# Gradient Clipping Helps in Non-Smooth Stochastic Optimization with Heavy-Tailed Noise

**Anonymous Author(s)**
Affiliation
Address
email

## Abstract

Thanks to their practical efficiency and random nature of the data, stochastic first-order methods are standard for training large-scale machine learning models. Random behavior may cause a particular run of an algorithm to result in a highly suboptimal objective value, whereas theoretical guarantees are usually proved for the expectation of the objective value. Thus, it is essential to theoretically guarantee that algorithms provide small objective residual with high probability. Existing methods for non-smooth stochastic convex optimization have complexity bounds with the dependence on the confidence level that is either negative-power or logarithmic but under an additional assumption of sub-Gaussian (light-tailed) noise distribution that may not hold in practice, e.g., in several NLP tasks. In our paper, we resolve this issue and derive the first high-probability convergence results with logarithmic dependence on the confidence level for non-smooth convex stochastic optimization problems with non-sub-Gaussian (heavy-tailed) noise. To derive our results, we propose novel stepsize rules for two stochastic methods with gradient clipping. Moreover, our analysis works for generalized smooth objectives with Hölder-continuous gradients, and for both methods, we provide an extension for strongly convex problems. Finally, our results imply that the first (accelerated) method we consider also has optimal iteration and oracle complexity in all the regimes, and the second one is optimal in the non-smooth setting.

## 1 Introduction

Stochastic first-order optimization methods like SGD [33], Adam [21], and their various modifications are extremely popular in solving a number of different optimization problems, especially those appearing in statistics [37], machine learning, and deep learning [14]. The success of these methods in real-world applications motivates the researchers to investigate theoretical properties for the methods and to develop new ones with better convergence guarantees. Typically, stochastic methods are analyzed in terms of the convergence in expectation (see [13, 25, 16] and references therein), whereas high-probability complexity results are established much rarely. However, as illustrated in [15], guarantees in terms of the convergence in expectation have much worse correlation with the real behavior of the methods than high-probability convergence guarantees when the noise in the stochastic gradients has *heavy-tailed distribution*.

Recent studies [36, 35, 42] show that in several popular problems such as training BERT [38] on `Wikipedia` dataset the noise in the stochastic gradients is heavy-tailed. Moreover, in [42], the authors justify empirically that in such cases SGD works significantly worse than clipped-SGD [31] and Adam. Therefore, it is important to theoretically study the methods' convergence when the noise is heavy-tailed.

For convex and strongly convex problems with Lipschitz continuous gradient, i.e., smooth convex and strongly convex problems, this question was properly addressed in [26, 3, 15] where the first high-probability complexity bounds with logarithmic dependence on the confidence level were derived for the stochastic problems with heavy-tailed noise. However, a number of practically important problems are non-smooth *on the whole space* [41, 23]. For example, in deep neural network training, the loss function often grows polynomially fast when the norm of the network's weights goes to infinity. Moreover, non-smoothness of the activation functions such as ReLU or loss functions such as hinge loss implies the non-smoothness of the whole problem. While being well-motivated by practical applications, the existing high-probability convergence guarantees for stochastic first-order methods applied to solve non-smooth convex optimization problems with heavy-tailed noise depend on the negative power of the confidence level that dramatically increases the number of iterations required to obtain high accuracy of the solution with probability close to one. Such a discrepancy in the theory between algorithms for stochastic smooth and non-smooth problems leads us to the natural question: *is it possible to obtain high-probability complexity bounds with logarithmic dependence on the confidence level for **non-smooth** convex stochastic problems with heavy-tailed noise?* In this paper, we give a positive answer to this question. To achieve this we focus on gradient clipping methods [31, 11, 24, 23, 41, 42].

## 1.1 Preliminaries

Before we describe our contributions in detail, we formally state the considered setup.

**Stochastic optimization.** We focus on the following problem

$$\min_{x \in \mathbb{R}^n} f(x), \quad f(x) = \mathbb{E}_\xi \left[ f(x, \xi) \right], \tag{1}$$

where $f(x)$ is a convex but possibly non-smooth function. Next, we assume that at each point $x \in \mathbb{R}^n$ we have an access to the unbiased estimator $\nabla f(x, \xi)$ of $\nabla f(x)$ with uniformly bounded variance

$$\mathbb{E}_\xi[\nabla f(x, \xi)] = \nabla f(x), \quad \mathbb{E}_\xi \left[ \|\nabla f(x, \xi) - \nabla f(x)\|_2^2 \right] \leq \sigma^2, \quad \sigma > 0. \tag{2}$$

This assumption on the stochastic oracle is widely used in stochastic optimization literature [12, 13, 20, 22, 27]. We emphasize that we do not assume that the stochastic gradients have so-called "light tails" [22], i.e., sub-Gaussian noise distribution meaning that $\mathbb{P}\{\|\nabla f(x, \xi) - \nabla f(x)\|_2 > b\} \leq 2\exp(-b^2/(2\sigma^2))$ for all $b > 0$.

**Level of smoothness.** Finally, we assume that function $f$ has $(\nu, M_\nu)$-Hölder continuous gradients on a compact set $Q \subseteq \mathbb{R}^n$ for some $\nu \in [0, 1]$, $M_\nu > 0$ meaning that

$$\|\nabla f(x) - \nabla f(y)\|_2 \leq M_\nu \|x - y\|_2^\nu \quad \forall x, y \in Q. \tag{3}$$

When $\nu = 1$ inequality (3) implies $M_1$-smoothness of $f$, and when $\nu = 0$ we have that $\nabla f(x)$ has bounded variation which is equivalent to being uniformly bounded. Moreover, when $\nu = 0$ differentiability of $f$ is not needed, and one can assume uniform boundedness of the subgradients of $f$. Linear regression in the case when the noise has generalized Gaussian distribution (Example 4.4 from [2]) serves as a natural example of the situation with $\nu \in (0, 1)$. Moreover, when (3) holds for $\nu = 0$ and $\nu = 1$ simultaneously then it holds for all $\nu \in [0, 1]$ with $M_\nu \leq M_0^{1-\nu} M_1^\nu$ [29]. As we show in our results, the set $Q$ should contain the ball centered at the solution $x^*$ of (1) with radius $2R_0 = 2\|x^0 - x^*\|_2$, where $x^0$ is a starting point of the method, i.e., our analysis does not require (3) to hold on $\mathbb{R}^n$.

**High-probability convergence.** For a given accuracy $\varepsilon > 0$ and confidence level $\beta \in (0, 1)$ we are interested in finding $\varepsilon$-solutions of problem (1) with probability at least $1 - \beta$, i.e., such $\widehat{x}$ that $\mathbb{P}\{f(\widehat{x}) - f(x^*) \leq \varepsilon\} \geq 1 - \beta$. For brevity, we will call such (in general, random) points $\widehat{x}$ as $(\varepsilon, \beta)$-solution of (1). Moreover, by high-probability complexity of a stochastic method $\mathcal{M}$ we mean the sufficient number of oracle calls, i.e., number of $\nabla f(x, \xi)$ computations, needed to guarantee that the output of $\mathcal{M}$ is an $(\varepsilon, \beta)$-solution of (1).

Table 1: Summary of known and new high-probability complexity bounds for solving (1) with $f$ being **convex** and having $(\nu, M_\nu)$-Hölder continuous gradients. Columns: "Ref." = reference, "Complexity" = high-probability complexity ($\varepsilon$ – accuracy, $\beta$ – confidence level, numerical constants and logarithmic factors are omitted), "HT" = heavy-tailed noise, "UD" = unbounded domain, "HCC" = Hölder continuity of the gradient is required only on the compact set. The results labeled by ♣ are obtained from the convergence guarantees in expectation via Markov's inequality. Negative-power dependencies on the confidence level $\beta$ are colored in red.

| Method | Ref. | Complexity | $\nu$ | HT? | UD? | HCC? |
|---|---|---|---|---|---|---|
| SGD | [27] | $\max\left\{\frac{M_0^2 R_0^2}{\varepsilon^2}, \frac{\sigma^2 R_0^2}{\varepsilon^2}\right\}$ | 0 | ✗ | ✓ | ✗ |
| AC-SA | [12, 22] | $\max\left\{\sqrt{\frac{M_1 R_0^2}{\varepsilon}}, \frac{\sigma^2 R_0^2}{\varepsilon^2}\right\}$ | 1 | ✗ | ✓ | ✗ |
| SIGMA | [6] | $\max\left\{\frac{M_\nu^{\frac{2}{1+3\nu}} R_0^{\frac{2(1+\nu)}{1+3\nu}}}{\varepsilon^{\frac{2}{1+3\nu}}}, \frac{\sigma^2 R_0^2}{\varepsilon^2}\right\}$ | $[0,1]$ | ✗ | ✓ | ✗ |
| SGD | [27]♣ | $\max\left\{\frac{M_0^2 R_0^2}{\beta^2 \varepsilon^2}, \frac{\sigma^2 R_0^2}{\beta^2 \varepsilon^2}\right\}$ | 0 | ✓ | ✗ | ✗ |
| AC-SA | [12, 22]♣ | $\max\left\{\sqrt{\frac{M_1 R_0^2}{\beta\varepsilon}}, \frac{\sigma^2 R_0^2}{\beta^2 \varepsilon^2}\right\}$ | 1 | ✓ | ✓ | ✗ |
| SIGMA | [6]♣ | $\max\left\{\frac{M_\nu^{\frac{2}{1+3\nu}} R_0^{\frac{2(1+\nu)}{1+3\nu}}}{\beta^{\frac{2}{1+3\nu}} \varepsilon^{\frac{2}{1+3\nu}}}, \frac{\sigma^2 R_0^2}{\beta^2 \varepsilon^2}\right\}$ | $[0,1]$ | ✓ | ✓ | ✗ |
| clipped-SSTM | [15] | $\max\left\{\sqrt{\frac{M_1 R_0^2}{\varepsilon}}, \frac{\sigma^2 R_0^2}{\varepsilon^2}\right\}$ | 1 | ✓ | ✓ | ✗ |
| clipped-SGD | [15] | $\max\left\{\frac{M_1 R_0^2}{\varepsilon}, \frac{\sigma^2 R_0^2}{\varepsilon^2}\right\}$ | 1 | ✓ | ✓ | ✗ |
| clipped-SSTM | Thm. 2.2 | $\max\left\{\frac{M_\nu^{\frac{2}{1+3\nu}} R_0^{\frac{2(1+\nu)}{1+3\nu}}}{\varepsilon^{\frac{2}{1+3\nu}}}, \frac{\sigma^2 R_0^2}{\varepsilon^2}\right\}$ | $[0,1]$ | ✓ | ✓ | ✓ |
| clipped-SGD | Thm. 3.1 | $\max\left\{\frac{M_\nu^{\frac{2}{1+\nu}} R_0^2}{\varepsilon^{\frac{2}{1+\nu}}}, \frac{\sigma^2 R_0^2}{\varepsilon^2}\right\}$ | $[0,1]$ | ✓ | ✓ | ✓ |

## 1.2 Contributions

- We propose novel stepsize rules for clipped-SSTM [15] to handle the problems with Hölder continuous gradients and derive high-probability complexity guarantees for convex stochastic optimization problems without using "light tails" assumption, i.e., we prove that our version of clipped-SSTM

$$\mathcal{O}\left(\max\left\{D \ln^{\frac{2(1+\nu)}{1+3\nu}}\frac{D}{\beta}, \frac{\sigma^2 R_0^2}{\varepsilon^2}\ln\frac{D}{\beta}\right\}\right), \quad D = \frac{M_\nu^{\frac{2}{1+3\nu}} R_0^{\frac{2(1+\nu)}{1+3\nu}}}{\varepsilon^{\frac{2}{1+3\nu}}}$$

high-probability complexity. Unlike all previous high-probability complexity results in this setup with $\nu < 1$ (see Tbl. 1), our result depends only logarithmically on the confidence level $\beta$ that is highly important when $\beta$ is small. Moreover, up to the difference in logarithmic factors the derived complexity guarantees meet the known lower bounds [22, 18] obtained for the problems with light-tailed noise. In particular, when $\nu = 1$ we recover accelerated convergence rate [30, 22]. That is, neglecting the logarithmic factors our results are unimprovable and, surprisingly coincide with the best-known results in the "light-tailed case".

- We derive the first high-probability complexity bounds for clipped-SGD when the objective functions is convex with $(\nu, M_\nu)$-Hölder continuous gradient and the noise is heavy tailed., i.e., we derive

$$\mathcal{O}\left(\max\left\{D^2, \max\left\{D^{1+\nu}, \frac{\sigma^2 R_0^2}{\varepsilon^2}\right\}\ln\frac{D^2 + D^{1+\nu}}{\beta}\right\}\right), \quad D = \frac{M_\nu^{\frac{1}{1+\nu}} R_0}{\varepsilon^{\frac{1}{1+\nu}}}$$

high-probability complexity bound. Interestingly, when $\nu = 0$ the derived bound for clipped-SGD has better dependence on the logarithms than the corresponding one for clipped-SSTM. Moreover, neglecting the dependence on $\varepsilon$ under the logarithm, our bound for clipped-SGD has the same

Table 2: Summary of known and new high-probability complexity bounds for solving (1) with $f$ being $\mu$-**strongly convex** and having $(\nu, M_\nu)$-Hölder continuous gradients. Columns: "Ref." = reference, "Complexity" = high-probability complexity ($\varepsilon$ – accuracy, $\beta$ – confidence level, numerical constants and logarithmic factors are omitted), "HT" = heavy-tailed noise, "UD" = unbounded domain, "HCC" = Hölder continuity of the gradient is required only on the compact set. The results labeled by ♣ are obtained from the convergence guarantees in expectation via Markov's inequality. Negative-power dependencies on the confidence level $\beta$ are colored in red.

| Method | Ref. | Complexity | $\nu$ | HT? | UD? | HCC? |
|---|---|---|---|---|---|---|
| SGD | [27] | $\max\left\{\frac{M_0^2}{\mu\varepsilon}, \frac{\sigma^2}{\mu\varepsilon}\right\}$ | 0 | ✗ | ✓ | ✗ |
| AC-SA | [12, 22] | $\max\left\{\sqrt{\frac{M_1}{\mu}}, \frac{\sigma^2}{\mu\varepsilon}\right\}$ | 1 | ✗ | ✓ | ✗ |
| SIGMA | [6] | $\max\left\{\hat{N}, \frac{\sigma^2}{\mu\varepsilon}\right\},$ $\hat{N} = \left(\frac{M_\nu}{\mu R_0^{1-\nu}}\right)^{\frac{2}{1+3\nu}} + \left(\frac{M_\nu^2}{\mu^{1+\nu}\varepsilon^{1-\nu}}\right)^{\frac{1}{1+3\nu}}$ | $[0,1]$ | ✗ | ✓ | ✗ |
| SGD | [27]♣ | $\max\left\{\frac{M_0^2}{\mu\beta\varepsilon}, \frac{\sigma^2}{\mu\beta\varepsilon}\right\}$ | 0 | ✓ | ✗ | ✗ |
| AC-SA | [12, 22]♣ | $\max\left\{\sqrt{\frac{M_1}{\mu}}, \frac{\sigma^2}{\mu\beta\varepsilon}\right\}$ | 1 | ✓ | ✓ | ✗ |
| SIGMA | [6]♣ | $\max\left\{\hat{N}, \frac{\sigma^2}{\mu\hat{\varepsilon}}\right\}, \hat{\varepsilon} = \beta\varepsilon,$ $\hat{N} = \left(\frac{M_\nu}{\mu R_0^{1-\nu}}\right)^{\frac{2}{1+3\nu}} + \left(\frac{M_\nu^2}{\mu^{1+\nu}\hat{\varepsilon}^{1-\nu}}\right)^{\frac{1}{1+3\nu}}$ | $[0,1]$ | ✓ | ✓ | ✗ |
| R-clipped-SSTM | [15] | $\max\left\{\sqrt{\frac{M_1}{\mu}}, \frac{\sigma^2}{\mu\varepsilon^2}\right\}$ | 1 | ✓ | ✓ | ✗ |
| R-clipped-SGD | [15] | $\max\left\{\frac{M_1}{\mu}, \frac{\sigma^2}{\mu\varepsilon^2}\right\}$ | 1 | ✓ | ✓ | ✗ |
| R-clipped-SSTM | Thm. 2.1 | $\max\left\{\hat{N}, \frac{\sigma^2}{\mu\varepsilon}\right\},$ $\hat{N} = \left(\frac{M_\nu}{\mu R_0^{1-\nu}}\right)^{\frac{2}{1+3\nu}} + \left(\frac{M_\nu^2}{\mu^{1+\nu}\varepsilon^{1-\nu}}\right)^{\frac{1}{1+3\nu}}$ | $[0,1]$ | ✓ | ✓ | ✓ |
| R-clipped-SGD | Thm. 3.2 | $\max\left\{\frac{M_\nu^{\frac{2}{1+\nu}}}{\mu^{\frac{2}{1+\nu}} R_0^{\frac{2(1-\nu)}{1+\nu}}}, \frac{M_\nu^{\frac{2}{1+\nu}}}{\mu\varepsilon^{\frac{1-\nu}{1+\nu}}}, \frac{\sigma^2}{\mu\varepsilon}\right\}$ | $[0,1]$ | ✓ | ✓ | ✓ |

dependence on the confidence level as the tightest known result in this case under the "light tails" assumption [17].

- Using restarts technique we extend the obtained results for clipped-SSTM and clipped-SGD to the strongly convex case (see Tbl. 2). As in the convex case, the obtained results are superior to all previous known results in the general setup we consider.

- As one of the key contributions of this work, we emphasize that in our theoretical results it is sufficient to assume Hölder continuity of the gradients of $f$ only on the ball with radius $2R_0 = 2\|x^0 - x^*\|_2$ and centered at a solution of the problem. This makes our results applicable to much larger class of problems than functions with Hölder continuous gradients on $\mathbb{R}^n$, e.g., our analysis works even for polynomially growing objectives.

- To test the performance of the considered methods we conduct several numerical experiments on image classification and NLP tasks, and observe that 1) clipped-SSTM and clipped-SGD show a comparable performance with SGD on the image classification task, when the noise distribution is almost sub-Gaussian, 2) converge much faster than SGD on the NLP task, when the noise distribution is heavy-tailed, and 3) clipped-SSTM achieves a comparable performance with Adam on the NLP task enjoying both the best known theoretical guarantees and good practical performance.

## 1.3 Related work

**Light-tailed noise.** The theory of high-probability complexity bounds for convex stochastic optimization with light-tailed noise is well-developed. Lower bounds and optimal methods for the problems with $(\nu, M_\nu)$-Hölder continuous gradients are obtained in [27] for $\nu = 0$, and in [12] for $\nu = 1$. Up to the logarithmic dependencies these high-probability convergence bounds coincide with

4

119 the corresponding results for the convergence in expectation (see first two rows of Tbl. 1) While not
120 being directly derived in the literature, the lower bound for the case when $\nu \in (0, 1)$ can be obtained
121 as a combination of lower bounds in the deterministic [28, 18] and smooth stochastic settings [12].
122 The corresponding optimal methods are analyzed in [4, 6] through the lens of inexact oracle.

123 **Heavy-tailed noise.** Unlike in the "light-tailed" case, the first theoretical guarantees with reasonable
124 dependence on both the accuracy $\varepsilon$ and the confidence level $\beta$ appeared just recently. In [26], the
125 first such results without acceleration [30] were derived for Mirror Descent with special truncation
126 technique for smooth ($\nu = 1$) convex problems on a bounded domain, and then were accelerated and
127 extended in [15]. For the strongly convex problems the first accelerated high-probability convergence
128 guarantees were obtained in [3] for the special method called proxBoost requiring solving auxiliary
129 nontrivial problems at each iteration. These bounds were tightened in [15].

130 In contrast, for the case when $\nu < 1$ and, in particular, when $\nu = 0$ the best-known high-probability
131 complexity bounds suffer from the negative-power dependence on the confidence level $\beta$, i.e., have
132 a factor $1/\beta^\alpha$ for some $\alpha > 0$, that affects the convergence rate dramatically for small enough
133 $\beta$. Without additional assumptions on the tails these results are obtained via Markov's inequality
134 $\mathbb{P}\{f(x) - f(x^*) > \varepsilon\} < \mathbb{E}[f(x)-f(x^*)]/\varepsilon$ from the guarantees for the convergence in expectation to
135 the accuracy $\varepsilon\beta$, see the results labeled by ♣ in Tbl. 1. Under an additional assumption on noise
136 tails that $\mathbb{P}\{\|\nabla f(x,\xi) - \nabla f(x)\|_2^2 > s\sigma^2\} = O(s^{-\alpha})$ for $\alpha > 2$ these results can be tightened [10]
137 when $\nu = 0$ as $O\left(M_0^2 R_0^2 \max\left\{\ln(\beta^{-1})/\varepsilon^2, (1/\beta\varepsilon^\alpha)^{2/(3\alpha-2)}\right\}\right)$ without removing the negative-power
138 dependence on the confidence level $\beta$. Different stepsize policies allow to change the last term in
139 max to $\beta^{-\frac{1}{2\alpha-1}}\varepsilon^{-\frac{2\alpha}{2\alpha-1}}$ without removing the negative-power dependence on $\beta$.

140 **Gradient clipping.** The methods based on gradient clipping [31] and normalization [19] are popular
141 in different machine learning and deep learning tasks due to their robustness in practice to the noise
142 in the stochastic gradients and rapid changes of the objective function [14]. In [41, 23], clipped-GD
143 and clipped-SGD are theoretically studied in applications to non-smooth problems that can grow
144 polynomially fast when $\|x - x^*\|_2 \to \infty$ showing the superiority of gradient clipping methods
145 to the methods without clipping. The results from [41] are obtained for non-convex problems
146 with almost surely bounded noise, and in [23], the authors derive the stability and expectation
147 convergence guarantees for strongly convex under assumption that the central $p$-th moment of the
148 stochastic gradient is bounded for $p \geq 2$. Since the authors of [23] do not provide convergence
149 guarantees with explicit dependencies on all important parameters of the problem it complicates direct
150 comparison with our results. Nevertheless, convergence guarantees from [23] are sub-linear and are
151 given for the convergence in expectation, and, as a consequence, the corresponding high-probability
152 convergence results obtained via Markov's inequality also suffer from negative-power dependence on
153 the confidence level. Next, in [42], the authors establish several expectation convergence guarantees
154 for clipped-SGD and prove their optimality in the non-convex case under assumption that the central
155 $\alpha$-moment of the stochastic gradient is uniformly bounded, where $\alpha \in (1, 2]$. It turns out that
156 clipped-SGD is able to converge even when $\alpha < 2$, whereas vanilla SGD can diverge in this setting.

## 2 Clipped Stochastic Similar Triangles Method

158 In this section, we propose a novel variation of Clipped Stochastic Similar Triangles Method [15]
159 adjusted to the class of objectives with Hölder continuous gradients (clipped-SSTM, see Alg. 1).

160 The method is based on the clipping of the stochastic gradients:

$$\text{clip}(\nabla f(x,\boldsymbol{\xi}),\lambda) = \min\left\{1, \frac{\lambda}{\|\nabla f(x,\boldsymbol{\xi})\|_2}\right\}\nabla f(x,\boldsymbol{\xi}) \tag{4}$$

161 where $\nabla f(x,\boldsymbol{\xi}) = \frac{1}{m}\sum_{i=1}^m \nabla f(x,\xi_i)$ is a mini-batched stochastic gradient. Gradient clipping
162 ensures that the resulting vector has a norm bounded by the clipping level $\lambda$. Since the clipped
163 stochastic gradient cannot have arbitrary large norm, the clipping helps to avoid unstable behavior of
164 the method when the noise is heavy-tailed and the clipping level $\lambda$ is properly adjusted.

165 However, unlike the stochastic gradient, clipped stochastic gradient is a *biased* estimate of $\nabla f(x)$:
166 the smaller the clipping level the larger the bias. The biasedness of the clipped stochastic gradient

---

**Algorithm 1** Clipped Stochastic Similar Triangles Method (clipped-SSTM): case $\nu \in [0, 1]$

---

**Input:** starting point $x^0$, number of iterations $N$, batchsizes $\{m_k\}_{k=1}^N$, stepsize parameter $\alpha$, clipping parameter $B$, Hölder exponent $\nu \in [0, 1]$.

1: Set $A_0 = \alpha_0 = 0$, $y^0 = z^0 = x^0$
2: **for** $k = 0, \ldots, N - 1$ **do**
3:     Set $\alpha_{k+1} = \alpha(k + 1)^{\frac{2\nu}{1+\nu}}$, $A_{k+1} = A_k + \alpha_{k+1}$, $\lambda_{k+1} = \frac{B}{\alpha_{k+1}}$
4:     $x^{k+1} = (A_k y^k + \alpha_{k+1} z^k)/A_{k+1}$
5:     Draw mini-batch $m_k$ of fresh i.i.d. samples $\xi_1^k, \ldots, \xi_{m_k}^k$ and compute $\nabla f(x^{k+1}, \boldsymbol{\xi}^k) = \frac{1}{m_k} \sum_{i=1}^{m_k} \nabla f(x^{k+1}, \xi_i^k)$
6:     Compute $\widetilde{\nabla} f(x^{k+1}, \boldsymbol{\xi}^k) = \text{clip}(\nabla f(x^{k+1}, \boldsymbol{\xi}^k), \lambda_{k+1})$ using (4)
7:     $z^{k+1} = z^k - \alpha_{k+1} \widetilde{\nabla} f(x^{k+1}, \boldsymbol{\xi}^k)$
8:     $y^{k+1} = (A_k y^k + \alpha_{k+1} z^{k+1})/A_{k+1}$
9: **end for**
**Output:** $y^N$

---

complicates the analysis of the method. On the other hand, to circumvent the negative effect of the heavy-tailed noise on the high-probability convergence one should choose $\lambda$ to be not too large. Therefore, the question on the appropriate choice of the clipping level is highly non-trivial.

Fortunately, there exists a simple but insightful observation that helps us to obtain the right formula for the clipping level $\lambda_k$ in clipped-SSTM: if $\lambda_k$ is chosen in such a way that $\|\nabla f(x^k)\|_2 \leq \lambda_k/2$ with high probability, then for the realizations $\nabla f(x^{k+1}, \boldsymbol{\xi}^k)$ of the mini-batched stochastic gradient such that $\|\nabla f(x^{k+1}, \boldsymbol{\xi}^k) - \nabla f(x^{k+1})\|_2 \leq \lambda_k/2$ the clipping is an identity operator. Next, if the probability mass of such realizations is big enough then the bias of the clipped stochastic gradient is properly bounded that helps to derive needed convergence guarantees. It turns out that the choice $\lambda_k \sim 1/\alpha_k$ ensures the method convergence with needed rate and high enough probability.

Guided by this observation we derive the precise expressions for all the parameters of clipped-SSTM and derive high-probability complexity bounds for the method. Below we provide a simplified version of the main result for clipped-SSTM in the convex case. The complete formulation and the full proof of the theorem are deferred to Appendix B.1 (see Thm. B.1).

