently outperforms all these methods in all the attacks that we have implemented.⁹ In particular, for attacks

 $^{^7784 \}times 25 = 19,600$ weights between the input and the first layer, 25 bias terms (one for each node in the hidden layer), $25 \times 10 = 250$ weights between the first layer and the output layer, and 10 bias terms (one for each node in the output layer).

⁸Note that changing the direction while keeping the norm same is among the worst attacks as the corrupt gradients cannot be filtered out just based on their norms.

⁹We found out that the Bulyan defense mechanism is significantly slower than all other mechanisms. Due to this, we only implemented this for the Baruch-attack, which was specifically designed against Krum/Bulyan algorithms. Since a basic building block of Bulyan is Krum, and Krum performs the worst among all the mechanisms that we implemented, we do not expect Bulyan

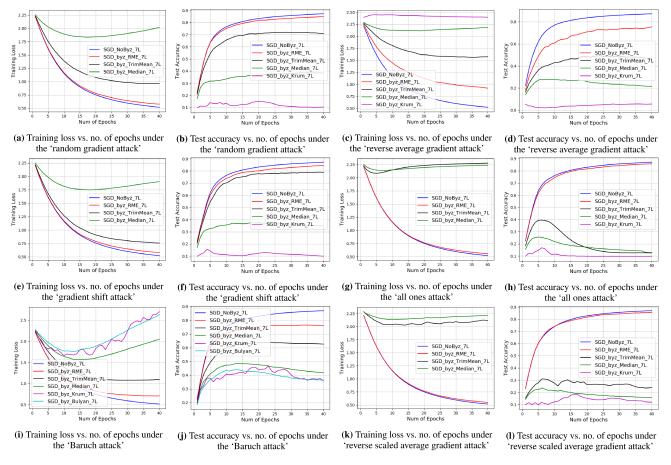


Figure 1. We compare the performance of our method (red) against four methods for robust gradient aggregation, namely, coordinate-wise trimmed-mean (black), coordinate-wise median (green), Krum (magenta), and Bulyan (cyan) under several adversarial attacks, and plot training loss and test accuracy against number of epochs. The plot in blue corresponds to running Algorithm 1 with no adversaries and no decoding. In the legends, 7L denotes that we are taking H = 7 local iterations. See also Footnotes 9, 10.

(i), (iii), (iv), (vi), our method (with adversaries) achieves *similar* performance for both training loss and test accuracy as that of running SGD with local iterations but *without* any adversaries and defense mechanism at the server; and for attacks (ii), (v), the performance difference (test accuracy) is around 0.1 at epoch 40, which is still significantly better than all other methods. ¹⁰ This conforms to the inadequacy of using these methods in our setting, as described in Section 3. Note that the experiments presented in (Xie et al., 2019a; Yin et al., 2018) only implemented a benign 'label-flipping' attack, which is a data poisoning attack. This is not a dynamic attack as, unlike gradient attacks, it does not adapt to the learning process over iterations. In contrast, in

to perform significantly better than Krum in other attacks as well – note that both Krum and Bulyan are the worst performing defense mechanisms against the Baruch-attack.

¹⁰We plot the Krum performance in the training loss vs. number of epochs figures only for the attacks (ii), (v); because in all other attacks, the Krum training loss became very high (above 100) even before epoch 40 and would have prevented observing other methods' performance if we had plotted it.

all our attacks, corrupt clients send adversarial gradients in *every iteration*, making them significantly more malicious than just flipping the labels. As we have mentioned in the related work (on page 2), and we want to emphasize again, that though (Xie et al., 2019a) also studied the same problem as ours, but employed 'coordinate-wise trimmed mean' for robust gradient aggregation, their convergence bound, in our opinion, are vacuous, as the sub-optimality gap in their bounds *always* scales linearly with the diameter of the parameter space. As far as we know, ours is the first theoretical result that combines Byzantine-resilience with local iterations for high-dimensional distributed training on heterogeneous datasets with good empirical performance.

Acknowledgments

Deepesh Data would like to thank Navjot Singh for his generous help with setting up the experiments. This work was supported in part by NSF grants #1740047, #2007714, and UC-NL grant LFR-18-548554.

References

- Alistarh, D., Allen-Zhu, Z., and Li, J. Byzantine stochastic gradient descent. In *Neural Information Processing Systems (NeurIPS)*, pp. 4618–4628, 2018.
- Baruch, G., Baruch, M., and Goldberg, Y. A little is enough: Circumventing defenses for distributed learning. In *Neural Information Processing Systems (NeurIPS)*, pp. 8632–8642, 2019.
- Basu, D., Data, D., Karakus, C., and Diggavi, S. N. Qsparse-local-sgd: Distributed SGD with quantization, sparsification and local computations. In *NeurIPS*, pp. 14668–14679, 2019.
- Blanchard, P., Mhamdi, E. M. E., Guerraoui, R., and Stainer, J. Machine learning with adversaries: Byzantine tolerant gradient descent. In *NIPS*, pp. 119–129, 2017.
- Bottou, L. Large-scale machine learning with stochastic gradient descent. In *COMPSTAT*, pp. 177–186, 2010.
- Charikar, M., Steinhardt, J., and Valiant, G. Learning from untrusted data. In *STOC*, pp. 47–60, 2017.
- Chen, L., Wang, H., Charles, Z. B., and Papailiopoulos, D. S. DRACO: byzantine-resilient distributed training via redundant gradients. In *International Conference on Machine Learning (ICML)*, pp. 902–911, 2018.
- Chen, Y., Su, L., and Xu, J. Distributed statistical machine learning in adversarial settings: Byzantine gradient descent. *POMACS*, 1(2):44:1–44:25, 2017.
- Data, D. and Diggavi, S. N. Byzantine-tolerant distributed coordinate descent. In *ISIT*, pp. 2724–2728, 2019.
- Data, D. and Diggavi, S. N. Byzantine-resilient SGD in high dimensions on heterogeneous data. *CoRR*, abs/2005.07866, 2020a. URL https://arxiv.org/abs/2005.07866. Preliminary version appeared in IEEE ISIT.
- Data, D. and Diggavi, S. N. On byzantine-resilient highdimensional stochastic gradient descent. In *IEEE International Symposium on Information Theory (ISIT)*, pp. 2628–2633. IEEE, 2020b.
- Data, D., Song, L., and Diggavi, S. N. Data encoding methods for byzantine-resilient distributed optimization. In *ISIT*, pp. 2719–2723, 2019.
- Data, D., Song, L., and Diggavi, S. N. Data encoding for byzantine-resilient distributed optimization. *IEEE Transactions on Information Theory*, 67(2):1117–1140, 2021. doi: 10.1109/TIT.2020.3035868.

- Dean, J., Corrado, G., Monga, R., Chen, K., Devin, M., Le,
 Q. V., Mao, M. Z., Ranzato, M., Senior, A. W., Tucker,
 P. A., Yang, K., and Ng, A. Y. Large scale distributed
 deep networks. In *Neural Information Processing Systems* (NIPS), pp. 1232–1240, 2012.
- Diakonikolas, I. and Kane, D. M. Recent advances in algorithmic high-dimensional robust statistics. *CoRR*, abs/1911.05911, 2019.
- Diakonikolas, I., Kamath, G., Kane, D., Li, J., Moitra, A., and Stewart, A. Robust estimators in high-dimensions without the computational intractability. *SIAM J. Comput.*, 48(2):742–864, 2019.
- Dong, Y., Hopkins, S. B., and Li, J. Quantum entropy scoring for fast robust mean estimation and improved outlier detection. In Wallach, H. M., Larochelle, H., Beygelzimer, A., d'Alché-Buc, F., Fox, E. B., and Garnett, R. (eds.), Neural Information Processing Systems (NeurIPS), pp. 6065–6075, 2019.
- Ghosh, A., Hong, J., Yin, D., and Ramchandran, K. Robust federated learning in a heterogeneous environment. *CoRR*, abs/1906.06629, 2019. URL http://arxiv.org/abs/1906.06629.
- Haddadpour, F. and Mahdavi, M. On the convergence of local descent methods in federated learning. *CoRR*, abs/1910.14425, 2019. URL http://arxiv.org/abs/1910.14425.
- Haddadpour, F., Kamani, M. M., Mahdavi, M., and Cadambe, V. R. Local SGD with periodic averaging: Tighter analysis and adaptive synchronization. In *Neural Information Processing Systems (NeurIPS)*, pp. 11080–11092, 2019.
- Kairouz, P. et al. Advances and open problems in federated learning. *CoRR*, abs/1912.04977, 2019.
- Karimireddy, S. P., Kale, S., Mohri, M., Reddi, S. J., Stich, S. U., and Suresh, A. T. SCAFFOLD: stochastic controlled averaging for federated learning. In *International Conference on Machine Learning (ICML)*, pp. 5132– 5143, 2020.
- Khaled, A., Mishchenko, K., and Richtárik, P. Tighter theory for local SGD on identical and heterogeneous data. In Chiappa, S. and Calandra, R. (eds.), *International Conference on Artificial Intelligence and Statistics (AISTATS)*, pp. 4519–4529, 2020.
- Konecný, J. Stochastic, distributed and federated optimization for machine learning. *CoRR*, abs/1707.01155, 2017.
- Konecný, J., McMahan, H. B., Ramage, D., and Richtárik, P. Federated optimization: Distributed machine learning for on-device intelligence. *CoRR*, abs/1610.02527, 2016.

- Lai, K. A., Rao, A. B., and Vempala, S. S. Agnostic estimation of mean and covariance. In *FOCS*, pp. 665–674, 2016.
- Li, J. Robustness in Machine Learning (CSE 599-M); Lecture 5 Efficient filtering from spectral signatures, 2019. URL https://jerryzli.github.io/robust-ml-fall19.html.
- Li, L., Xu, W., Chen, T., Giannakis, G. B., and Ling, Q. RSA: byzantine-robust stochastic aggregation methods for distributed learning from heterogeneous datasets. In *Conference on Artificial Intelligence (AAAI)*, pp. 1544– 1551, 2019a.
- Li, X., Yang, W., Wang, S., and Zhang, Z. Communication efficient decentralized training with multiple local updates. *CoRR*, abs/1910.09126, 2019b.
- Li, X., Huang, K., Yang, W., Wang, S., and Zhang, Z. On the convergence of fedavg on non-iid data. In *International Conference on Learning Representations (ICLR)*, 2020. URL https://openreview.net/forum?id=HJxNAnVtDS.
- McMahan, B., Moore, E., Ramage, D., Hampson, S., and y Arcas, B. A. Communication-efficient learning of deep networks from decentralized data. In *International Conference on Artificial Intelligence and Statistics (AISTATS)*, pp. 1273–1282, 2017.
- Mhamdi, E. M. E., Guerraoui, R., and Rouault, S. The hidden vulnerability of distributed learning in byzantium. In *International Conference on Machine Learning (ICML)*, pp. 3518–3527, 2018.
- Mohri, M., Sivek, G., and Suresh, A. T. Agnostic federated learning. In *International Conference on Machine Learning (ICML)*, pp. 4615–4625, 2019.
- Rajput, S., Wang, H., Charles, Z. B., and Papailiopoulos, D. S. DETOX: A redundancy-based framework for faster and more robust gradient aggregation. In *NeurIPS*, pp. 10320–10330, 2019.
- Sahu, A. K., Li, T., Sanjabi, M., Zaheer, M., Talwalkar, A., and Smith, V. Federated optimization in heterogeneous networks. In *Conference on Machine Learning and Systems (MLSys)*, 2020. URL http://arxiv.org/abs/1812. 06127.
- Steinhardt, J., Charikar, M., and Valiant, G. Resilience: A criterion for learning in the presence of arbitrary outliers. In *ITCS*, pp. 45:1–45:21, 2018.
- Su, L. and Xu, J. Securing distributed gradient descent in high dimensional statistical learning. *POMACS*, 3(1): 12:1–12:41, 2019.
- Vershynin, R. Introduction to the non-asymptotic analysis of random matrices. *CoRR*, abs/1011.3027, 2010.

- Xie, C., Koyejo, O., and Gupta, I. SLSGD: secure and efficient distributed on-device machine learning. In *Machine Learning and Knowledge Discovery in Databases European Conference, ECML PKDD, Proceedings, Part II*, pp. 213–228, 2019a.
- Xie, C., Koyejo, S., and Gupta, I. Zeno: Distributed stochastic gradient descent with suspicion-based faulttolerance. In *International Conference on Machine Learn*ing (ICML), pp. 6893–6901, 2019b.
- Yin, D., Chen, Y., Ramchandran, K., and Bartlett, P. L. Byzantine-robust distributed learning: Towards optimal statistical rates. In *ICML*, pp. 5636–5645, 2018.
- Yin, D., Chen, Y., Ramchandran, K., and Bartlett, P. L. Defending against saddle point attack in byzantine-robust distributed learning. In *ICML*, pp. 7074–7084, 2019.
- Yu, H., Jin, R., and Yang, S. On the linear speedup analysis of communication efficient momentum SGD for distributed non-convex optimization. In *ICML*, pp. 7184–7193, 2019a.
- 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 *Conference on Artificial Intelligence (AAAI)*, pp. 5693–5700, 2019b.

Supplementary Material

A. Complete Proof of Theorem 2

Let $t_k, t_{k+1} \in \mathcal{I}_T$ be any two consecutive synchronization indices. For $i \in \mathcal{K}_{t_k}$ corresponding to an honest client, let $Y_i^{t_k}, Y_i^{t_k+1}, \dots, Y_i^{t_{k+1}-1}$ be a sequence of $(t_{k+1}-t_k) \leq H$ (dependent) random variables, where, for any $t \in [t_k: t_{k+1}-1]$, the random variable Y_i^t is distributed as

$$Y_i^t \sim \text{Unif}\Big(\mathcal{F}_i^{\otimes b}\big(\boldsymbol{x}_i^t\big(\boldsymbol{x}_i^{t_k}, Y_i^{t_k}, \dots, Y_i^{t-1}\big)\big)\Big). \tag{13}$$

Here, Y_i^t corresponds to the stochastic sampling of mini-batch gradients from the set $\mathcal{F}_i^{\otimes b} \big(\boldsymbol{x}_i^t \big(\boldsymbol{x}_i^{t_k}, Y_i^{t_k}, \dots, Y_i^{t-1} \big) \big)$, which itself depends on the local parameters $\boldsymbol{x}_i^{t_k}$ (which is a deterministic quantity) at the last synchronization index and the past realizations of $Y_i^{t_k}, \dots, Y_i^{t-1}$. This is because the evolution of local parameters \boldsymbol{x}_i^t depends on $\boldsymbol{x}_i^{t_k}$ and the choice of gradients in between time indices t_k and t-1. Now define $Y_i := \sum_{t=t_k}^{t_{k+1}-1} Y_i^t$; and let p_i be the distribution of Y_i . This is the distribution p_i we will take when using Lemma 1.

Claim 1. For any honest client $i \in \mathcal{K}_{t_k}$, we have $\mathbb{E}||Y_i - \mathbb{E}[Y_i]||^2 \leq \frac{H^2\sigma^2}{b}$, where expectation is taken over sampling stochastic gradients by client i between the synchronization indices t_k and t_{k+1} .

Proof. Take an arbitrary honest client $i \in \mathcal{K}_{t_k}$.

$$\mathbb{E}\|Y_i - \mathbb{E}[Y_i]\|^2 = \mathbb{E}\left\|\sum_{t=t_k}^{t_{k+1}-1} \left(Y_i^t - \mathbb{E}[Y_i^t]\right)\right\|^2 \overset{\text{(a)}}{\leq} \left(t_{k+1} - t_k\right) \sum_{t=t_k}^{t_{k+1}-1} \mathbb{E}\|Y_i^t - \mathbb{E}[Y_i^t]\|^2 \overset{\text{(b)}}{\leq} \frac{H^2 \sigma^2}{b},$$

where (a) follows from the Jensen's inequality; in (b) we used $(t_{k+1} - t_k) \le H$ and that $\mathbb{E}||Y_i^t - \mathbb{E}[Y_i^t]||^2 \le \frac{\sigma^2}{b}$ for all $j \in [H]$, which follows from the explanation below:

$$\mathbb{E}\|Y_i^t - \mathbb{E}[Y_i^t]\|^2 = \sum_{\boldsymbol{y}_i^{t_k}, \dots, \boldsymbol{y}_i^{t-1}} \Pr\left[Y_i^j = \boldsymbol{y}_i^j, j \in [t_k : t-1]\right] \times \mathbb{E}\left[\|Y_i^t - \mathbb{E}[Y_i^t]\|^2 \mid Y_i^j = \boldsymbol{y}_i^j, j \in [t_k : t-1]\right]$$

$$\stackrel{\text{(c)}}{\leq} \sum_{\boldsymbol{y}_i^{t_k}, \dots, \boldsymbol{y}_i^{t-1}} \Pr\left[Y_i^j = \boldsymbol{y}_i^j, j \in [t_k : t-1]\right] \cdot \frac{\sigma^2}{b}$$

$$= \frac{\sigma^2}{b}.$$

Note that $Y_i^t \sim \mathrm{Unif}\Big(\mathcal{F}_i^{\otimes b}\big(\boldsymbol{x}_i^t\big(\boldsymbol{x}_i^{t_k},Y_i^{t_k},\dots,Y_i^{t-1}\big)\big)\Big)$. So, when we fix the values $Y_i^{t_k} = \boldsymbol{y}_i^{t_k},\dots,Y_i^{t-1} = \boldsymbol{y}_i^{t-1}$, the parameter vector $\boldsymbol{x}_i^t\big(\boldsymbol{x}_i^{t_k},Y_i^{t_k},\dots,Y_i^{t-1}\big)$ becomes a deterministic quantity. Now we can use the variance bound (5) in order to bound $\mathbb{E}\left[\|Y_i^t - \mathbb{E}[Y_i^t]\|^2 \mid Y_i^j = \boldsymbol{y}_i^j, j \in [t_k:t-1]\right] \leq \frac{\sigma^2}{b}$. This is what we used in (c).

It is easy to see that the hypothesis of Lemma 1 is satisfied with $\mu_i = \mathbb{E}[Y_i], \sigma_{p_i}^2 = \frac{H^2 \sigma^2}{b}$ for all honest clients $i \in \mathcal{K}_{t_k}$ (note that p_i is the distribution of Y_i):

$$\mathbb{E}_{\boldsymbol{y}_i \sim p_i}[\langle \boldsymbol{y}_i - \mathbb{E}[\boldsymbol{y}_i], \boldsymbol{v} \rangle^2] \overset{(\mathrm{d})}{\leq} \mathbb{E}[\|\boldsymbol{y}_i - \mathbb{E}_{\boldsymbol{y}_i \sim p_i}[\boldsymbol{y}_i]\|^2] \cdot \|\boldsymbol{v}\|^2 \overset{(\mathrm{e})}{\leq} \frac{H^2 \sigma^2}{h},$$

where (d) follows from the Cauchy-Schwarz inequality and (e) follows from Claim 1 and $||v|| \leq 1$.

We are given K different (summations of H) gradients, out of which at least $(1-\epsilon)K$ are according to the correct distribution. By considering only the uncorrupted gradients (i.e., taking $m=(1-\epsilon)K$), we have from Lemma 1 that there exists a

subset $S \subseteq \mathcal{K}_{t_k}$ of K gradients of size $(1 - \epsilon')(1 - \epsilon)K \ge (1 - (\epsilon + \epsilon'))K \ge \frac{2K}{3}$ (where in the last inequality we used $(\epsilon + \epsilon') \le \frac{1}{3}$) that satisfies

$$\lambda_{\max}\left(\frac{1}{|\mathcal{S}|}\sum_{i\in\mathcal{S}}\left(\boldsymbol{y}_{i}-\mathbb{E}[\boldsymbol{y}_{i}]\right)\left(\boldsymbol{y}_{i}-\mathbb{E}[\boldsymbol{y}_{i}]\right)^{T}\right) \leq \frac{4H^{2}\sigma^{2}}{b\epsilon'}\left(1+\frac{3d}{2K}\right). \tag{14}$$

Note that (14) bounds the deviation of the points in S from their respective means $\mathbb{E}[y_i]$. However, in (7), we need to bound the deviation of the points in S from their sample mean $\frac{1}{|S|} \sum_{i \in S} y_i$. As it turns out, due to our use of local iterations, this will require a non-trivial amount of technical work.

