

MOMENTUM AND ERROR FEEDBACK FOR CLIPPING WITH FAST RATES AND DIFFERENTIAL PRIVACY

Anonymous authors

Paper under double-blind review

ABSTRACT

Strong Differential Privacy (DP) and Optimization guarantees are two desirable properties for a method in Federated Learning (FL). However, existing algorithms do not achieve both properties at once: they either have optimal DP guarantees but rely on restrictive assumptions such as bounded gradients/bounded data heterogeneity, or they ensure strong optimization performance but lack DP guarantees. To address this gap in the literature, we propose and analyze a new method called Clip21-SGDM based on a novel combination of clipping, heavy-ball momentum, and Error Feedback. In particular, for non-convex smooth distributed problems with clients having arbitrarily heterogeneous data, we prove that Clip21-SGDM has optimal convergence rate and also optimal (local-)DP neighborhood. Our numerical experiments on non-convex logistic regression and training of neural networks highlight the superiority of Clip21-SGDM over baselines in terms of the optimization performance for a given DP-budget.

1 INTRODUCTION

Federated Learning (Konečný et al., 2016; McMahan et al., 2017a) is a modern training paradigm where multiple (possibly heterogeneous) clients aim to jointly train a machine learning model without sacrificing the privacy of their own data. This setup presents several noticeable challenges in terms of algorithm design affecting different aspects of training, including communication efficiency, partial participation of clients, data heterogeneity, security, and privacy (Kairouz et al., 2021; Wang et al., 2021). As a result, numerous optimization methods for Federated Learning (FL) have been introduced in recent years. However, despite extensive research in the field, achieving both strong optimization convergence and robust differential privacy (DP) guarantees (Dwork et al., 2014) simultaneously in an FL algorithm remains challenging due to the conflicting nature of these objectives. Indeed, most of the results in the field of DP are obtained by adding noise (e.g. Gaussian noise) to the method’s update (Abadi et al., 2016; Chen et al., 2020) in order to protect the client’s data that could be potentially reconstructed from the updates. Unfortunately, this approach results in less accurate updates, which negatively affects the convergence. Moreover, to ensure DP, this mechanism should be applied to the method with bounded updates, which is typically achieved via *gradient clipping* (Pascanu et al., 2013).

Further complicating the issue, naïve distributed Clipped Gradient Descent (Clip-GD) is not guaranteed to converge (Khirirat et al., 2023) when clients have heterogeneous data (even in the absence of any additive DP-noise), which is a common scenario in FL. To address this issue Khirirat et al. (2023) apply the EF21 mechanism – originally developed by Richtárik et al. (2021) for contractive compression operators to improve the standard Error Feedback (Seide et al., 2014) – to Clip-GD, resulting in a method known as Clip21-GD. Khirirat et al. (2023) show that in contrast to Clip-GD, Clip21-GD converges with $\mathcal{O}(1/T)$ rate for smooth non-convex problems with arbitrary heterogeneous data on clients. However, their analysis is limited to the case of full-batched gradients and does not work with DP-noise. This leads us to the natural question:

Is it possible to design a method that combines both strong optimization performance and DP guarantees in a stochastic setting?

Our contribution. In this paper, we provide a positive answer to the above question by introducing a new method, named Clip21-SGDM, which incorporates clipping, error feedback and heavy-ball

momentum (Polyak, 1964) in a novel way. For smooth non-convex distributed optimization problems, we show that Clip21-SGDM (i) converges with optimal $\mathcal{O}(1/T)$ rate when the workers compute full gradients, (ii) converges with optimal $\tilde{\mathcal{O}}(1/\sqrt{nT})$ high-probability convergence rate when the workers use stochastic gradients with sub-Gaussian noise, and (iii) has optimal local DP-error when DP-noise is added to the clients' updates. We also prove that Clip21-SGD is not guaranteed to converge in the stochastic case, underscoring the need for changes in the algorithm. Our experiments on logistic regression and neural networks highlight the robustness of Clip21-SGDM to the choice of clipping level and indicate Clip21-SGDM's superiority over Clip-SGD and Clip21-SGD in terms of optimization performance for a given DP-budget.

1.1 PROBLEM FORMULATION AND ASSUMPTIONS

We consider the optimization problem of the form

$$\min_{x \in \mathbb{R}^d} \left[f(x) := \frac{1}{n} \sum_{i=1}^n f_i(x) \right] \quad (1)$$

that typically appears in many machine learning applications and is standard for Federated Learning. Here x denotes the parameters of a model, f_i represents the loss associated with the local dataset \mathcal{D}_i of worker $i \in [n]$, and f is an average loss across all workers participating in the training process.

We make two main assumptions on the problem. The first one is smoothness, which is standard for non-convex optimization (Carmon et al., 2020; Danilova et al., 2022). In addition, we also assume that $f(x)$ is uniformly lower bounded since otherwise, problem (1) is intractable.

Assumption 1. We assume that each individual loss function f_i is L -smooth, i.e., for any $x, y \in \mathbb{R}^d$ and $i \in [n]$ we have

$$\|\nabla f_i(x) - \nabla f_i(y)\| \leq L\|x - y\|. \quad (2)$$

Moreover, we assume that $f^* := \inf_{x \in \mathbb{R}^d} f(x) > -\infty$.

We also note that our analysis can be easily generalized to the case when L depends on f_i .

Next, since computation of the full gradients is expensive in many practical applications, it is natural to consider the case when clients compute stochastic gradients. We make the following assumption on the stochastic noise of these gradients.

Assumption 2. We assume that each worker i has access to a σ -sub-Gaussian unbiased estimator $\nabla f_i(x, \xi)$ of a local gradient $\nabla f_i(x)$, i.e., for some¹ $\sigma \geq 0$ and any $x \in \mathbb{R}^d$ and $\forall i \in [n]$ we have

$$\mathbb{E}[\nabla f_i(x, \xi)] = \nabla f_i(x), \quad \mathbb{E}[\exp(\|\theta_i^t\|^2/\sigma^2)] \leq \exp(1), \quad (3)$$

where ξ denotes the source of the stochasticity and $\theta_i := \nabla f_i(x, \xi) - \nabla f_i(x)$.

Although this assumption is stronger than bounded variance, it is standard for the high-probability² analysis of SGD-type methods with polylogarithmic dependence on the confidence level (Nemirovski et al., 2009; Ghadimi & Lan, 2012). The second part of (3) is equivalent to $\Pr(\|\theta_i^t\| \geq b) \leq 2 \exp(-b^2/(2\sigma^2))$ up to a constant factor in σ^2 (Vershynin, 2018). We also note that it is possible to show high-probability bounds for SGD-type methods with polylogarithmic dependence on the confidence level when the noise has sub-Weibull tails (Madden et al., 2024), i.e., the noise can be even heavier but it affects the polylogarithmic factors.

Finally, we provide two important definitions for this work. The first one is the definition of the clipping operator, which is a non-linear map from \mathbb{R}^d to \mathbb{R}^d parameterized by the clipping threshold/level $\tau > 0$ and defined as

$$\text{clip}_\tau(x) := \begin{cases} \frac{\tau}{\|x\|}x, & \text{if } \|x\| > \tau, \\ x, & \text{if } \|x\| \leq \tau. \end{cases} \quad (4)$$

Next, we will use the following classical definition of (ϵ, δ) -Differential Privacy, which introduces plausible deniability into the output of a learning algorithm.

¹For simplicity, we define $0/0 := 0$. Then, (3) with $\sigma = 0$ implies $\nabla f_i(x, \xi) = \nabla f_i(x)$ almost surely.

²We elaborate on the reasons why we focus on high-probability analysis in Section 3.2.

Definition 1 ((ϵ, δ) -Differential Privacy (Dwork et al., 2014)). A randomized method $\mathcal{M} : \mathcal{D} \rightarrow \mathcal{R}$ satisfies (ϵ, δ) -Differential Privacy ((ϵ, δ) -DP) if for any adjacent $D, D' \in \mathcal{D}$ (e.g., if D and D' are datasets, then the adjacency means that D and D' differ in 1 sample) and for any $S \subseteq \mathcal{R}$

$$\Pr(\mathcal{M}(D) \in S) \leq e^\epsilon \Pr(\mathcal{M}(D') \in S) + \delta. \quad (5)$$

In this definition, the smaller ϵ, δ are, the more private the method is. Intuitively, if inequality (5) holds with small values of ϵ and δ , it becomes difficult to infer the specific data point that differs between two similar datasets based solely on the output of \mathcal{M} .

1.2 RELATED WORK

Differential Privacy. The most common approach to obtaining DP guarantees is to clip each client’s update, i.e., by bounding their ℓ_2 norm, and adding a calibrated amount of Gaussian noise to each update or the average. This is typically sufficient to obscure the influence of any single client (McMahan et al., 2017b). Commonly, two scenarios of the DP model are considered: *the central model* and *the local model*. In the first setting, central privacy, a trusted server collects updates and adds noise only before updating the server-side model. This ensures that client data remains private from external parties. In the second setting, local privacy, client data is protected even from the server by clipping and adding noise to updates locally before sending them to the server, ensuring privacy from both the server and other clients (Kasiviswanathan et al., 2011; Allouah et al., 2024). The local privacy setting offers stronger privacy against untrusted servers but results in poorer learning performance due to the need for more noise to obscure individual updates (Chan et al., 2012; Duchi et al., 2018). This can be improved by using a secure shuffler (Erlingsson et al., 2019; Balle et al., 2019), which permutes updates, or a secure aggregator (Bonawitz et al., 2017), which sums updates before sending them to the server. These methods anonymize updates and enhance privacy while maintaining reasonable learning performance, even without a fully trusted server. Finally, (Chaudhuri et al., 2022; Hegazy et al., 2024) show that when DP is required, one can also achieve compression of updates for free.

In this work, we adopt the local DP model by injecting Gaussian noise into each client’s update. However, the average noise can also be viewed as noise added to the average update. Therefore, Clip21-SGDM is compatible with all the aforementioned techniques and can also be applied to the central DP model with a smaller amount of noise.

Distributed methods with clipping. In the single-node regime, Clip-SGD has been analyzed under various assumptions by many authors (Zhang et al., 2020b;c;a; Gorbunov et al., 2020a; Cutkosky & Mehta, 2021; Sadiiev et al., 2023; Liu et al., 2023). Of course, these results can be generalized to the multi-node case if clipping is applied to the aggregated (e.g. averaged) vector, although mini-batching requires a refined analysis when the noise is heavy-tailed (Kornilov et al., 2024). However, to get DP, clipping has to be applied to the vectors communicated by clients to the server. In this regime, Clip-SGD is not guaranteed to converge even without any stochastic noise in the gradients (Chen et al., 2020; Khirirat et al., 2023). There exist several approaches to bypass this limitation that can be split into two lines of work. The first one relies on explicit or implicit assumptions about bounded heterogeneity. More precisely, Liu et al. (2022) analyze a version of Local-SGD/FedAvg (Mangasarian, 1995; McMahan et al., 2017a) with gradient clipping for homogeneous data case assuming that the stochastic gradients have symmetric distribution around their mean and Wei et al. (2020) consider Local-SGD with clipping of the models and analyze its convergence under bounded heterogeneity assumption. Moreover, the boundedness of the stochastic gradient is another assumption used in the literature but it implies the boundedness of gradients’ heterogeneity of clients as well. This assumption is used in numerous works, including: i) Zhang et al. (2022) in the analysis of a version of FedAvg with clipping of model difference (also empirically studied by Geyer et al. (2017)), ii) Noble et al. (2022) who propose and analyze a version of SCAFFOLD (Karimireddy et al., 2020) with gradient clipping (DP-SCAFFOLD), iii) Li & Chi (2023) who propose and analyze a version of BEER (Li et al., 2021) with gradient clipping (PORTER) under bounded gradient and/or bounded data heterogeneity assumption, and iv) Allouah et al. (2024) who study a version of Gossip-SGD (Nedic & Ozdaglar, 2009) with gradient clipping (DECOR). Although most of the mentioned works have rigorous DP guarantees, the corresponding methods are not guaranteed to converge for arbitrary heterogeneous problems.

The second line of work focuses on the clipping of shifted (stochastic) gradient. In particular, [Khirirat et al. \(2023\)](#) proposed and analyzed Clip21-GD, which is based on the application of EF21 ([Richtárik et al., 2021](#)) to the clipping operator, and [Gorbunov et al. \(2024\)](#) develop and analyze methods that apply clipping to the difference of stochastic gradients and learnable shift – an idea that was initially proposed by [Mishchenko et al. \(2019\)](#) to handle data heterogeneity in the Distributed Learning with unbiased communication compression. However, the analysis from ([Khirirat et al., 2023](#)) is limited to the noiseless regime, i.e., full-batched gradients are computed on workers, and both of the mentioned works do not provide³ DP guarantees. We also note that clipping of gradient differences is helpful in tolerating Byzantine attacks in the partial participation regime ([Malinovsky et al., 2023](#)).

Error Feedback. Error Feedback (EF) ([Seide et al., 2014](#)) is a popular technique for incorporating communication compression into Distributed/Federated Learning. However, for non-convex smooth problems, the existing analysis of EF is provided either for the single-node case or relies on restrictive assumptions such as boundedness of the gradient/compression error or boundedness of the data heterogeneity (gradient dissimilarity) ([Stich et al., 2018](#); [Stich & Karimireddy, 2019](#); [Karimireddy et al., 2019](#); [Koloskova et al., 2019](#); [Beznosikov et al., 2023](#); [Tang et al., 2019](#); [Xie et al., 2020](#); [Sahu et al., 2021](#)). Moreover, the convergence bounds for EF also depend on the data heterogeneity, which is not an artifact of the analysis as illustrated in the experiments on strongly convex problems [Gorbunov et al. \(2020b\)](#). [Richtárik et al. \(2021\)](#) address this limitation and propose a new version of Error Feedback called EF21. However, the existing analysis of EF21-SGD requires the usage of large batch sizes to achieve any predefined accuracy ([Fatkhullin et al., 2021](#)). It turns out that the large batch size requirement is unavoidable for EF21-SGD to converge, but this issue can be fixed using momentum ([Fatkhullin et al., 2024](#)). Momentum is also helpful in the decentralized extensions of Error Feedback ([Yau & Wai, 2022](#); [Huang et al., 2023](#); [Islamov et al., 2024](#)).

2 NON-CONVERGENCE OF Clip-SGD AND Clip21-SGD

We start with a discussion of the key limitations of Clip-SGD (Algorithm 1) and Clip21-SGD (Algorithm 2) – their potential non-convergence.

Algorithm 1 Clip-SGD ([Abadi et al., 2016](#))

Input: $x^0 \in \mathbb{R}^d$, stepsize $\gamma > 0$, clipping parameter $\tau > 0$

```

1: for  $t = 0, \dots, T - 1$  do
2:   for  $i = 1, \dots, n$  in parallel do
3:      $g_i^t = \text{clip}_\tau(\nabla f_i(x^t, \xi_i^t))$ 
4:   end for
5:    $g^t = \frac{1}{n} \sum_{i=1}^n g_i^t$ 
6:    $x^{t+1} = x^t - \gamma g^t$ 
7: end for

```

Algorithm 2 Clip21-SGD ([Khirirat et al., 2023](#))

Input: $x^0, g^0 \in \mathbb{R}^d$, stepsize $\gamma > 0$, clipping parameter $\tau > 0$

```

1: Initialize  $g_i^0 = g^0$  for all  $i \in [n]$ 
2: for  $t = 0, \dots, T - 1$  do
3:    $x^{t+1} = x^t - \gamma g^t$ 
4:   for  $i = 1, \dots, n$  in parallel do
5:      $c_i^{t+1} = \text{clip}_\tau(\nabla f_i(x^{t+1}, \xi_i^{t+1}) - g_i^t)$ 
6:      $g_i^{t+1} = g_i^t + c_i^{t+1}$ 
7:   end for
8:    $g^{t+1} = g^t + \frac{1}{n} \sum_{i=1}^n c_i^{t+1}$ 
9: end for

```

We start by restating the example from ([Chen et al., 2020](#)) illustrating the potential non-convergence of Clip-SGD even when full gradients are computed on clients (Clip-GD).

Example 1 (Non-Convergence of Clip-GD ([Chen et al., 2020](#))). Let $n = 2$, $d = 1$, and $f_1(x) = \frac{1}{2}(x - 3)^2$, $f_2(x) = \frac{1}{2}(x + 3)^2$ in problem (1) having a unique solution $x^* = 0$. Consider Clip-GD with $\tau = 1$ applied to this problem. If for some t_0 we have $x^{t_0} \in [-2, 2]$ in Clip-GD, then $g^t = 0$ and $x^t = x^{t_0}$ for any $t \geq t_0$, which can be seen via direct calculations. In particular, for any $x^0 \in [-2, 2]$, the method does not move away from x^0 .

³The proof of the DP guarantee by [Khirirat et al. \(2023\)](#) relies on the condition for some $C > 1$ and $\nu, \sigma_\omega \geq 0$ that implies $\min\{\nu^2, \sigma_\omega^2\} \geq C \max\{\nu^2, \sigma_\omega^2\}$. The latter one holds if and only if $\nu = \sigma_\omega = 0$, which means that no noise is added to the method since σ_ω^2 is the variance of DP-noise.

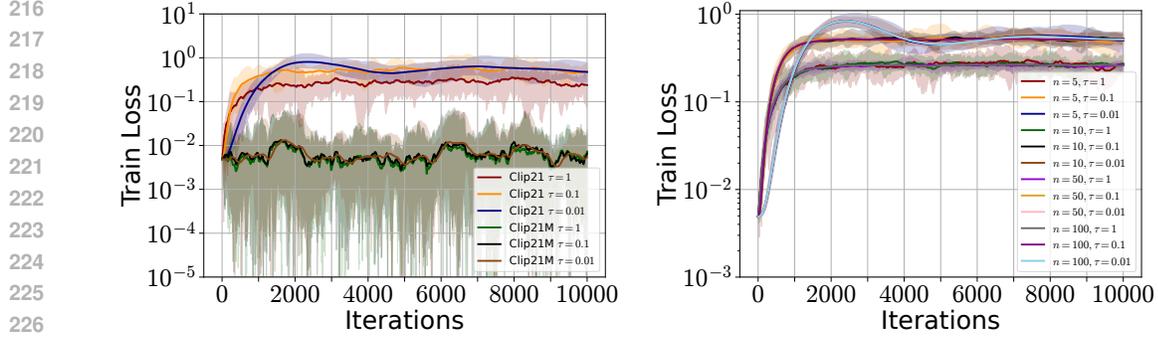


Figure 1: **Left:** behavior of stochastic Clip21-SGD and Clip21-SGDM without DP noise (see Algorithm 3) initialized at $x^0 = (0, -0.07)^\top$, with stepsize $\gamma = 1/\sqrt{T}$ where $T = 10^4$, i.e., close to the solution and small stepsize. We observe that Clip21-SGD escapes the good neighborhood of the solution for the problem from Theorem 1 with $n = 1$, $L = 2$, $\sigma = 5$, and varying $\tau \in \{1, 0.1, 0.01\}$. In contrast, Clip21-SGDM remains stable around the solution. **Right:** convergence of Clip21-SGD does not improve with the increase of n for the same problem.

To address the non-convergence of Clip-GD, Khirirat et al. (2023) propose Clip21-GD that applies the clipping operator to the difference between $\nabla f_i(x^{t+1})$ and the shift g_i^t , which is designed to approximate $\nabla f_i(x^t)$. In the deterministic case, this strategy ensures that after a certain number of steps, clipping turns off on all clients since $\|\nabla f_i(x^{t+1}) - g_i^t\|$ becomes smaller than τ for all $i \in [n]$ eventually. However, when workers compute stochastic gradients instead of the full gradients, Clip21-SGD can be non-convergent as well. To illustrate this, we consider the ideal version of Clip21-SGD with stochastic gradients, i.e., instead of g_i^t , we use $\nabla f_i(x^{t+1})$ as a shift:

$$x^{t+1} = x^t - \gamma g^t, \quad g^t = \frac{1}{n} \sum_{i=1}^n g_i^t, \quad (6)$$

$$g_i^{t+1} = \nabla f_i(x^{t+1}) + \text{clip}_\tau(\nabla f_i(x^{t+1}, \xi_i^{t+1}) - \nabla f_i(x^{t+1})) \quad (7)$$

The next theorem shows that even this (ideal) version of stochastic Clip21-SGD fails to converge even for a simple quadratic problem with sub-Gaussian noise.

Theorem 1. *Let $L, \sigma > 0$, $0 < \gamma \leq 1/L$, $n = 1$. There exists a convex, L -smooth problem, clipping parameter $\tau < 3\sigma\sqrt{3}/10$, and an unbiased stochastic gradient satisfying Assumption 2 such that the method (6) is run with a stepsize γ and clipping parameter τ , then for all $x^0 \in \{(0, x_{(2)}^0) \in \mathbb{R}^2 \mid x_{(2)}^0 < 0\}$ we have*

$$\mathbb{E} [\|\nabla f(x^T)\|^2] \geq \frac{1}{2} \min \left\{ \|\nabla f(x^0)\|^2, \frac{\tau^2}{45} \right\}.$$

Moreover, fix $0 < \varepsilon < L/\sqrt{2}$ and $x^0 = (0, -1)^\top$. Let the sub-Gaussian variance of stochastic gradients is bounded by σ^2/B where B is a batch size. If $B < 27\sigma^2/(60\varepsilon^2)$ and $\tau \geq \varepsilon/(3\sqrt{10})$, then we have $\mathbb{E} [\|\nabla f(x^T)\|^2] > \varepsilon^2$ for all $T > 0$.

We also illustrate the above result with simple numerical experiments reported in Figure 1. The left figure shows that Clip21-SGD diverges from the initial function sub-optimality level while the right one demonstrates non-improvement with the number of workers n — one of the desired properties of algorithms for FL.

3 Clip21-SGDM: NEW METHOD AND THEORETICAL RESULTS

This section introduces Clip21-SGDM (Algorithm 3), a novel distributed method with clipping that can be viewed as an enhanced version of Clip21-SGD, integrating momentum and DP-noise. That is, to control the noise coming from the stochastic gradients, we introduce momentum buffers $\{v_i^t\}_{i \in [n]}$ on the clients and $\text{clip} \{v_i^{t+1} - g_i^t\}_{i \in [n]}$ in contrast to the stochastic version of Clip21-SGD that

Algorithm 3 Clip21-SGDM

```

270 1: Input:  $x^0, g^0, v^0 \in \mathbb{R}^d$  (by default  $g^0 = v^0 = 0$ ), momentum parameter  $\beta \in (0, 1]$ , stepsize
271  $\gamma > 0$ , clipping parameter  $\tau > 0$ , DP-variance parameter  $\sigma_\omega^2 \geq 0$ 
272 2: Set  $g_i^0 = g^0$  and  $v_i^0 = v^0$  for all  $i \in [n]$ 
273 3: for  $t = 0, \dots, T - 1$  do
274 4:    $x^{t+1} = x^t - \gamma g^t$ 
275 5:   for  $i = 1, \dots, n$  do
276 6:      $v_i^{t+1} = (1 - \beta)v_i^t + \beta \nabla f_i(x^{t+1}, \xi_i^{t+1})$ 
277 7:      $\omega_i^{t+1} \sim \mathcal{N}(0, \sigma_\omega^2 \mathbf{I})$  only for DP version
278 8:      $c_i^{t+1} = \text{clip}_\tau(v_i^{t+1} - g_i^t) + \omega_i^{t+1}$ 
279 9:      $g_i^{t+1} = g_i^t + c_i^{t+1} - \omega_i^{t+1} = g_i^t + \text{clip}_\tau(v_i^{t+1} - g_i^t)$ 
280 10:   end for
281 11:    $g^{t+1} = g^t + \frac{1}{n} \sum_{i=1}^n c_i^{t+1}$ 
282 12: end for

```

applies clipping to potentially noisier vectors $\{\nabla f_i(x^{t+1}, \xi_i^{t+1}) - g_i^t\}_{i \in [n]}$. Moreover, similarly to Clip21-SGD – which can be seen as EF21 (Richtárik et al., 2021) where the compression operator is replaced by clipping – Clip21-SGDM can also be interpreted as EF21M (Fatkhullin et al., 2024) with the same replacement. However, both EF21 and EF21M rely on the contractiveness property of the compression operator $\mathcal{C}(x)$, i.e., the (randomized) mapping $\mathcal{C} : \mathbb{R}^d \rightarrow \mathbb{R}^d$ should satisfy

$$\mathbb{E} [\|\mathcal{C}(x) - x\|^2] \leq (1 - \nu)\|x\|^2 \quad \text{for some } \nu \in (0, 1], \quad (8)$$

where the expectation is w.r.t. the randomness of \mathcal{C} . As shown and discussed by Khirirat et al. (2023), clipping satisfies a condition that resembles (8) namely

$$\|\text{clip}_\tau(x) - x\|^2 \leq \begin{cases} 0, & \text{if } \|x\| \leq \tau, \\ \left(1 - \frac{\tau}{\|x\|}\right)^2 \|x\|^2, & \text{if } \|x\| > \tau, \end{cases} \quad (9)$$

but there is a significant difference: if $\|x\| > \tau$, the contraction factor is dependent of x and can be arbitrarily close to 1. To circumvent this issue, Khirirat et al. (2023) prove via induction that for all iterates of Clip21-SGD, the vectors $\nabla f_i(x^{t+1}) - g_i^t$ have norms bounded by some constant depending on the starting point. We show that a similar statement holds for Clip21-SGDM when the clients compute full-batched gradients and no DP-noise is added, and we start our analysis with this important special case. We also present the results in the stochastic case with and without DP noise.

3.1 ANALYSIS IN THE DETERMINISTIC CASE

The next result derives a convergence rate for Clip21-SGDM when $\nabla f_i(x^{t+1}, \xi_i^{t+1}) \equiv \nabla f_i(x^t)$ almost surely, i.e., Assumption 2 holds with $\sigma = 0$.

Theorem 2 (Simplified). *Let Assumptions 1 and 2 with $\sigma = 0$ hold. Let $B := \max_i \|\nabla f_i(x^0)\| > 3\tau$ and $\Delta \geq f(x^0) - f^*$. Then there exists a stepsize $\gamma \leq 1/12L$ and momentum parameter $\beta = 4L\gamma$ such that the iterates of Clip21-SGDM (Algorithm 3) converge with the rate*

$$\frac{1}{T} \sum_{t=0}^{T-1} \|\nabla f(x^t)\|^2 \leq \mathcal{O}\left(\frac{L\Delta}{T}\right). \quad (10)$$

Moreover, after at most $2B/\tau$ iterations, the clipping will eventually be turned off for all workers.

Proof sketch. The proof of Theorem 2 (and all following ones) relies on the same Lyapunov function that is used by Fatkhullin et al. (2024) in the analysis of EF21M:

$$\Phi^t := f(x^t) - f^* + \frac{\gamma}{\eta} \frac{1}{n} \sum_{i=1}^n \|g_i^t - v_i^t\|^2 + \frac{4\gamma\beta}{\eta^2} \frac{1}{n} \sum_{i=1}^n \|v_i^t - \nabla f_i(x^t)\|^2 + \frac{\gamma}{\beta} \|v^t - \nabla f(x^t)\|^2. \quad (11)$$

In the definition of Φ^t , the only parameter that was not introduced earlier in the paper is η , and it hides the main technical difficulty of the proof. That is, by induction we prove that $\|v_i^{t+1} -$

$g_i^t\| \leq \tau/\eta$ for some η defined in the proof. This bound is essential in deriving a descent of each term in the Lyapunov function. In view of (9) and (8), this allows us to consider clipping as a contractive compression operator for vectors $v_i^{t+1} - g_i^t$ generated by the method, and also this allows us to use the same Lyapunov function as in the analysis of EF21M. We defer the detailed proof to Appendix D. \square

The above result establishes a $\mathcal{O}(1/T)$ convergence rate that is optimal for non-convex smooth first-order optimization (Carmon et al., 2020; 2021). This result matches the one obtained by Khirirat et al. (2023), and, in particular, similarly to Clip21-SGD, Clip21-SGDM turns off clipping on each client after a finite number of steps t satisfying $\|v_i^{t+1} - g_i^t\| \leq \tau$. We also emphasize that Theorem 2 holds without bounded heterogeneity/gradient assumption. In contrast, even with bounded heterogeneity/gradient assumption, many existing convergence results in the non-convex case (Liu et al., 2022; Zhang et al., 2022; Li & Chi, 2023; Allouah et al., 2024) do not recover the $\mathcal{O}(1/T)$ rate in the noiseless regime.

3.2 ANALYSIS IN THE STOCHASTIC CASE WITHOUT DP-NOISE

Next, we turn to the stochastic setting where each worker has access to local gradient estimators satisfying Assumption 2. For simplicity, we first consider the case when no DP noise is added.

Theorem 3. *Let Assumptions 1 and 2 hold and $\alpha \in (0, 1)$. Let $\tilde{B} := \max_i \|\nabla f_i(x^0)\| > 3\tau$ and $\Delta \geq \Phi^0$. Then there exists a stepsize γ and momentum parameter β such that the iterates of Clip21-SGDM (Algorithm 3) satisfy with probability at least $1 - \alpha$*

$$\frac{1}{T} \sum_{t=0}^{T-1} \|\nabla f(x^t)\|^2 \leq \tilde{\mathcal{O}} \left(\frac{L\Delta}{T} + \frac{\sigma(\sqrt{L\Delta} + \tilde{B} + \sigma)}{\sqrt{Tn}} \right), \quad (12)$$

where $\tilde{\mathcal{O}}$ hides constant and logarithmic factors, and higher order terms that decrease in T .

Proof sketch. The core of the proof is similar to the one of Theorem 2. However, in contrast to the deterministic case, the vectors $v_i^{t+1} - g_i^t$ are stochastic, meaning that under Assumption 2, they can have arbitrarily large norms. Therefore, we focus on the high-probability analysis and prove by induction that the vectors $v_i^{t+1} - g_i^t$ are bounded *with high probability*, meaning that clipping can be seen as a contractive compressor with high probability for the vectors $v_i^{t+1} - g_i^t$ generated by the method. The proof is also based on a refined estimation of sums of martingale difference sequences; see the details in Appendix G. \square

This result demonstrates that Clip21-SGDM achieves an optimal $\mathcal{O}(1/\sqrt{nT})$ (Arjevani et al., 2023) rate in the stochastic setting. In contrast to the previous works establishing similar rates (Liu et al., 2022; Noble et al., 2022; Allouah et al., 2024), our result does not rely on the boundedness of the gradients or data heterogeneity. Moreover, when $\sigma = 0$ (no stochastic noise), the rate from (12) becomes $\mathcal{O}(1/T)$, recovering the one given by Theorem 2.

3.3 ANALYSIS IN THE STOCHASTIC CASE WITH DP-NOISE

Finally, we provide the convergence result for Clip21-SGDM with DP-noise.

Theorem 4. *Let Assumptions 1 and 2 hold and $\alpha \in (0, 1)$. Let $\Delta \geq \Phi^0$. Then there exists a stepsize γ and momentum parameter β such that the iterates of Clip21-SGDM (Algorithm 3) with the DP-noise variance σ_ω^2 with probability at least $1 - \alpha$ satisfy*

$$\frac{1}{T} \sum_{t=0}^{T-1} \|\nabla f(x^t)\|^2 \leq \tilde{\mathcal{O}} \left(\frac{L\Delta\sqrt{d}\sigma_\omega}{\sqrt{Tn}\tau} + \frac{(L\Delta)^{1/6}\sigma_\omega^{5/3}}{T^{1/6}n^{5/6}} + \frac{(L\Delta)^{4/9}\sigma_\omega^{5/9}d^{5/18}\sigma_\omega^{5/9}}{T^{4/9}n^{5/9}} \right), \quad (13)$$

where $\tilde{\mathcal{O}}$ hides constant and logarithmic factors, and higher order terms decreasing in T .

In the special case of local Differential Privacy, the noise level has to be chosen in a specific way. In this setting, we obtain the following privacy-utility trade-off.

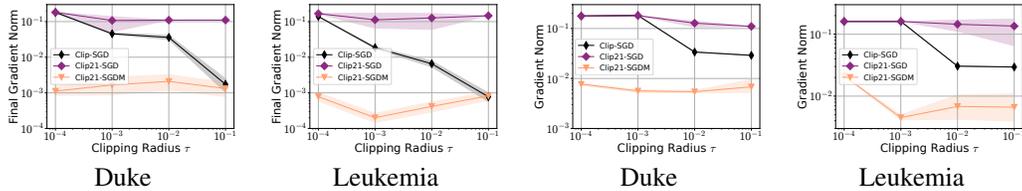


Figure 2: Comparison of tuned Clip-SGD, Clip21-SGD, and Clip21-SGDM on logistic regression with non-convex regularization for various clipping radii τ with mini-batch (**two left**) and Gaussian-added (**two right**) stochastic gradients. The final gradient norm is averaged over the last 100 iterations.

Corollary 1. Let Assumptions 1 and 2 hold and $\alpha \in (0, 1)$. Let $\Delta \geq \Phi^0$ and σ_ω be chosen as $\sigma_\omega = \Theta\left(\frac{\tau}{\varepsilon} \sqrt{T \log \frac{1}{\delta}}\right)$. Then there exists a stepsize γ and momentum parameter β such that the iterates of Clip21-SGDM (Algorithm 3) with probability at least $1 - \alpha$ satisfy local (ε, δ) -DP and

$$\frac{1}{T} \sum_{t=0}^{T-1} \|\nabla f(x^t)\|^2 \leq \tilde{O}\left(\frac{L\Delta\sqrt{d}}{\sqrt{ne}}\right), \quad (14)$$

where \tilde{O} hides constant and logarithmic factors, and terms decreasing in T .

To obtain local (ε, δ) -DP guarantees we follow Theorem 1 in (Abadi et al., 2016). This privacy-utility trade-off matches the known lower bound for locally private algorithms (Duchi et al., 2018). Overall, Theorems 2 and 3 and Corollary 1 indicate that Clip21-SGDM achieves optimal convergence rates in both deterministic and stochastic regime, and also has an optimal privacy-utility trade-off. These results are derived without assuming the boundedness of the gradients/data heterogeneity.

4 EXPERIMENTS

In this section, we provide an empirical evaluation of the proposed algorithm against baselines such as Clip21-SGD (Khirirat et al., 2023) and Clip-SGD. The learning rate and momentum (for Clip21-SGDM) are tuned in all experiments. We refer to Appendix H for the detailed description of tuning.

4.1 STOCHASTIC SETTING

First, we test the convergence of Clip-SGD, Clip21-SGD, and the proposed Clip21-SGDM algorithms with stochastic gradients for various clipping radii τ on several workloads. These results demonstrate the significance of using the momentum technique to achieve better performance.

4.1.1 NON-CONVEX LOGISTIC REGRESSION

We demonstrate the performance of all algorithms without adding noise for privacy but with stochastic gradients. We consider two cases: adding Gaussian noise to full local gradient $\nabla f_i(x)$ and mini-batch stochastic gradient. We conduct experiments on logistic regression with non-convex regularization, namely, $f_i(x) = \frac{1}{m} \sum_{j=1}^m \log(1 + \exp(-b_{ij} a_{ij}^\top x)) + \lambda \sum_{l=1}^d \frac{x_l^2}{1+x_l^2}$ which is a typical problem considered in previous works (Khirirat et al., 2023; Li & Chi, 2023). We use the Duke and Leukemia LibSVM (Chang & Lin, 2011) datasets.

We plot the gradient norm averaged across the last 100 iterations and 3 different runs in Figure 2. The results demonstrate the resilience of Clip21-SGDM to the choice of the clipping radius τ : it achieves a smaller or similar gradient norm compared to two other algorithms over all values of τ . This is especially visible when the clipping radius τ is small. These experimental findings align with the theoretical results presented in this work. Besides, the convergence plots are presented in Figure 7. The results demonstrate faster convergence for Clip21-SGDM than that of Clip21-SGD and Clip-SGD.

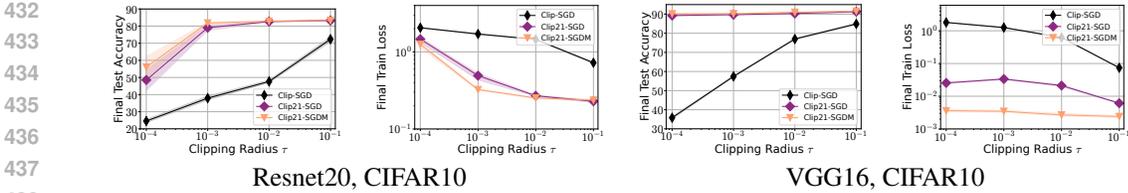


Figure 3: Comparison of tuned Clip-SGD, Clip21-SGD, and Clip21-SGDM on training Resnet20 (two left) and VGG16 (two right) models on CIFAR10 dataset where the clipping is applied globally.

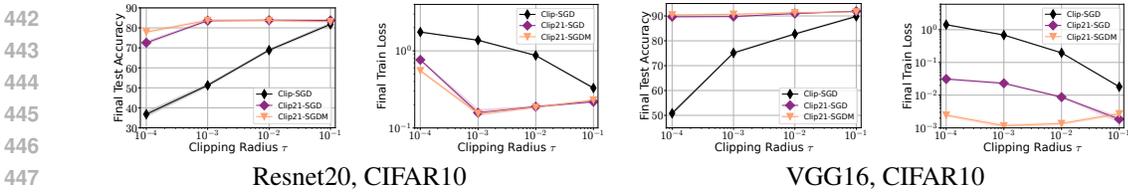


Figure 4: Comparison of tuned Clip-SGD, Clip21-SGD, and Clip21-SGDM on training Resnet20 (two left) and VGG16 (two right) models on CIFAR10 dataset where the clipping is applied layer-wise.

4.1.2 TRAINING RESNET20 AND VGG16

Next, we conduct experiments in training Resnet20 (He et al., 2016) and VGG16 (Simonyan & Zisserman, 2014) models on CIFAR10 dataset (Krizhevsky et al., 2009)⁴. The results are averaged across 3 different random seeds and shown in Figure 3 (the clipping operator is applied on all weights simultaneously) and Figure 4 (the clipping operator is applied layer-wise). We plot the test accuracy and train loss at the last point of the training. The results show that the performance of Clip-SGD consistently deteriorates as the clipping radius τ decreases, while Clip21-SGD and Clip21-SGDM are more stable to the changes of τ . Moreover, Clip21-SGDM outperforms Clip21-SGD for small values of τ reaching smaller train loss and larger test accuracy that supports the theoretical claims of this paper. For the convergence curves we refer to Figures 8 to 11.

4.2 ADDING GAUSSIAN NOISE FOR DP

In the second set of experiments, we test the performance of algorithms with additive Gaussian noise to preserve privacy. Since DP noise variance σ_ω typically scales with the clipping radius τ (e.g., see Corollary 1), we conduct the following set of experiments: we fix a noise-clipping ratio from $\{0.1, 1.0, 10\}$ for logistic regression and from $\{0.1, 0.3, 1.0, 3.0, 10.0\}$ for neural networks, and find such τ that gives the lowest final gradient norm, train loss, or test accuracy depending on the considered workload. The high values of the noise-clipping ratio correspond to stronger DP guarantees, while low values stand for weaker DP guarantees.

4.2.1 NON-CONVEX LOGISTIC REGRESSION

We provide the convergence results for non-convex logistic regression in Figure 5 where the gradient norm is averaged over the last 100 iterations and 5 random seeds. We demonstrate that Clip21-SGDM can achieve a smaller gradient norm for all values of the noise-clipping ratio than Clip-SGD. Besides, the performance of Clip21-SGD does not improve even if the noise-clipping ratio is small, demonstrating the importance of the use of momentum.

4.2.2 TRAINING NEURAL NETWORKS WITH DP NOISE

Next, we conduct experiments on training CNN and MLP models on MNIST dataset (Deng, 2012) varying the noise-clipping ratio. We highlight that it is a standard experiment setting considered in the literature on differential privacy (Papernot et al., 2020; Li & Chi, 2023; Allouah et al., 2024). The performance results are reported in Figure 6. We observe that no algorithm outperforms others

⁴We use the code base from (Horváth & Richtárik, 2020) with small modifications.