**Theorem 2.1.** *Assume that function $f$ is convex and its gradient satisfy* (3) *with $\nu \in [0, 1]$, $M_\nu > 0$ on $Q = B_{2R_0} = \{x \in \mathbb{R}^n \mid \|x - x^*\|_2 \leq 2R_0\}$, where $R_0 \geq \|x^0 - x^*\|_2$. Then there exist such a choice of parameters that* clipped-SSTM *achieves $f(y^N) - f(x^*) \leq \varepsilon$ with probability at least $1 - \beta$ after $\mathcal{O}\left(D \ln^{\frac{2(1+\nu)}{1+3\nu}} \frac{D}{\beta}\right)$ iterations with $D = \frac{M_\nu^{\frac{2}{1+3\nu}} R_0^{\frac{2(1+\nu)}{1+3\nu}}}{\varepsilon^{\frac{2}{1+3\nu}}}$ and requires*

$$\mathcal{O}\left(\max\left\{D \ln^{\frac{2(1+\nu)}{1+3\nu}} \frac{D}{\beta}, \frac{\sigma^2 R_0^2}{\varepsilon^2} \ln \frac{D}{\beta}\right\}\right) \text{ oracle calls.} \tag{5}$$

The obtained result has only logarithmic dependence on the confidence level $\beta$ and optimal dependence on the accuracy $\varepsilon$ up to logarithmic factors [22, 18] for all $\nu \in [0, 1]$. Moreover, we emphasize that our result does not require $f$ to have $(\nu, M_\nu)$-Hölder continuous gradient on the whole space. This is because we prove that for the proposed choice of parameters the iterates of clipped-SSTM stay inside the ball $B_{2R_0} = \{x \in \mathbb{R}^n \mid \|x - x^*\|_2 \leq 2R_0\}$ with probability at least $1 - \beta$, and, as a consequence, Hölder continuity of the gradient is required only inside this ball. In particular, this means that the better starting point leads not only to the reduction of $R_0$, but also it can reduce $M_\nu$. Moreover, our result is applicable to much wider class of functions than the standard class of functions with Hölder continuous gradients in $\mathbb{R}^n$, e.g., to the problems with polynomial growth.

For the strongly convex problems, we consider restarted version of Alg. 1 (R-clipped-SSTM, see Alg. 2) and derive high-probability complexity result for this version. Below we provide a simplified version of the result. The complete formulation and the full proof of the theorem are deferred to Appendix B.2 (see Thm. B.2).

**Theorem 2.2.** *Assume that function $f$ is $\mu$-strongly convex and its gradient satisfy* (3) *with $\nu \in [0, 1]$, $M_\nu > 0$ on $Q = B_{2R_0} = \{x \in \mathbb{R}^n \mid \|x - x^*\|_2 \leq 2R_0\}$, where $R_0 \geq \|x^0 - x^*\|_2$. Then there exist*

---
**Algorithm 2** Restarted clipped-SSTM (R-clipped-SSTM): case $\nu \in [0, 1]$
---
**Input:** starting point $x^0$, number of restarts $\tau$, number of steps of clipped-SSTM in restarts $\{N_t\}_{t=1}^{\tau}$, batchsizes $\{m_k^1\}_{k=1}^{N_1-1}, \{m_k^2\}_{k=1}^{N_2-1}, \ldots, \{m_k^{\tau}\}_{k=1}^{N_{\tau}-1}$, stepsize parameters $\{\alpha^t\}_{t=1}^{\tau}$, clipping parameters $\{B_t\}_{t=1}^{\tau}$, Hölder exponent $\nu \in [0, 1]$.
  1: $\hat{x}^0 = x^0$
  2: **for** $t = 1, \ldots, \tau$ **do**
  3:    Run clipped-SSTM (Alg. 1) for $N_t$ iterations with batchsizes $\{m_k^t\}_{k=1}^{N_t-1}$, stepsize parameter $\alpha_t$, clipping parameter $B_t$, and starting point $\hat{x}^{t-1}$. Define the output of clipped-SSTM by $\hat{x}^t$.
  4: **end for**
**Output:** $\hat{x}^{\tau}$
---

such a choice of parameters that R-clipped-SSTM achieves $f(\hat{x}^{\tau}) - f(x^*) \leq \varepsilon$ with probability at least $1 - \beta$ after

$$\hat{N} = O\left(D \ln^{\frac{2(1+\nu)}{1+3\nu}} \frac{D}{\beta}\right), \quad D = \max\left\{\left(\frac{M_\nu}{\mu R_0^{1-\nu}}\right)^{\frac{2}{1+3\nu}} \ln \frac{\mu R_0^2}{\varepsilon}, \left(\frac{M_\nu^2}{\mu^{1+\nu}\varepsilon^{1-\nu}}\right)^{\frac{1}{1+3\nu}}\right\} \quad (6)$$

iterations of Alg. 1 in total and requires

$$O\left(\max\left\{D \ln^{\frac{2(1+\nu)}{1+3\nu}} \frac{D}{\beta}, \frac{\sigma^2}{\mu\varepsilon} \ln \frac{D}{\beta}\right\}\right) \text{ oracle calls.} \quad (7)$$

Again, the obtained result has only logarithmic dependence on the confidence level $\beta$ and, as our result in the convex case, it has optimal dependence on the accuracy $\varepsilon$ up to logarithmic factors depending on $\beta$ [22, 18] for all $\nu \in [0, 1]$.

## 3   SGD with clipping

In this section, we present a new variant of clipped-SGD [31] properly adjusted to the class of objectives with $(\nu, M_\nu)$-Hölder continuous gradients (see Alg. 3).

---
**Algorithm 3** Clipped Stochastic Gradient Descent (clipped-SGD): case $\nu \in [0, 1]$
---
**Input:** starting point $x^0$, number of iterations $N$, batchsize $m$, stepsize $\gamma$, clipping parameter $B > 0$.
  1: **for** $k = 0, \ldots, N - 1$ **do**
  2:    Draw mini-batch of $m$ fresh i.i.d. samples $\xi_1^k, \ldots, \xi_m^k$ and compute $\nabla f(x^{k+1}, \boldsymbol{\xi}^k) = \frac{1}{m}\sum_{i=1}^{m} \nabla f(x^{k+1}, \xi_i^k)$
  3:    Compute $\widetilde{\nabla} f(x^k, \boldsymbol{\xi}^k) = \text{clip}(\nabla f(x^k, \boldsymbol{\xi}^k), \lambda)$ using (4) with $\lambda = {}^B\!/\!\gamma$
  4:    $x^{k+1} = x^k - \gamma \widetilde{\nabla} f(x^k, \boldsymbol{\xi}^k)$
  5: **end for**
**Output:** $\bar{x}^N = \frac{1}{N}\sum_{k=0}^{N-1} x^k$
---

We emphasize that as for clipped-SSTM we use clipping level $\lambda$ inversely proportional to the stepsize $\gamma$. Below we provide a simplified version of the main result for clipped-SGD in the convex case. The complete formulation and the full proof of the theorem are deferred to Appendix C.1 (see Thm. C.1).

**Theorem 3.1.** *Assume that function $f$ is convex and its gradient satisfy (3) with $\nu \in [0, 1]$, $M_\nu > 0$ on $Q = B_{2R_0} = \{x \in \mathbb{R}^n \mid \|x - x^*\|_2 \leq 2R_0\}$, where $R_0 \geq \|x^0 - x^*\|_2$. Then there exist such a choice of parameters that clipped-SGD achieves $f(\bar{x}^N) - f(x^*) \leq \varepsilon$ with probability at least $1 - \beta$ after*

$$\mathcal{O}\left(\max\left\{D^2, D^{1+\nu} \ln \frac{D^2 + D^{1+\nu}}{\beta}\right\}\right), \quad D = \frac{M_\nu^{\frac{1}{1+\nu}} R_0}{\varepsilon^{\frac{1}{1+\nu}}} \quad (8)$$

*iterations and requires*

$$\mathcal{O}\left(\max\left\{D^2, \max\left\{D^{1+\nu}, \frac{\sigma^2 R_0^2}{\varepsilon^2}\right\} \ln \frac{D^2 + D^{1+\nu}}{\beta}\right\}\right) \text{ oracle calls.} \quad (9)$$

7

As all our results in the paper, this result for clipped-SGD has two important features: 1) the dependence on the confidence level $\beta$ is logarithmic and 2) Hölder continuity is required only on the ball $B_{2R_0}$ centered at the solution. Moreover, up to the difference in the expressions under the logarithm the dependence on $\varepsilon$ in the result for clipped-SGD is the same as in the tightest known results for non-accelerated SGD-type methods [4, 17]. Finally, we emphasize that for $\nu < 1$ the logarithmic factors appearing in the complexity bound for clipped-SSTM are worse than the corresponding factor in the complexity bound for clipped-SGD. Therefore, clipped-SGD has the best known high-probability complexity results in the case when $\nu = 0$ and $f$ is convex.

For the strongly convex problems, we consider restarted version of Alg. 3 (R-clipped-SGD, see Alg. 4) and derive high-probability complexity result for this version. Below we provide a simplified

---

**Algorithm 4** Restarted clipped-SGD (R-clipped-SGD): case $\nu \in [0, 1]$

---

**Input:** starting point $x^0$, number of restarts $\tau$, number of steps of clipped-SGD in restarts $\{N_t\}_{t=1}^\tau$, batchsizes $\{m_t\}_{k=1}^\tau$, stepsizes $\{\gamma_t\}_{t=1}^\tau$, clipping parameters $\{B_t\}_{t=1}^\tau$
  1: $\hat{x}^0 = x^0$
  2: **for** $t = 1, \ldots, \tau$ **do**
  3:     Run clipped-SGD (Alg. 3) for $N_t$ iterations with batchsize $m_t$, stepsize $\gamma_t$, clipping parameter $B_t$, and starting point $\hat{x}^{t-1}$. Define the output of clipped-SGD by $\hat{x}^t$.
  4: **end for**
**Output:** $\hat{x}^\tau$

---

version of the result. The complete formulation and the full proof of the theorem are deferred to Appendix C.2 (see Thm. C.2).

**Theorem 3.2.** *Assume that function $f$ is $\mu$-strongly convex and its gradient satisfy (3) with $\nu \in [0, 1]$, $M_\nu > 0$ on $Q = B_{2R_0} = \{x \in \mathbb{R}^n \mid \|x - x^*\|_2 \leq 2R_0\}$, where $R_0 \geq \|x^0 - x^*\|_2$. Then there exist such a choice of parameters that R-clipped-SGD achieves $f(\bar{x}^N) - f(x^*) \leq \varepsilon$ with probability at least $1 - \beta$ after*

$$\mathcal{O}\left(\max\left\{D_1^{\frac{2}{1+\nu}}\ln\frac{\mu R_0^2}{\varepsilon}, D_2^{\frac{2}{1+\nu}}, \max\left\{D_1 \ln\frac{\mu R_0^2}{\varepsilon}, D_2\right\}\ln\frac{D}{\beta}\right\}\right)$$

*iterations of Alg. 3 in total and requires*

$$\mathcal{O}\left(\max\left\{D_1^{\frac{2}{1+\nu}}\ln\frac{\mu R_0^2}{\varepsilon}, D_2^{\frac{2}{1+\nu}}, \max\left\{D_1 \ln\frac{\mu R_0^2}{\varepsilon}, D_2, \frac{\sigma^2}{\mu\varepsilon}\right\}\ln\frac{D}{\beta}\right\}\right) \text{ oracle calls, where}$$

$$D_1 = \frac{M_\nu}{\mu R_0^{1-\nu}}, \quad D_2 = \frac{M_\nu}{\mu^{\frac{1+\nu}{2}}\varepsilon^{\frac{1-\nu}{2}}}, \quad D = (D_1^{\frac{2}{1+\nu}} + D_1)\ln\frac{\mu R_0^2}{\varepsilon} + D_2 + D_2^{\frac{2}{1+\nu}}.$$

As in the convex case, for $\nu < 1$ the log factors appearing in the complexity bound for R-clipped-SSTM are worse than the corresponding factor in the bound for R-clipped-SGD. Thus, R-clipped-SGD has the best known high-probability complexity results for strongly convex $f$ and $\nu = 0$.

## 4  Numerical experiments

We tested the performance of the methods on the following problems:

- BERT fine-tuning on CoLA dataset [39]. We use pretrained BERT from Transformers library [40] (bert-base-uncased) and freeze all layers except the last two linear ones.

- ResNet-18 training on ImageNet-100 (first 100 classes of ImageNet [34]).

First, we study the noise distribution for both problem as follows: at the starting point we sample large enough number of batched stochastic gradients $\nabla f(x^0, \boldsymbol{\xi}_1), \ldots, \nabla f(x^0, \boldsymbol{\xi}_K)$ with batchsize 32 and plot the histograms for $\|\nabla f(x^0, \boldsymbol{\xi}_1) - \nabla f(x^0)\|_2, \ldots, \|\nabla f(x^0, \boldsymbol{\xi}_K) - \nabla f(x^0)\|_2$, see Fig. 1. As one can see, the noise distribution for BERT + CoLA is substantially non-sub-Gaussian, whereas the distribution for ResNet-18 + Imagenet-100 is almost Gaussian.

Next, we compared 4 different optimizers on these problems: Adam, SGD (with Momentum), clipped-SGD (with Momentum and coordinate-wise clipping) and clipped-SSTM (with norm-clipping and $\nu = 1$). The results are presented in Fig. 2. We observed that the noise distributions do
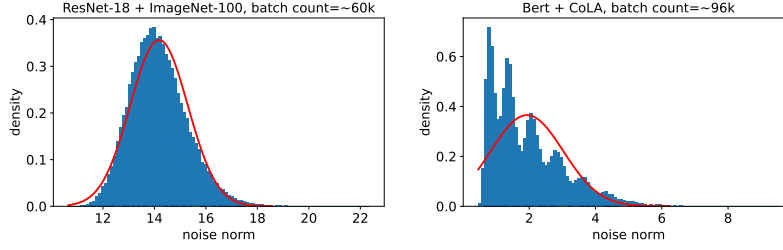
Figure 1: Noise distribution of the stochastic gradients for `ResNet-18` on `ImageNet-100` and BERT fine-tuning on the `CoLA` dataset before the training. Red lines: probability density functions with means and variances empirically estimated by the samples. Batch count is the total number of samples used to build a histogram.

not change significantly along the trajectories of the considered methods, see Appendix D. During the hyper-parameters search we compared different batchsizes, emulated via gradient accumulation (thus we compare methods with different batchsizes by the number of base batches used). The base batchsize was 32 for both problems, stepsizes and clipping levels were tuned. One can find additional details regarding our experiments in Appendix D.
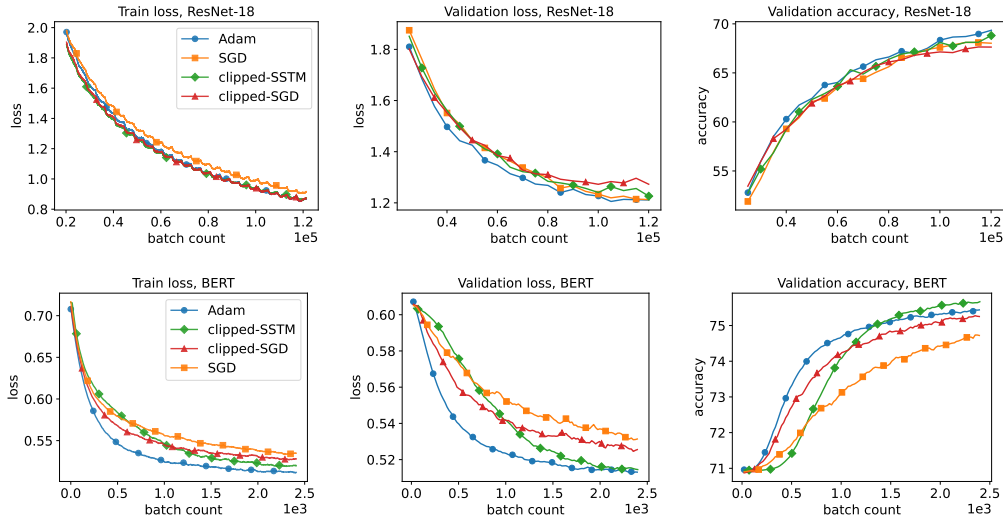


Figure 2: Train and validation loss + accuracy for different optimizers on both problems. Here, "batch count" denotes the total number of used stochastic gradients.

**Image classification.** On `ResNet-18` + `ImageNet-100` task, SGD performs relatively well, and even ties with Adam (with batchsize of $4 \times 32$) in validation loss. clipped-SSTM (with batchsize of $2 \times 32$) also ties with Adam and clipped-SGD is not far from them. The results were averaged from 5 different launches (with different starting points/weight initializations). Since the noise distribution is almost Gaussian even vanilla SGD performs well, i.e., gradient clipping is not required. At the same time, the clipping does not slow down the convergence significantly.

**Text classification.** On `BERT` + `CoLA` task, when the noise distribution is heavy-tailed, the methods with clipping outperform SGD by a large margin. This result is in good correspondence with the derived high-probability complexity bounds for clipped-SGD, clipped-SSTM and the best-known ones for SGD. Moreover, clipped-SSTM (with batchsize of $8 \times 32$) achieves the same loss on validation as Adam, and has better accuracy. These results were averaged from 5 different train-val splits and 20 launches (with different starting points/weight initializations) for each of the splits, 100 launches in total.

9

## References

[1] George Bennett. Probability inequalities for the sum of independent random variables. *Journal of the American Statistical Association*, 57(297):33–45, 1962.

[2] Caroline Chaux, Patrick L Combettes, Jean-Christophe Pesquet, and Valérie R Wajs. A variational formulation for frame-based inverse problems. *Inverse Problems*, 23(4):1495–1518, jun 2007.

[3] Damek Davis, Dmitriy Drusvyatskiy, Lin Xiao, and Junyu Zhang. From low probability to high confidence in stochastic convex optimization. *Journal of Machine Learning Research*, 22(49):1–38, 2021.

[4] Olivier Devolder. *Exactness, inexactness and stochasticity in first-order methods for large-scale convex optimization*. PhD thesis, PhD thesis, 2013.

[5] Olivier Devolder, François Glineur, and Yurii Nesterov. First-order methods of smooth convex optimization with inexact oracle. *Mathematical Programming*, 146(1):37–75, 2014.

[6] Pavel Dvurechensky and Alexander Gasnikov. Stochastic intermediate gradient method for convex problems with stochastic inexact oracle. *Journal of Optimization Theory and Applications*, 171(1):121–145, 2016.

[7] Kacha Dzhaparidze and JH Van Zanten. On bernstein-type inequalities for martingales. *Stochastic processes and their applications*, 93(1):109–117, 2001.

[8] David A Freedman et al. On tail probabilities for martingales. *the Annals of Probability*, 3(1):100–118, 1975.

[9] Alexander Gasnikov and Yurii Nesterov. Universal fast gradient method for stochastic composit optimization problems. *arXiv:1604.05275*, 2016.

[10] Alexander Vladimirovich Gasnikov, Yu E Nesterov, and Vladimir Grigor'evich Spokoiny. On the efficiency of a randomized mirror descent algorithm in online optimization problems. *Computational Mathematics and Mathematical Physics*, 55(4):580–596, 2015.

[11] Jonas Gehring, Michael Auli, David Grangier, Denis Yarats, and Yann N Dauphin. Convolutional sequence to sequence learning. In *Proceedings of the 34th International Conference on Machine Learning-Volume 70*, pages 1243–1252. JMLR. org, 2017.

[12] Saeed Ghadimi and Guanghui Lan. Optimal stochastic approximation algorithms for strongly convex stochastic composite optimization i: A generic algorithmic framework. *SIAM Journal on Optimization*, 22(4):1469–1492, 2012.

[13] Saeed Ghadimi and Guanghui Lan. Stochastic first-and zeroth-order methods for nonconvex stochastic programming. *SIAM Journal on Optimization*, 23(4):2341–2368, 2013.

[14] Ian Goodfellow, Yoshua Bengio, and Aaron Courville. *Deep Learning*. MIT Press, 2016. http://www.deeplearningbook.org.

[15] Eduard Gorbunov, Marina Danilova, and Alexander Gasnikov. Stochastic optimization with heavy-tailed noise via accelerated gradient clipping. In H. Larochelle, M. Ranzato, R. Hadsell, M. F. Balcan, and H. Lin, editors, *Advances in Neural Information Processing Systems*, volume 33, pages 15042–15053. Curran Associates, Inc., 2020.

[16] Robert Mansel Gower, Nicolas Loizou, Xun Qian, Alibek Sailanbayev, Egor Shulgin, and Peter Richtárik. Sgd: General analysis and improved rates. In *International Conference on Machine Learning*, pages 5200–5209, 2019.

[17] Vincent Guigues, Anatoli Juditsky, and Arkadi Nemirovski. Non-asymptotic confidence bounds for the optimal value of a stochastic program. *Optimization Methods and Software*, 32(5):1033–1058, 2017.

[18] Cristóbal Guzmán and Arkadi Nemirovski. On lower complexity bounds for large-scale smooth convex optimization. *Journal of Complexity*, 31(1):1–14, 2015.

[19] Elad Hazan, Kfir Levy, and Shai Shalev-Shwartz. Beyond convexity: Stochastic quasi-convex optimization. In *Advances in Neural Information Processing Systems*, pages 1594–1602, 2015.

[20] Anatoli Juditsky, Arkadi Nemirovski, et al. First order methods for nonsmooth convex large-scale optimization, i: general purpose methods. *Optimization for Machine Learning*, pages 121–148, 2011.

[21] Diederik P Kingma and Jimmy Ba. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980*, 2014.

[22] Guanghui Lan. An optimal method for stochastic composite optimization. *Mathematical Programming*, 133(1-2):365–397, 2012.

[23] Vien V Mai and Mikael Johansson. Stability and convergence of stochastic gradient clipping: Beyond lipschitz continuity and smoothness. *arXiv preprint arXiv:2102.06489*, 2021.

[24] Aditya Krishna Menon, Ankit Singh Rawat, Sashank J Reddi, and Sanjiv Kumar. Can gradient clipping mitigate label noise? In *International Conference on Learning Representations*, 2020.

[25] Eric Moulines and Francis R Bach. Non-asymptotic analysis of stochastic approximation algorithms for machine learning. In *Advances in Neural Information Processing Systems*, pages 451–459, 2011.

[26] Aleksandr Viktorovich Nazin, AS Nemirovsky, Aleksandr Borisovich Tsybakov, and AB Juditsky. Algorithms of robust stochastic optimization based on mirror descent method. *Automation and Remote Control*, 80(9):1607–1627, 2019.

[27] Arkadi Nemirovski, Anatoli Juditsky, Guanghui Lan, and Alexander Shapiro. Robust stochastic approximation approach to stochastic programming. *SIAM Journal on optimization*, 19(4):1574–1609, 2009.

[28] Arkadi Semenovich Nemirovsky and David Borisovich Yudin. Problem complexity and method efficiency in optimization. 1983.

[29] Yu Nesterov. Universal gradient methods for convex optimization problems. *Mathematical Programming*, 152(1-2):381–404, 2015.

[30] Yurii E Nesterov. A method for solving the convex programming problem with convergence rate o (1/kˆ 2). In *Dokl. akad. nauk Sssr*, volume 269, pages 543–547, 1983.

[31] Razvan Pascanu, Tomas Mikolov, and Yoshua Bengio. On the difficulty of training recurrent neural networks. In *International conference on machine learning*, pages 1310–1318, 2013.

[32] Adam Paszke, Sam Gross, Francisco Massa, Adam Lerer, James Bradbury, Gregory Chanan, Trevor Killeen, Zeming Lin, Natalia Gimelshein, Luca Antiga, Alban Desmaison, Andreas Kopf, Edward Yang, Zachary DeVito, Martin Raison, Alykhan Tejani, Sasank Chilamkurthy, Benoit Steiner, Lu Fang, Junjie Bai, and Soumith Chintala. Pytorch: An imperative style, high-performance deep learning library. In H. Wallach, H. Larochelle, A. Beygelzimer, F. d'Alché-Buc, E. Fox, and R. Garnett, editors, *Advances in Neural Information Processing Systems 32*, pages 8024–8035. Curran Associates, Inc., 2019.

[33] Herbert Robbins and Sutton Monro. A stochastic approximation method. *The annals of mathematical statistics*, pages 400–407, 1951.

[34] Olga Russakovsky, Jia Deng, Hao Su, Jonathan Krause, Sanjeev Satheesh, Sean Ma, Zhiheng Huang, Andrej Karpathy, Aditya Khosla, Michael Bernstein, Alexander C. Berg, and Li Fei-Fei. ImageNet Large Scale Visual Recognition Challenge. *International Journal of Computer Vision (IJCV)*, 115(3):211–252, 2015.

[35] Umut Şimşekli, Mert Gürbüzbalaban, Thanh Huy Nguyen, Gaël Richard, and Levent Sagun. On the heavy-tailed theory of stochastic gradient descent for deep neural networks. *arXiv preprint arXiv:1912.00018*, 2019.

[36] Umut Simsekli, Levent Sagun, and Mert Gurbuzbalaban. A tail-index analysis of stochastic gradient noise in deep neural networks. *arXiv preprint arXiv:1901.06053*, 2019.

[37] Vladimir Spokoiny et al. Parametric estimation. finite sample theory. *The Annals of Statistics*, 40(6):2877–2909, 2012.

[38] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. In *Advances in neural information processing systems*, pages 5998–6008, 2017.

[39] Alex Warstadt, Amanpreet Singh, and Samuel R Bowman. Neural network acceptability judgments. *arXiv preprint arXiv:1805.12471*, 2018.

[40] Thomas Wolf, Lysandre Debut, Victor Sanh, Julien Chaumond, Clement Delangue, Anthony Moi, Pierric Cistac, Tim Rault, Rémi Louf, Morgan Funtowicz, Joe Davison, Sam Shleifer, Patrick von Platen, Clara Ma, Yacine Jernite, Julien Plu, Canwen Xu, Teven Le Scao, Sylvain Gugger, Mariama Drame, Quentin Lhoest, and Alexander M. Rush. Transformers: State-of-the-art natural language processing. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing: System Demonstrations*, pages 38–45, Online, October 2020. Association for Computational Linguistics.

[41] Jingzhao Zhang, Tianxing He, Suvrit Sra, and Ali Jadbabaie. Why gradient clipping accelerates training: A theoretical justification for adaptivity. In *International Conference on Learning Representations*, 2020.

[42] Jingzhao Zhang, Sai Praneeth Karimireddy, Andreas Veit, Seungyeon Kim, Sashank Reddi, Sanjiv Kumar, and Suvrit Sra. Why are adaptive methods good for attention models? In H. Larochelle, M. Ranzato, R. Hadsell, M. F. Balcan, and H. Lin, editors, *Advances in Neural Information Processing Systems*, volume 33, pages 15383–15393. Curran Associates, Inc., 2020.

# Checklist

1. For all authors...

    (a) Do the main claims made in the abstract and introduction accurately reflect the paper's contributions and scope? [Yes]

    (b) Did you describe the limitations of your work? [Yes] Section 1.1 describes all assumptions that we use

    (c) Did you discuss any potential negative societal impacts of your work? [No] Our results are primarily theoretical, therefore, such a discussion is not applicable.

    (d) Have you read the ethics review guidelines and ensured that your paper conforms to them? [Yes]

2. If you are including theoretical results...

    (a) Did you state the full set of assumptions of all theoretical results? [Yes] Section 1.1 describes all assumptions that we use.