From the alternate definition of the largest eigenvalue of symmetric matrices $\mathbf{A} \in \mathbb{R}^{d \times d}$, we have

$$\lambda_{\max}(\mathbf{A}) = \sup_{\boldsymbol{v} \in \mathbb{R}^d, ||\boldsymbol{v}|| = 1} \boldsymbol{v}^T \mathbf{A} \boldsymbol{v}. \tag{15}$$

Applying this with $\mathbf{A} = \frac{1}{|\mathcal{S}|} \sum_{i \in \mathcal{S}} (\boldsymbol{y}_i - \mathbb{E}[\boldsymbol{y}_i]) (\boldsymbol{y}_i - \mathbb{E}[\boldsymbol{y}_i])^T$, we can equivalently write (14) as

$$\sup_{\boldsymbol{v} \in \mathbb{R}^d: \|\boldsymbol{v}\| = 1} \left(\frac{1}{|\mathcal{S}|} \sum_{i \in \mathcal{S}} \langle \boldsymbol{y}_i - \mathbb{E}[\boldsymbol{y}_i], \boldsymbol{v} \rangle^2 \right) \le \widehat{\sigma}_0^2 := \frac{4H^2 \sigma^2}{b\epsilon'} \left(1 + \frac{3d}{2K} \right). \tag{16}$$

Define $y_{\mathcal{S}} := \frac{1}{|\mathcal{S}|} \sum_{i \in \mathcal{S}} y_i$ to be the sample mean of the points in \mathcal{S} . Take an arbitrary $v \in \mathbb{R}^d$ such that ||v|| = 1.

$$\begin{split} \frac{1}{|\mathcal{S}|} \sum_{i \in \mathcal{S}} \langle \boldsymbol{y}_i - \boldsymbol{y}_{\mathcal{S}}, \boldsymbol{v} \rangle^2 &= \frac{1}{|\mathcal{S}|} \sum_{i \in \mathcal{S}} \left[\langle \boldsymbol{y}_i - \mathbb{E}[\boldsymbol{y}_i], \boldsymbol{v} \rangle + \langle \mathbb{E}[\boldsymbol{y}_i] - \boldsymbol{y}_{\mathcal{S}}, \boldsymbol{v} \rangle \right]^2 \\ &\leq \frac{2}{|\mathcal{S}|} \sum_{i \in \mathcal{S}} \langle \boldsymbol{y}_i - \mathbb{E}[\boldsymbol{y}_i], \boldsymbol{v} \rangle^2 + \frac{2}{|\mathcal{S}|} \sum_{i \in \mathcal{S}} \langle \mathbb{E}[\boldsymbol{y}_i] - \boldsymbol{y}_{\mathcal{S}}, \boldsymbol{v} \rangle^2 \qquad \text{(using } (a+b)^2 \leq 2a^2 + 2b^2) \end{split}$$

Using (16) to bound the first term, we get

$$\leq 2\widehat{\sigma}_{0}^{2} + \frac{2}{|\mathcal{S}|} \sum_{i \in \mathcal{S}} \left\langle \mathbb{E}[\boldsymbol{y}_{i}] - \frac{1}{|\mathcal{S}|} \sum_{j \in \mathcal{S}} \boldsymbol{y}_{j}, \boldsymbol{v} \right\rangle^{2}$$

$$= 2\widehat{\sigma}_{0}^{2} + \frac{2}{|\mathcal{S}|} \sum_{i \in \mathcal{S}} \left[\frac{1}{|\mathcal{S}|} \sum_{j \in \mathcal{S}} \langle \boldsymbol{y}_{j} - \mathbb{E}[\boldsymbol{y}_{i}], \boldsymbol{v} \rangle \right]^{2}$$

$$\leq 2\widehat{\sigma}_{0}^{2} + \frac{2}{|\mathcal{S}|} \sum_{i \in \mathcal{S}} \frac{1}{|\mathcal{S}|} \sum_{j \in \mathcal{S}} \langle \boldsymbol{y}_{j} - \mathbb{E}[\boldsymbol{y}_{i}], \boldsymbol{v} \rangle^{2} \qquad \text{(using the Jensen's inequality)}$$

$$= 2\widehat{\sigma}_{0}^{2} + \frac{2}{|\mathcal{S}|} \sum_{i \in \mathcal{S}} \frac{1}{|\mathcal{S}|} \sum_{j \in \mathcal{S}} \left[\langle \boldsymbol{y}_{j} - \mathbb{E}[\boldsymbol{y}_{j}], \boldsymbol{v} \rangle + \langle \mathbb{E}[\boldsymbol{y}_{j}] - \mathbb{E}[\boldsymbol{y}_{i}], \boldsymbol{v} \rangle \right]^{2}$$

$$\leq 2\widehat{\sigma}_{0}^{2} + \frac{2}{|\mathcal{S}|} \sum_{i \in \mathcal{S}} \frac{2}{|\mathcal{S}|} \sum_{j \in \mathcal{S}} \langle \boldsymbol{y}_{j} - \mathbb{E}[\boldsymbol{y}_{j}], \boldsymbol{v} \rangle^{2} + \frac{2}{|\mathcal{S}|} \sum_{i \in \mathcal{S}} \frac{2}{|\mathcal{S}|} \sum_{j \in \mathcal{S}} \langle \mathbb{E}[\boldsymbol{y}_{j}] - \mathbb{E}[\boldsymbol{y}_{i}], \boldsymbol{v} \rangle^{2}$$

$$\leq 2\widehat{\sigma}_{0}^{2} + \frac{4}{|\mathcal{S}|} \sum_{j \in \mathcal{S}} \langle \boldsymbol{y}_{j} - \mathbb{E}[\boldsymbol{y}_{j}], \boldsymbol{v} \rangle^{2} + \frac{4}{|\mathcal{S}|} \sum_{i \in \mathcal{S}} \frac{1}{|\mathcal{S}|} \sum_{j \in \mathcal{S}} \|\mathbb{E}[\boldsymbol{y}_{j}] - \mathbb{E}[\boldsymbol{y}_{i}]\|^{2}$$

$$\qquad \qquad \text{(using the Cauchy-Schwarz inequality and that } \|\boldsymbol{v}\| \leq 1$$

$$\leq 6\widehat{\sigma}_{0}^{2} + \frac{4}{|\mathcal{S}|} \sum_{i \in \mathcal{S}} \frac{1}{|\mathcal{S}|} \sum_{i \in \mathcal{S}} \|\mathbb{E}[\boldsymbol{y}_{j}] - \mathbb{E}[\boldsymbol{y}_{i}]\|^{2}$$

$$(17)$$

Claim 2. For any $r, s \in \mathcal{K}_{t_k}$, we have

$$\|\mathbb{E}[\boldsymbol{y}_r] - \mathbb{E}[\boldsymbol{y}_s]\|^2 \le H \sum_{t=t_k}^{t_{k+1}-1} \left(6\kappa^2 + 3L^2\mathbb{E}\|\boldsymbol{x}_r^t - \boldsymbol{x}_s^t\|^2\right),$$
 (18)

where expectations in $\mathbb{E}[\mathbf{y}_r]$ and $\mathbb{E}[\mathbf{y}_s]$ are taken over sampling stochastic gradients between the synchronization indices t_k, \ldots, t_{k+1} by client r and client s, respectively.

Proof. Note that we can equivalently write $\mathbb{E}[y_r] = \mathbb{E}[Y_r]$ and $\mathbb{E}[y_s] = \mathbb{E}[Y_s]$.

$$\|\mathbb{E}[Y_r] - \mathbb{E}[Y_s]\|^2 = \|\mathbb{E}[Y_r] - \mathbb{E}[Y_s]\|^2$$

$$= \left\| \sum_{t=t_k}^{t_{k+1}-1} \left(\mathbb{E}[Y_r^t] - \mathbb{E}[Y_s^t] \right) \right\|^2$$

$$\leq (t_{k+1} - t_k) \sum_{t=t_k}^{t_{k+1}-1} \left\| \mathbb{E}[Y_r^t] - \mathbb{E}[Y_s^t] \right\|^2$$
(19)

By definition of Y_s^t from (13), we have $Y_s^t \sim \mathrm{Unif}\Big(\mathcal{F}_s^{\otimes b}\big(\boldsymbol{x}_s^t\big(\boldsymbol{x}_s^{t_k},Y_s^{t_k},\ldots,Y_s^{t-1}\big)\Big)\Big)$, which implies using (4) that $\mathbb{E}[Y_s^t] = \mathbb{E}\left[\nabla F_s\big(\boldsymbol{x}_s^t\big(\boldsymbol{x}_s^{t_k},Y_s^{t_k},\ldots,Y_s^{t-1}\big)\right)\right]$, where on the RHS, expectation is taken over $(Y_s^{t_k},\ldots,Y_s^{t-1})$. To make the notation less cluttered, in the following, for any $s \in \mathcal{K}_{t_k}$, we write \boldsymbol{x}_s^t to denote $\boldsymbol{x}_s^t\big(\boldsymbol{x}_s^{t_k},Y_s^{t_k},\ldots,Y_s^{t-1}\big)$ with the understanding that expectation is always taken over the sampling of stochastic gradients between t_k and t_{k+1} . With these substitutions, the t'th term from (20) can be written as:

$$\|\mathbb{E}[Y_{r}^{t}] - \mathbb{E}[Y_{s}^{t}]\|^{2} = \|\mathbb{E}\left[\nabla F_{r}(\boldsymbol{x}_{r}^{t}) - \nabla F_{s}(\boldsymbol{x}_{s}^{t})\right]\|^{2}$$

$$\stackrel{\text{(a)}}{\leq} \mathbb{E}\left\|\nabla F_{r}\left(\boldsymbol{x}_{r}^{t}\right) - \nabla F_{s}\left(\boldsymbol{x}_{s}^{t}\right)\right\|^{2}$$

$$\stackrel{\text{(b)}}{\leq} 3\mathbb{E}\left\|\nabla F_{r}\left(\boldsymbol{x}_{r}^{t}\right) - \nabla F\left(\boldsymbol{x}_{s}^{t}\right)\right\|^{2} + 3\mathbb{E}\left\|\nabla F_{s}\left(\boldsymbol{x}_{s}^{t}\right) - \nabla F\left(\boldsymbol{x}_{s}^{t}\right)\right\|^{2}$$

$$+ 3\mathbb{E}\left\|\nabla F\left(\boldsymbol{x}_{r}^{t}\right) - \nabla F\left(\boldsymbol{x}_{s}^{t}\right)\right\|^{2}$$

$$\stackrel{\text{(c)}}{\leq} 6\kappa^{2} + 3L^{2}\mathbb{E}\|\boldsymbol{x}_{r}^{t} - \boldsymbol{x}_{s}^{t}\|^{2}.$$

$$(21)$$

Here, (a) and (b) both follow from the Jensen's inequality. (c) used the gradient dissimilarity bound from (6) to bound the first two terms¹¹ and L-Lipschitzness of ∇F to bound the last term. Substituting the bound from (21) back in (20) and using $(t_{k+1}-t_k) \leq H$ proves Claim 2.

Using the bound from (18) in (17) gives

$$\frac{1}{|\mathcal{S}|} \sum_{i \in \mathcal{S}} \langle \boldsymbol{y}_i - \boldsymbol{y}_{\mathcal{S}}, \boldsymbol{v} \rangle^2 \leq 6\widehat{\sigma}_0^2 + \frac{4}{|\mathcal{S}|} \sum_{i \in \mathcal{S}} \frac{1}{|\mathcal{S}|} \sum_{j \in \mathcal{S}} H \sum_{t=t_k}^{t_{k+1}-1} \left(6\kappa^2 + 3L^2 \mathbb{E} \|\boldsymbol{x}_r^t - \boldsymbol{x}_s^t\|^2 \right)$$

$$= 6\widehat{\sigma}_0^2 + 24H^2\kappa^2 + \frac{12HL^2}{|\mathcal{S}|} \sum_{i \in \mathcal{S}} \frac{1}{|\mathcal{S}|} \sum_{j \in \mathcal{S}} \sum_{t=t_k}^{t_{k+1}-1} \mathbb{E} \|\boldsymbol{x}_r^t - \boldsymbol{x}_s^t\|^2 \tag{22}$$

Now we bound the last term of (22), which is the drift in local parameters at different clients in between any two synchronization indices.

Lemma 3. For any $r, s \in \mathcal{K}_{t_k}$, if $\eta \leq \frac{1}{8HL}$, we have

$$\sum_{t=t_k}^{t_{k+1}-1} \mathbb{E} \left\| \boldsymbol{x}_r^t - \boldsymbol{x}_s^t \right\|^2 \le 7H^3 \eta^2 \left(\frac{\sigma^2}{b} + 3\kappa^2 \right), \tag{23}$$

where expectation is taken over sampling stochastic gradients at clients r, s between the synchronization indices t_k and t_{k+1} .

¹¹Note that though \boldsymbol{x}_r^t 's are random quantities, we can still bound $\mathbb{E} \|\nabla F_r(\boldsymbol{x}_r^t) - \nabla F_s(\boldsymbol{x}_s^t)\|^2 \le \kappa^2$ because the gradient dissimilarity bound (6) holds uniformly over the entire domain.

Proof. For any $t \in [t_k : t_{k+1} - 1]$ and $r, s \in \mathcal{K}_{t_k}$, define $D_{r,s}^t = \mathbb{E} \| \boldsymbol{x}_r^t - \boldsymbol{x}_s^t \|^2$. Note that at synchronization time t_k , all clients in the active set \mathcal{K}_{t_k} have the same parameters, i.e., $\boldsymbol{x}_r^{t_k} = \boldsymbol{x}^{t_k}$ for every $r \in \mathcal{K}_{t_k}$.

$$D_{r,s}^{t} = \mathbb{E} \left\| \boldsymbol{x}_{r}^{t} - \boldsymbol{x}_{s}^{t} \right\|^{2} = \mathbb{E} \left\| \left(\boldsymbol{x}_{r}^{t_{k}} - \eta \sum_{j=t_{k}}^{t-1} \boldsymbol{g}_{r}(\boldsymbol{x}_{r}^{j}) \right) - \left(\boldsymbol{x}_{s}^{t_{k}} - \eta \sum_{j=t_{k}}^{t-1} \boldsymbol{g}_{s}(\boldsymbol{x}_{s}^{j}) \right) \right\|^{2}$$

$$= \eta^{2} \mathbb{E} \left\| \sum_{j=t_{k}}^{t-1} \left(\boldsymbol{g}_{r}(\boldsymbol{x}_{r}^{j}) - \boldsymbol{g}_{s}(\boldsymbol{x}_{s}^{j}) \right) \right\|^{2}$$

$$\leq \eta^{2} (t - t_{k}) \sum_{j=t_{k}}^{t-1} \mathbb{E} \left\| \boldsymbol{g}_{r}(\boldsymbol{x}_{r}^{j}) - \boldsymbol{g}_{s}(\boldsymbol{x}_{s}^{j}) \right\|^{2}$$

$$\leq \eta^{2} H \sum_{j=t_{k}}^{t-1} \left(3 \mathbb{E} \left\| \boldsymbol{g}_{r}(\boldsymbol{x}_{r}^{j}) - \nabla F_{r}(\boldsymbol{x}_{r}^{j}) \right\|^{2} + 3 \mathbb{E} \left\| \boldsymbol{g}_{s}(\boldsymbol{x}_{s}^{j}) - \nabla F_{s}(\boldsymbol{x}_{s}^{j}) \right\|^{2}$$

$$+ 3 \mathbb{E} \left\| \nabla F_{r}(\boldsymbol{x}_{r}^{j}) - \nabla F_{s}(\boldsymbol{x}_{s}^{j}) \right\|^{2} \right)$$

$$(24)$$

To bound the first and the second terms we use the variance bound from (5).¹² We can bound the third term in the same way as we bounded it in (20) and obtained (21). This gives

$$\begin{split} D_{r,s}^{t} & \leq \eta^{2} H \sum_{j=t_{k}}^{t-1} \left(\frac{6\sigma^{2}}{b} + 18\kappa^{2} + 9L^{2} \mathbb{E} \| \boldsymbol{x}_{r}^{j} - \boldsymbol{x}_{s}^{j} \|^{2} \right) \\ & \leq \frac{6H^{2}\sigma^{2}\eta^{2}}{b} + 18H^{2}\eta^{2}\kappa^{2} + 9L^{2}H\eta^{2} \sum_{j=t_{k}}^{t-1} D_{r,s}^{j} \end{split} \tag{Since } D_{r,s}^{j} = \mathbb{E} \left\| \boldsymbol{x}_{r}^{j} - \boldsymbol{x}_{s}^{j} \right\|^{2}) \end{split}$$

Taking summation from $t = t_k$ to $t_{k+1} - 1$ gives

$$\begin{split} \sum_{t=t_k}^{t_{k+1}-1} D_{r,s}^t &\leq \sum_{t=t_k}^{t_{k+1}-1} \left(\frac{6H^2\sigma^2\eta^2}{b} + 18H^2\eta^2\kappa^2 + 9L^2H\eta^2 \sum_{j=t_k}^{t-1} D_{r,s}^j \right) \\ &\leq \frac{6H^3\sigma^2\eta^2}{b} + 18H^3\eta^2\kappa^2 + 9L^2H^2\eta^2 \sum_{t=t_k}^{t_{k+1}-1} D_{r,s}^t. \end{split}$$

After rearranging terms, we get

$$(1 - 9L^2H^2\eta^2) \sum_{t=t_k}^{t_{k+1}-1} D_{r,s}^t \le \frac{6H^3\sigma^2\eta^2}{b} + 18H^3\eta^2\kappa^2.$$
 (25)

If we take $\eta \leq \frac{1}{8HL}$, we get $\left(1-9\eta^2L^2H^2\right) \geq \frac{6}{7}$. Substituting this in the LHS of (25) yields $\sum_{t=t_k}^{t_{k+1}-1} D_{r,s}^t \leq \frac{7H^3\sigma^2\eta^2}{b} + 21H^3\eta^2\kappa^2$, which proves Lemma 3.

Substituting the bound from (23) for the last term in (22) gives

$$\frac{1}{|\mathcal{S}|} \sum_{i \in \mathcal{S}} \langle \boldsymbol{y}_i - \boldsymbol{y}_{\mathcal{S}}, \boldsymbol{v} \rangle^2 \le 6\widehat{\sigma}_0^2 + 24H^2\kappa^2 + \frac{12HL^2}{|\mathcal{S}|} \sum_{i \in \mathcal{S}} \frac{1}{|\mathcal{S}|} \sum_{j \in \mathcal{S}} \left(7H^3\eta^2 \left(\frac{\sigma^2}{b} + 3\kappa^2 \right) \right)$$

$$= 6\widehat{\sigma}_0^2 + 24H^2\kappa^2 + 84H^4L^2\eta^2 \left(\frac{\sigma^2}{b} + 3\kappa^2 \right)$$

¹²Note that \boldsymbol{x}_r^j 's are random quantities, however, since the variance bound (5) holds uniformly over the entire domain, we can bound $\mathbb{E} \|\boldsymbol{g}_r(\boldsymbol{x}_r^j) - \nabla F_r(\boldsymbol{x}_r^j)\|^2 \leq \frac{\sigma^2}{b}$.

$$\leq 6\hat{\sigma}_{0}^{2} + 28H^{2}\kappa^{2} + \frac{21H^{2}\sigma^{2}}{16b} \qquad (\text{Using } \eta \leq \frac{1}{8LH})$$

$$\leq \frac{24H^{2}\sigma^{2}}{b\epsilon'} \left(1 + \frac{3d}{2K} \right) + \frac{21H^{2}\sigma^{2}}{16b} + 28H^{2}\kappa^{2} \qquad (\text{Since } \hat{\sigma}_{0}^{2} = \frac{4H^{2}\sigma^{2}}{b\epsilon'} \left(1 + \frac{3d}{2K} \right))$$

$$\leq \frac{25H^{2}\sigma^{2}}{b\epsilon'} \left(1 + \frac{3d}{2K} \right) + 28H^{2}\kappa^{2}. \qquad (26)$$

In the last inequality we used $\frac{21}{16} \leq \frac{1}{\epsilon'} \leq \frac{1}{\epsilon'} \left(1 + \frac{3d}{2K}\right)$, where the first inequality follows because $\epsilon' \leq \frac{1}{3}$. Note that (26) holds for every unit vector $\boldsymbol{v} \in \mathbb{R}^d$. Using this and substituting $\boldsymbol{g}_{i,\text{accu}}^{t_k,t_{k+1}} = \boldsymbol{y}_i, \boldsymbol{g}_{\mathcal{S},\text{accu}}^{t_k,t_{k+1}} = \boldsymbol{y}_{\mathcal{S}}$ in (26), we get

$$\sup_{\boldsymbol{v} \in \mathbb{R}^d: \|\boldsymbol{v}\| = 1} \frac{1}{|\mathcal{S}|} \sum_{i \in \mathcal{S}} \left\langle \boldsymbol{g}_{i, \text{accu}}^{t_k, t_{k+1}} - \boldsymbol{g}_{\mathcal{S}, \text{accu}}^{t_k, t_{k+1}}, \boldsymbol{v} \right\rangle^2 \leq \frac{25H^2\sigma^2}{b\epsilon'} \left(1 + \frac{3d}{2K}\right) + 28H^2\kappa^2.$$

This, in view of the alternate definition of the largest eigenvalue given in (15), is equivalent to (7), which proves Theorem 2.

B. Convergence Proof of the Strongly-Convex Part of Theorem 1

Let $\mathcal{I}_T := \{t_1, t_2, \dots, t_k, \dots\}$ with $t_1 = 0$ be the set of synchronization indices at which server selects a subset $\mathcal{K} \subseteq [R]$ of K clients and sends the current global model parameters to them. Upon receiving that, clients in \mathcal{K} performs local SGD steps based on their own local datasets until the next synchronization index, at which they send their local model parameters to the server. When server has received the updates from clients, it applies the outlier-filtering procedure RAGE (see Algorithm 1) to robustly estimate the average of the uncorrupted accumulated gradients and then updates the global model parameters. We assume that $H = \max_{i>1} (t_{i+1} - t_i)$.

At any iteration $t \in [T]$, let $\mathcal{K}_t \subseteq [R]$ denote the set of clients that are active at time t. Let $\boldsymbol{x}^t := \frac{1}{K} \sum_{r \in \mathcal{K}_t} \boldsymbol{x}_r^t$ denote the average parameter vector of the clients in the active set \mathcal{K}_t . Note that, for any $t_i \in \mathcal{I}_T$, the clients in \mathcal{K}_{t_i} remain active at all time indices t such that $t \in [t_i : t_{i+1} - 1]$.

In the following, we denote the decoded gradient at the server at any synchronization time t_{i+1} by $\hat{g}_{\text{accu}}^{t_i,t_{i+1}}$, which is an estimate of the average of the accumulated gradients between time t_i and t_{i+1} of the honest clients in \mathcal{K}_{t_i} , as in Theorem 2. From Algorithm 1, we can write the parameter update rule for the global model at the synchronization indices as:

$$\boldsymbol{x}^{t_{i+1}} = \boldsymbol{x}^{t_i} - \eta \widehat{\boldsymbol{g}}_{\text{accu}}^{t_i, t_{i+1}}.$$

Note that at any synchronization index $t_i \in \mathcal{I}_T$, when server selects a subset \mathcal{K}_{t_i} of clients and sends the global parameter vector \boldsymbol{x}^{t_i} , all clients in \mathcal{K}_{t_i} set their local model parameters to be equal to the global model parameters, i.e., $\boldsymbol{x}_r^{t_i} = \boldsymbol{x}^{t_i}$ holds for every $r \in \mathcal{K}_{t_i}$.

Now we proceed with proving the strongly-convex part of Theorem 1.