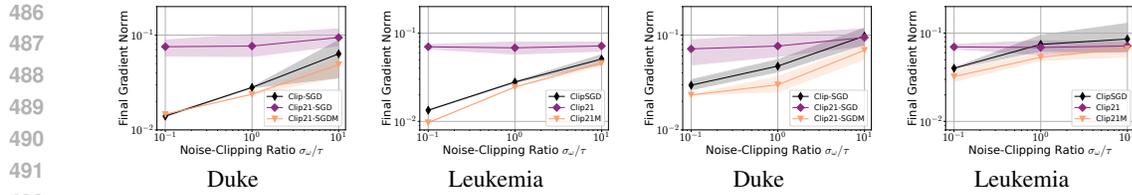


Figure 5: Comparison of tuned Clip-SGD, Clip21-SGD, Clip21-SGDM with mini-batch (**two left**) and Gaussian-added (**two right**) stochastic gradients and with additional DP-noise with variance σ_ω and varying noise-clipping ratio σ_ω/τ on non-convex logistic regression with non-convex regularization. The final gradient norm is averaged over the last 100 iterations.

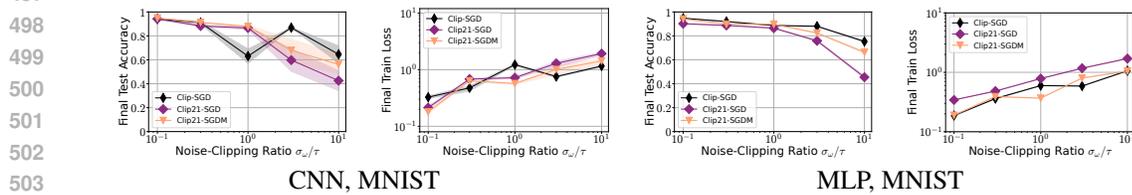


Figure 6: Comparison of tuned Clip-SGD, Clip21-SGD, and Clip21-SGDM on training CNN (**two left**) and MLP (**two right**) models on MNIST dataset varying the noise-clipping ratio where the clipping is applied globally.

across all values of the noise-clipping ratio in terms of the train loss. However, Clip-SGD typically attains smaller train loss than Clip21-SGDM for a large value of the noise-clipping ratio while Clip21-SGDM achieves smaller train loss Clip-SGD when that ratio is small.

5 CONCLUSION AND FUTURE WORK

In this work, we introduced a new method called Clip21-SGDM and proved that it achieves an optimal convergence rate and optimal privacy-utility trade-off without assuming boundedness of the gradients or boundedness of the data heterogeneity. Notably, several interesting directions remain unexplored. The first one is related to the generalization of the derived results to the case when stochastic gradients have heavy-tailed noise. Next, it would be interesting to study AdaGrad/Adam-type (Streeter & McMahan, 2010; Duchi et al., 2011; Kingma & Ba, 2014) versions of Clip21-SGDM due to their practical superiority over SGD in solving Deep Learning problems. Finally, it is important to extend the current analysis of Clip21-SGDM to the case when generalized smoothness is satisfied (Zhang et al., 2020b).

REFERENCES

- Martin Abadi, Andy Chu, Ian Goodfellow, H Brendan McMahan, Ilya Mironov, Kunal Talwar, and Li Zhang. Deep learning with differential privacy. In *Proceedings of the 2016 ACM SIGSAC conference on computer and communications security*, 2016. (Cited on pages 1, 4, and 8)
- Youssef Allouah, Anastasia Koloskova, Aymane El Firdoussi, Martin Jaggi, and Rachid Guerraoui. The privacy power of correlated noise in decentralized learning. *arXiv preprint arXiv:2405.01031*, 2024. (Cited on pages 3, 7, and 9)
- Yossi Arjevani, Yair Carmon, John C Duchi, Dylan J Foster, Nathan Srebro, and Blake Woodworth. Lower bounds for non-convex stochastic optimization. *Mathematical Programming*, 2023. (Cited on page 7)
- Borja Balle, James Bell, Adrià Gascón, and Kobbi Nissim. The privacy blanket of the shuffle model. In *Advances in Cryptology—CRYPTO 2019: 39th Annual International Cryptology Conference, Santa Barbara, CA, USA, August 18–22, 2019, Proceedings, Part II 39*, 2019. (Cited on page 3)

- 540 Aleksandr Beznosikov, Samuel Horváth, Peter Richtárik, and Mher Safaryan. On biased compres-
541 sion for distributed learning. *Journal of Machine Learning Research*, 2023. (Cited on page 4)
542
- 543 Keith Bonawitz, Vladimir Ivanov, Ben Kreuter, Antonio Marcedone, H Brendan McMahan, Sarvar
544 Patel, Daniel Ramage, Aaron Segal, and Karn Seth. Practical secure aggregation for privacy-
545 preserving machine learning. In *proceedings of the 2017 ACM SIGSAC Conference on Computer*
546 *and Communications Security*, 2017. (Cited on page 3)
- 547 Yair Carmon, John C Duchi, Oliver Hinder, and Aaron Sidford. Lower bounds for finding stationary
548 points i. *Mathematical Programming*, 2020. (Cited on pages 2 and 7)
549
- 550 Yair Carmon, John C Duchi, Oliver Hinder, and Aaron Sidford. Lower bounds for finding stationary
551 points ii: first-order methods. *Mathematical Programming*, 2021. (Cited on page 7)
- 552 TH Hubert Chan, Elaine Shi, and Dawn Song. Optimal lower bound for differentially private multi-
553 party aggregation. In *European Symposium on Algorithms*, 2012. (Cited on page 3)
- 554 Chih-Chung Chang and Chih-Jen Lin. Libsvm: A library for support vector machines. *ACM trans-*
555 *actions on intelligent systems and technology (TIST)*, 2011. (Cited on page 8)
556
- 557 Kamalika Chaudhuri, Chuan Guo, and Mike Rabbat. Privacy-aware compression for federated data
558 analysis. In *Uncertainty in Artificial Intelligence*, 2022. (Cited on page 3)
- 559 Xiangyi Chen, Steven Z Wu, and Mingyi Hong. Understanding gradient clipping in private sgd:
560 A geometric perspective. *Advances in Neural Information Processing Systems*, 2020. (Cited on
561 pages 1, 3, and 4)
562
- 563 Ashok Cutkosky and Harsh Mehta. High-probability bounds for non-convex stochastic optimization
564 with heavy tails. *Advances in Neural Information Processing Systems*, 2021. (Cited on page 3)
- 565 Marina Danilova, Pavel Dvurechensky, Alexander Gasnikov, Eduard Gorbunov, Sergey Guminov,
566 Dmitry Kamzolov, and Innokentiy Shibaev. Recent theoretical advances in non-convex optimiza-
567 tion. In *High-Dimensional Optimization and Probability: With a View Towards Data Science*, pp.
568 79–163. Springer, 2022. (Cited on page 2)
- 569 Li Deng. The mnist database of handwritten digit images for machine learning research. *IEEE*
570 *Signal Processing Magazine*, 2012. (Cited on page 9)
571
- 572 John Duchi, Elad Hazan, and Yoram Singer. Adaptive subgradient methods for online learning and
573 stochastic optimization. *Journal of machine learning research*, 2011. (Cited on page 10)
- 574 John C Duchi, Michael I Jordan, and Martin J Wainwright. Minimax optimal procedures for locally
575 private estimation. *Journal of the American Statistical Association*, 2018. (Cited on pages 3 and 8)
576
- 577 Cynthia Dwork, Aaron Roth, et al. The algorithmic foundations of differential privacy. *Foundations*
578 *and Trends® in Theoretical Computer Science*, 2014. (Cited on pages 1 and 3)
- 579 Úlfar Erlingsson, Vitaly Feldman, Ilya Mironov, Ananth Raghunathan, Kunal Talwar, and
580 Abhradeep Thakurta. Amplification by shuffling: From local to central differential privacy via
581 anonymity. In *Proceedings of the Thirtieth Annual ACM-SIAM Symposium on Discrete Algo-*
582 *rithms*, 2019. (Cited on page 3)
- 583 Ilyas Fatkhullin, Igor Sokolov, Eduard Gorbunov, Zhize Li, and Peter Richtárik. Ef21 with
584 bells & whistles: Practical algorithmic extensions of modern error feedback. *arXiv preprint*
585 *arXiv:2110.03294*, 2021. (Cited on page 4)
586
- 587 Ilyas Fatkhullin, Alexander Tyurin, and Peter Richtárik. Momentum provably improves error feed-
588 back! *Advances in Neural Information Processing Systems*, 2024. (Cited on pages 4 and 6)
- 589 Robin C Geyer, Tassilo Klein, and Moin Nabi. Differentially private federated learning: A client
590 level perspective. *arXiv preprint arXiv:1712.07557*, 2017. (Cited on page 3)
591
- 592 Saeed Ghadimi and Guanghui Lan. Optimal stochastic approximation algorithms for strongly con-
593 vex stochastic composite optimization i: A generic algorithmic framework. *SIAM Journal on*
Optimization, 2012. (Cited on page 2)

- 594 Eduard Gorbunov, Darina Dvinskikh, and Alexander Gasnikov. Optimal decentralized distributed
595 algorithms for stochastic convex optimization. *arXiv preprint arXiv:1911.07363*, 2019. (Cited on
596 page 17)
- 597 Eduard Gorbunov, Marina Danilova, and Alexander Gasnikov. Stochastic optimization with heavy-
598 tailed noise via accelerated gradient clipping. *Advances in Neural Information Processing Sys-*
599 *tems*, 2020a. (Cited on page 3)
- 600 Eduard Gorbunov, Dmitry Kovalev, Dmitry Makarenko, and Peter Richtárik. Linearly converging
601 error compensated sgd. *Advances in Neural Information Processing Systems*, 2020b. (Cited on
602 page 4)
- 603 Eduard Gorbunov, Abdurakhmon Sadiev, Marina Danilova, Samuel Horváth, Gauthier Gidel, Pavel
604 Dvurechensky, Alexander Gasnikov, and Peter Richtárik. High-probability convergence for com-
605 posite and distributed stochastic minimization and variational inequalities with heavy-tailed noise.
606 In *Proceedings of the 41st International Conference on Machine Learning*, 2024. (Cited on page 4)
- 607 Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recog-
608 nition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2016.
609 (Cited on page 9)
- 610 Mahmoud Hegazy, Rémi Leluc, Cheuk Ting Li, and Aymeric Dieuleveut. Compression with exact
611 error distribution for federated learning. In *International Conference on Artificial Intelligence
612 and Statistics*, 2024. (Cited on page 3)
- 613 Samuel Horváth and Peter Richtárik. A better alternative to error feedback for communication-
614 efficient distributed learning. *arXiv preprint arXiv:2006.11077*, 2020. (Cited on page 9)
- 615 Xinmeng Huang, Ping Li, and Xiaoyun Li. Stochastic controlled averaging for federated learning
616 with communication compression. *arXiv preprint arXiv:2308.08165*, 2023. (Cited on page 4)
- 617 Rustem Islamov, Yuan Gao, and Sebastian U Stich. Near optimal decentralized optimization with
618 compression and momentum tracking. *arXiv preprint arXiv:2405.20114*, 2024. (Cited on page 4)
- 619 Peter Kairouz, H Brendan McMahan, Brendan Avent, Aurélien Bellet, Mehdi Bennis, Arjun Nitin
620 Bhagoji, Kallista Bonawitz, Zachary Charles, Graham Cormode, Rachel Cummings, et al. Ad-
621 vances and open problems in federated learning. *Foundations and trends® in machine learning*,
622 2021. (Cited on page 1)
- 623 Sai Praneeth Karimireddy, Quentin Rebjock, Sebastian Stich, and Martin Jaggi. Error feedback
624 fixes signsgd and other gradient compression schemes. In *International Conference on Machine
625 Learning*, 2019. (Cited on page 4)
- 626 Sai Praneeth Karimireddy, Satyen Kale, Mehryar Mohri, Sashank Reddi, Sebastian Stich, and
627 Ananda Theertha Suresh. Scaffold: Stochastic controlled averaging for federated learning. In
628 *International conference on machine learning*, 2020. (Cited on page 3)
- 629 Shiva Prasad Kasiviswanathan, Homin K Lee, Kobbi Nissim, Sofya Raskhodnikova, and Adam
630 Smith. What can we learn privately? *SIAM Journal on Computing*, 2011. (Cited on page 3)
- 631 Sarit Khirirat, Eduard Gorbunov, Samuel Horváth, Rustem Islamov, Fakhri Karray, and Peter
632 Richtárik. Clip21: Error feedback for gradient clipping. *arXiv preprint arXiv:2305.18929*, 2023.
633 (Cited on pages 1, 3, 4, 5, 6, 7, 8, and 17)
- 634 Diederik P Kingma and Jimmy Ba. Adam: a method for stochastic optimization. *arXiv preprint
635 arXiv:1412.6980*, 2014. (Cited on page 10)
- 636 Anastasiia Koloskova, Tao Lin, Sebastian Urban Stich, and Martin Jaggi. Decentralized deep learn-
637 ing with arbitrary communication compression. In *Proceedings of the 8th International Confer-
638 ence on Learning Representations*, 2019. (Cited on page 4)
- 639 Jakub Konečný, H. Brendan McMahan, Felix Yu, Peter Richtárik, Ananda Theertha Suresh, and
640 Dave Bacon. Federated learning: Strategies for improving communication efficiency. In *NIPS
641 Private Multi-Party Machine Learning Workshop*, 2016. (Cited on page 1)

- 648 Nikita Kornilov, Ohad Shamir, Aleksandr Lobanov, Darina Dvinskikh, Alexander Gasnikov, Inno-
649 kentyi Shibaev, Eduard Gorbunov, and Samuel Horváth. Accelerated zeroth-order method for
650 non-smooth stochastic convex optimization problem with infinite variance. *Advances in Neural*
651 *Information Processing Systems*, 2024. (Cited on page 3)
- 652 Alex Krizhevsky, Vinod Nair, and Geoffrey Hinton. Learning multiple layers of features from tiny
653 images. *Scientific Report*, 2009. (Cited on page 9)
- 654
655 Boyue Li and Yuejie Chi. Convergence and privacy of decentralized nonconvex optimization with
656 gradient clipping and communication compression. *arXiv preprint arXiv:2305.09896*, 2023.
657 (Cited on pages 3, 7, 8, and 9)
- 658
659 Zhize Li, Hongyan Bao, Xiangliang Zhang, and Peter Richtárik. Page: A simple and optimal prob-
660 abilistic gradient estimator for nonconvex optimization. In *International conference on machine*
661 *learning*, 2021. (Cited on pages 3 and 17)
- 662
663 Mingrui Liu, Zhenxun Zhuang, Yunwen Lei, and Chunyang Liao. A communication-efficient dis-
664 tributed gradient clipping algorithm for training deep neural networks. *Advances in Neural Infor-*
665 *mation Processing Systems*, 2022. (Cited on pages 3 and 7)
- 666
667 Zijian Liu, Ta Duy Nguyen, Thien Hang Nguyen, Alina Ene, and Huy Nguyen. High probability
668 convergence of stochastic gradient methods. In *International Conference on Machine Learning*,
2023. (Cited on page 3)
- 669
670 Liam Madden, Emiliano Dall’Anese, and Stephen Becker. High probability convergence bounds
671 for non-convex stochastic gradient descent with sub-weibull noise. *Journal of Machine Learning*
672 *Research*, 2024. (Cited on page 2)
- 673
674 Grigory Malinovsky, Peter Richtárik, Samuel Horváth, and Eduard Gorbunov. Byzantine robustness
675 and partial participation can be achieved simultaneously: Just clip gradient differences. *arXiv*
preprint arXiv:2311.14127, 2023. (Cited on page 4)
- 676
677 LO Mangasarian. Parallel gradient distribution in unconstrained optimization. *SIAM Journal on*
678 *Control and Optimization*, 1995. (Cited on page 3)
- 679
680 Brendan McMahan, Eider Moore, Daniel Ramage, Seth Hampson, and Blaise Aguera y Arcas.
681 Communication-efficient learning of deep networks from decentralized data. In *Artificial intelli-*
gence and statistics, 2017a. (Cited on pages 1 and 3)
- 682
683 H Brendan McMahan, Daniel Ramage, Kunal Talwar, and Li Zhang. Learning differentially private
684 recurrent language models. *arXiv preprint arXiv:1710.06963*, 2017b. (Cited on page 3)
- 685
686 Konstantin Mishchenko, Eduard Gorbunov, Martin Takáč, and Peter Richtárik. Distributed learning
687 with compressed gradient differences. *arXiv preprint arXiv:1901.09269*, 2019. (Cited on page 4)
- 688
689 Angelia Nedic and Asuman Ozdaglar. Distributed subgradient methods for multi-agent optimiza-
690 tion. *IEEE Transactions on Automatic Control*, 2009. (Cited on page 3)
- 691
692 Arkadi Nemirovski, Anatoli Juditsky, Guanghui Lan, and Alexander Shapiro. Robust stochastic
693 approximation approach to stochastic programming. *SIAM Journal on optimization*, 2009. (Cited
694 on page 2)
- 695
696 Maxence Noble, Aurélien Bellet, and Aymeric Dieuleveut. Differentially private federated learn-
697 ing on heterogeneous data. In *Proceedings of The 25th International Conference on Artificial*
Intelligence and Statistics, 2022. (Cited on pages 3 and 7)
- 698
699 Francesco Orabona. A modern introduction to online learning. *arXiv preprint arXiv:1912.13213*,
2019. (Cited on page 17)
- 700
701 Nicolas Papernot, Steve Chien, Shuang Song, Abhradeep Thakurta, and Ulfar Erlingsson. Making
the shoe fit: Architectures, initializations, and tuning for learning with privacy. 2020. (Cited on
page 9)

- 702 Razvan Pascanu, Tomas Mikolov, and Yoshua Bengio. On the difficulty of training recurrent neural
703 networks. In *Proceedings of the 30th International Conference on International Conference on*
704 *Machine Learning-Volume 28*, 2013. (Cited on page 1)
705
- 706 Boris T Polyak. Some methods of speeding up the convergence of iteration methods. *Ussr compu-*
707 *tational mathematics and mathematical physics*, 1964. (Cited on page 2)
708
- 709 Peter Richtárik, Igor Sokolov, and Ilyas Fatkhullin. Ef21: A new, simpler, theoretically better, and
710 practically faster error feedback. In *Advances in Neural Information Processing Systems*, 2021.
711 (Cited on pages 1, 4, 6, 25, and 64)
- 712 Abdurakhmon Sadiev, Marina Danilova, Eduard Gorbunov, Samuel Horváth, Gauthier Gidel, Pavel
713 Dvurechensky, Alexander Gasnikov, and Peter Richtárik. High-probability bounds for stochas-
714 tic optimization and variational inequalities: the case of unbounded variance. In *International*
715 *Conference on Machine Learning*, 2023. (Cited on page 3)
- 716 Atal Sahu, Aritra Dutta, Ahmed M Abdelmoniem, Trambak Banerjee, Marco Canini, and Panos
717 Kalnis. Rethinking gradient sparsification as total error minimization. *Advances in Neural Infor-*
718 *mation Processing Systems*, 2021. (Cited on page 4)
719
- 720 Frank Seide, Hao Fu, Jasha Droppo, Gang Li, and Dong Yu. 1-bit stochastic gradient descent and
721 its application to data-parallel distributed training of speech dnns. In *Interspeech*, 2014. (Cited on
722 pages 1 and 4)
- 723 Karen Simonyan and Andrew Zisserman. Very deep convolutional networks for large-scale image
724 recognition. *arXiv preprint arXiv:1409.1556*, 2014. (Cited on page 9)
725
- 726 Sebastian U Stich and Sai Praneeth Karimireddy. The error-feedback framework: Better rates for
727 sgd with delayed gradients and compressed communication. *arXiv preprint arXiv:1909.05350*,
728 2019. (Cited on page 4)
- 729 Sebastian U Stich, Jean-Baptiste Cordonnier, and Martin Jaggi. Sparsified sgd with memory. *Ad-*
730 *vances in neural information processing systems*, 2018. (Cited on page 4)
731
- 732 Matthew Streeter and H Brendan McMahan. Less regret via online conditioning. *arXiv preprint*
733 *arXiv:1002.4862*, 2010. (Cited on page 10)
734
- 735 Hanlin Tang, Chen Yu, Xiangru Lian, Tong Zhang, and Ji Liu. Doublesqueeze: Parallel stochastic
736 gradient descent with double-pass error-compensated compression. In *International Conference*
737 *on Machine Learning*, 2019. (Cited on page 4)
- 738 Roman Vershynin. *High-dimensional probability: An introduction with applications in data science*.
739 Cambridge University Press, 2018. (Cited on page 2)
740
- 741 Jianyu Wang, Zachary Charles, Zheng Xu, Gauri Joshi, H Brendan McMahan, Maruan Al-Shedivat,
742 Galen Andrew, Salman Avestimehr, Katharine Daly, Deepesh Data, et al. A field guide to feder-
743 ated optimization. *arXiv preprint arXiv:2107.06917*, 2021. (Cited on page 1)
- 744 Kang Wei, Jun Li, Ming Ding, Chuan Ma, Howard H Yang, Farhad Farokhi, Shi Jin, Tony QS Quek,
745 and H Vincent Poor. Federated learning with differential privacy: Algorithms and performance
746 analysis. *IEEE transactions on information forensics and security*, 2020. (Cited on page 3)
747
- 748 Cong Xie, Shuai Zheng, Sanmi Koyejo, Indranil Gupta, Mu Li, and Haibin Lin. Cser:
749 Communication-efficient sgd with error reset. *Advances in Neural Information Processing Sys-*
750 *tems*, 2020. (Cited on page 4)
- 751 Chung-Yiu Yau and Hoi-To Wai. Docom: Compressed decentralized optimization with near-optimal
752 sample complexity. *arXiv preprint arXiv:2202.00255*, 2022. (Cited on page 4)
753
- 754 Bohang Zhang, Jikai Jin, Cong Fang, and Liwei Wang. Improved analysis of clipping algorithms for
755 non-convex optimization. In *Advances in Neural Information Processing Systems*, 2020a. (Cited
on page 3)

756 Jingzhao Zhang, Tianxing He, Suvrit Sra, and Ali Jadbabaie. Why gradient clipping accelerates
757 training: A theoretical justification for adaptivity. In *International Conference on Learning Rep-*
758 *resentations*, 2020b. (Cited on pages 3 and 10)

759
760 Jingzhao Zhang, Sai Praneeth Karimireddy, Andreas Veit, Seungyeon Kim, Sashank Reddi, Sanjiv
761 Kumar, and Suvrit Sra. Why are adaptive methods good for attention models? In *Advances in*
762 *Neural Information Processing Systems*, 2020c. (Cited on page 3)

763
764 Xinwei Zhang, Xiangyi Chen, Mingyi Hong, Zhiwei Steven Wu, and Jinfeng Yi. Understanding
765 clipping for federated learning: Convergence and client-level differential privacy. In *International*
766 *Conference on Machine Learning, ICML 2022*, 2022. (Cited on pages 3 and 7)

767
768
769
770
771
772
773
774
775
776
777
778
779
780
781
782
783
784
785
786
787
788
789
790
791
792
793
794
795
796
797
798
799
800
801
802
803
804
805
806
807
808
809

810	CONTENTS	
811		
812	1 Introduction	1
813		
814	1.1 Problem Formulation and Assumptions	2
815	1.2 Related Work	3
816		
817	2 Non-Convergence of Clip-SGD and Clip21-SGD	4
818		
819	3 Clip21-SGDM: New Method and Theoretical Results	5
820		
821	3.1 Analysis in the Deterministic Case	6
822	3.2 Analysis in the Stochastic Case without DP-Noise	7
823	3.3 Analysis in the Stochastic Case with DP-Noise	7
824		
825		
826	4 Experiments	8
827		
828	4.1 Stochastic Setting	8
829	4.1.1 Non-convex Logistic Regression	8
830	4.1.2 Training Resnet20 and VGG16	9
831	4.2 Adding Gaussian Noise for DP	9
832	4.2.1 Non-convex Logistic Regression	9
833	4.2.2 Training Neural Networks with DP Noise	9
834		
835		
836		
837	5 Conclusion and Future Work	10
838		
839	A Notation	17
840		
841	B Useful Lemmas	17
842		
843	C Proof of Theorem 1	18
844		
845	D Proof of Theorem 2	19
846		
847	E Proof of Theorem 4	26
848		
849	F Proof of Corollary 1	48
850		
851	G Proof of Theorem 3	48
852		
853	H Experiments Details and More	65
854		
855	H.1 Experiments with Logistic Regression	65
856	H.1.1 Stochastic Setting Varying Clipping Radius	65
857	H.1.2 Stochastic Setting with additive DP Noise	65
858	H.2 Experiments with Neural Networks	66
859	H.2.1 Varying Clipping Radius τ	66
860	H.2.2 Adding Additive DP Noise	67
861		
862		
863		

864 A NOTATION

865 For shortness, in all proofs, we use the following notation

$$866 \delta^t := f(x^t) - f^*, \quad \tilde{V}^t := \frac{1}{n} \sum_{i=1}^n \|g_i^t - v_i^t\|^2,$$

$$867 \tilde{P}^t := \frac{1}{n} \sum_{i=1}^n \|v_i^t - \nabla f_i(x^t)\|^2, \quad P^t := \|v^t - \nabla f(x^t)\|^2,$$

$$868 R^t := \|x^{t+1} - x^t\|^2.$$

869 We additionally denote $\eta_i^t := \frac{\tau}{\|v_i^t - g_i^{t-1}\|}$ and $\eta := \frac{\tau}{B}$ where B is defined in each section (it is
870 different in deterministic and stochastic settings). Besides, we define $\mathcal{I}_t := \{i \in [n] \mid \|v_i^t - g_i^{t-1}\| >$
871 $\tau\}$.

872 We denote $\theta_i^t := \nabla f_i(x^t, \xi_i^t) - \nabla f_i(x^t)$. From Assumption 2, we have that θ_i^t is zero-centered
873 σ -sub-Gaussian random vector conditioned at x^t , namely

$$874 \mathbb{E}[\theta_i^t \mid x^t] = 0, \quad \Pr(\|\theta_i^t\| > b) \leq 2 \exp\left(-\frac{b^2}{2\sigma^2}\right) \quad \forall b > 0. \quad (15)$$

875 Moreover, we define an average of θ_i^t as $\theta^t := \frac{1}{n} \sum_{i=1}^n \theta_i^t$.

876 B USEFUL LEMMAS

877 **Lemma 1** (Lemma C.3 in (Gorbunov et al., 2019)). Let $\{\xi_k\}_{k=1}^N$ be the sequence of random vectors
878 with values in \mathbb{R}^n such that

$$879 \mathbb{E}[\xi_k \mid \xi_{k-1}, \dots, \xi_1] = 0 \text{ almost surely, } \forall k \in \{1, \dots, N\},$$

880 and set $S_N := \sum_{k=1}^N \xi_k$. Assume that the sequence $\{\xi_k\}_{k=1}^N$ are sub-Gaussian, i.e.

$$881 \mathbb{E}[\exp(\|\xi_k\|^2/\sigma_k^2 \mid \xi_{k-1}, \dots, \xi_1)] \leq \exp(1) \text{ almost surely, } \forall k \in \{1, \dots, N\},$$

882 where $\sigma_2, \dots, \sigma_N$ are some positive numbers. Then for all $\gamma \geq 0$

$$883 \Pr\left(\|S_N\| \geq (\sqrt{2} + 2\gamma) \sqrt{\sum_{k=1}^N \sigma_k^2}\right) \leq \exp(-\gamma^2/3). \quad (16)$$

884 **Lemma 2** (Modification of Lemma 1 in (Li et al., 2021)). Let $\delta^t = f(x^t) - f^*$, $x^{t+1} = x^t - \gamma g^t$,
885 and the stepsize $\gamma \leq \frac{1}{2L}$. Then

$$886 \delta^{t+1} \leq \delta^t - \frac{\gamma}{2} \|\nabla f(x^t)\|^2 - \frac{1}{4\gamma} \|x^{t+1} - x^t\|^2 + \gamma \frac{1}{n} \sum_{i=1}^n \|g_i^t - v_i^t\|^2 + \gamma \|v^t - \nabla f(x^t)\|^2. \quad (17)$$

887 **Lemma 3** (Lemma 4.1 in (Khirirat et al., 2023)). The clipping operator satisfies for any $x \in \mathbb{R}^d$

$$888 \|\text{clip}_\tau(x) - x\| \leq \max\{\|x\| - \tau, 0\}. \quad (18)$$

889 **Lemma 4** (Property of smooth functions). Let $\phi: \mathbb{R}^d \rightarrow \mathbb{R}$ be L -smooth and lower bounded by
890 $\phi^* \in \mathbb{R}$, i.e. $\phi(x) \geq \phi^*$ for any $x \in \mathbb{R}^d$. Then we have

$$891 \|\nabla \phi(x)\|^2 \leq 2L(\phi(x) - \phi^*). \quad (19)$$

892 *Proof.* It is a standard property of smooth functions. We refer to Theorem 4.23 of (Orabona, 2019).
893 \square

C PROOF OF THEOREM 1

Proof. The case $n = 1$. Let us consider the problem $f(x) = \frac{L}{2}\|x\|^2$. Let vectors $\{z_j\}_{j=1}^3$ are defined as

$$z_1 = \begin{pmatrix} 3 \\ 0 \end{pmatrix} \sqrt{\frac{3\sigma^2}{100}}, \quad z_2 = \begin{pmatrix} 0 \\ 4 \end{pmatrix} \sqrt{\frac{3\sigma^2}{100}}, \quad z_3 = \begin{pmatrix} -3 \\ -4 \end{pmatrix} \sqrt{\frac{3\sigma^2}{100}}.$$

Note that we have

$$\|z_1\|^2 = \frac{27\sigma^2}{100}, \quad \|z_2\|^2 = \frac{24\sigma^2}{50}, \quad \|z_3\|^2 = \frac{3\sigma^2}{4},$$

meaning that $\tau < \|z_i\|$ for all $i \in [3]$. We define the stochastic gradient as $\nabla f(x^t, \xi^t) = \nabla f(x^t) + \xi^t = Lx^t + \xi^t$ where ξ^t is picked uniformly at random from $\{z_1, z_2, z_3\}$. Simple calculations verify that Assumption 2 holds for such noise. Next, the update rule of the method (6) in the case $n = 1$ is

$$x^{t+1} = x^t - \gamma g^t = x^t - \gamma(\nabla f(x^t) + \text{clip}_\tau(\nabla f(x^t, \xi^t) - \nabla f(x^t))) = x^t - L\gamma x^t - \gamma \text{clip}_\tau(\xi^t).$$

Since $\tau < \|z_i\|$ for any $i \in \{1, 2, 3\}$ clipping is always active and we have

$$\begin{aligned} \mathbb{E}[\text{clip}_\tau(\xi^t)] &= \frac{1}{3} \text{clip}_\tau(z_1) + \frac{1}{3} \text{clip}_\tau(z_2) + \frac{1}{3} \text{clip}_\tau(z_3) \\ &= \frac{1}{3} \frac{\tau}{\|z_1\|} z_1 + \frac{1}{3} \frac{\tau}{\|z_2\|} z_2 + \frac{1}{3} \frac{\tau}{\|z_3\|} z_3 \\ &= \frac{1}{3} \frac{\tau}{\frac{3\sqrt{3}\sigma}{10}} \frac{\sigma\sqrt{3}}{10} \begin{pmatrix} 3 \\ 0 \end{pmatrix} + \frac{1}{3} \frac{\tau}{\frac{4\sqrt{3}\sigma}{10}} \frac{\sigma\sqrt{3}}{10} \begin{pmatrix} 0 \\ 4 \end{pmatrix} + \frac{1}{3} \frac{\tau}{\frac{5\sqrt{3}\sigma}{10}} \frac{\sigma\sqrt{3}}{10} \begin{pmatrix} -3 \\ -4 \end{pmatrix} \\ &= \frac{\tau}{9} \begin{pmatrix} 3 \\ 0 \end{pmatrix} + \frac{\tau}{12} \begin{pmatrix} 0 \\ 4 \end{pmatrix} + \frac{\tau}{15} \begin{pmatrix} -3 \\ -4 \end{pmatrix} \\ &= \underbrace{\frac{\tau}{15} \begin{pmatrix} 2 \\ 1 \end{pmatrix}}_{:=h}. \end{aligned}$$

Thus, we obtain

$$\begin{aligned} \mathbb{E}[x^T] &= (1 - L\gamma)\mathbb{E}[x^{T-1}] - \gamma\mathbb{E}[\text{clip}_\tau(\xi^t)] \\ &= (1 - L\gamma)\mathbb{E}[x^{T-1}] - \gamma h \\ &= (1 - L\gamma)^T x^0 - \gamma h \sum_{t=0}^{T-1} (1 - L\gamma)^{T-1-t} \\ &= (1 - L\gamma)^T \begin{pmatrix} 0 \\ x_{(2)}^0 \end{pmatrix} - \frac{\tau\gamma}{15} \begin{pmatrix} 2 \\ 1 \end{pmatrix} \frac{1 - (1 - L\gamma)^T}{1 - (1 - L\gamma)} \\ &= (1 - L\gamma)^T \begin{pmatrix} 0 \\ x_{(2)}^0 \end{pmatrix} - \frac{\tau}{15L} \begin{pmatrix} 2 \\ 1 \end{pmatrix} (1 - (1 - L\gamma)^T). \end{aligned}$$

Therefore, since $x_{(2)}^0 < 0$ we have

$$\begin{aligned} \mathbb{E}[\|\nabla f(x^T)\|^2] &= \mathbb{E}[\|Lx^T\|^2] \\ &= \mathbb{E}[\|Lx^T\|^2] + \mathbb{E}[\|Lx^T - \mathbb{E}[Lx^T]\|^2] \\ &\geq \mathbb{E}[\|Lx^T\|^2] \\ &= \frac{4\tau^2}{165} \left(1 - (1 - L\gamma)^T\right)^2 + L^2 \left((1 - L\gamma)^T x_{(2)}^0 - \frac{\tau}{15L} \left(1 - (1 - L\gamma)^T\right) \right)^2 \\ &\geq \frac{4\tau^2}{165} \left(1 - (1 - L\gamma)^T\right)^2 + (1 - L\gamma)^{2T} \|Lx^0\|^2 + \frac{\tau^2}{165} (1 - (1 - L\gamma)^T)^2 \\ &= \frac{\tau^2}{45} \left(1 - (1 - L\gamma)^T\right)^2 + (1 - L\gamma)^{2T} \|\nabla f(x^0)\|^2. \end{aligned}$$

Note that the function $a(1-x)^2 + x^2b \geq \frac{ab}{a+b}$. Applying this result for $a = \frac{\tau^2}{45}$, $b = \|\nabla f(x^0)\|^2$, and $x = (1-L\gamma)^T$ we get

$$\mathbb{E} [\|\nabla f(x^T)\|^2] \geq \frac{\frac{\tau^2}{45} \|\nabla f(x^0)\|^2}{\frac{\tau^2}{45} + \|\nabla f(x^0)\|^2} \geq \frac{1}{2} \min \left\{ \|\nabla f(x^0)\|^2, \frac{\tau^2}{45} \right\}.$$

The case $n > 1$. If $n > 1$ then we can consider a similar example where each client is quadratic $\frac{L}{2}\|x\|^2$ and the stochastic gradient is constructed as $\nabla f_i(x^t, \xi_i^t) = \nabla f_i(x^t) + \xi_i^t = Lx^t + \xi_i^t$ where ξ_i^t is sampled uniformly at random from vectors $\{z_1, z_2, z_3\}$ such that

$$z_1 = \begin{pmatrix} 3 \\ 0 \end{pmatrix} \sqrt{\frac{3\sigma^2}{100B}}, \quad z_2 = \begin{pmatrix} 0 \\ 4 \end{pmatrix} \sqrt{\frac{3\sigma^2}{100B}}, \quad z_3 = \begin{pmatrix} -3 \\ -4 \end{pmatrix} \sqrt{\frac{3\sigma^2}{100B}}.$$

Then, Assumption 2 is satisfied with σ^2/B . Therefore, if $x_{(2)}^0 = -1$, $\varepsilon < \frac{L}{\sqrt{2}}$, and $\tau \geq \frac{\varepsilon}{3\sqrt{10}}$, this implies that $B \leq \frac{243\sigma^2}{5\varepsilon^2} < \frac{27\sigma^2}{50\tau^2}$, and

$$\mathbb{E} [\|\nabla f(x^T)\|^2] \geq \frac{1}{2} \min \left\{ \|\nabla f(x^0)\|^2, \frac{\tau^2}{45} \right\} \geq \varepsilon^2.$$

□

D PROOF OF THEOREM 2

Lemma 5. Let each f_i be L -smooth. Then we have the following inequality

$$\|v_i^{t+1} - g_i^t\| \leq \max \{0, \|v_i^t - g_i^{t-1}\| - \tau\} + \beta L \gamma \|g^t\| + \beta \|\nabla f_i(x^t) - v_i^t\|. \quad (20)$$

Proof. We have

$$\begin{aligned} \|v_i^{t+1} - g_i^t\| &\stackrel{(i)}{=} \|(1-\beta)v_i^t + \beta \nabla f_i(x^{t+1}) - g_i^t\| \\ &\stackrel{(ii)}{\leq} \|v_i^t - g_i^t\| + \beta \|\nabla f_i(x^{t+1}) - v_i^t\| \\ &\stackrel{(iii)}{\leq} \max \{0, \|v_i^t - g_i^{t-1}\| - \tau\} + \beta \|\nabla f_i(x^{t+1}) - \nabla f_i(x^t)\| + \beta \|\nabla f_i(x^t) - v_i^t\| \\ &\stackrel{(iv)}{\leq} \max \{0, \|v_i^t - g_i^{t-1}\| - \tau\} + \beta L \|x^{t+1} - x^t\| + \beta \|\nabla f_i(x^t) - v_i^t\| \\ &\stackrel{(v)}{=} \max \{0, \|v_i^t - g_i^{t-1}\| - \tau\} + \beta L \gamma \|g^t\| + \beta \|\nabla f_i(x^t) - v_i^t\|, \end{aligned}$$

where (i) follows from the update rule of v_i^t in deterministic case; (ii) from triangle inequality; (iii) from the update rule of g_i^t , properties of clipping from Lemma 3, and triangle inequality; (iv) from L -smoothness of f_i ; (v) from the update rule of x^t . □

Lemma 6. Let each f_i be L -smooth and $\Delta \geq \Phi^0$. Assume that the following inequalities hold

1. $\|g^{t-1}\| \leq \sqrt{64L\Delta} + 3(B - \tau)$;
2. $\|\nabla f_i(x^{t-1}) - v_i^{t-1}\| \leq \sqrt{4L\Delta} + \frac{3}{2}(B - \tau)$;
3. $\|v_i^t - g_i^{t-1}\| \leq B$;
4. $\gamma \leq \frac{1}{12L}$;
5. $0 \leq \beta \leq \frac{1}{2}$;
6. $\Phi^t \leq \Delta$

Then we have

$$\|g^t\| \leq \sqrt{64L\Delta} + 3(B - \tau). \quad (21)$$

1026 *Proof.* We have

$$\begin{aligned}
1027 \quad \|g^t\| &\stackrel{(i)}{=} \left\| \frac{1}{n} \sum_{i=1}^n g_i^{t-1} + \text{clip}_\tau(v_i^t - g_i^{t-1}) \right\| \\
1028 \quad &= \left\| \frac{1}{n} \sum_{i=1}^n \nabla f_i(x^t) + (v_i^t - \nabla f_i(x^t)) + \text{clip}_\tau(v_i^t - g_i^{t-1}) - (v_i^t - g_i^{t-1}) \right\| \\
1029 \quad &\stackrel{(ii)}{\leq} \|\nabla f(x^t)\| + \frac{1}{n} \sum_{i=1}^n \|v_i^t - \nabla f_i(x^t)\| + \frac{1}{n} \sum_{i=1}^n \max\{0, \|v_i^t - g_i^{t-1}\| - \tau\}, \\
1030 \quad &
\end{aligned}$$

1031 where (i) follows from the update rule g_i^t ; (ii) from triangle inequality and clipping properties from
1032 Lemma 3. We continue to bound $\|g^t\|$ as follows

$$\begin{aligned}
1033 \quad \|g^t\| &\stackrel{(i)}{\leq} \|\nabla f(x^{t-1})\| + \|\nabla f(x^t) - \nabla f(x^{t-1})\| + \frac{1}{n} \sum_{i=1}^n (1-\beta) \|v_i^{t-1} - \nabla f_i(x^t)\| + B - \tau \\
1034 \quad &\stackrel{(ii)}{\leq} \|\nabla f(x^{t-1})\| + \|\nabla f(x^t) - \nabla f(x^{t-1})\| + \frac{1}{n} \sum_{i=1}^n (1-\beta) \|v_i^{t-1} - \nabla f_i(x^{t-1})\| \\
1035 \quad &+ \frac{1}{n} \sum_{i=1}^n (1-\beta) \|\nabla f_i(x^t) - \nabla f_i(x^{t-1})\| + B - \tau \\
1036 \quad &\stackrel{(iii)}{\leq} \sqrt{2L(f(x^{t-1}) - f^*)} + L\gamma(2-\beta)\|g^{t-1}\| + (1-\beta)\frac{1}{n} \sum_{i=1}^n \|v_i^{t-1} - \nabla f_i(x^{t-1})\| + B - \tau \\
1037 \quad &\stackrel{(iv)}{\leq} \sqrt{2L\Phi^t} + 2L\gamma\|g^{t-1}\| + (1-\beta)\frac{1}{n} \sum_{i=1}^n \|v_i^{t-1} - \nabla f_i(x^{t-1})\| + B - \tau \\
1038 \quad &\stackrel{(v)}{\leq} \sqrt{2L\Delta} + 2L\gamma \left(\sqrt{64L\Delta} + 3(B-\tau) \right) + \left(\sqrt{4L\Delta} + \frac{3}{2}(B-\tau) \right) + B - \tau \\
1039 \quad &= \left(\sqrt{2} + 16L\gamma + 2 \right) \sqrt{L\Delta} + (6L\gamma + 1 + 3/2)(B-\tau), \\
1040 \quad &
\end{aligned}$$

1041 where (i) follows from triangle inequality and update of v_i^t , and assumption 3 in the statement of the
1042 lemma; (ii) from triangle inequality; (iii) from properties of smooth function from Lemma 4 and
1043 update rule of x^t ; (iv) from the definition of Φ^t ; (v) from assumption 1, 2, and 6 in the statement of
1044 the lemma. Since $\gamma \leq \frac{1}{12L} \leq \frac{6-\sqrt{2}}{16L}$, then $16L\gamma + \sqrt{2} + 2 \leq 8$, and $\gamma \leq \frac{1}{12L}$, then $6L\gamma + 5/2 \leq 3$. \square

1045 **Lemma 7.** Let each f_i be L -smooth and $\Delta \geq \Phi^0$. Let the following inequalities hold

- 1046 1. $4L\gamma = \beta$ and $\gamma \leq \frac{1}{4L}$;
- 1047 2. $\|\nabla f_i(x^{t-1}) - v_i^{t-1}\| \leq \sqrt{4L\Delta} + \frac{3}{2}(B-\tau)$;
- 1048 3. $\|g^{t-1}\| \leq \sqrt{64L\Delta} + 3(B-\tau)$.