    (b) Did you include complete proofs of all theoretical results? [Yes] Appendix B and C include the complete proofs of all the results we derive.

3. If you ran experiments...

    (a) Did you include the code, data, and instructions needed to reproduce the main experimental results (either in the supplemental material or as a URL)? [Yes] See our code in the supplementary material.

    (b) Did you specify all the training details (e.g., data splits, hyperparameters, how they were chosen)? [Yes] See Appendix D.

    (c) Did you report error bars (e.g., with respect to the random seed after running experiments multiple times)? [No] Instead of it, we show the averaged trajectories of the methods' convergence.

    (d) Did you include the total amount of compute and the type of resources used (e.g., type of GPUs, internal cluster, or cloud provider)? [Yes] See Appendix D.

4. If you are using existing assets (e.g., code, data, models) or curating/releasing new assets...

(a) If your work uses existing assets, did you cite the creators? [Yes]

(b) Did you mention the license of the assets? [No] We use only publicly available resources.

(c) Did you include any new assets either in the supplemental material or as a URL? [No]

(d) Did you discuss whether and how consent was obtained from people whose data you're using/curating? [No] We use only publicly available resources.

(e) Did you discuss whether the data you are using/curating contains personally identifiable information or offensive content? [N/A]

5. If you used crowdsourcing or conducted research with human subjects...

(a) Did you include the full text of instructions given to participants and screenshots, if applicable? [N/A]

(b) Did you describe any potential participant risks, with links to Institutional Review Board (IRB) approvals, if applicable? [N/A]

(c) Did you include the estimated hourly wage paid to participants and the total amount spent on participant compensation? [N/A]

# Contents

# A  Basic facts, technical lemmas, and auxiliary results

## A.1  Notation, missing definitions, and useful inequalities

**Notation and missing definitions.**  We use standard notation for stochastic optimization. For all $x \in \mathbb{R}^n$ we use $\|x\|_2 = \sqrt{\langle x, x \rangle}$ to denote standard Euclidean norm, where $\langle x, y \rangle = x_1 y_1 + x_2 y_2 + \ldots + x_n y_n$, $x = (x_1, \ldots, x_n)^\top$, $x = (x_1, \ldots, x_n)^\top \in \mathbb{R}^n$. Next, we use $\mathbb{E}[\xi]$ and $\mathbb{E}[\xi \mid \eta]$ to denote expectation of $\xi$ and expectation of $\xi$ conditioned on $\eta$ respectively. In some places of the paper, we also use $\mathbb{E}_\xi[\cdot]$ to denote conditional expectation taken w.r.t. the randomness coming from $\xi$. The probability of event $E$ is defined as $\mathbb{P}\{E\}$.

Finally, we use a standard definition of differentiable strongly convex function.

**Definition A.1.** *Differentiable function $f : Q \subseteq \mathbb{R}^n \to \mathbb{R}$ is called $\mu$-strongly convex for some $\mu \geq 0$ if for all $x, y \in Q$*

$$f(y) \geq f(x) + \langle \nabla f(x), y - x \rangle + \frac{\mu}{2}\|y - x\|_2^2.$$

*When $\mu = 0$ function $f$ is called convex.*

**Useful inequalities.**  For all $a, b \in \mathbb{R}^n$ and $\lambda > 0$

$$|\langle a, b \rangle| \leq \frac{\|a\|_2^2}{2\lambda} + \frac{\lambda \|b\|_2^2}{2}, \tag{10}$$

$$\|a + b\|_2^2 \leq 2\|a\|_2^2 + 2\|b\|_2^2, \tag{11}$$

$$\langle a, b \rangle = \frac{1}{2}\left(\|a + b\|_2^2 - \|a\|_2^2 - \|b\|_2^2\right). \tag{12}$$

## A.2  Auxiliary lemmas

**Lemma A.1** ([5, 29])**.** *Let $f$ be $(\nu, M_\nu)$-Hölder continuous on $Q \subseteq \mathbb{R}^n$. Then for all $x, y \in Q$ and for all $\delta > 0$*

$$f(y) \leq f(x) + \langle \nabla f(x), y - x \rangle + \frac{M_\nu}{1 + \nu}\|x - y\|_2^{1+\nu}, \tag{13}$$

$$f(y) \leq f(x) + \langle \nabla f(x), y - x \rangle + \frac{L(\delta, \nu)}{2}\|x - y\|_2^2 + \frac{\delta}{2}, \quad L(\delta, \nu) = \left(\frac{1}{\delta}\right)^{\frac{1-\nu}{1+\nu}} M_\nu^{\frac{2}{1+\nu}}. \tag{14}$$

**Lemma A.2** (Bernstein inequality for martingale differences [1, 7, 8])**.** *Let the sequence of random variables $\{X_i\}_{i \geq 1}$ form a martingale difference sequence, i.e. $\mathbb{E}[X_i \mid X_{i-1}, \ldots, X_1] = 0$ for all $i \geq 1$. Assume that conditional variances $\sigma_i^2 \stackrel{def}{=} \mathbb{E}[X_i^2 \mid X_{i-1}, \ldots, X_1]$ exist and are bounded and assume also that there exists deterministic constant $c > 0$ such that $\|X_i\|_2 \leq c$ almost surely for all $i \geq 1$. Then for all $b > 0$, $F > 0$ and $n \geq 1$*

$$\mathbb{P}\left\{\left|\sum_{i=1}^n X_i\right| > b \text{ and } \sum_{i=1}^n \sigma_i^2 \leq F\right\} \leq 2\exp\left(-\frac{b^2}{2F + {2cb}/{3}}\right). \tag{15}$$

## A.3  Technical lemmas

**Lemma A.3.** *Let sequences $\{\alpha_k\}_{k \geq 0}$ and $\{A_k\}_{k \geq 0}$ satisfy*

$$\alpha_0 = A_0 = 0, \quad \alpha_{k+1} = \frac{(k+1)^{\frac{2\nu}{1+\nu}}(\varepsilon/2)^{\frac{1-\nu}{1+\nu}}}{2^{\frac{2\nu}{1+\nu}} a M_\nu^{\frac{2}{1+\nu}}}, \quad A_{k+1} = A_k + \alpha_{k+1}, \quad a, \varepsilon, M_\nu > 0, \ \nu \in [0, 1] \tag{16}$$

*for all $k \geq 0$. Then for all $k \geq 0$ we have*

$$A_k \geq a L_k \alpha_k^2, \quad A_k \geq \frac{k^{\frac{1+3\nu}{1+\nu}}(\varepsilon/2)^{\frac{1-\nu}{1+\nu}}}{2^{\frac{1+3\nu}{1+\nu}} a M_\nu^{\frac{2}{1+\nu}}}, \tag{17}$$

*where $L_0 = 0$ and for $k > 0$*

$$L_k = \left(\frac{2A_k}{\alpha_k \varepsilon}\right)^{\frac{1-\nu}{1+\nu}} M_\nu^{\frac{2}{1+\nu}}. \tag{18}$$

*Moreover, for all $k \geq 0$*

$$A_k \leq \frac{k^{\frac{1+3\nu}{1+\nu}} (\varepsilon/2)^{\frac{1-\nu}{1+\nu}}}{2^{\frac{2\nu}{1+\nu}} a M_\nu^{\frac{2}{1+\nu}}}. \tag{19}$$

*Proof.* We start with deriving the second inequality from (17). The proof goes by induction. For $k = 0$ the inequality holds. Next, we assume that it holds for all $k \leq K$. Then,

$$A_{K+1} = A_K + \alpha_{K+1} \geq \frac{K^{\frac{1+3\nu}{1+\nu}} (\varepsilon/2)^{\frac{1-\nu}{1+\nu}}}{2^{\frac{1+3\nu}{1+\nu}} a M_\nu^{\frac{2}{1+\nu}}} + \frac{(K+1)^{\frac{2\nu}{1+\nu}} (\varepsilon/2)^{\frac{1-\nu}{1+\nu}}}{2^{\frac{2\nu}{1+\nu}} a M_\nu^{\frac{2}{1+\nu}}}.$$

Let us estimate the right-hand side of the previous inequality. We want to show that

$$\frac{K^{\frac{1+3\nu}{1+\nu}} (\varepsilon/2)^{\frac{1-\nu}{1+\nu}}}{2^{\frac{1+3\nu}{1+\nu}} a M_\nu^{\frac{2}{1+\nu}}} + \frac{(K+1)^{\frac{2\nu}{1+\nu}} (\varepsilon/2)^{\frac{1-\nu}{1+\nu}}}{2^{\frac{2\nu}{1+\nu}} a M_\nu^{\frac{2}{1+\nu}}} \geq \frac{(K+1)^{\frac{1+3\nu}{1+\nu}} (\varepsilon/2)^{\frac{1-\nu}{1+\nu}}}{2^{\frac{1+3\nu}{1+\nu}} a M_\nu^{\frac{2}{1+\nu}}}$$

that is equivalent to the inequality:

$$\frac{K^{\frac{1+3\nu}{1+\nu}}}{2} + (K+1)^{\frac{2\nu}{1+\nu}} \geq \frac{(K+1)^{\frac{1+3\nu}{1+\nu}}}{2} \iff \frac{K^{\frac{1+3\nu}{1+\nu}}}{2} \geq \frac{(K+1)^{\frac{2\nu}{1+\nu}}(K-1)}{2}.$$

If $K = 1$, it trivially holds. If $K > 1$, it is equivalent to

$$\frac{K}{K-1} \geq \left(\frac{K+1}{K}\right)^{2-\frac{2}{1+\nu}}.$$

Since $2 - \frac{2}{1+\nu}$ is monotonically increasing function for $\nu \in [0,1]$ we have that

$$\left(\frac{K+1}{K}\right)^{2-\frac{2}{1+\nu}} \leq \frac{K+1}{K} \leq \frac{K}{K-1}.$$

That is, the second inequality in (17) holds for $k = K + 1$, and, as a consequence, it holds for all $k \geq 0$. Next, we derive the first part of (17). For $k = 0$ it trivially holds. For $k > 0$ we consider cases $\nu = 0$ and $\nu > 0$ separately. When $\nu = 0$ the inequality is equivalent to

$$1 \geq \frac{2a\alpha_k M_0^2}{\varepsilon}, \text{ where } \frac{2a\alpha_k M_0^2}{\varepsilon} \overset{(16)}{=} 1,$$

i.e., we have $A_k = aL_k \alpha_k^2$ for all $k \geq 0$. When $\nu > 0$ the first inequality in (17) is equivalent to

$$A_k \geq a^{\frac{1+\nu}{2\nu}} \alpha_k^{\frac{1+3\nu}{2\nu}} (\varepsilon/2)^{-\frac{1-\nu}{2\nu}} M_\nu^{\frac{1}{\nu}} \overset{(16)}{\iff} A_k \geq \frac{k^{\frac{1+3\nu}{1+\nu}} (\varepsilon/2)^{\frac{1-\nu}{1+\nu}}}{2^{\frac{1+3\nu}{1+\nu}} a M_\nu^{\frac{2}{1+\nu}}},$$

where the last inequality coincides with the second inequality from (17) that we derived earlier in the proof.

To finish the proof it remains to derive (19). Again, the proof goes by induction. For $k = 0$ inequality (19) is trivial. Next, we assume that it holds for all $k \leq K$. Then,

$$A_{K+1} = A_K + \alpha_{K+1} \leq \frac{K^{\frac{1+3\nu}{1+\nu}} (\varepsilon/2)^{\frac{1-\nu}{1+\nu}}}{2^{\frac{2\nu}{1+\nu}} a M_\nu^{\frac{2}{1+\nu}}} + \frac{(K+1)^{\frac{2\nu}{1+\nu}} (\varepsilon/2)^{\frac{1-\nu}{1+\nu}}}{2^{\frac{2\nu}{1+\nu}} a M_\nu^{\frac{2}{1+\nu}}}.$$

Let us estimate the right-hand side of the previous inequality. We want to show that

$$\frac{K^{\frac{1+3\nu}{1+\nu}} (\varepsilon/2)^{\frac{1-\nu}{1+\nu}}}{2^{\frac{2\nu}{1+\nu}} a M_\nu^{\frac{2}{1+\nu}}} + \frac{(K+1)^{\frac{2\nu}{1+\nu}} (\varepsilon/2)^{\frac{1-\nu}{1+\nu}}}{2^{\frac{2\nu}{1+\nu}} a M_\nu^{\frac{2}{1+\nu}}} \leq \frac{(K+1)^{\frac{1+3\nu}{1+\nu}} (\varepsilon/2)^{\frac{1-\nu}{1+\nu}}}{2^{\frac{2\nu}{1+\nu}} a M_\nu^{\frac{2}{1+\nu}}}$$

that is equivalent to the inequality:

$$K^{\frac{1+3\nu}{1+\nu}} + (K+1)^{\frac{2\nu}{1+\nu}} \leq (K+1)^{\frac{1+3\nu}{1+\nu}}.$$

This inequality holds due to

$$K^{\frac{1+3\nu}{1+\nu}} \leq (K+1)^{\frac{2\nu}{1+\nu}} K.$$

That is, (19) holds for $k = K + 1$, and, as a consequence, it holds for all $k \geq 0$. $\qquad\square$

**Lemma A.4.** *Let $f$ have Hölder continuous gradients on $Q \subseteq \mathbb{R}^n$ for some $\nu \in [0,1]$ with constant $M_\nu > 0$, be convex and $x^* \in Q$ be some minimum of $f(x)$ on $\mathbb{R}^n$. Then, for all $x \in \mathbb{R}^n$*

$$\|\nabla f(x)\|_2 \leq \left(\frac{1+\nu}{\nu}\right)^{\frac{\nu}{1+\nu}} M_\nu^{\frac{1}{1+\nu}} \left(f(x) - f(x^*)\right)^{\frac{\nu}{1+\nu}}, \tag{20}$$

*where for $\nu = 0$ we use $\left[\left(\frac{1+\nu}{\nu}\right)^{\frac{\nu}{1+\nu}}\right]_{\nu=0} := \lim_{\nu \to 0} \left(\frac{1+\nu}{\nu}\right)^{\frac{\nu}{1+\nu}} = 1$.*

*Proof.* For $\nu = 0$ inequality (20) follows from (3) and[1] $\nabla f(x^*) = 0$. When $\nu > 0$ for arbitrary point $x \in Q$ we consider the point $y = x - \alpha \nabla f(x)$, where $\alpha = \left(\frac{\|\nabla f(x)\|_2^{1-\nu}}{M_\nu}\right)^{\frac{1}{\nu}}$. Since $x^* \in Q$ and $f$ is convex one can easily show that $y \in Q$. For the pair of points $x, y$ we apply (13) and get

$$
\begin{aligned}
f(y) &\leq f(x) + \langle \nabla f(x), y - x\rangle + \frac{M_\nu}{1+\nu}\|x-y\|_2^{1+\nu} \\
&= f(x) - \alpha\|\nabla f(x)\|^2 + \frac{\alpha^{\nu+1}M_\nu}{1+\nu}\|\nabla f(x)\|_2^{1+\nu} \\
&= f(x) - \frac{\|\nabla f(x)\|_2^{\frac{1+\nu}{\nu}}}{M_\nu^{\frac{1}{\nu}}} + \frac{\|\nabla f(x)\|_2^{\frac{1+\nu}{\nu}}}{(1+\nu)M_\nu^{\frac{1}{\nu}}} = f(x) - \frac{\nu\|\nabla f(x)\|_2^{\frac{1+\nu}{\nu}}}{(1+\nu)M_\nu^{\frac{1}{\nu}}}
\end{aligned}
$$

implying

$$\|\nabla f(x)\|_2 \leq \left(\frac{1+\nu}{\nu}\right)^{\frac{\nu}{1+\nu}} M_\nu^{\frac{1}{1+\nu}} \left(f(x) - f(y)\right)^{\frac{\nu}{1+\nu}} \leq \left(\frac{1+\nu}{\nu}\right)^{\frac{\nu}{1+\nu}} M_\nu^{\frac{1}{1+\nu}} \left(f(x) - f(x^*)\right)^{\frac{\nu}{1+\nu}}.$$

□

**Lemma A.5.** *Let $f$ have Hölder continuous gradients on $Q \subseteq \mathbb{R}^n$ for some $\nu \in [0,1]$ with constant $M_\nu > 0$, be convex and $x^* \in Q$ be some minimum of $f(x)$ on $\mathbb{R}^n$. Then, for all $x \in \mathbb{R}^n$ and all $\delta > 0$,*

$$\|\nabla f(x)\|_2^2 \leq 2\left(\frac{1}{\delta}\right)^{\frac{1-\nu}{1+\nu}} M_\nu^{\frac{2}{1+\nu}} \left(f(x) - f(x^*)\right) + \delta^{\frac{2\nu}{1+\nu}} M_\nu^{\frac{2}{1+\nu}}. \tag{21}$$

*Proof.* For a given $\delta > 0$ we consider an arbitrary point $x \in Q$ and $y = x - \frac{1}{L(\delta,\nu)}\nabla f(x)$, where $L(\delta,\nu) = \left(\frac{1}{\delta}\right)^{\frac{1-\nu}{1+\nu}} M_\nu^{\frac{2}{1+\nu}}$. Since $x^* \in Q$ and $f$ is convex one can easily show that $y \in Q$. For the pair of points $x, y$ we apply (14) and get

$$
\begin{aligned}
f(y) &\leq f(x) + \langle \nabla f(x), y - x\rangle + \frac{L(\delta,\nu)}{2}\|x-y\|_2^2 + \frac{\delta}{2} \\
&= f(x) - \frac{1}{2L(\delta,\nu)}\|x-y\|_2^2 + \frac{\delta}{2}
\end{aligned}
$$

implying

$$
\begin{aligned}
\|\nabla f(x)\|_2^2 &\leq 2L(\delta,\nu)\left(f(x) - f(y)\right) + \delta L(\delta,\nu) \\
&\leq 2\left(\frac{1}{\delta}\right)^{\frac{1-\nu}{1+\nu}} M_\nu^{\frac{2}{1+\nu}} \left(f(x) - f(x^*)\right) + \delta^{\frac{2\nu}{1+\nu}} M_\nu^{\frac{2}{1+\nu}}.
\end{aligned}
$$

□

---

[1] When $f$ is not differentiable, we use subgradients. In this case, 0 belongs to the subdifferential of $f$ at the point $x^*$ and we take it as $\nabla f(x^*)$.

# B Clipped Similar Triangles Method: missing details and proofs

## B.1 Convergence in the convex case

In this section, we provide the full proof of Thm. 2.1 together with complete statement of the result.

### B.1.1 Two lemmas

The analysis of clipped-SSTM consists of 3 main steps. The first one is an "optimization lemma" – a modification of a standard lemma for Similar Triangles Method (see [9] and Lemma F.4 from [15]). This result helps to estimate the progress of the method after $N$ iterations.

**Lemma B.1.** *Let $f$ be a convex function with a minimum at some[2] point $x^*$, its gradient be $(\nu, M_\nu)$-Hölder continuous on a ball $B_{3R_0}(x^*)$, where $R_0 \geq \|x^0 - x^*\|_2$, and let stepsize parameter $a$ satisfy $a \geq 1$. If $x^k, y^k, z^k \in B_{3R_0}(x^*)$ for all $k = 0, 1, \ldots, N$, $N \geq 0$, then after $N$ iterations of* clipped-SSTM *for all $z \in \mathbb{R}^n$ we have*

$$
\begin{aligned}
A_N \left( f(y^N) - f(z) \right) &\leq \frac{1}{2}\|z^0 - z\|_2^2 - \frac{1}{2}\|z^N - z\|_2^2 + \sum_{k=0}^{N-1} \alpha_{k+1} \left\langle \theta_{k+1}, z - z^k \right\rangle \\
&\quad + \sum_{k=0}^{N-1} \alpha_{k+1}^2 \|\theta_{k+1}\|_2^2 + \sum_{k=0}^{N-1} \alpha_{k+1}^2 \left\langle \theta_{k+1}, \nabla f(x^{k+1}) \right\rangle + \frac{A_N \varepsilon}{4}
\end{aligned}
\tag{22}
$$

$$
\theta_{k+1} \overset{def}{=} \widetilde{\nabla} f(x^{k+1}, \boldsymbol{\xi}^k) - \nabla f(x^{k+1}).
\tag{23}
$$

*Proof.* Consider an arbitrary $k \in \{0, 1, \ldots, N - 1\}$. Using $z^{k+1} = z^k - \alpha_{k+1}\widetilde{\nabla} f(x^{k+1}, \boldsymbol{\xi}^k)$ we get that for all $z \in \mathbb{R}^n$

$$
\begin{aligned}
\alpha_{k+1} \left\langle \widetilde{\nabla} f(x^{k+1}, \boldsymbol{\xi}^k), z^k - z \right\rangle &= \alpha_{k+1} \left\langle \widetilde{\nabla} f(x^{k+1}, \boldsymbol{\xi}^k), z^k - z^{k+1} \right\rangle \\
&\quad + \alpha_{k+1} \left\langle \widetilde{\nabla} f(x^{k+1}, \boldsymbol{\xi}^k), z^{k+1} - z \right\rangle \\
&= \alpha_{k+1} \left\langle \widetilde{\nabla} f(x^{k+1}, \boldsymbol{\xi}^k), z^k - z^{k+1} \right\rangle + \left\langle z^{k+1} - z^k, z - z^{k+1} \right\rangle \\
&\overset{(12)}{=} \alpha_{k+1} \left\langle \widetilde{\nabla} f(x^{k+1}, \boldsymbol{\xi}^k), z^k - z^{k+1} \right\rangle - \frac{1}{2}\|z^k - z^{k+1}\|_2^2 \\
&\quad + \frac{1}{2}\|z^k - z\|_2^2 - \frac{1}{2}\|z^{k+1} - z\|_2^2.
\end{aligned}
\tag{24}
$$

Next, we notice that

$$
y^{k+1} = \frac{A_k y^k + \alpha_{k+1} z^{k+1}}{A_{k+1}} = \frac{A_k y^k + \alpha_{k+1} z^k}{A_{k+1}} + \frac{\alpha_{k+1}}{A_{k+1}} \left( z^{k+1} - z^k \right) = x^{k+1} + \frac{\alpha_{k+1}}{A_{k+1}} \left( z^{k+1} - z^k \right)
\tag{25}
$$

---

[2]Our proofs are valid for any solution $x^*$ and, for example, one can take as $x^*$ the closest solution to the starting point $x^0$.

532 implying

$$
\begin{aligned}
\alpha_{k+1}\left\langle\widetilde{\nabla}f(x^{k+1},\boldsymbol{\xi}^k),z^k-z\right\rangle
\overset{(23),(24)}{\leq}\;&\alpha_{k+1}\left\langle\nabla f(x^{k+1}),z^k-z^{k+1}\right\rangle-\frac{1}{2}\|z^k-z^{k+1}\|_2^2\\
&+\alpha_{k+1}\left\langle\theta_{k+1},z^k-z^{k+1}\right\rangle+\frac{1}{2}\|z^k-z\|_2^2-\frac{1}{2}\|z^{k+1}-z\|_2^2\\
\overset{(25)}{=}\;&A_{k+1}\left\langle\nabla f(x^{k+1}),x^{k+1}-y^{k+1}\right\rangle-\frac{1}{2}\|z^k-z^{k+1}\|_2^2\\
&+\alpha_{k+1}\left\langle\theta_{k+1},z^k-z^{k+1}\right\rangle+\frac{1}{2}\|z^k-z\|_2^2-\frac{1}{2}\|z^{k+1}-z\|_2^2\\
\overset{(14)}{\leq}\;&A_{k+1}\left(f(x^{k+1})-f(y^{k+1})\right)+\frac{A_{k+1}L_{k+1}}{2}\|x^{k+1}-y^{k+1}\|_2^2\\
&+\frac{\alpha_{k+1}\varepsilon}{4}-\frac{1}{2}\|z^k-z^{k+1}\|_2^2+\alpha_{k+1}\left\langle\theta_{k+1},z^k-z^{k+1}\right\rangle\\
&+\frac{1}{2}\|z^k-z\|_2^2-\frac{1}{2}\|z^{k+1}-z\|_2^2\\
\overset{(25)}{=}\;&A_{k+1}\left(f(x^{k+1})-f(y^{k+1})\right)+\frac{1}{2}\left(\frac{\alpha_{k+1}^2L_{k+1}}{A_{k+1}}-1\right)\|z^k-z^{k+1}\|_2^2\\
&+\alpha_{k+1}\left\langle\theta_{k+1},z^k-z^{k+1}\right\rangle+\frac{1}{2}\|z^k-z\|_2^2-\frac{1}{2}\|z^{k+1}-z\|_2^2+\frac{\alpha_{k+1}\varepsilon}{4},
\end{aligned}
$$

533 where in the third inequality we used $x^{k+1},y^{k+1}\in B_{3R_0}(x^*)$ and (14) with $\delta=\frac{\alpha_{k+1}}{2A_{k+1}}\varepsilon$ and
534 $L(\delta,\nu)=L_{k+1}=\left(\frac{2A_{k+1}}{\varepsilon\alpha_{k+1}}\right)^{\frac{1-\nu}{1+\nu}}M_\nu^{\frac{2}{1+\nu}}$. Since $A_{k+1}\geq aL_{k+1}\alpha_{k+1}^2$ (Lemma A.3) and $a\geq 1$ we
535 can continue our derivations:

$$
\begin{aligned}
\alpha_{k+1}\left\langle\widetilde{\nabla}f(x^{k+1},\boldsymbol{\xi}^k),z^k-z\right\rangle\;\leq\;&A_{k+1}\left(f(x^{k+1})-f(y^{k+1})\right)+\alpha_{k+1}\left\langle\theta_{k+1},z^k-z^{k+1}\right\rangle\\
&+\frac{1}{2}\|z^k-z\|_2^2-\frac{1}{2}\|z^{k+1}-z\|_2^2+\frac{\alpha_{k+1}\varepsilon}{4}.
\end{aligned}\tag{26}
$$

536 Next, due to convexity of $f$ we have

$$
\begin{aligned}
\left\langle\widetilde{\nabla}f(x^{k+1},\boldsymbol{\xi}^k),y^k-x^{k+1}\right\rangle\;\overset{(23)}{=}\;&\left\langle\nabla f(x^{k+1}),y^k-x^{k+1}\right\rangle+\left\langle\theta_{k+1},y^k-x^{k+1}\right\rangle\\
\leq\;&f(y^k)-f(x^{k+1})+\left\langle\theta_{k+1},y^k-x^{k+1}\right\rangle.
\end{aligned}\tag{27}
$$