First we derive a recurrence relation for the synchronization indices and then later we extend the proof to all indices. Consider the (i+1)'st synchronization index $t_{i+1} \in \mathcal{I}_T$.

$$\boldsymbol{x}^{t_{i+1}} = \boldsymbol{x}^{t_i} - \eta \widehat{\boldsymbol{g}}_{\text{accu}}^{t_i, t_{i+1}}$$

$$= \boldsymbol{x}^{t_i} - \eta \frac{1}{K} \sum_{r \in \mathcal{K}_{t_i}} \sum_{t=t_i}^{t_{i+1}-1} \nabla F_r(\boldsymbol{x}_r^t) - \eta \left(\widehat{\boldsymbol{g}}_{\text{accu}}^{t_i, t_{i+1}} - \frac{1}{K} \sum_{r \in \mathcal{K}_{t_i}} \sum_{t=t_i}^{t_{i+1}-1} \nabla F_r(\boldsymbol{x}_r^t) \right)$$

For simplicity of notation, define $\mathcal{E} \triangleq \left(\widehat{\boldsymbol{g}}_{\mathrm{accu}}^{t_i,t_{i+1}} - \frac{1}{K}\sum_{r \in \mathcal{K}_{t_i}}\sum_{t=t_i}^{t_{i+1}-1}\nabla F_r(\boldsymbol{x}_r^t)\right)$. Substituting this in the above and using $\boldsymbol{x}^{t_i} = \frac{1}{K}\sum_{r \in \mathcal{K}_{t_i}} \boldsymbol{x}_r^{t_i}$ gives

$$\boldsymbol{x}^{t_{i+1}} = \frac{1}{K} \sum_{r \in \mathcal{K}_{t_i}} \boldsymbol{x}_r^{t_i} - \eta \frac{1}{K} \sum_{r \in \mathcal{K}_{t_i}} \sum_{t=t_i}^{t_{i+1}-1} \nabla F_r(\boldsymbol{x}_r^t) - \eta \mathcal{E}$$

$$= \frac{1}{K} \sum_{r \in \mathcal{K}_{t_i}} \left(\boldsymbol{x}_r^{t_i} - \eta \sum_{t=t_i}^{t_{i+1}-1} \nabla F_r(\boldsymbol{x}_r^t) \right) - \eta \mathcal{E}$$

$$= \frac{1}{K} \sum_{r \in \mathcal{K}_{t_i}} \left(\boldsymbol{x}_r^{t_{i+1}-1} - \eta \nabla F_r(\boldsymbol{x}_r^{t_{i+1}-1}) \right) - \eta \mathcal{E}$$

$$= \boldsymbol{x}^{t_{i+1}-1} - \eta \frac{1}{K} \sum_{r \in \mathcal{K}_{t_i}} \nabla F_r(\boldsymbol{x}_r^{t_{i+1}-1}) - \eta \mathcal{E}$$

$$= \boldsymbol{x}^{t_{i+1}-1} - \eta \nabla F(\boldsymbol{x}^{t_{i+1}-1}) + \eta \frac{1}{K} \sum_{r \in \mathcal{K}_{t_i}} \left(\nabla F(\boldsymbol{x}^{t_{i+1}-1}) - \nabla F_r(\boldsymbol{x}_r^{t_{i+1}-1}) \right) - \eta \mathcal{E}$$

$$(27)$$

Subtracting x^* from both sides gives:

$$\boldsymbol{x}^{t_{i+1}} - \boldsymbol{x}^* = \underbrace{\boldsymbol{x}^{t_{i+1}-1} - \boldsymbol{x}^* - \eta \nabla F(\boldsymbol{x}^{t_{i+1}-1})}_{=: \boldsymbol{u}} + \eta \underbrace{\frac{1}{K} \sum_{r \in \mathcal{K}_{t_i}} \left(\nabla F(\boldsymbol{x}^{t_{i+1}-1}) - \nabla F_r(\boldsymbol{x}^{t_{i+1}-1}_r) \right) - \eta \mathcal{E}}_{=: \boldsymbol{v}}$$
(28)

This gives $x^{t_{i+1}} - x^* = u + \eta(v - \mathcal{E})$. Taking norm on both sides and then squaring gives

$$\|x^{t_{i+1}} - x^*\|^2 = \|u\|^2 + \eta^2 \|v - \mathcal{E}\|^2 + 2\eta \langle u, v - \mathcal{E} \rangle$$
(29)

Now we use a simple but powerful trick on inner-products together with the inequality $2\langle a, b \rangle \leq ||a||^2 + ||b||^2$ and get:

$$2\eta \langle \boldsymbol{u}, \boldsymbol{v} - \mathcal{E} \rangle = 2 \left\langle \sqrt{\frac{\eta \mu}{2}} \boldsymbol{u}, \sqrt{\frac{2\eta}{\mu}} (\boldsymbol{v} - \mathcal{E}) \right\rangle \le \frac{\eta \mu}{2} \|\boldsymbol{u}\|^2 + \frac{2\eta}{\mu} \|\boldsymbol{v} - \mathcal{E}\|^2$$
(30)

Substituting this back in (29) gives

$$\begin{split} \left\| \boldsymbol{x}^{t_{i+1}} - \boldsymbol{x}^* \right\|^2 & \leq \left(1 + \frac{\eta \mu}{2} \right) \| \boldsymbol{u} \|^2 + \eta \left(\eta + \frac{2}{\mu} \right) \| \boldsymbol{v} - \mathcal{E} \|^2 \\ & \leq \left(1 + \frac{\eta \mu}{2} \right) \| \boldsymbol{u} \|^2 + 2\eta \left(\eta + \frac{2}{\mu} \right) \| \boldsymbol{v} \|^2 + 2\eta \left(\eta + \frac{2}{\mu} \right) \| \mathcal{E} \|^2 \end{split}$$

Substituting the values of u, v, \mathcal{E} and taking expectation w.r.t. the stochastic sampling of gradients by clients in \mathcal{K}_{t_i} between iterations t_i and t_{i+1} (while conditioning on the past) gives:

$$\mathbb{E} \left\| \boldsymbol{x}^{t_{i+1}} - \boldsymbol{x}^* \right\|^2 \le \left(1 + \frac{\mu \eta}{2} \right) \mathbb{E} \left\| \boldsymbol{x}^{t_{i+1}-1} - \eta \nabla F(\boldsymbol{x}^{t_{i+1}-1}) - \boldsymbol{x}^* \right\|^2$$

$$+ 2\eta \left(\eta + \frac{2}{\mu} \right) \mathbb{E} \left\| \frac{1}{K} \sum_{r \in \mathcal{K}_{t_i}} \left(\nabla F(\boldsymbol{x}^{t_{i+1}-1}) - \nabla F_r(\boldsymbol{x}^{t_{i+1}-1}) \right) \right\|^2$$

$$+ 2\eta \left(\eta + \frac{2}{\mu} \right) \mathbb{E} \left\| \widehat{\boldsymbol{g}}_{\text{accu}}^{t_i, t_{i+1}} - \frac{1}{K} \sum_{r \in \mathcal{K}_{t_i}} \sum_{t=t_i}^{t_{i+1}-1} \nabla F_r(\boldsymbol{x}_r^t) \right\|^2$$

$$(31)$$

Now we bound each of the three terms on the RHS of (31) separately in Claim 3, Claim 4, and Claim 5, respectively.

Claim 3. For $\eta < \frac{1}{L}$, we have

$$\mathbb{E} \| \boldsymbol{x}^{t_{i+1}-1} - \eta \nabla F(\boldsymbol{x}^{t_{i+1}-1}) - \boldsymbol{x}^* \|^2 \le (1 - \mu \eta) \mathbb{E} \| \boldsymbol{x}^{t_{i+1}-1} - \boldsymbol{x}^* \|^2.$$
(32)

Proof. Expand the LHS.

$$\mathbb{E} \|\boldsymbol{x}^{t_{i+1}-1} - \boldsymbol{x}^* - \eta \nabla F(\boldsymbol{x}^{t_{i+1}-1})\|^2 = \mathbb{E} \|\boldsymbol{x}^{t_{i+1}-1} - \boldsymbol{x}^*\|^2 + \eta^2 \mathbb{E} \|\nabla F(\boldsymbol{x}^{t_{i+1}-1})\|^2$$

$$+2\eta \mathbb{E}\left\langle \boldsymbol{x}^* - \boldsymbol{x}^{t_{i+1}-1}, \nabla F(\boldsymbol{x}^{t_{i+1}-1})\right\rangle \tag{33}$$

We can bound the second term on the RHS using L-smoothness of F, which implies that $\|\nabla F(\boldsymbol{x})\|^2 \leq 2L(F(\boldsymbol{x}) - F(\boldsymbol{x}^*))$ holds for every $\boldsymbol{x} \in \mathbb{R}^d$; see Fact 1 on page 23. We can bound the third term on the RHS using μ -strong convexity of F as follows: $\langle \boldsymbol{x}^* - \boldsymbol{x}^{t_{i+1}-1}, \nabla F(\boldsymbol{x}^{t_{i+1}-1}) \rangle \leq F(\boldsymbol{x}^*) - F(\boldsymbol{x}^{t_{i+1}-1}) - \frac{\mu}{2} \|\boldsymbol{x}^{t_{i+1}-1} - \boldsymbol{x}^*\|^2$. Substituting these back in (33) gives:

$$\mathbb{E} \| \boldsymbol{x}^{t_{i+1}-1} - \boldsymbol{x}^* - \eta \nabla F(\boldsymbol{x}^{t_{i+1}-1}) \|^2 \le (1 - \mu \eta) \, \mathbb{E} \| \boldsymbol{x}^{t_{i+1}-1} - \boldsymbol{x}^* \|^2 - 2\eta (1 - \eta L) \mathbb{E} \left(F(\boldsymbol{x}^{t_{i+1}-1}) - F(\boldsymbol{x}^*) \right)$$
(34)

Since $\eta < \frac{1}{L}$, we have $(1 - \eta L) > 0$. We also have $F(x^{t_{i+1}-1}) \ge F(x^*)$. Using these together, we can ignore the last term in the RHS of (34). This proves Claim 3.

Claim 4. For $\eta \leq \frac{1}{8HL}$, we have

$$\mathbb{E}\left\|\frac{1}{K}\sum_{r\in\mathcal{K}_{t_i}}\left(\nabla F_r(\boldsymbol{x}_r^{t_{i+1}-1}) - \nabla F(\boldsymbol{x}^{t_{i+1}-1})\right)\right\|^2 \le 2\kappa^2 + \frac{7H}{32}\left(\frac{\sigma^2}{b} + 3\kappa^2\right). \tag{35}$$

Proof. By definition, we have $x^{t_{i+1}-1} = \frac{1}{K} \sum_{r \in \mathcal{K}_{t_i}} x^{t_{i+1}-1}$.

$$\mathbb{E} \left\| \frac{1}{K} \sum_{r \in \mathcal{K}_{t_{i}}} \left(\nabla F_{r}(\boldsymbol{x}_{r}^{t_{i+1}-1}) - \nabla F(\boldsymbol{x}^{t_{i+1}-1}) \right) \right\|^{2} \leq \frac{1}{K} \sum_{r \in \mathcal{K}_{t_{i}}} \mathbb{E} \left\| \nabla F_{r}(\boldsymbol{x}_{r}^{t_{i+1}-1}) - \nabla F(\boldsymbol{x}^{t_{i+1}-1}) \right\|^{2} \\
\leq \frac{2}{K} \sum_{r \in \mathcal{K}_{t_{i}}} \left(\mathbb{E} \left\| \nabla F_{r}(\boldsymbol{x}_{r}^{t_{i+1}-1}) - \nabla F(\boldsymbol{x}_{r}^{t_{i+1}-1}) \right\|^{2} + \mathbb{E} \left\| \nabla F(\boldsymbol{x}_{r}^{t_{i+1}-1}) - \nabla F(\boldsymbol{x}^{t_{i+1}-1}) \right\|^{2} \right) \\
\leq \frac{2}{K} \sum_{r \in \mathcal{K}_{t_{i}}} \left(\kappa^{2} + L^{2} \mathbb{E} \left\| \boldsymbol{x}_{r}^{t_{i+1}-1} - \boldsymbol{x}^{t_{i+1}-1} \right\|^{2} \right) \\
= 2\kappa^{2} + \frac{2L^{2}}{K} \sum_{r \in \mathcal{K}_{t_{i}}} \mathbb{E} \left\| \boldsymbol{x}_{r}^{t_{i+1}-1} - \frac{1}{K} \sum_{s \in \mathcal{K}_{t_{i}}} \boldsymbol{x}_{s}^{t_{i+1}-1} \right\|^{2} \\
\leq 2\kappa^{2} + \frac{2L^{2}}{K} \sum_{r \in \mathcal{K}_{t_{i}}} \frac{1}{K} \sum_{s \in \mathcal{K}_{t_{i}}} \mathbb{E} \left\| \boldsymbol{x}_{r}^{t_{i+1}-1} - \boldsymbol{x}_{s}^{t_{i+1}-1} \right\|^{2} \\
\leq 2\kappa^{2} + \frac{2L^{2}}{K} \sum_{r \in \mathcal{K}_{t_{i}}} \frac{1}{K} \sum_{s \in \mathcal{K}_{t_{i}}} \left(7H^{3}\eta^{2} \left(\frac{\sigma^{2}}{b} + 3\kappa^{2} \right) \right) \\
\leq 2\kappa^{2} + 14L^{2}H^{3}\eta^{2} \left(\frac{\sigma^{2}}{b} + 3\kappa^{2} \right) \stackrel{\text{(c)}}{\leq} 2\kappa^{2} + \frac{7H}{32} \left(\frac{\sigma^{2}}{b} + 3\kappa^{2} \right) \right)$$
(36)

In (a) we used the gradient dissimilarity bound from (6) to bound the first term and L-Lipschitz gradient property of F to bound the second term. For (b), note that we have already bounded $\sum_{t=t_i}^{t_{i+1}-1} \mathbb{E} \| \boldsymbol{x}_r^t - \boldsymbol{x}_s^t \|^2 \leq 7H^3\eta^2 \left(\frac{\sigma^2}{b} + 3\kappa^2 \right)$ in (23) in Lemma 3. Since each term in the summation is trivially bounded by the same quantity, which we used in (b) to bound $\mathbb{E} \left\| \boldsymbol{x}_r^{t_{i+1}-1} - \boldsymbol{x}_s^{t_{i+1}-1} \right\|^2 \leq 7H^3\eta^2 \left(\frac{\sigma^2}{b} + 3\kappa^2 \right)$. In (c) we used $\eta \leq \frac{1}{8HL}$.

Claim 5. If $\eta \leq \frac{1}{8HL}$, then with probability at least $1 - \exp\left(-\frac{\epsilon'^2(1-\epsilon)K}{16}\right)$, we have

$$\mathbb{E}\left\|\widehat{g}_{\text{accu}}^{t_{i},t_{i+1}} - \frac{1}{K} \sum_{r \in \mathcal{K}_{t_{i}}} \sum_{t=t_{i}}^{t_{i+1}-1} \nabla F_{r}(\boldsymbol{x}_{r}^{t})\right\|^{2} \leq 3\Upsilon^{2} + \frac{8H^{2}\sigma^{2}}{b} + 30H^{2}\kappa^{2},\tag{37}$$

where $\Upsilon^2 = \mathcal{O}\left(\sigma_0^2(\epsilon + \epsilon')\right)$ and $\sigma_0^2 = \frac{25H^2\sigma^2}{b\epsilon'}\left(1 + \frac{3d}{2K}\right) + 28H^2\kappa^2$.

Proof. Let $S \subseteq \mathcal{K}_{t_i}$ denote the subset of honest clients of size $(1-(\epsilon+\epsilon'))K$, whose average accumulated gradient between time t_i and t_{i+1} that server approximates at time t_{i+1} in Theorem 2. Let the average accumulated gradient be denoted by $\mathbf{g}_{S,\mathrm{accu}}^{t_i,t_{i+1}} = \frac{1}{|S|} \sum_{r \in S} \mathbf{g}_{r,\mathrm{accu}}^{t_i,t_{i+1}}$, where $\mathbf{g}_{r,\mathrm{accu}}^{t_i,t_{i+1}} = \sum_{t=t_i}^{t_{i+1}-1} \mathbf{g}_r(\mathbf{x}_r^t)$, and server approximates it by $\widehat{\mathbf{g}}_{\mathrm{accu}}^{t_i,t_{i+1}}$. Note that S exists with probability at least $1-\exp\left(-\frac{\epsilon'^2(1-\epsilon)K}{16}\right)$. To make the notation less cluttered, for every $r \in \mathcal{K}_{t_i}$, define $\nabla F_r^{t_i,t_{i+1}} := \sum_{t=t_i}^{t_{i+1}-1} \nabla F_r(\mathbf{x}_r^t)$.

$$\mathbb{E} \left\| \widehat{g}_{\text{accu}}^{t_{i},t_{i+1}} - \frac{1}{K} \sum_{r \in \mathcal{K}_{t_{i}}} \nabla F_{r}^{t_{i},t_{i+1}} \right\|^{2} \leq 3\mathbb{E} \left\| \widehat{g}_{\text{accu}}^{t_{i},t_{i+1}} - \frac{1}{|\mathcal{S}|} \sum_{r \in \mathcal{S}} g_{r,\text{accu}}^{t_{i},t_{i+1}} \right\|^{2} \\
3\mathbb{E} \left\| \frac{1}{|\mathcal{S}|} \sum_{r \in \mathcal{S}} g_{r,\text{accu}}^{t_{i},t_{i+1}} - \frac{1}{|\mathcal{S}|} \sum_{r \in \mathcal{S}} \nabla F_{r}^{t_{i},t_{i+1}} \right\|^{2} \\
+ 3\mathbb{E} \left\| \frac{1}{|\mathcal{S}|} \sum_{r \in \mathcal{S}} \nabla F_{r}^{t_{i},t_{i+1}} - \frac{1}{K} \sum_{s \in \mathcal{K}_{t_{i}}} \nabla F_{s}^{t_{i},t_{i+1}} \right\|^{2} \tag{38}$$

Now we bound each term on the RHS of (38).

Bounding the first term on the RHS of (38). We can bound this using the second part of Theorem 2 as follows (note that given the first part of Theorem 2 is satisfied, the second part provides deterministic approximation guarantees, which implies that it also holds in expectation):

$$\mathbb{E}\left\|\widehat{g}_{\text{accu}}^{t_{i},t_{i+1}} - \frac{1}{|\mathcal{S}|} \sum_{r \in \mathcal{S}} g_{r,\text{accu}}^{t_{i},t_{i+1}} \right\|^{2} \leq \Upsilon^{2},\tag{39}$$

where $\Upsilon^2 = \mathcal{O}\left(\sigma_0^2(\epsilon + \epsilon')\right)$ and $\sigma_0^2 = \frac{25H^2\sigma^2}{b\epsilon'}\left(1 + \frac{3d}{2K}\right) + 28H^2\kappa^2$.

Bounding the second term on the RHS of (38). We can bound this using the variance bound (5).

$$\mathbb{E} \left\| \frac{1}{|\mathcal{S}|} \sum_{r \in \mathcal{S}} \left(\boldsymbol{g}_{r, \text{accu}}^{t_i, t_{i+1}} - \nabla F_r^{t_i, t_{i+1}} \right) \right\|^2 = \mathbb{E} \left\| \sum_{t=t_i}^{t_{i+1}-1} \frac{1}{|\mathcal{S}|} \sum_{r \in \mathcal{S}} \left(\boldsymbol{g}_r(\boldsymbol{x}_r^t) - \nabla F_r(\boldsymbol{x}_r^t) \right) \right\|^2 \\
\stackrel{\text{(a)}}{\leq} \left(t_{i+1} - t_i \right) \sum_{t=t_i}^{t_{i+1}-1} \mathbb{E} \left\| \frac{1}{|\mathcal{S}|} \sum_{r \in \mathcal{S}} \left(\boldsymbol{g}_r(\boldsymbol{x}_r^t) - \nabla F_r(\boldsymbol{x}_r^t) \right) \right\|^2 \\
\stackrel{\text{(b)}}{\leq} H \sum_{t=t_i}^{t_{i+1}-1} \frac{1}{|\mathcal{S}|^2} \mathbb{E} \left\| \sum_{r \in \mathcal{S}} \left(\boldsymbol{g}_r(\boldsymbol{x}_r^t) - \nabla F_r(\boldsymbol{x}_r^t) \right) \right\|^2 \\
\stackrel{\text{(c)}}{=} H \sum_{t=t_i}^{t_{i+1}-1} \frac{1}{|\mathcal{S}|^2} \sum_{r \in \mathcal{S}} \mathbb{E} \left\| \boldsymbol{g}_r(\boldsymbol{x}_r^t) - \nabla F_r(\boldsymbol{x}_r^t) \right\|^2 \\
\stackrel{\text{(d)}}{\leq} H \sum_{t=t_i}^{t_{i+1}-1} \frac{1}{|\mathcal{S}|} \frac{\sigma^2}{b} \stackrel{\text{(e)}}{\leq} \frac{4H^2 \sigma^2}{3bK}. \tag{40}$$

In (a) we used the Jensen's inequality. In (b) used $|t_{i+1}-t_i| \leq H$. In (c) we used (4) (which states that $\mathbb{E}[\boldsymbol{g}_r(\boldsymbol{x})] = \nabla F_r(\boldsymbol{x})$ holds for every honest client $r \in [R]$ and $\boldsymbol{x} \in \mathbb{R}^d$) together with that the stochastic gradients at different clients are sampled independently, and then we used the fact that the variance of independent random variables is equal to the sum of the variances. Note that $\operatorname{Var}(\boldsymbol{g}_r(\boldsymbol{x}_r^t)) = \mathbb{E} \|\boldsymbol{g}_r(\boldsymbol{x}_r^t) - \nabla F_r(\boldsymbol{x}_r^t)\|^2$. In (d) we used the variance bound (5). In (e) we used $|\mathcal{S}| \geq (1 - (\epsilon + \epsilon'))K \geq \frac{2K}{3}$, where the last inequality uses $(\epsilon + \epsilon') \leq \frac{1}{3}$.

Bounding the third term on the RHS of (38).

$$\mathbb{E} \left\| \frac{1}{|\mathcal{S}|} \sum_{r \in \mathcal{S}} \nabla F_r^{t_i, t_{i+1}} - \frac{1}{K} \sum_{s \in \mathcal{K}_{t_i}} \nabla F_s^{t_i, t_{i+1}} \right\|^2 = \mathbb{E} \left\| \sum_{t=t_i}^{t_{i+1}-1} \left(\frac{1}{|\mathcal{S}|} \sum_{r \in \mathcal{S}} \nabla F_r(\boldsymbol{x}_r^t) - \frac{1}{K} \sum_{s \in \mathcal{K}_{t_i}} \nabla F_s(\boldsymbol{x}_s^t) \right) \right\|^2 \\
\stackrel{\text{(a)}}{\leq} H \sum_{t=t_i}^{t_{i+1}-1} \mathbb{E} \left\| \frac{1}{|\mathcal{S}|} \sum_{r \in \mathcal{S}} \nabla F_r(\boldsymbol{x}_r^t) - \frac{1}{K} \sum_{s \in \mathcal{K}_{t_i}} \nabla F_s(\boldsymbol{x}_s^t) \right\|^2 \tag{41}$$

In (a), first we used the Jensen's inequality and then substituted $|t_{i+1} - t_i| \leq H$. In order to bound (41), it suffices to bound $\mathbb{E}\left\|\frac{1}{|\mathcal{S}|}\sum_{r\in\mathcal{S}}\nabla F_r(\boldsymbol{x}_r^t) - \frac{1}{K}\sum_{s\in\mathcal{K}_{t_i}}\nabla F_s(\boldsymbol{x}_s^t)\right\|^2$ for every $t\in[t_i:t_{i+1}-1]$. We bound this in the following. Take an arbitrary $t\in[t_i:t_{i+1}-1]$.

$$\begin{split} & \mathbb{E} \left\| \frac{1}{|\mathcal{S}|} \sum_{r \in \mathcal{S}} \nabla F_r(\boldsymbol{x}_r^t) - \frac{1}{K} \sum_{s \in \mathcal{K}_{t_i}} \nabla F_s(\boldsymbol{x}_s^t) \right\|^2 \leq 3\mathbb{E} \left\| \frac{1}{|\mathcal{S}|} \sum_{r \in \mathcal{S}} \left(\nabla F_r(\boldsymbol{x}_r^t) - \nabla F(\boldsymbol{x}_r^t) \right) \right\|^2 \\ & + 3\mathbb{E} \left\| \frac{1}{|\mathcal{S}|} \sum_{r \in \mathcal{S}} \nabla F(\boldsymbol{x}_r^t) - \frac{1}{K} \sum_{s \in \mathcal{K}_{t_i}} \nabla F(\boldsymbol{x}_s^t) \right\|^2 + 3\mathbb{E} \left\| \frac{1}{K} \sum_{s \in \mathcal{K}_{t_i}} \left(\nabla F(\boldsymbol{x}_s^t) - \nabla F_s(\boldsymbol{x}_s^t) \right) \right\|^2 \\ & \leq \frac{3}{|\mathcal{S}|} \sum_{r \in \mathcal{S}} \mathbb{E} \left\| \nabla F_r(\boldsymbol{x}_r^t) - \nabla F(\boldsymbol{x}_r^t) \right\|^2 + \frac{3}{K} \sum_{s \in \mathcal{K}_{t_i}} \mathbb{E} \left\| \nabla F(\boldsymbol{x}_s^t) - \nabla F_r(\boldsymbol{x}_r^t) \right\|^2 \\ & + 3\mathbb{E} \left\| \frac{1}{|\mathcal{S}|} \sum_{r \in \mathcal{S}} \left(\nabla F(\boldsymbol{x}_r^t) - \nabla F(\boldsymbol{x}^t) \right) - \frac{1}{K} \sum_{s \in \mathcal{K}_{t_i}} \left(\nabla F(\boldsymbol{x}_s^t) - \nabla F(\boldsymbol{x}^t) \right) \right\|^2 \\ & \leq 3\kappa^2 + 3\kappa^2 + 6\mathbb{E} \left\| \frac{1}{|\mathcal{S}|} \sum_{r \in \mathcal{S}} \nabla F(\boldsymbol{x}_r^t) - \nabla F(\boldsymbol{x}^t) \right\|^2 + 6\mathbb{E} \left\| \frac{1}{K} \sum_{s \in \mathcal{K}_{t_i}} \left(\nabla F(\boldsymbol{x}_s^t) - \nabla F(\boldsymbol{x}^t) \right) \right\|^2 \\ & \leq 6\kappa^2 + \frac{6}{|\mathcal{S}|} \sum_{r \in \mathcal{S}} \mathbb{E} \left\| \nabla F(\boldsymbol{x}_r^t) - \nabla F(\boldsymbol{x}^t) \right\|^2 + \frac{6}{K} \sum_{s \in \mathcal{K}_{t_i}} \mathbb{E} \left\| \nabla F(\boldsymbol{x}_s^t) - \nabla F(\boldsymbol{x}^t) \right\|^2 \\ & \leq 6\kappa^2 + \frac{6L^2}{|\mathcal{S}|} \sum_{r \in \mathcal{S}} \mathbb{E} \left\| \boldsymbol{x}_r^t - \boldsymbol{x}_s^t \right\|^2 + \frac{6L^2}{K} \sum_{r \in \mathcal{K}_{t_i}} \mathbb{E} \left\| \boldsymbol{x}_r^t - \frac{1}{K} \sum_{s \in \mathcal{K}_{t_i}} \boldsymbol{x}_s^t \right\|^2 \\ & \leq 6\kappa^2 + \frac{6L^2}{|\mathcal{S}|} \sum_{r \in \mathcal{S}} \frac{1}{K} \sum_{s \in \mathcal{K}_{t_i}} \mathbb{E} \left\| \boldsymbol{x}_r^t - \boldsymbol{x}_s^t \right\|^2 + \frac{6L^2}{K} \sum_{r \in \mathcal{K}_{t_i}} \frac{1}{K} \sum_{s \in \mathcal{K}_{t_i}} \mathbb{E} \left\| \boldsymbol{x}_r^t - \boldsymbol{x}_s^t \right\|^2 \\ & \leq 6\kappa^2 + \frac{6L^2}{|\mathcal{S}|} \sum_{r \in \mathcal{S}} \frac{1}{K} \sum_{s \in \mathcal{K}_{t_i}} \mathbb{E} \left\| \boldsymbol{x}_r^t - \boldsymbol{x}_s^t \right\|^2 + \frac{6L^2}{K} \sum_{r \in \mathcal{K}_{t_i}} \frac{1}{K} \sum_{s \in \mathcal{K}_{t_i}} \mathbb{E} \left\| \boldsymbol{x}_r^t - \boldsymbol{x}_s^t \right\|^2 \end{aligned}$$