1049 Then we have

$$1050 \quad \|\nabla f_i(x^t) - v_i^t\| \leq \sqrt{4L\Delta} + \frac{3}{2}(B-\tau) \quad \forall i \in [n]. \quad (22)$$

1051 *Proof.* We have

$$\begin{aligned}
1052 \quad \|\nabla f_i(x^t) - v_i^t\| &\stackrel{(i)}{=} \|\nabla f_i(x^t) - (1-\beta)v_i^{t-1} - \beta\nabla f_i(x^t)\| \\
1053 \quad &= (1-\beta)\|\nabla f_i(x^t) - v_i^{t-1}\| \\
1054 \quad &\stackrel{(ii)}{\leq} (1-\beta)L\gamma\|g^{t-1}\| + (1-\beta)\|\nabla f_i(x^{t-1}) - v_i^{t-1}\| \\
1055 \quad &\stackrel{(iii)}{\leq} L\gamma \left(\sqrt{64L\Delta} + 3(B-\tau) \right) + (1-\beta) \left(\sqrt{4L\Delta} + \frac{3}{2}(B-\tau) \right) \\
1056 \quad &= (8L\gamma + 2(1-\beta))\sqrt{L\Delta} + (3L\gamma + 3(1-\beta)/2)(B-\tau), \\
1057 \quad &
\end{aligned}$$

where (i) follows from the update rule of v_i^t ; (ii) from triangle inequality, smoothness, and update of x^t ; (iii) from assumption 2-3 of the statement of the lemma. Since $4L\gamma = \beta$, then $4L\gamma + 2(1-\beta) = 2$ and $3L\gamma + 3^{(1-\beta)/2} = 3L\gamma + \frac{3}{2}(1-4L\gamma) \leq \frac{3}{2}$. \square

Lemma 8. Let each f_i be L -smooth, $\Delta \geq \Phi^0$ and $i \in \mathcal{I}_t$. Let the following inequalities hold

1. $\beta = 4L\gamma$ and $\beta \leq \frac{1}{2}$;
2. $\gamma \leq \frac{\tau}{48L\sqrt{L\Delta}}$;
3. $\gamma \leq \frac{\tau}{30L(B-\tau)}$;
4. $\|g^t\| \leq \sqrt{64L\Delta} + 3(B-\tau)$;
5. $\|\nabla f_i(x^t) - v_i^t\| \leq \sqrt{4L\Delta} + \frac{3}{2}(B-\tau)$.

Then

$$\|v_i^{t+1} - g_i^t\| \leq \|v_i^t - g_i^{t-1}\| - \frac{\tau}{2}. \quad (23)$$

Proof. Since $i \in \mathcal{I}_t$, then $\|v_i^t - g_i^{t-1}\| > \tau$, thus from Lemma 5 we have

$$\begin{aligned} \|v_i^{t+1} - g_i^t\| &\leq \|v_i^t - g_i^{t-1}\| - \tau + \beta L\gamma \|g^t\| + \beta \|\nabla f_i(x^t) - v_i^t\| \\ &\stackrel{(i)}{\leq} \|v_i^t - g_i^{t-1}\| - \tau + \frac{1}{2}L\gamma \left(\sqrt{64L\Delta} + 3(B-\tau) \right) + \beta \left(\sqrt{4L\Delta} + \frac{3}{2}(B-\tau) \right) \\ &= \|v_i^t - g_i^{t-1}\| - \tau + (4L\gamma + 2\beta)\sqrt{L\Delta} + (3L\gamma/2 + 3\beta/2)(B-\tau), \end{aligned}$$

where (i) follows from assumptions 4-5 of the statement of the lemma. Since $\beta = 4L\gamma$, we have

$$\|v_i^{t+1} - g_i^t\| \leq \|v_i^t - g_i^{t-1}\| - \tau + 12L\gamma\sqrt{L\Delta} + \frac{15}{2}L\gamma(B-\tau).$$

Since $\gamma \leq \frac{\tau}{48L\sqrt{L\Delta}}$, then $12L\gamma\sqrt{L\Delta} \leq \frac{\tau}{4}$, and since $\gamma \leq \frac{\tau}{30L(B-\tau)}$, then $\frac{15}{2}L\gamma(B-\tau) \leq \frac{\tau}{4}$. Therefore, we have

$$\|v_i^{t+1} - g_i^t\| \leq \|v_i^t - g_i^{t-1}\| - \frac{\tau}{2}. \quad \square$$

Lemma 9. Let each f_i be L -smooth. Then \tilde{P}^t decreases as

$$\tilde{P}^{t+1} \leq (1-\beta)\tilde{P}^t + \frac{3L^2}{\beta}R^t. \quad (24)$$

Proof. We have

$$\begin{aligned} \|v_i^{t+1} - \nabla f_i(x^{t+1})\|^2 &\stackrel{(i)}{\leq} \|(1-\beta)v_i^t + \beta\nabla f_i(x^{t+1}) - \nabla f_i(x^{t+1})\|^2 \\ &= (1-\beta)^2 \|\nabla f_i(x^{t+1}) - v_i^t\|^2 \\ &\stackrel{(ii)}{\leq} (1-\beta)^2(1+\beta/2) \|v_i^t - \nabla f_i(x^t)\|^2 \\ &\quad + (1-\beta)^2(1+2/\beta) \|\nabla f_i(x^t) - \nabla f_i(x^{t+1})\|^2 \\ &\stackrel{(iii)}{\leq} (1-\beta) \|v_i^t - \nabla f_i(x^t)\|^2 + \frac{3L^2}{\beta} \|x^t - x^{t+1}\|^2, \end{aligned}$$

where (i) follows from the update rule of v_i^t ; (ii) from the inequality $\|a+b\|^2 \leq (1+\beta/2)\|a\|^2 + (1+2/\beta)\|b\|^2$; (iii) from smoothness. Averaging the inequalities above across $i \in [n]$, we get the statement of the lemma. \square

Similarly, we can get the descent of P^t .

Lemma 10. Let each f_i be L -smooth. Then P^t decreases as

$$P^{t+1} \leq (1 - \beta)P^t + \frac{3L^2}{\beta}R^t. \quad (25)$$

Now we present the descent of \tilde{V}^t .

Lemma 11. Let each f_i be L -smooth. Let $\|v_i^t - g_i^{t-1}\| \leq B$ for all $i \in [n]$. Then

$$\|g_i^t - v_i^t\|^2 \leq (1 - \eta)\|g_i^{t-1} - v_i^{t-1}\|^2 + \frac{4\beta^2}{\eta}\|v_i^{t-1} - \nabla f_i(x^{t-1})\|^2 + \frac{4\beta^2 L^2}{\eta}R^{t-1}.$$

Proof. Since $\|v_i^t - g_i^{t-1}\| \leq B$, we have $\eta_i^t \geq \eta$. Thus, we have

$$\begin{aligned} \|g_i^t - v_i^t\|^2 &\stackrel{(i)}{=} \|g_i^{t-1} + \text{clip}_\tau(v_i^t - g_i^{t-1}) - v_i^t\|^2 \\ &\stackrel{(ii)}{\leq} (1 - \eta_i^t)^2 \|g_i^{t-1} - v_i^t\|^2 \\ &\stackrel{(iii)}{\leq} (1 - \eta)^2 \|g_i^{t-1} - v_i^t\|^2 \\ &\stackrel{(iv)}{=} (1 - \eta)^2 \|g_i^{t-1} - (1 - \beta)v_i^{t-1} - \beta \nabla f_i(x^t)\|^2 \\ &\stackrel{(v)}{\leq} (1 - \eta)^2 (1 + \rho) \|g_i^{t-1} - v_i^{t-1}\|^2 + (1 - \eta)^2 (1 + \rho^{-1}) \beta^2 \|v_i^{t-1} - \nabla f_i(x^t)\|^2 \\ &\stackrel{(vi)}{\leq} (1 - \eta)^2 (1 + \rho) \|g_i^{t-1} - v_i^{t-1}\|^2 + 2(1 - \eta)^2 (1 + \rho^{-1}) \beta^2 \|v_i^{t-1} - \nabla f_i(x^{t-1})\|^2 \\ &\quad + 2(1 - \eta)^2 (1 + \rho^{-1}) \beta^2 L^2 \|x^{t-1} - x^t\|^2, \end{aligned}$$

where (i) follows from the update rule of g_i^t ; (ii) from properties of clipping from Lemma 3; (iii) from the fact that $\eta_i^t \geq \eta$; (iv) from the update rule of v_i^t ; (v) from the inequality $\|a + b\|^2 \leq (1 + r/2)\|a\|^2 + (1 + 2/r)\|b\|^2$ for any positive r ; (vi) from the inequality $\|a + b\|^2 \leq (1 + r/2)\|a\|^2 + (1 + 2/r)\|b\|^2$ for any positive r and smoothness. If we choose $\rho = \eta/2$, we get

$$\|g_i^t - v_i^t\|^2 \leq (1 - \eta)\|g_i^{t-1} - v_i^{t-1}\|^2 + \frac{4\beta^2}{\eta}\|v_i^{t-1} - \nabla f_i(x^{t-1})\|^2 + \frac{4\beta^2 L^2}{\eta}R^{t-1}.$$

□

Theorem 5 (Full statement of Theorem 2). *Let Assumptions 1 holds. Let $B := \max_i \|\nabla f_i(x^0)\| > 3\tau$ and $\Delta \geq \Phi^0$. Assume the following inequalities hold*

1. **stepsize restrictions:** $\gamma \leq \frac{1}{12L}$, $\gamma \leq \frac{\tau}{48L\sqrt{L\Delta}}$, $\gamma \leq \frac{\tau}{30L(B-\tau)}$, and

$$\frac{2}{3} - \frac{16\beta^2 L^2}{\eta^2} \gamma^2 - \frac{48L^2}{\eta^2} \gamma^2 \geq 0;$$

2. **momentum restrictions:** $\beta = 4L\gamma \leq \frac{1}{2}$.

Then the Lyapunov function decreases as

$$\Phi^{t+1} \leq \Phi^t - \frac{\gamma}{2} \|\nabla f(x^t)\|^2,$$

therefore we obtain

$$\frac{1}{T} \sum_{t=0}^{T-1} \|\nabla f(x^t)\|^2 \leq \frac{2\Delta}{T} = \mathcal{O}\left(\frac{1}{T}\right). \quad (26)$$

Moreover, after at most $\frac{2B}{\tau}$ iterations, the clipping operator will be turned off for all workers.

Proof. We prove the main theorem by induction. The conventional choice is

$$\nabla f_i(x^{-1}) = v_i^{-1} = g_i^{-1} = 0, \quad \Phi^{-1} = +\infty.$$

We will show that

1. the Lyapunov function decreases as $\Phi^{t+1} \leq \Phi^t - \frac{\gamma}{2} \|\nabla f(x^t)\|^2$;
2. $\|g^t\| \leq \sqrt{64L\Delta} + 3(B - \tau)$;
3. $\|v_i^t - \nabla f_i(x^t)\| \leq \sqrt{4L\Delta} + \frac{3}{2}(B - \tau)$;
4. and $\|v_i^t - g_i^{t-1}\| \leq B - \frac{t\tau}{2}$.

First, we prove that the base of induction holds.

Base of induction.

1. $\|v_i^0 - g_i^{-1}\| = \|v_i^0\| = \beta \|\nabla f_i(x^0)\| \leq \frac{1}{2}B \leq B$ holds;
2. $g^0 = \frac{1}{n} \sum_{i=1}^n (g_i^{-1} + \text{clip}_\tau(v_i^0 - g_i^{-1})) = \frac{1}{n} \sum_{i=1}^n \text{clip}_\tau(\beta \nabla f_i(x^0))$. Therefore, we have

$$\begin{aligned} \|g^0\| &\leq \left\| \frac{1}{n} \sum_{i=1}^n \beta \nabla f_i(x^0) + (\text{clip}_\tau(\beta \nabla f_i(x^0)) - \beta \nabla f_i(x^0)) \right\| \\ &\leq \beta \|\nabla f(x^0)\| + \frac{1}{n} \sum_{i=1}^n \max\{0, \beta \|\nabla f_i(x^0)\| - \tau\} \\ &\leq \beta \sqrt{2L(f(x^0) - f^*)} + B - \tau \\ &\leq \sqrt{64L\Delta} + 3(B - \tau). \end{aligned}$$

3. We have

$$\begin{aligned} \|v_i^0 - \nabla f_i(x^0)\| &= \|\beta \nabla f_i(x^0) - \nabla f_i(x^0)\| \\ &\leq (1 - \beta)B \\ &\leq \sqrt{4L\Delta} + \frac{3}{2}(B - \tau) \end{aligned}$$

4. $\Phi^0 \leq \Phi^{-1} - \frac{\gamma}{2} \|\nabla f(x^{-1})\|^2 = \Phi^{-1}$ holds.

Transition of induction. Assume that for K we have that for all $t \in [0, K]$

1. $\Phi^t \leq \Delta$;
2. $\|g^t\| \leq \sqrt{64L\Delta} + 3(B - \tau)$;
3. $\|v_i^t - \nabla f_i(x^t)\| \leq \sqrt{4L\Delta} + \frac{3}{2}(B - \tau)$;
4. $\|v_i^t - g_i^{t-1}\| \leq B$ for $i \in \mathcal{I}_t$.

CASE $|\mathcal{I}_{K+1}| > 0$. Since all requirements of Lemma 8 are satisfied at iteration K we get for all $i \in \mathcal{I}_{K+1}$

$$\|v_i^{K+1} - g_i^K\| \leq \|v_i^K - g_i^{K-1}\| - \frac{\tau}{2} \leq B - \frac{\tau}{2}.$$

Similarly due to the assumption of induction, from Lemma 6 we get that

$$\|g^{K+1}\| \leq \sqrt{64L\Delta} + 3(B - \tau),$$

and from Lemma 7

$$\|\nabla f_i(x^{K+1}) - v_i^{K+1}\| \leq \sqrt{4L\Delta} + \frac{3}{2}(B - \tau).$$

This means that steps 1-3 in the assumption of the induction are also verified for step $K + 1$.

The remaining part is the descent of the Lyapunov function. For \tilde{V}^{K+1} we have Lemma 11 since $\|v_i^{K+1} - g_i^K\| \leq B - \frac{\tau}{2}$

$$\tilde{V}^{K+1} \leq (1 - \eta)\tilde{V}^K + \frac{4\beta^2}{\eta}\tilde{P}^K + \frac{4\beta^2L^2}{\eta}R^K.$$

Combining this result with the claims of Lemmas 2, 9 and 10 we get

$$\begin{aligned}
\Phi^{K+1} &= \delta^{K+1} + \frac{\gamma}{\eta} \tilde{V}^{K+1} + \frac{4\gamma\beta}{\eta^2} \tilde{P}^{K+1} + \frac{\gamma}{\beta} P^{K+1} \\
&\leq \delta^K - \frac{\gamma}{2} \|\nabla f(x^K)\|^2 - \frac{1}{4\gamma} R^K + \gamma \tilde{V}^K + \gamma P^K \\
&\quad + \frac{\gamma}{\eta} \left((1-\eta) \tilde{V}^K + \frac{4\beta^2}{\eta} \tilde{P}^K + \frac{4\beta^2 L^2}{\eta} R^K \right) \\
&\quad + \frac{4\gamma\beta}{\eta^2} \left((1-\beta) \tilde{P}^K + \frac{3L^2}{\beta} R^K \right) \\
&\quad + \frac{\gamma}{\beta} \left((1-\beta) P^K + \frac{3L^2}{\beta} R^K \right) \\
&= \delta^K - \frac{\gamma}{2} \|\nabla f(x^K)\|^2 + \frac{\gamma}{\eta} \tilde{V}^K (1-\eta+\eta) + \frac{4\gamma\beta}{\eta^2} \tilde{P}^{t*} (1-\beta+\beta) \\
&\quad + \frac{\gamma}{\beta} P^K (1-\beta+\beta) - \frac{1}{4\gamma} \left(1 - \frac{16\beta^2 L^2}{\eta^2} \gamma^2 - \frac{48L^2}{\eta^2} \gamma^2 - \frac{12L^2}{\beta^2} \gamma^2 \right) R^K \\
&\leq \Phi^K - \frac{\gamma}{2} \|\nabla f(x^K)\|^2 - \frac{1}{4\gamma} \left(1 - \frac{16\beta^2 L^2}{\eta^2} \gamma^2 - \frac{48L^2}{\eta^2} \gamma^2 - \frac{12L^2}{\beta^2} \gamma^2 \right) R^K.
\end{aligned}$$

Since we choose $\beta^2 = 64L^2\gamma^2$, then $-\frac{1}{\beta^2} = -\frac{1}{64L^2\gamma^2}$ and $-\frac{12L^2}{\beta^2}\gamma^2 = -\frac{12L^2}{64L^2\gamma^2}\gamma^2 \geq -\frac{1}{3}$. Therefore,

$$1 - \frac{16\beta^2 L^2}{\eta^2} \gamma^2 - \frac{48L^2}{\eta^2} \gamma^2 - \frac{12L^2}{\beta^2} \gamma^2 \geq \frac{2}{3} - \frac{16\beta^2 L^2}{\eta^2} \gamma^2 - \frac{48L^2}{\eta^2} \gamma^2 \geq 0,$$

by the choice of γ . Thus, we get

$$\Phi^{K+1} \leq \Phi^K - \frac{\gamma}{2} \|\nabla f(x^K)\|^2.$$

In particular, this implies $\Phi^{K+1} \leq \Phi^K \leq \Delta$.

CASE $|\mathcal{I}_{K+1}| = 0$. In this case $\eta_i^{K+1} = 1$ for all $i \in [n]$, i.e. $\text{clip}_\tau(v_i^{K+1} - g_i^K) = v_i^{K+1} - g_i^K$ that leads to $g_i^{K+1} = v_i^{K+1}$. Thus, $\tilde{V}^{K+1} = 0$. We can perform similar steps as before for Φ^{K+1} and get less restrictive inequality

$$\Phi^{K+1} \leq \Phi^K - \frac{\gamma}{2} \|\nabla f(x^K)\|^2 - \frac{1}{4\gamma} \left(1 - \frac{48L^2}{\eta^2} \gamma^2 - \frac{12L^2}{\beta^2} \gamma^2 \right) R^K.$$

Again, $1 - \frac{48L^2}{\eta^2} \gamma^2 - \frac{12L^2}{\beta^2} \gamma^2 \geq \frac{2}{3} - \frac{48L^2}{\eta^2} \gamma^2 \geq 0$ which is satisfied by the choice of γ .

We conclude that in both cases the Lyapunov function decreases as $\Phi^{K+1} \leq \Phi^K - \frac{\gamma}{2} \|\nabla f(x^K)\|^2$, and consequently, $\Phi^{K+1} \leq \Delta$. This finalizes the induction step. Therefore, we can guarantee that for all iterations $t \in [0, T-1]$ we have

$$\Phi^{t+1} \leq \Phi^t - \frac{\gamma}{2} \|\nabla f(x^t)\|^2 \Rightarrow \frac{1}{T} \sum_{t=0}^{T-1} \|\nabla f(x^t)\|^2 \leq \frac{2\Delta}{\gamma T}.$$

Moreover, the proof shows that the clipping operator will be eventually turned off since $\|v_i^t - g_i^{t-1}\| \leq B - \frac{t\tau}{2}$, i.e. after at most $\frac{2B}{\tau}$ iterations. \square

Remark 1. With $v_i^{-1} = g_i^{-1} = 0$ we have

$$\begin{aligned}
\Phi^0 &= \delta^0 + \frac{\gamma}{\eta} \frac{1}{n} \sum_{i=1}^n \|\text{clip}_\tau(\beta \nabla f_i(x^0)) - \beta \nabla f_i(x^0)\|^2 + \frac{4\gamma\beta}{\eta^2} \frac{1}{n} (1-\beta)^2 \sum_{i=1}^n \|\nabla f_i(x^0)\|^2 \\
&\quad + \frac{\gamma}{\beta} (1-\beta)^2 \|\nabla f(x^0)\|^2 \\
&\leq \delta^0 + \frac{\gamma}{\eta} \frac{1}{n} \sum_{i=1}^n \max\{(\|\nabla f(x^0)\| - \tau)^2, 0\} + \frac{16L\gamma^2}{\eta^2} \frac{1}{n} (1-\beta)^2 \sum_{i=1}^n \|\nabla f_i(x^0)\|^2 \\
&\quad + \frac{1}{4L} (1-\beta)^2 \|\nabla f(x^0)\|^2 \\
&\leq \frac{3}{2} \delta^0 + \frac{\gamma}{\eta} \frac{1}{n} \sum_{i=1}^n \max\{(\|\nabla f_i(x^0)\| - \tau)^2, 0\} + \frac{16L\gamma^2}{\eta^2} \frac{1}{n} (1-\beta)^2 \sum_{i=1}^n \|\nabla f_i(x^0)\|^2.
\end{aligned}$$

We have the stepsize restriction

$$\frac{2}{3} - \frac{64L^4\gamma^2}{\eta^2} - \frac{48L^2\gamma^2}{\eta^2} \geq 0. \tag{27}$$

For inequality of the form $a\gamma^2 + b\gamma \leq 1$ the stepsize restriction of the form $\gamma \leq \frac{1}{\sqrt{a+b}}$ is tight up to a constant factor 2, i.e. $\frac{2}{\sqrt{a+b}}$ does not satisfy the inequality (see Lemma 5 in (Richtárik et al., 2021)). Using this lemma in our case we get that the stepsize satisfying Equation (27) should also satisfy

$$L^2\gamma^2 \leq 2 \cdot \frac{\eta}{72/\eta + 4\sqrt{6}}.$$

This implies that $L^2\gamma^2 \leq \frac{\eta}{4\sqrt{6}}$ and $L^2\gamma^2 \leq \frac{\eta^2}{72}$. Consequently, it also satisfies $\frac{\gamma}{\eta} \leq \frac{1}{6L\sqrt{2}}$ (from the last inequality). Therefore, we have

$$\begin{aligned}
\Phi^0 &\leq \frac{3}{2} \delta^0 + \frac{1}{6L\sqrt{2}} \frac{1}{n} \sum_{i=1}^n \max\{(\|\nabla f_i(x^0)\| - \tau)^2, 0\} + \frac{2}{9L} \frac{1}{n} \sum_{i=1}^n \|\nabla f_i(x^0)\|^2 \\
&\leq \frac{3}{2} \delta^0 + \left(\frac{1}{6L\sqrt{2}} + \frac{2}{9L} \right) \frac{1}{n} \sum_{i=1}^n \|\nabla f_i(x^0)\|^2,
\end{aligned}$$

which is independent of τ , and can be use as a bound for Δ .

E PROOF OF THEOREM 4

We define constants a , b , and c as follows that will be used later in the proofs:

$$\begin{aligned}
 a &:= \left(\sqrt{2} + 2\sqrt{3 \log \frac{6(T+1)}{\alpha}} \right) \sqrt{d} \sigma_\omega \sqrt{T/n}, \\
 b^2 &:= 2\sigma^2 \log \left(\frac{6(T+1)n}{\alpha} \right), \\
 c^2 &:= \left(\sqrt{2} + 2\sqrt{3 \log \frac{6(T+1)}{\alpha}} \right)^2 \sigma^2,
 \end{aligned} \tag{28}$$

where T is the number of iterations, n is the number of workers, d is the dimension of the problem, σ is from Assumption 2, $\alpha \in (0, 1)$ is a constant, and σ_ω is the variance of DP noise.

Lemma 12. Let each f_i be L -smooth. Then we have the following inequality with probability 1

$$\|v_i^{t+1} - g_i^t\| \leq \max\{0, \|v_i^t - g_i^{t-1}\| - \tau\} + \beta L \gamma \|g^t\| + \beta \|\nabla f_i(x^t) - v_i^t\| + \beta \|\theta_i^{t+1}\|. \tag{29}$$

Proof. We have

$$\begin{aligned}
 \|v_i^{t+1} - g_i^t\| &\stackrel{(i)}{=} \|(1 - \beta)v_i^t + \beta \nabla f_i(x^{t+1}, \xi_i^{t+1}) - g_i^t\| \\
 &\stackrel{(ii)}{\leq} \|v_i^t - g_i^t\| + \beta \|\nabla f_i(x^{t+1}, \xi_i^{t+1}) - v_i^t\| \\
 &\stackrel{(iii)}{=} \|v_i^t - \text{clip}_\tau(v_i^t - g_i^{t-1}) - g_i^{t-1}\| + \beta \|\nabla f_i(x^{t+1}, \xi_i^{t+1}) - v_i^t\| \\
 &\stackrel{(iv)}{\leq} \max\{0, \|v_i^t - g_i^{t-1}\| - \tau\} + \beta \|\nabla f_i(x^{t+1}, \xi_i^{t+1}) - \nabla f_i(x^{t+1})\| \\
 &\quad + \beta \|\nabla f_i(x^{t+1}) - \nabla f_i(x^t)\| + \beta \|\nabla f_i(x^t) - v_i^t\| \\
 &\stackrel{(v)}{\leq} \max\{0, \|v_i^t - g_i^{t-1}\| - \tau\} + \beta L \|x^{t+1} - x^t\| + \beta \|\nabla f_i(x^t) - v_i^t\| + \beta \|\theta_i^{t+1}\| \\
 &\stackrel{(vi)}{=} \max\{0, \|v_i^t - g_i^{t-1}\| - \tau\} + \beta L \gamma \|g^t\| + \beta \|\nabla f_i(x^t) - v_i^t\| + \beta \|\theta_i^{t+1}\|,
 \end{aligned}$$

where (i) follows from the update rule of v_i^t ; (ii) from triangle inequality; (iii) from the update rule of g_i^t ; (iv) from the properties of the clipping operator from Lemma 3 and triangle inequality; (v) from smoothness; (vi) from the update rule of x^t . \square

Let us choose $p \in [0.2, 0.8]$. With this choice we have $3x^{1-p} \geq 4x$ for any $x \in (0, 1/12]$.

Lemma 13. Let each f_i be L -smooth and $\Delta \geq \Phi^0$. Assume that the following inequalities hold

1. $g^0 = \frac{1}{n} \sum_{i=1}^n g_i^0$;
2. $\|g^{t-1}\| \leq \sqrt{64L\Delta} + 3(B - \tau) + 3b + 3a$;
3. $\|\nabla f_i(x^{t-1}) - v_i^{t-1}\| \leq \sqrt{4L\Delta} + \frac{3}{2}(B - \tau) + 2b + (L\gamma)^p a$ for all $i \in [n]$;
4. $\|v_i^t - g_i^{t-1}\| \leq B$ for all $i \in [n]$;
5. $\gamma \leq \frac{1}{12L}$;
6. $\|\theta_i^t\| \leq b$ for all $i \in [n]$;
7. $\left\| \frac{1}{n} \sum_{l=1}^t \sum_{i=1}^n \omega_i^l \right\| \leq a$;
8. $1 \geq \beta \geq 4L\gamma$;
9. $\Phi^{t-1} \leq 2\Delta$.

1404 Then we have
1405
1406

$$1407 \quad \|g^t\| \leq \sqrt{64L\Delta} + 3(B - \tau) + 3b + 3a. \quad (30)$$

1408
1409
1410
1411
1412
1413

1414 *Proof.* We start as follows
1415
1416

$$1417 \quad \|g^t\| \stackrel{(i)}{=} \left\| g^{t-1} + \frac{1}{n} \sum_{i=1}^n \text{clip}_\tau(v_i^t - g_i^{t-1}) + \frac{1}{n} \sum_{i=1}^n \omega_i^t \right\|$$

$$1418 \quad = \left\| \frac{1}{n} \sum_{i=1}^n [\nabla f_i(x^t) + (v_i^t - \nabla f_i(x^t)) + \text{clip}_\tau(v_i^t - g_i^{t-1}) - (v_i^t - g_i^{t-1})] \right.$$

$$1419 \quad \left. + g^{t-1} - \frac{1}{n} \sum_{i=1}^n g_i^{t-1} + \frac{1}{n} \sum_{i=1}^n \omega_i^t \right\|$$

$$1420 \quad \stackrel{(ii)}{\leq} \|\nabla f(x^t)\| + \frac{1}{n} \sum_{i=1}^n \|v_i^t - \nabla f_i(x^t)\| + \frac{1}{n} \sum_{i=1}^n \max\{0, \|v_i^t - g_i^{t-1}\| - \tau\}$$

$$1421 \quad + \left\| g^{t-2} + \frac{1}{n} \sum_{i=1}^n [\text{clip}_\tau(v_i^{t-1} - g_i^{t-2}) + \omega_i^{t-1}] - \frac{1}{n} \sum_{i=1}^n [g_i^{t-2} + \text{clip}_\tau(v_i^{t-1} - g_i^{t-2})] \right.$$

$$1422 \quad \left. + \frac{1}{n} \sum_{i=1}^n \omega_i^t \right\|$$

1423
1424
1425
1426
1427
1428
1429
1430
1431
1432
1433
1434
1435
1436

1437 where (i) follows from the update rule of g^t ; (ii) from the triangle inequality and the properties of
1438 the clipping operator from Lemma 3. Cancelling terms inside the norm in the last term above we
1439 obtain
1440

$$1441 \quad \|g^t\| \stackrel{(i)}{=} \|\nabla f(x^t)\| + \frac{1}{n} \sum_{i=1}^n \|v_i^t - \nabla f_i(x^t)\| + \frac{1}{n} \sum_{i=1}^n \max\{0, \|v_i^t - g_i^{t-1}\| - \tau\}$$

$$1442 \quad + \left\| g^{t-2} - \frac{1}{n} \sum_{i=1}^n g_i^{t-2} + \frac{1}{n} \sum_{l=t-1}^t \sum_{i=1}^n \omega_i^l \right\|$$

$$1443 \quad \stackrel{(ii)}{=} \|\nabla f(x^t)\| + \frac{1}{n} \sum_{i=1}^n \|v_i^t - \nabla f_i(x^t)\| + \frac{1}{n} \sum_{i=1}^n \max\{0, \|v_i^t - g_i^{t-1}\| - \tau\}$$

$$1444 \quad + \left\| \frac{1}{n} \sum_{l=1}^t \sum_{i=1}^n \omega_i^l \right\|,$$

1445
1446
1447
1448
1449
1450
1451
1452
1453
1454

1455 where (ii) follows from performing similar steps as in (i) and having in mind assumption 1 from
1456 the statement of the lemma.
1457

We continue to bound $\|g^t\|$ in the following way

1458
 1459
 1460
 1461
 1462
 1463
 1464
 1465
 1466
 1467
 1468
 1469
 1470
 1471
 1472
 1473
 1474
 1475
 1476
 1477
 1478
 1479
 1480
 1481
 1482
 1483
 1484
 1485
 1486
 1487
 1488
 1489
 1490
 1491
 1492
 1493
 1494
 1495
 1496
 1497
 1498
 1499
 1500
 1501
 1502
 1503
 1504
 1505
 1506
 1507
 1508
 1509
 1510
 1511

$$\begin{aligned}
 \|g^t\| &\stackrel{(i)}{\leq} \|\nabla f(x^{t-1})\| + \|\nabla f(x^t) - \nabla f(x^{t-1})\| + \frac{1}{n} \sum_{i=1}^n \|(1-\beta)v_i^{t-1} + \beta\nabla f_i(x^t, \xi_i^t) - \nabla f_i(x^t)\| \\
 &\quad + B - \tau + \left\| \frac{1}{n} \sum_{l=1}^t \sum_{i=1}^n \omega_i^l \right\| \\
 &\stackrel{(ii)}{\leq} \|\nabla f(x^{t-1})\| + \|\nabla f(x^t) - \nabla f(x^{t-1})\| + \frac{1}{n} \sum_{i=1}^n \|(1-\beta)v_i^{t-1} + \beta\nabla f_i(x^t) - \nabla f_i(x^t)\| \\
 &\quad + \frac{\beta}{n} \sum_{i=1}^n \|\theta_i^t\| + B - \tau + \left\| \frac{1}{n} \sum_{l=1}^t \sum_{i=1}^n \omega_i^l \right\| \\
 &\stackrel{(iii)}{\leq} \|\nabla f(x^{t-1})\| + \|\nabla f(x^t) - \nabla f(x^{t-1})\| + \frac{1}{n} \sum_{i=1}^n (1-\beta) \|v_i^{t-1} - \nabla f_i(x^{t-1})\| \\
 &\quad + \frac{1}{n} \sum_{i=1}^n (1-\beta) \|\nabla f_i(x^t) - \nabla f_i(x^{t-1})\| + \frac{\beta}{n} \sum_{i=1}^n \|\theta_i^t\| + B - \tau + \left\| \frac{1}{n} \sum_{l=1}^t \sum_{i=1}^n \omega_i^l \right\|,
 \end{aligned}$$

where (i) follows from triangle inequality, the update rule of v_i^t , and properties of the clipping operator from Lemma 3; (ii) and (iii) from triangle inequality. Using smoothness of f we continue

$$\begin{aligned}
 \|g^t\| &\stackrel{(i)}{\leq} \sqrt{2L(f(x^{t-1}) - f^*)} + L\gamma(2-\beta)\|g^{t-1}\| + (1-\beta)\frac{1}{n} \sum_{i=1}^n \|v_i^{t-1} - \nabla f_i(x^{t-1})\| \\
 &\quad + \frac{\beta}{n} \sum_{i=1}^n \|\theta_i^t\| + B - \tau + \frac{1}{n} \sum_{i=1}^n \|\omega_i^t\| \\
 &\leq \sqrt{2L\Phi^{t-1}} + 2L\gamma\|g^{t-1}\| + (1-\beta)\frac{1}{n} \sum_{i=1}^n \|v_i^{t-1} - \nabla f_i(x^{t-1})\| \\
 &\quad + \frac{\beta}{n} \sum_{i=1}^n \|\theta_i^t\| + B - \tau + \left\| \frac{1}{n} \sum_{l=1}^t \sum_{i=1}^n \omega_i^l \right\| \\
 &\stackrel{(ii)}{\leq} \sqrt{4L\Delta} + 2L\gamma \left(\sqrt{64L\Delta} + 3(B-\tau) + 3b + 3a \right) \\
 &\quad + (1-\beta) \left(\sqrt{4L\Delta} + \frac{3}{2}(B-\tau) + 2b + (L\gamma)^p a \right) + B - \tau + \beta b + a \\
 &\leq (16L\gamma + 4) \sqrt{L\Delta} + (6L\gamma + 1 + 3/2)(B-\tau) + b(6L\gamma + 2(1-\beta) + \beta) \\
 &\quad + a(6L\gamma + (L\gamma)^p(1-\beta) + 1),
 \end{aligned}$$

where (i) follows from Lemma 4 and smoothness; from (ii) from assumptions 2-4 in the statement of the lemma.

The claim of the lemma comes by noticing that since $\gamma \leq \frac{1}{12L} < \frac{1}{4L}$, then $16L\gamma + 4 \leq 8$. Moreover, $6L\gamma + 1 + 3/2 \leq 1/2 + 1 + 3/2 = 3$. Next, we have that

$$6L\gamma + 2(1-\beta) + \beta \leq 3 \Leftrightarrow 6L\gamma \leq 1 + \beta,$$

which is satisfied if $12L\gamma \leq 1$, and

$$6L\gamma + (L\gamma)^p(1-\beta) + 1 \leq \frac{1}{2} + (1/12)^p + 1 \leq \frac{1}{2} + 1 + 1 < 3,$$

where the last inequality holds for any $p \in [0.2, 0.8]$ since $\beta \leq 1$ and $L\gamma \leq 1/12$. \square

Lemma 14. Let each f_i is L -smooth and $\Delta \geq \Phi^0$. Assume the following inequalities hold

1. $\gamma \leq \frac{1}{12L}$;

- 1512 2. $3(L\gamma)^{1-p} = \max\{4L\gamma, 3(L\gamma)^{1-p}\} \leq \beta \leq 1^5$;
 1513
 1514 3. $\|\nabla f_i(x^{t-1}) - v_i^{t-1}\| \leq \sqrt{4L\Delta} + \frac{3}{2}(B - \tau) + 2b + (L\gamma)^p a$;
 1515
 1516 4. $\|\theta_i^t\| \leq b$;
 1517
 1518 5. $\|g^{t-1}\| \leq \sqrt{64L\Delta} + 3(B - \tau) + 3b + 3a$.