537 By definition of $x^{k+1}$ we have $x^{k+1}=\frac{A_ky^k+\alpha_{k+1}z^k}{A_{k+1}}$ implying

$$
\alpha_{k+1}\left(x^{k+1}-z^k\right)=A_k\left(y^k-x^{k+1}\right)\tag{28}
$$

19

since $A_{k+1} = A_k + \alpha_{k+1}$. Putting all together we derive that

$$
\begin{aligned}
\alpha_{k+1}\left\langle \widetilde{\nabla} f(x^{k+1}, \boldsymbol{\xi}^k), x^{k+1} - z \right\rangle \quad &= \quad \alpha_{k+1}\left\langle \widetilde{\nabla} f(x^{k+1}, \boldsymbol{\xi}^k), x^{k+1} - z^k \right\rangle \\
&\quad + \alpha_{k+1}\left\langle \widetilde{\nabla} f(x^{k+1}, \boldsymbol{\xi}^k), z^k - z \right\rangle \\
&\overset{(28)}{=} \quad A_k\left\langle \widetilde{\nabla} f(x^{k+1}, \boldsymbol{\xi}^k), y^k - x^{k+1} \right\rangle \\
&\quad + \alpha_{k+1}\left\langle \widetilde{\nabla} f(x^{k+1}, \boldsymbol{\xi}^k), z^k - z \right\rangle \\
&\overset{(27),(26)}{\leq} \quad A_k\left( f(y^k) - f(x^{k+1}) \right) + A_k\left\langle \theta_{k+1}, y^k - x^{k+1} \right\rangle \\
&\quad + A_{k+1}\left( f(x^{k+1}) - f(y^{k+1}) \right) + \alpha_{k+1}\left\langle \theta_{k+1}, z^k - z^{k+1} \right\rangle \\
&\quad + \frac{1}{2}\|z^k - z\|_2^2 - \frac{1}{2}\|z^{k+1} - z\|_2^2 + \frac{\alpha_{k+1}\varepsilon}{4} \\
&\overset{(28)}{=} \quad A_k f(y^k) - A_{k+1} f(y^{k+1}) + \alpha_{k+1}\left\langle \theta_{k+1}, x^{k+1} - z^k \right\rangle \\
&\quad + \alpha_{k+1} f(x^{k+1}) + \alpha_{k+1}\left\langle \theta_{k+1}, z^k - z^{k+1} \right\rangle \\
&\quad + \frac{1}{2}\|z^k - z\|_2^2 - \frac{1}{2}\|z^{k+1} - z\|_2^2 + \frac{\alpha_{k+1}\varepsilon}{4} \\
&= \quad A_k f(y^k) - A_{k+1} f(y^{k+1}) + \alpha_{k+1} f(x^{k+1}) \\
&\quad + \alpha_{k+1}\left\langle \theta_{k+1}, x^{k+1} - z^{k+1} \right\rangle \\
&\quad + \frac{1}{2}\|z^k - z\|_2^2 - \frac{1}{2}\|z^{k+1} - z\|_2^2 + \frac{\alpha_{k+1}\varepsilon}{4}.
\end{aligned}
$$

Rearranging the terms we get

$$
\begin{aligned}
A_{k+1} f(y^{k+1}) - A_k f(y^k) \quad &\leq \quad \alpha_{k+1}\left( f(x^{k+1}) + \left\langle \widetilde{\nabla} f(x^{k+1}, \boldsymbol{\xi}^k), z - x^{k+1} \right\rangle \right) + \frac{1}{2}\|z^k - z\|_2^2 \\
&\quad - \frac{1}{2}\|z^{k+1} - z\|_2^2 + \alpha_{k+1}\left\langle \theta_{k+1}, x^{k+1} - z^{k+1} \right\rangle + \frac{\alpha_{k+1}\varepsilon}{4} \\
&\overset{(23)}{=} \quad \alpha_{k+1}\left( f(x^{k+1}) + \left\langle \nabla f(x^{k+1}), z - x^{k+1} \right\rangle \right) \\
&\quad + \alpha_{k+1}\left\langle \theta_{k+1}, z - x^{k+1} \right\rangle + \frac{1}{2}\|z^k - z\|_2^2 - \frac{1}{2}\|z^{k+1} - z\|_2^2 \\
&\quad + \alpha_{k+1}\left\langle \theta_{k+1}, x^{k+1} - z^{k+1} \right\rangle + \frac{\alpha_{k+1}\varepsilon}{4} \\
&\leq \quad \alpha_{k+1} f(z) + \frac{1}{2}\|z^k - z\|_2^2 - \frac{1}{2}\|z^{k+1} - z\|_2^2 + \alpha_{k+1}\left\langle \theta_{k+1}, z - z^{k+1} \right\rangle + \frac{\alpha_{k+1}\varepsilon}{4}
\end{aligned}
$$

where in the last inequality we use the convexity of $f$. Taking into account $A_0 = \alpha_0 = 0$ and $A_N = \sum_{k=0}^{N-1} \alpha_{k+1}$ we sum up these inequalities for $k = 0, \dots, N-1$ and get

$$
\begin{aligned}
A_N f(y^N) \quad &\leq \quad A_N f(z) + \frac{1}{2}\|z^0 - z\|_2^2 - \frac{1}{2}\|z^N - z\|_2^2 + \sum_{k=0}^{N-1} \alpha_{k+1}\left\langle \theta_{k+1}, z - z^{k+1} \right\rangle + \frac{A_N\varepsilon}{4} \\
&= \quad A_N f(z) + \frac{1}{2}\|z^0 - z\|_2^2 - \frac{1}{2}\|z^N - z\|_2^2 + \sum_{k=0}^{N-1} \alpha_{k+1}\left\langle \theta_{k+1}, z - z^k \right\rangle \\
&\quad + \sum_{k=0}^{N-1} \alpha_{k+1}^2\left\langle \theta_{k+1}, \widetilde{\nabla} f(x^{k+1}, \boldsymbol{\xi}^k) \right\rangle + \frac{A_N\varepsilon}{4} \\
&\overset{(23)}{=} \quad A_N f(z) + \frac{1}{2}\|z^0 - z\|_2^2 - \frac{1}{2}\|z^N - z\|_2^2 + \sum_{k=0}^{N-1} \alpha_{k+1}\left\langle \theta_{k+1}, z - z^k \right\rangle \\
&\quad + \sum_{k=0}^{N-1} \alpha_{k+1}^2\|\theta_{k+1}\|_2^2 + \sum_{k=0}^{N-1} \alpha_{k+1}^2\left\langle \theta_{k+1}, \nabla f(x^{k+1}) \right\rangle + \frac{A_N\varepsilon}{4}
\end{aligned}
$$

that concludes the proof. $\qquad\square$

From Lemma A.3 we know that

$$A_N \sim \frac{N^{\frac{1+3\nu}{1+\nu}} \varepsilon^{\frac{1-\nu}{1+\nu}}}{M_\nu^{\frac{2}{1+\nu}}}.$$

Therefore, in view of Lemma B.1 (inequality (22) with $z = x^*$), to derive the desired complexity bound from Thm. 2.1 it is sufficient to show that

$$\sum_{k=0}^{N-1} \alpha_{k+1} \left\langle \theta_{k+1}, z - z^k \right\rangle + \sum_{k=0}^{N-1} \alpha_{k+1}^2 \left\| \theta_{k+1} \right\|_2^2 + \sum_{k=0}^{N-1} \alpha_{k+1}^2 \left\langle \theta_{k+1}, \nabla f(x^{k+1}) \right\rangle + \frac{A_N \varepsilon}{4} \lesssim R_0^2.$$

with probability at least $1 - \beta$. One possible way to achieve this goal is to apply some concentration inequality to these three sums. Since we use clipped stochastic gradients, under a proper choice of the clipping parameter, random vector $\theta_{k+1} = \widetilde{\nabla} f(x^{k+1}, \boldsymbol{\xi}^k) - \nabla f(x^{k+1})$ is bounded in $\ell_2$-norm by $2\lambda_{k+1}$ with high probability as well. Taking into account the assumption on the stochastic gradients (see (2)), it is natural to apply Bernstein's inequality (see Lemma A.2). Despite the seeming simplicity, this part of the proof is the trickiest one.

First of all, it is useful to derive tight enough upper bounds for bias, variance and distortion of $\widetilde{\nabla} f(x^{k+1}, \boldsymbol{\xi}^k)$ – this is the second step of the whole proof. Fortunately, Lemma F.5 from [15] does exactly what we need in our proof and holds without any changes.

**Lemma B.2** (Lemma F.5 from [15].)**.** *For all $k \geq 0$ the following inequality holds:*

$$\left\| \widetilde{\nabla} f(x^{k+1}, \boldsymbol{\xi}^k) - \mathbb{E}_{\boldsymbol{\xi}^k} \left[ \widetilde{\nabla} f(x^{k+1}, \boldsymbol{\xi}^k) \right] \right\|_2 \leq 2\lambda_{k+1}. \tag{29}$$

*Moreover, if $\|\nabla f(x^{k+1})\|_2 \leq \frac{\lambda_{k+1}}{2}$ for some $k \geq 0$, then for this $k$ we have:*

$$\left\| \mathbb{E}_{\boldsymbol{\xi}^k} \left[ \widetilde{\nabla} f(x^{k+1}, \boldsymbol{\xi}^k) \right] - \nabla f(x^{k+1}) \right\|_2 \leq \frac{4\sigma^2}{m_k \lambda_{k+1}}, \tag{30}$$

$$\mathbb{E}_{\boldsymbol{\xi}^k} \left[ \left\| \widetilde{\nabla} f(x^{k+1}, \boldsymbol{\xi}^k) - \nabla f(x^{k+1}) \right\|_2^2 \right] \leq \frac{18\sigma^2}{m_k}, \tag{31}$$

$$\mathbb{E}_{\boldsymbol{\xi}^k} \left[ \left\| \widetilde{\nabla} f(x^{k+1}, \boldsymbol{\xi}^k) - \mathbb{E}_{\boldsymbol{\xi}^k} \left[ \widetilde{\nabla} f(x^{k+1}, \boldsymbol{\xi}^k) \right] \right\|_2^2 \right] \leq \frac{18\sigma^2}{m_k}. \tag{32}$$

### B.1.2 Proof of Theorem 2.1

The final, third, step of the proof is consists of providing explicit formulas and bounds for the parameters of the method and derivation of the desired result using induction and Bernstein's inequality. Below we provide the complete statement of Thm. 2.1.

**Theorem B.1.** *Assume that function $f$ is convex, achieves minimum value at some[3] $x^*$ , and the gradients of $f$ satisfy (3) with $\nu \in [0, 1]$, $M_\nu > 0$ on $B_{3R_0}(x^*)$, where $R_0 \geq \|x^0 - x^*\|_2$. Then for all $\beta \in (0, 1)$ and $N \geq 1$ such that*

$$\ln \frac{4N}{\beta} \geq 2 \tag{33}$$

*we have that after $N$ iterations of clipped-SSTM with*

$$\alpha = \frac{(\varepsilon/2)^{\frac{1-\nu}{1+\nu}}}{2^{\frac{2\nu}{1+\nu}} a M_\nu^{\frac{2}{1+\nu}}}, \quad m_k = \max \left\{ 1, \frac{20736 N \sigma^2 \alpha_{k+1}^2 \ln \frac{4N}{\beta}}{C^2 R_0^2} \right\}, \tag{34}$$

$$B = \frac{C R_0}{16 \ln \frac{4N}{\beta}}, \quad a \geq 16384 \ln^2 \frac{4N}{\beta}, \tag{35}$$

$$\varepsilon^{\frac{1-\nu}{1+\nu}} \leq \frac{a C M_\nu^{\frac{1-\nu}{1+\nu}} R_0^{1-\nu}}{16 \ln \frac{4N}{\beta}}, \quad \varepsilon \leq \frac{2^{\frac{1+\nu}{2}} a^{\frac{1+\nu}{2}} C^{1+\nu} R_0^{1+\nu} M_\nu}{100^{\frac{1+3\nu}{2}}}, \tag{36}$$

---

[3]Our proofs are valid for any solution $x^*$ and, for example, one can take as $x^*$ the closest solution to the starting point $x^0$.

$$\varepsilon^{\frac{1-\nu}{1+3\nu}} \le \min\left\{ \frac{a^{\frac{2+3\nu-\nu^2}{2(1+3\nu)}}}{2^{2+4\nu+\frac{3+8\nu-5\nu^2-6\nu^3}{(1+\nu)(1+3\nu)}}\ln\frac{4N}{\beta}}, \frac{a^{\frac{(1+\nu)^2}{1+3\nu}}}{2^{4+7\nu+\frac{2+7\nu+2\nu^2-3\nu^3}{(1+\nu)(1+3\nu)}}\ln^{1+\nu}\frac{4N}{\beta}} \right\} C^{\frac{1-\nu^2}{1+3\nu}} R_0^{\frac{1-\nu^2}{1+3\nu}} M_\nu^{\frac{1-\nu}{1+3\nu}} \tag{37}$$

568  *with probability at least $1-\beta$*

$$f(y^N) - f(x^*) \le \frac{4aC^2 R_0^2 M_\nu^{\frac{2}{1+\nu}}}{N^{\frac{1+3\nu}{1+\nu}} \varepsilon^{\frac{1-\nu}{1+\nu}}}, \tag{38}$$

569  *where*

$$N = \left\lceil \frac{2^{\frac{1+\nu}{1+3\nu}} a^{\frac{1+\nu}{1+3\nu}} C^{\frac{2(1+\nu)}{1+3\nu}} R_0^{\frac{2(1+\nu)}{1+3\nu}} M_\nu^{\frac{2}{1+3\nu}}}{\varepsilon^{\frac{2}{1+3\nu}}} \right\rceil + 1, \quad C = \sqrt{7}. \tag{39}$$

570  *In other words, if we choose $a = 16384\ln^2\frac{4N}{\beta}$, then the method achieves $f(y^N) - f(x^*) \le \varepsilon$ with*

571  *probability at least $1-\beta$ after $O\left( \frac{M_\nu^{\frac{2}{1+3\nu}} R_0^{\frac{2(1+\nu)}{1+3\nu}}}{\varepsilon^{\frac{2}{1+3\nu}}} \ln^{\frac{2(1+\nu)}{1+3\nu}} \frac{M_\nu^{\frac{2}{1+3\nu}} R_0^{\frac{2(1+\nu)}{1+3\nu}}}{\varepsilon^{\frac{2}{1+3\nu}}\beta} \right)$ iterations and requires*

$$O\left( \max\left\{ \frac{M_\nu^{\frac{2}{1+3\nu}} R_0^{\frac{2(1+\nu)}{1+3\nu}}}{\varepsilon^{\frac{2}{1+3\nu}}} \ln^{\frac{2(1+\nu)}{1+3\nu}} \frac{M_\nu^{\frac{2}{1+3\nu}} R_0^{\frac{2(1+\nu)}{1+3\nu}}}{\varepsilon^{\frac{2}{1+3\nu}}\beta}, \frac{\sigma^2 R_0^2}{\varepsilon^2} \ln \frac{M_\nu^{\frac{2}{1+3\nu}} R_0^{\frac{2(1+\nu)}{1+3\nu}}}{\varepsilon^{\frac{2}{1+3\nu}}\beta} \right\} \right) \text{ oracle calls.} \tag{40}$$

572  *Proof.* First of all, we notice that for each $k \ge 0$ iterates $x^{k+1}, z^k, y^k$ lie in the ball $B_{\widetilde{R}_k}(x^*)$, where

573  $R_k = \|z^k - x^*\|_2$, $\widetilde{R}_0 = R_0$, $\widetilde{R}_{k+1} = \max\{\widetilde{R}_k, R_{k+1}\}$. We prove it using induction. Since $y^0 =$

574  $z^0 = x^0$, $\widetilde{R}_0 = R_0 \ge \|z^0 - x^*\|_2$ and $x^1 = \frac{A_0 y^0 + \alpha_1 z^0}{A_1} = z^0$ we have that $x^1, z^0, y^0 \in B_{\widetilde{R}_0}(x^*)$.

575  Next, assume that $x^l, z^{l-1}, y^{l-1} \in B_{\widetilde{R}_{l-1}}(x^*)$ for some $l \ge 1$. By definitions of $R_l$ and $\widetilde{R}_l$ we have

576  that $z^l \in B_{R_l}(x^*) \subseteq B_{\widetilde{R}_l}(x^*)$. Since $y^l$ is a convex combination of $y^{l-1} \in B_{\widetilde{R}_{l-1}}(x^*) \subseteq B_{\widetilde{R}_l}(x^*)$,

577  $z^l \in B_{\widetilde{R}_l}(x^*)$ and $B_{\widetilde{R}_l}(x^*)$ is a convex set we conclude that $y^l \in B_{\widetilde{R}_l}(x^*)$. Finally, since $x^{l+1}$ is a

578  convex combination of $y^l$ and $z^l$ we have that $x^{l+1}$ lies in $B_{\widetilde{R}_l}(x^*)$ as well.

579  Next, our goal is to prove via induction that for all $k = 0, 1, \ldots, N$ with probability at least $1 - \frac{k\beta}{N}$

580  the following statement holds: inequalities

$$\begin{aligned} R_t^2 &\le R_0^2 + 2\sum_{l=0}^{t-1}\alpha_{l+1}\langle\theta_{l+1}, x^* - z^l\rangle + 2\sum_{l=0}^{t-1}\alpha_{l+1}^2\langle\theta_{l+1}, \nabla f(x^{l+1})\rangle \\ &\quad + 2\sum_{l=0}^{t-1}\alpha_{k+1}^2\|\theta_{l+1}\|_2^2 + \frac{A_N\varepsilon}{2} \\ &\le C^2 R_0^2 \end{aligned} \tag{41}$$

581  hold for $t = 0, 1, \ldots, k$ simultaneously where $C$ is defined in (39). Let $E_k$ denote the probabilistic

582  event that this statement holds. Then, our goal is to show that $\mathbb{P}\{E_k\} \ge 1 - \frac{k\beta}{N}$ for all $k = 0, 1, \ldots, N$.

583  For $t = 0$ inequality (41) holds with probability 1 since $C \ge 1$, hence $\mathbb{P}\{E_0\} = 1$. Next, assume

584  that for some $k = T - 1 \le N - 1$ we have $\mathbb{P}\{E_k\} = \mathbb{P}\{E_{T-1}\} \ge 1 - \frac{(T-1)\beta}{N}$. Let us prove that

585  $\mathbb{P}\{E_T\} \ge 1 - \frac{T\beta}{N}$. First of all, since $R_{T-1}$ implies $R_t \le CR_0$ for all $t = 0, 1, \ldots, T-1$ we have

586  that $\widetilde{R}_{T-1} \le CR_0$, and, as a consequence, $z^{T-1} \in B_{CR_0}(x^*)$. Therefore, probability event $E_{T-1}$

587  implies

$$\begin{aligned} \|z^T - x^*\|_2 &= \|z^{T-1} - x^* - \alpha_T\widetilde{\nabla}f(x^T, \boldsymbol{\xi}^{T-1})\|_2 \le \|z^{T-1} - x^*\|_2 + \alpha_T\|\widetilde{\nabla}f(x^T, \boldsymbol{\xi}^{T-1})\|_2 \\ &\le CR_0 + \alpha_T\lambda_T = \left(1 + \frac{1}{16\ln\frac{4N}{\beta}}\right)CR_0 \overset{(33),(39)}{\le} \left(1 + \frac{1}{32}\right)\sqrt{7}R_0 \le 3R_0, \end{aligned}$$

22

hence $\widetilde{R}_T \leq 3R_0$. Then, one can apply Lemma B.1 and get that probability event $E_{T-1}$ implies

$$
\begin{aligned}
A_t \left( f(y^t) - f(x^*) \right) \leq \ & \frac{1}{2}\|z^0 - x^*\|_2^2 - \frac{1}{2}\|z^t - x^*\|_2^2 + \sum_{k=0}^{t-1} \alpha_{k+1} \left\langle \theta_{k+1}, x^* - z^k \right\rangle \\
& + \sum_{k=0}^{t-1} \alpha_{k+1}^2 \|\theta_{k+1}\|_2^2 + \sum_{k=0}^{t-1} \alpha_{k+1}^2 \left\langle \theta_{k+1}, \nabla f(x^{k+1}) \right\rangle + \frac{A_t \varepsilon}{4}, (42)
\end{aligned}
$$

$$
\theta_{k+1} \overset{\text{def}}{=} \widetilde{\nabla} f(x^{k+1}, \boldsymbol{\xi}^k) - \nabla f(x^{k+1}) \tag{43}
$$

for all $t = 0, 1, \ldots, T-1, T$. Taking into account that $f(y^t) - f(x^*) \geq 0$ for all $y^t$ we derive that probability event $E_{T-1}$ implies

$$
R_t^2 \leq R_0^2 + 2\sum_{l=0}^{t-1} \alpha_{l+1} \left\langle \theta_{l+1}, x^* - z^l \right\rangle + 2\sum_{l=0}^{t-1} \alpha_{l+1}^2 \left\langle \theta_{l+1}, \nabla f(x^{l+1}) \right\rangle + 2\sum_{l=0}^{t-1} \alpha_{l+1}^2 \|\theta_{l+1}\|_2^2 + \frac{A_t \varepsilon}{2}. \tag{44}
$$

for all $t = 0, 1, \ldots, T$.

The rest of the proof is based on the refined analysis of inequality (44). First of all, when $\nu = 0$ from (14) for all $t \geq 0$ we have

$$
\left\| \nabla f(x^{t+1}) \right\|_2 \leq M_0 = \frac{16 M_0 B \ln \frac{4N}{\beta}}{C R_0} \leq \frac{a M_0^2 B}{\varepsilon} = \frac{B}{2\alpha_{t+1}} = \frac{\lambda_{t+1}}{2}
$$

where we use $B = \frac{C R_0}{16 \ln \frac{4N}{\beta}}$ and $\varepsilon \leq \frac{a C M_0 R_0}{16 \ln \frac{4N}{\beta}}$. Next, we prove that $\|\nabla f(x^{t+1})\|_2 \leq \frac{\lambda_{t+1}}{2}$ when $\nu > 0$. For $t = 0$ we have

$$
\left\| \nabla f(x^1) \right\|_2 = \|\nabla f(z^0)\|_2 \overset{(3)}{\leq} M_\nu \|z^0 - x^*\|_2^\nu \leq M_\nu R_0^\nu = \frac{16 \varepsilon^{\frac{1-\nu}{1+\nu}} \ln \frac{4N}{\beta}}{a C M_\nu^{\frac{1-\nu}{1+\nu}} R_0^{1-\nu}} \leq \frac{B}{2\alpha_1} = \frac{\lambda_1}{2}
$$

since $\varepsilon^{\frac{1-\nu}{1+\nu}} \leq \frac{a C M_\nu^{\frac{1-\nu}{1+\nu}} R_0^{1-\nu}}{16 \ln \frac{4N}{\beta}}$. For $0 < t \leq T-1$ probability event $E_{T-1}$ implies

$$
\begin{aligned}
\|\nabla f(x^{t+1})\|_2 \ \leq \ & \|\nabla f(x^{t+1}) - \nabla f(y^t)\|_2 + \|\nabla f(y^t)\|_2 \\
\overset{(3)}{\leq} \ & M_\nu \|x^{t+1} - y^t\|_2^\nu + \left(\frac{1+\nu}{\nu}\right)^{\frac{\nu}{1+\nu}} M_\nu^{\frac{1}{1+\nu}} \left( f(y^t) - f(x^*) \right)^{\frac{\nu}{1+\nu}} \\
\overset{(28),(41)}{\leq} \ & M_\nu \left(\frac{\alpha_{t+1}}{A_t}\right)^\nu \|x^{t+1} - z^t\|_2^\nu + \left(\frac{1+\nu}{\nu}\right)^{\frac{\nu}{1+\nu}} M_\nu^{\frac{1}{1+\nu}} \left(\frac{C^2 R_0^2}{2A_t}\right)^{\frac{\nu}{1+\nu}} \\
= \ & \frac{\lambda_{t+1}}{2} \Bigg( \underbrace{\frac{2M_\nu}{\lambda_{t+1}} \left(\frac{\alpha_{t+1}}{A_t}\right)^\nu \|x^{t+1} - z^t\|_2^\nu}_{D_1} + \underbrace{\left(\frac{1+\nu}{\nu}\right)^{\frac{\nu}{1+\nu}} \frac{2M_\nu^{\frac{1}{1+\nu}}}{\lambda_{t+1}} \left(\frac{C^2 R_0^2}{2A_t}\right)^{\frac{\nu}{1+\nu}}}_{D_2} \Bigg).
\end{aligned}
$$

23

Next, we show that $D_1 + D_2 \leq 1$. Using the definition of $\lambda_{t+1}$, triangle inequality $\|x^{t+1} - z^t\|_2 \leq \|x^{t+1} - x^*\|_2 + \|z^t - x^*\|_2 \leq 2CR_0$, and lower bound (17) for $A_t$ (see Lemma A.3) we derive

$$
\begin{aligned}
D_1 &= \frac{2^{\nu+4} M_\nu \alpha_{t+1}^{1+\nu} \ln \frac{4N}{\beta}}{C^{1-\nu} R_0^{1-\nu} A_t^\nu} = \frac{2^{\nu+4} M_\nu (t+1)^{2\nu} (\varepsilon/2)^{1-\nu} \ln \frac{4N}{\beta}}{2^{2\nu} a^{1+\nu} C^{1-\nu} R_0^{1-\nu} M_\nu^2 A_t^\nu} \\[2mm]
&\overset{(17)}{\leq} \frac{2^3 (t+1)^{2\nu} \varepsilon^{1-\nu} \ln \frac{4N}{\beta}}{a^{1+\nu} C^{1-\nu} R_0^{1-\nu} M_\nu} \cdot \frac{2^{\frac{(1+3\nu)\nu}{1+\nu}} a^\nu M_\nu^{\frac{2\nu}{1+\nu}}}{t^{\frac{(1+3\nu)\nu}{1+\nu}} (\varepsilon/2)^{\frac{\nu(1-\nu)}{1+\nu}}} \\[2mm]
&= \frac{(t+1)^{2\nu}}{t^{\frac{\nu(1+3\nu)}{1+\nu}}} \cdot \frac{2^{3+2\nu} \varepsilon^{\frac{1-\nu}{1+\nu}} \ln \frac{4N}{\beta}}{a M_\nu^{\frac{1-\nu}{1+\nu}} C^{1-\nu} R_0^{1-\nu}} \leq \frac{2^{3+4\nu} t^{\frac{\nu(1-\nu)}{1+\nu}} \varepsilon^{\frac{1-\nu}{1+\nu}} \ln \frac{4N}{\beta}}{a M_\nu^{\frac{1-\nu}{1+\nu}} C^{1-\nu} R_0^{1-\nu}} \\[2mm]
&\overset{(39)}{\leq} \frac{2^{3+4\nu} \varepsilon^{\frac{1-\nu}{1+\nu}} \ln \frac{4N}{\beta}}{a M_\nu^{\frac{1-\nu}{1+\nu}} C^{1-\nu} R_0^{1-\nu}} \cdot \frac{2^{\frac{2\nu(1-\nu)(1+2\nu)}{(1+\nu)(1+3\nu)}} a^{\frac{\nu(1-\nu)}{1+3\nu}} C^{\frac{2\nu(1-\nu)}{1+3\nu}} R_0^{\frac{2\nu(1-\nu)}{1+3\nu}} M_\nu^{\frac{2\nu(1-\nu)}{(1+\nu)(1+3\nu)}}}{\varepsilon^{\frac{2\nu(1-\nu)}{(1+\nu)(1+3\nu)}}} \\[2mm]
&= \frac{2^{3+4\nu+\frac{2\nu(1-\nu)(1+2\nu)}{(1+\nu)(1+3\nu)}} \varepsilon^{\frac{1-\nu}{1+3\nu}} \ln \frac{4N}{\beta}}{a^{\frac{(1+\nu)^2}{1+3\nu}} M_\nu^{\frac{1-\nu}{1+3\nu}} C^{\frac{(1-\nu)(1+\nu)}{1+3\nu}} R_0^{\frac{(1-\nu)(1+\nu)}{1+3\nu}}} \overset{(37)}{\leq} \frac{1}{2^{\frac{3+6\nu-7\nu^2-2\nu^3}{(1+\nu)(1+3\nu)}} a^{\frac{\nu}{2}}}.
\end{aligned}
$$