Substituting this back in (41) gives:

$$\mathbb{E} \left\| \frac{1}{|\mathcal{S}|} \sum_{r \in \mathcal{S}} \nabla F_r^{t_i, t_{i+1}} - \frac{1}{K} \sum_{s \in \mathcal{K}_{t_i}} \nabla F_s^{t_i, t_{i+1}} \right\|^2 \le H \sum_{t=t_i}^{t_{i+1}-1} 6\kappa^2$$

$$+ H \sum_{t=t_i}^{t_{i+1}-1} \left(\frac{6L^2}{|\mathcal{S}|} \sum_{r \in \mathcal{S}} \frac{1}{K} \sum_{s \in \mathcal{K}_{t_i}} \mathbb{E} \left\| \boldsymbol{x}_r^t - \boldsymbol{x}_s^t \right\|^2 + \frac{6L^2}{K} \sum_{r \in \mathcal{K}_{t_i}} \frac{1}{K} \sum_{s \in \mathcal{K}_{t_i}} \mathbb{E} \left\| \boldsymbol{x}_r^t - \boldsymbol{x}_s^t \right\|^2 \right)$$

$$\stackrel{\text{(a)}}{\le} 6H^2 \kappa^2 + 6HL^2 \left(7H^3 \eta^2 \left(\frac{\sigma^2}{b} + 3\kappa^2 \right) \right) + 6HL^2 \left(7H^3 \eta^2 \left(\frac{\sigma^2}{b} + 3\kappa^2 \right) \right)$$

$$= 6H^2 \kappa^2 + 84L^2 H^4 \eta^2 \left(\frac{\sigma^2}{b} + 3\kappa^2 \right)$$

$$\stackrel{\text{(b)}}{\leq} 10H^2\kappa^2 + \frac{21H^2\sigma^2}{16b}.\tag{42}$$

In (a) we used $t_{i+1} - t_i \le H$ and the bound $\sum_{t=t_i}^{t_{i+1}-1} \mathbb{E} \| \boldsymbol{x}_r^t - \boldsymbol{x}_s^t \|^2 \le 7H^3\eta^2 \left(\frac{\sigma^2}{b} + 3\kappa^2 \right)$, which holds when $\eta \le \frac{1}{8HL}$; we have already shown this in (23) in Lemma 3. In (b) we used $\eta \le \frac{1}{8HL}$.

Substituting the bounds from (39), (40), (42) in (38) gives

$$\mathbb{E}\left\|\widehat{g}_{\text{accu}}^{t_{i},t_{i+1}} - \frac{1}{K} \sum_{r \in \mathcal{K}_{t_{i}}} \nabla F_{r}^{t_{i},t_{i+1}} \right\|^{2} \leq 3\Upsilon^{2} + \frac{4H^{2}\sigma^{2}}{bK} + 3\left(10H^{2}\kappa^{2} + \frac{21H^{2}\sigma^{2}}{16b}\right)$$

$$\leq 3\Upsilon^{2} + \frac{4H^{2}\sigma^{2}}{bK} + 30H^{2}\kappa^{2} + \frac{4H^{2}\sigma^{2}}{b}$$

$$= 3\Upsilon^{2} + \frac{8H^{2}\sigma^{2}}{b} + 30H^{2}\kappa^{2},$$

where $\Upsilon^2 = \mathcal{O}\left(\sigma_0^2(\epsilon + \epsilon')\right)$ and $\sigma_0^2 = \frac{25H^2\sigma^2}{b\epsilon'}\left(1 + \frac{3d}{2K}\right) + 28H^2\kappa^2$.

This completes the proof of Claim 5.

Using the bounds from (32), (35), (37) in (31) and using $\left(1 + \frac{\mu\eta}{2}\right)\left(1 - \mu\eta\right) \leq \left(1 - \frac{\mu\eta}{2}\right)$ for the first term gives

$$\mathbb{E} \|\boldsymbol{x}^{t_{i+1}} - \boldsymbol{x}^*\|^2 \le \left(1 - \frac{\mu\eta}{2}\right) \mathbb{E} \|\boldsymbol{x}^{t_{i+1}-1} - \boldsymbol{x}^*\|^2 + 2\eta \left(\eta + \frac{2}{\mu}\right) \left(2\kappa^2 + \frac{7H}{32} \left(\frac{\sigma^2}{b} + 3\kappa^2\right)\right) + 2\eta \left(\eta + \frac{2}{\mu}\right) \left(3\Upsilon^2 + \frac{8H^2\sigma^2}{b} + 30H^2\kappa^2\right) \le \left(1 - \frac{\mu\eta}{2}\right) \mathbb{E} \|\boldsymbol{x}^{t_{i+1}-1} - \boldsymbol{x}^*\|^2 + \frac{6\eta}{\mu} \left(3\Upsilon^2 + \frac{9H^2\sigma^2}{b} + 33H^2\kappa^2\right),$$
(43)

where $\Upsilon^2 = \mathcal{O}\left(\sigma_0^2(\epsilon + \epsilon')\right)$ and $\sigma_0^2 = \frac{25H^2\sigma^2}{b\epsilon'}\left(1 + \frac{3d}{2K}\right) + 28H^2\kappa^2$. In the last inequality (43) we used $\eta \leq \frac{1}{8LH} \leq \frac{1}{L} \leq \frac{1}{\mu}$, which implies $(\eta + \frac{2}{\mu}) \leq \frac{3}{\mu}$. Note that (43) holds with probability at least $1 - \exp\left(-\frac{\epsilon'^2(1-\epsilon)K}{16}\right)$.

Note that above recurrence in (43) holds only at the synchronization indices $t_i \in \mathcal{I}_T$ for $i = 1, 2, 3, \ldots$. However, in order to establish a recurrence that we can use to prove convergence, we need to show a recurrence relation for all t. Now we give a recurrence at non-synchronization indices.

Take an arbitrary $t \in [T]$ and let $t_i \in \mathcal{I}_T$ be such that $t \in [t_i : t_{i+1} - 1]$; when $H \geq 2$, such t's exist. Note that $\boldsymbol{x}^t = \frac{1}{K} \sum_{r \in \mathcal{K}_{t_i}} \boldsymbol{x}_r^t$.

$$\boldsymbol{x}^{t+1} = \boldsymbol{x}^{t} - \eta \frac{1}{K} \sum_{r \in \mathcal{K}_{t_{i}}} \boldsymbol{g}_{r}(\boldsymbol{x}_{r}^{t})$$

$$= \boldsymbol{x}^{t} - \eta \frac{1}{K} \sum_{r \in \mathcal{K}_{t_{i}}} \nabla F_{r}(\boldsymbol{x}_{r}^{t}) - \eta \left(\frac{1}{K} \sum_{r \in \mathcal{K}_{t_{i}}} \boldsymbol{g}_{r}(\boldsymbol{x}_{r}^{t}) - \frac{1}{K} \sum_{r \in \mathcal{K}_{t_{i}}} \nabla F_{r}(\boldsymbol{x}_{r}^{t}) \right)$$

$$= \boldsymbol{x}^{t} - \eta \nabla F(\boldsymbol{x}^{t}) + \frac{\eta}{K} \sum_{r \in \mathcal{K}_{t_{i}}} \left(\nabla F(\boldsymbol{x}^{t}) - \nabla F_{r}(\boldsymbol{x}_{r}^{t}) \right) - \frac{\eta}{K} \sum_{r \in \mathcal{K}_{t_{i}}} \left(\boldsymbol{g}_{r}(\boldsymbol{x}_{r}^{t}) - \nabla F_{r}(\boldsymbol{x}_{r}^{t}) \right)$$

$$(44)$$

Now, subtracting x^* from both sides and following the same steps as in from (28) to (31), we get (in the following, expectation is taken w.r.t. the stochastic sampling of gradients at the t'th iteration while conditioning on the past):

$$\mathbb{E} \|\boldsymbol{x}^{t+1} - \boldsymbol{x}^*\|^2 \le \left(1 + \frac{\mu\eta}{2}\right) \mathbb{E} \|\boldsymbol{x}^t - \boldsymbol{x}^* - \eta \nabla F(\boldsymbol{x}^t)\|^2 + 2\eta \left(\eta + \frac{2}{\mu}\right) \mathbb{E} \left\|\frac{1}{K} \sum_{r \in \mathcal{K}_{t_i}} \left(\nabla F(\boldsymbol{x}^t) - \nabla F_r(\boldsymbol{x}_r^t)\right)\right\|^2$$

$$+2\eta \left(\eta + \frac{2}{\mu}\right) \mathbb{E} \left\| \frac{1}{K} \sum_{r \in \mathcal{K}_{t_i}} \left(\boldsymbol{g}_r(\boldsymbol{x}_r^t) - \nabla F_r(\boldsymbol{x}_r^t) \right) \right\|^2$$
(45)

We can bound the first and the second terms on the RHS of (45) using (32) and (35), respectively, as $\mathbb{E} \| \boldsymbol{x}^t - \eta \nabla F(\boldsymbol{x}^t) - \boldsymbol{x}^* \|^2 \le (1 - \mu \eta) \mathbb{E} \| \boldsymbol{x}^t - \boldsymbol{x}^* \|^2$ and $\mathbb{E} \left\| \frac{1}{K} \sum_{r \in \mathcal{K}_{t_i}} (\nabla F(\boldsymbol{x}^t) - \nabla F_r(\boldsymbol{x}_r^t)) \right\|^2 \le 2\kappa^2 + \frac{7H}{32} \left(\frac{\sigma^2}{b} + 3\kappa^2 \right)$. To bound the third term on the RHS of (45), we use the fact that variance of the sum of independent random variables is equal to the sum of the variances and that clients sample stochastic gradients $\boldsymbol{g}_r(\boldsymbol{x}_r^t)$ independent of each other; using this fact and (5), we can bound $\mathbb{E} \left\| \frac{1}{K} \sum_{r \in \mathcal{K}_{t_i}} (\boldsymbol{g}_r(\boldsymbol{x}_r^t) - \nabla F_r(\boldsymbol{x}_r^t)) \right\|^2 \le \frac{\sigma^2}{bK}$. Substituting these in (45) and using $\left(1 + \frac{\mu \eta}{2}\right) (1 - \mu \eta) \le \left(1 - \frac{\mu \eta}{2}\right)$ for the first term and $\left(\eta + \frac{2}{\mu}\right) \le \frac{3}{\mu}$ (which follows because $\eta \le \frac{1}{8HL} \le \frac{1}{L} \le \frac{1}{\mu}$) give

$$\mathbb{E} \left\| \boldsymbol{x}^{t+1} - \boldsymbol{x}^* \right\|^2 \le \left(1 - \frac{\mu \eta}{2} \right) \mathbb{E} \left\| \boldsymbol{x}^t - \boldsymbol{x}^* \right\|^2 + \frac{6\eta}{\mu} \left(2\kappa^2 + \frac{7H}{32} \left(\frac{\sigma^2}{b} + 3\kappa^2 \right) + \frac{\sigma^2}{bK} \right)$$

$$\le \left(1 - \frac{\mu \eta}{2} \right) \mathbb{E} \left\| \boldsymbol{x}^t - \boldsymbol{x}^* \right\|^2 + \frac{6\eta}{\mu} \left(3H\kappa^2 + \frac{2H\sigma^2}{b} \right)$$
(46)

Note that (46) holds with probability 1.

Now we have a recurrence at the synchronization indices given in (43) and at non-synchronization indices given in (46). Let $\alpha=\left(1-\frac{\mu\eta}{2}\right),\,\beta_1=\left(3\varUpsilon^2+\frac{9H^2\sigma^2}{b}+33H^2\kappa^2\right)$, and $\beta_2=\left(3H\kappa^2+\frac{2H\sigma^2}{b}\right)$. Substituting these and using (43) for the synchronization indices and (46) for the rest of the indices, we get:

$$\mathbb{E} \|\boldsymbol{x}^{T} - \boldsymbol{x}^{*}\|^{2} \leq \alpha^{T} \|\boldsymbol{x}^{0} - \boldsymbol{x}^{*}\|^{2} + \frac{6\eta}{\mu} \left(\sum_{i=0}^{T/H} \sum_{j=1}^{H-1} \alpha^{iH+j} \beta_{2} + \sum_{i=0}^{T/H} \alpha^{iH} \beta_{1} \right)$$

$$\leq \alpha^{T} \|\boldsymbol{x}^{0} - \boldsymbol{x}^{*}\|^{2} + \frac{6\eta}{\mu} \left(\sum_{i=0}^{\infty} \alpha^{i} \beta_{2} + \sum_{i=0}^{\infty} \alpha^{iH} \beta_{1} \right)$$

$$= \alpha^{T} \|\boldsymbol{x}^{0} - \boldsymbol{x}^{*}\|^{2} + \frac{6\eta}{\mu} \left(\frac{1}{1 - \alpha} \beta_{2} + \frac{1}{1 - \alpha^{H}} \beta_{1} \right)$$

$$(48)$$

Since $\alpha=\left(1-\frac{\mu\eta}{2}\right)$, we have $\alpha^H=\left(1-\frac{\mu\eta}{2}\right)^H\overset{\text{(a)}}{\leq}\exp(-\frac{\mu\eta H}{2})\overset{\text{(b)}}{\leq}1-\frac{\mu\eta H}{2}+\left(\frac{\mu\eta H}{2}\right)^2\overset{\text{(c)}}{\leq}1-\frac{\mu\eta H}{2}+\frac{1}{16}\frac{\mu\eta H}{2}=1-\frac{15}{16}\frac{\mu\eta H}{2}.$ In (a) we used the inequality $(1-\frac{1}{x})^x\leq\frac{1}{e}$ which holds for any x>0; in (b) we used $\exp(-x)\leq 1-x+x^2$ which holds for any $x\geq0$; in (c) we used $\eta\leq\frac{1}{8HL}$ and $\mu\leq L$, which together imply $\frac{\mu\eta H}{2}\leq\frac{1}{16}.$ Substituting these in (48) gives

$$\mathbb{E} \|\boldsymbol{x}^{T} - \boldsymbol{x}^{*}\|^{2} \leq \left(1 - \frac{\mu\eta}{2}\right)^{T} \|\boldsymbol{x}^{0} - \boldsymbol{x}^{*}\|^{2} + \frac{6\eta}{\mu} \left(\frac{2}{\mu\eta}\beta_{2} + \frac{32}{15\mu\eta H}\beta_{1}\right)
\leq \left(1 - \frac{\mu\eta}{2}\right)^{T} \|\boldsymbol{x}^{0} - \boldsymbol{x}^{*}\|^{2} + \frac{6 \times 32}{15\mu^{2}} \left(\frac{15}{16}\beta_{2} + \frac{1}{H}\beta_{1}\right)
\leq \left(1 - \frac{\mu\eta}{2}\right)^{T} \|\boldsymbol{x}^{0} - \boldsymbol{x}^{*}\|^{2} + \frac{13}{\mu^{2}} \left(\frac{3\Upsilon^{2}}{H} + \frac{11H\sigma^{2}}{b} + 36H\kappa^{2}\right)$$
(49)

Note that the last term on the RHS of (49) is independent of η , which together with the dependence of η on the first term implies that bigger the η , faster the convergence. Since we need $\eta \leq \frac{1}{8HL}$ for Claim 4 and Claim 5 to hold, we choose $\eta = \frac{1}{8HL}$. Substituting this in (49) yields the convergence rate in the strongly-convex part of Theorem 1.

Error probability analysis. Note that (43) holds with probability at least $1 - \exp\left(-\frac{\epsilon'^2(1-\epsilon)K}{16}\right)$ and (46) holds with probability 1. Since to arrive at (47) (which leads to our final bound (49)), we used (43) $\frac{T}{H}$ times and (46) $\left(T - \frac{T}{H}\right)$ times; as a consequence, by union bound, we have that (49) holds with probability at least $1 - \frac{T}{H} \exp\left(-\frac{\epsilon'^2(1-\epsilon)K}{16}\right)$, which is at least $(1-\delta)$, for any $\delta > 0$, provided we run our algorithm for at most $T \le \delta H \exp\left(\frac{\epsilon'^2(1-\epsilon)K}{16}\right)$ iterations.

This concludes the proof of the strongly-convex part of Theorem 1.

Fact 1. Let $F: \mathbb{R}^d \to \mathbb{R}$ be an L-smooth function with a global minimizer x^* . Then, for every $x \in \mathbb{R}^d$, we have

$$\|\nabla F(\boldsymbol{x})\|^2 \le 2L(F(\boldsymbol{x}) - F(\boldsymbol{x}^*)).$$

Proof. By definition of *L*-smoothness, we have $F(y) \leq F(x) + \langle \nabla F(x), y - x \rangle + \frac{L}{2} ||y - x||^2$ holds for every $x, y \in \mathbb{R}^d$. Fix an arbitrary $x \in \mathbb{R}^d$ and take infimum over y on both sides:

$$\inf_{\boldsymbol{y}} F(\boldsymbol{y}) \leq \inf_{\boldsymbol{y}} \left(F(\boldsymbol{x}) + \langle \nabla F(\boldsymbol{x}), \boldsymbol{y} - \boldsymbol{x} \rangle + \frac{L}{2} \|\boldsymbol{y} - \boldsymbol{x}\|^{2} \right)$$

$$\stackrel{\text{(a)}}{=} \inf_{\boldsymbol{v}: \|\boldsymbol{v}\| = 1} \inf_{\boldsymbol{t}} \left(F(\boldsymbol{x}) + t \langle \nabla F(\boldsymbol{x}), \boldsymbol{v} \rangle + \frac{Lt^{2}}{2} \right)$$

$$\stackrel{\text{(b)}}{=} \inf_{\boldsymbol{v}: \|\boldsymbol{v}\| = 1} \left(F(\boldsymbol{x}) - \frac{1}{2L} \langle \nabla F(\boldsymbol{x}), \boldsymbol{v} \rangle^{2} \right)$$

$$\stackrel{\text{(c)}}{=} \left(F(\boldsymbol{x}) - \frac{1}{2L} \|\nabla F(\boldsymbol{x})\|^{2} \right)$$

The value of t that minimizes the RHS of (a) is $t = -\frac{1}{L} \langle \nabla F(\boldsymbol{x}), \boldsymbol{v} \rangle$, this implies (b); (c) follows from the Cauchy-Schwarz inequality: $\langle \boldsymbol{u}, \boldsymbol{v} \rangle \leq \|\boldsymbol{u}\| \|\boldsymbol{v}\|$, where equality is achieved whenever $\boldsymbol{u} = \boldsymbol{v}$. Now, substituting $\inf_{\boldsymbol{y}} F(\boldsymbol{y}) = F(\boldsymbol{x}^*)$ yields the result.

C. Convergence Proof of the Non-Convex Part of Theorem 1

Let $\mathcal{K}_t \subseteq [R]$ denote the subset of clients of size $|\mathcal{K}_t| = K$ sampled at the t'th iteration. For any $t \in [t_i : t_{i+1} - 1]$, let $x^t = \frac{1}{K} \sum_{k \in \mathcal{K}_{t,i}} x_k^t$ denote the average of the local parameters of clients in the sampling set \mathcal{K}_{t_i} .

Similar to the proof given in Appendix B for the strongly-convex part of Theorem 1, here also, first we derive a recurrence for the synchronization indices and then for non-synchronization indices.