1519 Then we have

$$1520 \|\nabla f_i(x^t) - v_i^t\| \leq \sqrt{4L\Delta} + \frac{3}{2}(B - \tau) + 2b + (L\gamma)^p a. \quad (31)$$

1522 *Proof.* We have

$$1524 \begin{aligned} 1525 \|\nabla f_i(x^t) - v_i^t\| &\stackrel{(i)}{=} \|\nabla f_i(x^t) - (1 - \beta)v_i^{t-1} - \beta\nabla f_i(x^t, \xi_i^t)\| \\ 1526 &\stackrel{(ii)}{\leq} (1 - \beta)\|\nabla f_i(x^t) - v_i^{t-1}\| + \beta\|\nabla f_i(x^t) - \nabla f_i(x^t, \xi_i^t)\| \\ 1527 &\stackrel{(iii)}{\leq} (1 - \beta)L\gamma\|g^{t-1}\| + (1 - \beta)\|\nabla f_i(x^{t-1}) - v_i^{t-1}\| + \beta\|\theta_i^t\| \\ 1528 &\stackrel{(iv)}{\leq} (1 - \beta)L\gamma\left(\sqrt{64L\Delta} + 3(B - \tau) + 3b + 3a\right) \\ 1529 &\quad + (1 - \beta)\left(\sqrt{4L\Delta} + \frac{3}{2}(B - \tau) + 2b + (L\gamma)^p a\right) + \beta b \\ 1530 &= (8L\gamma + 2(1 - \beta))\sqrt{L\Delta} + (3L\gamma + 3^{(1-\beta)/2})(B - \tau) \\ 1531 &\quad + (3L\gamma + 2(1 - \beta) + \beta)b + (3L\gamma + (L\gamma)^p(1 - \beta))a, \end{aligned}$$

1532 where (i) follows from the update rule of v_i^t ; (ii) from the triangle inequality; (iii) from triangle
 1533 inequality, smoothness, and the update rule of x^t ; (iv) from assumptions 2-4 of the lemma. Since
 1534 $\beta = 6L\gamma$, then

$$1540 \begin{aligned} 1541 8L\gamma + 2(1 - \beta) &\leq 2 \Leftrightarrow 4L\gamma \leq \beta, \\ 1542 3L\gamma + \frac{3}{2}(1 - \beta) &\leq \frac{3}{2} \Leftrightarrow 2L\gamma \leq \beta, \\ 1543 3L\gamma + 2(1 - \beta) + \beta &\leq 2 \Leftrightarrow 3L\gamma \leq \beta, \\ 1544 3L\gamma + (L\gamma)^p(1 - \beta) &\leq (L\gamma)^p \Leftrightarrow 3L\gamma \leq (L\gamma)^p \beta \Leftrightarrow 3(L\gamma)^{1-p} \leq \beta, \end{aligned}$$

1545 where the last inequalities in each line hold by the choice of β . \square

1546 **Lemma 15.** Let each f_i be L -smooth, $\Delta \geq \Phi^0$, and $i \in \mathcal{I}_t$. Let the following inequalities hold

- 1550 1. $12L\gamma \leq 1$;
 1551 2. $1 \geq \beta \geq \max\{4L\gamma, 3(L\gamma)^{1-p}\} = 3(L\gamma)^{1-p}$;
 1552 3. $\beta \leq \frac{\tau}{32\sqrt{L\Delta}}$;
 1553 4. $\beta \leq \frac{\tau}{18(B - \tau)}$;
 1554 5. $\beta \leq \frac{\tau}{30b}$;
 1555 6. $\beta \leq \left(\frac{\tau}{48a}\right)^{1-p}$;
 1556 7. $\|g^t\| \leq \sqrt{64L\Delta} + 3(B - \tau) + 3b + 3a$;
 1557 8. $\|\theta_i^{t+1}\| \leq b$;
 1558 9. $\|\nabla f_i(x^t) - v_i^t\| \leq \sqrt{4L\Delta} + \frac{3}{2}(B - \tau) + 2b + (L\gamma)^p a$.

1565 ⁵For $p \in [1/5, 0.8]$ we have $3x^{1-p} \geq 4x$ for any $x \in [0, 1/12]$.

1566 Then

$$1567 \quad \|v_i^{t+1} - g_i^t\| \leq \|v_i^t - g_i^{t-1}\| - \frac{\tau}{2}. \quad (32)$$

1570 *Proof.* Since $i \in \mathcal{I}_t$, then $\|v_i^t - g_i^{t-1}\| > \tau$, thus from Lemma 12 we have

$$1571 \quad \|v_i^{t+1} - g_i^t\| \leq \|v_i^t - g_i^{t-1}\| - \tau + \beta L\gamma \|g^t\| + \beta \|\nabla f_i(x^t) - v_i^t\| + \beta \|\theta_i^{t+1}\|$$

$$1572 \quad \stackrel{(i)}{\leq} \|v_i^t - g_i^{t-1}\| - \tau + L\gamma \left(\sqrt{64L\Delta} + 3(B - \tau) + 3b + 3a \right)$$

$$1573 \quad + \beta \left(\sqrt{4L\Delta} + \frac{3}{2}(B - \tau) + 2b + (L\gamma)^p a \right) + \beta b$$

$$1574 \quad = \|v_i^t - g_i^{t-1}\| - \tau + (8L\gamma + 2\beta)\sqrt{L\Delta} + (3L\gamma + 3\beta/2)(B - \tau) + (3L\gamma + 2\beta)b$$

$$1575 \quad + (3L\gamma + (L\gamma)^p)\beta a,$$

1580 where (i) follows from assumptions 7-9 of the lemma. Since $4L\gamma \leq \beta$ we have

$$1581 \quad (8L\gamma + 2\beta)\sqrt{L\Delta} \leq 4\beta\sqrt{L\Delta} \leq \frac{\tau}{8},$$

1582 where $\beta \leq \frac{\tau}{32\sqrt{L\Delta}}$. Since $4L\gamma \leq \beta$ we have

$$1583 \quad \left(3L\gamma + \frac{3\beta}{2} \right) (B - \tau) \leq \frac{9}{4}\beta(B - \tau) \leq \frac{\tau}{8},$$

1584 where $\beta \leq \frac{\tau}{18(B - \tau)}$. Since $4L\gamma \leq \beta$ we have

$$1585 \quad (3L\gamma + 3\beta)b \leq \frac{15}{4}\beta b \leq \frac{\tau}{8},$$

1586 where $\beta \leq \frac{\tau}{30b}$. Since $3(L\gamma)^{1-p} \leq \beta$ we have

$$1587 \quad (3L\gamma + (L\gamma)^p)\beta a \leq 6(\beta/3)^{\frac{1}{1-p}} a \leq \frac{\tau}{8},$$

1588 where $\beta \leq \left(\frac{\tau}{48a}\right)^{1-p}$. Thus we have

$$1589 \quad \|v_i^{t+1} - g_i^t\| \leq \|v_i^t - g_i^{t-1}\| - \tau + 4 \cdot \frac{\tau}{8} = \|v_i^t - g_i^{t-1}\| - \frac{\tau}{2}.$$

1590 □

1600 **Lemma 16.** Let $\|\theta_i^{t+1}\| \leq b$ for all $i \in [n]$. Let each f_i be L -smooth. Then \tilde{P}^t decreases as

$$1601 \quad \tilde{P}^{t+1} \leq (1 - \beta)\tilde{P}^t + \frac{3L^2}{\beta}R^t + \beta^2b^2 + \frac{2}{n}\beta(1 - \beta) \sum_{i=1}^n \langle v_i^t - \nabla f_i(x^{t+1}), \theta_i^{t+1} \rangle. \quad (33)$$

1602 *Proof.* We have

$$1603 \quad \|v_i^{t+1} - \nabla f_i(x^{t+1})\|^2 \stackrel{(i)}{=} \|(1 - \beta)v_i^t + \beta\nabla f_i(x^{t+1}, \xi_i^{t+1}) - \nabla f_i(x^{t+1})\|^2$$

$$1604 \quad = \|(1 - \beta)(v_i^t - \nabla f_i(x^{t+1})) + \beta(\nabla f_i(x^{t+1}, \xi_i^{t+1}) - \nabla f_i(x^{t+1}))\|^2$$

$$1605 \quad = (1 - \beta)^2 \|v_i^t - \nabla f_i(x^{t+1})\|^2 + \beta^2 \|\theta_i^{t+1}\|^2$$

$$1606 \quad + 2\beta(1 - \beta) \langle v_i^t - \nabla f_i(x^{t+1}), \theta_i^{t+1} \rangle$$

$$1607 \quad \stackrel{(ii)}{\leq} (1 - \beta)^2 (1 + \beta/2) \|v_i^t - \nabla f_i(x^t)\|^2$$

$$1608 \quad + (1 - \beta)^2 (1 + 2/\beta) \|\nabla f_i(x^t) - \nabla f_i(x^{t+1})\|^2 + \beta^2 b^2$$

$$1609 \quad + 2\beta(1 - \beta) \langle v_i^t - \nabla f_i(x^{t+1}), \theta_i^{t+1} \rangle$$

$$1610 \quad \stackrel{(iii)}{\leq} (1 - \beta) \|v_i^t - \nabla f_i(x^t)\|^2 + \frac{3L^2}{\beta} \|x^t - x^{t+1}\|^2 + \beta^2 b^2$$

$$1611 \quad + 2\beta(1 - \beta) \langle v_i^t - \nabla f_i(x^{t+1}), \theta_i^{t+1} \rangle,$$

where (i) follows from the update rule of v_i^t ; (ii) from $\|x+y\|^2 \leq (1+r)\|x\|^2 + (1+r^{-1})\|y\|^2$ for any $x, y \in \mathbb{R}^d$ and $r > 0$; (iii) from the smoothness and inequalities $(1-\beta)^2(1+\beta/2) \leq (1-\beta)$ and $(1-\beta)^2(1+2/\beta) \leq 3/\beta$. Averaging the inequalities above across all $i \in [n]$ we get the lemma statement. \square

Similarly, we can get the descent of P^t .

Lemma 17. Let $\|\theta^{t+1}\| \leq \frac{c}{\sqrt{n}}$, and each f_i be L -smooth. Then P^t decreases as

$$P^{t+1} \leq (1-\beta)P^t + \frac{3L^2}{\beta}R^t + \beta^2b^2 + 2\beta(1-\beta)\langle v^t - \nabla f(x^{t+1}), \theta^{t+1} \rangle.$$

Proof. For shortness, we denote $\nabla f(x^t, \xi^t) := \frac{1}{n} \sum_{i=1}^n \nabla f(x^t, \xi_i^t)$ and $\theta^t := \frac{1}{n} \sum_{i=1}^n \nabla f_i(x^t, \xi^t) - \nabla f(x^t)$. Then we have

$$\begin{aligned} \|v^{t+1} - \nabla f(x^t)\|^2 &\stackrel{(i)}{=} \|(1-\beta)v^t + \beta\nabla f(x^{t+1}, \xi^{t+1}) - \nabla f(x^{t+1})\|^2 \\ &= \|(1-\beta)(v_i^t - \nabla f_i(x^{t+1})) + \beta(\nabla f(x^{t+1}, \xi^{t+1}) - \nabla f(x^{t+1}))\|^2 \\ &= (1-\beta)^2\|v^t - \nabla f(x^{t+1})\|^2 + \beta^2\|\theta^{t+1}\|^2 \\ &\quad + 2\beta(1-\beta)\langle v^t - \nabla f(x^{t+1}), \theta^{t+1} \rangle \\ &\stackrel{(ii)}{\leq} (1-\beta)^2(1+\beta/2)\|v^t - \nabla f(x^t)\|^2 \\ &\quad + (1-\beta)^2(1+2/\beta)\|\nabla f_i(x^t) - \nabla f_i(x^{t+1})\|^2 + \beta^2\frac{c^2}{n} \\ &\quad + 2\beta(1-\beta)\langle v^t - \nabla f(x^{t+1}), \theta_i^{t+1} \rangle \\ &\stackrel{(iii)}{\leq} (1-\beta)\|v^t - \nabla f(x^t)\|^2 + \frac{3L^2}{\beta}\|x^t - x^{t+1}\|^2 + \beta^2\frac{c^2}{n} \\ &\quad + 2\beta(1-\beta)\langle v^t - \nabla f(x^{t+1}), \theta^{t+1} \rangle, \end{aligned}$$

where (i) follows from the update rule of v_i^t ; (ii) from $\|x+y\|^2 \leq (1+r)\|x\|^2 + (1+r^{-1})\|y\|^2$ for any $x, y \in \mathbb{R}^d$ and $r > 0$; (iii) from the smoothness and inequalities $(1-\beta)^2(1+\beta/2) \leq (1-\beta)$ and $(1-\beta)^2(1+2/\beta) \leq 3/\beta$. \square

Now we present the descent of \tilde{V}^t .

Lemma 18. Let $\|\theta_i^t\| \leq b$ for all $i \in [n]$, each f_i be L -smooth, and $\|v_i^t - g_i^{t-1}\| \leq B$ for all $i \in [n]$. Then

$$\begin{aligned} \|g_i^t - v_i^t\|^2 &\leq (1-\eta)\|g_i^{t-1} - v_i^{t-1}\|^2 + \frac{4\beta^2}{\eta}\|v_i^{t-1} - \nabla f_i(x^{t-1})\|^2 + \frac{4\beta^2L^2}{\eta}R^{t-1} + \beta^2b^2 \quad (34) \\ &\quad + 2(1-\eta)^2\beta\langle (g_i^{t-1} - v_i^{t-1}) + \beta(v_i^{t-1} - \nabla f_i(x^{t-1})), \theta_i^t \rangle \\ &\quad + 2(1-\eta)^2\beta\langle \beta(\nabla f_i(x^{t-1}) - \nabla f_i(x^t)), \theta_i^t \rangle. \end{aligned}$$

Moreover, averaging the inequalities across all $i \in [n]$ we get

$$\begin{aligned} \tilde{V}^t &\leq (1-\eta)\tilde{V}^{t-1} + \frac{4\beta^2}{\eta}\tilde{P}^{t-1} + \frac{4\beta^2L^2}{\eta}R^{t-1} + \beta^2b^2 \quad (35) \\ &\quad + \frac{2}{n}(1-\eta)^2\beta\sum_{i=1}^n\langle (g_i^{t-1} - v_i^{t-1}) + \beta(v_i^{t-1} - \nabla f_i(x^{t-1})) + \beta(\nabla f_i(x^{t-1}) - \nabla f_i(x^t)), \theta_i^t \rangle. \end{aligned}$$

Proof. Since $\|v_i^t - g_i^{t-1}\| \leq B$, we have $\eta_i^t \geq \eta \in (0, 1)$. Thus, we have

$$\begin{aligned} \|g_i^t - v_i^t\|^2 &\stackrel{(i)}{=} \|g_i^{t-1} + \text{clip}_\tau(v_i^t - g_i^{t-1}) - v_i^t\|^2 = \|g_i^{t-1} - v_i^t - (v_i^t - g_i^{t-1}) \cdot \tau / \|v_i^t - g_i^{t-1}\|\|^2 \\ &\leq (1-\eta_i^t)^2\|g_i^{t-1} - v_i^t\|^2 \leq (1-\eta)^2\|g_i^{t-1} - v_i^t\|^2, \end{aligned}$$

where (i) follows from the update rule of g_i^t . We can rewrite RHS in the above inequality as follows using the update rule of v_i^t

$$\begin{aligned} \|g_i^t - v_i^t\|^2 &\leq (1-\eta)^2 \|g_i^{t-1} - (1-\beta)v_i^{t-1} - \beta\nabla f_i(x^t, \xi_i^t)\|^2 \\ &= (1-\eta)^2 \|g_i^{t-1} - (1-\beta)v_i^{t-1} - \beta\nabla f_i(x^t) - \beta\theta_i^t\|^2 \\ &= (1-\eta)^2 \|g_i^{t-1} - (1-\beta)v_i^{t-1} - \beta\nabla f_i(x^t)\|^2 + (1-\eta)^2 \beta^2 \|\theta_i^t\|^2 \\ &\quad + 2(1-\eta)^2 \beta \langle g_i^{t-1} - (1-\beta)v_i^{t-1} - \beta\nabla f_i(x^t), \theta_i^t \rangle \end{aligned}$$

Therefore, we have

$$\begin{aligned} \|g_i^t - v_i^t\|^2 &\stackrel{(i)}{\leq} (1-\eta)^2 \|g_i^{t-1} - (1-\beta)v_i^{t-1} - \beta\nabla f_i(x^t)\|^2 + \beta^2 b^2 \\ &\quad + 2(1-\eta)^2 \beta \langle g_i^{t-1} - (1-\beta)v_i^{t-1} - \beta\nabla f_i(x^t), \theta_i^t \rangle \\ &\stackrel{(ii)}{\leq} (1-\eta)^2 (1+\rho) \|g_i^{t-1} - v_i^{t-1}\|^2 + (1-\eta)^2 (1+\rho^{-1}) \beta^2 \|v_i^{t-1} - \nabla f_i(x^t)\|^2 \\ &\quad + \beta^2 b^2 + 2(1-\eta)^2 \beta \langle g_i^{t-1} - (1-\beta)v_i^{t-1} - \beta\nabla f_i(x^t), \theta_i^t \rangle \\ &\stackrel{(iii)}{\leq} (1-\eta)^2 (1+\rho) \|g_i^{t-1} - v_i^{t-1}\|^2 + 2(1-\eta)^2 (1+\rho^{-1}) \beta^2 \|v_i^{t-1} - \nabla f_i(x^{t-1})\|^2 \\ &\quad + 2(1-\eta)^2 (1+\rho^{-1}) \beta^2 L^2 \|x^{t-1} - x^t\|^2 + \beta^2 b^2 \\ &\quad + 2(1-\eta)^2 \beta \langle (g_i^{t-1} - v_i^{t-1}) + \beta(v_i^{t-1} - \nabla f_i(x^{t-1})), \theta_i^t \rangle \\ &\quad + 2(1-\eta)^2 \beta \langle \beta(\nabla f_i(x^{t-1}) - \nabla f_i(x^t)), \theta_i^t \rangle, \end{aligned}$$

where (i) follows from the assumption of the lemma; (ii) from the inequality $\|x+y\|^2 \leq (1+r)\|x\|^2 + (1+r^{-1})\|y\|^2$ for any $x, y \in \mathbb{R}^d$ and $r > 0$; from $\|x+y\|^2 \leq 2\|x\|^2 + 2\|y\|^2$ for any $x, y \in \mathbb{R}^d$ and smoothness.

If we choose $\rho = \eta/2$, we get the final bound

$$\begin{aligned} \|g_i^t - v_i^t\|^2 &\leq (1-\eta) \|g_i^{t-1} - v_i^{t-1}\|^2 + \frac{4\beta^2}{\eta} \|v_i^{t-1} - \nabla f_i(x^{t-1})\|^2 \\ &\quad + \frac{4\beta^2 L^2}{\eta} R^{t-1} + \beta^2 b^2 + 2(1-\eta)^2 \beta \langle (g_i^{t-1} - v_i^{t-1}) + \beta(v_i^{t-1} - \nabla f_i(x^{t-1})), \theta_i^t \rangle \\ &\quad + 2(1-\eta)^2 \beta \langle \beta(\nabla f_i(x^{t-1}) - \nabla f_i(x^t)), \theta_i^t \rangle \end{aligned}$$

□

Theorem 6 (Full statement of Theorem 4). *Let Assumptions 1 and 2 hold, $B := \max_i \{\|\nabla f_i(x^0)\|\} + b > \tau$, probability constant $\alpha \in (0, 1)$, constants a, b , and c be defined as in (28), $p = 0.8$, and $\Delta \geq \Phi^0$. Let us run Algorithm 3 for T iterations with DP noise with variance σ_ω . Assume the following inequalities hold*

1. stepsize restrictions:

- i) $12L\gamma \leq 1$;
- ii)

$$\frac{2}{3} - \frac{16\beta^2 L^2}{\eta^2} \gamma^2 - \frac{48L^2}{\eta^2} \gamma^2 \geq 0;$$

2. momentum restrictions:

- i) $1 \geq \beta \geq \max\{4L\gamma, 3L\gamma\}^{1-p} = 3(L\gamma)^{1-p}$;
- ii) $\beta \leq \frac{\tau}{32\sqrt{L\Delta}}$;
- iii) $\beta \leq \frac{\tau}{18(B-\tau)}$;
- iv) $\beta \leq \frac{\tau}{30b}$;
- v) $\beta \leq \left(\frac{3\tau}{40a}\right)^{1-p}$;

1728 v_i) and momentum restrictions defined in (38), (39), (40), (41), (42), (44), (43), and (45);

1729
1730 Then with probability $1 - \alpha$ we have

$$1731 \frac{1}{T} \sum_{t=0}^{T-1} \|\nabla f(x^t)\|^2 \leq \tilde{\mathcal{O}} \left(\frac{L^2 \Delta \sqrt{d} \sigma_\omega}{\sqrt{T} n \tau} + \frac{(L\Delta)^{\frac{1}{6}} \sigma^{\frac{5}{3}}}{T^{\frac{1}{6}} n^{\frac{5}{6}}} + \frac{(L\Delta)^{\frac{4}{9}} \sigma^{\frac{5}{9}} d^{\frac{5}{18}} \sigma_\omega^{\frac{5}{9}}}{T^{\frac{4}{9}} n^{\frac{5}{9}}} \right),$$

1732 where $\tilde{\mathcal{O}}$ hides constant and logarithmic factors, and higher order terms decreasing in T .

1733 *Proof.* We prove the main theorem by induction. The conventional choice $\nabla f_i(x^{-1}, \xi_i^{-1}) = v_i^{-1} =$
1734 $g_i^{-1} = 0, \Phi^{-1} = \Phi^0$.

1735 Let us define an event E_t for each $t \in \{0, \dots, T\}$ such that the following inequalities hold for all
1736 $k \in \{0, \dots, t\}$

- 1737 1. $\|v_i^k - g_i^{k-1}\| \leq B$ for $i \in \mathcal{I}_k$;
- 1738 2. $\|g^k\| \leq \sqrt{64L\Delta} + 3(B - \tau) + 3b + 3a$;
- 1739 3. $\|v_i^k - \nabla f_i(x^k)\| \leq \sqrt{4L\Delta} + \frac{3}{2}(B - \tau) + 2b + (L\gamma)^p a$;
- 1740 4. $\|\theta_i^k\| \leq b$ for all $i \in [n]$ and $\|\theta^k\| \leq \frac{c}{\sqrt{n}}$;
- 1741 5. $\left\| \frac{1}{n} \sum_{l=1}^{k+1} \sum_{i=1}^n \omega_i^l \right\| \leq a$;
- 1742 6. $\Phi^k \leq 2\Delta$;
- 1743 7.

$$1744 \Delta \geq \frac{2\gamma\beta}{n\eta} (1 - \eta)^2 \sum_{l=0}^{k-1} \sum_{i=1}^n \langle (g_i^l - v_i^l) + \beta(v_i^l - \nabla f_i(x^l)) + \beta(\nabla f_i(x^l) - \nabla f_i(x^{l+1})), \theta_i^l \rangle$$

$$1745 + \frac{8\gamma\beta^2}{n\eta^2} (1 - \beta) \sum_{l=0}^{k-1} \sum_{i=1}^n \langle v_i^l - \nabla f_i(x^l), \theta_i^{l+1} \rangle + 2\gamma(1 - \beta) \sum_{l=0}^{k-1} \langle v^l - \nabla f(x^l), \theta^{l+1} \rangle$$

$$1746 + \frac{8\gamma\beta^2}{n\eta^2} (1 - \beta) \sum_{l=0}^{k-1} \sum_{i=1}^n \langle \nabla f_i(x^l) - \nabla f_i(x^{l+1}), \theta_i^{l+1} \rangle$$

$$1747 + 2\gamma(1 - \beta) \sum_{l=0}^{k-1} \langle \nabla f(x^l) - \nabla f(x^{l+1}), \theta^{l+1} \rangle.$$

1748 Denote the events Θ_i^t, Θ^t and N^{t+1} as

$$1749 \Theta_i^t := \{\|\theta_i^t\| \geq b\}, \quad \Theta^t := \{\|\theta^t\| \geq \frac{c}{\sqrt{n}}\}, \quad \text{and} \quad N^{t+1} := \left\{ \left\| \frac{1}{n} \sum_{l=1}^t \sum_{i=1}^n \omega_i^l \right\| \geq a \right\} \quad (36)$$

1750 respectively. From Assumption 2 we have

$$1751 \Pr(\Theta_i^t) \leq 2 \exp\left(-\frac{b^2}{2\sigma^2}\right) = \frac{\alpha}{6(T+1)n}$$

1752 where the last equality is by definition of b^2 . Therefore, $\Pr(\bar{\Theta}_i^t) \geq 1 - \frac{\alpha}{6(T+1)n}$.

1753 Besides, notice that the constant c in (28) can be viewed as

$$1754 c = (\sqrt{2} + 2b_3)\sigma \quad \text{where} \quad b_3^2 = 3 \log \frac{6(T+1)}{\alpha}.$$

Now we can use Lemma 1 to bound $\Pr(\Theta^t)$. Since all θ_i^t are independent σ -sub-Gaussian random vectors, then we have

$$\Pr\left(\left\|\sum_{i=1}^n \theta_i^t\right\| \geq c\sqrt{n}\right) = \Pr\left(\|\theta^t\| \geq \frac{c}{\sqrt{n}}\right) \leq \exp(-b_3^2/3) = \frac{\alpha}{6(T+1)}.$$

We also use Lemma 1 to bound $\Pr(N^t)$. Indeed, since all ω_i^l are independent Gaussian random vectors, then we have

$$\Pr\left(\left\|\sum_{l=1}^t \sum_{i=1}^n \omega_i^l\right\| \geq (\sqrt{2} + 2b_2) \sqrt{\sum_{l=1}^t \sum_{i=1}^n \sigma_\omega^2 d}\right) \leq \exp(-b_2^2/3) = \frac{\alpha}{6(T+1)}.$$

with $b_2^2 = 3 \log\left(\frac{4(T+1)}{\alpha}\right)$.

This implies that

$$\Pr\left(\left\|\frac{1}{n} \sum_{l=1}^t \sum_{i=1}^n \omega_i^l\right\| \geq a\right) \leq \frac{\alpha}{6(T+1)}$$

due to the choice of a from (28):

$$a = (\sqrt{2} + 2b_2)\sigma_\omega \sqrt{d} \sqrt{T/n} \quad \text{where} \quad b_2^2 = 3 \log \frac{6(T+1)}{\alpha}.$$

Note that with this choice of a we have that the above is true for any $t \in \{1, T\}$, i.e. $\Pr(N^t) \geq 1 - \frac{\alpha}{6(T+1)}$ for all $t \in \{1, T\}$.

Now we prove that $\Pr(E_t) \geq 1 - \frac{\alpha(t+1)}{T+1}$ for all $t \in \{0, \dots, T-1\}$. First, we show that the base of induction holds.

Base of induction.

1. $\|v_i^0 - g_i^{-1}\| = \|v_i^0\| = \beta \|\nabla f_i(x^0, \xi_i^0)\| = \beta \|\theta_i^0\| + \beta \|\nabla f_i(x^0)\| \leq \frac{1}{2}b + \frac{1}{2}B \leq \frac{1}{2}B + \frac{1}{2}B = B$ holds with probability $1 - \frac{\alpha}{6(T+1)}$. Indeed, we have

$$\Pr(\Theta_i^0) \leq 2 \exp\left(-\frac{b^2}{2\sigma^2}\right) = \frac{\alpha}{6(T+1)n}.$$

Therefore, we have

$$\Pr\left(\bigcap_{i=1}^n \bar{\Theta}_i^0\right) = 1 - \Pr\left(\bigcup_{i=1}^n \Theta_i^0\right) \geq 1 - \sum_{i=1}^n \Pr(\Theta_i^0) = 1 - n \frac{\alpha}{6(T+1)n} = 1 - \frac{\alpha}{6(T+1)}.$$

Moreover, by concentration inequality we have

$$\Pr(\Theta^0) \leq \frac{\alpha}{6(T+1)}.$$

This means that the probability of the event that each $\left\|\frac{1}{n} \sum_{l=1}^1 \sum_{i=1}^n \omega_i^l\right\| \leq a$, $\|\theta_i^0\| \leq b$, and $\|\theta^0\| \leq \frac{c}{\sqrt{n}}$, and is at least

$$1 - \frac{\alpha}{6(T+1)} - n \frac{\alpha}{6n(T+1)} - \frac{\alpha}{6(T+1)} = 1 - \frac{\alpha}{2(T+1)}.$$

2. We have already shown that

$$\Pr\left(\left\|\frac{1}{n} \sum_{i=1}^n \omega_i^1\right\| \geq a\right) \leq \frac{\alpha}{6(T+1)},$$

implying that $\left\|\frac{1}{n} \sum_{i=1}^n \omega_i^1\right\| \leq a$ with probability at least $1 - \frac{\alpha}{6(T+1)}$.

3. $g^0 = \frac{1}{n} \sum_{i=1}^n (g_i^{-1} + \text{clip}_\tau(v_i^0 - g_i^{-1})) = \frac{1}{n} \sum_{i=1}^n \text{clip}_\tau(\beta \nabla f_i(x^0, \xi_i^0))$. Therefore, we have

$$\begin{aligned}
\|g^0\| &\leq \left\| \frac{1}{n} \sum_{i=1}^n \beta \nabla f_i(x^0) + \beta \theta_i^0 + (\text{clip}_\tau(\beta \nabla f_i(x^0, \xi_i^0)) - \beta \nabla f_i(x^0, \xi_i^0)) \right\| \\
&\leq \beta \|\nabla f(x^0)\| + \frac{\beta}{n} \sum_{i=1}^n \|\theta_i^0\| + \frac{1}{n} \sum_{i=1}^n \max\{0, \beta \|\nabla f_i(x^0, \xi_i^0)\| - \tau\} \\
&\leq \beta \sqrt{2L(f(x^0) - f(x^*))} + \frac{\beta}{n} \sum_{i=1}^n \|\theta_i^0\| + \frac{1}{n} \sum_{i=1}^n \max\{0, \beta \|\nabla f_i(x^0)\| + \beta \|\theta_i^0\| - \tau\} \\
&\leq \frac{1}{2} \sqrt{2L\Phi^0} + \frac{2\beta}{n} \sum_{i=1}^n \|\theta_i^0\| + \frac{\beta}{n} \sum_{i=1}^n \|\nabla f_i(x^0)\| - \tau \\
&\leq \sqrt{64L\Delta} + 2\beta b + \beta B - \tau \\
&\leq \sqrt{64L\Delta} + \frac{3}{2}B - \tau + b \leq \sqrt{64L\Delta} + 3(B - \tau) + 2b + (L\gamma)^p a.
\end{aligned}$$

The inequalities above again hold in $\cap_{i=1}^n \bar{\Theta}_i^0$, i.e. with probability at least $1 - \frac{\alpha}{6(T+1)}$.

4. We have

$$\|v_i^0 - \nabla f_i(x^0)\| = \|\nabla f_i(x^0, \xi_i^0) - \nabla f_i(x^0)\| \leq b.$$

This bound holds with probability at least $1 - \frac{\alpha}{6(T+1)}$ because it holds in $\cap_{i=1}^n \bar{\Theta}_i^0$.

5. Inequalities 5 obviously also hold, as $\Phi^0 \leq 2\Phi^0 \leq 2\Delta$ by the choice of Δ .

Therefore, we conclude that the inequalities 1-7 hold with a probability at least

$$\begin{aligned}
\Pr\left(\Theta^0 \cap \left(\cap_{i=1}^n \bar{\Theta}_i^0\right) \cap \bar{N}^t\right) &\geq 1 - \Pr(\Theta^0) - \sum_{i=1}^n \Pr(\Theta_i^0) - \Pr(N^0) \\
&\geq 1 - \frac{\alpha}{6(T+1)} - n \cdot \frac{\alpha}{6n(T+1)} - \frac{\alpha}{6(T+1)} \\
&= 1 - \frac{\alpha}{2(T+1)} > 1 - \frac{\alpha}{T+1},
\end{aligned}$$

i.e. $\Pr(E_0) \geq 1 - \frac{\alpha}{T+1}$ holds. This is the base of the induction.

Transition step of induction.

CASE $|\mathcal{I}_{K+1}| > 0$. Assume that all events $\bar{\Theta}^{K+1}$, $\bar{\Theta}_i^{K+1}$ and \bar{N}^{K+1} take place, i.e. $\|\theta_i^{K+1}\| \leq b$, $\|\theta^{K+1}\| \leq \frac{c}{\sqrt{n}}$ for all $i \in [n]$ and $\left\| \frac{1}{n} \sum_{l=1}^{t+1} \sum_{i=1}^n \omega_i^l \right\| \leq a$. For that we need to work in $\bar{\Theta}^{K+1} \cap \left(\cap_{i=1}^n \bar{\Theta}_i^{K+1}\right) \cap \bar{N}^{K+1}$. Then, by the assumptions of the induction, from Lemma 15 we get for all $i \in \mathcal{I}_{K+1}$

$$\|v_i^{K+1} - g_i^K\| \leq \|v_i^K - g_i^{K-1}\| - \frac{\tau}{2} \leq B - \frac{\tau}{2}.$$

Therefore, from Lemma 13 we get that

$$\|g^{K+1}\| \leq \sqrt{64L\Delta} + 3(B - \tau) + 3b + 3a,$$

and from Lemma 14

$$\|\nabla f_i(x^{K+1}) - v_i^{K+1}\| \leq \sqrt{4L\Delta} + \frac{3}{2}(B - \tau) + 2b + (L\gamma)^p a.$$

This means that 1-5 in the induction assumption are also verified for the step $K + 1$.

Since we have for all $t \in \{0, \dots, K+1\}$ that inequalities 1-5 are verified, then we can write for each $t \in \{0, K\}$ by Lemmas 2 and 16 to 18 the following

$$\begin{aligned}
\Phi^{t+1} &= \delta^{t+1} + \frac{\gamma}{\eta} \tilde{V}^{t+1} + \frac{4\gamma\beta}{\eta^2} \tilde{P}^{t+1} + \frac{\gamma}{\beta} P^{t+1} \\
&\leq \delta^t - \frac{\gamma}{2} \|\nabla f(x^t)\|^2 - \frac{1}{4\gamma} R^t + \gamma \tilde{V}^t + \gamma P^t \\
&\quad + \frac{\gamma}{\eta} \left((1-\eta) \tilde{V}^t + \frac{4\beta^2}{\eta} \tilde{P}^t + \frac{4\beta^2 L^2}{\eta} R^t + \beta^2 b^2 \right. \\
&\quad \left. + \frac{2}{n} \beta (1-\eta^2) \sum_{i=1}^n \langle (g_i^t - v_i^t) + \beta(v_i^t - \nabla f_i(x^t)) + \beta(\nabla f_i(x^t) - \nabla f_i(x^{t+1})), \theta_i^{t+1} \rangle \right) \\
&\quad + \frac{4\gamma\beta}{\eta^2} \left((1-\beta) \tilde{P}^t + \frac{3L^2}{\beta} R^t + \beta^2 b^2 + \frac{2}{n} \beta (1-\beta) \sum_{i=1}^n \langle v_i^t - \nabla f_i(x^{t+1}), \theta_i^{t+1} \rangle \right) \\
&\quad + \frac{\gamma}{\beta} \left((1-\beta) P^t + \frac{3L^2}{\beta} R^t + \beta^2 \frac{c^2}{n} + 2\beta(1-\beta) \langle v^t - \nabla f(x^{t+1}), \theta^{t+1} \rangle \right)
\end{aligned}$$

Rearranging terms we get

$$\begin{aligned}
\Phi^{t+1} &\leq \delta^t - \frac{\gamma}{2} \|\nabla f(x^t)\|^2 + \frac{\gamma}{\eta} \tilde{V}^t (\eta + 1 - \eta) + \frac{4\gamma\beta}{\eta^2} \tilde{P}^t (\beta + 1 - \beta) + \frac{\gamma}{\beta} P^t (\beta + 1 - \beta) \\
&\quad - \frac{1}{4\gamma} R^t \left(1 - \frac{16L^2\beta^2}{\eta^2} \gamma^2 - \frac{48L^2}{\eta^2} \gamma^2 - \frac{12L^2}{\beta^2} \gamma^2 \right) + b^2 \left(\frac{\beta^2\gamma}{\eta} + \frac{4\gamma\beta^3}{\eta^2} \right) + c^2 \frac{\gamma\beta}{n} \\
&\quad + \frac{2\gamma\beta}{n\eta} (1-\eta)^2 \sum_{i=1}^n \langle (g_i^t - v_i^t) + \beta(v_i^t - \nabla f_i(x^t)) + \beta(\nabla f_i(x^t) - \nabla f_i(x^{t+1})), \theta_i^{t+1} \rangle \\
&\quad + \frac{8\gamma\beta^2}{n\eta^2} (1-\beta) \sum_{i=1}^n \langle v_i^t - \nabla f_i(x^t), \theta_i^{t+1} \rangle + 2\gamma(1-\beta) \langle v^t - \nabla f(x^t), \theta^{t+1} \rangle \\
&\quad + \frac{8\gamma\beta^2}{n\eta^2} (1-\beta) \sum_{i=1}^n \langle \nabla f_i(x^t) - \nabla f_i(x^{t+1}), \theta_i^{t+1} \rangle \\
&\quad + 2\gamma(1-\beta) \langle \nabla f(x^t) - \nabla f(x^{t+1}), \theta^{t+1} \rangle.
\end{aligned}$$

Using stepsize restriction (v_i) we get rid of the term with R^t and obtain

$$\begin{aligned}
\Phi^{t+1} &\leq \Phi^t - \frac{\gamma}{2} \|\nabla f(x^t)\|^2 + b^2 \left(\frac{\beta^2\gamma}{\eta} + \frac{4\gamma\beta^3}{\eta^2} \right) + c^2 \frac{\gamma\beta}{n} \\
&\quad + \frac{2\gamma\beta}{n\eta} (1-\eta)^2 \sum_{i=1}^n \langle (g_i^t - v_i^t) + \beta(v_i^t - \nabla f_i(x^t)) + \beta(\nabla f_i(x^t) - \nabla f_i(x^{t+1})), \theta_i^{t+1} \rangle \\
&\quad + \frac{8\gamma\beta^2}{n\eta^2} (1-\beta) \sum_{i=1}^n \langle v_i^t - \nabla f_i(x^t), \theta_i^{t+1} \rangle + 2\gamma(1-\beta) \langle v^t - \nabla f(x^t), \theta^{t+1} \rangle \\
&\quad + \frac{8\gamma\beta^2}{n\eta^2} (1-\beta) \sum_{i=1}^n \langle \nabla f_i(x^t) - \nabla f_i(x^{t+1}), \theta_i^{t+1} \rangle \\
&\quad + 2\gamma(1-\beta) \langle \nabla f(x^t) - \nabla f(x^{t+1}), \theta^{t+1} \rangle.
\end{aligned}$$

Now we sum all the inequalities above for $t \in \{0, \dots, K\}$ and get

$$\begin{aligned}
\Phi^{K+1} &\leq \Phi^0 - \frac{\gamma}{2} \sum_{t=0}^K \|\nabla f(x^t)\|^2 + Kb^2 \left(\frac{\beta^2 \gamma}{\eta} + \frac{4\gamma\beta^3}{\eta^2} \right) + Kc^2 \frac{\gamma\beta}{n} \\
&\quad + \frac{2\gamma\beta}{n\eta} (1-\eta)^2 \sum_{t=0}^K \sum_{i=1}^n \langle (g_i^t - v_i^t) + \beta(v_i^t - \nabla f_i(x^t)) + \beta(\nabla f_i(x^t) - \nabla f_i(x^{t+1})), \theta_i^{t+1} \rangle \\
&\quad + \frac{8\gamma\beta^2}{n\eta^2} (1-\beta) \sum_{t=0}^K \sum_{i=1}^n \langle v_i^t - \nabla f_i(x^t), \theta_i^{t+1} \rangle + 2\gamma(1-\beta) \sum_{t=0}^K \langle v^t - \nabla f(x^t), \theta^{t+1} \rangle \\
&\quad + \frac{8\gamma\beta^2}{n\eta^2} (1-\beta) \sum_{t=0}^K \sum_{i=1}^n \langle \nabla f_i(x^t) - \nabla f_i(x^{t+1}), \theta_i^{t+1} \rangle \\
&\quad + 2\gamma(1-\beta) \sum_{t=0}^K \langle \nabla f(x^t) - \nabla f(x^{t+1}), \theta^{t+1} \rangle. \tag{37}
\end{aligned}$$