Applying the same inequalities and $\left(\frac{1+\nu}{\nu}\right)^{\frac{\nu}{1+\nu}} \leq 2$ we estimate $D_2$:

$$
\begin{aligned}
D_2 &= \left(\frac{1+\nu}{\nu}\right)^{\frac{\nu}{1+\nu}} \frac{2^{4-\frac{\nu}{1+\nu}} M_\nu^{\frac{1}{1+\nu}} \alpha_{t+1} \ln \frac{4N}{\beta}}{C^{\frac{1-\nu}{1+\nu}} R_0^{\frac{1-\nu}{1+\nu}} A_t^{\frac{\nu}{1+\nu}}} \leq 2 \cdot \frac{2^{4-\frac{\nu}{1+\nu}} M_\nu^{\frac{1}{1+\nu}} \ln \frac{4N}{\beta}}{C^{\frac{1-\nu}{1+\nu}} R_0^{\frac{1-\nu}{1+\nu}} A_t^{\frac{\nu}{1+\nu}}} \cdot \frac{(t+1)^{\frac{2\nu}{1+\nu}} (\varepsilon/2)^{\frac{1-\nu}{1+\nu}}}{2^{\frac{2\nu}{1+\nu}} a M_\nu^{\frac{2}{1+\nu}}} \\[2mm]
&\leq \frac{2^{4-\frac{\nu}{1+\nu}} \cdot 2^{\frac{2\nu}{1+\nu}} t^{\frac{2\nu}{1+\nu}} \varepsilon^{\frac{1-\nu}{1+\nu}} \ln \frac{4N}{\beta}}{a C^{\frac{1-\nu}{1+\nu}} R_0^{\frac{1-\nu}{1+\nu}} M_\nu^{\frac{1}{1+\nu}} A_t^{\frac{\nu}{1+\nu}}} \\[2mm]
&\overset{(17)}{\leq} \frac{2^{4+\frac{\nu}{1+\nu}} t^{\frac{2\nu}{1+\nu}} \varepsilon^{\frac{1-\nu}{1+\nu}} \ln \frac{4N}{\beta}}{a C^{\frac{1-\nu}{1+\nu}} R_0^{\frac{1-\nu}{1+\nu}} M_\nu^{\frac{1}{1+\nu}}} \cdot \frac{2^{\frac{\nu(1+3\nu)}{(1+\nu)^2}} a^{\frac{\nu}{1+\nu}} M_\nu^{\frac{2\nu}{(1+\nu)^2}}}{t^{\frac{\nu(1+3\nu)}{(1+\nu)^2}} (\varepsilon/2)^{\frac{\nu(1-\nu)}{(1+\nu)^2}}} \\[2mm]
&= \frac{2^{4+\frac{3\nu}{1+\nu}} t^{\frac{\nu(1-\nu)}{(1+\nu)^2}} \varepsilon^{\frac{1-\nu}{(1+\nu)^2}} \ln \frac{4N}{\beta}}{a^{\frac{1}{1+\nu}} C^{\frac{1-\nu}{1+\nu}} R_0^{\frac{1-\nu}{1+\nu}} M_\nu^{\frac{1-\nu}{(1+\nu)^2}}} \\[2mm]
&\overset{(39)}{\leq} \frac{2^{4+\frac{3\nu}{1+\nu}} \varepsilon^{\frac{1-\nu}{(1+\nu)^2}} \ln \frac{4N}{\beta}}{a^{\frac{1}{1+\nu}} C^{\frac{1-\nu}{1+\nu}} R_0^{\frac{1-\nu}{1+\nu}} M_\nu^{\frac{1-\nu}{(1+\nu)^2}}} \cdot \frac{2^{\frac{2\nu(1+2\nu)(1-\nu)}{(1+\nu)^2(1+3\nu)}} a^{\frac{\nu(1-\nu)}{(1+\nu)(1+3\nu)}} C^{\frac{2\nu(1-\nu)}{(1+\nu)(1+3\nu)}} R_0^{\frac{2\nu(1-\nu)}{(1+\nu)(1+3\nu)}} M_\nu^{\frac{2\nu(1-\nu)}{(1+\nu)^2(1+3\nu)}}}{\varepsilon^{\frac{2\nu(1-\nu)}{(1+\nu)^2(1+3\nu)}}} \\[2mm]
&= \frac{2^{4+\frac{3\nu}{1+\nu}+\frac{2\nu(1+2\nu)(1-\nu)}{(1+\nu)^2(1+3\nu)}} \varepsilon^{\frac{1-\nu}{(1+\nu)(1+3\nu)}} \ln \frac{4N}{\beta}}{a^{\frac{1+\nu}{1+3\nu}} C^{\frac{1-\nu}{1+3\nu}} R_0^{\frac{1-\nu}{1+3\nu}} M_\nu^{\frac{1-\nu}{(1+\nu)(1+3\nu)}}} \overset{(37)}{\leq} \frac{1}{2^{\frac{2+5\nu+\nu^3}{(1+\nu)^2(1+3\nu)}}}.
\end{aligned}
$$

Combining the upper bounds for $D_1$ and $D_2$ we get

$$
D_1 + D_2 \leq \frac{1}{2^{\frac{3+6\nu-7\nu^2-2\nu^3}{(1+\nu)(1+3\nu)}} a^{\frac{\nu}{2}}} + \frac{1}{2^{\frac{2+5\nu+\nu^3}{(1+\nu)^2(1+3\nu)}}}.
$$

Since $\frac{2+5\nu+\nu^3}{(1+\nu)^2(1+3\nu)}$ is a decreasing function of $\nu$ for $\nu \in [0,1]$ we continue as

$$
D_1 + D_2 \leq \frac{1}{2^{\frac{3+6\nu-7\nu^2-2\nu^3}{(1+\nu)(1+3\nu)}} a^{\frac{\nu}{2}}} + \frac{1}{\sqrt{2}}.
$$

Next, we use $a \geq 16384 \ln^2 \frac{4N}{\beta} \geq 2^{10}$ and obtain

$$
D_1 + D_2 \leq \frac{1}{2^{\frac{3+11\nu+13\nu^2+13\nu^3}{(1+\nu)(1+3\nu)}}} + \frac{1}{\sqrt{2}}.
$$

24

One can numerically verify that $\frac{1}{2^{\frac{3+11\nu+13\nu^2+13\nu^3}{(1+\nu)(1+3\nu)}}} + \frac{1}{\sqrt{2}}$ is smaller than 1 for $\nu \in [0,1]$. Putting all together we conclude that probability event $E_{T-1}$ implies

$$\|\nabla f(x^{t+1})\|_2 \leq \frac{\lambda_{t+1}}{2} \tag{45}$$

for all $t = 0, 1, \ldots, T-1$. Having inequality (45) in hand we show in the rest of the proof that (41) holds for $t = T$ with large enough probability. First of all, we introduce new random variables:

$$\eta_l = \begin{cases} x^* - z^l, & \text{if } \|x^* - z^l\|_2 \leq CR_0, \\ 0, & \text{otherwise}, \end{cases} \quad \text{and} \quad \zeta_l = \begin{cases} \nabla f(x^{l+1}), & \text{if } \|\nabla f(x^{l+1})\|_2 \leq \frac{B}{2\alpha_{l+1}}, \\ 0, & \text{otherwise}, \end{cases} \tag{46}$$

for $l = 0, 1, \ldots T-1$. Note that these random variables are bounded with probability 1, i.e. with probability 1 we have

$$\|\eta_l\|_2 \leq CR_0 \quad \text{and} \quad \|\zeta_l\|_2 \leq \frac{B}{2\alpha_{l+1}}. \tag{47}$$

Secondly, we use the introduced notation and get that $E_{T-1}$ implies

$$R_T^2 \overset{(44),(41),(45),(46)}{\leq} R_0^2 + 2\sum_{l=0}^{T-1} \alpha_{l+1} \langle \theta_{l+1}, \eta_l \rangle + 2\sum_{l=0}^{T-1} \alpha_{l+1}^2 \|\theta_{l+1}\|_2^2 + 2\sum_{l=0}^{T-1} \alpha_{l+1}^2 \langle \theta_{l+1}, \zeta_l \rangle + \frac{A_N \varepsilon}{2}$$

$$= R_0^2 + \sum_{l=0}^{T-1} \alpha_{l+1} \langle \theta_{l+1}, 2\eta_l + 2\alpha_{l+1}\zeta_l \rangle + 2\sum_{l=0}^{T-1} \alpha_{l+1}^2 \|\theta_{l+1}\|_2^2 + \frac{A_N \varepsilon}{2}.$$

Finally, we do some preliminaries in order to apply Bernstein's inequality (see Lemma A.2) and obtain that $E_{T-1}$ implies

$$R_T^2 \overset{(11)}{\leq} R_0^2 + \underbrace{\sum_{l=0}^{T-1} \alpha_{l+1} \langle \theta_{l+1}^u, 2\eta_l + 2\alpha_{l+1}\zeta_l \rangle}_{\textcircled{1}} + \underbrace{\sum_{l=0}^{T-1} \alpha_{l+1} \langle \theta_{l+1}^b, 2\eta_l + 2\alpha_{l+1}\zeta_l \rangle}_{\textcircled{2}}$$

$$+ \underbrace{\sum_{l=0}^{T-1} 4\alpha_{l+1}^2 \left( \|\theta_{l+1}^u\|_2^2 - \mathbb{E}_{\boldsymbol{\xi}^l} \left[ \|\theta_{l+1}^u\|_2^2 \right] \right)}_{\textcircled{3}} + \underbrace{\sum_{l=0}^{T-1} 4\alpha_{l+1}^2 \mathbb{E}_{\boldsymbol{\xi}^l} \left[ \|\theta_{l+1}^u\|_2^2 \right]}_{\textcircled{4}}$$

$$+ \underbrace{\sum_{l=0}^{T-1} 4\alpha_{l+1}^2 \|\theta_{l+1}^b\|_2^2}_{\textcircled{5}} + \frac{A_N \varepsilon}{2} \tag{48}$$

where we introduce new notations:

$$\theta_{l+1}^u \overset{\text{def}}{=} \widetilde{\nabla} f(x^{l+1}, \boldsymbol{\xi}^l) - \mathbb{E}_{\boldsymbol{\xi}^l} \left[ \widetilde{\nabla} f(x^{l+1}, \boldsymbol{\xi}^l) \right], \quad \theta_{l+1}^b \overset{\text{def}}{=} \mathbb{E}_{\boldsymbol{\xi}^l} \left[ \widetilde{\nabla} f(x^{l+1}, \boldsymbol{\xi}^l) \right] - \nabla f(x^{l+1}), \tag{49}$$

$$\theta_{l+1} \overset{(23)}{=} \theta_{l+1}^u + \theta_{l+1}^b.$$

It remains to provide tight upper bounds for $\textcircled{1}$, $\textcircled{2}$, $\textcircled{3}$, $\textcircled{4}$ and $\textcircled{5}$, i.e. in the remaining part of the proof we show that $\textcircled{1} + \textcircled{2} + \textcircled{3} + \textcircled{4} + \textcircled{5} \leq \delta C^2 R_0^2$ for some $\delta < 1$.

**Upper bound for $\textcircled{1}$.** First of all, since $\mathbb{E}_{\boldsymbol{\xi}^l}[\theta_{l+1}^u] = 0$ summands in $\textcircled{1}$ are conditionally unbiased:

$$\mathbb{E}_{\boldsymbol{\xi}^l} \left[ \alpha_{l+1} \langle \theta_{l+1}^u, 2\eta_l + 2\alpha_{l+1}\zeta_l \rangle \right] = 0.$$

Secondly, these summands are bounded with probability 1:

$$\left| \alpha_{l+1} \langle \theta_{l+1}^u, 2\eta_l + 2\alpha_{l+1}\zeta_l \rangle \right| \leq \alpha_{l+1} \|\theta_{l+1}^u\|_2 \|2\eta_l + 2\alpha_{l+1}\zeta_l\|_2$$

$$\overset{(29),(47)}{\leq} 2\alpha_{l+1} \lambda_{l+1} (2CR_0 + B) = 2B(2CR_0 + B)$$

$$= \left( 1 + \frac{1}{32 \ln \frac{4N}{\beta}} \right) \frac{C^2 R_0^2}{4 \ln \frac{4N}{\beta}} \overset{(33)}{\leq} \left( 1 + \frac{1}{64} \right) \frac{C^2 R_0^2}{4 \ln \frac{4N}{\beta}}.$$

Finally, one can bound conditional variances $\sigma_l^2 \overset{\text{def}}{=} \mathbb{E}_{\boldsymbol{\xi}^l} \left[ \alpha_{l+1}^2 \left\langle \theta_{l+1}^u, 2\eta_l + 2\alpha_{l+1}\zeta_l \right\rangle^2 \right]$ in the following way:

$$
\begin{aligned}
\sigma_l^2 &\leq \mathbb{E}_{\boldsymbol{\xi}^l} \left[ \alpha_{l+1}^2 \left\| \theta_{l+1}^u \right\|_2^2 \| 2\eta_l + 2\alpha_{l+1}\zeta_l \|_2^2 \right] \\
&\overset{(47)}{\leq} \alpha_{l+1}^2 \mathbb{E}_{\boldsymbol{\xi}^l} \left[ \left\| \theta_{l+1}^u \right\|_2^2 \right] (2CR_0 + B)^2 = 4\alpha_{l+1}^2 \mathbb{E}_{\boldsymbol{\xi}^l} \left[ \left\| \theta_{l+1}^u \right\|_2^2 \right] \left( 1 + \frac{1}{32 \ln \frac{4N}{\beta}} \right)^2 C^2 R_0^2 \\
&\overset{(33)}{\leq} 4\alpha_{l+1}^2 \mathbb{E}_{\boldsymbol{\xi}^l} \left[ \left\| \theta_{l+1}^u \right\|_2^2 \right] \left( 1 + \frac{1}{64} \right)^2 C^2 R_0^2.
\end{aligned}
\tag{50}
$$

In other words, sequence $\left\{ \alpha_{l+1} \left\langle \theta_{l+1}^u, 2\eta_l + 2\alpha_{l+1}\zeta_l \right\rangle \right\}_{l \geq 0}$ is a bounded martingale difference sequence with bounded conditional variances $\{\sigma_l^2\}_{l \geq 0}$. Therefore, we can apply Bernstein's inequality, i.e. we apply Lemma A.2 with $X_l = \alpha_{l+1} \left\langle \theta_{l+1}^u, 2\eta_l + 2\alpha_{l+1}\zeta_l \right\rangle$, $c = \left( 1 + \frac{1}{64} \right) \frac{C^2 R_0^2}{4 \ln \frac{4N}{\beta}}$ and $F = \frac{c^2 \ln \frac{4N}{\beta}}{18}$ and get that for all $b > 0$

$$
\mathbb{P} \left\{ \left| \sum_{l=0}^{T-1} X_l \right| > b \text{ and } \sum_{l=0}^{T-1} \sigma_l^2 \leq F \right\} \leq 2 \exp \left( -\frac{b^2}{2F + 2cb/3} \right)
$$

or, equivalently, with probability at least $1 - 2 \exp \left( -\frac{b^2}{2F + 2cb/3} \right)$

$$
\text{either } \sum_{l=0}^{T-1} \sigma_l^2 > F \quad \text{or} \quad \underbrace{\left| \sum_{l=0}^{T-1} X_l \right|}_{|\textcircled{1}|} \leq b.
$$

The choice of $F$ will be clarified below. Let us now choose $b$ in such a way that $2 \exp \left( -\frac{b^2}{2F + 2cb/3} \right) = \frac{\beta}{2N}$. This implies that $b$ is the positive root of the quadratic equation

$$
b^2 - \frac{2c \ln \frac{4N}{\beta}}{3} b - 2F \ln \frac{4N}{\beta} = 0,
$$

hence

$$
\begin{aligned}
b &= \frac{c \ln \frac{4N}{\beta}}{3} + \sqrt{\frac{c^2 \ln^2 \frac{4N}{\beta}}{9} + 2F \ln \frac{4N}{\beta}} \leq \frac{c \ln \frac{4N}{\beta}}{3} + \sqrt{\frac{2c^2 \ln^2 \frac{4N}{\beta}}{9}} \\
&= \frac{1 + \sqrt{2}}{3} c \ln \frac{4N}{\beta} \leq c \ln \frac{4N}{\beta} = \left( 1 + \frac{1}{64} \right) \frac{C^2 R_0^2}{4} = \left( \frac{1}{4} + \frac{1}{256} \right) C^2 R_0^2.
\end{aligned}
$$

That is, with probability at least $1 - \frac{\beta}{2N}$

$$
\underbrace{\text{either } \sum_{l=0}^{T-1} \sigma_l^2 > F \quad \text{or} \quad |\textcircled{1}| \leq \left( \frac{1}{4} + \frac{1}{256} \right) C^2 R_0^2}_{\text{probability event } E_{\textcircled{1}}}.
$$

Next, we notice that probability event $E_{T-1}$ implies that

$$
\begin{aligned}
\sum_{l=0}^{T-1} \sigma_l^2 &\overset{(50)}{\leq} 4 \left( 1 + \frac{1}{64} \right)^2 C^2 R_0^2 \sum_{l=0}^{T-1} \alpha_{l+1}^2 \mathbb{E}_{\boldsymbol{\xi}^l} \left[ \left\| \theta_{l+1}^u \right\|_2^2 \right] \\
&\overset{(32),(45)}{\leq} 72 \left( 1 + \frac{1}{64} \right)^2 \sigma^2 C^2 R_0^2 \sum_{l=0}^{T-1} \frac{\alpha_{l+1}^2}{m_l} \\
&\overset{(34)}{\leq} \frac{\left( 1 + \frac{1}{64} \right)^2 C^4 R_0^4}{288 \ln \frac{4N}{\beta}} \sum_{l=0}^{T-1} \frac{1}{N} \\
&\overset{T \leq N}{\leq} \frac{\left( 1 + \frac{1}{64} \right)^2 C^4 R_0^4}{288 \ln \frac{4N}{\beta}} = \frac{c^2 \ln \frac{4N}{\beta}}{18} = F.
\end{aligned}
$$

26

**Upper bound for ②.** The probability event $E_{T-1}$ implies

$$
\begin{aligned}
\alpha_{l+1}\left\langle\theta_{l+1}^b, 2\eta_l + 2\alpha_{l+1}\zeta_l\right\rangle
&\leq \alpha_{l+1}\left\|\theta_{l+1}^b\right\|_2\left\|2\eta_l + 2\alpha_{l+1}\zeta_l\right\|_2 \\
&\overset{(30),(47)}{\leq} \alpha_{l+1}\cdot\frac{4\sigma^2}{m_l\lambda_{l+1}}(2CR_0 + B) \\
&= \frac{4\sigma^2\alpha_{l+1}^2}{m_l}\left(1+\frac{2CR_0}{B}\right) = \frac{4\sigma^2\alpha_{l+1}^2\left(1+32\ln\frac{4N}{\beta}\right)}{m_l} \\
&\overset{(34)}{\leq} \frac{4\left(\frac{1}{\ln\frac{4N}{\beta}}+32\right)C^2R_0^2}{20736N} \overset{(33)}{\leq} \frac{11C^2R_0^2}{1728N}.
\end{aligned}
$$

This implies that

$$
② = \sum_{l=0}^{T-1}\alpha_{l+1}\left\langle\theta_{l+1}^b, 2\eta_l + 2\alpha_{l+1}\zeta_l\right\rangle \overset{T\leq N}{\leq} \frac{11C^2R_0^2}{1728}.
$$

**Upper bound for ③.** We derive the upper bound for ③ using the same technique as for ①. First of all, we notice that the summands in ③ are conditionally unbiased:

$$
\mathbb{E}_{\boldsymbol{\xi}^l}\left[4\alpha_{l+1}^2\left(\|\theta_{l+1}^u\|_2^2 - \mathbb{E}_{\boldsymbol{\xi}^l}\left[\|\theta_{l+1}^u\|_2^2\right]\right)\right] = 0.
$$

Secondly, the summands are bounded with probability 1:

$$
\begin{aligned}
\left|4\alpha_{l+1}^2\left(\|\theta_{l+1}^u\|_2^2 - \mathbb{E}_{\boldsymbol{\xi}^l}\left[\|\theta_{l+1}^u\|_2^2\right]\right)\right|
&\leq 4\alpha_{l+1}^2\left(\|\theta_{l+1}^u\|_2^2 + \mathbb{E}_{\boldsymbol{\xi}^l}\left[\|\theta_{l+1}^u\|_2^2\right]\right) \\
&\overset{(29)}{\leq} 4\alpha_{l+1}^2\left(4\lambda_{l+1}^2 + 4\lambda_{l+1}^2\right) \\
&= 32B^2 = \frac{C^2R_0^2}{8\ln^2\frac{4N}{\beta}} \overset{(33)}{\leq} \frac{C^2R_0^2}{16\ln\frac{4N}{\beta}} \overset{\text{def}}{=} c_1. \quad (51)
\end{aligned}
$$

Finally, one can bound conditional variances $\hat{\sigma}_l^2 \overset{\text{def}}{=} \mathbb{E}_{\boldsymbol{\xi}^l}\left[\left|4\alpha_{l+1}^2\left(\|\theta_{l+1}^u\|_2^2 - \mathbb{E}_{\boldsymbol{\xi}^l}\left[\|\theta_{l+1}^u\|_2^2\right]\right)\right|^2\right]$ in the following way:

$$
\begin{aligned}
\hat{\sigma}_l^2 &\overset{(51)}{\leq} c_1\mathbb{E}_{\boldsymbol{\xi}^l}\left[\left|4\alpha_{l+1}^2\left(\|\theta_{l+1}^u\|_2^2 - \mathbb{E}_{\boldsymbol{\xi}^l}\left[\|\theta_{l+1}^u\|_2^2\right]\right)\right|\right] \\
&\leq 4c_1\alpha_{l+1}^2\mathbb{E}_{\boldsymbol{\xi}^l}\left[\|\theta_{l+1}^u\|_2^2 + \mathbb{E}_{\boldsymbol{\xi}^l}\left[\|\theta_{l+1}^u\|_2^2\right]\right] = 8c_1\alpha_{l+1}^2\mathbb{E}_{\boldsymbol{\xi}^l}\left[\|\theta_{l+1}^u\|_2^2\right]. \quad (52)
\end{aligned}
$$

In other words, sequence $\left\{4\alpha_{l+1}^2\left(\|\theta_{l+1}^u\|_2^2 - \mathbb{E}_{\boldsymbol{\xi}^l}\left[\|\theta_{l+1}^u\|_2^2\right]\right)\right\}_{l\geq 0}$ is bounded martingale difference sequence with bounded conditional variances $\{\hat{\sigma}_l^2\}_{l\geq 0}$. Therefore, we can apply Bernstein's inequality, i.e. we apply Lemma A.2 with $X_l = \hat{X}_l = 4\alpha_{l+1}^2\left(\|\theta_{l+1}^u\|_2^2 - \mathbb{E}_{\boldsymbol{\xi}^l}\left[\|\theta_{l+1}^u\|_2^2\right]\right)$, $c = c_1 = \frac{C^2R_0^2}{16\ln\frac{4N}{\beta}}$ and $F = F_1 = \frac{c_1^2\ln\frac{4N}{\beta}}{18}$ and get that for all $b > 0$

$$
\mathbb{P}\left\{\left|\sum_{l=0}^{T-1}\hat{X}_l\right| > b \text{ and } \sum_{l=0}^{T-1}\hat{\sigma}_l^2 \leq F_1\right\} \leq 2\exp\left(-\frac{b^2}{2F_1 + 2c_1b/3}\right)
$$

or, equivalently, with probability at least $1 - 2\exp\left(-\frac{b^2}{2F_1 + 2c_1b/3}\right)$

$$
\text{either } \sum_{l=0}^{T-1}\hat{\sigma}_l^2 > F_1 \quad \text{or} \quad \underbrace{\left|\sum_{l=0}^{T-1}\hat{X}_l\right|}_{|③|} \leq b.
$$

As in our derivations of the upper bound for ① we choose such $b$ that $2\exp\left(-\frac{b^2}{2F_1 + 2c_1b/3}\right) = \frac{\beta}{2N}$, i.e.,

$$
b = \frac{c_1\ln\frac{4N}{\beta}}{3} + \sqrt{\frac{c_1^2\ln^2\frac{4N}{\beta}}{9} + 2F_1\ln\frac{4N}{\beta}} \leq \frac{1+\sqrt{2}}{3}c_1\ln\frac{4N}{\beta} \leq \frac{C^2R_0^2}{16}.
$$

That is, with probability at least $1 - \frac{\beta}{2N}$

$$\underbrace{\text{either} \sum_{l=0}^{T-1} \hat{\sigma}_l^2 > F_1 \quad \text{or} \quad |\text{③}| \le \frac{C^2 R_0^2}{16}}_{\text{probability event } E_\text{③}}.$$

Next, we notice that probability event $E_{T-1}$ implies that

$$\sum_{l=0}^{T-1} \hat{\sigma}_l^2 \overset{(52)}{\le} 8c_1 \sum_{l=0}^{T-1} \alpha_{l+1}^2 \mathbb{E}_{\boldsymbol{\xi}^l}\left[\left\|\theta_{l+1}^u\right\|_2^2\right]$$

$$\overset{(32),(45)}{\le} \frac{9\sigma^2 C^2 R_0^2}{\ln\frac{4N}{\beta}} \sum_{l=0}^{T-1} \frac{\alpha_{l+1}^2}{m_l} \overset{(34)}{\le} \frac{C^4 R_0^4}{2304 \ln^2\frac{4N}{\beta}} \sum_{l=0}^{T-1} \frac{1}{N}$$

$$\overset{T\le N}{\le} \frac{C^4 R_0^4}{2304 \ln^2\frac{4N}{\beta}} \overset{(33)}{\le} \frac{C^4 R_0^4}{4608 \ln\frac{4N}{\beta}} = \frac{c_1^2 \ln\frac{4N}{\beta}}{18} = F_1.$$

**Upper bound for ④.** The probability event $E_{T-1}$ implies

$$\text{④} = \sum_{l=0}^{T-1} 4\alpha_{l+1}^2 \mathbb{E}_{\boldsymbol{\xi}^l}\left[\|\theta_{l+1}^u\|_2^2\right] \overset{(32),(45)}{\le} \sum_{l=0}^{T-1} \frac{72\alpha_{l+1}^2 \sigma^2}{m_l} \overset{(34)}{\le} \sum_{l=0}^{T-1} \frac{C^2 R_0^2}{288N \ln\frac{4N}{\beta}}$$

$$\overset{T\le N}{\le} \frac{C^2 R_0^2}{288 \ln\frac{4N}{\beta}} \overset{(33)}{\le} \frac{C^2 R_0^2}{576}.$$