For the synchronization indices $t_1, t_2, \ldots, t_k, \ldots \in \mathcal{I}_T$, from (27), we have

$$\mathbf{x}^{t_{i+1}} = \mathbf{x}^{t_{i+1}-1} - \eta \nabla F(\mathbf{x}^{t_{i+1}-1}) + \eta C$$
(50)

where

$$C = \frac{1}{K} \sum_{r \in \mathcal{K}_{t_i}} \left(\nabla F(\boldsymbol{x}^{t_{i+1}-1}) - \nabla F_r(\boldsymbol{x}_r^{t_{i+1}-1}) \right) - \left(\widehat{\boldsymbol{g}}_{\text{accu}}^{t_i, t_{i+1}} - \frac{1}{K} \sum_{r \in \mathcal{K}_{t_i}} \sum_{t=t_i}^{t_{i+1}-1} \nabla F_r(\boldsymbol{x}_r^t) \right).$$
 (51)

Now, using the definition of L-smoothness in (50), we have

$$F(\boldsymbol{x}^{t_{i+1}}) \leq F(\boldsymbol{x}^{t_{i+1}-1}) + \left\langle \nabla F(\boldsymbol{x}^{t_{i+1}-1}), \boldsymbol{x}^{t_{i+1}} - \boldsymbol{x}^{t_{i+1}-1} \right\rangle + \frac{L}{2} \left\| \boldsymbol{x}^{t_{i+1}} - \boldsymbol{x}^{t_{i+1}-1} \right\|^{2}$$

$$= F(\boldsymbol{x}^{t_{i+1}-1}) - \eta \left\langle \nabla F(\boldsymbol{x}^{t_{i+1}-1}), \nabla F(\boldsymbol{x}^{t_{i+1}-1}) - C \right\rangle + \frac{\eta^{2}L}{2} \left\| \nabla F(\boldsymbol{x}^{t_{i+1}-1}) - C \right\|^{2}$$

$$= F(\boldsymbol{x}^{t_{i+1}-1}) - \eta \left\| \nabla F(\boldsymbol{x}^{t_{i+1}-1}) \right\|^{2} + \eta \left\langle \nabla F(\boldsymbol{x}^{t_{i+1}-1}), C \right\rangle + \frac{\eta^{2}L}{2} \left\| \nabla F(\boldsymbol{x}^{t_{i+1}-1}) - C \right\|^{2}$$

$$\stackrel{\text{(a)}}{\leq} F(\boldsymbol{x}^{t_{i+1}-1}) - \eta \left\| \nabla F(\boldsymbol{x}^{t_{i+1}-1}) \right\|^{2} + \eta \left(\frac{\left\| \nabla F(\boldsymbol{x}^{t_{i+1}-1}) \right\|^{2}}{4} + \left\| C \right\|^{2} \right)$$

$$+ \frac{\eta^{2}L}{2} \left\| \nabla F(\boldsymbol{x}^{t_{i+1}-1}) - C \right\|^{2}$$

$$\stackrel{\text{(b)}}{\leq} F(\boldsymbol{x}^{t_{i+1}-1}) - \frac{3\eta}{4} \left\| \nabla F(\boldsymbol{x}^{t_{i+1}-1}) \right\|^{2} + \eta \|C\|^{2} + \eta^{2}L \left(\left\| \nabla F(\boldsymbol{x}^{t_{i+1}-1}) \right\|^{2} + \|C\|^{2} \right)$$

$$= F(\boldsymbol{x}^{t_{i+1}-1}) - \eta \left(\frac{3}{4} - \eta L\right) \left\|\nabla F(\boldsymbol{x}^{t_{i+1}-1})\right\|^2 + \eta \left(1 + \eta L\right) \|C\|^2$$
(52)

In (a), we used the inequality $2\langle \boldsymbol{a}, \boldsymbol{b} \rangle \leq \tau \|\boldsymbol{a}\|^2 + \frac{1}{\tau} \|\boldsymbol{b}\|^2$, which holds for every $\tau > 0$, and we used $\tau = \frac{1}{2}$ in (a). In (b), we used the inequality $\|\boldsymbol{a} + \boldsymbol{b}\|^2 \leq 2(\|\boldsymbol{a}\|^2 + \|\boldsymbol{b}\|^2)$. For $\eta \leq \frac{1}{8HL} \leq \frac{1}{8L}$, we have $(3/4 - \eta L) \geq 1/2$ and $(1 + \eta L) \leq \frac{9}{8}$. Substituting these in (52) and taking expectation w.r.t. the stochastic sampling of gradients at clients in \mathcal{K}_{i_t} between iterations t_i and t_{i+1} (while conditioning on the past) gives:

$$\mathbb{E}[F(\boldsymbol{x}^{t_{i+1}})] \leq \mathbb{E}[F(\boldsymbol{x}^{t_{i+1}-1})] - \frac{\eta}{2}\mathbb{E}\left\|\nabla F(\boldsymbol{x}^{t_{i+1}-1})\right\|^2 + \frac{9\eta}{8}\mathbb{E}\|C\|^2.$$
 (53)

Now we bound $\mathbb{E}||C||^2$. Substituting the value of C from (51) gives:

$$\mathbb{E}\|C\|^{2} \leq 2\mathbb{E}\left\|\frac{1}{K}\sum_{r\in\mathcal{K}_{t_{i}}}\left(\nabla F(\boldsymbol{x}^{t_{i+1}-1}) - \nabla F_{r}(\boldsymbol{x}_{r}^{t_{i+1}-1})\right)\right\|^{2} + 2\mathbb{E}\left\|\widehat{\boldsymbol{g}}_{\text{accu}}^{t_{i},t_{i+1}} - \frac{1}{K}\sum_{r\in\mathcal{K}_{t_{i}}}\sum_{t=t_{i}}^{t_{i+1}-1}\nabla F_{r}(\boldsymbol{x}_{r}^{t})\right\|^{2} \\
\leq 2\left(2\kappa^{2} + \frac{7H}{32}\left(\frac{\sigma^{2}}{b} + 3\kappa^{2}\right)\right) + 2\left(3\boldsymbol{\Upsilon}^{2} + \frac{8H^{2}\sigma^{2}}{b} + 30H^{2}\kappa^{2}\right) \\
\leq 2\left(3\boldsymbol{\Upsilon}^{2} + \frac{9H^{2}\sigma^{2}}{b} + 33H^{2}\kappa^{2}\right) \tag{54}$$

Here, the first inequality used $\|a + b\|^2 \le 2(\|a\|^2 + \|b\|^2)$ and the second inequality used the bounds from (35) and (37). Substituting the bound from (54) into (53) gives

$$\mathbb{E}[F(\boldsymbol{x}^{t_{i+1}})] \le \mathbb{E}[F(\boldsymbol{x}^{t_{i+1}-1})] - \frac{\eta}{2}\mathbb{E}\left\|\nabla F(\boldsymbol{x}^{t_{i+1}-1})\right\|^2 + \frac{9\eta}{4}\left(3\Upsilon^2 + \frac{9H^2\sigma^2}{b} + 33H^2\kappa^2\right)$$
(55)

where $\Upsilon^2 = \mathcal{O}\left(\sigma_0^2(\epsilon+\epsilon')\right)$ and $\sigma_0^2 = \frac{25H^2\sigma^2}{b\epsilon'}\left(1+\frac{3d}{2K}\right) + 28H^2\kappa^2$. Note that (55) holds with probability at least $1-\exp\left(-\frac{\epsilon'^2(1-\epsilon)K}{16}\right)$.

Note that the above recurrence in (55) holds only at the synchronization indices $t_i \in \mathcal{I}_T$ for $i = 1, 2, 3, \ldots$ Now we give a recurrence at non-synchronization indices.

We have done a similar calculation in the proof of the strongly-convex part of Theorem 1 in Appendix B.

Take an arbitrary $t \in [T]$ and let $t_i \in \mathcal{I}_T$ be such that $t \in [t_i : t_{i+1} - 1]$; when $H \geq 2$, such t's exist. Note that $\boldsymbol{x}^t = \frac{1}{K} \sum_{r \in \mathcal{K}_{t_i}} \boldsymbol{x}_r^t$.

From (44), we have $\mathbf{x}^{t+1} = \mathbf{x}^t - \eta \nabla F(\mathbf{x}^t) + \eta D$, where

$$D = \frac{1}{K} \sum_{r \in \mathcal{K}_{t_i}} \left(\nabla F(\boldsymbol{x}^t) - \nabla F_r(\boldsymbol{x}_r^t) \right) - \frac{1}{K} \sum_{r \in \mathcal{K}_{t_i}} \left(\boldsymbol{g}_r(\boldsymbol{x}_r^t) - \nabla F_r(\boldsymbol{x}_r^t) \right).$$

Using L-smoothness of F, and then performing similar algebraic manipulations that we used in order to arrive at (53), we get:

$$\mathbb{E}[F(\boldsymbol{x}^{t+1})] \le \mathbb{E}[F(\boldsymbol{x}^t)] - \frac{\eta}{2} \mathbb{E} \left\| \nabla F(\boldsymbol{x}^t) \right\|^2 + \frac{9\eta}{8} \mathbb{E} \|D\|^2$$
(56)

Now we bound $\mathbb{E}||D||^2$:

$$\mathbb{E}\|D\|^{2} \leq 2\mathbb{E}\left\|\frac{1}{K}\sum_{r \in \mathcal{K}_{t_{i}}}\left(\nabla F(\boldsymbol{x}^{t}) - \nabla F_{r}(\boldsymbol{x}_{r}^{t})\right)\right\|^{2} + 2\mathbb{E}\left\|\frac{1}{K}\sum_{r \in \mathcal{K}_{t_{i}}}\left(\boldsymbol{g}_{r}(\boldsymbol{x}_{r}^{t}) - \nabla F_{r}(\boldsymbol{x}_{r}^{t})\right)\right\|^{2} \\
\leq 2\left(2\kappa^{2} + \frac{7H}{32}\left(\frac{\sigma^{2}}{b} + 3\kappa^{2}\right) + \frac{\sigma^{2}}{bK}\right)$$

$$\leq 2\left(3H\kappa^2 + \frac{2H\sigma^2}{b}\right) \tag{57}$$

Here, the second inequality used the same bounds on both the quantities on the RHS of the first inequality that we used to go from (45) to (46).

Substituting the bound on $\mathbb{E}||D||^2$ from (57) into (56) gives

$$\mathbb{E}[F(\boldsymbol{x}^{t+1})] \le \mathbb{E}[F(\boldsymbol{x}^t)] - \frac{\eta}{2} \mathbb{E} \left\| \nabla F(\boldsymbol{x}^t) \right\|^2 + \frac{9\eta}{4} \left(3H\kappa^2 + \frac{2H\sigma^2}{b} \right)$$
 (58)

Note that (58) holds with probability 1.

Now we have a recurrence at synchronization indices given in (55) and at non-synchronization indices given in (58). Adding (55) and (58) from t = 0 to T (use (55) for the synchronization indices and (58) for the rest of the indices) gives:

$$\sum_{t=0}^{T} \mathbb{E}[F(\boldsymbol{x}^{t+1})] \leq \sum_{t=0}^{T} \mathbb{E}[F(\boldsymbol{x}^{t})] - \frac{\eta}{2} \sum_{t=0}^{T} \mathbb{E} \left\| \nabla F(\boldsymbol{x}^{t}) \right\|^{2} + \frac{9\eta}{4} \left[\frac{T}{H} \left(3\Upsilon^{2} + \frac{9H^{2}\sigma^{2}}{b} + 33H^{2}\kappa^{2} \right) + \left(T - \frac{T}{H} \right) \left(3H\kappa^{2} + \frac{2H\sigma^{2}}{b} \right) \right]$$

$$(59)$$

We can simplifying the constant term in the RHS of (59) as follows:

$$\begin{split} \frac{1}{H} \left(3\varUpsilon^2 + \frac{9H^2\sigma^2}{b} + 33H^2\kappa^2 \right) + \left(1 - \frac{1}{H} \right) \left(3H\kappa^2 + \frac{2H\sigma^2}{b} \right) \\ & \leq \frac{1}{H} \left(3\varUpsilon^2 + \frac{9H^2\sigma^2}{b} + 33H^2\kappa^2 \right) + \left(3H\kappa^2 + \frac{2H\sigma^2}{b} \right) \\ & \leq \frac{3\varUpsilon^2}{H} + \frac{11H\sigma^2}{b} + 36H\kappa^2 \end{split}$$

Substituting this in (59) and then rearranging, we get:

$$\frac{1}{T} \sum_{t=0}^{T} \mathbb{E} \left\| \nabla F(\boldsymbol{x}^{t}) \right\|^{2} \leq \frac{2}{\eta T} \left[\mathbb{E}[F(\boldsymbol{x}^{0})] - \mathbb{E}[F(\boldsymbol{x}^{T+1})] \right] + \frac{9}{2} \left(\frac{3\Upsilon^{2}}{H} + \frac{11H\sigma^{2}}{b} + 36H\kappa^{2} \right)$$
(60)

Note that the last term in (60) is a constant. So, it would be best to take the step-size η to be as large as possible such that it satisfies $\eta \leq \frac{1}{8HL}$. We take $\eta = \frac{1}{8HL}$. Substituting this in (60) and using $F(\boldsymbol{x}^{T+1}) \geq F(\boldsymbol{x}^*)$ gives

$$\frac{1}{T} \sum_{t=0}^{T} \mathbb{E} \left\| \nabla F(\boldsymbol{x}^{t}) \right\|^{2} \leq \frac{16HL}{T} \left[\mathbb{E}[F(\boldsymbol{x}^{0})] - \mathbb{E}[F(\boldsymbol{x}^{*})] \right] + \frac{9}{2} \left(\frac{3\Upsilon^{2}}{H} + \frac{11H\sigma^{2}}{b} + 36H\kappa^{2} \right), \tag{61}$$

where $\Upsilon^2 = \mathcal{O}\left(\sigma_0^2(\epsilon + \epsilon')\right)$ and $\sigma_0^2 = \frac{25H^2\sigma^2}{b\epsilon'}\left(1 + \frac{3d}{2K}\right) + 28H^2\kappa^2$. Note that (61) is the convergence rate in the non-convex part of Theorem 1.

Error probability analysis. Note that (55) holds with probability at least $1 - \exp\left(-\frac{\epsilon'^2(1-\epsilon)K}{16}\right)$ and (58) holds with probability 1. Since to arrive at (59) (which leads to our final bound (61)), we used (55) $\frac{T}{H}$ times and (58) $\left(T - \frac{T}{H}\right)$ times; as a consequence, by union bound, we have that (61) holds with probability at least $1 - \frac{T}{H} \exp\left(-\frac{\epsilon'^2(1-\epsilon)K}{16}\right)$, which is at least $(1-\delta)$, for any $\delta > 0$, provided we run our algorithm for at most $T \le \delta H \exp\left(\frac{\epsilon'^2(1-\epsilon)K}{16}\right)$ iterations.

This concludes the proof of the non-convex part of Theorem 1.

D. Results on Full-Batch Local Gradient Descent

In this section, we focus on the case when in each local iteration clients compute *full-batch* gradients (instead of computing mini-batch stochastic gradients) in Algorithm 1. Our main result for full-batch gradient descent with local iteration is given below:

Theorem 4 (Full-Batch Local Gradient Descent). In the same setting as that of Theorem 1, except for that we running Algorithm 1 with a fixed step-size $\eta = \frac{1}{5HL}$, and in any iteration, instead of sampling mini-batch stochastic gradients, every honest client takes full-batch gradients from their local datasets. If $\epsilon \leq \frac{1}{3}$, then with probability 1, the sequence of average iterates $\{x^t = \frac{1}{K} \sum_{r \in \mathcal{K}_+} x_t^r : t \in [0:T]\}$ satisfy the following convergence guarantees:

• **Strongly-convex:** If F is L-smooth for $L \ge 0$ and μ -strongly convex for $\mu > 0$, we get:

$$\|\boldsymbol{x}^T - \boldsymbol{x}^*\|^2 \le \left(1 - \frac{\mu}{10HL}\right)^T \|\boldsymbol{x}^0 - \boldsymbol{x}^*\|^2 + \frac{14}{\mu^2} \Gamma_{GD}.$$
 (62)

• Non-convex: If F is L-smooth for $L \ge 0$, we get:

$$\frac{1}{T} \sum_{t=0}^{T} \|\nabla F(\mathbf{x}^{t})\|^{2} \le \frac{10HL}{T} \left[F(\mathbf{x}^{0}) - F(\mathbf{x}^{*}) \right] + \frac{24}{5} \Gamma_{GD}.$$
 (63)

In (62), (63), $\Gamma_{GD} = \frac{2\Upsilon_{GD}^2}{H} + 25H\kappa^2$, where $\Upsilon_{GD} = \mathcal{O}(H\kappa\sqrt{\epsilon})$.

The rest of this section is devoted to proving Theorem 4.

Note that the robust accumulated gradient estimation (RAGE) result of Theorem 2 (which is for stochastic gradients) is one of the main ingredients behind the convergence analyses of Theorem 1. So, in order to prove Theorem 4, first we need to show a RAGE result for full-batch gradients. Note that we can obtain such a result by substituting $\sigma=0$ in both the parts of Theorem 2; however, this would give a loose bound on the approximation error in the second part. In the following, we get a tighter bound (both for RAGE and the convergence rates in Theorem 4) by working directly with full-batch gradients. To get a RAGE result for full-batch gradients, we do a much simplified analysis than what we did before to prove Theorem 2, and the resulting result is stated and proved below in Theorem 5.

Note that, in order to prove Theorem 2, we showed an existence of a subset \mathcal{S} of honest clients (from the set \mathcal{K} of clients who communicate with the server) from whom the accumulated stochastic gradients are well-concentrated, as stated in form of a matrix concentration bound (7) in Theorem 2. It turns out that for full-batch gradients, an analogous result can be proven directly (as there is no randomness due to stochastic gradients); and below we provide such a result. Note that Theorem 2 is a probabilistic statement, where we show that with high probability, there exists a large subset $\mathcal{S} \subseteq \mathcal{K}$ of honest clients whose stochastic accumulated gradients are well-concentrated. In contrast, in the following result, we can deterministically take the set of *all* honest clients in \mathcal{K} to be that subset for which we can directly show the concentration.

First we setup the notation to state our main result on RAGE for full-batch gradients. Let $\mathcal{K}_t\subseteq [R]$ denote the subset of clients of size K that are active at any time $t\in [0:T]$. Let Algorithm 1 generate a sequence of iterates $\{\boldsymbol{x}_r^t:t\in [0:T],r\in\mathcal{K}_t\}$ when run with a fixed step-size η satisfying $\eta\leq \frac{1}{5HL}$ while minimizing a global objective function $F:\mathbb{R}^d\to\mathbb{R}$, where in any iteration, instead of sampling mini-batch stochastic gradients, every honest client takes full-batch gradients from their local datasets. Take any two consecutive synchronization indices $t_k,t_{k+1}\in\mathcal{I}_T$. Note that $|t_{k+1}-t_k|\leq H$. For an honest client $r\in\mathcal{K}_{t_k}$, let $\nabla F_{r,\mathrm{accu}}^{t_k,t_{k+1}}:=\sum_{t=t_k}^{t_{k+1}-1}\nabla F_r(\boldsymbol{x}_r^t)$ denote the sum of local full-batch gradients taken by client r between time t_k and t_{k+1} . Note that at iteration t_{k+1} , every honest client $r\in\mathcal{K}_{t_k}$ reports its local parameters $\boldsymbol{x}_r^{t_{k+1}}$ to the server, from which server can compute $\nabla F_{r,\mathrm{accu}}^{t_k,t_{k+1}}$, whereas, corrupt clients may report arbitrary and adversarially chosen vectors in \mathbb{R}^d . The goal of the server is to produce an estimate $\nabla \widehat{F}_{\mathrm{accu}}^{t_k,t_{k+1}}$ of the average accumulated gradients from honest clients as best as possible.

Theorem 5 (Robust Accumulated Gradient Estimation for Full-Batch Gradient Descent). Suppose an ϵ fraction of clients who communicate with the server are corrupt. In the setting and notation described above, suppose we are given $K \leq R$ accumulated full-batch gradients $\nabla \widetilde{F}_{r,\text{accu}}^{t_k,t_{k+1}}$, $r \in \mathcal{K}_{t_k}$ in \mathbb{R}^d , where $\nabla \widetilde{F}_{r,\text{accu}}^{t_k,t_{k+1}} = \nabla F_{r,\text{accu}}^{t_k,t_{k+1}}$ if the r'th client is honest, otherwise can be arbitrary. Let $S \subseteq \mathcal{K}_{t_k}$ be the subset of all honest clients in \mathcal{K}_{t_k} and $\nabla F_{S,\text{accu}}^{t_k,t_{k+1}} := \frac{1}{|S|} \sum_{i \in S} \nabla F_{i,\text{accu}}^{t_k,t_{k+1}}$ be the sample average of uncorrupted full-batch gradients. If $\epsilon \leq \frac{1}{3}$, then with probability 1, we can find an estimate $\nabla \widehat{F}_{\text{accu}}^{t_k,t_{k+1}}$ of $\nabla F_{S,\text{accu}}^{t_k,t_{k+1}}$ in polynomial-time, such that $\left\|\nabla \widehat{F}_{\text{accu}}^{t_k,t_{k+1}} - \nabla F_{S,\text{accu}}^{t_k,t_{k+1}}\right\| \leq \mathcal{O}(H\kappa\sqrt{\epsilon})$.

Proof. First we prove that

$$\lambda_{\max} \left(\frac{1}{|\mathcal{S}|} \sum_{i \in \mathcal{S}} \left(\nabla F_{i, \text{accu}}^{t_k, t_{k+1}} - \nabla F_{\mathcal{S}, \text{accu}}^{t_k, t_{k+1}} \right) \left(\nabla F_{i, \text{accu}}^{t_k, t_{k+1}} - \nabla F_{\mathcal{S}, \text{accu}}^{t_k, t_{k+1}} \right)^T \right) \le 11 H^2 \kappa^2.$$
 (64)

In view of the alternate characterization the largest eigenvalue given in (15), this is equivalent to showing

$$\sup_{\boldsymbol{v} \in \mathbb{R}^d: \|\boldsymbol{v}\| = 1} \frac{1}{|\mathcal{S}|} \sum_{i \in \mathcal{S}} \left\langle \nabla F_{i, \text{accu}}^{t_k, t_{k+1}} - \nabla F_{\mathcal{S}, \text{accu}}^{t_k, t_{k+1}}, \boldsymbol{v} \right\rangle^2 \le 11 H^2 \kappa^2, \tag{65}$$

which we prove below. Define $F_{\text{accu}}^{t_k,t_{k+1}} := \sum_{t=t_k}^{t_{k+1}-1} F(\boldsymbol{x}^t)$, where $\boldsymbol{x}^t = \frac{1}{K} \sum_{r \in \mathcal{K}_{t_k}} \boldsymbol{x}_r^t$ for any $t \in [t_k : t_{k+1}-1]$. Take an arbitrary unit vector $\boldsymbol{v} \in \mathbb{R}^d$.

$$\begin{split} &\frac{1}{|\mathcal{S}|} \sum_{i \in \mathcal{S}} \left\langle \nabla F_{i, \text{accu}}^{t_{k, t, k_{k+1}}} - \nabla F_{\mathcal{S}, \text{accu}}^{t_{k, t_{k+1}}}, \mathbf{v} \right\rangle^{2} \\ &= \frac{1}{|\mathcal{S}|} \sum_{i \in \mathcal{S}} \left[\left\langle \nabla F_{i, \text{accu}}^{t_{k, t_{k+1}}} - \nabla F_{\text{accu}}^{t_{k, t_{k+1}}} + \nabla F_{\text{accu}}^{t_{k, t_{k+1}}} - \nabla F_{\mathcal{S}, \text{accu}}^{t_{k, t_{k+1}}}, \mathbf{v} \right\rangle^{2} \\ &\leq \frac{2}{|\mathcal{S}|} \sum_{i \in \mathcal{S}} \left\langle \nabla F_{i, \text{accu}}^{t_{k, t_{k+1}}} - \nabla F_{\text{accu}}^{t_{k, t_{k+1}}}, \mathbf{v} \right\rangle^{2} + \frac{2}{|\mathcal{S}|} \sum_{i \in \mathcal{S}} \left\langle \nabla F_{\mathcal{S}, \text{accu}}^{t_{k, t_{k+1}}}, \mathbf{v} \right\rangle^{2} \\ &= \frac{2}{|\mathcal{S}|} \sum_{i \in \mathcal{S}} \left\langle \nabla F_{i, \text{accu}}^{t_{k, t_{k+1}}} - \nabla F_{\text{accu}}^{t_{k, t_{k+1}}}, \mathbf{v} \right\rangle^{2} + 2 \left\langle \nabla F_{\mathcal{S}, \text{accu}}^{t_{k, t_{k+1}}} - \nabla F_{\text{accu}}^{t_{k, t_{k+1}}}, \mathbf{v} \right\rangle^{2} \\ &= \frac{2}{|\mathcal{S}|} \sum_{i \in \mathcal{S}} \left\langle \nabla F_{i, \text{accu}}^{t_{k, t_{k+1}}} - \nabla F_{\text{accu}}^{t_{k, t_{k+1}}}, \mathbf{v} \right\rangle^{2} + 2 \left\langle \nabla F_{\mathcal{S}, \text{accu}}^{t_{k, t_{k+1}}} - \nabla F_{\text{accu}}^{t_{k, t_{k+1}}}, \mathbf{v} \right\rangle^{2} \\ &= \frac{2}{|\mathcal{S}|} \sum_{i \in \mathcal{S}} \left\langle \nabla F_{i, \text{accu}}^{t_{k, t_{k+1}}} - \nabla F_{\text{accu}}^{t_{k, t_{k+1}}}, \mathbf{v} \right\rangle^{2} + 2 \left\langle \nabla F_{\mathcal{S}, \text{accu}}^{t_{k, t_{k+1}}} - \nabla F_{\text{accu}}^{t_{k, t_{k+1}}}, \mathbf{v} \right\rangle^{2} \\ &= \frac{2}{|\mathcal{S}|} \sum_{i \in \mathcal{S}} \left\langle \nabla F_{i, \text{accu}}^{t_{k, t_{k+1}}} - \nabla F_{\text{accu}}^{t_{k, t_{k+1}}}, \mathbf{v} \right\rangle^{2} + 2 \left\langle \nabla F_{i, \text{accu}}^{t_{k, t_{k+1}}} - \nabla F_{\text{accu}}^{t_{k, t_{k+1}}}, \mathbf{v} \right\rangle^{2} \\ &\leq \frac{2}{|\mathcal{S}|} \sum_{i \in \mathcal{S}} \left\langle \nabla F_{i, \text{accu}}^{t_{k, t_{k+1}}} - \nabla F_{\text{accu}}^{t_{k, t_{k+1}}}, \mathbf{v} \right\rangle^{2} + \frac{2}{|\mathcal{S}|} \sum_{i \in \mathcal{S}} \left\langle \nabla F_{i, \text{accu}}^{t_{k, t_{k+1}}}, \mathbf{v} \right\rangle^{2} \\ &= \frac{4}{|\mathcal{S}|} \sum_{i \in \mathcal{S}} \left\langle \nabla F_{i, \text{accu}}^{t_{k, t_{k+1}}} - \nabla F_{\text{accu}}^{t_{k, t_{k+1}}}, \mathbf{v} \right\rangle^{2} \\ &\leq \frac{4}{|\mathcal{S}|} \sum_{i \in \mathcal{S}} \left\| \nabla F_{i, \text{accu}}^{t_{k, t_{k+1}}} - \nabla F_{\text{accu}}^{t_{k, t_{k+1}}}, \mathbf{v} \right\rangle^{2} \\ &\leq \frac{4}{|\mathcal{S}|} \sum_{i \in \mathcal{S}} \left\| \nabla F_{i, \text{accu}}^{t_{k, t_{k+1}}} - \nabla F_{\text{accu}}^{t_{k, t_{k+1}}}, \mathbf{v} \right\|^{2} \right\} \\ &\leq \frac{4}{|\mathcal{S}|} \sum_{i \in \mathcal{S}} \left\| \nabla F_{i, \text{accu}}^{t_{k, t_{k+1}}} - \nabla F_{\text{accu}}^{t_{k, t_{k+1}}}, \mathbf{v} \right\|^{2} \right\|^{2} \\ &\leq \frac{4}{|\mathcal{S}|} \sum_{i \in \mathcal{S}} \left\| \nabla F_{i, \text{accu}}^{t_{k, t_{k+1}}}$$

The last inequality follows from the Jensen's inequality. In (a) we used (6) to bound $\|\nabla F_i(\boldsymbol{x}_i^t) - \nabla F(\boldsymbol{x}_i^t)\|^2 \le \kappa^2$ and L-Lipschitz gradient property of F to bound $\|\nabla F(\boldsymbol{x}_i^t) - \nabla F(\boldsymbol{x}^t)\| \le L\|\boldsymbol{x}_i^t - \boldsymbol{x}^t\|$.