Rearranging terms we get

$$\begin{aligned}
\frac{\gamma}{2} \sum_{t=0}^K \|\nabla f(x^t)\|^2 &\leq \Phi^0 - \Phi^{K+1} + Kb^2 \left(\frac{\beta^2 \gamma}{\eta} + \frac{4\gamma\beta^3}{\eta^2} \right) + Kc^2 \frac{\gamma\beta}{n} \\
&\quad + \frac{2\gamma\beta}{n\eta} (1-\eta)^2 \sum_{t=0}^K \sum_{i=1}^n \langle (g_i^t - v_i^t) + \beta(v_i^t - \nabla f_i(x^t)) + \beta(\nabla f_i(x^t) - \nabla f_i(x^{t+1})), \theta_i^{t+1} \rangle \\
&\quad + \frac{8\gamma\beta^2}{n\eta^2} (1-\beta) \sum_{t=0}^K \sum_{i=1}^n \langle v_i^t - \nabla f_i(x^t), \theta_i^{t+1} \rangle + 2\gamma(1-\beta) \sum_{t=0}^K \langle v^t - \nabla f(x^t), \theta^{t+1} \rangle \\
&\quad + \frac{8\gamma\beta^2}{n\eta^2} (1-\beta) \sum_{t=0}^K \sum_{i=1}^n \langle \nabla f_i(x^t) - \nabla f_i(x^{t+1}), \theta_i^{t+1} \rangle \\
&\quad + 2\gamma(1-\beta) \sum_{t=0}^K \langle \nabla f(x^t) - \nabla f(x^{t+1}), \theta^{t+1} \rangle.
\end{aligned}$$

Taking into account that $\frac{\gamma}{2} \sum_{t=0}^K \|\nabla f(x^t)\|^2 \geq 0$, we get that the event $E_K \cap \left(\bigcap_{i=1}^n \bar{\Theta}_i^{K+1} \right) \cap \bar{N}^t$ implies

$$\begin{aligned}
\Phi^{K+1} &\leq \Phi^0 + Kb^2 \left(\frac{\beta^2 \gamma}{\eta} + \frac{4\gamma\beta^3}{\eta^2} \right) + Kc^2 \frac{\gamma\beta}{n} \\
&\quad + \frac{2\gamma\beta}{n\eta} (1-\eta)^2 \sum_{t=0}^K \sum_{i=1}^n \langle (g_i^t - v_i^t) + \beta(v_i^t - \nabla f_i(x^t)) + \beta(\nabla f_i(x^t) - \nabla f_i(x^{t+1})), \theta_i^{t+1} \rangle \\
&\quad + \frac{8\gamma\beta^2}{n\eta^2} (1-\beta) \sum_{t=0}^K \sum_{i=1}^n \langle v_i^t - \nabla f_i(x^t), \theta_i^{t+1} \rangle + \frac{2\gamma(1-\beta)}{n} \sum_{t=0}^K \sum_{i=1}^n \langle v^t - \nabla f(x^t), \theta_i^{t+1} \rangle \\
&\quad + \frac{8\gamma\beta^2}{n\eta^2} (1-\beta) \sum_{t=0}^K \sum_{i=1}^n \langle \nabla f_i(x^t) - \nabla f_i(x^{t+1}), \theta_i^{t+1} \rangle \\
&\quad + \frac{2\gamma(1-\beta)}{n} \sum_{t=0}^K \sum_{i=1}^n \langle \nabla f(x^t) - \nabla f(x^{t+1}), \theta_i^{t+1} \rangle.
\end{aligned}$$

Now we define the following random vectors:

$$\zeta_{1,i}^t := \begin{cases} g_i^t - v_i^t, & \text{if } \|g_i^t - v_i^t\| \leq B, \\ 0, & \text{otherwise} \end{cases},$$

$$\zeta_{2,t}^t := \begin{cases} v_i^t - \nabla f_i(x^t), & \text{if } \|v_i^t - \nabla f_i(x^t)\| \leq \sqrt{4L\Delta} + \frac{3}{2}(B - \tau) + 2b + (L\gamma)^p a, \\ 0, & \text{otherwise} \end{cases},$$

$$\zeta_{3,i}^t := \begin{cases} \nabla f_i(x^t) - \nabla f_i(x^{t+1}), & \text{if } \|\nabla f_i(x^t) - \nabla f_i(x^{t+1})\| \leq L\gamma \left(\sqrt{64L\Delta} + 3(B - \tau) + 3b + 3a \right), \\ 0, & \text{otherwise} \end{cases},$$

$$\zeta_4^t := \begin{cases} v^t - \nabla f(x^t), & \text{if } \|v^t - \nabla f(x^t)\| \leq \sqrt{4L\Delta} + \frac{3}{2}(B - \tau) + 2b + (L\gamma)^p a, \\ 0, & \text{otherwise} \end{cases},$$

$$\zeta_5^t := \begin{cases} \nabla f(x^t) - \nabla f(x^{t+1}), & \text{if } \|\nabla f(x^t) - \nabla f(x^{t+1})\| \leq L\gamma \left(\sqrt{64L\Delta} + 3(B - \tau) + 3b + 3a \right), \\ 0, & \text{otherwise} \end{cases}.$$

By definition, all introduced random vectors $\zeta_{l,i}^t, l \in [3], i \in [n], \zeta_{4,5}^t$ are bounded with probability

1. Moreover, by the definition of Φ^t and definition of E_t we get that the event $E_K \cap \bar{\Theta}^{K+1} \cap \left(\bigcap_{i=1}^n \bar{\Theta}_i^{K+1} \right) \cap \bar{N}^{K+1}$ implies

$$\begin{aligned} \zeta_{1,i}^t &= g_i^t - v_i^t, & \zeta_{2,i}^t &= v_i^t - \nabla f_i(x^t), & \zeta_{3,i}^t &= \nabla f_i(x^t) - \nabla f_i(x^{t+1}), \\ \zeta_4^t &= v^t - \nabla f(x^t), & \zeta_5^t &= \nabla f(x^t) - \nabla f(x^{t+1}). \end{aligned}$$

Therefore, the event $E_K \cap \bar{\Theta}^{K+1} \cap \left(\bigcap_{i=1}^n \bar{\Theta}_i^{K+1} \right) \cap \bar{N}^{K+1}$ implies

$$\begin{aligned} \Phi^{K+1} &\leq \Phi^0 + \underbrace{Kb^2 \left(\frac{\beta^2\gamma}{\eta} + \frac{4\gamma\beta^3}{\eta^2} \right)}_{\textcircled{1}} + \underbrace{Kc^2 \frac{\gamma\beta}{n} + \frac{2\gamma\beta}{n\eta} (1-\eta)^2 \sum_{t=0}^K \sum_{i=1}^n \langle \zeta_{1,i}^t, \theta_i^{t+1} \rangle}_{\textcircled{2}} \\ &+ \underbrace{\frac{2\gamma\beta^2}{n\eta} (1-\eta)^2 \sum_{t=0}^K \sum_{i=1}^n \langle \zeta_{2,i}^t, \theta_i^{t+1} \rangle}_{\textcircled{3}} + \underbrace{\frac{2\gamma\beta^2}{n\eta} (1-\eta)^2 \sum_{t=0}^K \sum_{i=1}^n \langle \zeta_{3,i}^t, \theta_i^{t+1} \rangle}_{\textcircled{4}} \\ &+ \underbrace{\frac{8\gamma\beta^2}{n\eta^2} (1-\beta) \sum_{t=0}^K \sum_{i=1}^n \langle \zeta_{2,i}^t, \theta_i^{t+1} \rangle}_{\textcircled{5}} + \underbrace{\frac{2\gamma(1-\beta)}{n} \sum_{t=0}^K \sum_{i=1}^n \langle \zeta_4^t, \theta_i^{t+1} \rangle}_{\textcircled{6}} \\ &+ \underbrace{\frac{8\gamma\beta^2}{n\eta^2} (1-\beta) \sum_{t=0}^K \sum_{i=1}^n \langle \zeta_{3,i}^t, \theta_i^{t+1} \rangle}_{\textcircled{7}} + \underbrace{\frac{2\gamma(1-\beta)}{n} \sum_{t=0}^K \sum_{i=1}^n \langle \zeta_5^t, \theta_i^{t+1} \rangle}_{\textcircled{8}}. \end{aligned}$$

BOUND OF THE TERM ①. For the term ① we have

$$Kb^2 \left(\frac{\beta^2\gamma}{\eta} + \frac{4\gamma\beta^3}{\eta^2} \right) + Kc^2 \frac{\gamma\beta}{n} \leq Kb^2 \left(\frac{9(\beta/3)^{\frac{3-2p}{1-p}}}{L\eta} + \frac{108(\beta/3)^{\frac{4-3p}{1-p}}}{L\eta^2} \right) + 3Kc^2 \frac{(\beta/3)^{\frac{2-p}{1-p}}}{Ln}.$$

By choosing γ such that

$$\beta \leq \min \left\{ 3 \left(\frac{L\Delta\eta}{216Tb^2} \right)^{\frac{1-p}{3-2p}}, 3 \left(\frac{L\Delta\eta^2}{48Tb^2} \right)^{\frac{1-p}{4-3p}}, 3 \left(\frac{L\Delta n}{72Tc^2} \right)^{\frac{1-p}{2-p}} \right\} \quad (38)$$

we get that

$$Kb^2 \left(\frac{\beta^2\gamma}{\eta} + \frac{4\gamma\beta^3}{\eta^2} \right) + Kc^2 \frac{\gamma\beta}{n} \leq 3 \cdot \frac{\Delta}{24} = \frac{\Delta}{8}.$$

This bound holds with probability 1. Note that the worst dependency in the restriction on β in T is $\mathcal{O}(1/T^{\frac{1-p}{2-p}})$ that comes from the last term in min.

2052 BOUND OF THE TERM ②. For term ②, let us enumerate random variables as

$$2053 \langle \zeta_{1,1}^0, \theta_1^1 \rangle, \dots, \langle \zeta_{1,n}^0, \theta_n^1 \rangle, \langle \zeta_{1,1}^1, \theta_1^2 \rangle, \dots, \langle \zeta_{1,n}^1, \theta_n^2 \rangle, \dots, \langle \zeta_{1,1}^K, \theta_1^{K+1} \rangle, \dots, \langle \zeta_{1,n}^K, \theta_n^{K+1} \rangle,$$

2055 i.e. first by index i , then by index t . Then we have that the event $E_K \cap \left(\bigcap_{i=1}^n \bar{\Theta}_i^{K+1} \right)$ implies

$$2057 \mathbb{E} \left[\frac{2\gamma\beta}{n\eta} (1-\eta)^2 \langle \zeta_{1,i}^l, \theta_i^{l+1} \rangle \mid \langle \zeta_{1,i-1}^l, \theta_{i-1}^{l+1} \rangle, \dots, \langle \zeta_{1,1}^l, \theta_1^{l+1} \rangle, \dots, \langle \zeta_{1,1}^0, \theta_1^1 \rangle \right] = 0,$$

2059 because $\{\theta_i^{l+1}\}_{i=1}^n$ are independent. Let

$$2061 \sigma_2^2 := \frac{4\gamma^2\beta^2}{n^2\eta^2} \cdot B^2 \cdot \sigma^2.$$

2063 Since θ_i^{l+1} is σ -sub-Gaussian random vector, we have

$$\begin{aligned} 2064 & \mathbb{E} \left[\exp \left(\left| \frac{1}{\sigma_2^2} \frac{4\gamma^2\beta^2}{n^2\eta^2} (1-\eta)^4 \langle \zeta_{1,i}^l, \theta_i^{l+1} \rangle^2 \right) \mid l, i-1 \right] \\ 2065 & \leq \mathbb{E} \left[\exp \left(\frac{1}{\sigma_1^2} \frac{4\gamma^2\beta^2}{n^2\eta^2} \|\zeta_{1,i}^l\|^2 \cdot \|\theta_i^{l+1}\|^2 \right) \mid l, i-1 \right] \\ 2066 & \leq \mathbb{E} \left[\exp \left(\frac{1}{\sigma_2^2} \frac{4\gamma^2\beta^2}{n^2\eta^2} \cdot B^2 \|\theta_i^{l+1}\|^2 \right) \mid l, i-1 \right] \\ 2067 & \leq \mathbb{E} \left[\exp \left(\frac{n^2\eta^2}{4\gamma^2\beta^2 \cdot B^2 \cdot \sigma^2} \frac{4\gamma^2\beta^2}{n^2\eta^2} \cdot B^2 \|\theta_i^{l+1}\|^2 \right) \mid l, i-1 \right] \\ 2068 & = \mathbb{E} \left[\exp \left(\frac{\|\theta_i^{l+1}\|^2}{\sigma^2} \mid l, i-1 \right) \right] \leq \exp(1). \end{aligned}$$

2074 Here $\mathbb{E}[\cdot \mid l, i-1]$ means

$$2075 \mathbb{E}[\cdot \mid \langle \zeta_{1,i-1}^l, \theta_{i-1}^{l+1} \rangle, \dots, \langle \zeta_{1,1}^l, \theta_1^{l+1} \rangle, \dots, \langle \zeta_{1,1}^0, \theta_1^1 \rangle] = 0,$$

2077 Therefore, we have by Lemma 1 with $\sigma_k^2 \equiv \sigma_2^2$ that

$$\begin{aligned} 2081 & \Pr \left(\frac{2\gamma\beta}{n\eta} (1-\eta)^2 \left\| \sum_{t=0}^K \sum_{i=1}^n \langle \zeta_{1,i}^t, \theta_i^{t+1} \rangle \right\| \geq (\sqrt{2} + \sqrt{2}b_1) \sqrt{\sum_{t=0}^K \sum_{i=1}^n \frac{4B^2\gamma^2\beta^2\sigma^2}{n^2\eta^2}} \right) \\ 2082 & \leq \exp(-b_1^2/3) \\ 2083 & = \frac{\alpha}{14(T+1)} \end{aligned}$$

2084 with $b_1^2 = 3 \log \left(\frac{14(T+1)}{\alpha} \right)$. Note that

$$\begin{aligned} 2089 & (\sqrt{2} + \sqrt{2}b_1) \sqrt{\sum_{t=0}^K \sum_{i=1}^n \frac{4B^2\gamma^2\beta^2\sigma^2}{n^2\eta^2}} \leq (\sqrt{2} + \sqrt{2}b_1) \sqrt{\sum_{t=0}^K \sum_{i=1}^n \frac{36B^2(\beta/3)^{\frac{4-p}{1-p}} \sigma^2}{L^2 n^2 \eta^2}} \\ 2090 & = (\sqrt{2} + \sqrt{2}b_1) \frac{6B(\beta/3)^{\frac{2-p}{1-p}} \sigma}{Ln\eta} \sqrt{(K+1)n} \\ 2091 & \leq \frac{\Delta}{8}, \end{aligned}$$

2092 because we choose

$$2093 \beta \leq 3 \left(\frac{L\Delta\sqrt{n\eta}}{48\sqrt{2}(1+b_1)B\sigma\sqrt{T}} \right)^{\frac{1-p}{2-p}}, \quad \text{and} \quad K+1 \leq T. \quad (39)$$

2094 This implies that

$$2095 \Pr \left(\frac{2\gamma\beta}{n\eta} (1-\eta)^2 \left\| \sum_{t=0}^K \sum_{i=1}^n \langle \zeta_{1,i}^t, \theta_i^{t+1} \rangle \right\| \geq \frac{\Delta}{8} \right) \leq \frac{\alpha}{14(T+1)}$$

2096 with this choice of momentum parameter. The dependency on T is $\tilde{O}(1/T^{\frac{1-p}{2(2-p)}})$.

2106 BOUND OF THE TERM ③. The bound in this case is similar to the previous one. Let

$$2108 \quad \sigma_3^2 := \frac{4\gamma^2\beta^4}{n^2\eta^2} \cdot \left(\sqrt{4L\Delta} + \frac{3}{2}(B - \tau) + 2b + (L\gamma)^p a \right)^2 \cdot \sigma^2.$$

$$2110 \quad \mathbb{E} \left[\exp \left(\left| \frac{1}{\sigma_3^2} \frac{4\gamma^2\beta^4}{n^2\eta^2} (1 - \eta)^4 \langle \zeta_{2,i}^l, \theta_i^{l+1} \rangle^2 \right| \right) \mid l, i - 1 \right]$$

$$2112 \quad \leq \mathbb{E} \left[\exp \left(\frac{1}{\sigma_3^2} \frac{4\gamma^2\beta^4}{n^2\eta^2} \|\zeta_{2,i}^l\|^2 \cdot \|\theta_i^{l+1}\|^2 \right) \right]$$

$$2114 \quad \leq \mathbb{E} \left[\exp \left(\frac{1}{\sigma_3^2} \frac{4\gamma^2\beta^4}{n^2\eta^2} \cdot \left(\sqrt{4L\Delta} + \frac{3}{2}(B - \tau) + 2b + (L\gamma)^p a \right)^2 \cdot \|\theta_i^{l+1}\|^2 \right) \mid l, i - 1 \right]$$

$$2116 \quad \leq \mathbb{E} \left[\exp \left(\left[\frac{4\gamma^2\beta^4}{n^2\eta^2} \cdot \left(\sqrt{4L\Delta} + \frac{3}{2}(B - \tau) + 2b + (L\gamma)^p a \right)^2 \cdot \sigma^2 \right]^{-1} \right) \right]$$

$$2118 \quad \frac{4\gamma^2\beta^4}{n^2\eta^2} \cdot \left(\sqrt{4L\Delta} + \frac{3}{2}(B - \tau) + 2b + (L\gamma)^p a \right)^2 \cdot \|\theta_i^{l+1}\|^2 \mid l, i - 1 \Big]$$

$$2120 \quad = \mathbb{E} \left[\exp \left(\frac{\|\theta_i^{l+1}\|^2}{\sigma^2} \right) \mid l, i - 1 \right] \leq \exp(1).$$

2122 Therefore, we have by Lemma 1 that

$$2124 \quad \Pr \left[\frac{2\gamma\beta^2}{n\eta} (1 - \eta)^2 \left\| \sum_{t=0}^K \sum_{i=1}^n \langle \zeta_{2,i}^t, \theta_i^{t+1} \rangle \right\| \right]$$

$$2126 \quad \geq (\sqrt{2} + \sqrt{2}b_1) \sqrt{\sum_{t=0}^K \sum_{i=1}^n \frac{4\gamma^2\beta^4\sigma^2}{n^2\eta^2} \cdot \left(\sqrt{4L\Delta} + \frac{3}{2}(B - \tau) + 2b + (L\gamma)^p a \right)^2}$$

$$2128 \quad \leq \exp(-b_1^2/3) = \frac{\alpha}{14(T+1)},$$

2130 Note that

$$2132 \quad (\sqrt{2} + \sqrt{2}b_1) \sqrt{(K+1)n} \frac{2\gamma\beta^2\sigma}{\eta n} \left(\sqrt{4L\Delta} + \frac{3}{2}(B - \tau) + 2b + (L\gamma)^p a \right)$$

$$2134 \quad \leq (\sqrt{2} + \sqrt{2}b_1) \sqrt{(K+1)n} \frac{18(\beta/3)^{\frac{3-2p}{1-p}} \sigma}{L\eta n} \left(\sqrt{4L\Delta} + \frac{3}{2}(B - \tau) + 2b + (\beta/3)^{\frac{1-p}{1-p}} a \right)$$

$$2136 \quad \leq \frac{\Delta}{8}.$$

2138 because we choose

$$2140 \quad \beta \leq \min \left\{ 3 \left(\frac{L\Delta\eta\sqrt{n}}{288\sqrt{2}(1+b_1)\sigma\sqrt{T} \left(\sqrt{4L\Delta} + \frac{3}{2}(B - \tau) + 2b \right)} \right)^{\frac{1-p}{3-2p}} \right.$$

$$2142 \quad \left. 3 \left(\frac{L\Delta\eta\sqrt{n}}{288\sqrt{2}(1+b_1)\sqrt{T}\sigma a} \right)^{\frac{1-p}{3-p}} \right\},$$

2144 and $K+1 \leq T$.

(40)

2146 This implies

$$2148 \quad \Pr \left(\frac{2\gamma\beta^2}{n\eta} (1 - \eta)^2 \left\| \sum_{t=0}^K \sum_{i=1}^n \langle \zeta_{2,i}^t, \theta_i^{t+1} \rangle \right\| \geq \frac{\Delta}{8} \right) \leq \frac{\alpha}{14(T+1)}.$$

2150 Note that the worst dependency w.r.t. T is $\tilde{O}(1/T^{\frac{3(1-p)}{2(3-p)}})$ since $a \approx \sigma_\omega \sqrt{T} \sim T$.

2160 BOUND OF THE TERM ④. The bound in this case is similar to the previous one. Let

$$2161 \sigma_4^2 := \frac{4L^2\gamma^4\beta^4}{n^2\eta^2} \left(\sqrt{64L\Delta} + 3(B - \tau) + 3b + 3a \right)^2 \cdot \sigma^2.$$

2162 Then we have

$$2163 \begin{aligned} & \mathbb{E} \left[\exp \left(\left| \frac{1}{\sigma_4^2} \frac{4\gamma^2\beta^4}{n^2\eta^2} (1 - \eta)^4 \langle \zeta_{3,i}^l, \theta_i^{l+1} \rangle^2 \right| \right) \mid l, i - 1 \right] \\ & \leq \mathbb{E} \left[\exp \left(\frac{1}{\sigma_4^2} \frac{4\gamma^2\beta^4}{n^2\eta^2} \|\zeta_{3,i}^l\|^2 \cdot \|\theta_i^{l+1}\|^2 \right) \mid l, i - 1 \right] \\ & \leq \mathbb{E} \left[\exp \left(\frac{1}{\sigma_4^2} \frac{4\gamma^2\beta^4}{n^2\eta^2} \cdot L^2\gamma^2 \left(\sqrt{64L\Delta} + 3(B - \tau) + 3b + 3a \right)^2 \cdot \|\theta_i^{l+1}\|^2 \right) \mid l, i - 1 \right] \\ & \leq \mathbb{E} \left[\exp \left(\left[\frac{4L^2\gamma^4\beta^4}{n^2\eta^2} \left(\sqrt{64L\Delta} + 3(B - \tau) + 3b + 3a \right)^2 \cdot \sigma^2 \right]^{-1} \right. \right. \\ & \quad \left. \left. \frac{4L^2\gamma^4\beta^4}{n^2\eta^2} \cdot \frac{16\beta^4}{81} \left(\sqrt{64L\Delta} + 3(B - \tau) + 3b + 3a \right)^2 \cdot \|\theta_i^{l+1}\|^2 \right) \mid l, i - 1 \right] \\ & = \mathbb{E} \left[\exp \left(\frac{\|\theta_i^{l+1}\|^2}{\sigma^2} \right) \right] \leq \exp(1). \end{aligned}$$

2175 Therefore, we have by Lemma 1 that

$$2176 \begin{aligned} & \Pr \left(\frac{2\gamma\beta^2}{n\eta} (1 - \eta)^2 \left\| \sum_{t=0}^K \sum_{i=1}^n \langle \zeta_{3,i}^t, \theta_i^{t+1} \rangle \right\| \right. \\ & \quad \left. \geq (\sqrt{2} + \sqrt{2}b_1) \sqrt{\sum_{t=0}^K \sum_{i=1}^n \frac{4L^2\gamma^4\beta^4\sigma^2}{n^2\eta^2} \cdot \left(\sqrt{64L\Delta} + 3(B - \tau + b) + 3a \right)^2} \right) \\ & \leq \exp(-b_1^2/3) = \frac{\alpha}{14(T+1)}. \end{aligned}$$

2180 Note that

$$2181 \begin{aligned} & (\sqrt{2} + \sqrt{2}b_1) \sqrt{(K+1)n} \frac{2L\gamma^2\beta^2\sigma}{\eta n} \left(\sqrt{64L\Delta} + 3(B - \tau + b) + 3a \right) \\ & \leq \sqrt{2}(1 + b_1) \sqrt{(K+1)n} \frac{18(\beta/3)^{\frac{4-2p}{1-p}} \sigma}{L\eta n} \left(\sqrt{64L\Delta} + 3(B - \tau + b) + 3a \right) \\ & \leq \frac{\Delta}{8}. \end{aligned}$$

2185 because we choose

$$2186 \begin{aligned} & \beta \leq \min \left\{ 3 \left(\frac{L\Delta\eta\sqrt{n}}{288\sqrt{2}(1+b_1)\sigma\sqrt{T} \left(\sqrt{64L\Delta} + 3(B - \tau + b) \right)} \right)^{\frac{1-p}{4-2p}} \right. \\ & \quad \left. 3 \left(\frac{L\Delta\eta\sqrt{n}}{288\sqrt{2}(1+b_1)\sigma\sqrt{Ta}} \right)^{\frac{1-p}{4-2p}} \right\}, \\ & \text{and } K + 1 \leq T. \end{aligned} \tag{41}$$

2190 This implies

$$2191 \Pr \left(\frac{2\gamma\beta^2}{n\eta} (1 - \eta)^2 \left\| \sum_{t=0}^K \sum_{i=1}^n \langle \zeta_{2,i}^t, \theta_i^{t+1} \rangle \right\| \geq \frac{\Delta}{8} \right) \leq \frac{\alpha}{14(T+1)},$$

2192 Note that the worst dependency w.r.t. T is $\tilde{O}(1/T^{\frac{3(1-p)}{4(2-p)}})$ since $a \sim T$.

2214 BOUND OF THE TERM ⑤. The bound in this case is similar to the previous one. Let

$$2215 \sigma_5^2 := \frac{64\gamma^2\beta^4}{n^2\eta^4} \cdot \left(\sqrt{4L\Delta} + \frac{3}{2}(B - \tau) + 2b + (L\gamma)^p a \right)^2 \cdot \sigma^2.$$

2218 Then we have

$$\begin{aligned} 2219 & \mathbb{E} \left[\exp \left(\left| \frac{1}{\sigma_5^2} \frac{64\gamma^2\beta^4}{n^2\eta^4} (1 - \beta)^2 \langle \zeta_{2,i}^l, \theta_i^{l+1} \rangle^2 \right| \right) \mid l, i - 1 \right] \\ 2220 & \leq \mathbb{E} \left[\exp \left(\frac{1}{\sigma_5^2} \frac{64\gamma^2\beta^4}{n^2\eta^4} \|\zeta_{2,i}^l\|^2 \cdot \|\theta_i^{l+1}\|^2 \right) \mid l, i - 1 \right] \\ 2221 & \leq \mathbb{E} \left[\exp \left(\frac{1}{\sigma_5^2} \frac{64\gamma^2\beta^4}{n^2\eta^4} \cdot \left(\sqrt{4L\Delta} + \frac{3}{2}(B - \tau) + 2b + (L\gamma)^p a \right)^2 \cdot \|\theta_i^{l+1}\|^2 \right) \mid l, i - 1 \right] \\ 2222 & = \mathbb{E} \left[\exp \left(\left[\frac{64\gamma^2\beta^4}{81L^2n^2\eta^4} \cdot \left(\sqrt{4L\Delta} + \frac{3}{2}(B - \tau) + 2b + (L\gamma)^p a \right)^2 \cdot \sigma^2 \right]^{-1} \right. \right. \\ 2223 & \quad \left. \left. \frac{64\gamma^2\beta^4}{n^2\eta^4} \cdot \left(\sqrt{4L\Delta} + \frac{3}{2}(B - \tau) + 2b + (L\gamma)^p a \right)^2 \cdot \|\theta_i^{l+1}\|^2 \right) \mid l, i - 1 \right] \\ 2224 & = \mathbb{E} \left[\exp \left(\frac{\|\theta_i^{l+1}\|^2}{\sigma^2} \right) \mid l, i - 1 \right] \leq \exp(1). \end{aligned}$$

2236 Therefore, we have by Lemma 1 that

$$\begin{aligned} 2237 & \Pr \left[\frac{8\gamma\beta^2}{n\eta^2} (1 - \beta) \left\| \sum_{t=0}^K \sum_{i=1}^n \langle \zeta_{2,i}^t, \theta_i^{t+1} \rangle \right\| \right. \\ 2238 & \quad \left. \geq (\sqrt{2} + \sqrt{2}b_1) \sqrt{\sum_{t=0}^K \sum_{i=1}^n \frac{64\gamma^2\beta^4\sigma^2}{n^2\eta^4} \left(\sqrt{4L\Delta} + \frac{3}{2}(B - \tau) + 2b + (L\gamma)^p a \right)^2} \right] \\ 2239 & \leq \exp(-b_1^2/3) = \frac{\alpha}{14(T+1)}. \end{aligned}$$

2246 Note that

$$\begin{aligned} 2247 & (\sqrt{2} + \sqrt{2}b_1) \sqrt{(K+1)n} \frac{8\gamma\beta^2\sigma}{n\eta^2} \left(\sqrt{4L\Delta} + \frac{3}{2}(B - \tau) + 2b + (L\gamma)^p a \right) \\ 2248 & \leq (\sqrt{2} + \sqrt{2}b_1) \sqrt{(K+1)n} \frac{72(\beta/3)^{\frac{3-2p}{1-p}} \sigma}{Ln\eta^2} \left(\sqrt{4L\Delta} + \frac{3}{2}(B - \tau) + 2b + (\beta/3)^{\frac{p}{1-p}} a \right) \\ 2249 & \leq \frac{\Delta}{8} \end{aligned}$$

2254 because we choose

$$\begin{aligned} 2255 & \beta \leq \min \left\{ 3 \left(\frac{L\Delta\eta^2\sqrt{n}}{1152\sqrt{2}(1+b_1)\sigma\sqrt{T} \left(\sqrt{4L\Delta} + \frac{3}{2}(B - \tau) + 2b \right)} \right)^{\frac{1-p}{3-2p}} \right. \\ 2256 & \quad \left. 3 \left(\frac{L\Delta\eta^2\sqrt{n}}{1152\sqrt{2}(1+b_1)\sigma\sqrt{T}a} \right)^{\frac{1-p}{3-p}} \right\}, \\ 2257 & \text{and } K+1 \leq T. \end{aligned} \tag{42}$$

2263 This implies

$$2264 \Pr \left(\frac{8\gamma\beta^2}{n\eta^2} (1 - \beta) \left\| \sum_{t=0}^K \sum_{i=1}^n \langle \zeta_{2,i}^t, \theta_i^{t+1} \rangle \right\| \geq \frac{\Delta}{8} \right) \leq \frac{\alpha}{14(T+1)}.$$

2266 Note that the worst dependency w.r.t. T is $\tilde{O}(1/T^{\frac{3(1-p)}{2(3-p)}})$ since $a \sim T$.

2268 BOUND OF THE TERM \mathcal{O} . The bound in this case is similar to the previous one. Let

$$2270 \sigma_7^2 := \frac{64\gamma^2\beta^4}{n^2\eta^4} \left(\sqrt{64L\Delta} + 3(B - \tau + b) + 3a \right)^2 \cdot \sigma^2.$$

2272 Then we have

$$\begin{aligned} 2273 & \mathbb{E} \left[\exp \left(\left| \frac{1}{\sigma_7^2} \frac{64L^2\gamma^4\beta^4}{n^2\eta^4} (1 - \beta)^2 \langle \zeta_{3,i}^l, \theta_i^{l+1} \rangle^2 \right| \right) \mid l, i - 1 \right] \\ 2274 & \leq \mathbb{E} \left[\exp \left(\frac{1}{\sigma_7^2} \frac{64\gamma^2\beta^4}{n^2\eta^4} \|\zeta_{3,i}^l\|^2 \cdot \|\theta_i^{l+1}\|^2 \right) \mid l, i - 1 \right] \\ 2275 & \leq \mathbb{E} \left[\exp \left(\frac{64\gamma^2\beta^4}{n^2\eta^4} \cdot L^2\gamma^2 \left(\sqrt{64L\Delta} + 3(B - \tau + b) + 3a \right)^2 \cdot \|\theta_i^{l+1}\|^2 \right) \mid l, i - 1 \right] \\ 2276 & \leq \mathbb{E} \left[\exp \left(\left[\frac{64L^2\gamma^4\beta^4}{n^2\eta^4} \left(\sqrt{64L\Delta} + 3(B - \tau + b) + 3a \right)^2 \cdot \sigma^2 \right]^{-1} \right. \right. \\ 2277 & \quad \left. \left. \frac{64L^2\gamma^4\beta^4}{n^2\eta^4} \left(\sqrt{64L\Delta} + 3(B - \tau + b) + 3a \right)^2 \cdot \|\theta_i^{l+1}\|^2 \right) \mid l, i - 1 \right] \\ 2278 & = \mathbb{E} \left[\exp \left(\frac{\|\theta_i^{l+1}\|^2}{\sigma^2} \right) \mid l, i - 1 \right] \leq \exp(1). \end{aligned}$$

2288 Therefore, we have by Lemma 1 that

$$\begin{aligned} 2290 & \Pr \left[\frac{8\gamma\beta^2}{n\eta^2} (1 - \beta) \left\| \sum_{t=0}^K \sum_{i=1}^n \langle \zeta_{3,i}^t, \theta_i^{t+1} \rangle \right\| \geq \right. \\ 2291 & \quad \left. (\sqrt{2} + \sqrt{2}b_1) \sqrt{\sum_{t=0}^K \sum_{i=1}^n \frac{64L^2\gamma^4\beta^4\sigma^2}{n^2\eta^4} \cdot \left(\sqrt{64L\Delta} + 3(B - \tau + b) + 3a \right)^2} \right] \\ 2292 & \leq \exp(-b_1^2/3) = \frac{\alpha}{14(T+1)}. \end{aligned}$$

2299 Note that

$$\begin{aligned} 2300 & (\sqrt{2} + \sqrt{2}b_1) \sqrt{(K+1)n} \frac{8L\gamma^2\beta^2\sigma}{\eta^2n} \left(\sqrt{64L\Delta} + 3(B - \tau + b) + 3a \right) \\ 2301 & \leq (\sqrt{2} + \sqrt{2}b_1) \sqrt{(K+1)n} \frac{72(\beta/3)^{\frac{4-2p}{1-p}} \sigma}{L\eta^2n} \left(\sqrt{64L\Delta} + 3(B - \tau + b) + 3a \right) \\ 2302 & \leq \frac{\Delta}{8} \end{aligned}$$

2307 because we choose

$$\begin{aligned} 2308 & \beta \leq \min \left\{ \left(\frac{L\Delta\eta^2\sqrt{n}}{1152\sqrt{2}(1+b_1)\sigma\sqrt{T} \left(\sqrt{64L\Delta} + 3(B - \tau + b) \right)} \right)^{\frac{1-p}{4-2p}} \right. \\ 2309 & \quad \left. \left(\frac{L\Delta\eta^2\sqrt{n}}{3456\sqrt{2}(1+b_1)\sigma\sqrt{Ta}} \right)^{\frac{1-p}{4-2p}} \right\} \\ 2310 & \text{and } K+1 \leq T. \end{aligned} \tag{43}$$

2317 This implies

$$2318 \Pr \left(\frac{8\gamma\beta^2}{n\eta^2} (1 - \beta) \left\| \sum_{t=0}^K \sum_{i=1}^n \langle \zeta_{3,i}^t, \theta_i^{t+1} \rangle \right\| \geq \frac{\Delta}{8} \right) \leq \frac{\alpha}{14(T+1)}.$$

2320 Note that the worst dependency w.r.t. T is $\tilde{\mathcal{O}}(1/T^{\frac{3(1-p)}{4(2-p)}})$.

2322 BOUND OF THE TERM ⑥. The bound in this case is similar to the previous one. Let

$$2323 \sigma_6^2 := \frac{4\gamma^2}{n^2} \left(\sqrt{4L\Delta} + \frac{3}{2}(B - \tau) + 2b + (L\gamma)^p a \right)^2 \cdot \sigma^2.$$

2326 Then we have

$$2327 \mathbb{E} \left[\exp \left(\left| \frac{1}{\sigma_6^2} \frac{4\gamma^2}{n^2} (1 - \beta)^2 \langle \zeta_4^l, \theta_i^{l+1} \rangle^2 \right| \mid l, i - 1 \right) \right]$$

$$2328 \leq \mathbb{E} \left[\exp \left(\frac{1}{\sigma_6^2} \frac{4\gamma^2}{n^2} \|\zeta_4^l\|^2 \cdot \|\theta_i^{l+1}\|^2 \mid l, i - 1 \right) \right]$$

$$2329 \leq \mathbb{E} \left[\exp \left(\frac{1}{\sigma_6^2} \frac{4\gamma^2}{n^2} \left(\sqrt{4L\Delta} + \frac{3}{2}(B - \tau) + 2b + (L\gamma)^p a \right)^2 \cdot \|\theta_i^{l+1}\|^2 \mid l, i - 1 \right) \right]$$

$$2330 \leq \mathbb{E} \left[\exp \left(\left[\frac{4\gamma^2}{n^2} \left(\sqrt{4L\Delta} + \frac{3}{2}(B - \tau) + 2b + (L\gamma)^p a \right)^2 \cdot \sigma^2 \right]^{-1} \right. \right.$$

$$2331 \left. \frac{4\gamma^2}{n^2} \left(\sqrt{4L\Delta} + \frac{3}{2}(B - \tau) + 2b + (L\gamma)^p a \right)^2 \cdot \|\theta_i^{l+1}\|^2 \mid l, i - 1 \right)$$

$$2332 = \mathbb{E} \left[\exp \left(\frac{\|\theta_i^{l+1}\|^2}{\sigma^2} \mid l, i - 1 \right) \right] \leq \exp(1).$$

2344 Therefore, we have by Lemma 1 that

$$2345 \Pr \left[\frac{2\gamma(1 - \beta)}{n} \left\| \sum_{t=0}^K \sum_{i=1}^n \langle \zeta_{4,i}^t, \theta_i^{t+1} \rangle \right\| \right]$$

$$2346 \geq (\sqrt{2} + \sqrt{2}b_1) \sqrt{\sum_{t=0}^K \sum_{i=1}^n \frac{4\gamma^2}{n^2} \sigma^2 \cdot \left(\sqrt{4L\Delta} + \frac{3}{2}(B - \tau) + 2b + (L\gamma)^p a \right)^2}$$

$$2347 \leq \exp(-b_1^2/3) = \frac{\alpha}{14(T + 1)},$$

2353 Note that

$$2354 (\sqrt{2} + \sqrt{2}b_1) \sqrt{(K + 1)n} \cdot \frac{2\gamma}{n} \sigma \left(\sqrt{4L\Delta} + \frac{3}{2}(B - \tau) + 2b + (L\gamma)^p a \right)$$

$$2355 \leq (\sqrt{2} + \sqrt{2}b_1) \sqrt{(K + 1)n} \cdot \frac{2(\beta/3)^{\frac{1}{1-p}}}{Ln} \sigma \left(\sqrt{4L\Delta} + \frac{3}{2}(B - \tau) + 2b + (\beta/3)^{\frac{1}{1-p}} a \right)$$

$$2356 \leq \frac{\Delta}{8}$$

2362 because we choose

$$2363 \beta \leq \min \left\{ 3 \left(\frac{L\Delta\sqrt{n}}{32\sqrt{2}(1 + b_1)\sigma\sqrt{T} \left(\sqrt{4L\Delta} + \frac{3}{2}(B - \tau) + 2b \right)} \right)^{1-p} \right.$$

$$2364 \left. \left(\frac{L\Delta\sqrt{n}}{32\sqrt{2}(1 + b_1)\sigma\sqrt{T}a} \right)^{\frac{1-p}{1+p}} \right\},$$

$$2365 \text{ and } K + 1 \leq T. \tag{44}$$

2371 This implies

$$2372 \Pr \left(\frac{2\gamma(1 - \beta)}{n} \left\| \sum_{t=0}^K \sum_{i=1}^n \langle \zeta_{4,i}^t, \theta_i^{t+1} \rangle \right\| \geq \frac{\Delta}{8} \right) \leq \frac{\alpha}{14(T + 1)}.$$