**Upper bound for ⑤.** Again, we use corollaries of probability event $E_{T-1}$:

$$\text{⑤} = \sum_{l=0}^{T-1} 4\alpha_{l+1}^2 \|\theta_{l+1}^b\|_2^2 \overset{(30),(45)}{\le} \sum_{l=0}^{T-1} \frac{64\alpha_{l+1}^2 \sigma^4}{m_l^2 \lambda_{l+1}^2} = \frac{64\sigma^4}{B^2} \sum_{l=0}^{T-1} \frac{\alpha_{l+1}^4}{m_l^2}$$

$$\overset{(34),(35)}{\le} \frac{256 \cdot 64\sigma^4 \ln^2\frac{4N}{\beta}}{C^2 R_0^2} \sum_{l=0}^{T-1} \frac{C^4 R_0^4}{20736^2 N^2 \sigma^4 \ln^2\frac{4N}{\beta}} \overset{T\le N}{\le} \frac{C^2 R_0^2}{26244}.$$

Now we summarize all bounds that we have: probability event $E_{T-1}$ implies

$$R_T^2 \overset{(44)}{\le} R_0^2 + 2\sum_{l=0}^{T-1} \alpha_{l+1}\left\langle\theta_{l+1}, x^* - z^l\right\rangle + 2\sum_{l=0}^{k-1} \alpha_{l+1}^2\left\langle\theta_{l+1}, \nabla f(x^{l+1})\right\rangle + 2\sum_{l=0}^{T-1} \alpha_{l+1}^2\|\theta_{l+1}\|_2^2 + \frac{A_N\varepsilon}{2}$$

$$\overset{(48)}{\le} R_0^2 + \text{①} + \text{②} + \text{③} + \text{④} + \text{⑤} + \frac{A_N\varepsilon}{2},$$

$$\text{②} \le \frac{11C^2 R_0^2}{1728}, \quad \text{④} \le \frac{CR_0^2}{576}, \quad \text{⑤} \le \frac{C^2 R_0^2}{26244},$$

$$\sum_{l=0}^{T-1} \sigma_l^2 \le F, \quad \sum_{l=0}^{T-1} \hat{\sigma}_l^2 \le F_1$$

and

$$\mathbb{P}\{E_{T-1}\} \ge 1 - \frac{(T-1)\beta}{N}, \quad \mathbb{P}\{E_\text{①}\} \ge 1 - \frac{\beta}{2N}, \quad \mathbb{P}\{E_\text{③}\} \ge 1 - \frac{\beta}{2N},$$

where

$$E_\text{①} = \left\{\text{either} \sum_{l=0}^{T-1} \sigma_l^2 > F \quad \text{or} \quad |\text{①}| \le \left(\frac{1}{4} + \frac{1}{256}\right) C^2 R_0^2\right\},$$

$$E_\text{③} = \left\{\text{either} \sum_{l=0}^{T-1} \hat{\sigma}_l^2 > F_1 \quad \text{or} \quad |\text{③}| \le \frac{C^2 R_0^2}{16}\right\}.$$

Moreover, since $N \overset{(39)}{\leq} \dfrac{2^{\frac{1+\nu}{1+3\nu}} a^{\frac{1+\nu}{1+3\nu}} C^{\frac{2(1+\nu)}{1+3\nu}} R_0^{\frac{2(1+\nu)}{1+3\nu}} M_\nu^{\frac{2}{1+3\nu}}}{\varepsilon^{\frac{2}{1+3\nu}}} + 1$ and $\varepsilon \overset{(36)}{\leq} \dfrac{2^{\frac{1+\nu}{2}} a^{\frac{1+\nu}{2}} C^{1+\nu} R_0^{1+\nu} M_\nu}{100^{\frac{1+3\nu}{2}}}$ we have

$$
\begin{aligned}
\frac{A_N \varepsilon}{2} &\overset{(19)}{\leq} \frac{N^{\frac{1+3\nu}{1+\nu}} \varepsilon^{\frac{2}{1+\nu}}}{4aM_\nu^{\frac{2}{1+\nu}}} \overset{(39)}{\leq} \left( \frac{2^{\frac{1+\nu}{1+3\nu}} a^{\frac{1+\nu}{1+3\nu}} C^{\frac{2(1+\nu)}{1+3\nu}} R_0^{\frac{2(1+\nu)}{1+3\nu}} M_\nu^{\frac{2}{1+3\nu}}}{\varepsilon^{\frac{2}{1+3\nu}}} + 1 \right)^{\frac{1+3\nu}{1+\nu}} \frac{\varepsilon^{\frac{2}{1+\nu}}}{4aM_\nu^{\frac{2}{1+\nu}}} \\
&\overset{(36)}{\leq} \left( \frac{101}{100} \right)^{\frac{1+3\nu}{1+\nu}} \frac{C^2 R_0^2}{2} \leq \frac{10201 C^2 R_0^2}{20000}.
\end{aligned}
$$

Taking into account these inequalities we get that probability event $E_{T-1} \cap E_{①} \cap E_{③}$ implies

$$
\begin{aligned}
R_T^2 &\overset{(44)}{\leq} R_0^2 + 2\sum_{l=0}^{T-1} \alpha_{l+1} \left\langle \theta_{l+1}, x^* - z^l \right\rangle + 2\sum_{l=0}^{k-1} \alpha_{l+1}^2 \left\langle \theta_{l+1}, \nabla f(x^{l+1}) \right\rangle + 2\sum_{l=0}^{T-1} \alpha_{l+1}^2 \|\theta_{l+1}\|_2^2 + \frac{A_N \varepsilon}{2} \\
&\leq \left( 1 + \left( \frac{1}{4} + \frac{1}{256} + \frac{11}{1728} + \frac{1}{16} + \frac{1}{576} + \frac{1}{26244} + \frac{10201}{20000} \right) C^2 \right) R_0^2 \\
&\overset{(39)}{\leq} C^2 R_0^2. \tag{53}
\end{aligned}
$$

Moreover, using union bound we derive

$$
\mathbb{P}\left\{ E_{T-1} \cap E_{①} \cap E_{③} \right\} = 1 - \mathbb{P}\left\{ \overline{E}_{T-1} \cup \overline{E}_{①} \cup \overline{E}_{③} \right\} \geq 1 - \frac{T\beta}{N}. \tag{54}
$$

That is, by definition of $E_T$ and $E_{T-1}$ we have proved that

$$
\mathbb{P}\{E_T\} \overset{(53)}{\geq} \mathbb{P}\left\{ E_{T-1} \cap E_{①} \cap E_{③} \right\} \overset{(54)}{\geq} 1 - \frac{T\beta}{N},
$$

which implies that for all $k = 0, 1, \dots, N$ we have $\mathbb{P}\{E_k\} \geq 1 - \frac{k\beta}{N}$. Then, for $k = N$ we have that with probability at least $1 - \beta$

$$
\begin{aligned}
A_N \left( f(y^N) - f(x^*) \right) &\overset{(42)}{\leq} \frac{1}{2}\|z^0 - z\|_2^2 - \frac{1}{2}\|z^N - z\|_2^2 + \sum_{k=0}^{N-1} \alpha_{k+1} \left\langle \theta_{k+1}, z - z^k \right\rangle \\
&\quad + \sum_{k=0}^{N-1} \alpha_{k+1}^2 \|\theta_{k+1}\|_2^2 + \sum_{k=0}^{N-1} \alpha_{k+1}^2 \left\langle \theta_{k+1}, \nabla f(x^{k+1}) \right\rangle + \frac{A_N \varepsilon}{4} \\
&\overset{(41)}{\leq} \frac{C^2 R_0^2}{2}.
\end{aligned}
$$

Since $A_N \overset{(17)}{\geq} \dfrac{N^{\frac{1+3\nu}{1+\nu}} (\varepsilon/2)^{\frac{1-\nu}{1+\nu}}}{2^{\frac{1+3\nu}{1+\nu}} aM_\nu^{\frac{2}{1+\nu}}}$ we get that with probability at least $1 - \beta$

$$
f(y^N) - f(x^*) \leq \frac{4aC^2 R_0^2 M_\nu^{\frac{2}{1+\nu}}}{N^{\frac{1+3\nu}{1+\nu}} \varepsilon^{\frac{1-\nu}{1+\nu}}}.
$$

29

In other words, `clipped-SSTM` with $a = 16384 \ln^2 \frac{4N}{\beta}$ achieves $f(y^N) - f(x^*) \leq \varepsilon$ with probability at least $1 - \beta$ after $\mathcal{O}\left( \frac{M_\nu^{\frac{2}{1+3\nu}} R_0^{\frac{2(1+\nu)}{1+3\nu}}}{\varepsilon^{\frac{2}{1+3\nu}}} \ln^{\frac{2(1+\nu)}{1+3\nu}} \frac{M_\nu^{\frac{2}{1+3\nu}} R_0^{\frac{2(1+\nu)}{1+3\nu}}}{\varepsilon^{\frac{2}{1+3\nu}} \beta} \right)$ iterations and requires

$$
\begin{aligned}
\sum_{k=0}^{N-1} m_k &\stackrel{(34)}{=} \sum_{k=0}^{N-1} \mathcal{O}\left( \max\left\{ 1, \frac{\sigma^2 \alpha_{k+1}^2 N \ln \frac{N}{\beta}}{R_0^2} \right\} \right) \\
&= \mathcal{O}\left( \max\left\{ N, \sum_{k=0}^{N-1} \frac{\sigma^2 (k+1)^{\frac{4\nu}{1+\nu}} \varepsilon^{\frac{2(1-\nu)}{1+\nu}} N \ln \frac{N}{\beta}}{M_\nu^{\frac{4}{1+\nu}} R_0^2 a^2} \right\} \right) \\
&\stackrel{(35)}{=} \mathcal{O}\left( \max\left\{ N, \frac{\sigma^2 \varepsilon^{\frac{2(1-\nu)}{1+\nu}} N^{\frac{2(1+3\nu)}{1+\nu}}}{M_\nu^{\frac{4}{1+\nu}} R_0^2 \ln^3 \frac{N}{\beta}} \right\} \right) \\
&= \mathcal{O}\left( \max\left\{ \frac{M_\nu^{\frac{2}{1+3\nu}} R_0^{\frac{2(1+\nu)}{1+3\nu}}}{\varepsilon^{\frac{2}{1+3\nu}}} \ln^{\frac{2(1+\nu)}{1+3\nu}} \frac{M_\nu^{\frac{2}{1+3\nu}} R_0^{\frac{2(1+\nu)}{1+3\nu}}}{\varepsilon^{\frac{2}{1+3\nu}} \beta}, \frac{\sigma^2 R_0^2}{\varepsilon^2} \ln \frac{M_\nu^{\frac{2}{1+3\nu}} R_0^{\frac{2(1+\nu)}{1+3\nu}}}{\varepsilon^{\frac{2}{1+3\nu}} \beta} \right\} \right).
\end{aligned}
$$

oracle calls. $\qquad\square$

### B.1.3  On the batchsizes and numerical constants

The obtained complexity result is discussed in details in the main part of the paper. Here we discuss the choice of the parameters. For convenience, we provide all assumptions from Thm. B.1 on the parameters below:

$$
\ln \frac{4N}{\beta} \geq 2 \tag{55}
$$

$$
\alpha = \frac{(\varepsilon/2)^{\frac{1-\nu}{1+\nu}}}{2^{\frac{2\nu}{1+\nu}} a M_\nu^{\frac{2}{1+\nu}}}, \quad m_k = \max\left\{ 1, \frac{20736 N \sigma^2 \alpha_{k+1}^2 \ln \frac{4N}{\beta}}{C^2 R_0^2} \right\}, \tag{56}
$$

$$
B = \frac{C R_0}{16 \ln \frac{4N}{\beta}}, \quad a \geq 16384 \ln^2 \frac{4N}{\beta}, \tag{57}
$$

$$
\varepsilon^{\frac{1-\nu}{1+\nu}} \leq \frac{a C M_\nu^{\frac{1-\nu}{1+\nu}} R_0^{1-\nu}}{16 \ln \frac{4N}{\beta}}, \quad \varepsilon \leq \frac{2^{\frac{1+\nu}{2}} a^{\frac{1+\nu}{2}} C^{1+\nu} R_0^{1+\nu} M_\nu}{100^{\frac{1+3\nu}{2}}}, \tag{58}
$$

$$
\varepsilon^{\frac{1-\nu}{1+3\nu}} \leq \min\left\{ \frac{a^{\frac{2+3\nu-\nu^2}{2(1+3\nu)}}}{2^{2+4\nu+\frac{3+8\nu-5\nu^2-6\nu^3}{(1+\nu)(1+3\nu)}} \ln \frac{4N}{\beta}}, \frac{a^{\frac{(1+\nu)^2}{1+3\nu}}}{2^{4+7\nu+\frac{2+7\nu+2\nu^2-3\nu^3}{(1+\nu)(1+3\nu)}} \ln^{1+\nu} \frac{4N}{\beta}} \right\} C^{\frac{1-\nu^2}{1+3\nu}} R_0^{\frac{1-\nu^2}{1+3\nu}} M_\nu^{\frac{1-\nu}{1+3\nu}} \tag{59}
$$

$$
N = \left\lceil \frac{2^{\frac{1+\nu}{1+3\nu}} a^{\frac{1+\nu}{1+3\nu}} C^{\frac{2(1+\nu)}{1+3\nu}} R_0^{\frac{2(1+\nu)}{1+3\nu}} M_\nu^{\frac{2}{1+3\nu}}}{\varepsilon^{\frac{2}{1+3\nu}}} \right\rceil + 1, \quad C = \sqrt{7}. \tag{60}
$$

We emphasize that (55), (58), and (59) are not restrictive at all since the target accuracy $\varepsilon$ and confidence level $\beta$ are often chosen to be small enough, whereas $a$ can be made large enough.

Next, one can notice that the assumptions on parameter $a$ and batchsize $m_k$ contain huge numerical constants (see (56)-(57)) that results in large numerical constants in the expression for the number of iterations $N$ and the total number of oracle calls required to guarantee accuracy $\varepsilon$ of the solution. However, for the sake of simplicity of the proofs, we do not try to provide an analysis with optimal or near-optimal dependence on the numerical constants. Moreover, the main goal in this paper is to derive improved high-probability complexity guarantees in terms of $\mathcal{O}(\cdot)$-notation – such guarantees are insensitive to numerical constants by definition.

30

Finally, (56) implies that the batchsize at iteration $k$ is

$$m_k = \Theta\left(\max\left\{1, \frac{N\sigma^2(k+1)^{\frac{4\nu}{1+\nu}}\varepsilon^{\frac{2(1-\nu)}{1+\nu}}\ln\frac{N}{\beta}}{a^2 M_\nu^{\frac{4}{1+\nu}}R_0^2}\right\}\right)$$

meaning that for $k \sim N$ and $a = \mathcal{O}\left(\ln^2\frac{N}{\beta}\right)$ we have that the second term in the maximum is proportional to $N^{\frac{1+5\nu}{1+\nu}}\varepsilon^{\frac{2(1-\nu)}{1+\nu}}$. When $\nu$ is close to $1$ and $\sigma^2 \gg 0$ it implies that $m_k$ is huge for big enough $k$ making the method completely impractical. Fortunately, this issue can be easily solved without sacrificing the oracle complexity of the method: it is sufficient to choose large enough $a$.

**Corollary B.1.** *Let the assumptions of Thm. B.1 hold and*

$$a = \max\left\{16384\ln^2\frac{4N}{\beta}, \frac{5184^{\frac{1+3\nu}{1+\nu}}\cdot 2^{\frac{2(1+5\nu)(1+2\nu)}{(1+\nu)^2}}\sigma^{\frac{2(1+3\nu)}{1+\nu}}C^{\frac{4\nu}{1+\nu}}R_0^{\frac{4\nu}{1+\nu}}\ln^{\frac{1+3\nu}{1+\nu}}\frac{4N}{\beta}}{M_\nu^{\frac{2}{1+\nu}}\varepsilon^{\frac{6\nu}{1+\nu}}}\right\}. \tag{61}$$

*Then for all $k = 0, 1, \ldots, N-1$ we have $m_k = 1$ and to achieve $f(y^N) - f(x^*) \le \varepsilon$ with probability at least $1 - \beta$* clipped-SSTM *requires*

$$\mathcal{O}\left(\max\left\{\frac{M_\nu^{\frac{2}{1+3\nu}}R_0^{\frac{2(1+\nu)}{1+3\nu}}}{\varepsilon^{\frac{2}{1+3\nu}}}\ln^{\frac{2(1+\nu)}{1+3\nu}}\frac{M_\nu^{\frac{2}{1+3\nu}}R_0^{\frac{2(1+\nu)}{1+3\nu}}}{\varepsilon^{\frac{2}{1+3\nu}}\beta}, \frac{\sigma^2 R_0^2}{\varepsilon^2}\ln\frac{\sigma^2 R_0^2}{\varepsilon^2\beta}\right\}\right) \tag{62}$$

*iterations/oracle calls.*

*Proof.* We start with showing that for the new choice of $a$ we have $m_k = 1$ for all $k = 0, 1, \ldots, N-1$. Indeed, using the assumptions on the parameters from Thm. B.1 we derive

$$\begin{aligned}
m_k &= \max\left\{1, \frac{20736N\sigma^2\alpha_{k+1}^2\ln\frac{4N}{\beta}}{C^2 R_0^2}\right\} = \max\left\{1, \frac{5184N\sigma^2(k+1)^{\frac{4\nu}{1+\nu}}\varepsilon^{\frac{2(1-\nu)}{1+\nu}}}{a^2 M_\nu^{\frac{4}{1+\nu}}C^2 R_0^2}\right\} \\
&\overset{k<N}{\le} \max\left\{1, \frac{5184\sigma^2 N^{\frac{1+5\nu}{1+\nu}}\varepsilon^{\frac{2(1-\nu)}{1+\nu}}}{a^2 M_\nu^{\frac{4}{1+\nu}}C^2 R_0^2}\right\} \\
&\overset{(39)}{\le} \max\left\{1, \frac{5184\cdot 2^{\frac{2(1+5\nu)(1+2\nu)}{(1+\nu)(1+3\nu)}}\sigma^2 C^{\frac{4\nu}{1+3\nu}}R_0^{\frac{4\nu}{1+3\nu}}\ln\frac{4N}{\beta}}{a^{\frac{1+\nu}{1+3\nu}}M_\nu^{\frac{2}{1+3\nu}}\varepsilon^{\frac{6\nu}{1+3\nu}}}\right\} \overset{(61)}{\le} 1.
\end{aligned}$$

That is, with the choice of the stepsize parameter $a$ as in (61) the method uses unit batchsizes at each iteration. Therefore, iteration and oracle complexities coincide in this case. Next, we consider two possible situations.

1. If $a = 16384\ln^2\frac{4N}{\beta}$, then

$$\begin{aligned}
N &\overset{(39)}{=} \left\lceil\frac{2^{\frac{1+\nu}{1+3\nu}}a^{\frac{1+\nu}{1+3\nu}}C^{\frac{2(1+\nu)}{1+3\nu}}R_0^{\frac{2(1+\nu)}{1+3\nu}}M_\nu^{\frac{2}{1+3\nu}}}{\varepsilon^{\frac{2}{1+3\nu}}}\right\rceil + 1 = \mathcal{O}\left(\frac{M_\nu^{\frac{2}{1+3\nu}}R_0^{\frac{2(1+\nu)}{1+3\nu}}}{\varepsilon^{\frac{2}{1+3\nu}}}\ln^{\frac{2(1+\nu)}{1+3\nu}}\frac{N}{\beta}\right) \\
&= \mathcal{O}\left(\frac{M_\nu^{\frac{2}{1+3\nu}}R_0^{\frac{2(1+\nu)}{1+3\nu}}}{\varepsilon^{\frac{2}{1+3\nu}}}\ln^{\frac{2(1+\nu)}{1+3\nu}}\frac{M_\nu^{\frac{2}{1+3\nu}}R_0^{\frac{2(1+\nu)}{1+3\nu}}}{\varepsilon^{\frac{2}{1+3\nu}}\beta}\right).
\end{aligned}$$

2. If $a = \frac{5184^{\frac{1+3\nu}{1+\nu}}\cdot 2^{\frac{2(1+5\nu)(1+2\nu)}{(1+\nu)^2}}\sigma^{\frac{2(1+3\nu)}{1+\nu}}C^{\frac{4\nu}{1+\nu}}R_0^{\frac{4\nu}{1+\nu}}\ln^{\frac{1+3\nu}{1+\nu}}\frac{4N}{\beta}}{M_\nu^{\frac{2}{1+\nu}}\varepsilon^{\frac{6\nu}{1+\nu}}}$, then

$$\begin{aligned}
N &\overset{(39)}{=} \left\lceil\frac{2^{\frac{1+\nu}{1+3\nu}}a^{\frac{1+\nu}{1+3\nu}}C^{\frac{2(1+\nu)}{1+3\nu}}R_0^{\frac{2(1+\nu)}{1+3\nu}}M_\nu^{\frac{2}{1+3\nu}}}{\varepsilon^{\frac{2}{1+3\nu}}}\right\rceil + 1 \\
&= \mathcal{O}\left(\frac{M_\nu^{\frac{2}{1+3\nu}}R_0^{\frac{2(1+\nu)}{1+3\nu}}}{\varepsilon^{\frac{2}{1+3\nu}}}\cdot\frac{\sigma^2 R_0^{\frac{4\nu}{1+3\nu}}\ln\frac{4N}{\beta}}{M_\nu^{\frac{2}{1+3\nu}}\varepsilon^{\frac{6\nu}{1+3\nu}}}\right) = \mathcal{O}\left(\frac{\sigma^2 R_0^2}{\varepsilon^2}\ln\frac{\sigma^2 R_0^2}{\varepsilon^2\beta}\right).
\end{aligned}$$

31

696    Putting all together we derive (62).                 □

## B.2    Convergence in the strongly convex case

698 In this section, we provide the full proof of Thm. 2.2 together with complete statement of the result.
699 Note that due to strong convexity the solution $x^*$ is unique.

700 **Theorem B.2.** *Assume that function $f$ is $\mu$-strongly convex and its gradients satisfy (3) with $\nu \in [0, 1]$,*
701 $M_\nu > 0$ *on* $Q = B_{3R_0} = \{x \in \mathbb{R}^n \mid \|x - x^*\|_2 \le 3R_0\}$, *where* $R_0 \ge \|x^0 - x^*\|_2$. *Let* $\varepsilon > 0$,
702 $\beta \in (0, 1)$ *and for* $t = 1, \ldots, \tau$

$$N_t = \left\lceil \frac{2^{\frac{1+\nu}{1+3\nu}} a_t^{\frac{1+\nu}{1+3\nu}} C^{\frac{2(1+\nu)}{1+3\nu}} R_0^{\frac{2(1+\nu)}{1+3\nu}} M_\nu^{\frac{2}{1+3\nu}}}{2^{\frac{(1+\nu)(t-1)}{1+3\nu}} \varepsilon_t^{\frac{2}{1+3\nu}}} \right\rceil + 1, \quad \varepsilon_t = \frac{\mu R_0^2}{2^{t+1}}, \tag{63}$$

703

$$\tau = \left\lceil \log_2 \frac{\mu R_0}{2\varepsilon} \right\rceil, \quad \ln \frac{4N_t \tau}{\beta} \ge 2, \quad C = \sqrt{7}, \tag{64}$$

704

$$\alpha^t = \frac{(\varepsilon_t/2)^{\frac{1-\nu}{1+\nu}}}{2^{\frac{2\nu}{1+\nu}} a_t M_\nu^{\frac{2}{1+\nu}}}, \quad m_k^t = \max \left\{ 1, \frac{20736 \cdot 2^{t-1} N_t \sigma^2 (\alpha_{k+1}^t)^2 \ln \frac{4N_t \tau}{\beta}}{C^2 R_0^2} \right\}, \tag{65}$$

705

$$\alpha_{k+1}^t = \alpha^t (k+1)^{\frac{2\nu}{1+\nu}}, \quad B = \frac{CR_0}{16 \ln \frac{4N_t \tau}{\beta}}, \quad a_t = 16384 \ln^2 \frac{4N_t \tau}{\beta}, \tag{66}$$

706

$$\varepsilon_t^{\frac{1-\nu}{1+\nu}} \le \frac{a_t C M_\nu^{\frac{1-\nu}{1+\nu}} R_0^{1-\nu}}{16 \cdot 2^{\frac{(1-\nu)(t-1)}{2}} \ln \frac{4N_t \tau}{\beta}}, \quad \varepsilon_t \le \frac{2^{\frac{1+\nu}{2}} a_t^{\frac{1+\nu}{2}} C^{1+\nu} R_0^{1+\nu} M_\nu}{100^{\frac{1+3\nu}{2}} \cdot 2^{\frac{(1+\nu)(t-1)}{2}}}, \tag{67}$$

707

$$\varepsilon_t^{\frac{1-\nu}{1+3\nu}} \le \min \left\{ \frac{a_t^{\frac{2+3\nu-\nu^2}{2(1+3\nu)}}}{2^{2+4\nu+\frac{3+8\nu-5\nu^2-6\nu^3}{(1+\nu)(1+3\nu)}} \ln \frac{4N_t \tau}{\beta}}, \frac{a_t^{\frac{(1+\nu)^2}{1+3\nu}}}{2^{4+7\nu+\frac{2+7\nu+2\nu^2-3\nu^3}{(1+\nu)(1+3\nu)}} \ln^{1+\nu} \frac{4N_t \tau}{\beta}} \right\} \frac{C^{\frac{1-\nu^2}{1+3\nu}} R_0^{\frac{1-\nu^2}{1+3\nu}} M_\nu^{\frac{1-\nu}{1+3\nu}}}{2^{\frac{(1-\nu^2)(t-1)}{2(1+3\nu)}}}. \tag{68}$$

708 *Then, after $\tau$ restarts* R-clipped-SSTM *produces $\hat{x}^\tau$ such that with probability at least $1 - \beta$*

$$f(\hat{x}^\tau) - f(x^*) \le \varepsilon. \tag{69}$$

709 *That is, to achieve (69) with probability at least $1 - \beta$ the method requires*

$$\hat{N} = \mathcal{O}\left( \max \left\{ \left( \frac{M_\nu}{\mu R_0^{1-\nu}} \right)^{\frac{2}{1+3\nu}} \ln \frac{\mu R_0^2}{\varepsilon}, \left( \frac{M_\nu^2}{\mu^{1+\nu} \varepsilon^{1-\nu}} \right)^{\frac{1}{1+3\nu}} \right\} \ln^{\frac{2(1+\nu)}{1+3\nu}} \frac{M_\nu^{\frac{2}{1+3\nu}} \ln \frac{\mu R_0^2}{\varepsilon}}{\mu^{\frac{1+\nu}{1+3\nu}} \varepsilon^{\frac{1-\nu}{1+3\nu}} \beta} \right) \tag{70}$$

710 *iterations of Alg. 1 and*

$$\mathcal{O}\left( \max \left\{ \hat{N}, \frac{\sigma^2}{\mu \varepsilon} \ln \frac{M_\nu^{\frac{2}{1+3\nu}} \ln \frac{\mu R_0^2}{\varepsilon}}{\mu^{\frac{1+\nu}{1+3\nu}} \varepsilon^{\frac{1-\nu}{1+3\nu}} \beta} \right\} \right) \quad \text{oracle calls.} \tag{71}$$