Now we bound the last term of (66).

Lemma 4. For any $r, s \in \mathcal{K}_{t_k}$, if $\eta \leq \frac{1}{5HL}$, we have

$$\sum_{t=t_k}^{t_{k+1}-1} \|\boldsymbol{x}_r^t - \boldsymbol{x}_s^t\|^2 \le 7\eta^2 H^3 \kappa^2.$$
 (67)

Proof. Note that we have shown a similar result (but, in expectation) in Lemma 3 (on page 14), which is for stochastic gradients. We will simplify that proof to prove Lemma 4, which is for full-batch deterministic gradients.

Take an arbitrary $t \in [t_k : t_{k+1} - 1]$. Following the proof of Lemma 3 until (24) and removing the factor of 3 inside the summation (the factor of 3 appeared because we applied the Jensen's inequality earlier to separate the deterministic gradient term and the stochastic gradient terms) would give

$$\left\|\boldsymbol{x}_{r}^{t} - \boldsymbol{x}_{s}^{t}\right\|^{2} \leq \eta^{2} H \sum_{j=t_{k}}^{t-1} \left\|\nabla F_{r}(\boldsymbol{x}_{r}^{j}) - \nabla F_{s}(\boldsymbol{x}_{s}^{j})\right\|^{2}.$$
(68)

Following the remaining proof of Lemma 3 from (24) until the end and substituting $\sigma = 0$ gives the desired result.

Substituting the bound from (67) into (66) gives

$$\begin{split} \frac{1}{|\mathcal{S}|} \sum_{i \in \mathcal{S}} \left\langle \nabla F_{i,\text{accu}}^{t_k,t_{k+1}} - \nabla F_{\mathcal{S},\text{accu}}^{t_k,t_{k+1}}, \boldsymbol{v} \right\rangle^2 &\leq 8H^2 \kappa^2 + 56H^4 L^2 \eta^2 \kappa^2 \\ &\leq 8H^2 \kappa^2 + \frac{56}{25} H^2 \kappa^2 \\ &\leq 11H^2 \kappa^2. \end{split} \tag{Substituting } \eta \leq \frac{1}{5HL}) \end{split}$$

Note that (69) holds for an arbitrary unit vector $v \in \mathbb{R}^d$, implying that (65) holds true. Since (65) and (64) are equivalent, we have thus shown (64).

Now apply Theorem 3 with \mathcal{S} being the set of all honest clients, and $\boldsymbol{g}_{i,\mathrm{accu}}^{t_k,t_{k+1}} = \nabla F_{i,\mathrm{accu}}^{t_k,t_{k+1}}, \ \boldsymbol{g}_{\mathcal{S},\mathrm{accu}}^{t_k,t_{k+1}} = \nabla F_{\mathcal{S},\mathrm{accu}}^{t_k,t_{k+1}}$ of $\nabla F_{\mathcal{S},\mathrm{accu}}^{t_k,t_{k+1}}$ of $\nabla F_{\mathcal{S},\mathrm{accu}}^{t_k,t_{k+1}}$ in polynomial-time, such that $\left\| \nabla \widehat{F}_{\mathrm{accu}}^{t_k,t_{k+1}} - \nabla F_{\mathcal{S},\mathrm{accu}}^{t_k,t_{k+1}} \right\| \leq \mathcal{O}\left(H\kappa\sqrt{\epsilon}\right)$ holds with probability 1.

Theorem 4 can be proved with appropriate modifications in the proof of Theorem 1, and for completeness, we prove it below.

D.1. Convergence Proof of the Strongly-Convex Part of Theorem 4

Let $\mathcal{K}_t \subseteq [R]$ denote the subset of clients of size $|\mathcal{K}_t| = K$ that are active at the t'th iteration. For any $t \in [t_i : t_{i+1} - 1]$, let $x^t = \frac{1}{K} \sum_{k \in \mathcal{K}_t} x_k^t$ denote the average of the local parameters of clients in the sampling set \mathcal{K}_{t_i} .

Following the proof of the strongly-convex part of Theorem 1 given in Appendix B until (31) gives

$$\|\boldsymbol{x}^{t_{i+1}} - \boldsymbol{x}^*\|^2 \le \left(1 + \frac{\mu\eta}{2}\right) \|\boldsymbol{x}^{t_{i+1}-1} - \eta\nabla F(\boldsymbol{x}^{t_{i+1}-1}) - \boldsymbol{x}^*\|^2$$

$$+ 2\eta \left(\eta + \frac{2}{\mu}\right) \left\|\frac{1}{K} \sum_{r \in \mathcal{K}_{t_i}} \left(\nabla F(\boldsymbol{x}^{t_{i+1}-1}) - \nabla F_r(\boldsymbol{x}_r^{t_{i+1}-1})\right)\right\|^2$$

$$+ 2\eta \left(\eta + \frac{2}{\mu}\right) \left\|\widehat{F}_{\text{accu}}^{t_i, t_{i+1}} - \frac{1}{K} \sum_{r \in \mathcal{K}_{t_i}} \sum_{t=t_i}^{t_{i+1}-1} \nabla F_r(\boldsymbol{x}_r^t)\right\|^2$$

$$(70)$$

We have already bounded the first term in Claim 3 (on page 17) by

$$\|\boldsymbol{x}^{t_{i+1}} - \eta \nabla F(\boldsymbol{x}^{t_{i+1}-1}) - \boldsymbol{x}^*\|^2 \le (1 - \eta \mu) \|\boldsymbol{x}^{t_{i+1}-1} - \boldsymbol{x}^*\|^2.$$
 (71)

In order to bound the second term, we follow the proof of Claim 4 exactly until (36), and then to bound $\|x_r^{t_{i+1}-1} - x_s^{t_{i+1}-1}\|^2$ for every $r, s \in \mathcal{K}_{t_i}$, we use the bound from (67) in Lemma 4 and use $\eta \leq \frac{1}{5HL}$, which gives

$$\left\| \frac{1}{K} \sum_{r \in \mathcal{K}_{t_i}} \left(\nabla F_r(\boldsymbol{x}^{t_{i+1}-1}) - \nabla F_r(\boldsymbol{x}_r^{t_{i+1}-1}) \right) \right\|^2 \le 3H\kappa^2. \tag{72}$$

To bound the third term in the RHS of (70), we can simplify the proof of Claim 5: Firstly, note that with full-batch gradients, the variance σ^2 becomes zero; secondly, as shown in Theorem 5, the robust estimation of accumulated gradients holds with probability 1. Following the proof of Claim 5 with these changes and using $\eta \leq \frac{1}{5HL}$, we get

$$\left\| \widehat{F}_{\text{accu}}^{t_i, t_{i+1}} - \frac{1}{K} \sum_{r \in \mathcal{K}_{t_i}} \sum_{t=t_i}^{t_{i+1}-1} \nabla F_r(\boldsymbol{x}_r^t) \right\|^2 \le 2\Upsilon_{\text{GD}}^2 + 20H^2\kappa^2, \tag{73}$$

where $\Upsilon_{\rm GD}=\mathcal{O}\left(H\kappa\sqrt{\epsilon}\right)$. Substituting all these bounds from (71)-(73) into (70) and simplifying further using $\left(1+\frac{\mu\eta}{2}\right)\left(1-\mu\eta\right)\leq \left(1-\frac{\mu\eta}{2}\right)$ and $\left(\eta+\frac{2}{\mu}\right)\leq \frac{3}{\mu}$ gives

$$\|\boldsymbol{x}^{t_{i+1}} - \boldsymbol{x}^*\|^2 \le \left(1 - \frac{\mu\eta}{2}\right) \|\boldsymbol{x}^{t_{i+1}-1} - \boldsymbol{x}^*\|^2 + \frac{6\eta}{\mu} \left(2\Upsilon_{GD}^2 + 23H^2\kappa^2\right)$$
 (74)

Note that (74) gives a recurrence at the synchronization indices. Now we give a recurrence at non-synchronization indices. Take an arbitrary $t \in [T]$ and let $t_i \in \mathcal{I}_T$ be such that $t \in [t_i : t_{i+1} - 1]$; when $H \ge 2$, such t's exist. Following the steps that we used to arrive at (45), we get the following (note that the last term on the RHS of (45) is zero, as $g_r(\boldsymbol{x}_r^t) = \nabla F_r(\boldsymbol{x}_r^t)$ holds for every $r \in [R]$ and $t \in [T]$; this will also save us the factor of 2 in the previous term as we don't have to use the Jensen's inequality to get to (45)):

$$\left\| \boldsymbol{x}^{t+1} - \boldsymbol{x}^* \right\|^2 \le \left(1 + \frac{\mu \eta}{2} \right) \left\| \boldsymbol{x}^t - \boldsymbol{x}^* - \eta \nabla F(\boldsymbol{x}^t) \right\|^2 + \eta \left(\eta + \frac{2}{\mu} \right) \left\| \frac{1}{K} \sum_{r \in \mathcal{K}} \left(\nabla F(\boldsymbol{x}^t) - \nabla F_r(\boldsymbol{x}_r^t) \right) \right\|^2$$
(75)

Substituting the bounds from (71) and (72) into (75) and simplifying the coefficients as above, we get

$$\|\boldsymbol{x}^{t+1} - \boldsymbol{x}^*\|^2 \le \left(1 - \frac{\mu\eta}{2}\right) \|\boldsymbol{x}^t - \boldsymbol{x}^*\|^2 + \frac{3\eta}{\mu} (3H\kappa^2)$$
 (76)

Now we have a recurrence at the synchronization indices given in (74) and at non-synchronization indices given in (76). Let $\alpha = \left(1 - \frac{\mu\eta}{2}\right)$, $\beta_1 = \left(2\Upsilon_{GD}^2 + 23H^2\kappa^2\right)$, and $\beta_2 = \left(\frac{3}{2}H\kappa^2\right)$. Following the same steps that we used to arrive at (48) gives:

$$\|\boldsymbol{x}^{T} - \boldsymbol{x}^{*}\|^{2} \le \alpha^{T} \|\boldsymbol{x}^{0} - \boldsymbol{x}^{*}\|^{2} + \frac{6\eta}{\mu} \left(\frac{1}{1 - \alpha}\beta_{2} + \frac{1}{1 - \alpha^{H}}\beta_{1}\right)$$
 (77)

Since $\alpha=\left(1-\frac{\mu\eta}{2}\right)$, we have $\alpha^H=\left(1-\frac{\mu\eta}{2}\right)^H\overset{\text{(a)}}{\leq}\exp(-\frac{\mu\eta H}{2})\overset{\text{(b)}}{\leq}1-\frac{\mu\eta H}{2}+\left(\frac{\mu\eta H}{2}\right)^2\overset{\text{(c)}}{\leq}1-\frac{\mu\eta H}{2}+\frac{1}{10}\frac{\mu\eta H}{2}=1-\frac{9}{10}\frac{\mu\eta H}{2}$. In (a) we used the inequality $(1-\frac{1}{x})^x\leq\frac{1}{e}$ which holds for any x>0; in (b) we used $\exp(-x)\leq 1-x+x^2$ which holds for any $x\geq0$; in (c) we used $\eta\leq\frac{1}{5HL}$ and $\mu\leq L$, which imply $\frac{\mu\eta H}{2}\leq\frac{1}{10}$. Substituting these in (77) gives

$$\|\boldsymbol{x}^{T} - \boldsymbol{x}^{*}\|^{2} \leq \left(1 - \frac{\mu\eta}{2}\right)^{T} \|\boldsymbol{x}^{0} - \boldsymbol{x}^{*}\|^{2} + \frac{6\eta}{\mu} \left(\frac{2}{\mu\eta}\beta_{2} + \frac{20}{9\mu\eta H}\beta_{1}\right)$$

$$\leq \left(1 - \frac{\mu\eta}{2}\right)^{T} \|\boldsymbol{x}^{0} - \boldsymbol{x}^{*}\|^{2} + \frac{6 \times 20}{9\mu^{2}} \left(\frac{9}{10}\beta_{2} + \frac{1}{H}\beta_{1}\right)$$

$$\leq \left(1 - \frac{\mu\eta}{2}\right)^{T} \|\boldsymbol{x}^{0} - \boldsymbol{x}^{*}\|^{2} + \frac{14}{\mu^{2}} \left(\frac{2\Upsilon_{GD}^{2}}{H} + 25H\kappa^{2}\right), \tag{78}$$

where $\Upsilon_{GD} = \mathcal{O}(H\kappa\sqrt{\epsilon})$. Substituting the value of $\eta = \frac{1}{5HL}$ yields the convergence rate (62) in the strongly-convex part of Theorem 4. Note that (78) holds with probability 1.

D.2. Convergence Proof of the Non-Convex Part of Theorem 4

Following the proof of the non-convex part of Theorem 1 given in Appendix C until (53) and using $\eta \leq \frac{1}{5HL}$ gives:

$$F(\boldsymbol{x}^{t_{i+1}}) \le F(\boldsymbol{x}^{t_{i+1}-1}) - \frac{\eta}{2} \left\| \nabla F(\boldsymbol{x}^{t_{i+1}-1}) \right\|^2 + \frac{6\eta}{5} \|C\|^2, \tag{79}$$

where
$$C = \frac{1}{K} \sum_{r \in \mathcal{K}_{t_i}} \left(\nabla F(\boldsymbol{x}^{t_{i+1}-1}) - \nabla F_r(\boldsymbol{x}^{t_{i+1}-1}_r) \right) - \left(\widehat{F}_{\text{accu}}^{t_i, t_{i+1}} - \frac{1}{K} \sum_{r \in \mathcal{K}_{t_i}} \sum_{t=t_i}^{t_{i+1}-1} \nabla F_r(\boldsymbol{x}^t_r) \right).$$

Using the bounds from (72) and (73), together with the Jensen's inequality, we can bound $||C||^2$ as follows:

$$||C||^2 \le 2(3H\kappa^2) + 2(2\Upsilon_{GD}^2 + 20H^2\kappa^2) \le 2(2\Upsilon_{GD}^2 + 23H^2\kappa^2)$$
(80)

Substituting the bound from (80) into (79) gives:

$$F(\boldsymbol{x}^{t_{i+1}}) \le F(\boldsymbol{x}^{t_{i+1}-1}) - \frac{\eta}{2} \left\| \nabla F(\boldsymbol{x}^{t_{i+1}-1}) \right\|^2 + \frac{12\eta}{5} \left(2\Upsilon_{GD}^2 + 23H^2\kappa^2 \right), \tag{81}$$

where $\Upsilon_{GD} = \mathcal{O}(H\kappa\sqrt{\epsilon})$.

Note that above recurrence in (81) holds only at the synchronization indices. Now we give a recurrence at non-synchronization indices.

We have done a similar calculations in the non-convex part of Theorem 1 in Appendix C.

Take an arbitrary $t \in [T]$ and let $t_i \in \mathcal{I}_T$ be such that $t \in [t_i : t_{i+1} - 1]$; when $H \ge 2$, such t's exist. Following the same steps until (56) and using $\eta \le \frac{1}{5HL}$ gives:

$$F(\boldsymbol{x}^{t+1}) \le F(\boldsymbol{x}^t) - \frac{\eta}{2} \|\nabla F(\boldsymbol{x}^t)\|^2 + \frac{6\eta}{5} \|D\|^2,$$
 (82)

where $D = \frac{1}{K} \sum_{r \in \mathcal{K}_{t_i}} (\nabla F(\boldsymbol{x}^t) - \nabla F_r(\boldsymbol{x}_r^t)).$

Using the bound from (72), we have $||D||^2 < 3H\kappa^2$. Substituting this in (82) gives:

$$F(\boldsymbol{x}^{t+1}) \le F(\boldsymbol{x}^t) - \frac{\eta}{2} \left\| \nabla F(\boldsymbol{x}^t) \right\|^2 + \frac{6\eta}{5} (3H\kappa^2)$$
(83)

Now we have a recurrence at the synchronization indices given in (81) and at non-synchronization indices given in (83). Adding (81) and (83) from t = 0 to T (use (81) for the synchronization indices and (83) for the rest of the indices) gives:

$$\sum_{t=0}^{T} F(\boldsymbol{x}^{t+1}) \le \sum_{t=0}^{T} F(\boldsymbol{x}^{t}) - \frac{\eta}{2} \sum_{t=0}^{T} \left\| \nabla F(\boldsymbol{x}^{t}) \right\|^{2} + \frac{12\eta}{5} \left[\frac{T}{H} \left(2\Upsilon_{\text{GD}}^{2} + 23H^{2}\kappa^{2} \right) + \left(T - \frac{T}{H} \right) \left(\frac{3}{2}H\kappa^{2} \right) \right]$$
(84)

After rearranging and simplifying the last constant terms, we get:

$$\frac{1}{T} \sum_{t=0}^{T} \left\| \nabla F(\boldsymbol{x}^{t}) \right\|^{2} \leq \frac{2}{\eta T} \left[F(\boldsymbol{x}^{0}) - F(\boldsymbol{x}^{T+1}) \right] + \frac{24}{5} \left(\frac{2 \Upsilon_{\text{GD}}^{2}}{H} + 25 H \kappa^{2} \right)$$
(85)

Note that the last term in (85) is a constant. So, it would be best to take the step-size η to be as large as possible such that it satisfies $\eta \leq \frac{1}{5HL}$. We take $\eta = \frac{1}{5HL}$. Substituting this in (85) and using $F(\boldsymbol{x}^{T+1}) \geq F(\boldsymbol{x}^*)$ gives

$$\frac{1}{T} \sum_{t=0}^{T} \left\| \nabla F(\boldsymbol{x}^{t}) \right\|^{2} \leq \frac{10HL}{T} \left[F(\boldsymbol{x}^{0}) - F(\boldsymbol{x}^{*}) \right] + \frac{24}{5} \left(\frac{2\Upsilon_{\text{GD}}^{2}}{H} + 25H\kappa^{2} \right), \tag{86}$$

where $\Upsilon_{GD} = \mathcal{O}(H\kappa\sqrt{\epsilon})$. This yields the convergence rate (63) in the non-convex part of Theorem 4. Note that (86) holds with probability 1.

This concludes the proof of Theorem 4.

E. Bounding Local Variances and Gradient Dissimilarity in the Statistical Heterogeneous Model

In this section, we bound the gradient dissimilarity κ^2 (from (6)) and local variance σ^2 (from (2)) in the statistical model in heterogeneous setting, where different clients may have local data generated from potentially different distributions. The purpose of this section is to provide upper bounds on κ and σ in the statistical model.

Let q_1,q_2,\ldots,q_R denote the R probability distributions from which the local data samples at the clients are drawn. Specifically, the data samples at any client r are drawn from q_r in an i.i.d. fashion and independently from other clients. For $r \in [R]$, let \mathcal{Q}_r denote the alphabet over which q_r is distributed. For $r \in [R]$, let $f_r : \mathcal{Q}_r \times \mathcal{C} \to \mathbb{R}$ denote the local loss function at client r, where $f_r(\boldsymbol{z}, \boldsymbol{x})$ is the loss associated with the sample $\boldsymbol{z} \in \mathcal{Q}_r$ w.r.t. the model parameters $\boldsymbol{x} \in \mathcal{C}$ and $\mathcal{C} \subset \mathbb{R}^d$ is a bounded subset of \mathbb{R}^d . Linear regression is a classic example of this, where, if $\boldsymbol{z} = (\boldsymbol{w}, y)$ denote the pair of a feature vector $\boldsymbol{w} \in \mathbb{R}^d$ and the response $\boldsymbol{y} \in \mathbb{R}$, then $f_r(\boldsymbol{z}, \boldsymbol{x}) = \frac{1}{2}(\langle \boldsymbol{w}, \boldsymbol{x} \rangle - \boldsymbol{y})^2$. For each client $r \in [R]$, we assume that for any fixed $\boldsymbol{z} \in \mathcal{Q}_r$, the local loss function $f_r(\boldsymbol{z}, \boldsymbol{x})$ is L-smooth w.r.t. \boldsymbol{x} , i.e., for any $\boldsymbol{z} \in \mathcal{Q}_r$, we have $\|\nabla f_r(\boldsymbol{z}, \boldsymbol{x}) - \nabla f_r(\boldsymbol{z}, \boldsymbol{y})\| \le L\|\boldsymbol{x} - \boldsymbol{y}\|, \forall \boldsymbol{x}, \boldsymbol{y} \in \mathcal{C}$.