2373 Note that the worst dependency w.r.t. T is $\tilde{O}(1/T^{\frac{3(1-p)}{2(1+p)}})$ since $a \sim T$.

2376 BOUND OF THE TERM ⑧. The bound in this case is similar to the previous one. Let

$$2377 \sigma_8^2 := \frac{4L^2\gamma^4}{n^2} \cdot \left(\sqrt{64L\Delta} + 3(B - \tau + b) + 3a \right)^2 \cdot \sigma^2.$$

2379 Then we have

$$2380 \mathbb{E} \left[\exp \left(\left| \frac{1}{\sigma_8^2} \frac{4\gamma^2}{n^2} (1 - \beta)^2 \langle \zeta_5^l, \theta_i^{l+1} \rangle^2 \right| \mid l, i - 1 \right) \right]$$

$$2381 \leq \mathbb{E} \left[\exp \left(\frac{1}{\sigma_8^2} \frac{4\gamma^2}{n^2} \|\zeta_5^l\|^2 \cdot \|\theta_i^{l+1}\|^2 \mid l, i - 1 \right) \right]$$

$$2382 \leq \mathbb{E} \left[\exp \left(\frac{1}{\sigma_8^2} \frac{4\gamma^2}{n^2} L^2 \gamma^2 \left(\sqrt{64L\Delta} + 3(B - \tau + b) + 3a \right) \cdot \|\theta_i^{l+1}\|^2 \mid l, i - 1 \right) \right].$$

2383 Since θ_i^{l+1} is sub-Gaussian with parameter σ^2 , then we can continue the chain of inequalities above

$$2384 \text{ using the definition of } \sigma_8^2$$

$$2385 \mathbb{E} \left[\exp \left(\left[\frac{4L^2\gamma^4}{n^2} \cdot \left(\sqrt{64L\Delta} + 3(B - \tau + b) + 3a \right)^2 \cdot \sigma^2 \right]^{-1} \right. \right.$$

$$2386 \left. \left. \frac{4L^2\gamma^4}{n^2} \cdot \left(\sqrt{64L\Delta} + 3(B - \tau + b) + 3a \right)^2 \cdot \|\theta_i^{l+1}\|^2 \right) \mid l, i - 1 \right]$$

$$2387 = \mathbb{E} \left[\exp \left(\frac{\|\theta_i^{l+1}\|^2}{\sigma^2} \right) \right] \leq \exp(1).$$

2388 Therefore, we have by Lemma 1 that

$$2389 \Pr \left[\frac{2\gamma(1 - \beta)}{n} \left\| \sum_{t=0}^K \sum_{i=1}^n \langle \zeta_{5,i}^t, \theta^{t+1} \rangle \right\| \right]$$

$$2390 \geq (\sqrt{2} + \sqrt{2}b_1) \sqrt{\sum_{t=0}^K \sum_{i=1}^n \frac{4L^2\gamma^4}{n^2} \sigma^2 \cdot \left(\sqrt{64L\Delta} + 3(B - \tau + b) + 3a \right)^2}$$

$$2391 \leq \exp(-b_1^2/3) = \frac{\alpha}{14(T + 1)}.$$

2392 Note that

$$2393 (\sqrt{2} + \sqrt{2}b_1) \sqrt{(K + 1)n} \cdot \frac{2L\gamma^2}{n} \sigma \left(\sqrt{64L\Delta} + 3(B - \tau + b) + 3a \right)$$

$$2394 \leq (\sqrt{2} + \sqrt{2}b_1) \sqrt{(K + 1)n} \cdot \frac{2(\beta/3)^{\frac{2}{1-p}}}{Ln} \sigma \left(\sqrt{64L\Delta} + 3(B - \tau + b) + 3a \right)$$

$$2395 \leq \frac{\Delta}{8}$$

2396 because we choose

$$2397 \beta \leq \min \left\{ \left(\frac{L\Delta\sqrt{n}}{32\sqrt{2}(1 + b_1)\sigma\sqrt{T} \left(\sqrt{64L\Delta} + 3(B - \tau + b) \right)} \right)^{\frac{1-p}{2}} \right.$$

$$2398 \left. \left(\frac{L\Delta\sqrt{n}}{96\sqrt{2}(1 + b_1)\sigma\sqrt{Ta}} \right)^{\frac{1-p}{2}} \right\}, \quad (45)$$

$$2399 \text{ and } K + 1 \leq T. \quad (46)$$

2400 This implies

$$2401 \Pr \left(2\gamma(1 - \beta) \left\| \sum_{t=0}^K \sum_{i=1}^n \langle \zeta_{5,i}^t, \theta^{t+1} \rangle \right\| \geq \frac{\Delta}{8} \right) \leq \frac{\alpha}{14(T + 1)}.$$

2402 Note that the worst dependency w.r.t T is $\tilde{O}(1/T^{\frac{3(1-p)}{4}})$.

Final probability. Therefore, the probability event

$$\Omega := E_K \cap \bar{\Theta}^{K+1} \cap \left(\bigcap_{i=1}^n \bar{\Theta}_i^{K+1} \right) \cap \bar{N}^{K+1} \cap E_{\textcircled{1}} \cap E_{\textcircled{2}} \cap E_{\textcircled{3}} \cap E_{\textcircled{4}} \cap E_{\textcircled{5}} \cap E_{\textcircled{6}} \cap E_{\textcircled{7}} \cap E_{\textcircled{8}},$$

where each $E_{\textcircled{1}}-E_{\textcircled{8}}$ denotes that each of 1-8-th terms is smaller than $\frac{\Delta}{8}$ implies that

$$\textcircled{1} + \textcircled{2} + \textcircled{3} + \textcircled{4} + \textcircled{5} + \textcircled{6} + \textcircled{7} + \textcircled{8} \leq 8 \cdot \frac{\Phi^0}{8} = \Delta,$$

i.e. 7 in the induction assumption holds. Moreover, this also implies that

$$\Phi^{K+1} \leq \Phi^0 + \Delta \leq \Delta + \Delta = 2\Delta,$$

i.e. 6 in the induction assumption holds. The probability $\Pr(E_{K+1})$ can be lower bounded as follows

$$\begin{aligned} \Pr(E_{K+1}) &\geq \Pr(\Omega) \\ &= \Pr\left(E_K \cap \bar{\Theta}^{K+1} \cap \left(\bigcap_{i=1}^n \bar{\Theta}_i^{K+1}\right) \cap \bar{N}^{K+1} \cap E_{\textcircled{1}} \cap E_{\textcircled{2}} \cap E_{\textcircled{3}} \cap E_{\textcircled{4}} \cap E_{\textcircled{5}} \cap E_{\textcircled{6}} \right. \\ &\quad \left. \cap E_{\textcircled{7}} \cap E_{\textcircled{8}}\right) \\ &= 1 - \Pr\left(\bar{E}_K \cup \Theta^{K+1} \cup \left(\bigcup_{i=1}^n \Theta_i^{K+1}\right) \cup N^{K+1} \cup \bar{E}_{\textcircled{1}} \cup \bar{E}_{\textcircled{2}} \cup \bar{E}_{\textcircled{3}} \cup \bar{E}_{\textcircled{4}} \cup \bar{E}_{\textcircled{5}} \cup \bar{E}_{\textcircled{6}} \right. \\ &\quad \left. \cup \bar{E}_{\textcircled{7}} \cup \bar{E}_{\textcircled{8}}\right) \\ &\geq 1 - \Pr(\bar{E}_K) - \Pr(\Theta^{K+1}) - \sum_{i=1}^n \Pr(\Theta_i^{K+1}) - \Pr(N^{K+1}) - \Pr(\bar{E}_{\textcircled{1}}) - \Pr(\bar{E}_{\textcircled{2}}) \\ &\quad - \Pr(\bar{E}_{\textcircled{3}}) - \Pr(\bar{E}_{\textcircled{4}}) - \Pr(\bar{E}_{\textcircled{5}}) - \Pr(\bar{E}_{\textcircled{6}}) - \Pr(\bar{E}_{\textcircled{7}}) - \Pr(\bar{E}_{\textcircled{8}}) \\ &\geq 1 - \frac{\alpha(K+1)}{T+1} - \frac{\alpha}{6(T+1)} - \sum_{i=1}^n \frac{\alpha}{6n(T+1)} - \frac{\alpha}{6(T+1)} - 0 - 7 \cdot \frac{\alpha}{14(T+1)} \\ &= 1 - \frac{\alpha(K+2)}{T+1}. \end{aligned}$$

This finalizes the transition step of induction. The result of the theorem follows by setting $K = T - 1$. Indeed, from (37) we obtain

$$\frac{\gamma}{2} \sum_{t=0}^K \|\nabla f(x^t)\|^2 \leq \Phi^0 - \Phi^{K+1} + \Delta \leq 2\Delta \Rightarrow \frac{1}{T} \sum_{t=0}^{T-1} \|\nabla f(x^t)\|^2 \leq \frac{4\Delta}{\gamma T}. \quad (47)$$

Final rate. Now we have the following restrictions on the momentum parameter in terms of dependency on T from each bound of terms 1-8 correspondingly

$$\begin{aligned} \beta \leq \tilde{\mathcal{O}} &\left(\underbrace{\left(\frac{L\Delta n}{T\sigma^2}\right)^{\frac{1-p}{2-p}}}_{\text{from term 1}}, \underbrace{\left(\frac{L\Delta\sqrt{n}\eta}{B\sigma\sqrt{T}}\right)^{\frac{1-p}{2-p}}}_{\text{from term 2}}, \underbrace{\left(\frac{L\Delta\sqrt{n}\eta}{\sqrt{T}\sigma a}\right)^{\frac{1-p}{3-p}}}_{\text{from term 3}}, \underbrace{\left(\frac{L\Delta\eta\sqrt{n}}{\sigma\sqrt{T}a}\right)^{\frac{1-p}{2(2-p)}}}_{\text{from term 4}}, \right. \\ &\left. \underbrace{\left(\frac{L\Delta\eta^2\sqrt{n}}{\sigma\sqrt{T}a}\right)^{\frac{1-p}{3-p}}}_{\text{from term 5}}, \underbrace{\left(\frac{L\Delta\eta^2\sqrt{n}}{\sigma\sqrt{T}a}\right)^{\frac{1-p}{2(2-p)}}}_{\text{from term 7}}, \underbrace{\left(\frac{L\Delta\sqrt{n}}{\sigma\sqrt{T}a}\right)^{\frac{1-p}{1+p}}}_{\text{from term 6}}, \underbrace{\left(\frac{L\Delta\sqrt{n}}{\sigma\sqrt{T}a}\right)^{\frac{1-p}{2}}}_{\text{from term 8}} \right). \end{aligned}$$

Now we need to understand which stepsize restrictions give the worst T complexity. We have

$$\begin{aligned} \gamma \leq \frac{1}{L} \tilde{\mathcal{O}} &\left(\underbrace{\left(\frac{L\Delta n}{T\sigma^2}\right)^{\frac{1}{2-p}}}_{\text{from term 1}}, \underbrace{\left(\frac{L\Delta\sqrt{n}\eta}{B\sigma\sqrt{T}}\right)^{\frac{1}{2-p}}}_{\text{from term 2}}, \underbrace{\left(\frac{L\Delta\sqrt{n}\eta}{\sqrt{T}\sigma a}\right)^{\frac{1}{3-p}}}_{\text{from term 3}}, \underbrace{\left(\frac{L\Delta\eta\sqrt{n}}{\sigma\sqrt{T}a}\right)^{\frac{1}{2(2-p)}}}_{\text{from term 4}}, \right. \\ &\left. \underbrace{\left(\frac{L\Delta\eta^2\sqrt{n}}{\sigma\sqrt{T}a}\right)^{\frac{1}{3-p}}}_{\text{from term 5}}, \underbrace{\left(\frac{L\Delta\eta^2\sqrt{n}}{\sigma\sqrt{T}a}\right)^{\frac{1}{2(2-p)}}}_{\text{from term 7}}, \underbrace{\left(\frac{L\Delta\sqrt{n}}{\sigma\sqrt{T}a}\right)^{\frac{1}{1+p}}}_{\text{from term 6}}, \underbrace{\left(\frac{L\Delta\sqrt{n}}{\sigma\sqrt{T}a}\right)^{\frac{1}{2}}}_{\text{from term 8}} \right). \quad (48) \end{aligned}$$

By (28) we get that

$$\gamma \leq \frac{1}{L} \tilde{\mathcal{O}} \left(\underbrace{\left(\frac{L\Delta n}{T\sigma^2} \right)^{\frac{1}{2-p}}}_{\text{from term 1}}, \underbrace{\left(\frac{L\Delta\sqrt{n}\eta}{B\sigma\sqrt{T}} \right)^{\frac{1}{2-p}}}_{\text{from term 2}}, \underbrace{\left(\frac{L\Delta n\eta}{T\sigma\sqrt{d}\sigma_\omega} \right)^{\frac{1}{3-p}}}_{\text{from term 3}}, \underbrace{\left(\frac{L\Delta\eta n}{\sigma T\sqrt{d}\sigma_\omega} \right)^{\frac{1}{2(2-p)}}}_{\text{from term 4}}, \right. \\ \left. \underbrace{\left(\frac{L\Delta\eta^2 n}{\sigma T\sqrt{d}\sigma_\omega} \right)^{\frac{1}{3-p}}}_{\text{from term 5}}, \underbrace{\left(\frac{L\Delta\eta^2 n}{\sigma T\sqrt{d}\sigma_\omega} \right)^{\frac{1}{2(2-p)}}}_{\text{from term 7}}, \underbrace{\left(\frac{L\Delta n}{\sigma T\sqrt{d}\sigma_\omega} \right)^{\frac{1}{1+p}}}_{\text{from term 6}}, \underbrace{\left(\frac{L\Delta n}{\sigma T\sqrt{d}\sigma_\omega} \right)^{\frac{1}{2}}}_{\text{from term 8}} \right). \quad (49)$$

If we choose $p = 0.8$, then the worst power of T comes from the term ① and equals to $\frac{1}{6}$. The second worst comes from the term ⑥ and equals to $\frac{4}{9}$. These two terms give the rate of the form

$$\tilde{\mathcal{O}} \left(\frac{L\Delta}{T} \left(\frac{T\sigma^2}{L\Delta n} \right)^{\frac{1}{2-p}} + \frac{L\Delta}{T} \left(\frac{\sigma T\sqrt{d}\sigma_\omega}{L\Delta n} \right)^{\frac{1}{1+p}} \right) \\ = \tilde{\mathcal{O}} \left(\frac{(L\Delta)^{\frac{1-p}{2-p}} \sigma^{\frac{2}{2-p}}}{T^{\frac{1-p}{2-p}} n^{\frac{1}{2-p}}} + \frac{(L\Delta)^{\frac{p}{1+p}} \sigma^{\frac{1}{1+p}} d^{\frac{1}{2(1+p)}} \sigma_\omega^{\frac{1}{1+p}}}{T^{\frac{p}{1+p}} n^{\frac{1}{1+p}}} \right) \\ \stackrel{p=0.8}{=} \tilde{\mathcal{O}} \left(\frac{(L\Delta)^{\frac{1}{6}} \sigma^{\frac{5}{3}}}{T^{\frac{1}{6}} n^{\frac{5}{6}}} + \frac{(L\Delta)^{\frac{4}{9}} \sigma^{\frac{5}{9}} d^{\frac{5}{18}} \sigma_\omega^{\frac{5}{9}}}{T^{\frac{4}{9}} n^{\frac{5}{9}}} \right).$$

Besides, we have the momentum restriction of the form $\beta \leq \left(\frac{\tau}{48a} \right)^{1-p}$ that translates to

$$\gamma \leq \tilde{\mathcal{O}} \left(\frac{\tau}{La} \right),$$

and therefore, gives an additional term in the rate of the form

$$\tilde{\mathcal{O}} \left(\frac{L\Delta}{T} \frac{\sqrt{d}\sqrt{T/n}\sigma_\omega}{\tau} \right) = \tilde{\mathcal{O}} \left(\frac{L\Delta\sqrt{d}\sigma_\omega}{\sqrt{Tn}\tau} \right).$$

To conclude, we obtain with probability at least $1 - \alpha$ that

$$\frac{1}{T} \sum_{t=0}^{T-1} \|\nabla f(x^t)\|^2 \leq \tilde{\mathcal{O}} \left(\frac{L\Delta\sqrt{d}\sigma_\omega}{\sqrt{Tn}\tau} + \frac{(L\Delta)^{\frac{1}{6}} \sigma^{\frac{5}{3}}}{T^{\frac{1}{6}} n^{\frac{5}{6}}} + \frac{(L\Delta)^{\frac{4}{9}} \sigma^{\frac{5}{9}} d^{\frac{5}{18}} \sigma_\omega^{\frac{5}{9}}}{T^{\frac{4}{9}} n^{\frac{5}{9}}} \right).$$

CASE $\mathcal{I}_{K+1} = 0$. This case is even easier. The only change will be with the term next to R^t . We will get

$$1 - \frac{48L^2}{\eta^2} \gamma^2 - \frac{12L^2}{\beta^2} \gamma^2 \geq \frac{2}{3} - \frac{48L^2}{\eta} \gamma^2 \geq 0$$

instead of

$$1 - \frac{16\beta^2 L^2}{\eta^2} \gamma^2 - \frac{48L^2}{\eta^2} \gamma^2 - \frac{12L^2}{\beta^2} \gamma^2 \geq 0$$

as in the previous case. This difference comes from Lemma 18 because $\tilde{V}^{K+1} = 0$. The rest is a repetition of the previous derivations. \square

F PROOF OF COROLLARY 1

Corollary 1. Let Assumptions 1 and 2 hold and $\alpha \in (0, 1)$. Let $\Delta \geq \Phi^0$ and σ_ω be chosen as $\sigma_\omega = \Theta\left(\frac{\tau}{\varepsilon}\sqrt{T \log \frac{1}{\delta}}\right)$. Then there exists a stepsize γ and momentum parameter β such that the iterates of Clip21-SGDM (Algorithm 3) with probability at least $1 - \alpha$ satisfy local (ε, δ) -DP and

$$\frac{1}{T} \sum_{t=0}^{T-1} \|\nabla f(x^t)\|^2 \leq \tilde{\mathcal{O}}\left(\frac{L\Delta\sqrt{d}}{\sqrt{n\varepsilon}}\right), \quad (14)$$

where $\tilde{\mathcal{O}}$ hides constant and logarithmic factors, and terms decreasing in T .

Proof. We need to plug in the value of σ_ω inside (13). Indeed, we have that

$$\begin{aligned} & \tilde{\mathcal{O}}\left(\frac{L\Delta\sqrt{d}}{\sqrt{Tn\tau}}\frac{\tau}{\varepsilon}\sqrt{T} + \frac{(L\Delta)^{1/6}\sigma^{5/3}}{T^{1/6}n^{5/6}} + \frac{(L\Delta)^{4/9}\sigma^{5/9}d^{5/18}}{T^{4/9}n^{5/9}}\frac{\tau}{\varepsilon}\sqrt{T}\right) \\ &= \tilde{\mathcal{O}}\left(\frac{L\Delta\sqrt{d}}{\sqrt{Tn\tau}}\frac{\tau}{\varepsilon}\sqrt{T} + \frac{(L\Delta)^{1/6}\sigma^{5/3}}{T^{1/6}n^{5/6}} + \frac{(L\Delta)^{4/9}\sigma^{5/9}d^{5/18}}{T^{4/9}n^{5/9}}\left(\frac{\tau}{\varepsilon}\sqrt{T}\right)^{5/9}\right) \\ &= \tilde{\mathcal{O}}\left(\frac{L\Delta\sqrt{d}}{\sqrt{n\varepsilon}} + \frac{(L\Delta)^{1/6}\sigma^{5/3}}{T^{1/6}n^{5/6}} + \frac{(L\Delta)^{4/9}\sigma^{5/9}d^{5/18}\tau^{5/9}}{T^{1/6}n^{5/9}\varepsilon^{5/9}}\right) / \end{aligned}$$

Leaving only the terms that do not improve with T we get the result. \square

G PROOF OF THEOREM 3

We highlight that the proof of Theorem 3 mainly follows that of Theorem 4. The main difference comes from the fact that stepsize and momentum restrictions become less demanding as in purely stochastic setting (without DP noise) $a = 0$. In particular, we can choose $p = 1$. Therefore, we only list the modified lemmas without the proofs.

Lemma 19. Let each f_i be L -smooth. Then we have the following inequality with probability 1

$$\|v_i^{t+1} - g_i^t\| \leq \max\{0, \|v_i^t - g_i^{t-1}\| - \tau\} + \beta L\gamma\|g^t\| + \beta\|\nabla f_i(x^t) - v_i^t\| + \beta\|\theta_i^{t+1}\|. \quad (50)$$

Lemma 20. Let each f_i be L -smooth and $\Delta \geq \Phi^0$. Assume that the following inequalities hold

1. $g^0 = \frac{1}{n} \sum_{i=1}^n g_i^0$;
2. $\|g^{t-1}\| \leq \sqrt{64L\Delta} + 3(B - \tau) + 3b$;
3. $\|\nabla f_i(x^{t-1}) - v_i^{t-1}\| \leq \sqrt{4L\Delta} + \frac{3}{2}(B - \tau) + 2b$ for all $i \in [n]$;
4. $\|v_i^t - g_i^{t-1}\| \leq B$ for all $i \in [n]$;
5. $\gamma \leq \frac{1}{12L}$;
6. $\|\theta_i^t\| \leq b$ for all $i \in [n]$;
7. $1 \geq \beta \geq 4L\gamma$;
8. $\Phi^{t-1} \leq 2\Delta$.

Then we have

$$\|g^t\| \leq \sqrt{64L\Delta} + 3(B - \tau) + 3b. \quad (51)$$

Lemma 21. Let each f_i is L -smooth and $\Delta \geq \Phi^0$. Assume the following inequalities hold

1. $\gamma \leq \frac{1}{12L}$;

- 2592 2. $4L\gamma \leq \beta \leq 1$;
 2593
 2594 3. $\|\nabla f_i(x^{t-1}) - v_i^{t-1}\| \leq \sqrt{4L\Delta} + \frac{3}{2}(B - \tau) + 2b$;
 2595
 2596 4. $\|\theta_i^t\| \leq b$;
 2597
 2598 5. $\|g^{t-1}\| \leq \sqrt{64L\Delta} + 3(B - \tau) + 3b$.

2599 Then we have

$$2600 \|\nabla f_i(x^t) - v_i^t\| \leq \sqrt{4L\Delta} + \frac{3}{2}(B - \tau) + 2b. \quad (52)$$

2601 **Lemma 22.** Let each f_i be L -smooth, $\Delta \geq \Phi^0$, and $i \in \mathcal{I}_t$. Let the following inequalities hold

- 2602 1. $12L\gamma \leq 1$;
 2603
 2604 2. $1 \geq \beta \geq 4L\gamma$;
 2605
 2606 3. $\beta \leq \frac{\tau}{32\sqrt{L\Delta}}$;
 2607
 2608 4. $\beta \leq \frac{\tau}{18(B-\tau)}$;
 2609
 2610 5. $\beta \leq \frac{\tau}{30b}$;
 2611
 2612 6. $\|g^t\| \leq \sqrt{64L\Delta} + 3(B - \tau) + 3b$;
 2613
 2614 7. $\|\theta_i^{t+1}\| \leq b$;
 2615
 2616 8. $\|\nabla f_i(x^t) - v_i^t\| \leq \sqrt{4L\Delta} + \frac{3}{2}(B - \tau) + 2b$.

2617 Then

$$2618 \|v_i^{t+1} - g_i^t\| \leq \|v_i^t - g_i^{t-1}\| - \frac{\tau}{2}. \quad (53)$$

2619 **Lemma 23.** Let $\|\theta_i^{t+1}\| \leq b$ for all $i \in [n]$. Let each f_i be L -smooth. Then \tilde{P}^t decreases as

$$2620 \tilde{P}^{t+1} \leq (1 - \beta)\tilde{P}^t + \frac{3L^2}{\beta}R^t + \beta^2b^2 + \frac{2}{n}\beta(1 - \beta) \sum_{i=1}^n \langle v_i^t - \nabla f_i(x^{t+1}), \theta_i^{t+1} \rangle. \quad (54)$$

2621 Similarly, we can get the descent of P^t .

2622 **Lemma 24.** Let $\|\theta^{t+1}\| \leq \frac{c}{\sqrt{n}}$, and each f_i be L -smooth. Then P^t decreases as

$$2623 P^{t+1} \leq (1 - \beta)P^t + \frac{3L^2}{\beta}R^t + \beta^2b^2 + 2\beta(1 - \beta) \langle v^t - \nabla f(x^{t+1}), \theta^{t+1} \rangle.$$

2624 Now we present the descent of \tilde{V}^t .

2625 **Lemma 25.** Let $\|\theta_i^t\| \leq b$ for all $i \in [n]$, each f_i be L -smooth, and $\|v_i^t - g_i^{t-1}\| \leq B$ for all $i \in [n]$. Then

$$2626 \begin{aligned} 2627 \|g_i^t - v_i^t\|^2 &\leq (1 - \eta)\|g_i^{t-1} - v_i^{t-1}\|^2 + \frac{4\beta^2}{\eta}\|v_i^{t-1} - \nabla f_i(x^{t-1})\|^2 + \frac{4\beta^2L^2}{\eta}R^{t-1} + \beta^2b^2 \quad (55) \\ 2628 &+ 2(1 - \eta)^2\beta \langle (g_i^{t-1} - v_i^{t-1}) + \beta(v_i^{t-1} - \nabla f_i(x^{t-1})), \theta_i^t \rangle \\ 2629 &+ 2(1 - \eta)^2\beta \langle \beta(\nabla f_i(x^{t-1}) - \nabla f_i(x^t)), \theta_i^t \rangle. \end{aligned}$$

2630 Moreover, averaging the inequalities across all $i \in [n]$ we get

$$2631 \begin{aligned} 2632 \tilde{V}^t &\leq (1 - \eta)\tilde{V}^{t-1} + \frac{4\beta^2}{\eta}\tilde{P}^{t-1} + \frac{4\beta^2L^2}{\eta}R^{t-1} + \beta^2b^2 \quad (56) \\ 2633 &+ \frac{2}{n}(1 - \eta)^2\beta \sum_{i=1}^n \langle (g_i^{t-1} - v_i^{t-1}) + \beta(v_i^{t-1} - \nabla f_i(x^{t-1})) + \beta(\nabla f_i(x^{t-1}) - \nabla f_i(x^t)), \theta_i^t \rangle. \end{aligned}$$

Theorem 7 (Full statement of Theorem 3). *Let Assumptions 1 and 2 hold, $B := \max_i \{\|\nabla f_i(x^0)\|\} + b > \tau$, probability constant $\alpha \in (0, 1)$, and $\Delta \geq \Phi^0$. Let us run Algorithm 3 for T iterations. Assume the following inequalities hold*

1. **stepsize restrictions:**

- i) $12L\gamma \leq 1$;
- ii)

$$\frac{2}{3} - \frac{16\beta^2 L^2}{\eta^2} \gamma^2 - \frac{48L^2}{\eta^2} \gamma^2 \geq 0;$$

2. **momentum restrictions:**

- i) $1 \geq \beta \geq 4L\gamma$;
- ii) $\beta \leq \frac{\tau}{32\sqrt{L\Delta}}$;
- iii) $\beta \leq \frac{\tau}{18(B-\tau)}$;
- iv) $\beta \leq \frac{\tau}{30b}$;
- v) *and momentum restrictions defined in (59), (60), (61), (62), (63), (65), (64), and (66);*

Then with probability $1 - \alpha$ we have

$$\frac{1}{T} \sum_{t=0}^{T-1} \|\nabla f(x^t)\|^2 \leq \tilde{\mathcal{O}} \left(\frac{\sigma(\sqrt{L\Delta} + B + \sigma)}{\sqrt{Tn}} \right),$$

where $\tilde{\mathcal{O}}$ hides constant and logarithmic factors, and higher order terms decreasing in T .

Proof. We prove the main theorem by induction. The conventional choice $\nabla f_i(x^{-1}, \xi_i^{-1}) = v_i^{-1} = g_i^{-1} = 0, \Phi^{-1} = \Phi^0$.

Let us define an event E_t for each $t \in \{0, \dots, T\}$ such that the following inequalities hold for all $k \in \{0, \dots, t\}$

1. $\|v_i^k - g_i^{k-1}\| \leq B$ for $i \in \mathcal{I}_k$;
2. $\|g^k\| \leq \sqrt{64L\Delta} + 3(B - \tau) + 3b$;
3. $\|v_i^k - \nabla f_i(x^k)\| \leq \sqrt{4L\Delta} + \frac{3}{2}(B - \tau) + 2b$;
4. $\|\theta_i^k\| \leq b$ for all $i \in [n]$ and $\|\theta^k\| \leq \frac{c}{\sqrt{n}}$;
5. $\Phi^k \leq 2\Delta$;
- 6.

$$\begin{aligned} \Delta &\geq \frac{2\gamma\beta}{n\eta} (1 - \eta)^2 \sum_{l=0}^{k-1} \sum_{i=1}^n \langle (g_i^l - v_i^l) + \beta(v_i^l - \nabla f_i(x^l)) + \beta(\nabla f_i(x^l) - \nabla f_i(x^{l+1})), \theta_i^l \rangle \\ &\quad + \frac{8\gamma\beta^2}{n\eta^2} (1 - \beta) \sum_{l=0}^{k-1} \sum_{i=1}^n \langle v_i^l - \nabla f_i(x^l), \theta_i^{l+1} \rangle + 2\gamma(1 - \beta) \sum_{l=0}^{k-1} \langle v^l - \nabla f(x^l), \theta^{l+1} \rangle \\ &\quad + \frac{8\gamma\beta^2}{n\eta^2} (1 - \beta) \sum_{l=0}^{k-1} \sum_{i=1}^n \langle \nabla f_i(x^l) - \nabla f_i(x^{l+1}), \theta_i^{l+1} \rangle \\ &\quad + 2\gamma(1 - \beta) \sum_{l=0}^{k-1} \langle \nabla f(x^l) - \nabla f(x^{l+1}), \theta^{l+1} \rangle. \end{aligned}$$

2700 Denote the events Θ_i^t and Θ^t as

$$2701 \Theta_i^t := \{\|\theta_i^t\| \geq b\}, \quad \text{and} \quad \Theta^t := \{\|\theta^t\| \geq \frac{c}{\sqrt{n}}\} \quad (57)$$

2702 respectively. From Assumption 2 we have

$$2703 \Pr(\Theta_i^t) \leq 2 \exp\left(-\frac{b^2}{2\sigma^2}\right) = \frac{\alpha}{6(T+1)n}$$

2704 where the last equality is by definition of b^2 . Therefore, $\Pr(\bar{\Theta}_i^t) \geq 1 - \frac{\alpha}{6(T+1)n}$.

2705 Besides, notice that the constant c in (28) can be viewed as

$$2706 c = (\sqrt{2} + 2b_3)\sigma \quad \text{where} \quad b_3^2 = 3 \log \frac{6(T+1)}{\alpha}.$$

2707 Now we can use Lemma 1 to bound $\Pr(\Theta^t)$. Since all θ_i^t are independent σ -sub-Gaussian random

2708 vectors, then we have

$$2709 \Pr\left(\left\|\sum_{i=1}^n \theta_i^t\right\| \geq c\sqrt{n}\right) = \Pr\left(\|\theta^t\| \geq \frac{c}{\sqrt{n}}\right) \leq \exp(-b_3^2/3) = \frac{\alpha}{6(T+1)}.$$

2710 Now we prove that $\Pr(E_t) \geq 1 - \frac{\alpha(t+1)}{T+1}$ for all $t \in \{0, \dots, T-1\}$. First, we show that the base of

2711 induction holds.

2712 Base of induction.

- 2713 1. $\|v_i^0 - g_i^{-1}\| = \|v_i^0\| = \beta \|\nabla f_i(x^0, \xi_i^0)\| = \beta \|\theta_i^0\| + \beta \|\nabla f_i(x^0)\| \leq \frac{1}{2}b + \frac{1}{2}B \leq \frac{1}{2}B + \frac{1}{2}B =$
 2714 B holds with probability $1 - \frac{\alpha}{6(T+1)}$. Indeed, we have

$$2715 \Pr(\Theta_i^0) \leq 2 \exp\left(-\frac{b^2}{2\sigma^2}\right) = \frac{\alpha}{6(T+1)n}.$$

2716 Therefore, we have

$$2717 \Pr\left(\bigcap_{i=1}^n \bar{\Theta}_i^0\right) = 1 - \Pr\left(\bigcup_{i=1}^n \Theta_i^0\right) \geq 1 - \sum_{i=1}^n \Pr(\Theta_i^0) = 1 - n \frac{\alpha}{6(T+1)n} = 1 - \frac{\alpha}{6(T+1)}.$$

2718 Moreover, by concentration inequality we have

$$2719 \Pr(\Theta^0) \leq \frac{\alpha}{6(T+1)}.$$

2720 This means that the probability of the event that each $\|\theta_i^0\| \leq b$, and $\|\theta^0\| \leq \frac{c}{\sqrt{n}}$, and is at

2721 least

$$2722 1 - n \frac{\alpha}{6n(T+1)} - \frac{\alpha}{6(T+1)} = 1 - \frac{\alpha}{3(T+1)}.$$

- 2723 2. $g^0 = \frac{1}{n} \sum_{i=1}^n (g_i^{-1} + \text{clip}_\tau(v_i^0 - g_i^{-1})) = \frac{1}{n} \sum_{i=1}^n \text{clip}_\tau(\beta \nabla f_i(x^0, \xi_i^0))$. Therefore, we have

$$2724 \|g^0\| \leq \left\| \frac{1}{n} \sum_{i=1}^n \beta \nabla f_i(x^0) + \beta \theta_i^0 + (\text{clip}_\tau(\beta \nabla f_i(x^0, \xi_i^0)) - \beta \nabla f_i(x^0, \xi_i^0)) \right\|$$

$$2725 \leq \beta \|\nabla f(x^0)\| + \frac{\beta}{n} \sum_{i=1}^n \|\theta_i^0\| + \frac{1}{n} \sum_{i=1}^n \max\{0, \beta \|\nabla f_i(x^0, \xi_i^0)\| - \tau\}$$

$$2726 \leq \beta \sqrt{2L(f(x^0) - f(x^*))} + \frac{\beta}{n} \sum_{i=1}^n \|\theta_i^0\| + \frac{1}{n} \sum_{i=1}^n \max\{0, \beta \|\nabla f_i(x^0)\| + \beta \|\theta_i^0\| - \tau\}$$

$$2727 \leq \frac{1}{2} \sqrt{2L\Phi^0} + \frac{2\beta}{n} \sum_{i=1}^n \|\theta_i^0\| + \frac{\beta}{n} \sum_{i=1}^n \|\nabla f_i(x^0)\| - \tau$$

$$2728 \leq \sqrt{64L\Delta} + 2\beta b + \beta B - \tau$$

$$2729 \leq \sqrt{64L\Delta} + \frac{3}{2}B - \tau + b \leq \sqrt{64L\Delta} + 3(B - \tau) + 2b.$$

2754 The inequalities above again hold in $\cap_{i=1}^n \bar{\Theta}_i^0$, i.e. with probability at least $1 - \frac{\alpha}{6(T+1)}$.
 2755

2756 3. We have

$$2757 \quad \|v_i^0 - \nabla f_i(x^0)\| = \|\nabla f_i(x^0, \xi_i^0) - \nabla f_i(x^0)\| \leq b.$$

2759 This bound holds with probability at least $1 - \frac{\alpha}{6(T+1)}$ because it holds in $\cap_{i=1}^n \bar{\Theta}_i^0$.
 2760

2761 4. Inequalities 5 obviously also hold, as $\Phi^0 \leq 2\Phi^0 \leq 2\Delta$ by the choice of Δ .
 2762

2763 Therefore, we conclude that the inequalities 1-7 hold with a probability at least
 2764

$$2765 \quad \Pr\left(\Theta^0 \cap \left(\cap_{i=1}^n \bar{\Theta}_i^0\right) \cap \bar{N}^t\right) \geq 1 - \Pr(\Theta^0) - \sum_{i=1}^n \Pr(\Theta_i^0) - \Pr(N^0)$$

$$2766 \quad \geq 1 - \frac{\alpha}{6(T+1)} - n \cdot \frac{\alpha}{6n(T+1)} - \frac{\alpha}{6(T+1)}$$

$$2767 \quad = 1 - \frac{\alpha}{2(T+1)} > 1 - \frac{\alpha}{T+1},$$

2770 i.e. $\Pr(E_0) \geq 1 - \frac{\alpha}{T+1}$ holds. This is the base of the induction.
 2771

2772 Transition step of induction.

2773 CASE $|\mathcal{I}_{K+1}| > 0$. Assume that all events $\bar{\Theta}^{K+1}$ and $\bar{\Theta}_i^{K+1}$ take place, i.e. $\|\theta_i^{K+1}\| \leq$
 2774 $b, \|\theta^{K+1}\| \leq \frac{c}{\sqrt{n}}$ for all $i \in [n]$. For that we need to work in $\bar{\Theta}^{K+1} \cap \left(\cap_{i=1}^n \bar{\Theta}_i^{K+1}\right)$. Then,
 2775 by the assumptions of the induction, from Lemma 22 we get for all $i \in \mathcal{I}_{K+1}$

$$2776 \quad \|v_i^{K+1} - g_i^K\| \leq \|v_i^K - g_i^{K-1}\| - \frac{\tau}{2} \leq B - \frac{\tau}{2}.$$

2777 Therefore, from Lemma 20 we get that

$$2778 \quad \|g^{K+1}\| \leq \sqrt{64L\Delta} + 3(B - \tau) + 3b,$$

2779 and from Lemma 21

$$2780 \quad \|\nabla f_i(x^{K+1}) - v_i^{K+1}\| \leq \sqrt{4L\Delta} + \frac{3}{2}(B - \tau) + 2b.$$

2781 This means that 1-3 in the induction assumption are also verified for the step $K + 1$.
 2782

2783 Since we have for all $t \in \{0, \dots, K + 1\}$ that inequalities 1-3 are verified, then we can write for
 2784 each $t \in \{0, K\}$ by Lemmas 2 and 23 to 25 the following
 2785

$$2786 \quad \Phi^{t+1} = \delta^{t+1} + \frac{\gamma}{\eta} \tilde{V}^{t+1} + \frac{4\gamma\beta}{\eta^2} \tilde{P}^{t+1} + \frac{\gamma}{\beta} P^{t+1}$$

$$2787 \quad \leq \delta^t - \frac{\gamma}{2} \|\nabla f(x^t)\|^2 - \frac{1}{4\gamma} R^t + \gamma \tilde{V}^t + \gamma P^t$$

$$2788 \quad + \frac{\gamma}{\eta} \left((1 - \eta) \tilde{V}^t + \frac{4\beta^2}{\eta} \tilde{P}^t + \frac{4\beta^2 L^2}{\eta} R^t + \beta^2 b^2 \right)$$

$$2789 \quad + \frac{2}{n} \beta (1 - \eta^2) \sum_{i=1}^n \left(\langle (g_i^t - v_i^t) + \beta(v_i^t - \nabla f_i(x^t)) + \beta(\nabla f_i(x^t) - \nabla f_i(x^{t+1})), \theta_i^{t+1} \rangle \right)$$

$$2790 \quad + \frac{4\gamma\beta}{\eta^2} \left((1 - \beta) \tilde{P}^t + \frac{3L^2}{\beta} R^t + \beta^2 b^2 + \frac{2}{n} \beta (1 - \beta) \sum_{i=1}^n \langle v_i^t - \nabla f_i(x^{t+1}), \theta_i^{t+1} \rangle \right)$$

$$2791 \quad + \frac{\gamma}{\beta} \left((1 - \beta) P^t + \frac{3L^2}{\beta} R^t + \beta^2 \frac{c^2}{n} + 2\beta(1 - \beta) \langle v^t - \nabla f(x^{t+1}), \theta^{t+1} \rangle \right)$$