711 *Proof.* Applying Thm. B.1, we obtain that with probability at least $1 - \frac{\beta}{\tau}$

$$f(\hat{x}^1) - f(x^*) \le \frac{\mu R_0^2}{4}.$$

712 Since $f$ is $\mu$-strongly convex we have

$$\frac{\mu \|\hat{x}^1 - x^*\|_2^2}{2} \le f(\hat{x}^1) - f(x^*).$$

713 Therefore, with probability at least $1 - \frac{\beta}{\tau}$

$$f(\hat{x}^1) - f(x^*) \le \frac{\mu R_0^2}{4}, \quad \|\hat{x}^1 - x^*\|_2^2 \le \frac{R_0^2}{2}.$$

32

From mathematical induction and the union bound for probability events it follows that inequalities

$$f(\hat{x}^t) - f(x^*) \leq \frac{\mu R_0^2}{2^{t+1}}, \quad \|\hat{x}^t - x^*\|_2^2 \leq \frac{R_0^2}{2^t}$$

hold simultaneously for $t = 1, \ldots, \tau$ with probability at least $1 - \beta$. In particular, it means that after $\tau = \left\lceil \log_2 \frac{\mu R_0^2}{\varepsilon} \right\rceil - 1$ restarts R-clipped-SSTM finds an $\varepsilon$-solution with probability at least $1 - \beta$. The total number of iterations $\hat{N}$ is

$$
\begin{aligned}
\sum_{t=1}^{\tau} N_t &= \mathcal{O}\left( \sum_{t=1}^{\tau} \frac{M_\nu^{\frac{2}{1+3\nu}} R_0^{\frac{2(1+\nu)}{1+3\nu}}}{2^{\frac{(1+\nu)t}{1+3\nu}} \varepsilon_t^{\frac{2}{1+3\nu}}} \ln^{\frac{2(1+\nu)}{1+3\nu}} \frac{M_\nu^{\frac{2}{1+3\nu}} R_0^{\frac{2(1+\nu)}{1+3\nu}} \tau}{2^{\frac{(1+\nu)t}{1+3\nu}} \varepsilon_t^{\frac{2}{1+3\nu}} \beta} \right) \\
&= \mathcal{O}\left( \sum_{t=1}^{\tau} \frac{M_\nu^{\frac{2}{1+3\nu}} R_0^{\frac{2(1+\nu)}{1+3\nu}} 2^{\frac{2t}{1+3\nu}}}{2^{\frac{(1+\nu)t}{1+3\nu}} \mu^{\frac{2}{1+3\nu}} R_0^{\frac{4}{1+3\nu}}} \ln^{\frac{2(1+\nu)}{1+3\nu}} \frac{M_\nu^{\frac{2}{1+3\nu}} R_0^{\frac{2(1+\nu)}{1+3\nu}} 2^{\frac{2t}{1+3\nu}} \tau}{2^{\frac{(1+\nu)t}{1+3\nu}} \mu^{\frac{2}{1+3\nu}} R_0^{\frac{4}{1+3\nu}} \beta} \right) \\
&= \mathcal{O}\left( \sum_{t=1}^{\tau} \frac{M_\nu^{\frac{2}{1+3\nu}} 2^{\frac{(1-\nu)t}{1+3\nu}}}{\mu^{\frac{2}{1+3\nu}} R_0^{\frac{2(1-\nu)}{1+3\nu}}} \ln^{\frac{2(1+\nu)}{1+3\nu}} \frac{M_\nu^{\frac{2}{1+3\nu}} 2^{\frac{(1-\nu)t}{1+3\nu}} \tau}{\mu^{\frac{2}{1+3\nu}} R_0^{\frac{2(1-\nu)}{1+3\nu}} \beta} \right) \\
&= \mathcal{O}\left( \frac{M_\nu^{\frac{2}{1+3\nu}} \max\left\{ \tau, 2^{\frac{(1-\nu)\tau}{1+3\nu}} \right\}}{\mu^{\frac{2}{1+3\nu}} R_0^{\frac{2(1-\nu)}{1+3\nu}}} \ln^{\frac{2(1+\nu)}{1+3\nu}} \frac{M_\nu^{\frac{2}{1+3\nu}} 2^{\frac{(1-\nu)\tau}{1+3\nu}} \tau}{\mu^{\frac{2}{1+3\nu}} R_0^{\frac{2(1-\nu)}{1+3\nu}} \beta} \right) \\
&= \mathcal{O}\left( \max\left\{ \left( \frac{M_\nu}{\mu R_0^{1-\nu}} \right)^{\frac{2}{1+3\nu}} \ln \frac{\mu R_0^2}{\varepsilon}, \left( \frac{M_\nu^2}{\mu^{1+\nu}\varepsilon^{1-\nu}} \right)^{\frac{1}{1+3\nu}} \right\} \ln^{\frac{2(1+\nu)}{1+3\nu}} \frac{M_\nu^{\frac{2}{1+3\nu}} \ln \frac{\mu R_0^2}{\varepsilon}}{\mu^{\frac{1+\nu}{1+3\nu}} \varepsilon^{\frac{1-\nu}{1+3\nu}} \beta} \right),
\end{aligned}
$$

and the total number of oracle calls equals

$$
\begin{aligned}
\sum_{t=1}^{\tau} \sum_{k=0}^{N_t-1} m_k^t &= \mathcal{O}\left( \max\left\{ \sum_{t=1}^{\tau} N_t, \sum_{t=1}^{\tau} \frac{\sigma^2 R_0^2}{2^t \varepsilon_t^2} \ln \frac{M_\nu^{\frac{2}{1+3\nu}} 2^{\frac{(1-\nu)t}{1+3\nu}} \tau}{\mu^{\frac{2}{1+3\nu}} R_0^{\frac{2(1-\nu)}{1+3\nu}} \beta} \right\} \right) \\
&= \mathcal{O}\left( \max\left\{ \hat{N}, \sum_{t=1}^{\tau} \frac{\sigma^2 \cdot 2^t}{\mu^2 R_0^2} \ln \frac{M_\nu^{\frac{2}{1+3\nu}} 2^{\frac{(1-\nu)\tau}{1+3\nu}} \tau}{\mu^{\frac{2}{1+3\nu}} R_0^{\frac{2(1-\nu)}{1+3\nu}} \beta} \right\} \right) \\
&= \mathcal{O}\left( \max\left\{ \hat{N}, \frac{\sigma^2}{\mu\varepsilon} \ln \frac{M_\nu^{\frac{2}{1+3\nu}} \ln \frac{\mu R_0^2}{\varepsilon}}{\mu^{\frac{1+\nu}{1+3\nu}} \varepsilon^{\frac{1-\nu}{1+3\nu}} \beta} \right\} \right).
\end{aligned}
$$

$\square$

One can also derive a similar result for R-clipped-SSTM when stepsize parameter $a$ is chosen as in Cor. B.1 for all restarts.

# C  SGD with clipping: missing details and proofs

## C.1  Convex case

In this section, we provide a full statement of Thm. 3.1 together with its proof. The proof is based on a similar idea as the proof of the complexity bounds for clipped-SSTM.

**Theorem C.1.** *Assume that function $f$ is convex, achieves its minimum at a point $x^*$, and its gradients satisfy (3) with $\nu \in [0, 1]$, $M_\nu$ on $Q = B_{7R_0} = \{x \in \mathbb{R}^n \mid \|x - x^*\|_2 \leq 7R_0\}$, where $R_0 \geq \|x^0 - x^*\|_2$. Then, for all $\beta \in (0, 1)$ and $N$ such that*

$$\ln \frac{4N}{\beta} \geq 2, \tag{72}$$

*we have that after $N$ iterations of* clipped-SGD *with*

$$\lambda = \frac{R_0}{\gamma \ln \frac{4N}{\beta}}, \quad m \geq \max\left\{1, \frac{81N\sigma^2}{\lambda^2 \ln \frac{4N}{\beta}}\right\} \tag{73}$$

*and stepsize*

$$\gamma \leq \min\left\{\frac{\varepsilon^{\frac{1-\nu}{1+\nu}}}{8M_\nu^{\frac{2}{1+\nu}}}, \frac{R_0}{\sqrt{2N}\varepsilon^{\frac{\nu}{1+\nu}}M_\nu^{\frac{1}{1+\nu}}}, \frac{R_0^{1-\nu}}{2C^\nu M_\nu \ln \frac{4N}{\beta}}\right\}, \tag{74}$$

*with probability at least $1 - \beta$ it holds that*

$$f(\bar{x}^N) - f(x^*) \leq \frac{C^2 R_0^2}{\gamma N}, \tag{75}$$

*where $\bar{x}^N = \frac{1}{N}\sum_{k=0}^{N-1} x^k$ and*

$$C = 7. \tag{76}$$

*In other words,* clipped-SGD *with $\gamma = \min\left\{\frac{\varepsilon^{\frac{1-\nu}{1+\nu}}}{8M_\nu^{\frac{2}{1+\nu}}}, \frac{R_0}{\sqrt{2N}\varepsilon^{\frac{\nu}{1+\nu}}M_\nu^{\frac{1}{1+\nu}}}, \frac{R_0^{1-\nu}}{2C^\nu M_\nu \ln \frac{4N}{\beta}}\right\}$ achieves $f(\bar{x}^N) - f(x^*) \leq \varepsilon$ with probability at least $1-\beta$ after $\mathcal{O}\left(\max\left\{\frac{M_\nu^{\frac{2}{1+\nu}} R_0^2}{\varepsilon^{\frac{2}{1+\nu}}}, \frac{M_\nu R_0^{1+\nu}}{\varepsilon} \ln \frac{M_\nu R_0^{1+\nu}}{\varepsilon\beta}\right\}\right)$ iterations and requires*

$$\mathcal{O}\left(\max\left\{\frac{M_\nu^{\frac{2}{1+\nu}} R_0^2}{\varepsilon^{\frac{2}{1+\nu}}}, \max\left\{\frac{M_\nu R_0^{1+\nu}}{\varepsilon}, \frac{\sigma^2 R_0^2}{\varepsilon^2}\right\} \ln \frac{M_\nu R_0^{1+\nu}}{\varepsilon\beta}\right\}\right) \tag{77}$$

*oracle calls.*

*Proof.* Since $f(x)$ is convex and its gradients satisfy (3), we get the following inequality under assumption that $x^k \in B_{7R_0}(x^*)$:

$$
\begin{aligned}
\|x^{k+1} - x^*\|_2^2 &= \|x^k - \gamma\widetilde{\nabla}f(x^k, \boldsymbol{\xi}^k) - x^*\|_2^2 \\
&= \|x^k - x^*\|_2^2 + \gamma^2\|\widetilde{\nabla}f(x^k, \boldsymbol{\xi}^k)\|_2^2 - 2\gamma\left\langle x^k - x^*, \widetilde{\nabla}f(x^k, \boldsymbol{\xi}^k)\right\rangle \\
&= \|x^k - x^*\|_2^2 + \gamma^2\|\nabla f(x^k) + \theta_k\|_2^2 - 2\gamma\left\langle x^k - x^*, \nabla f(x^k) + \theta_k\right\rangle \\
&\overset{(11)}{\leq} \|x^k - x^*\|_2^2 + 2\gamma^2\|\nabla f(x^k)\|_2^2 + 2\gamma^2\|\theta_k\|_2^2 - 2\gamma\left\langle x^k - x^*, \nabla f(x^k) + \theta_k\right\rangle \\
&\overset{(21)}{\leq} \|x^k - x^*\|_2^2 - 2\gamma\left(1 - 2\gamma\left(\frac{1}{\varepsilon}\right)^{\frac{1-\nu}{1+\nu}} M_\nu^{\frac{2}{1+\nu}}\right)\left(f(x^k) - f(x^*)\right) + 2\gamma^2\|\theta_k\|_2^2 \\
&\quad -2\gamma\left\langle x^k - x^*, \theta_k\right\rangle + 2\gamma^2\varepsilon^{\frac{2\nu}{1+\nu}}M_\nu^{\frac{2}{1+\nu}},
\end{aligned}
$$

where $\theta_k = \widetilde{\nabla}f(x^k, \boldsymbol{\xi}^k) - \nabla f(x^k)$ and the last inequality follows from the convexity of $f$. Using notation $R_k \overset{\text{def}}{=} \|x^k - x^*\|_2$, $k > 0$ we derive that for all $k \geq 0$

$$R_{k+1}^2 \leq R_k^2 - 2\gamma\left(1 - 2\gamma\left(\frac{1}{\varepsilon}\right)^{\frac{1-\nu}{1+\nu}} M_\nu^{\frac{2}{1+\nu}}\right)\left(f(x^k) - f(x^*)\right) + 2\gamma^2\|\theta_k\|_2^2 - 2\gamma\left\langle x^k - x^*, \theta_k\right\rangle + 2\gamma^2\varepsilon^{\frac{2\nu}{1+\nu}} M_\nu^{\frac{2}{1+\nu}}$$

under assumption that $x^k \in B_{7R_0}(x^*)$. Let us define $A = 2\gamma \left( 1 - 2\gamma \left( \frac{1}{\varepsilon} \right)^{\frac{1-\nu}{1+\nu}} M_\nu^{\frac{2}{1+\nu}} \right) \overset{(74)}{\geq} \gamma > 0$, then

$$A \left( f(x^k) - f(x^*) \right) \leq R_k^2 - R_{k+1}^2 + 2\gamma^2 \|\theta_k\|_2^2 - 2\gamma \left\langle x^k - x^*, \theta_k \right\rangle + 2\gamma^2 \varepsilon^{\frac{2\nu}{1+\nu}} M_\nu^{\frac{2}{1+\nu}}$$

under assumption that $x^k \in B_{7R_0}(x^*)$. Summing up these inequalities for $k = 0, \dots, N-1$, we obtain

$$
\begin{aligned}
\frac{A}{N} \sum_{k=0}^{N-1} \left[ f(x^k) - f(x^*) \right] &\leq \frac{1}{N} \sum_{k=0}^{N-1} \left( R_k^2 - R_{k+1}^2 \right) + 2\gamma^2 \varepsilon^{\frac{2\nu}{1+\nu}} M_\nu^{\frac{2}{1+\nu}} + \frac{2\gamma^2}{N} \sum_{k=0}^{N-1} \|\theta_k\|_2^2 \\
&\quad - \frac{2\gamma}{N} \sum_{k=0}^{N-1} \left\langle x^k - x^*, \theta_k \right\rangle \\
&= \frac{1}{N} \left( R_0^2 - R_N^2 \right) + 2\gamma^2 \varepsilon^{\frac{2\nu}{1+\nu}} M_\nu^{\frac{2}{1+\nu}} + \frac{2\gamma^2}{N} \sum_{k=0}^{N-1} \|\theta_k\|_2^2 \\
&\quad - \frac{2\gamma}{N} \sum_{k=0}^{N-1} \left\langle x^k - x^*, \theta_k \right\rangle
\end{aligned}
$$

under assumption that $x^k \in B_{7R_0}(x^*)$. Noticing that for $\bar{x}^N = \frac{1}{N} \sum_{k=0}^{N-1} x^k$ Jensen's inequality gives

$f(\bar{x}^N) = f \left( \frac{1}{N} \sum_{k=0}^{N-1} x^k \right) \leq \frac{1}{N} \sum_{k=0}^{N-1} f(x^k)$, we have

$$AN \left( f(\bar{x}^N) - f(x^*) \right) \leq R_0^2 - R_N^2 + 2\gamma^2 N \varepsilon^{\frac{2\nu}{1+\nu}} M_\nu^{\frac{2}{1+\nu}} + 2\gamma^2 \sum_{k=0}^{N-1} \|\theta_k\|_2^2 - 2\gamma \sum_{k=0}^{N-1} \left\langle x^k - x^*, \theta_k \right\rangle \tag{78}$$

under assumption that $x^k \in B_{7R_0}(x^*)$ for $k = 0, 1, \dots, N-1$. Taking into account that $f(\bar{x}^N) - f(x^*) \geq 0$ and changing the indices we get that for all $k = 0, 1, \dots, N$

$$R_k^2 \leq R_0^2 + 2\gamma^2 k \varepsilon^{\frac{2\nu}{1+\nu}} M_\nu^{\frac{2}{1+\nu}} + 2\gamma^2 \sum_{l=0}^{k-1} \|\theta_l\|_2^2 - 2\gamma \sum_{l=0}^{k-1} \left\langle x^l - x^*, \theta_k \right\rangle. \tag{79}$$

under assumption that $x^l \in B_{7R_0}(x^*)$ for $l = 0, 1, \dots, k-1$. The remaining part of the proof is based on the analysis of inequality (79). In particular, via induction we prove that for all $k = 0, 1, \dots, N$ with probability at least $1 - \frac{k\beta}{N}$ the following statement holds: inequalities

$$R_t^2 \overset{(79)}{\leq} R_0^2 + 2\gamma^2 t \varepsilon^{\frac{2\nu}{1+\nu}} M_\nu^{\frac{2}{1+\nu}} + 2\gamma^2 \sum_{l=0}^{t-1} \|\theta_k\|_2^2 - 2\gamma \sum_{l=0}^{t-1} \left\langle x^k - x^*, \theta_k \right\rangle \leq C^2 R_0^2 \tag{80}$$

hold for $t = 0, 1, \dots, k$ simultaneously where $C$ is defined in (76). Let us define the probability event when this statement holds as $E_k$. Then, our goal is to show that $\mathbb{P}\{E_k\} \geq 1 - \frac{k\beta}{N}$ for all $k = 0, 1, \dots, N$. For $t = 0$ inequality (80) holds with probability 1 since $C \geq 1$. Next, assume that for some $k = T - 1 \leq N - 1$ we have $\mathbb{P}\{E_k\} = \mathbb{P}\{E_{T-1}\} \geq 1 - \frac{(T-1)\beta}{N}$. Let us prove that $\mathbb{P}\{E_T\} \geq 1 - \frac{T\beta}{N}$. First of all, probability event $E_{T-1}$ implies that $x^t \in B_{7R_0}(x^*)$ for $t = 0, 1, \dots, T-1$, and, as a consequence, (79) holds for $k = T$. Since $\nabla f(x)$ is $(\nu, M_\nu)$-Hölder continuous on $B_{7R_0}(x^*)$, we have that probability event $E_{T-1}$ implies

$$\left\| \nabla f(x^t) \right\|_2 \overset{(3)}{\leq} M_\nu \|x^t - x^0\|^\nu \leq M_\nu C^\nu R_0^\nu \overset{(74)}{\leq} \frac{\lambda}{2} \tag{81}$$

for $t = 0, \dots, T - 1$. Next, we introduce new random variables:

$$\eta_l = \begin{cases} x^* - x^l, & \text{if } \|x^* - x^l\|_2 \leq C R_0, \\ 0, & \text{otherwise,} \end{cases} \tag{82}$$

35

for $l = 0, 1, \ldots T - 1$. Note that these random variables are bounded with probability 1, i.e. with probability 1 we have

$$\|\eta_l\|_2 \le C R_0. \tag{83}$$

Using the introduced notation, we obtain that $E_{T-1}$ implies

$$R_T^2 \overset{(73),(74),(79),(80),(82)}{\le} 2R_0^2 + 2\gamma \sum_{l=0}^{T-1} \langle \theta_l, \eta_l \rangle + 2\gamma^2 \sum_{l=0}^{T-1} \|\theta_l\|_2^2.$$

Finally, we do some preliminaries in order to apply Bernstein's inequality (see Lemma A.2) and obtain that $E_{T-1}$ implies

$$R_T^2 \overset{(11)}{\le} 2R_0^2 + \underbrace{2\gamma \sum_{l=0}^{T-1} \langle \theta_l^u, \eta_l \rangle}_{①} + \underbrace{2\gamma \sum_{l=0}^{T-1} \langle \theta_l^b, \eta_l \rangle}_{②} + \underbrace{4\gamma^2 \sum_{l=0}^{T-1} \left( \|\theta_l^u\|_2^2 - \mathbb{E}_{\boldsymbol{\xi}^l} \left[ \|\theta_l^u\|_2^2 \right] \right)}_{③}$$

$$+ \underbrace{4\gamma^2 \sum_{l=0}^{T-1} \mathbb{E}_{\boldsymbol{\xi}^l} \left[ \|\theta_l^u\|_2^2 \right]}_{④} + \underbrace{4\gamma^2 \sum_{l=0}^{T-1} \|\theta_l^b\|_2^2}_{⑤}, \tag{84}$$

where we introduce new notations:

$$\theta_l^u \overset{\text{def}}{=} \widetilde{\nabla} f(x^l, \boldsymbol{\xi}^l) - \mathbb{E}_{\boldsymbol{\xi}^l} \left[ \widetilde{\nabla} f(x^l, \boldsymbol{\xi}^l) \right], \quad \theta_l^b \overset{\text{def}}{=} \mathbb{E}_{\boldsymbol{\xi}^l} \left[ \widetilde{\nabla} f(x^l, \boldsymbol{\xi}^l) \right] - \nabla f(x^l), \tag{85}$$

$$\theta_l = \theta_l^u + \theta_l^b.$$

It remains to provide tight upper bounds for ①, ②, ③, ④ and ⑤, i.e. in the remaining part of the proof we show that $① + ② + ③ + ④ + ⑤ \le \delta C^2 R_0^2$ for some $\delta < 1$.

**Upper bound for ①.** First of all, since $\mathbb{E}_{\boldsymbol{\xi}^l}[\theta_l^u] = 0$ summands in ① are conditionally unbiased:

$$\mathbb{E}_{\boldsymbol{\xi}^l} \left[ 2\gamma \langle \theta_l^u, \eta_l \rangle \right] = 0.$$

Secondly, these summands are bounded with probability 1:

$$|2\gamma \langle \theta_l^u, \eta_l \rangle| \quad \le \quad 2\gamma \|\theta_l^u\|_2 \|\eta_l\|_2 \overset{(29),(83)}{\le} 4\gamma \lambda C R_0.$$

Finally, one can bound conditional variances $\sigma_l^2 \overset{\text{def}}{=} \mathbb{E}_{\boldsymbol{\xi}^l} \left[ 4\gamma^2 \langle \theta_l^u, \eta_l \rangle^2 \right]$ in the following way:

$$\sigma_l^2 \quad \le \quad \mathbb{E}_{\boldsymbol{\xi}^l} \left[ 4\gamma^2 \|\theta_l^u\|_2^2 \|\eta_l\|_2^2 \right] \overset{(83)}{\le} 4\gamma^2 (C R_0)^2 \mathbb{E}_{\boldsymbol{\xi}^l} \left[ \|\theta_l^u\|_2^2 \right].$$

In other words, sequence $\{2\gamma \langle \theta_l^u, \eta_l \rangle\}_{l \ge 0}$ is a bounded martingale difference sequence with bounded conditional variances $\{\sigma_l^2\}_{l \ge 0}$. Therefore, we can apply Bernstein's inequality, i.e., we apply Lemma A.2 with $X_l = 2\gamma \langle \theta_l^u, \eta_l \rangle$, $c = 4\gamma \lambda C R_0$ and $F = \frac{c^2 \ln \frac{4N}{\beta}}{6}$ and get that for all $b > 0$

$$\mathbb{P} \left\{ \left| \sum_{l=0}^{T-1} X_l \right| > b \text{ and } \sum_{l=0}^{T-1} \sigma_l^2 \le F \right\} \le 2 \exp \left( -\frac{b^2}{2F + 2cb/3} \right)$$

or, equivalently, with probability at least $1 - 2 \exp \left( -\frac{b^2}{2F + 2cb/3} \right)$

$$\text{either } \sum_{l=0}^{T-1} \sigma_l^2 > F \quad \text{or} \quad \underbrace{\left| \sum_{l=0}^{T-1} X_l \right|}_{|①|} \le b.$$

The choice of $F$ will be clarified further, let us now choose $b$ in such a way that $2 \exp \left( -\frac{b^2}{2F + 2cb/3} \right) = \frac{\beta}{2N}$. This implies that $b$ is the positive root of the quadratic equation

$$b^2 - \frac{2c \ln \frac{4N}{\beta}}{3} b - 2F \ln \frac{4N}{\beta} = 0,$$

36

hence

$$
\begin{aligned}
b &= \frac{c\ln\frac{4N}{\beta}}{3} + \sqrt{\frac{c^2\ln^2\frac{4N}{\beta}}{9} + 2F\ln\frac{4N}{\beta}} = \frac{c\ln\frac{4N}{\beta}}{3} + \sqrt{\frac{4c^2\ln^2\frac{4N}{\beta}}{9}} \\
&= c\ln\frac{4N}{\beta} = 4\gamma\lambda CR_0\ln\frac{4N}{\beta}.
\end{aligned}
$$

That is, with probability at least $1 - \frac{\beta}{2N}$

$$
\underbrace{\text{either } \sum_{l=0}^{T-1}\sigma_l^2 > F \quad\text{or}\quad |①| \leq 4\gamma\lambda CR_0\ln\frac{4N}{\beta}}_{\text{probability event } E_①}.
$$

Next, we notice that probability event $E_{T-1}$ implies that

$$
\begin{aligned}
\sum_{l=0}^{T-1}\sigma_l^2 &\leq 4\gamma^2(CR_0)^2\sum_{l=0}^{T-1}\mathbb{E}_{\boldsymbol{\xi}^l}\left[\|\theta_l^u\|_2^2\right] \overset{(32)}{\leq} 72\gamma^2(CR_0)^2\sigma^2\frac{T}{m} \\
&\overset{T\leq N}{\leq} 72\gamma^2(CR_0)^2\sigma^2\frac{N}{m} \leq \frac{c^2\ln\frac{4N}{\beta}}{6} = F,
\end{aligned}
$$

where the last inequality follows from $c = 4\gamma\lambda CR_0$ and simple arithmetic.