Let $\mu_r(\boldsymbol{x}) := \mathbb{E}_{\boldsymbol{z} \sim q_r}[f_r(\boldsymbol{z}, \boldsymbol{x})]$ denote the expected value of $f_r(\boldsymbol{z}, \boldsymbol{x})$, when \boldsymbol{z} is sampled from \mathcal{Q}_r according to q_r . For any $\boldsymbol{x} \in \mathcal{C}$, let $\mu(\boldsymbol{x}) := \frac{1}{R} \sum_{r=1}^R \mu_r(\boldsymbol{x})$ denote the average value of $\mu_r(\boldsymbol{x})$, $r \in [R]$.

We are given n_r i.i.d. samples $z_{r,1}, z_{r,2}, \ldots, z_{r,n_r}$ at the r'th client from q_r . Fix an arbitrary parameter vector $\boldsymbol{x} \in \mathcal{C}$. Let $\bar{f}_r(\boldsymbol{x}) := \frac{1}{n_r} \sum_{i=1}^{n_r} f_r(z_{r,i}, \boldsymbol{x})$ denote the average loss at client r on the n_r samples $z_{r,1}, \ldots, z_{r,n_r}$ w.r.t. \boldsymbol{x} . Let $\bar{f}(\boldsymbol{x}) := \frac{1}{R} \sum_{r=1}^{r} \bar{f}_r(\boldsymbol{x})$ denote the average loss across all clients. The analogues of (6) and (2) in this statistical heterogeneous model are the following:

$$\|\nabla \bar{f}_r(\boldsymbol{x}) - \nabla \bar{f}(\boldsymbol{x})\|^2 \le \kappa^2, \quad \forall \boldsymbol{x} \in \mathcal{C},$$
 (87)

$$\mathbb{E}_{i \in U[n_r]} \left\| \nabla f_r(\boldsymbol{z}_{r,i}, \boldsymbol{x}) - \nabla \bar{f}_r(\boldsymbol{x}) \right\|^2 \le \sigma^2, \quad \forall \boldsymbol{x} \in \mathcal{C}.$$
(88)

We need to find good upper bounds on κ and σ that hold for all $r \in [R]$, $x \in \mathcal{C}$ with high probability. We provide two bounds on κ , one when the local gradients at clients are assumed to be sub-exponential random vectors, and other when they are sub-Gaussian random vectors. We provide a bound on σ assuming that the local gradients are sub-Gaussian random vectors. These are standard assumptions on gradients in statistical models, where data at all clients are sampled from the *same* distribution in an i.i.d. fashion (Chen et al., 2017; Su & Xu, 2019; Yin et al., 2019), which is in contrast to our heterogeneous data setting, where data at different clients may be sampled from *different* distributions. Note that these works minimize the *population risk* with *full batch* gradient descent, whereas, we minimize the *empirical risk* with *stochastic* gradient descent. In particular, (Chen et al., 2017; Su & Xu, 2019) make sub-exponential gradient assumption and give convergence guarantees only for strong-convex objectives. On the other hand, (Yin et al., 2019) gives convergence guarantees for non-convex objectives, but under a stricter condition of sub-Gaussian distribution on gradients. In this paper, we provide convergence guarantees for both strongly-convex and non-convex objectives. Moreover, as opposed to (Chen et al., 2017; Su & Xu, 2019; Yin et al., 2019), our results are in a more general heterogeneous data model. Note that we need sub-Gaussian assumption only to bound the variance, which occurs because clients sample stochastic gradients. In case of full batch gradient descent, we only need sub-exponential assumption, as the variance is zero.

Now we state the distributional assumptions on local gradients. We defer the definitions of sub-exponential/sub-Gaussian random variables/vectors and their concentration inequalities that we will use in this section to Section E.3.

Assumption 3 (Sub-exponential local gradients). For every $x \in C$, the local gradient vectors at any client $r \in [R]$ are sub-exponential random vectors, i.e., there exist non-negative parameters (ν, α) such that

$$\sup_{\boldsymbol{v} \in \mathbb{R}^d: ||\boldsymbol{v}|| = 1} \mathbb{E}_{\boldsymbol{z} \sim q_r} \left[\exp\left(\lambda \left\langle \nabla f_r(\boldsymbol{z}, \boldsymbol{x}) - \nabla \mu_r(\boldsymbol{x}), \boldsymbol{v} \right\rangle \right) \right] \le \exp\left(\lambda^2 \nu^2 / 2\right), \quad \forall |\lambda| < \frac{1}{\alpha}.$$
 (89)

Assumption 4 (Sub-Gaussian local gradients). For every $x \in C$, the local gradient vectors at any client $r \in [R]$ are sub-Gaussian random vectors, i.e., there exists a non-negative parameter σ_{σ} such that

$$\sup_{\boldsymbol{v} \in \mathbb{R}^d: \|\boldsymbol{v}\| = 1} \mathbb{E}_{\boldsymbol{z} \sim q_r} \left[\exp \left(\lambda \left\langle \nabla f_r(\boldsymbol{z}, \boldsymbol{x}) - \nabla \mu_r(\boldsymbol{x}), \boldsymbol{v} \right\rangle \right) \right] \le \exp \left(\lambda^2 \sigma_g^2 / 2 \right), \quad \forall \lambda \in \mathbb{R}.$$
 (90)

Though, as stated above in both the assumptions, local gradients at all clients have the same parameters $((\nu, \alpha))$ for sub-exponential and σ_g for sub-Gaussian), this is without loss of generality. In case they have different parameters

 $((\nu_r, \alpha_r), r \in [R]$ for sub-exponential and $\sigma_r, r \in [R]$ for sub-Gaussian), we can take the final parameters to be the maximum of the respective local parameters – for sub-exponential, we can take $\nu = \max_{r \in [R]} \nu_r$ and $\alpha = \max_{r \in [R]} \alpha_r$, and for sub-Gaussian, we can take $\sigma_g = \max_{r \in [R]} \sigma_r$.

E.1. Bounding the gradient dissimilarity κ

In this section, we provide an upper bound on $\|\nabla \bar{f}_r(x) - \nabla \bar{f}(x)\|$.

$$\|\nabla \bar{f}_r(\boldsymbol{x}) - \nabla \bar{f}(\boldsymbol{x})\| \le \|\nabla \bar{f}_r(\boldsymbol{x}) - \nabla \mu_r(\boldsymbol{x})\| + \|\nabla \mu_r(\boldsymbol{x}) - \nabla \mu(\boldsymbol{x})\| + \|\nabla \bar{f}(\boldsymbol{x}) - \nabla \mu(\boldsymbol{x})\|$$

$$\le \|\nabla \bar{f}_r(\boldsymbol{x}) - \nabla \mu_r(\boldsymbol{x})\| + \|\nabla \mu_r(\boldsymbol{x}) - \nabla \mu(\boldsymbol{x})\| + \frac{1}{R} \sum_{r=1}^{R} \|\nabla \bar{f}_r(\boldsymbol{x}) - \nabla \mu_r(\boldsymbol{x})\|,$$
(91)

where for the third term, we used $\bar{f}(\boldsymbol{x}) = \frac{1}{R} \sum_{r=1}^{R} \bar{f}_r(\boldsymbol{x})$ and $\mu(\boldsymbol{x}) = \frac{1}{R} \sum_{r=1}^{R} \mu_r(\boldsymbol{x})$, and applied the triangle inequality. It follows from (91) that in order to bound $\|\nabla \bar{f}_r(\boldsymbol{x}) - \nabla \bar{f}(\boldsymbol{x})\|$ uniformly over $\boldsymbol{x} \in \mathcal{C}$, it suffices to bound $\|\nabla \mu_r(\boldsymbol{x}) - \nabla \mu_r(\boldsymbol{x})\|$ and $\|\nabla \bar{f}_r(\boldsymbol{x}) - \nabla \mu_r(\boldsymbol{x})\|$, $\forall r \in [R]$ uniformly over $\boldsymbol{x} \in \mathcal{C}$.

Bounding $\|\nabla \mu_r(x) - \nabla \mu(x)\|$. Note that $\nabla \mu_r(x) = \mathbb{E}_{z \sim q_r}[\nabla f_r(z, x)]$ is a property of the distribution q_r from which the data samples have been drawn and so is $\nabla \mu(x) = \frac{1}{R} \sum_{r=1}^{R} \nabla \mu_r(x)$ the property of q_1, \ldots, q_R . Note that $\|\nabla \mu_r(x) - \nabla \mu(x)\|$ captures heterogeneity among distributions through their expected values, and is equal to zero in the i.i.d. homogeneous data setting of (Chen et al., 2017; Su & Xu, 2019; Yin et al., 2018; 2019). In order to get a meaningful bound for κ , it is reasonable to assume that this heterogeneity is bounded. We assume a uniform bound on the $\|\nabla \mu_r(x) - \nabla \mu(x)\|$ for every $x \in \mathcal{C}$.

Assumption 5. For every client $r \in [R]$, the population mean of the local gradients has a uniformly bounded deviation from the population mean of the global gradient, i.e.,

$$\|\nabla \mu_r(\boldsymbol{x}) - \nabla \mu(\boldsymbol{x})\| \le \kappa_{\text{mean}}, \quad \forall \boldsymbol{x} \in \mathcal{C}.$$
 (92)

Bounding $\|\nabla \bar{f}_r(x) - \nabla \mu_r(x)\|$. Now we bound the difference between the sample mean and the true mean under both sub-exponential and sub-Gaussian distributional assumptions on local gradients.

Let $D = \max\{\|\boldsymbol{x} - \boldsymbol{x}'\| : \boldsymbol{x}, \boldsymbol{x}' \in \mathcal{C}\}$ be the diameter of \mathcal{C} . Note that \mathcal{C} is contained in $\mathcal{B}^d_{D/2}$, which is the Euclidean ball of radius $\frac{D}{2}$ in d dimensions that contains \mathcal{C} . Note that $D = \Omega(\sqrt{d})$, and we assume that D can grow at most polynomially in d

Below we state two lemmas, each of which uniformly bounds $\|\nabla \bar{f}_r(x) - \nabla \mu_r(x)\|$ over all $x \in \mathcal{C}$ under different distributional assumptions on gradients.

Lemma 5 (Sub-exponential gradients). Suppose Assumption 3 holds. Take an arbitrary $r \in [R]$. Let $n_r \in \mathbb{N}$ be sufficiently large such that $n_r = \Omega\left(d\log(n_r d)\right)$. Then, with probability at least $1 - \frac{1}{(1+n_r LD)^d}$, we have

$$\|\nabla \bar{f}_r(\boldsymbol{x}) - \nabla \mu_r(\boldsymbol{x})\| \le 3\nu \sqrt{\frac{8d \log(1 + n_r L D)}{n_r}}, \quad \forall x \in \mathcal{C}.$$
 (93)

Lemma 6 (Sub-Gaussian gradients). Suppose Assumption 4 holds. Take an arbitrary $r \in [R]$. For any $n_r \in \mathbb{N}$, with probability at least $1 - \frac{1}{(1+n_rLD)^d}$, we have

$$\|\nabla \bar{f}_r(\boldsymbol{x}) - \nabla \mu_r(\boldsymbol{x})\| \le 3\sigma_g \sqrt{\frac{8d\log(1 + n_r LD)}{n_r}}, \quad \forall x \in \mathcal{C}.$$
 (94)

We prove Lemma 5 in Appendix E.4 and Lemma 6 in Appendix E.5.

Now we state our main result on bounding the gradient dissimilarity, which we will prove with the help of the above two lemmas. For notational convenience, we state for the case when all clients have the same number of data samples.

Theorem 6 (Gradient dissimilarity). Suppose $n := n_r, \forall r \in [R]$, and Assumption 5 holds. Then, the gradient dissimilarity bound under different distributional assumptions is as follows:

1. Sub-exponential: Suppose Assumption 3 holds. Let $n \in \mathbb{N}$ be sufficiently large such that $n = \Omega(d \log(nd))$. Then, with probability at least $1 - \frac{R}{(1+nLD)^d}$, the following bound holds for all $r \in [R]$:

$$\|\nabla \bar{f}_r(\boldsymbol{x}) - \nabla \bar{f}(\boldsymbol{x})\| \le \kappa_{\text{mean}} + \mathcal{O}\left(\sqrt{\frac{d\log(nd)}{n}}\right), \quad \forall \boldsymbol{x} \in \mathcal{C}.$$
 (95)

2. Sub-Gaussian: Suppose Assumption 4 holds. For every $n \in \mathbb{N}$, with probability at least $1 - \frac{R}{(1+nLD)^d}$, the following bound holds for all $r \in [R]$:

$$\|\nabla \bar{f}_r(\boldsymbol{x}) - \nabla \bar{f}(\boldsymbol{x})\| \le \kappa_{\text{mean}} + \mathcal{O}\left(\sqrt{\frac{d\log(nd)}{n}}\right), \quad \forall \boldsymbol{x} \in \mathcal{C}.$$
 (96)

Remark 1. Note that under Assumption 3 (sub-exponential), the gradient dissimilarity bound (95) holds only when each client has sufficiently large number of samples $n = \Omega\left(d\log(nd)\right)$. On the other hand, under Assumption 4 (sub-Gaussian), the gradient dissimilarity bound (96) holds for every $n \in \mathbb{N}$.

Proof of Theorem 6. In order to prove Theorem 6, we need to show two bounds, one (stated in (95)) under the sub-exponential gradient assumption, and the other (stated in (96)) under the sub-Gaussian assumption. We can show (95) using Lemma 5 and (96) using Lemma 6. Here we only show (95); and (96) can be shown similarly.

Using Assumption 5 (i.e., $\|\nabla \mu_r(x) - \nabla \mu(x)\| \le \kappa_{\text{mean}}, \forall x \in \mathcal{C}$) in (91) gives

$$\|\nabla \bar{f}_r(\boldsymbol{x}) - \nabla \bar{f}(\boldsymbol{x})\| \le \|\nabla \bar{f}_r(\boldsymbol{x}) - \nabla \mu_r(\boldsymbol{x})\| + \kappa_{\text{mean}} + \frac{1}{R} \sum_{r=1}^{R} \|\nabla \bar{f}_r(\boldsymbol{x}) - \nabla \mu_r(\boldsymbol{x})\|.$$
 (97)

Note that (93) holds for any fixed client $r \in [R]$. By the union bound, we have that with probability at least $1 - \frac{R}{(1 + n_r LD)^d}$, for every $r \in [R]$, we have $\|\nabla \bar{f}_r(\boldsymbol{x}) - \nabla \mu_r(\boldsymbol{x})\| \le 3\nu \sqrt{\frac{8d \log(1 + n_r LD)}{n_r}}, \forall \boldsymbol{x} \in \mathcal{C}$.

Let $n_r = n, \forall r \in [R]$. Using these in (97), we get that with probability at least $1 - \frac{R}{(1 + n_r LD)^d}$, for every client $r \in [R]$, we have $\|\nabla \bar{f}_r(\boldsymbol{x}) - \nabla \bar{f}(\boldsymbol{x})\| \le \kappa_{\text{mean}} + \mathcal{O}\left(\sqrt{\frac{d \log(nd)}{n}}\right), \forall \boldsymbol{x} \in \mathcal{C}$, which proves (95). This completes the proof of Theorem 6.

E.2. Bounding the local variances

The local variance bound at the r'th client is $\mathbb{E}_{i \in U[n_r]} \left\| \nabla f_r(\boldsymbol{z}_{r,i}, \boldsymbol{x}) - \nabla \bar{f}_r(\boldsymbol{x}) \right\|^2 \leq \sigma^2$ (from (88)). We simplify the LHS:

$$\mathbb{E}_{i \in U[n_r]} \left\| \nabla f_r(\boldsymbol{z}_{r,i}, \boldsymbol{x}) - \nabla \bar{f}_r(\boldsymbol{x}) \right\|^2 \leq 2\mathbb{E}_{i \in U[n_r]} \left\| \nabla f_r(\boldsymbol{z}_{r,i}, \boldsymbol{x}) - \nabla \mu_r(\boldsymbol{x}) \right\|^2$$

$$+ 2\mathbb{E}_{i \in U[n_r]} \left\| \nabla \bar{f}_r(\boldsymbol{x}) - \nabla \mu_r(\boldsymbol{x}) \right\|^2$$

$$\stackrel{\text{(a)}}{=} 2 \left\| \nabla f_r(\boldsymbol{z}_{r,1}, \boldsymbol{x}) - \nabla \mu_r(\boldsymbol{x}) \right\|^2 + 2 \left\| \nabla \bar{f}_r(\boldsymbol{x}) - \nabla \mu_r(\boldsymbol{x}) \right\|^2$$

$$\stackrel{\text{(b)}}{\leq} 4 \left\| \nabla f_r(\boldsymbol{z}_{r,1}, \boldsymbol{x}) - \nabla \mu_r(\boldsymbol{x}) \right\|^2$$

$$(98)$$

For the first term on the RHS of (a), we used that $\mathbf{z}_{r,i}, i \in [n_r]$ are i.i.d., and the second term follows because it is independent of $i \in [n_r]$. Inequality (b) follows because $\|\nabla \bar{f}_r(\mathbf{z}) - \nabla \mu_r(\mathbf{z})\|^2 \le \|\nabla f_r(\mathbf{z}_{r,1}, \mathbf{z}) - \nabla \mu_r(\mathbf{z})\|^2$, since the average of i.i.d. samples gives tighter concentration in comparison to if we use just one sample.

Note that bounding $\|\nabla f_r(\boldsymbol{z}_{r,1}, \boldsymbol{x}) - \nabla \mu_r(\boldsymbol{x})\|$ is equivalent to bounding $\|\nabla f_r(\boldsymbol{z}, \boldsymbol{x}) - \nabla \mu_r(\boldsymbol{x})\|$ for a random $\boldsymbol{z} \sim q_r$. Now we provide a uniform bound on $\|\nabla f_r(\boldsymbol{z}, \boldsymbol{x}) - \nabla \mu_r(\boldsymbol{x})\|$ for a random $\boldsymbol{z} \sim q_r$ using the sub-Gaussian gradient assumption. Bounding $\|\nabla f_r(z,x) - \nabla \mu_r(x)\|$ for a random $z \sim q_r$. To bound this, we need sub-Gaussian assumption on local gradients (we can also bound this using sub-exponential assumption, but that will give a bound that scales as $\widetilde{\Omega}(d)$ as opposed to $\widetilde{\Omega}(\sqrt{d})$). Note that Lemma 6 holds for any $n_r \in \mathbb{N}$. In particular, it also holds for $n_r = 1$. So, under Assumption 4, with probability at least $1 - \frac{1}{(1+n_rLD)^d}$, we have

$$\|\nabla f_r(\boldsymbol{z}, \boldsymbol{x}) - \nabla \mu_r(\boldsymbol{x})\| \le 3\sigma_g \sqrt{8d \log(1 + LD)}, \quad \forall \boldsymbol{x} \in \mathcal{C},$$
 (99)

where $z \sim q_r$, and probability is over the randomness due to the sub-Gaussian distribution of local gradients. So, with probability at least $1 - \frac{1}{(1+n_rLD)^d}$, we have

$$\mathbb{E}_{i \in U[n_r]} \left\| \nabla f_r(\boldsymbol{z}_{r,i}, \boldsymbol{x}) - \nabla \bar{f}_r(\boldsymbol{x}) \right\|^2 \le 288\sigma_g^2 d \log(1 + LD), \quad \forall \boldsymbol{x} \in \mathcal{C}.$$
 (100)

Note that (100) holds for a fixed client $r \in [R]$. By taking the union bound over all clients $r \in [R]$ proves our variance bound, which we state below.

Theorem 7 (Variance bound). Suppose $n := n_r, \forall r \in [R]$, and Assumption 4 holds. Then, with probability at least $1 - \frac{R}{(1+nLD)^d}$, the following bound holds for all $r \in [R]$:

$$\mathbb{E}_{i \in_{U}[n]} \left\| \nabla f_r(\boldsymbol{z}_{r,i}, \boldsymbol{x}) - \nabla \bar{f}_r(\boldsymbol{x}) \right\|^2 \le \mathcal{O}\left(d \log(d)\right), \quad \forall \boldsymbol{x} \in \mathcal{C}.$$
(101)

Remark 2 (Sub-Gaussian vs. sub-exponential assumption). Note that, we needed sub-Gaussian assumption on local gradients because we wanted to uniformly bound $\mathbb{E}_{i \in [n_r]} \|\nabla f_r(z_{r,i}, x) - \nabla \mu_r(x)\|^2$, which is the case when we use only one data sample in each SGD iteration. In this paper, we use mini-batch SGD with a variable batch size (to control the approximation error of our solution; see the approximation error analysis in Section 2.2). So, when the batch-size b is sufficiently large and satisfies $b = \Omega(d \log(bd))$, we can work with the sub-exponential gradient assumption because the large batch size gives a concentration similar to sub-Gaussian. This would give a bound of $\mathcal{O}\left(\frac{d \log(bd)}{b}\right)$ on variance.

E.3. Definitions of sub-exponential/sub-Gaussian distributions and concentration inequalities

In this section, we give formal definitions of sub-exponential/sub-Gaussian random variables/vectors and the concentration inequalities for them that we will use later on to prove Lemma 5 and Lemma 6.

Definition 1 (Sub-exponential distribution). A random variable Z with mean $\mu = \mathbb{E}[Z]$ is sub-exponential if there are non-negative parameters (ν, α) such that

$$\mathbb{E}\left[\exp\left(\lambda(Z-\mu)\right)\right] \le \exp\left(\lambda^2\nu^2/2\right), \quad \forall |\lambda| < \frac{1}{\alpha}.$$

A random vector Z with mean $\mu = \mathbb{E}[Z]$ is sub-exponential if its projection on every unit vector is sub-exponential, i.e., there are non-negative parameters (ν, α) such that

$$\sup_{\boldsymbol{v} \in \mathbb{R}^d: \|\boldsymbol{v}\| = 1} \mathbb{E}\left[\exp\left(\lambda \langle Z - \mu, \boldsymbol{v} \rangle\right)\right] \leq \exp\left(\lambda^2 \nu^2 / 2\right), \qquad \forall |\lambda| < \frac{1}{\alpha}.$$

Now we state a concentration inequality for sums of independent sub-exponential random variables.

Fact 2 (Sub-exponential concentration inequality). Suppose X_1, X_2, \ldots, X_n are independent random variables, where for every $i \in [n]$, X_i is sub-exponential with parameters (ν_i, α_i) and mean μ_i . Then $\sum_{i=1}^n X_i$ is sub-exponential with parameters (ν, α) , where $\nu^2 = \sum_{i=1}^n \nu_i^2$ and $\alpha = \max_{1 \le i \le n} \alpha_i$. Moreover, we have

$$\Pr\left[\sum_{i=1}^{n} (X_i - \mu_i) \ge t\right] \le \exp\left(-\frac{1}{2}\min\left\{\frac{t^2}{\nu^2}, \frac{t}{\alpha}\right\}\right), \quad \forall t \ge 0$$
(102)

Definition 2 (Sub-Gaussian distribution). A random variable Z with mean $\mu = \mathbb{E}[Z]$ is sub-Gaussian if there is a non-negative parameter σ_g such that

$$\mathbb{E}\left[\exp\left(\lambda(Z-\mu)\right)\right] \leq \exp\left(\lambda^2 \sigma_{\rm g}^2/2\right), \qquad \forall \lambda \in \mathbb{R}.$$

A random vector Z with mean $\mu = \mathbb{E}[Z]$ is sub-Gaussian if its projection on every unit vector is sub-Gaussian, i.e., there is a non-negative parameter σ_g such that

$$\sup_{\boldsymbol{v} \in \mathbb{R}^d: \|\boldsymbol{v}\| = 1} \mathbb{E}\left[\exp\left(\lambda \langle Z - \mu, \boldsymbol{v} \rangle\right)\right] \leq \exp\left(\lambda^2 \sigma_g^2 / 2\right), \qquad \forall \lambda \in \mathbb{R}.$$

Now we state a concentration inequality for sums of independent sub-Gaussian random variables.