2808 Rearranging terms we get

$$\begin{aligned}
2811 \quad \Phi^{t+1} &\leq \delta^t - \frac{\gamma}{2} \|\nabla f(x^t)\|^2 + \frac{\gamma}{\eta} \tilde{V}^t (\eta + 1 - \eta) + \frac{4\gamma\beta}{\eta^2} \tilde{P}^t (\beta + 1 - \beta) + \frac{\gamma}{\beta} P^t (\beta + 1 - \beta) \\
2812 & - \frac{1}{4\gamma} R^t \left(1 - \frac{16L^2\beta^2}{\eta^2} \gamma^2 - \frac{48L^2}{\eta^2} \gamma^2 - \frac{12L^2}{\beta^2} \gamma^2 \right) + b^2 \left(\frac{\beta^2\gamma}{\eta} + \frac{4\gamma\beta^3}{\eta^2} \right) + c^2 \frac{\gamma\beta}{n} \\
2813 & + \frac{2\gamma\beta}{n\eta} (1 - \eta)^2 \sum_{i=1}^n \langle (g_i^t - v_i^t) + \beta(v_i^t - \nabla f_i(x^t)) + \beta(\nabla f_i(x^t) - \nabla f_i(x^{t+1})), \theta_i^{t+1} \rangle \\
2814 & + \frac{8\gamma\beta^2}{n\eta^2} (1 - \beta) \sum_{i=1}^n \langle v_i^t - \nabla f_i(x^t), \theta_i^{t+1} \rangle + 2\gamma(1 - \beta) \langle v^t - \nabla f(x^t), \theta^{t+1} \rangle \\
2815 & + \frac{8\gamma\beta^2}{n\eta^2} (1 - \beta) \sum_{i=1}^n \langle \nabla f_i(x^t) - \nabla f_i(x^{t+1}), \theta_i^{t+1} \rangle \\
2816 & + 2\gamma(1 - \beta) \langle \nabla f(x^t) - \nabla f(x^{t+1}), \theta^{t+1} \rangle.
\end{aligned}$$

2828 Using stepsize restriction (v_i) we get rid of the term with R^t and obtain

$$\begin{aligned}
2831 \quad \Phi^{t+1} &\leq \Phi^t - \frac{\gamma}{2} \|\nabla f(x^t)\|^2 + b^2 \left(\frac{\beta^2\gamma}{\eta} + \frac{4\gamma\beta^3}{\eta^2} \right) + c^2 \frac{\gamma\beta}{n} \\
2832 & + \frac{2\gamma\beta}{n\eta} (1 - \eta)^2 \sum_{i=1}^n \langle (g_i^t - v_i^t) + \beta(v_i^t - \nabla f_i(x^t)) + \beta(\nabla f_i(x^t) - \nabla f_i(x^{t+1})), \theta_i^{t+1} \rangle \\
2833 & + \frac{8\gamma\beta^2}{n\eta^2} (1 - \beta) \sum_{i=1}^n \langle v_i^t - \nabla f_i(x^t), \theta_i^{t+1} \rangle + 2\gamma(1 - \beta) \langle v^t - \nabla f(x^t), \theta^{t+1} \rangle \\
2834 & + \frac{8\gamma\beta^2}{n\eta^2} (1 - \beta) \sum_{i=1}^n \langle \nabla f_i(x^t) - \nabla f_i(x^{t+1}), \theta_i^{t+1} \rangle \\
2835 & + 2\gamma(1 - \beta) \langle \nabla f(x^t) - \nabla f(x^{t+1}), \theta^{t+1} \rangle.
\end{aligned}$$

2845 Now we sum all the inequalities above for $t \in \{0, \dots, K\}$ and get

$$\begin{aligned}
2849 \quad \Phi^{K+1} &\leq \Phi^0 - \frac{\gamma}{2} \sum_{t=0}^K \|\nabla f(x^t)\|^2 + Kb^2 \left(\frac{\beta^2\gamma}{\eta} + \frac{4\gamma\beta^3}{\eta^2} \right) + Kc^2 \frac{\gamma\beta}{n} \\
2850 & + \frac{2\gamma\beta}{n\eta} (1 - \eta)^2 \sum_{t=0}^K \sum_{i=1}^n \langle (g_i^t - v_i^t) + \beta(v_i^t - \nabla f_i(x^t)) + \beta(\nabla f_i(x^t) - \nabla f_i(x^{t+1})), \theta_i^{t+1} \rangle \\
2851 & + \frac{8\gamma\beta^2}{n\eta^2} (1 - \beta) \sum_{t=0}^K \sum_{i=1}^n \langle v_i^t - \nabla f_i(x^t), \theta_i^{t+1} \rangle + 2\gamma(1 - \beta) \sum_{t=0}^K \langle v^t - \nabla f(x^t), \theta^{t+1} \rangle \\
2852 & + \frac{8\gamma\beta^2}{n\eta^2} (1 - \beta) \sum_{t=0}^K \sum_{i=1}^n \langle \nabla f_i(x^t) - \nabla f_i(x^{t+1}), \theta_i^{t+1} \rangle \\
2853 & + 2\gamma(1 - \beta) \sum_{t=0}^K \langle \nabla f(x^t) - \nabla f(x^{t+1}), \theta^{t+1} \rangle.
\end{aligned} \tag{58}$$

2862 Rearranging terms we get

$$\begin{aligned}
2863 & \frac{\gamma}{2} \sum_{t=0}^K \|\nabla f(x^t)\|^2 \leq \Phi^0 - \Phi^{K+1} + Kb^2 \left(\frac{\beta^2 \gamma}{\eta} + \frac{4\gamma\beta^3}{\eta^2} \right) + Kc^2 \frac{\gamma\beta}{n} \\
2864 & + \frac{2\gamma\beta}{n\eta} (1-\eta)^2 \sum_{t=0}^K \sum_{i=1}^n \langle (g_i^t - v_i^t) + \beta(v_i^t - \nabla f_i(x^t)) + \beta(\nabla f_i(x^t) - \nabla f_i(x^{t+1})), \theta_i^{t+1} \rangle \\
2865 & + \frac{8\gamma\beta^2}{n\eta^2} (1-\beta) \sum_{t=0}^K \sum_{i=1}^n \langle v_i^t - \nabla f_i(x^t), \theta_i^{t+1} \rangle + 2\gamma(1-\beta) \sum_{t=0}^K \langle v^t - \nabla f(x^t), \theta^{t+1} \rangle \\
2866 & + \frac{8\gamma\beta^2}{n\eta^2} (1-\beta) \sum_{t=0}^K \sum_{i=1}^n \langle \nabla f_i(x^t) - \nabla f_i(x^{t+1}), \theta_i^{t+1} \rangle \\
2867 & + 2\gamma(1-\beta) \sum_{t=0}^K \langle \nabla f(x^t) - \nabla f(x^{t+1}), \theta^{t+1} \rangle.
\end{aligned}$$

2878 Taking into account that $\frac{\gamma}{2} \sum_{t=0}^K \|\nabla f(x^t)\|^2 \geq 0$, we get that the event $E_K \cap \left(\bigcap_{i=1}^n \bar{\Theta}_i^{K+1} \right) \cap \bar{N}^t$
2879 implies

$$\begin{aligned}
2881 & \Phi^{K+1} \leq \Phi^0 + Kb^2 \left(\frac{\beta^2 \gamma}{\eta} + \frac{4\gamma\beta^3}{\eta^2} \right) + Kc^2 \frac{\gamma\beta}{n} \\
2882 & + \frac{2\gamma\beta}{n\eta} (1-\eta)^2 \sum_{t=0}^K \sum_{i=1}^n \langle (g_i^t - v_i^t) + \beta(v_i^t - \nabla f_i(x^t)) + \beta(\nabla f_i(x^t) - \nabla f_i(x^{t+1})), \theta_i^{t+1} \rangle \\
2883 & + \frac{8\gamma\beta^2}{n\eta^2} (1-\beta) \sum_{t=0}^K \sum_{i=1}^n \langle v_i^t - \nabla f_i(x^t), \theta_i^{t+1} \rangle + \frac{2\gamma(1-\beta)}{n} \sum_{t=0}^K \sum_{i=1}^n \langle v^t - \nabla f(x^t), \theta_i^{t+1} \rangle \\
2884 & + \frac{8\gamma\beta^2}{n\eta^2} (1-\beta) \sum_{t=0}^K \sum_{i=1}^n \langle \nabla f_i(x^t) - \nabla f_i(x^{t+1}), \theta_i^{t+1} \rangle \\
2885 & + \frac{2\gamma(1-\beta)}{n} \sum_{t=0}^K \sum_{i=1}^n \langle \nabla f(x^t) - \nabla f(x^{t+1}), \theta_i^{t+1} \rangle.
\end{aligned}$$

2895 Now we define the following random vectors:

$$\begin{aligned}
2896 & \zeta_{1,i}^t := \begin{cases} g_i^t - v_i^t, & \text{if } \|g_i^t - v_i^t\| \leq B, \\ 0, & \text{otherwise} \end{cases}, \\
2897 & \zeta_{2,t}^t := \begin{cases} v_i^t - \nabla f_i(x^t), & \text{if } \|v_i^t - \nabla f_i(x^t)\| \leq \sqrt{4L\Delta} + \frac{3}{2}(B - \tau) + 2b, \\ 0, & \text{otherwise} \end{cases}, \\
2898 & \zeta_{3,i}^t := \begin{cases} \nabla f_i(x^t) - \nabla f_i(x^{t+1}), & \text{if } \|\nabla f_i(x^t) - \nabla f_i(x^{t+1})\| \leq L\gamma \left(\sqrt{64L\Delta} + 3(B - \tau) + 3b \right), \\ 0, & \text{otherwise} \end{cases}, \\
2899 & \zeta_4^t := \begin{cases} v^t - \nabla f(x^t), & \text{if } \|v^t - \nabla f(x^t)\| \leq \sqrt{4L\Delta} + \frac{3}{2}(B - \tau) + 2ba, \\ 0, & \text{otherwise} \end{cases}, \\
2900 & \zeta_5^t := \begin{cases} \nabla f(x^t) - \nabla f(x^{t+1}), & \text{if } \|\nabla f(x^t) - \nabla f(x^{t+1})\| \leq L\gamma \left(\sqrt{64L\Delta} + 3(B - \tau) + 3b \right), \\ 0, & \text{otherwise} \end{cases}.
\end{aligned}$$

2910 By definition, all introduced random vectors $\zeta_{l,i}^t, l \in [3], i \in [n], \zeta_{4,5}^t$ are bounded with probability

2911 1. Moreover, by the definition of Φ^t and definition of E_t we get that the event $E_K \cap \bar{\Theta}^{K+1} \cap$
2912 $\left(\bigcap_{i=1}^n \bar{\Theta}_i^{K+1} \right) \cap \bar{N}^{K+1}$ implies

$$\begin{aligned}
2913 & \zeta_{1,i}^t = g_i^t - v_i^t, \quad \zeta_{2,i}^t = v_i^t - \nabla f_i(x^t), \quad \zeta_{3,i}^t = \nabla f_i(x^t) - \nabla f_i(x^{t+1}), \\
2914 & \zeta_4^t = v^t - \nabla f(x^t), \quad \zeta_5^t = \nabla f(x^t) - \nabla f(x^{t+1}).
\end{aligned}$$

Therefore, the event $E_K \cap \bar{\Theta}^{K+1} \cap \left(\bigcap_{i=1}^n \bar{\Theta}_i^{K+1}\right) \cap \bar{N}^{K+1}$ implies

$$\begin{aligned}
\Phi^{K+1} &\leq \Phi^0 + \underbrace{Kb^2 \left(\frac{\beta^2 \gamma}{\eta} + \frac{4\gamma \beta^3}{\eta^2} \right) + Kc^2 \frac{\gamma \beta}{n}}_{\textcircled{1}} + \underbrace{\frac{2\gamma \beta}{n\eta} (1-\eta)^2 \sum_{t=0}^K \sum_{i=1}^n \langle \zeta_{1,i}^t, \theta_i^{t+1} \rangle}_{\textcircled{2}} \\
&+ \underbrace{\frac{2\gamma \beta^2}{n\eta} (1-\eta)^2 \sum_{t=0}^K \sum_{i=1}^n \langle \zeta_{2,i}^t, \theta_i^{t+1} \rangle}_{\textcircled{3}} + \underbrace{\frac{2\gamma \beta^2}{n\eta} (1-\eta)^2 \sum_{t=0}^K \sum_{i=1}^n \langle \zeta_{3,i}^t, \theta_i^{t+1} \rangle}_{\textcircled{4}} \\
&+ \underbrace{\frac{8\gamma \beta^2}{n\eta^2} (1-\beta) \sum_{t=0}^K \sum_{i=1}^n \langle \zeta_{2,i}^t, \theta_i^{t+1} \rangle}_{\textcircled{5}} + \underbrace{\frac{2\gamma(1-\beta)}{n} \sum_{t=0}^K \sum_{i=1}^n \langle \zeta_4, \theta_i^{t+1} \rangle}_{\textcircled{6}} \\
&+ \underbrace{\frac{8\gamma \beta^2}{n\eta^2} (1-\beta) \sum_{t=0}^K \sum_{i=1}^n \langle \zeta_{3,i}^t, \theta_i^{t+1} \rangle}_{\textcircled{7}} + \underbrace{\frac{2\gamma(1-\beta)}{n} \sum_{t=0}^K \sum_{i=1}^n \langle \zeta_5, \theta_i^{t+1} \rangle}_{\textcircled{8}}.
\end{aligned}$$

BOUND OF THE TERM ①. For the term ① we have

$$Kb^2 \left(\frac{\beta^2 \gamma}{\eta} + \frac{4\gamma \beta^3}{\eta^2} \right) + Kc^2 \frac{\gamma \beta}{n} \leq Kb^2 \left(\frac{\beta^3}{4L\eta} + \frac{\beta^4}{L\eta^2} \right) + Kc^2 \frac{\beta^2}{4Ln}.$$

By choosing γ such that

$$\beta \leq \min \left\{ \left(\frac{L\Delta\eta}{6Tb^2} \right)^{\frac{1}{3}}, \left(\frac{L\Delta\eta^2}{24Tb^2} \right)^{\frac{1}{4}}, \left(\frac{L\Delta n}{6Tc^2} \right)^{\frac{1}{2}} \right\} \quad (59)$$

we get that

$$Kb^2 \left(\frac{\beta^2 \gamma}{\eta} + \frac{4\gamma \beta^3}{\eta^2} \right) + Kc^2 \frac{\gamma \beta}{n} \leq 3 \cdot \frac{\Delta}{24} = \frac{\Delta}{8}.$$

This bound holds with probability 1. Note that the worst dependency in the restriction on β in T is $\mathcal{O}(1/T^{\frac{1}{2}})$ that comes from the last term in min.

BOUND OF THE TERM ②. For term ②, let us enumerate random variables as

$$\langle \zeta_{1,1}^0, \theta_1^0 \rangle, \dots, \langle \zeta_{1,n}^0, \theta_n^0 \rangle, \langle \zeta_{1,1}^1, \theta_1^2 \rangle, \dots, \langle \zeta_{1,n}^1, \theta_n^2 \rangle, \dots, \langle \zeta_{1,1}^K, \theta_1^{K+1} \rangle, \dots, \langle \zeta_{1,n}^K, \theta_n^{K+1} \rangle,$$

i.e. first by index i , then by index t . Then we have that the event $E_K \cap \left(\bigcap_{i=1}^n \bar{\Theta}_i^{K+1}\right)$ implies

$$\mathbb{E} \left[\frac{2\gamma \beta}{n\eta} (1-\eta)^2 \langle \zeta_{1,i}^l, \theta_i^{l+1} \rangle \mid \langle \zeta_{1,i-1}^l, \theta_{i-1}^{l+1} \rangle, \dots, \langle \zeta_{1,1}^l, \theta_1^{l+1} \rangle, \dots, \langle \zeta_{1,1}^0, \theta_1^1 \rangle \right] = 0,$$

because $\{\theta_i^{l+1}\}_{i=1}^n$ are independent. Let

$$\sigma_2^2 := \frac{4\gamma^2 \beta^2}{n^2 \eta^2} \cdot B^2 \cdot \sigma^2.$$

2970 Since θ_i^{l+1} is σ -sub-Gaussian random vector, we have

$$\begin{aligned}
2971 & \\
2972 & \mathbb{E} \left[\exp \left(\left| \frac{1}{\sigma_2^2} \frac{4\gamma^2\beta^2}{n^2\eta^2} (1-\eta)^4 \langle \zeta_{1,i}^l, \theta_i^{l+1} \rangle^2 \right) \mid l, i-1 \right] \\
2973 & \leq \mathbb{E} \left[\exp \left(\frac{1}{\sigma_1^2} \frac{4\gamma^2\beta^2}{n^2\eta^2} \|\zeta_{1,i}^l\|^2 \cdot \|\theta_i^{l+1}\|^2 \mid l, i-1 \right) \right] \\
2974 & \leq \mathbb{E} \left[\exp \left(\frac{1}{\sigma_2^2} \frac{4\gamma^2\beta^2}{n^2\eta^2} \cdot B^2 \|\theta_i^{l+1}\|^2 \mid l, i-1 \right) \right] \\
2975 & \leq \mathbb{E} \left[\exp \left(\frac{n^2\eta^2}{4\gamma^2\beta^2 \cdot B^2 \cdot \sigma^2} \frac{4\gamma^2\beta^2}{n^2\eta^2} \cdot B^2 \|\theta_i^{l+1}\|^2 \mid l, i-1 \right) \right] \\
2976 & = \mathbb{E} \left[\exp \left(\frac{\|\theta_i^{l+1}\|^2}{\sigma^2} \mid l, i-1 \right) \right] \leq \exp(1). \\
2977 & \\
2978 & \\
2979 & \\
2980 & \\
2981 & \\
2982 & \\
2983 & \\
2984 &
\end{aligned}$$

2985 Here $\mathbb{E}[\cdot \mid l, i-1]$ means

$$\begin{aligned}
2986 & \\
2987 & \mathbb{E}[\cdot \mid \langle \zeta_{1,i-1}^l, \theta_{i-1}^{l+1} \rangle, \dots, \langle \zeta_{1,1}^l, \theta_1^{l+1} \rangle, \dots, \langle \zeta_{1,1}^0, \theta_1^1 \rangle] = 0, \\
2988 & \\
2989 &
\end{aligned}$$

2989 Therefore, we have by Lemma 1 with $\sigma_k^2 \equiv \sigma_2^2$ that

$$\begin{aligned}
2990 & \\
2991 & \Pr \left(\frac{2\gamma\beta}{n\eta} (1-\eta)^2 \left\| \sum_{t=0}^K \sum_{i=1}^n \langle \zeta_{1,i}^t, \theta_i^{t+1} \rangle \right\| \geq (\sqrt{2} + \sqrt{2}b_1) \sqrt{\sum_{t=0}^K \sum_{i=1}^n \frac{4B^2\gamma^2\beta^2\sigma^2}{n^2\eta^2}} \right) \\
2992 & \leq \exp(-b_1^2/3) \\
2993 & = \frac{\alpha}{14(T+1)} \\
2994 & \\
2995 & \\
2996 & \\
2997 &
\end{aligned}$$

2998 with $b_1^2 = 3 \log \left(\frac{14(T+1)}{\alpha} \right)$. Note that

$$\begin{aligned}
3000 & \\
3001 & (\sqrt{2} + \sqrt{2}b_1) \sqrt{\sum_{t=0}^K \sum_{i=1}^n \frac{4B^2\gamma^2\beta^2\sigma^2}{n^2\eta^2}} \leq (\sqrt{2} + \sqrt{2}b_1) \sqrt{\sum_{t=0}^K \sum_{i=1}^n \frac{B^2\beta^4\sigma^2}{4L^2n^2\eta^2}} \\
3002 & = (\sqrt{2} + \sqrt{2}b_1) \frac{B\beta^2\sigma}{2Ln\eta} \sqrt{(K+1)n} \\
3003 & \leq \frac{\Delta}{8}, \\
3004 & \\
3005 & \\
3006 & \\
3007 & \\
3008 &
\end{aligned}$$

3008 because we choose

$$\begin{aligned}
3009 & \\
3010 & \beta \leq \left(\frac{L\Delta\sqrt{n}\eta}{4\sqrt{2}(1+b_1)B\sigma\sqrt{T}} \right)^{\frac{1}{2}}, \quad \text{and} \quad K+1 \leq T. \tag{60} \\
3011 & \\
3012 &
\end{aligned}$$

3013 This implies that

$$\begin{aligned}
3014 & \\
3015 & \Pr \left(\frac{2\gamma\beta}{n\eta} (1-\eta)^2 \left\| \sum_{t=0}^K \sum_{i=1}^n \langle \zeta_{1,i}^t, \theta_i^{t+1} \rangle \right\| \geq \frac{\Delta}{8} \right) \leq \frac{\alpha}{14(T+1)} \\
3016 & \\
3017 &
\end{aligned}$$

3018 with this choice of momentum parameter. The dependency on T is $\tilde{O}(1/T^{\frac{1}{4}})$.

3021 **BOUND OF THE TERM ③.** The bound in this case is similar to the previous one. Let

$$\begin{aligned}
3022 & \\
3023 & \sigma_3^2 := \frac{4\gamma^2\beta^4}{n^2\eta^2} \cdot \left(\sqrt{4L\Delta} + \frac{3}{2}(B-\tau) + 2b \right)^2 \cdot \sigma^2.
\end{aligned}$$

$$\begin{aligned}
& \mathbb{E} \left[\exp \left(\left| \frac{1}{\sigma_3^2} \frac{4\gamma^2\beta^4}{n^2\eta^2} (1-\eta)^4 \langle \zeta_{2,i}^l, \theta_i^{l+1} \rangle^2 \right| \mid l, i-1 \right) \right] \\
& \leq \mathbb{E} \left[\exp \left(\frac{1}{\sigma_3^2} \frac{4\gamma^2\beta^4}{n^2\eta^2} \|\zeta_{2,i}^l\|^2 \cdot \|\theta_i^{l+1}\|^2 \right) \right] \\
& \leq \mathbb{E} \left[\exp \left(\frac{1}{\sigma_3^2} \frac{4\gamma^2\beta^4}{n^2\eta^2} \cdot \left(\sqrt{4L\Delta} + \frac{3}{2}(B-\tau) + 2b \right)^2 \cdot \|\theta_i^{l+1}\|^2 \right) \mid l, i-1 \right] \\
& \leq \mathbb{E} \left[\exp \left(\left[\frac{4\gamma^2\beta^4}{n^2\eta^2} \cdot \left(\sqrt{4L\Delta} + \frac{3}{2}(B-\tau) + 2b \right)^2 \cdot \sigma^2 \right]^{-1} \cdot \frac{4\gamma^2\beta^4}{n^2\eta^2} \cdot \left(\sqrt{4L\Delta} + \frac{3}{2}(B-\tau) + 2b \right)^2 \cdot \|\theta_i^{l+1}\|^2 \right) \mid l, i-1 \right] \\
& = \mathbb{E} \left[\exp \left(\frac{\|\theta_i^{l+1}\|^2}{\sigma^2} \right) \mid l, i-1 \right] \leq \exp(1).
\end{aligned}$$

Therefore, we have by Lemma 1 that

$$\begin{aligned}
& \Pr \left[\frac{2\gamma\beta^2}{n\eta} (1-\eta)^2 \left\| \sum_{t=0}^K \sum_{i=1}^n \langle \zeta_{2,i}^t, \theta_i^{t+1} \rangle \right\| \right. \\
& \quad \left. \geq (\sqrt{2} + \sqrt{2}b_1) \sqrt{\sum_{t=0}^K \sum_{i=1}^n \frac{4\gamma^2\beta^4\sigma^2}{n^2\eta^2} \cdot \left(\sqrt{4L\Delta} + \frac{3}{2}(B-\tau) + 2b \right)^2} \right] \\
& \leq \exp(-b_1^2/3) = \frac{\alpha}{14(T+1)},
\end{aligned}$$

Note that

$$\begin{aligned}
& (\sqrt{2} + \sqrt{2}b_1) \sqrt{(K+1)n} \frac{2\gamma\beta^2\sigma}{\eta n} \left(\sqrt{4L\Delta} + \frac{3}{2}(B-\tau) + 2b \right) \\
& \leq (\sqrt{2} + \sqrt{2}b_1) \sqrt{(K+1)n} \frac{\beta^3\sigma}{2L\eta n} \left(\sqrt{4L\Delta} + \frac{3}{2}(B-\tau) + 2b \right) \\
& \leq \frac{\Delta}{8}.
\end{aligned}$$

because we choose

$$\beta \leq \left(\frac{L\Delta\eta\sqrt{n}}{4\sqrt{2}(1+b_1)\sigma\sqrt{T} \left(\sqrt{4L\Delta} + \frac{3}{2}(B-\tau) + 2b \right)} \right)^{\frac{1}{3}} \quad \text{and} \quad K+1 \leq T. \quad (61)$$

This implies

$$\Pr \left(\frac{2\gamma\beta^2}{n\eta} (1-\eta)^2 \left\| \sum_{t=0}^K \sum_{i=1}^n \langle \zeta_{2,i}^t, \theta_i^{t+1} \rangle \right\| \geq \frac{\Delta}{8} \right) \leq \frac{\alpha}{14(T+1)}.$$

Note that the worst dependency w.r.t. T is $\tilde{O}(1/T^{\frac{1}{6}})$.

BOUND OF THE TERM ④. The bound in this case is similar to the previous one. Let

$$\sigma_4^2 := \frac{4L^2\gamma^4\beta^4}{n^2\eta^2} \left(\sqrt{64L\Delta} + 3(B-\tau) + 3b \right)^2 \cdot \sigma^2.$$

3078 Then we have

$$\begin{aligned}
3080 & \mathbb{E} \left[\exp \left(\left| \frac{1}{\sigma_4^2} \frac{4\gamma^2\beta^4}{n^2\eta^2} (1-\eta)^4 \langle \zeta_{3,i}^l, \theta_i^{l+1} \rangle^2 \right| \mid l, i-1 \right) \right] \\
3082 & \leq \mathbb{E} \left[\exp \left(\frac{1}{\sigma_4^2} \frac{4\gamma^2\beta^4}{n^2\eta^2} \|\zeta_{3,i}^l\|^2 \cdot \|\theta_i^{l+1}\|^2 \mid l, i-1 \right) \right] \\
3084 & \leq \mathbb{E} \left[\exp \left(\frac{1}{\sigma_4^2} \frac{4\gamma^2\beta^4}{n^2\eta^2} \cdot L^2\gamma^2 \left(\sqrt{64L\Delta} + 3(B-\tau) + 3b + 3a \right)^2 \cdot \|\theta_i^{l+1}\|^2 \mid l, i-1 \right) \right] \\
3087 & \leq \mathbb{E} \left[\exp \left(\left[\frac{4L^2\gamma^4\beta^4}{n^2\eta^2} \left(\sqrt{64L\Delta} + 3(B-\tau) + 3b + 3a \right)^2 \cdot \sigma^2 \right]^{-1} \right. \right. \\
3088 & \quad \left. \left. \frac{4L^2\gamma^4\beta^4}{n^2\eta^2} \cdot \frac{16\beta^4}{81} \left(\sqrt{64L\Delta} + 3(B-\tau) + 3b + 3a \right)^2 \cdot \|\theta_i^{l+1}\|^2 \mid l, i-1 \right) \right] \\
3092 & = \mathbb{E} \left[\exp \left(\frac{\|\theta_i^{l+1}\|^2}{\sigma^2} \right) \right] \leq \exp(1).
\end{aligned}$$

3095 Therefore, we have by Lemma 1 that

$$\begin{aligned}
3097 & \Pr \left(\frac{2\gamma\beta^2}{n\eta} (1-\eta)^2 \left\| \sum_{t=0}^K \sum_{i=1}^n \langle \zeta_{3,i}^t, \theta_i^{t+1} \rangle \right\| \right. \\
3098 & \quad \left. \geq (\sqrt{2} + \sqrt{2}b_1) \sqrt{\sum_{t=0}^K \sum_{i=1}^n \frac{4L^2\gamma^4\beta^4\sigma^2}{n^2\eta^2} \cdot \left(\sqrt{64L\Delta} + 3(B-\tau+b) \right)^2} \right) \\
3104 & \leq \exp(-b_1^2/3) = \frac{\alpha}{14(T+1)}.
\end{aligned}$$

3106 Note that

$$\begin{aligned}
3108 & (\sqrt{2} + \sqrt{2}b_1) \sqrt{(K+1)n} \frac{2L\gamma^2\beta^2\sigma}{\eta n} \left(\sqrt{64L\Delta} + 3(B-\tau+b) + 3a \right) \\
3110 & \leq \sqrt{2}(1+b_1) \sqrt{(K+1)n} \frac{\beta^4\sigma}{8L\eta n} \left(\sqrt{64L\Delta} + 3(B-\tau+b) \right) \\
3113 & \leq \frac{\Delta}{8}.
\end{aligned}$$

3115 because we choose

$$3117 \beta \leq \left(\frac{L\Delta\eta\sqrt{n}}{\sqrt{2}(1+b_1)\sigma\sqrt{T} \left(\sqrt{64L\Delta} + 3(B-\tau+b) \right)} \right)^{\frac{1}{4}}, \quad \text{and } K+1 \leq T. \quad (62)$$

3121 This implies

$$3123 \Pr \left(\frac{2\gamma\beta^2}{n\eta} (1-\eta)^2 \left\| \sum_{t=0}^K \sum_{i=1}^n \langle \zeta_{2,i}^t, \theta_i^{t+1} \rangle \right\| \geq \frac{\Delta}{8} \right) \leq \frac{\alpha}{14(T+1)},$$

3126 Note that the worst dependency w.r.t. T is $\tilde{\mathcal{O}}(1/T^{\frac{1}{8}})$.

3128 **BOUND OF THE TERM ⑤.** The bound in this case is similar to the previous one. Let

$$3130 \sigma_5^2 := \frac{64\gamma^2\beta^4}{n^2\eta^4} \cdot \left(\sqrt{4L\Delta} + \frac{3}{2}(B-\tau) + 2b \right)^2 \cdot \sigma^2.$$

Then we have

$$\begin{aligned}
& \mathbb{E} \left[\exp \left(\left| \frac{1}{\sigma_5^2} \frac{64\gamma^2\beta^4}{n^2\eta^4} (1-\beta)^2 \langle \zeta_{2,i}^l, \theta_i^{l+1} \rangle^2 \right| \mid l, i-1 \right) \right] \\
& \leq \mathbb{E} \left[\exp \left(\frac{1}{\sigma_5^2} \frac{64\gamma^2\beta^4}{n^2\eta^4} \|\zeta_{2,i}^l\|^2 \cdot \|\theta_i^{l+1}\|^2 \mid l, i-1 \right) \right] \\
& \leq \mathbb{E} \left[\exp \left(\frac{1}{\sigma_5^2} \frac{64\gamma^2\beta^4}{n^2\eta^4} \cdot \left(\sqrt{4L\Delta} + \frac{3}{2}(B-\tau) + 2b \right)^2 \cdot \|\theta_i^{l+1}\|^2 \mid l, i-1 \right) \right] \\
& = \mathbb{E} \left[\exp \left(\left[\frac{64\gamma^2\beta^4}{81L^2n^2\eta^4} \cdot \left(\sqrt{4L\Delta} + \frac{3}{2}(B-\tau) + 2b \right)^2 \cdot \sigma^2 \right]^{-1} \right. \right. \\
& \quad \left. \left. \frac{64\gamma^2\beta^4}{n^2\eta^4} \cdot \left(\sqrt{4L\Delta} + \frac{3}{2}(B-\tau) + 2b \right)^2 \cdot \|\theta_i^{l+1}\|^2 \right) \mid l, i-1 \right] \\
& = \mathbb{E} \left[\exp \left(\frac{\|\theta_i^{l+1}\|^2}{\sigma^2} \right) \mid l, i-1 \right] \leq \exp(1).
\end{aligned}$$

Therefore, we have by Lemma 1 that

$$\begin{aligned}
& \Pr \left[\frac{8\gamma\beta^2}{n\eta^2} (1-\beta) \left\| \sum_{t=0}^K \sum_{i=1}^n \langle \zeta_{2,i}^t, \theta_i^{t+1} \rangle \right\| \right. \\
& \quad \left. \geq (\sqrt{2} + \sqrt{2}b_1) \sqrt{\sum_{t=0}^K \sum_{i=1}^n \frac{64\gamma^2\beta^4\sigma^2}{n^2\eta^4} \left(\sqrt{4L\Delta} + \frac{3}{2}(B-\tau) + 2b + (L\gamma)^pa \right)^2} \right] \\
& \leq \exp(-b_1^2/3) = \frac{\alpha}{14(T+1)}.
\end{aligned}$$

Note that

$$\begin{aligned}
& (\sqrt{2} + \sqrt{2}b_1) \sqrt{(K+1)n} \frac{8\gamma\beta^2\sigma}{n\eta^2} \left(\sqrt{4L\Delta} + \frac{3}{2}(B-\tau) + 2b \right) \\
& \leq (\sqrt{2} + \sqrt{2}b_1) \sqrt{(K+1)n} \frac{2\beta^3\sigma}{L\eta^2} \left(\sqrt{4L\Delta} + \frac{3}{2}(B-\tau) + 2b \right) \\
& \leq \frac{\Delta}{8}
\end{aligned}$$

because we choose

$$\beta \leq \left(\frac{L\Delta\eta^2\sqrt{n}}{16\sqrt{2}(1+b_1)\sigma\sqrt{T} \left(\sqrt{4L\Delta} + \frac{3}{2}(B-\tau) + 2b \right)} \right)^{\frac{1}{3}} \quad \text{and} \quad K+1 \leq T. \quad (63)$$

This implies

$$\Pr \left(\frac{8\gamma\beta^2}{n\eta^2} (1-\beta) \left\| \sum_{t=0}^K \sum_{i=1}^n \langle \zeta_{2,i}^t, \theta_i^{t+1} \rangle \right\| \geq \frac{\Delta}{8} \right) \leq \frac{\alpha}{14(T+1)}.$$

Note that the worst dependency w.r.t. T is $\tilde{\mathcal{O}}(1/T^{\frac{1}{6}})$.