Fact 3 (Sub-Gaussian concentration inequality). Suppose X_1, X_2, \ldots, X_n are independent random variables, where for every $i \in [n]$, X_i is sub-Gaussian with parameter $\sigma_i > 0$ and mean μ_i . Then $\sum_{i=1}^n X_i$ is sub-Gaussian with parameter $\sigma_g = \sqrt{\sum_{i=1}^n \sigma_i^2}$. Moreover, we have

$$\Pr\left[\sum_{i=1}^{n} (X_i - \mu_i) \ge t\right] \le \exp\left(-t^2/2\sigma_g^2\right), \qquad \forall t \ge 0.$$
(103)

E.4. Proof of Lemma 5 (sub-exponential gradients)

We prove Lemma 5 with the help of the following result, which holds for any fixed $x \in C$. Then we extend this bound to all $x \in C$ using an ϵ -net argument. These are standard calculations and have appeared in literature (Chen et al., 2017; Yin et al., 2019).

Lemma 7. Suppose Assumption 3 holds. Take an arbitrary $r \in [R]$. For any $\delta \in (0,1)$ and $n_r \in \mathbb{N}$, define $\Delta = \sqrt{2}\nu\sqrt{\frac{d\log 5 + \log(1/\delta)}{n_r}}$. If n_r is such that $\Delta \leq \frac{\nu^2}{\alpha}$, then, for any fixed $\mathbf{x} \in \mathcal{C}$, with probability at least $1 - \delta$, we have

$$\|\nabla \bar{f}_r(\boldsymbol{x}) - \nabla \mu_r(\boldsymbol{x})\| \le 2\sqrt{2}\nu\sqrt{\frac{d\log 5 + \log(1/\delta)}{n_r}},$$
 (104)

where randomness is due to the sub-exponential distribution of local gradients.

Proof. Let $\mathcal{B}^d = \{ \boldsymbol{v} \in \mathbb{R}^d : \|\boldsymbol{v}\| \le 1 \}$. Let $\mathcal{V} = \{ \boldsymbol{v}_1, \boldsymbol{v}_2, \dots, \boldsymbol{v}_{N_{1/2}} \}$ denote an $\frac{1}{2}$ -net of \mathcal{B}^d , which implies that for every $\boldsymbol{v} \in \mathcal{B}^d$, there exists a $\boldsymbol{v}' \in \mathcal{V}$ such that $\|\boldsymbol{v} - \boldsymbol{v}'\| \le \frac{1}{2}$. We have from (Vershynin, 2010, Lemma 5.2) that $N_{1/2} = |\mathcal{V}| \le 5^d$.

Fix an arbitrary $\boldsymbol{x} \in \mathcal{C}$. Note that there exists a $\boldsymbol{v}^* \in \mathcal{B}^d$ (namely, $\boldsymbol{v}^* = \frac{\nabla \bar{f}_r(\boldsymbol{x}) - \nabla \mu_r(\boldsymbol{x})}{\|\nabla \bar{f}_r(\boldsymbol{x}) - \nabla \mu_r(\boldsymbol{x})\|}$) such that $\|\nabla \bar{f}_r(\boldsymbol{x}) - \nabla \mu_r(\boldsymbol{x})\| = \langle \nabla \bar{f}_r(\boldsymbol{x}) - \nabla \mu_r(\boldsymbol{x}), \boldsymbol{v}^* \rangle$. By the property of \mathcal{V} , there exists an index $i^* \in [N_{1/2}]$ such that $\|\boldsymbol{v}^* - \boldsymbol{v}_{i^*}\| \leq \frac{1}{2}$. Now we bound $\|\nabla \bar{f}_r(\boldsymbol{x}) - \nabla \mu_r(\boldsymbol{x})\|$.

$$\begin{split} \left\| \nabla \bar{f}_r(\boldsymbol{x}) - \nabla \mu_r(\boldsymbol{x}) \right\| &= \left\langle \nabla \bar{f}_r(\boldsymbol{x}) - \nabla \mu_r(\boldsymbol{x}), \boldsymbol{v}^* \right\rangle \\ &= \left\langle \nabla \bar{f}_r(\boldsymbol{x}) - \nabla \mu_r(\boldsymbol{x}), \boldsymbol{v}_{i^*} \right\rangle + \left\langle \nabla \bar{f}_r(\boldsymbol{x}) - \nabla \mu_r(\boldsymbol{x}), \boldsymbol{v}^* - \boldsymbol{v}_{i^*} \right\rangle \\ &\leq \left\langle \nabla \bar{f}_r(\boldsymbol{x}) - \nabla \mu_r(\boldsymbol{x}), \boldsymbol{v}_{i^*} \right\rangle + \left\| \nabla \bar{f}_r(\boldsymbol{x}) - \nabla \mu_r(\boldsymbol{x}) \right\| \left\| \boldsymbol{v}^* - \boldsymbol{v}_{i^*} \right\| \\ &\leq \left\langle \nabla \bar{f}_r(\boldsymbol{x}) - \nabla \mu_r(\boldsymbol{x}), \boldsymbol{v}_{i^*} \right\rangle + \frac{1}{2} \left\| \nabla \bar{f}_r(\boldsymbol{x}) - \nabla \mu_r(\boldsymbol{x}) \right\| \\ &\leq \max_{\boldsymbol{v} \in \mathcal{V}} \left\langle \nabla \bar{f}_r(\boldsymbol{x}) - \nabla \mu_r(\boldsymbol{x}), \boldsymbol{v} \right\rangle + \frac{1}{2} \left\| \nabla \bar{f}_r(\boldsymbol{x}) - \nabla \mu_r(\boldsymbol{x}) \right\| \end{split}$$

By collecting similar terms together, we get

$$\|\nabla \bar{f}_r(\boldsymbol{x}) - \nabla \mu_r(\boldsymbol{x})\| \le 2 \max_{\boldsymbol{v} \in \mathcal{V}} \langle \nabla \bar{f}_r(\boldsymbol{x}) - \nabla \mu_r(\boldsymbol{x}), \boldsymbol{v} \rangle$$
(105)

Note that the RHS of (105) is a non-negative number (because LHS is). Note also that, since $\mathcal{V} \subset \mathcal{B}^d$, for every $\boldsymbol{v} \in \mathcal{V}$, we have $\|\boldsymbol{v}\| \leq 1$. This implies that $\max_{\boldsymbol{v} \in \mathcal{V}} \left\langle \nabla \bar{f}_r(\boldsymbol{x}) - \nabla \mu_r(\boldsymbol{x}), \boldsymbol{v} \right\rangle \leq \max_{\boldsymbol{v} \in \mathcal{V}} \left\langle \nabla \bar{f}_r(\boldsymbol{x}) - \nabla \mu_r(\boldsymbol{x}), \frac{\boldsymbol{v}}{\|\boldsymbol{v}\|} \right\rangle$. Using this in (105), we get

$$\left\|\nabla \bar{f}_r(\boldsymbol{x}) - \nabla \mu_r(\boldsymbol{x})\right\| \le 2 \max_{\boldsymbol{v} \in \mathcal{V}} \left\langle \nabla \bar{f}_r(\boldsymbol{x}) - \nabla \mu_r(\boldsymbol{x}), \frac{\boldsymbol{v}}{\|\boldsymbol{v}\|} \right\rangle. \tag{106}$$

Fix any $v \in \mathcal{V}$. It follows from Assumption 3 that $\left\langle \nabla f_r(\boldsymbol{z}, \boldsymbol{x}) - \nabla \mu_r(\boldsymbol{x}), \frac{v}{\|\boldsymbol{v}\|} \right\rangle$, where $\boldsymbol{z} \sim q_r$, is a sub-exponential random variable (with mean zero) with parameters (ν, α) . From Fact 2 (stated on page 34), we have that $\sum_{i=1}^{n_r} \left\langle \nabla f_r(\boldsymbol{z}_{r,i}, \boldsymbol{x}) - \nabla \mu_r(\boldsymbol{x}), \frac{v}{\|\boldsymbol{v}\|} \right\rangle$ (where $\boldsymbol{z}_{r,i} \sim q_r, i \in [n_r]$ are i.i.d.) is a sub-exponential random variable with parameters $(\sqrt{n_r}\nu, \alpha)$.

Now, apply the concentration bound from (102) with $t=n_r\Delta$. Substituting this and the parameters $(\sqrt{n_r}\nu,\alpha)$, the bound becomes $\exp(-\frac{1}{2}\min\{\frac{n_r^2\Delta^2}{n_r\nu^2},\frac{n_r\Delta}{\alpha}\})\stackrel{\text{(a)}}{=}\exp(-\frac{1}{2}\frac{n_r\Delta^2}{\nu^2})$, where (a) follows because $\Delta\leq\frac{\nu^2}{\alpha}$. This gives

$$\Pr\left[\sum_{i=1}^{n_r} \left\langle \nabla f_r(\boldsymbol{z}_{r,i}, \boldsymbol{x}) - \nabla \mu_r(\boldsymbol{x}), \frac{\boldsymbol{v}}{\|\boldsymbol{v}\|} \right\rangle \ge n_r \Delta\right] \le \exp\left(-\frac{n_r \Delta^2}{2\nu^2}\right). \tag{107}$$

Note that $\sum_{i=1}^{n_r} \left\langle \nabla f_r(\boldsymbol{z}_{r,i}, \boldsymbol{x}) - \nabla \mu_r(\boldsymbol{x}), \frac{\boldsymbol{v}}{\|\boldsymbol{v}\|} \right\rangle = n_r \left\langle \nabla \bar{f}_r(\boldsymbol{x}) - \nabla \mu_r(\boldsymbol{x}), \frac{\boldsymbol{v}}{\|\boldsymbol{v}\|} \right\rangle$. Using this in (107) yields

$$\Pr\left[\left\langle \nabla \bar{f}_r(\boldsymbol{x}) - \nabla \mu_r(\boldsymbol{x}), \frac{\boldsymbol{v}}{\|\boldsymbol{v}\|} \right\rangle \ge \Delta\right] \le \exp\left(-\frac{n_r \Delta^2}{2\nu^2}\right)$$
(108)

This implies that

$$\Pr\left[\max_{\boldsymbol{v}\in\mathcal{V}}\left\langle\nabla\bar{f}_r(\boldsymbol{x}) - \nabla\mu_r(\boldsymbol{x}), \frac{\boldsymbol{v}}{\|\boldsymbol{v}\|}\right\rangle \ge \Delta\right] \le \sum_{\boldsymbol{v}\in\mathcal{V}} \Pr\left[\left\langle\nabla\bar{f}_r(\boldsymbol{x}) - \nabla\mu_r(\boldsymbol{x}), \frac{\boldsymbol{v}}{\|\boldsymbol{v}\|}\right\rangle \ge \Delta\right]$$

$$\le |\mathcal{V}| \exp\left(-\frac{n_r\Delta^2}{2\nu^2}\right) \le 5^d \exp\left(-\frac{n_r\Delta^2}{2\nu^2}\right)$$

$$= \exp\left(-\frac{n_r\Delta^2}{2\nu^2} + d\log 5\right)$$
(109)

Together with (106), which implies that

$$\Pr\left[\left\|\nabla \bar{f}_r(\boldsymbol{x}) - \nabla \mu_r(\boldsymbol{x})\right\| \ge t\right] \le \Pr\left[2 \max_{\boldsymbol{v} \in \mathcal{V}} \left\langle \nabla \bar{f}_r(\boldsymbol{x}) - \nabla \mu_r(\boldsymbol{x}), \frac{\boldsymbol{v}}{\|\boldsymbol{v}\|} \right\rangle \ge t\right]$$

holds for every t > 0, (109) gives

$$\Pr\left[\left\|\nabla \bar{f}_r(\boldsymbol{x}) - \nabla \mu_r(\boldsymbol{x})\right\| \ge 2\Delta\right] \le \exp\left(-\frac{n_r \Delta^2}{2\nu^2} + d\log 5\right) \le \delta,\tag{110}$$

where in the last inequality we used $\Delta = \sqrt{2}\nu\sqrt{\frac{d\log 5 + \log(1/\delta)}{n_r}}$

This completes the proof of Lemma 7.

Proof of Lemma 5. We have from Lemma 7 that for each fixed $x \in \mathcal{C}$, with probability at least $1 - \delta$, we have

$$\left\|\nabla \bar{f}_r(\boldsymbol{x}) - \nabla \mu_r(\boldsymbol{x})\right\| \le 2\nu \sqrt{\frac{2d\log 5 + 2\log(1/\delta)}{n_r}}.$$
(111)

To extend this argument uniformly over the entire set C, we use another covering argument. Recall that D is the diameter of C. Note that C is contained in $\mathcal{B}^d_{D/2}$, which is the Euclidean ball of radius $\frac{D}{2}$ in d dimensions that contains C. For some $\delta_0 > 0$,

let
$$C_{\delta_0} = \{x_0, x_2, \dots, x_{N_{\delta_0}}\}$$
 be the δ_0 -net of C . It follows from (Vershynin, 2010, Lemma 5.2) that $N_{\delta_0} \leq \left(1 + \frac{D}{\delta_0}\right)^d$.

Applying the union bound in (111), we get that with probability at least $1 - \delta$, we have for all $x_i \in \mathcal{C}_{\delta_0}$,

$$\left\|\nabla \bar{f}_r(\boldsymbol{x}_i) - \nabla \mu_r(\boldsymbol{x}_i)\right\| \le 2\nu \sqrt{\frac{2d\log 5 + 2\log\left(\frac{N_{\delta_0}}{\delta}\right)}{n_r}}.$$
(112)

We want to bound $\|\nabla \bar{f}_r(\boldsymbol{x}) - \nabla \mu_r(\boldsymbol{x})\|$ for all $\boldsymbol{x} \in \mathcal{C}$. Take any $\boldsymbol{x} \in \mathcal{C}$. Since \mathcal{C}_{δ_0} is a δ_0 -net of \mathcal{C} , there exists an $\boldsymbol{x}' \in \mathcal{C}_{\delta_0}$ such that $\|\boldsymbol{x} - \boldsymbol{x}'\| \leq \delta_0$.

$$\|\nabla \bar{f}_{r}(\boldsymbol{x}) - \nabla \mu_{r}(\boldsymbol{x})\| = \|\nabla \bar{f}_{r}(\boldsymbol{x}) - \nabla \bar{f}_{r}(\boldsymbol{x}') + \nabla \bar{f}_{r}(\boldsymbol{x}') - \nabla \mu_{r}(\boldsymbol{x}) + \nabla \mu_{r}(\boldsymbol{x}') - \nabla \mu_{r}(\boldsymbol{x}')\|$$

$$\leq \underbrace{\|\nabla \bar{f}_{r}(\boldsymbol{x}) - \nabla \bar{f}_{r}(\boldsymbol{x}')\|}_{=: T_{1}} + \underbrace{\|\nabla \mu_{r}(\boldsymbol{x}) - \nabla \mu_{r}(\boldsymbol{x}')\|}_{=: T_{2}} + \|\nabla \bar{f}_{r}(\boldsymbol{x}') - \nabla \mu_{r}(\boldsymbol{x}')\|$$

$$(113)$$

Now we bound each term on the RHS of (113).

$$T_{1} = \left\| \frac{1}{n_{r}} \sum_{i=1}^{n_{r}} \left(\nabla f_{r}(\boldsymbol{z}_{r,i}, \boldsymbol{x}) - \nabla f_{r}(\boldsymbol{z}_{r,i}, \boldsymbol{x}') \right) \right\| \leq \frac{1}{n_{r}} \sum_{i=1}^{n_{r}} \left\| \nabla f_{r}(\boldsymbol{z}_{r,i}, \boldsymbol{x}) - \nabla f_{r}(\boldsymbol{z}_{r,i}, \boldsymbol{x}') \right\|$$

$$\leq L \|\boldsymbol{x} - \boldsymbol{x}'\| \leq L \delta_{0}$$

$$T_{2} = \left\| \mathbb{E}_{\boldsymbol{z} \sim q_{r}} \left[\nabla f_{r}(\boldsymbol{z}, \boldsymbol{x}) - \nabla f_{r}(\boldsymbol{z}, \boldsymbol{x};) \right] \right\| \leq \mathbb{E}_{\boldsymbol{z} \sim q_{r}} \left\| \nabla f_{r}(\boldsymbol{z}, \boldsymbol{x}) - \nabla f_{r}(\boldsymbol{z}, \boldsymbol{x};) \right\|$$

$$\leq \mathbb{E}_{\boldsymbol{z} \sim q_{r}} L \|\boldsymbol{x} - \boldsymbol{x}'\| \leq L \delta_{0}$$

Substituting the above bounds on T_1 , T_2 in (113) and bounding the third term of (113) using (112) gives

$$\|\nabla \bar{f}_r(\boldsymbol{x}) - \nabla \mu_r(\boldsymbol{x})\| \le 2L\delta_0 + 2\nu \sqrt{\frac{2d\log 5 + 2\log\left(\frac{N_{\delta_0}}{\delta}\right)}{n_r}}.$$
(114)

Note that $N_{\delta_0} \leq \left(1 + \frac{D}{\delta_0}\right)^d$. Take $\delta = 1/\left(1 + \frac{D}{\delta_0}\right)^d$. If we take $\delta_0 = \frac{1}{n_r L}$, which implies $\delta = \frac{1}{(1 + n_r LD)^d}$, we would get $2d\log 5 + 2\log\left(\frac{N_{\delta_0}}{\delta}\right) \leq 4d + 4d\log(1 + n_r LD) \leq 8d\log(1 + n_r LD)$. Substituting these in above gives

$$\left\|\nabla \bar{f}_r(\boldsymbol{x}) - \nabla \mu_r(\boldsymbol{x})\right\| \le \frac{2}{n_r} + \frac{2\nu}{\sqrt{n_r}} \sqrt{8d \log(1 + n_r LD)}.$$
 (115)

When $n_r \ge \frac{1}{2\nu^2 d \log(1+n_r LD)}$ (which is a very small number less than 1), with probability at least $1 - \frac{1}{(1+n_r LD)^d}$, we have

$$\|\nabla \bar{f}_r(\boldsymbol{x}) - \nabla \mu_r(\boldsymbol{x})\| \le 3\nu \sqrt{\frac{8d\log(1 + n_r LD)}{n_r}}, \quad \forall \boldsymbol{x} \in \mathcal{C}.$$
 (116)

Lower bound on n_r . Note that Lemma 7 requires $\Delta \leq \frac{\nu^2}{\alpha}$, where $\Delta = \sqrt{2}\nu\sqrt{\frac{d\log 5 + \log(1/\delta)}{n_r}}$. Substituting the value of $\delta = \frac{1}{(1+n_rLD)^d}$ gives $n_r \geq \frac{2\alpha^2}{\nu^2} \left(d\log 5 + d\log(1+n_rLD)\right)$, which is $\Omega(d\log(n_rLD))$ for constant α,ν . Treating the smoothness parameter L a constant, we get $n_r = \Omega(d\log(n_rd))$ to be requirement on the sample size at the r'th client for the bound in Lemma 5 to hold.

This completes the proof of Lemma 5.

E.5. Proof of Lemma 6 (sub-Gaussian gradients)

We prove Lemma 6 with the help of the following result, which holds for any fixed $x \in \mathcal{C}$.

Lemma 8. Suppose Assumption 4 holds. Take an arbitrary $r \in [R]$. For any $\delta \in (0,1)$ and $n_r \in \mathbb{N}$, with probability at least $1 - \delta$, we have for any fixed $x \in C$:

$$\|\nabla \bar{f}_r(\boldsymbol{x}) - \nabla \mu_r(\boldsymbol{x})\| \le 2\sqrt{2}\sigma_g \sqrt{\frac{d\log 5 + \log(1/\delta)}{n_r}},$$
 (117)

where randomness is due to the sub-Gaussian distribution of local gradients.

Proof. Follow the proof of Lemma 7 exactly until (106). Then instead of the sub-exponential assumption, use the sub-Gaussian assumption (Assumption 4) on local gradients. Then apply the concentration bound from (103) with $t = n_r \Delta$. This gives that for any fixed $v \in \mathcal{V}$ and any $\Delta \geq 0$, we have

$$\Pr\left[\left\langle \nabla \bar{f}_r(\boldsymbol{x}) - \nabla \mu_r(\boldsymbol{x}), \frac{\boldsymbol{v}}{\|\boldsymbol{v}\|} \right\rangle \ge \Delta\right] \le \exp\left(-\frac{n_r \Delta^2}{2\sigma_{\rm g}^2}\right). \tag{118}$$

Now following the proof of Lemma 7 from (108) to (110) gives

$$\Pr\left[\left\|\nabla \bar{f}_r(\boldsymbol{x}) - \nabla \mu_r(\boldsymbol{x})\right\| \ge 2\Delta\right] \le \exp\left(-\frac{n_r \Delta^2}{2\sigma_g^2} + d\log 5\right) \le \delta,\tag{119}$$

where in the last inequality we used
$$\Delta = \sqrt{2}\sigma_g\sqrt{\frac{d\log 5 + \log(1/\delta)}{n_r}}$$
.

We can extend the bound from Lemma 8 to all $x \in C$ (and prove Lemma 6) using an ϵ -net argument exactly in the same way as used in the proof of Lemma 5. So, to avoid repetition, we do not show this extension here.

F. Additional Experimental Details

There are some implementation issues about the decoding algorithm (as described in Algorithm 2) that could be important in the deployment of the algorithm. In the following, we describe these issues and also explain our approach in the implementation to address them.

- Note that the stopping criterion (see line 7) in our decoding algorithm described in Algorithm 2 requires the matrix concentration bound σ_0^2 that we show in Theorem 2 in terms of the SGD variance bound σ^2 (see (2)) and the bounded gradient dissimilarity κ^2 (see (6)). Since these are properties of the local datasets stored at clients, which is challenging to determine in a adversarial federated learning setting. In order to mitigate this, we observe two things:
 - 1. the only place where Algorithm 2 uses this matrix concentration bound is in the stopping criterion (in line 7); and
 - 2. in each iteration of the while loop, at least one sample gets its weight reduced to zero.

Since we know an upper bound on the fraction of corrupt samples, these two observations suggest replacing the stopping condition in line 7 with the condition that break the while loop when the number of samples whose weights become zero is more than the number of corrupt samples. This is what we used as a stopping criterion (in line 7) in our implementation of Algorithm 2.

• Note that each iteration of the while loop (line 7) of Algorithm 2 requires computing the principal eigenvector of the covariance matrix (line 8), which can be done using the singular value decomposition (SVD) algorithm. This, however, could be computationally expensive. To mitigate this, we choose uniformly at random 1024 coordinates from the all gradient vectors (same 1024 random coordinates from all the gradients), and run the decoding algorithm only on them. Suppose \mathcal{A} denotes the set of indices of the surviving gradients (i.e., whose weight are not zero when the filtering algorithm terminates), then we will discard all those full gradients whose indices are outside the set \mathcal{A} .

Furthermore, we observed performance boost when replacing the line 13 of Algorithm 2 (i.e., $\hat{g} = \sum_{i=1}^K \frac{w_i^{(t)}}{\|\boldsymbol{w}^{(t)}\|_1} \boldsymbol{g}_i$) with $\hat{g} = \sum_{i \in \mathcal{A}} \frac{1}{|\mathcal{A}|} \boldsymbol{g}_i$, where \mathcal{A} contains the identities of the surviving samples; in other words, we replaced the weighted average with the uniform average.