BOUND OF THE TERM ⑦. The bound in this case is similar to the previous one. Let

$$\sigma_7^2 := \frac{64\gamma^2\beta^4}{n^2\eta^4} \left(\sqrt{64L\Delta} + 3(B-\tau+b) \right)^2 \cdot \sigma^2.$$

3186 Then we have
3187

$$\begin{aligned}
3188 & \mathbb{E} \left[\exp \left(\left| \frac{1}{\sigma_7^2} \frac{64L^2\gamma^4\beta^4}{n^2\eta^4} (1-\beta)^2 \langle \zeta_{3,i}^l, \theta_i^{l+1} \rangle^2 \right| \mid l, i-1 \right) \right] \\
3189 & \leq \mathbb{E} \left[\exp \left(\frac{1}{\sigma_7^2} \frac{64\gamma^2\beta^4}{n^2\eta^4} \|\zeta_{3,i}^l\|^2 \cdot \|\theta_i^{l+1}\|^2 \mid l, i-1 \right) \right] \\
3190 & \leq \mathbb{E} \left[\exp \left(\frac{64\gamma^2\beta^4}{n^2\eta^4} \cdot L^2\gamma^2 \left(\sqrt{64L\Delta} + 3(B-\tau+b) \right)^2 \cdot \|\theta_i^{l+1}\|^2 \mid l, i-1 \right) \right] \\
3191 & \leq \mathbb{E} \left[\exp \left(\left[\frac{64L^2\gamma^4\beta^4}{n^2\eta^4} \left(\sqrt{64L\Delta} + 3(B-\tau+b) \right)^2 \cdot \sigma^2 \right]^{-1} \right. \right. \\
3192 & \quad \left. \left. \frac{64L^2\gamma^4\beta^4}{n^2\eta^4} \left(\sqrt{64L\Delta} + 3(B-\tau+b) \right)^2 \cdot \|\theta_i^{l+1}\|^2 \mid l, i-1 \right) \right] \\
3193 & = \mathbb{E} \left[\exp \left(\frac{\|\theta_i^{l+1}\|^2}{\sigma^2} \mid l, i-1 \right) \right] \leq \exp(1). \\
3194 & \\
3195 & \\
3196 & \\
3197 & \\
3198 & \\
3199 & \\
3200 & \\
3201 & \\
3202 & \\
3203 & \\
3204 &
\end{aligned}$$

3204 Therefore, we have by Lemma 1 that

$$\begin{aligned}
3205 & \Pr \left[\frac{8\gamma\beta^2}{n\eta^2} (1-\beta) \left\| \sum_{t=0}^K \sum_{i=1}^n \langle \zeta_{3,i}^t, \theta_i^{t+1} \rangle \right\| \geq \right. \\
3206 & \quad \left. (\sqrt{2} + \sqrt{2}b_1) \sqrt{\sum_{t=0}^K \sum_{i=1}^n \frac{64L^2\gamma^4\beta^4\sigma^2}{n^2\eta^4} \cdot \left(\sqrt{64L\Delta} + 3(B-\tau+b) \right)^2} \right] \\
3207 & \leq \exp(-b_1^2/3) = \frac{\alpha}{14(T+1)}. \\
3208 & \\
3209 & \\
3210 & \\
3211 & \\
3212 & \\
3213 &
\end{aligned}$$

3214 Note that

$$\begin{aligned}
3215 & (\sqrt{2} + \sqrt{2}b_1) \sqrt{(K+1)n} \frac{8L\gamma^2\beta^2\sigma}{\eta^2n} \left(\sqrt{64L\Delta} + 3(B-\tau+b) \right) \\
3216 & \leq (\sqrt{2} + \sqrt{2}b_1) \sqrt{(K+1)n} \frac{\beta^4\sigma}{2L\eta^2n} \left(\sqrt{64L\Delta} + 3(B-\tau+b) \right) \\
3217 & \leq \frac{\Delta}{8} \\
3218 & \\
3219 & \\
3220 & \\
3221 & \\
3222 &
\end{aligned}$$

3223 because we choose

$$3224 \beta \leq \left(\frac{L\Delta\eta^2\sqrt{n}}{4\sqrt{2}(1+b_1)\sigma\sqrt{T} \left(\sqrt{64L\Delta} + 3(B-\tau+b) \right)} \right)^{\frac{1}{4}} \quad \text{and} \quad K+1 \leq T. \quad (64)$$

3229 This implies

$$3230 \Pr \left(\frac{8\gamma\beta^2}{n\eta^2} (1-\beta) \left\| \sum_{t=0}^K \sum_{i=1}^n \langle \zeta_{3,i}^t, \theta_i^{t+1} \rangle \right\| \geq \frac{\Delta}{8} \right) \leq \frac{\alpha}{14(T+1)}.$$

3234 Note that the worst dependency w.r.t. T is $\tilde{\mathcal{O}}(1/T^{\frac{1}{8}})$.

3237 BOUND OF THE TERM ⑥. The bound in this case is similar to the previous one. Let

$$3238 \sigma_6^2 := \frac{4\gamma^2}{n^2} \left(\sqrt{4L\Delta} + \frac{3}{2}(B-\tau) + 2b \right)^2 \cdot \sigma^2.$$

3240 Then we have
3241

$$\begin{aligned}
3242 & \mathbb{E} \left[\exp \left(\left| \frac{1}{\sigma_6^2} \frac{4\gamma^2}{n^2} (1-\beta)^2 \langle \zeta_4^l, \theta_i^{l+1} \rangle^2 \right| \mid l, i-1 \right) \right] \\
3243 & \leq \mathbb{E} \left[\exp \left(\frac{1}{\sigma_6^2} \frac{4\gamma^2}{n^2} \|\zeta_4^l\|^2 \cdot \|\theta_i^{l+1}\|^2 \mid l, i-1 \right) \right] \\
3244 & \leq \mathbb{E} \left[\exp \left(\frac{1}{\sigma_6^2} \frac{4\gamma^2}{n^2} \left(\sqrt{4L\Delta} + \frac{3}{2}(B-\tau) + 2b \right)^2 \cdot \|\theta_i^{l+1}\|^2 \mid l, i-1 \right) \right] \\
3245 & \leq \mathbb{E} \left[\exp \left(\left[\frac{4\gamma^2}{n^2} \left(\sqrt{4L\Delta} + \frac{3}{2}(B-\tau) + 2b \right)^2 \cdot \sigma^2 \right]^{-1} \right. \right. \\
3246 & \quad \left. \left. \frac{4\gamma^2}{n^2} \left(\sqrt{4L\Delta} + \frac{3}{2}(B-\tau) + 2b \right)^2 \cdot \|\theta_i^{l+1}\|^2 \mid l, i-1 \right) \right] \\
3247 & = \mathbb{E} \left[\exp \left(\frac{\|\theta_i^{l+1}\|^2}{\sigma^2} \mid l, i-1 \right) \right] \leq \exp(1). \\
3248 & \\
3249 & \\
3250 & \\
3251 & \\
3252 & \\
3253 & \\
3254 & \\
3255 & \\
3256 & \\
3257 & \\
3258 &
\end{aligned}$$

3258 Therefore, we have by Lemma 1 that
3259

$$\begin{aligned}
3260 & \Pr \left[\frac{2\gamma(1-\beta)}{n} \left\| \sum_{t=0}^K \sum_{i=1}^n \langle \zeta_{4,i}^t, \theta_i^{t+1} \rangle \right\| \right] \\
3261 & \geq (\sqrt{2} + \sqrt{2}b_1) \sqrt{\sum_{t=0}^K \sum_{i=1}^n \frac{4\gamma^2}{n^2} \sigma^2 \cdot \left(\sqrt{4L\Delta} + \frac{3}{2}(B-\tau) + 2b \right)^2} \\
3262 & \leq \exp(-b_1^2/3) = \frac{\alpha}{14(T+1)}, \\
3263 & \\
3264 & \\
3265 & \\
3266 & \\
3267 & \\
3268 &
\end{aligned}$$

3269 Note that
3270

$$\begin{aligned}
3271 & (\sqrt{2} + \sqrt{2}b_1) \sqrt{(K+1)n} \cdot \frac{2\gamma}{n} \sigma \left(\sqrt{4L\Delta} + \frac{3}{2}(B-\tau) + 2b \right) \\
3272 & \leq (\sqrt{2} + \sqrt{2}b_1) \sqrt{(K+1)n} \cdot \frac{\beta\sigma}{2Ln} \left(\sqrt{4L\Delta} + \frac{3}{2}(B-\tau) + 2b \right) \\
3273 & \leq \frac{\Delta}{8} \\
3274 & \\
3275 & \\
3276 & \\
3277 &
\end{aligned}$$

3278 because we choose
3279

$$\beta \leq \left(\frac{L\Delta\sqrt{n}}{32\sqrt{2}(1+b_1)\sigma\sqrt{T} \left(\sqrt{4L\Delta} + \frac{3}{2}(B-\tau) + 2b \right)} \right) \quad \text{and} \quad K+1 \leq T. \quad (65)$$

3283 This implies
3284

$$\Pr \left(\frac{2\gamma(1-\beta)}{n} \left\| \sum_{t=0}^K \sum_{i=1}^n \langle \zeta_{4,i}^t, \theta_i^{t+1} \rangle \right\| \geq \frac{\Delta}{8} \right) \leq \frac{\alpha}{14(T+1)}.$$

3288 Note that the worst dependency w.r.t. T is $\tilde{\mathcal{O}}(1/T^{1/2})$.
3289

3290 **BOUND OF THE TERM ⑧.** The bound in this case is similar to the previous one. Let
3291

$$\sigma_8^2 := \frac{4L^2\gamma^4}{n^2} \cdot \left(\sqrt{64L\Delta} + 3(B-\tau+b) \right)^2 \cdot \sigma^2.$$

3292
3293

Then we have

$$\begin{aligned}
& \mathbb{E} \left[\exp \left(\left| \frac{1}{\sigma_8^2} \frac{4\gamma^2}{n^2} (1-\beta)^2 \langle \zeta_5^l, \theta_i^{l+1} \rangle^2 \right| \right) \mid l, i-1 \right] \\
& \leq \mathbb{E} \left[\exp \left(\frac{1}{\sigma_8^2} \frac{4\gamma^2}{n^2} \|\zeta_5^l\|^2 \cdot \|\theta_i^{l+1}\|^2 \right) \mid l, i-1 \right] \\
& \leq \mathbb{E} \left[\exp \left(\frac{1}{\sigma_8^2} \frac{4\gamma^2}{n^2} L^2 \gamma^2 \left(\sqrt{64L\Delta} + 3(B-\tau+b) \right) \cdot \|\theta_i^{l+1}\|^2 \right) \mid l, i-1 \right].
\end{aligned}$$

Since θ_i^{l+1} is sub-Gaussian with parameter σ^2 , then we can continue the chain of inequalities above using the definition of σ_8^2

$$\begin{aligned}
& \mathbb{E} \left[\exp \left(\left[\frac{4L^2\gamma^4}{n^2} \cdot \left(\sqrt{64L\Delta} + 3(B-\tau+b) + 3a \right)^2 \cdot \sigma^2 \right]^{-1} \right. \right. \\
& \quad \left. \left. \frac{4L^2\gamma^4}{n^2} \cdot \left(\sqrt{64L\Delta} + 3(B-\tau+b) \right)^2 \cdot \|\theta_i^{l+1}\|^2 \right) \mid l, i-1 \right] \\
& = \mathbb{E} \left[\exp \left(\frac{\|\theta_i^{l+1}\|^2}{\sigma^2} \right) \right] \leq \exp(1).
\end{aligned}$$

Therefore, we have by Lemma 1 that

$$\begin{aligned}
& \Pr \left[\frac{2\gamma(1-\beta)}{n} \left\| \sum_{t=0}^K \sum_{i=1}^n \langle \zeta_{5,i}^t, \theta^{t+1} \rangle \right\| \right. \\
& \geq (\sqrt{2} + \sqrt{2}b_1) \sqrt{\sum_{t=0}^K \sum_{i=1}^n \frac{4L^2\gamma^4}{n^2} \sigma^2 \cdot \left(\sqrt{64L\Delta} + 3(B-\tau+b) \right)^2} \\
& \left. \leq \exp(-b_1^2/3) = \frac{\alpha}{14(T+1)}. \right]
\end{aligned}$$

Note that

$$\begin{aligned}
& (\sqrt{2} + \sqrt{2}b_1) \sqrt{(K+1)n} \cdot \frac{2L\gamma^2}{n} \sigma \left(\sqrt{64L\Delta} + 3(B-\tau+b) \right) \\
& \leq (\sqrt{2} + \sqrt{2}b_1) \sqrt{(K+1)n} \cdot \frac{\beta^2}{2Ln} \sigma \left(\sqrt{64L\Delta} + 3(B-\tau+b) \right) \\
& \leq \frac{\Delta}{8}
\end{aligned}$$

because we choose

$$\beta \leq \left(\frac{L\Delta\sqrt{n}}{4\sqrt{2}(1+b_1)\sigma\sqrt{T} \left(\sqrt{64L\Delta} + 3(B-\tau+b) \right)} \right)^{\frac{1}{2}} \quad \text{and} \quad K+1 \leq T. \quad (66)$$

This implies

$$\Pr \left(2\gamma(1-\beta) \left\| \sum_{t=0}^K \sum_{i=1}^n \langle \zeta_{5,i}^t, \theta^{t+1} \rangle \right\| \geq \frac{\Delta}{8} \right) \leq \frac{\alpha}{14(T+1)}.$$

Note that the worst dependency w.r.t T is $\tilde{O}(1/T^{\frac{1}{4}})$.

Final probability. Therefore, the probability event

$$\Omega := E_K \cap \bar{\Theta}^{K+1} \cap \left(\bigcap_{i=1}^n \bar{\Theta}_i^{K+1} \right) \cap \bar{N}^{K+1} \cap E_{\text{①}} \cap E_{\text{②}} \cap E_{\text{③}} \cap E_{\text{④}} \cap E_{\text{⑤}} \cap E_{\text{⑥}} \cap E_{\text{⑦}} \cap E_{\text{⑧}},$$

where each $E_{\textcircled{1}}-E_{\textcircled{8}}$ denotes that each of 1-8-th terms is smaller than $\frac{\Delta}{8}$ implies that

$$\textcircled{1} + \textcircled{2} + \textcircled{3} + \textcircled{4} + \textcircled{5} + \textcircled{6} + \textcircled{7} + \textcircled{8} \leq 8 \cdot \frac{\Phi^0}{8} = \Delta,$$

i.e. 7 in the induction assumption holds. Moreover, this also implies that

$$\Phi^{K+1} \leq \Phi^0 + \Delta \leq \Delta + \Delta = 2\Delta,$$

i.e. 6 in the induction assumption holds. The probability $\Pr(E_{K+1})$ can be lower bounded as follows

$$\begin{aligned} \Pr(E_{K+1}) &\geq \Pr(\Omega) \\ &= \Pr\left(E_K \cap \bar{\Theta}^{K+1} \cap \left(\bigcap_{i=1}^n \bar{\Theta}_i^{K+1}\right) \cap \bar{N}^{K+1} \cap E_{\textcircled{1}} \cap E_{\textcircled{2}} \cap E_{\textcircled{3}} \cap E_{\textcircled{4}} \cap E_{\textcircled{5}} \cap E_{\textcircled{6}} \right. \\ &\quad \left. \cap E_{\textcircled{7}} \cap E_{\textcircled{8}}\right) \\ &= 1 - \Pr\left(\bar{E}_K \cup \Theta^{K+1} \cup \left(\bigcup_{i=1}^n \Theta_i^{K+1}\right) \cup N^{K+1} \cup \bar{E}_{\textcircled{1}} \cup \bar{E}_{\textcircled{2}} \cup \bar{E}_{\textcircled{3}} \cup \bar{E}_{\textcircled{4}} \cup \bar{E}_{\textcircled{5}} \cup \bar{E}_{\textcircled{6}} \right. \\ &\quad \left. \cup \bar{E}_{\textcircled{7}} \cup \bar{E}_{\textcircled{8}}\right) \\ &\geq 1 - \Pr(\bar{E}_K) - \Pr(\Theta^{K+1}) - \sum_{i=1}^n \Pr(\Theta_i^{K+1}) - \Pr(N^{K+1}) - \Pr(\bar{E}_{\textcircled{1}}) - \Pr(\bar{E}_{\textcircled{2}}) \\ &\quad - \Pr(\bar{E}_{\textcircled{3}}) - \Pr(\bar{E}_{\textcircled{4}}) - \Pr(\bar{E}_{\textcircled{5}}) - \Pr(\bar{E}_{\textcircled{6}}) - \Pr(\bar{E}_{\textcircled{7}}) - \Pr(\bar{E}_{\textcircled{8}}) \\ &\geq 1 - \frac{\alpha(K+1)}{T+1} - \frac{\alpha}{6(T+1)} - \sum_{i=1}^n \frac{\alpha}{6n(T+1)} - \frac{\alpha}{6(T+1)} - 0 - 7 \cdot \frac{\alpha}{14(T+1)} \\ &= 1 - \frac{\alpha(K+2)}{T+1}. \end{aligned}$$

This finalizes the transition step of induction. The result of the theorem follows by setting $K = T - 1$. Indeed, from (58) we obtain

$$\frac{\gamma}{2} \sum_{t=0}^K \|\nabla f(x^t)\|^2 \leq \Phi^0 - \Phi^{K+1} + \Delta \leq 2\Delta \Rightarrow \frac{1}{T} \sum_{t=0}^{T-1} \|\nabla f(x^t)\|^2 \leq \frac{4\Delta}{\gamma T}. \quad (67)$$

Final rate. Now we have the following restrictions on the momentum parameter in terms of dependency on T from each bound of terms 1-8 correspondingly

$$\begin{aligned} \beta \leq \tilde{\mathcal{O}} &\left(\underbrace{\left(\frac{L\Delta n}{T\sigma^2}\right)^{\frac{1}{2}}}_{\text{from term 1}}, \underbrace{\left(\frac{L\Delta\sqrt{n}\eta}{B\sigma\sqrt{T}}\right)^{\frac{1}{2}}}_{\text{from term 2}}, \underbrace{\left(\frac{L\Delta\sqrt{n}\eta}{\sigma\sqrt{T}(\sqrt{4L\Delta} + (B-\tau) + 2b)}\right)^{\frac{1}{3}}}_{\text{from term 3}}, \right. \\ &\underbrace{\left(\frac{L\Delta\eta\sqrt{n}}{\sigma\sqrt{T}(\sqrt{64L\Delta} + 3(B-\tau) + b)}\right)^{\frac{1}{4}}}_{\text{from term 4}}, \underbrace{\left(\frac{L\Delta\eta^2\sqrt{n}}{\sigma\sqrt{T}(\sqrt{4L\Delta} + 3/2(B-\tau) + 2b)}\right)^{\frac{1}{3}}}_{\text{from term 5}}, \\ &\underbrace{\left(\frac{L\Delta\eta^2\sqrt{n}}{\sigma\sqrt{T}(\sqrt{64L\Delta} + 3(B-\tau) + b)}\right)^{\frac{1-p}{2(2-p)}}}_{\text{from term 7}}, \underbrace{\left(\frac{L\Delta\sqrt{n}}{\sigma\sqrt{T}(\sqrt{4L\Delta} + 3/2(B-\tau) + 2b)}\right)}_{\text{from term 6}}, \\ &\left. \underbrace{\left(\frac{L\Delta\sqrt{n}}{\sigma\sqrt{T}(\sqrt{64L\Delta} + 3(B-\tau) + b)}\right)^{\frac{1}{2}}}_{\text{from term 8}} \right). \end{aligned}$$

Multiplying the above by $\frac{1}{4L}$ gives restrictions on γ . The worst dependency on T is given by $\textcircled{6}$ term and translates to the rate of the form

$$\tilde{\mathcal{O}}\left(\frac{L\Delta}{T} \frac{\sigma\sqrt{T}(\sqrt{L\Delta} + B + b)}{L\Delta\sqrt{n}}\right) = \tilde{\mathcal{O}}\left(\frac{\sigma(\sqrt{L\Delta} + B + b)}{\sqrt{Tn}}\right).$$

Therefore, with probability $1 - \alpha$ Clip21-SGDM converges as

$$\frac{1}{T} \sum_{t=0}^{T-1} \|\nabla f(x^t)\|^2 \leq \tilde{\mathcal{O}} \left(\frac{L\Delta \sigma \sqrt{T} (\sqrt{L\Delta} + B + b)}{T L\Delta \sqrt{n}} \right) = \tilde{\mathcal{O}} \left(\frac{\sigma (\sqrt{L\Delta} + B + \sigma)}{\sqrt{Tn}} \right), \quad (68)$$

where $\tilde{\mathcal{O}}$ hides constant and logarithmic factors, and higher order terms decreasing with T .

CASE $\mathcal{I}_{K+1} = 0$. This case is even easier. The only change will be with the term next to R^t . We will get

$$1 - \frac{48L^2}{\eta^2} \gamma^2 - \frac{12L^2}{\beta^2} \gamma^2 \geq \frac{2}{3} - \frac{48L^2}{\eta} \gamma^2 \geq 0$$

instead of

$$1 - \frac{16\beta^2 L^2}{\eta^2} \gamma^2 - \frac{48L^2}{\eta^2} \gamma^2 - \frac{12L^2}{\beta^2} \gamma^2 \geq 0$$

as in the previous case. This difference comes from Lemma 18 because $\tilde{V}^{K+1} = 0$. The rest is a repetition of the previous derivations. \square

Remark 2. With $v_i^{-1} = g_i^{-1} = 0$ we have

$$\begin{aligned} \Phi^0 &= F^0 + \frac{\gamma}{\eta} \frac{1}{n} \sum_{i=1}^n \|\text{clip}_\tau(\beta \nabla f_i(x^0, \xi_i^0)) - \beta \nabla f_i(x^0, \xi_i^0)\|^2 + \frac{4\gamma\beta}{\eta^2} \frac{1}{n} (1-\beta)^2 \sum_{i=1}^n \|\nabla f_i(x^0, \xi_i^0)\|^2 \\ &\quad + \frac{\gamma}{\beta} (1-\beta)^2 \|\nabla f(x^0, \xi^0)\|^2 \\ &\leq F^0 + \frac{\gamma}{\eta} \frac{1}{n} \sum_{i=1}^n \max\{(\|\nabla f(x^0, \xi_i^0)\| - \tau)^2, 0\} + \frac{16L\gamma^2}{\eta^2} \frac{1}{n} (1-\beta)^2 \sum_{i=1}^n \|\nabla f_i(x^0, \xi_i^0)\|^2 \\ &\quad + \frac{1}{4L} (1-\beta)^2 \|\nabla f(x^0, \xi^0)\|^2. \end{aligned}$$

We have the stepsize restriction

$$\frac{2}{3} - \frac{64L^4\gamma^2}{\eta^2} - \frac{48L^2\gamma^2}{\eta^2} \geq 0. \quad (69)$$

For inequality of the form $a\gamma^2 + b\gamma \leq 1$ the stepsize restriction of the form $\gamma \leq \frac{1}{\sqrt{a+b}}$ is tight up to a constant factor 2, i.e. $\frac{2}{\sqrt{a+b}}$ does not satisfy the inequality (see Lemma 5 in (Richtárik et al., 2021)). Using this lemma in our case we get that the stepsize satisfying Equation (69) should also satisfy

$$L^2\gamma^2 \leq 2 \cdot \frac{\eta}{72/\eta + 4\sqrt{6}}.$$

This implies that $L^2\gamma^2 \leq \frac{\eta}{4\sqrt{6}}$ and $L^2\gamma^2 \leq \frac{\eta^2}{72}$. Consequently, it also satisfies $\frac{\gamma}{\eta} \leq \frac{1}{6L\sqrt{2}}$ (from the last inequality). Therefore, we have

$$\begin{aligned} \Phi^0 &\leq F^0 + \frac{1}{6L\sqrt{2}} \frac{1}{n} \sum_{i=1}^n \max\{(\|\nabla f_i(x^0, \xi_i^0)\| - \tau)^2, 0\} + \frac{2}{9L} \frac{1}{n} \sum_{i=1}^n \|\nabla f_i(x^0, \xi_i^0)\|^2 \\ &\quad + \frac{1}{4L} (1-\beta)^2 \|\nabla f(x^0, \xi^0)\|^2 \\ &\leq F^0 + \frac{1}{6L\sqrt{2}} \frac{1}{n} \sum_{i=1}^n \max\{\|\nabla f_i(x^0, \xi_i^0)\|^2, 0\} + \frac{2}{9L} \frac{1}{n} \sum_{i=1}^n \|\nabla f_i(x^0, \xi_i^0)\|^2 \\ &\quad + \frac{1}{4L} (1-\beta)^2 \|\nabla f(x^0, \xi^0)\|^2 \end{aligned}$$

which is independent of τ , and can be use as a bound for Δ . Terms containing $\|\nabla f_i(x^0, \xi_i^0)\|^2$ can be bounded by B^2 and $\|\nabla f_i(x^0)\|^2$ with high probability, i.e. Δ is again independent of τ .

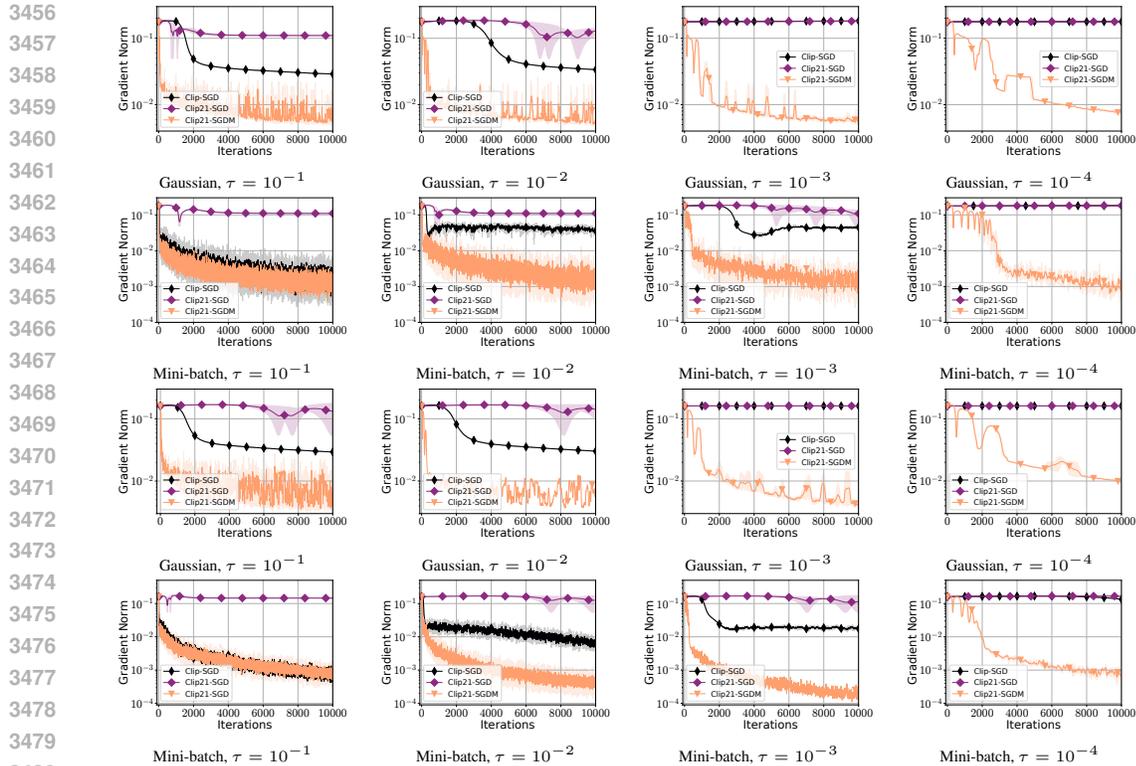


Figure 7: Comparison of Clip-SGD, Clip21-SGD, and Clip21-SGDM on logistic regression with non-convex regularization for various the clipping radii τ with mini-batch and Gaussian-added stochastic gradients on Duke (two first rows) and Leukemia (two last rows).

H EXPERIMENTS DETAILS AND MORE

H.1 EXPERIMENTS WITH LOGISTIC REGRESSION

H.1.1 STOCHASTIC SETTING VARYING CLIPPING RADIUS

We conduct experiments on non-convex logistic regression with regularization parameter $\lambda = 10^{-3}$ for 10^4 iterations. We use Duke and Leukemia datasets from LibSVM library and split the dataset into $n = 4$ equal parts. We normalize the row of the feature matrix to demonstrate the differences between algorithms. To simulate the stochastic gradients we either add centered Gaussian noise with variance $\sigma = 0.05$ for the Duke dataset and $\sigma = 0.1$ for the Leukemia dataset, or mini-batch gradients with batch-size of $\frac{1}{3}$ of the whole local dataset for Duke dataset and $\frac{1}{4}$ of the whole local dataset for Leukemia dataset. For Clip21-SGD and Clip-SGD algorithms, we tune the stepsize in $\{2^{-5}, \dots, 2^5\}$ and choose the one that gives the lowest final gradient norm in average across 3 random seeds. For Clip21-SGDM, we tune both the stepsize in $\{2^{-5}, \dots, 2^5\}$ and the momentum parameter in $\{0.1, 0.5, 0.9\}$ and choose the best pair of parameters similarly as before. For completeness, we report the convergence curves in Figure 7. We observe that Clip21-SGDM is more robust to the choice of the clipping radius τ while Clip-SGD converges well only for large enough τ . Besides, Clip21-SGD does not converge in all cases which is also highlighted by our theory in Theorem 1.

H.1.2 STOCHASTIC SETTING WITH ADDITIVE DP NOISE

We describe the setting in more detail for completeness. First, note that we use the same set of problem parameters as in Appendix H.1.1 such as n, λ, σ , and batch-size. Next, we fix a ratio between Gaussian DP noise variance σ_ω and the clipping parameter τ from $\{0.1, 1.0, 10.0\}$. For a given ratio, we tune Clip21-SGD and Clip-SGD algorithms across all possible pairs of the step-

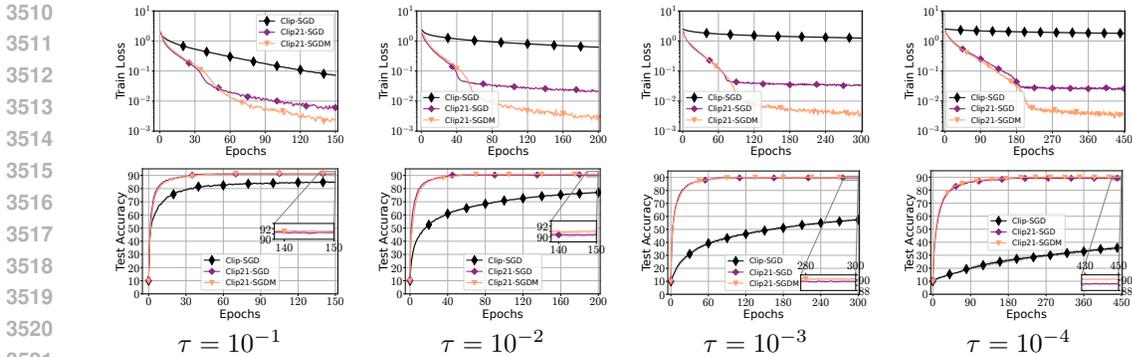


Figure 8: Comparison of Clip-SGD, Clip21-SGD, and Clip21-SGDM on training VGG16 model on CIFAR10 dataset where the clipping is applied globally.

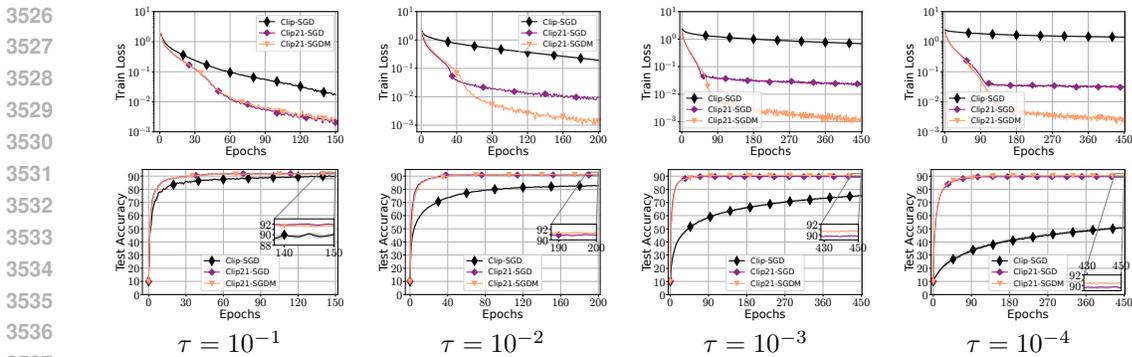


Figure 9: Comparison of Clip-SGD, Clip21-SGD, and Clip21-SGDM on training VGG16 model on CIFAR10 dataset where the clipping is applied layer-wise.

size and the clipping radius τ where the step-size γ is taken from $\{2^{-10}, \dots, 2^0\}$ and τ — from $\{10^{-4}, \dots, 10^0\}$, and choose the pair (γ, τ) that gives the smallest final gradient norm averaged over 3 runs. For Clip21-SGDM we perform the same grid search with an additional tuning of the momentum parameter $\beta \in \{0.1, 0.5, 0.9\}$. We report the last final gradient norm reached by each algorithm averaged over 3 runs.

H.2 EXPERIMENTS WITH NEURAL NETWORKS

H.2.1 VARYING CLIPPING RADIUS τ

Now we switch to the training of Resnet20 and VGG16 models on CIFAR10 dataset. For all algorithms, we do not use any techniques such as learning rate schedule, warm-up, or weight decay. However, we do tuning of the learning rate for Clip-SGD and Clip21-SGD from $\{10^{-3}, 10^{-2}, 10^{-1}, 10^0\}$ and choose the one that gives the highest test accuracy. For Clip21-SGDM we tune both the learning rate from $\{10^{-3}, 10^{-2}, 10^{-1}, 10^0\}$ and the momentum parameter from $\{0.1, 0.5, 0.9\}$ and choose the pair that reaches the highest test accuracy. The batch size for all algorithms is set to 32. We compare the performance of algorithms in two cases: when the clipping is applied globally on the whole model and layer-wise.

We observe in Figures 8 to 11 that the performance of Clip-SGD gets worsen once the clipping radius is small enough. For Clip21-SGDM is more robust to the choice of τ and can achieve smaller train loss and test accuracy even when τ is small.

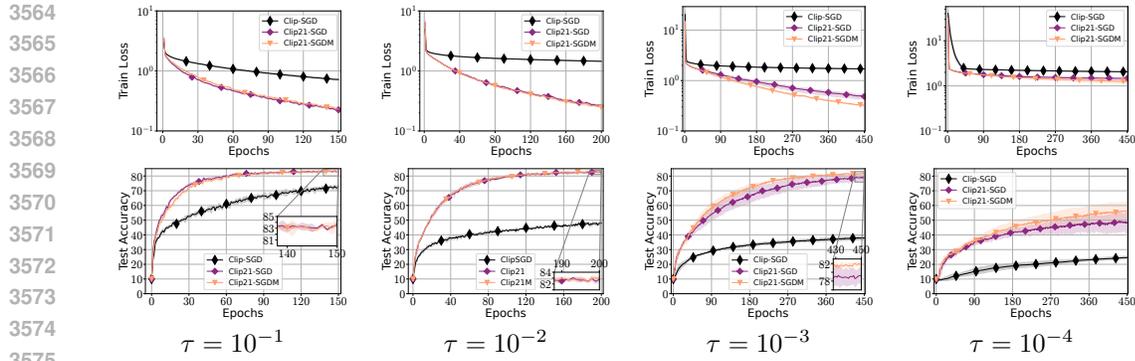


Figure 10: Comparison of Clip-SGD, Clip21-SGD, and Clip21-SGDM on training Resnet20 model on CIFAR10 dataset where the clipping is applied globally.

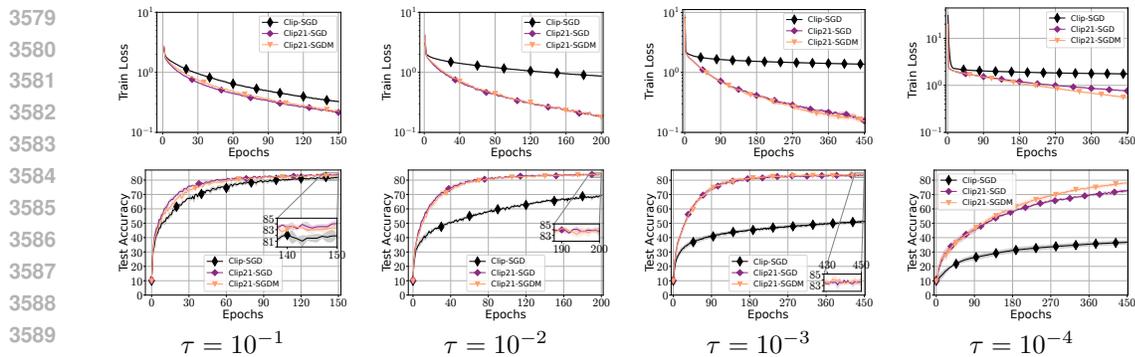


Figure 11: Comparison of Clip-SGD, Clip21-SGD, and Clip21-SGDM on training Resnet20 model on CIFAR10 dataset where the clipping is applied layer-wise.

H.2.2 ADDING ADDITIVE DP NOISE

We consider the same experiment section described in Appendix H.1.2 but now in the training of MLP and CNN models on MNIST dataset.

We use MLP model with 1 hidden layer of size 256 and Tanh activation function. CNN model has 2 convolution layers with 16 convolutions each and kernel size 5 with one max-pooling layer and Tanh activation function. We perform a grid search over the learning rate from $\{10^{-3}, \dots, 10^0\}$ and the clipping radius from $\{10^{-4}, \dots, 10^{-1}\}$. The aforementioned tuning is performed for each value of the noise-clipping ratio from $\{0.1, 0.3, 1.0, 3.0, 10.0\}$. The momentum parameter is tuned over $\{0.5, 0.1, 0.01\}$. We highlight that we do not use the techniques such as a learning rate scheduler although it might improve the performance of algorithms. The batch size for all algorithms is set to 64.

In Figures 12 to 15 we demonstrate that Clip-SGD and Clip-SGDM always outperform Clip21-SGD. However, there is no clear separation between Clip21-SGDM and Clip-SGD: in some cases, the latter has better performance, and in some cases — the former.

3618
 3619
 3620
 3621
 3622
 3623
 3624
 3625
 3626
 3627
 3628
 3629
 3630
 3631
 3632
 3633
 3634
 3635
 3636
 3637
 3638
 3639
 3640
 3641
 3642
 3643
 3644
 3645
 3646
 3647
 3648
 3649
 3650
 3651
 3652
 3653
 3654
 3655
 3656
 3657
 3658
 3659
 3660
 3661
 3662
 3663
 3664
 3665
 3666
 3667
 3668
 3669
 3670
 3671

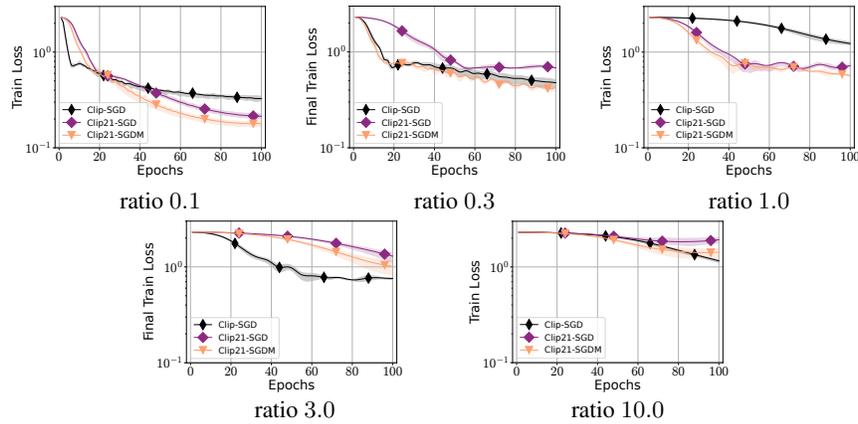


Figure 12: Comparison of Clip-SGD, Clip21-SGD, and Clip21-SGDM on training CNN model on MNIST dataset varying the noise-clipping ratio.

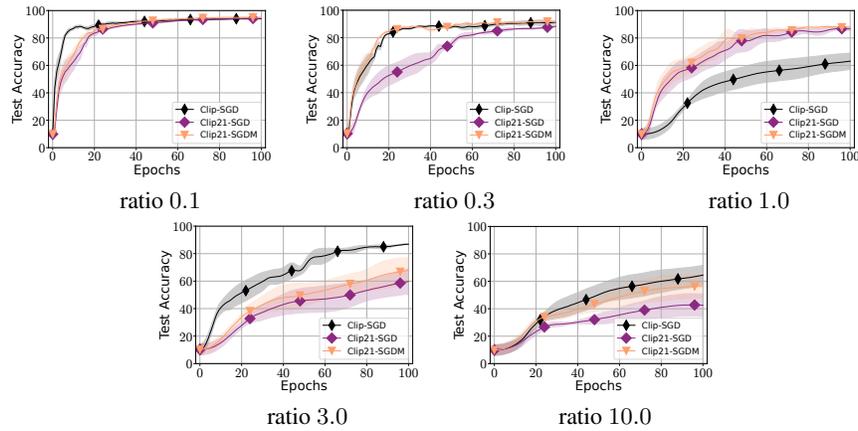


Figure 13: Comparison of Clip-SGD, Clip21-SGD, and Clip21-SGDM on training CNN model on MNIST dataset varying the noise-clipping ratio.

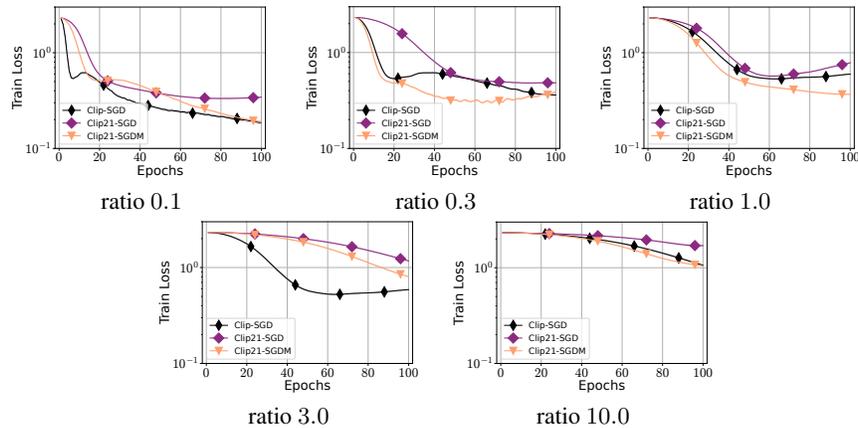


Figure 14: Comparison of Clip-SGD, Clip21-SGD, and Clip21-SGDM on training MLP model on MNIST dataset varying the noise-clipping ratio.

3672
 3673
 3674
 3675
 3676
 3677
 3678
 3679
 3680
 3681
 3682
 3683
 3684
 3685
 3686
 3687
 3688
 3689
 3690
 3691
 3692
 3693
 3694
 3695
 3696
 3697
 3698
 3699
 3700
 3701
 3702
 3703
 3704
 3705
 3706
 3707
 3708
 3709
 3710
 3711
 3712
 3713
 3714
 3715
 3716
 3717
 3718
 3719
 3720
 3721
 3722
 3723
 3724
 3725

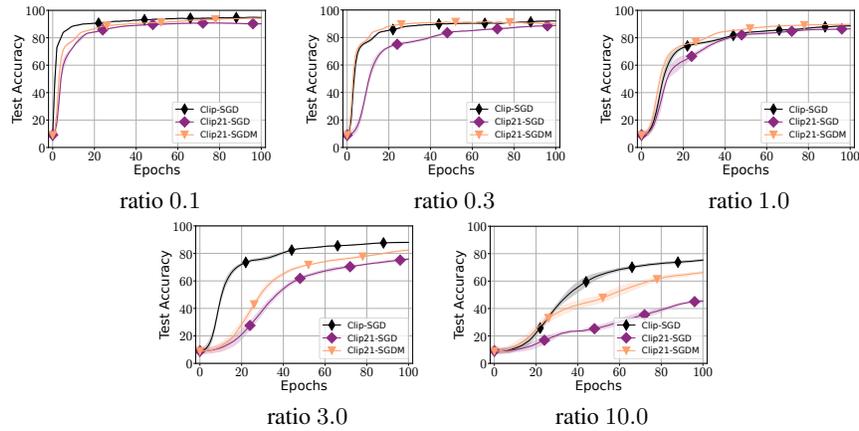


Figure 15: Comparison of Clip-SGD, Clip21-SGD, and Clip21-SGDM on training MLP model on MNIST dataset varying the noise-clipping ratio.