**Upper bound for ②.** First of all, we notice that probability event $E_{T-1}$ implies

$$
2\gamma\left\langle\theta_l^b,\eta_l\right\rangle \leq 2\gamma\left\|\theta_l^b\right\|_2\|\eta_l\|_2 \overset{(30),(83)}{\leq} 2\gamma\frac{4\sigma^2}{m\lambda}CR_0 = \frac{8\gamma\sigma^2CR_0}{m\lambda}.
$$

This implies that

$$
② = 2\gamma\sum_{l=0}^{T-1}\left\langle\theta_l^b,\eta_l\right\rangle \overset{T\leq N}{\leq} \frac{8\gamma\sigma^2CR_0N}{m\lambda} \overset{(73)}{\leq} \frac{8}{81}\lambda\gamma CR_0\ln\frac{4N}{\beta}.
$$

**Upper bound for ③.** We derive the upper bound for ③ using the same technique as for ①. First of all, we notice that the summands in ③ are conditionally unbiased:

$$
\mathbb{E}_{\boldsymbol{\xi}^l}\left[4\gamma^2\left(\|\theta_l^u\|_2^2 - \mathbb{E}_{\boldsymbol{\xi}^l}\left[\|\theta_l^u\|_2^2\right]\right)\right] = 0.
$$

Secondly, the summands are bounded with probability 1:

$$
\begin{aligned}
\left|4\gamma^2\left(\|\theta_l^u\|_2^2 - \mathbb{E}_{\boldsymbol{\xi}^l}\left[\|\theta_l^u\|_2^2\right]\right)\right| &\leq 4\gamma^2\left(\|\theta_l^u\|_2^2 + \mathbb{E}_{\boldsymbol{\xi}^l}\left[\|\theta_l^u\|_2^2\right]\right) \overset{(29)}{\leq} 4\gamma^2\left(4\lambda^2 + 4\lambda^2\right) \\
&= 32\gamma^2\lambda^2 \overset{\text{def}}{=} c_1. \tag{86}
\end{aligned}
$$

Finally, one can bound conditional variances $\hat{\sigma}_l^2 \overset{\text{def}}{=} \mathbb{E}_{\boldsymbol{\xi}^l}\left[\left|4\gamma^2\left(\|\theta_l^u\|_2^2 - \mathbb{E}_{\boldsymbol{\xi}^l}\left[\|\theta_l^u\|_2^2\right]\right)\right|^2\right]$ in the following way:

$$
\begin{aligned}
\hat{\sigma}_l^2 &\overset{(86)}{\leq} c_1\mathbb{E}_{\boldsymbol{\xi}^l}\left[\left|4\gamma^2\left(\|\theta_l^u\|_2^2 - \mathbb{E}_{\boldsymbol{\xi}^l}\left[\|\theta_l^u\|_2^2\right]\right)\right|\right] \\
&\leq 4\gamma^2c_1\mathbb{E}_{\boldsymbol{\xi}^l}\left[\|\theta_l^u\|_2^2 + \mathbb{E}_{\boldsymbol{\xi}^l}\left[\|\theta_l^u\|_2^2\right]\right] = 8\gamma^2c_1\mathbb{E}_{\boldsymbol{\xi}^l}\left[\|\theta_l^u\|_2^2\right]. \tag{87}
\end{aligned}
$$

In other words, sequence $\left\{4\gamma^2\left(\|\theta_l^u\|_2^2 - \mathbb{E}_{\boldsymbol{\xi}^l}\left[\|\theta_l^u\|_2^2\right]\right)\right\}_{l\geq 0}$ is a bounded martingale difference sequence with bounded conditional variances $\{\hat{\sigma}_l^2\}_{l\geq 0}$. Therefore, we can apply Bernstein's inequality, i.e. we apply Lemma A.2 with $X_l = \hat{X}_l = 4\gamma^2\left(\|\theta_l^u\|_2^2 - \mathbb{E}_{\boldsymbol{\xi}^l}\left[\|\theta_l^u\|_2^2\right]\right)$, $c = c_1 = 32\gamma^2\lambda^2$ and $F = F_1 = \frac{c_1^2\ln\frac{4N}{\beta}}{18}$ and get that for all $b > 0$

$$
\mathbb{P}\left\{\left|\sum_{l=0}^{T-1}\hat{X}_l\right| > b \text{ and } \sum_{l=0}^{T-1}\hat{\sigma}_l^2 \leq F_1\right\} \leq 2\exp\left(-\frac{b^2}{2F_1 + {2c_1b}/{3}}\right)
$$

37

or, equivalently, with probability at least $1 - 2\exp\left(-\frac{b^2}{2F_1 + 2c_1 b/3}\right)$

$$\text{either } \sum_{l=0}^{T-1} \hat{\sigma}_l^2 > F_1 \quad \text{or} \quad \underbrace{\left|\sum_{l=0}^{T-1} \hat{X}_l\right|}_{|③|} \leq b.$$

As in our derivations of the upper bound for ① we choose such $b$ that $2\exp\left(-\frac{b^2}{2F_1 + 2c_1 b/3}\right) = \frac{\beta}{2N}$, i.e.,

$$b = \frac{c_1 \ln \frac{4N}{\beta}}{3} + \sqrt{\frac{c_1^2 \ln^2 \frac{4N}{\beta}}{9} + 2F_1 \ln \frac{4N}{\beta}} \leq c_1 \ln \frac{4N}{\beta} = 32\gamma^2 \lambda^2 \ln \frac{4N}{\beta}.$$

That is, with probability at least $1 - \frac{\beta}{2N}$

$$\underbrace{\text{either } \sum_{l=0}^{T-1} \hat{\sigma}_l^2 > F_1 \quad \text{or} \quad |③| \leq 32\gamma^2 \lambda^2 \ln \frac{4N}{\beta}}_{\text{probability event } E_③}.$$

Next, we notice that probability event $E_{T-1}$ implies that

$$\sum_{l=0}^{T-1} \hat{\sigma}_l^2 \overset{(87)}{\leq} 8\gamma^2 c_1 \sum_{l=0}^{T-1} \mathbb{E}_{\boldsymbol{\xi}^l}\left[\|\theta_l^u\|_2^2\right] \overset{(32)}{\leq} 144\gamma^2 c_1 \sigma^2 \frac{T}{m}$$

$$\overset{T \leq N}{\leq} 144\gamma^2 c_1 \sigma^2 \frac{N}{m} = \frac{c_1^2 \ln \frac{4N}{\beta}}{18} \leq F_1.$$

**Upper bound for ④.** The probability event $E_{T-1}$ implies

$$④ = 4\gamma^2 \sum_{l=0}^{T-1} \mathbb{E}_{\boldsymbol{\xi}^l}\left[\|\theta_l^u\|_2^2\right] \overset{(32)}{\leq} 72\gamma^2 \sigma^2 \sum_{l=0}^{T-1} \frac{1}{m} \overset{T \leq N}{\leq} \frac{72\gamma^2 N \sigma^2}{m} \overset{(73)}{\leq} \frac{8}{9}\lambda^2 \gamma^2 \ln \frac{4N}{\beta}.$$

**Upper bound for ⑤.** Again, we use corollaries of probability event $E_{T-1}$:

$$⑤ = 4\gamma^2 \sum_{l=0}^{T-1} \|\theta_l^b\|_2^2 \overset{(30)}{\leq} 64\gamma^2 \sigma^4 \frac{T}{m^2 \lambda^2} \overset{T \leq N}{\leq} 64\gamma^2 \sigma^4 \frac{N}{m^2 \lambda^2} \overset{(73)}{\leq} \frac{64}{6561} \frac{\lambda^2 \gamma^2 \ln^2 \frac{4N}{\beta}}{N}.$$

Now we summarize all bound that we have: probability event $E_{T-1}$ implies

$$R_T^2 \overset{(79)}{\leq} 2R_0^2 + 2\gamma^2 \sum_{l=0}^{T-1} \|\theta_l\|_2^2 - 2\gamma \sum_{l=0}^{T-1} \langle x^l - x^*, \theta_l \rangle$$

$$\overset{(84)}{\leq} 2R_0^2 + ① + ② + ③ + ④ + ⑤,$$

$$② \leq \frac{8}{81}\lambda\gamma C R_0 \ln \frac{4N}{\beta}, \quad ④ \leq \frac{8}{9}\lambda^2 \gamma^2 \ln \frac{4N}{\beta}, \quad ⑤ \leq \frac{64}{6561} \frac{\lambda^2 \gamma^2 \ln^2 \frac{4N}{\beta}}{N},$$

$$\sum_{l=0}^{T-1} \sigma_l^2 \leq F, \quad \sum_{l=0}^{T-1} \hat{\sigma}_l^2 \leq F_1$$

and

$$\mathbb{P}\{E_{T-1}\} \geq 1 - \frac{(T-1)\beta}{N}, \quad \mathbb{P}\{E_①\} \geq 1 - \frac{\beta}{2N}, \quad \mathbb{P}\{E_③\} \geq 1 - \frac{\beta}{2N},$$

where

$$E_① = \left\{\text{either } \sum_{l=0}^{T-1} \sigma_l^2 > F \quad \text{or} \quad |①| \leq 4\gamma\lambda C R_0 \ln \frac{4N}{\beta}\right\},$$

$$E_③ = \left\{\text{either } \sum_{l=0}^{T-1} \hat{\sigma}_l^2 > F_1 \quad \text{or} \quad |③| \leq 32\gamma^2 \lambda^2 \ln \frac{4N}{\beta}\right\}.$$

Taking into account these inequalities and our assumptions on $\lambda$ and $\gamma$ (see (73) and (74)) we get that probability event $E_{T-1} \cap E_① \cap E_③$ implies

$$
\begin{aligned}
R_T^2 &\overset{(79)}{\leq} 2R_0^2 + 2\gamma^2 \sum_{l=0}^{T-1} \|\theta_l\|_2^2 - 2\gamma \sum_{l=0}^{T-1} \langle x^l - x^*, \theta_l \rangle \\
&\leq 2R_0^2 + \left( \frac{4}{7} + \frac{8}{567} + \frac{16}{49} + \frac{4}{441} + \frac{64}{321489} \right) C^2 R_0^2 \overset{(76)}{\leq} C^2 R_0^2.
\end{aligned} \tag{88}
$$

Moreover, using union bound we derive

$$
\mathbb{P}\left\{ E_{T-1} \cap E_① \cap E_③ \right\} = 1 - \mathbb{P}\left\{ \overline{E}_{T-1} \cup \overline{E}_① \cup \overline{E}_③ \right\} \geq 1 - \frac{T\beta}{N}. \tag{89}
$$

That is, by definition of $E_T$ and $E_{T-1}$ we have proved that

$$
\mathbb{P}\{E_T\} \overset{(88)}{\geq} \mathbb{P}\left\{ E_{T-1} \cap E_① \cap E_③ \right\} \overset{(89)}{\geq} 1 - \frac{T\beta}{N},
$$

which implies that for all $k = 0, 1, \ldots, N$ we have $\mathbb{P}\{E_k\} \geq 1 - \frac{k\beta}{N}$. Then, for $k = N$ we have that with probability at least $1 - \beta$

$$
AN \left( f(\bar{x}^N) - f(x^*) \right) \overset{(78)}{\leq} 2R_0^2 + 2\gamma^2 \sum_{k=0}^{N-1} \|\theta_k\|_2^2 - 2\gamma \sum_{k=0}^{N-1} \langle x^k - x^*, \theta_k \rangle \overset{(80)}{\leq} C^2 R_0^2.
$$

Since $A = 2\gamma \left( 1 - 2\gamma \left( \frac{1}{\varepsilon} \right)^{\frac{1-\nu}{1+\nu}} M_\nu^{\frac{2}{1+\nu}} \right) \overset{(74)}{\geq} \gamma$ we get that with probability at least $1 - \beta$

$$
f(\bar{x}^N) - f(x^*) \leq \frac{C^2 R_0^2}{AN} = \frac{C^2 R_0^2}{\gamma N}.
$$

When

$$
\gamma = \min \left\{ \frac{\varepsilon^{\frac{1-\nu}{1+\nu}}}{8 M_\nu^{\frac{2}{1+\nu}}}, \frac{R_0}{\sqrt{2N} \varepsilon^{\frac{\nu}{1+\nu}} M_\nu^{\frac{1}{1+\nu}}}, \frac{R_0^{1-\nu}}{2C^\nu M_\nu \ln \frac{4N}{\beta}} \right\}
$$

we have that with probability at least $1 - \beta$

$$
f(\bar{x}^N) - f(x^*) \leq \max \left\{ \frac{8C^2 M_\nu^{\frac{2}{1+\nu}} R_0^2}{\varepsilon^{\frac{1-\nu}{1+\nu}} N}, \frac{\sqrt{2} C^2 M_\nu^{\frac{1}{1+\nu}} R_0 \varepsilon^{\frac{\nu}{1+\nu}}}{\sqrt{N}}, \frac{2C^{2+\nu} M_\nu R_0^{1+\nu} \ln \frac{4N}{\beta}}{N} \right\}.
$$

Next, we estimate the iteration and oracle complexities of the method and consider 3 possible situations.

1. If $\gamma = \frac{\varepsilon^{\frac{1-\nu}{1+\nu}}}{8 M_\nu^{\frac{2}{1+\nu}}}$, then with probability at least $1 - \beta$

$$
f(\bar{x}^N) - f(x^*) \leq \frac{8C^2 M_\nu^{\frac{2}{1+\nu}} R_0^2}{\varepsilon^{\frac{1-\nu}{1+\nu}} N}.
$$

In other words, clipped-SGD achieves $f(\bar{x}^N) - f(x^*) \leq \varepsilon$ with probability at least $1 - \beta$ after

$$
\mathcal{O}\left( \frac{M_\nu^{\frac{2}{1+\nu}} R_0^2}{\varepsilon^{\frac{2}{1+\nu}}} \right)
$$

iterations and requires

$$
\begin{aligned}
Nm &\overset{(73)}{=} \mathcal{O}\left( \max \left\{ N, \frac{N^2 \sigma^2 \gamma^2 \ln \frac{N}{\beta}}{R_0^2} \right\} \right) = \mathcal{O}\left( \max \left\{ N, \frac{N^2 \varepsilon^{\frac{2(1-\nu)}{1+\nu}} \sigma^2 \ln \frac{N}{\beta}}{M_\nu^{\frac{4}{1+\nu}} R_0^2} \right\} \right) \\
&= \mathcal{O}\left( \max \left\{ \frac{M_\nu^{\frac{2}{1+\nu}} R_0^2}{\varepsilon^{\frac{2}{1+\nu}}}, \frac{\sigma^2 R_0^2}{\varepsilon^2} \ln \frac{M_\nu^{\frac{2}{1+\nu}} R_0^2}{\varepsilon^{\frac{2}{1+\nu}} \beta} \right\} \right)
\end{aligned}
$$

oracle calls.

39

2. If $\gamma = \frac{R_0}{\sqrt{2N}\varepsilon^{\frac{\nu}{1+\nu}}M_\nu^{\frac{1}{1+\nu}}}$, then with probability at least $1 - \beta$

$$f(\bar{x}^N) - f(x^*) \leq \frac{\sqrt{2}C^2 M_\nu^{\frac{1}{1+\nu}} R_0 \varepsilon^{\frac{\nu}{1+\nu}}}{\sqrt{N}}.$$

In other words, clipped-SGD achieves $f(\bar{x}^N) - f(x^*) \leq \varepsilon$ with probability at least $1 - \beta$ after

$$\mathcal{O}\left(\frac{M_\nu^{\frac{2}{1+\nu}} R_0^2}{\varepsilon^{\frac{2}{1+\nu}}}\right)$$

iterations and requires

$$Nm \overset{(73)}{=} \mathcal{O}\left(\max\left\{N, \frac{N^2\sigma^2\gamma^2 \ln\frac{N}{\beta}}{R_0^2}\right\}\right) = \mathcal{O}\left(\max\left\{N, \frac{N\sigma^2 \ln\frac{N}{\beta}}{\varepsilon^{\frac{2\nu}{1+\nu}} M_\nu^{\frac{2}{1+\nu}}}\right\}\right)$$

$$= \mathcal{O}\left(\max\left\{\frac{M_\nu^{\frac{2}{1+\nu}} R_0^2}{\varepsilon^{\frac{2}{1+\nu}}}, \frac{\sigma^2 R_0^2}{\varepsilon^2} \ln\frac{M_\nu^{\frac{2}{1+\nu}} R_0^2}{\varepsilon^{\frac{2}{1+\nu}}\beta}\right\}\right)$$

oracle calls.

3. If $\gamma = \frac{R_0^{1-\nu}}{2C^\nu M_\nu \ln\frac{4N}{\beta}}$, then with probability at least $1 - \beta$

$$f(\bar{x}^N) - f(x^*) \leq \frac{2C^{2+\nu} M_\nu R_0^{1+\nu} \ln\frac{4N}{\beta}}{N}.$$

In other words, clipped-SGD achieves $f(\bar{x}^N) - f(x^*) \leq \varepsilon$ with probability at least $1 - \beta$ after

$$\mathcal{O}\left(\frac{M_\nu R_0^{1+\nu} \ln\frac{M_\nu R_0^{1+\nu}}{\varepsilon\beta}}{\varepsilon}\right)$$

iterations and requires

$$Nm \overset{(73)}{=} \mathcal{O}\left(\max\left\{N, \frac{N^2\sigma^2\gamma^2 \ln\frac{N}{\beta}}{R_0^2}\right\}\right) = \mathcal{O}\left(\max\left\{N, \frac{N^2\sigma^2}{M_\nu^2 R_0^{2\nu} \ln\frac{N}{\beta}}\right\}\right)$$

$$= \mathcal{O}\left(\max\left\{\frac{M_\nu R_0^{1+\nu}}{\varepsilon}, \frac{\sigma^2 R_0^2}{\varepsilon^2}\right\} \ln\frac{M_\nu R_0^{1+\nu}}{\varepsilon\beta}\right)$$

oracle calls.

Putting all together and noticing that $\ln\frac{M_\nu^{\frac{2}{1+\nu}} R_0^2}{\varepsilon^{\frac{2}{1+\nu}}\beta} = \mathcal{O}\left(\ln\frac{M_\nu R_0^{1+\nu}}{\varepsilon\beta}\right)$ we get the desired result. $\quad\square$

As for clipped-SSTM it is possible to get rid of using large batchsizes without sacrificing the oracle complexity via a proper choice of $\gamma$, i.e., it is sufficient to choose

$$\gamma = \min\left\{\frac{\varepsilon^{\frac{1-\nu}{1+\nu}}}{8M_\nu^{\frac{2}{1+\nu}}}, \frac{R_0}{\sqrt{2N}\varepsilon^{\frac{\nu}{1+\nu}}M_\nu^{\frac{1}{1+\nu}}}, \frac{R_0^{1-\nu}}{2C^\nu M_\nu \ln\frac{4N}{\beta}}, \frac{R_0}{9\sigma N \ln\frac{4N}{\beta}}\right\}.$$

## C.2   Strongly convex case

In this section, we provide a full statement of Thm. 3.2 together with its proof. Note that due to strong convexity the solution $x^*$ is unique.

40

**Theorem C.2.** *Assume that function $f$ is $\mu$-strongly convex and its gradients satisfy* (3) *with $\nu \in [0,1]$, $M_\nu > 0$ on $Q = B_{2R_0} = \{x \in \mathbb{R}^n \mid \|x - x^*\|_2 \le 2R_0\}$, where $R_0 \ge \|x^0 - x^*\|_2$. Let $\varepsilon > 0$, $\beta \in (0,1)$, and for all $t = 1, \ldots, \tau$*

$$N_t = \max \left\{ \frac{2C^4 M_\nu^{\frac{2}{1+\nu}} R_0^2}{2^t \varepsilon_t^{\frac{2}{1+\nu}}}, \frac{4C^{2+\nu} M_\nu R_0^{1+\nu} \ln \frac{16 C^{2+\nu} M_\nu R_0^{1+\nu}}{2^{\frac{(1+\nu)t}{2}} \varepsilon_t \beta}}{2^{\frac{(1+\nu)t}{2}} \varepsilon_t} \right\}, \quad \varepsilon_t = \frac{\mu R_0^2}{2^{t+1}},$$

$$\lambda_t = \frac{R_0}{2^{\frac{t}{2}} \gamma_t \ln \frac{4 N_t \tau}{\beta}}, \quad m_t \ge \max \left\{ 1, \frac{81 N_t \sigma^2}{\lambda_t^2 \ln \frac{4 N_t \tau}{\beta}} \right\}, \quad \ln \frac{4 N_t \tau}{\beta} \ge 2,$$

$$\gamma_t = \min \left\{ \frac{\varepsilon_t^{\frac{1-\nu}{1+\nu}}}{8 M_\nu^{\frac{2}{1+\nu}}}, \frac{R_0}{2^{\frac{t}{2}} \sqrt{2 N_t} \varepsilon_t^{\frac{\nu}{1+\nu}} M_\nu^{\frac{1}{1+\nu}}}, \frac{R_0^{1-\nu}}{2^{1 + \frac{(1-\nu)t}{2}} C^\nu M_\nu \ln \frac{4 N_t \tau}{\beta}} \right\}.$$

*Then R-clipped-SGD achieves $f(\bar{x}^\tau) - f(x^*) \le \varepsilon$ with probability at least $1 - \beta$ after*

$$\mathcal{O} \left( \max \left\{ D_1^{\frac{2}{1+\nu}} \ln \frac{\mu R_0^2}{\varepsilon}, D_2^{\frac{2}{1+\nu}}, \max \left\{ D_1 \ln \frac{\mu R_0^2}{\varepsilon}, D_2 \right\} \ln \frac{D}{\beta} \right\} \right)$$

*iterations of Alg. 3 in total and requires*

$$\mathcal{O} \left( \max \left\{ D_1^{\frac{2}{1+\nu}} \ln \frac{\mu R_0^2}{\varepsilon}, D_2^{\frac{2}{1+\nu}}, \max \left\{ D_1 \ln \frac{\mu R_0^2}{\varepsilon}, D_2, \frac{\sigma^2}{\mu \varepsilon} \right\} \ln \frac{D}{\beta} \right\} \right) \tag{90}$$

*oracle calls, where*

$$D_1 = \frac{M_\nu}{\mu R_0^{1-\nu}}, \quad D_2 = \frac{M_\nu}{\mu^{\frac{1+\nu}{2}} \varepsilon^{\frac{1-\nu}{2}}}, \quad D = D_2 \ln \frac{\mu R_0^2}{\varepsilon}.$$

*Proof.* Applying Thm. C.1, we obtain that with probability at least $1 - \frac{\beta}{\tau}$

$$f(\hat{x}^1) - f(x^*) \le \frac{\mu R_0^2}{4}.$$

Since $f$ is $\mu$-strongly convex we have

$$\frac{\mu \|\hat{x}^1 - x^*\|_2^2}{2} \le f(\hat{x}^1) - f(x^*).$$

Therefore, with probability at least $1 - \frac{\beta}{\tau}$

$$f(\hat{x}^1) - f(x^*) \le \frac{\mu R_0^2}{4}, \quad \|\hat{x}^1 - x^*\|_2^2 \le \frac{R_0^2}{2}.$$

From mathematical induction and the union bound for probability events it follows that inequalities

$$f(\hat{x}^t) - f(x^*) \le \frac{\mu R_0^2}{2^{t+1}}, \quad \|\hat{x}^t - x^*\|_2^2 \le \frac{R_0^2}{2^t}$$

hold simultaneously for $t = 1, \ldots, \tau$ with probability at least $1 - \beta$. In particular, it means that after $\tau = \left\lceil \log_2 \frac{\mu R_0^2}{\varepsilon} \right\rceil - 1$ restarts R-clipped-SGD finds an $\varepsilon$-solution with probability at least $1 - \beta$. The total number of iterations $\hat{N}$ is

$$
\begin{aligned}
\sum_{t=1}^{\tau} N_t &= \mathcal{O} \left( \sum_{t=1}^{\tau} \max \left\{ \frac{M_\nu^{\frac{2}{1+\nu}} R_0^2}{2^t \varepsilon_t^{\frac{2}{1+\nu}}}, \frac{M_\nu R_0^{1+\nu}}{2^{\frac{(1+\nu)t}{2}} \varepsilon_t} \ln \frac{M_\nu R_0^{1+\nu} \tau}{2^{\frac{(1+\nu)t}{2}} \varepsilon_t \beta} \right\} \right) \\
&= \mathcal{O} \left( \sum_{t=1}^{\tau} \max \left\{ \frac{M_\nu^{\frac{2}{1+\nu}} \cdot 2^{\frac{(1-\nu)t}{1+\nu}}}{\mu^{\frac{2}{1+\nu}} R_0^{\frac{2(1-\nu)}{1+\nu}}}, \frac{M_\nu \cdot 2^{\frac{(1-\nu)t}{2}}}{\mu R_0^{1-\nu}} \ln \frac{M_\nu \cdot 2^{\frac{(1-\nu)\tau}{2}} \tau}{\mu R_0^{1-\nu} \beta} \right\} \right) \\
&= \mathcal{O} \left( \max \left\{ \frac{M_\nu^{\frac{2}{1+\nu}}}{\mu^{\frac{2}{1+\nu}} R_0^{\frac{2(1-\nu)}{1+\nu}}}, \frac{M_\nu}{\mu R_0^{1-\nu}} \ln \frac{M_\nu \ln \frac{\mu R_0^2}{\varepsilon}}{\mu^{\frac{1+\nu}{2}} \varepsilon^{\frac{1-\nu}{2}} \beta} \right\} \cdot \max \left\{ \ln \frac{\mu R_0^2}{\varepsilon}, \left( \frac{\mu R_0^2}{\varepsilon} \right)^{\frac{1-\nu}{2}} \right\} \right) \\
&= \mathcal{O} \left( \max \left\{ D_1^{\frac{2}{1+\nu}} \ln \frac{\mu R_0^2}{\varepsilon}, D_2^{\frac{2}{1+\nu}}, \max \left\{ D_1 \ln \frac{\mu R_0^2}{\varepsilon}, D_2 \right\} \ln \frac{D}{\beta} \right\} \right),
\end{aligned}
$$

41

850  where

$$D_1 = \frac{M_\nu}{\mu R_0^{1-\nu}}, \quad D_2 = \frac{M_\nu}{\mu^{\frac{1+\nu}{2}} \varepsilon^{\frac{1-\nu}{2}}}, \quad D = D_2 \ln \frac{\mu R_0^2}{\varepsilon}.$$

851  Finally, the total number of oracle calls equals

$$
\begin{aligned}
\sum_{t=1}^{\tau} \sum_{k=0}^{N_t-1} m_k^t &= \mathcal{O}\left( \max\left\{ \sum_{t=1}^{\tau} N_t, \sum_{t=1}^{\tau} \frac{\sigma^2 R_0^2}{2^t \varepsilon_t^2} \ln \frac{M_\nu R_0^{1+\nu} \tau}{2^{\frac{(1+\nu)t}{2}} \varepsilon_t \beta} \right\} \right) \\
&= \mathcal{O}\left( \max\left\{ \hat{N}, \sum_{t=1}^{\tau} \frac{\sigma^2 \cdot 2^t}{\mu^2 R_0^2} \ln \frac{D}{\beta} \right\} \right) = \mathcal{O}\left( \max\left\{ \hat{N}, \frac{\sigma^2}{\mu \varepsilon} \ln \frac{D}{\beta} \right\} \right).
\end{aligned}
$$

852  $\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad$ □

# D  Additional experimental details

## D.1  Main experiment hyper-parameters

In our experiments, we use standard implementations of Adam and SGD from PyTorch [32], we write only the parameters we changed from the default.

To conduct these experiments we used Nvidia RTX 2070s. The longest experiment (evolution of the noise distribution for image classification task) took 53 hours (we iterated several times over train dataset to build better histogram, see Appendix D.3).

### D.1.1  Image classification

For ResNet-18 + ImageNet-100 the parameters of the methods were chosen as follows:

- Adam: $lr = 1e - 3$ and a batchsize of $4 \times 32$

- SGD: $lr = 1e - 2$, $momentum = 0.9$ and a batchsize of $32$

- clipped-SSTM: $\nu = 1$, stepsize parameter $\alpha = 1e-3$ (in code we use separately $lr = 1e-2$ and $L = 10$ and $\alpha = \frac{lr}{L}$), norm clipping with clipping parameter $B = 1$ and a batchsize of $2 \times 32$. We also upper bounded the ratio $A_k/A_{k+1}$ by $0.99$ (see $a\_k\_ratio\_upper\_bound$ parameter in code).

- clipped-SGD: $lr = 5e - 2$, $momentum = 0.9$, coordinate-wise clipping with clipping parameter $B = 0.1$ and a batchsize of $32$

The main two parameters that we grid-searched were $lr$ and batchsize. For both of them we used logarithmic grid (i.e. for $lr$ we used $1e - 5, 2e - 5, 5e - 5, 1e - 4, \dots, 1e - 2, 2e - 2, 5e - 2$ for Adam). Batchsize was chosen from $32, 2 \cdot 32, 4 \cdot 32$ and $8 \cdot 32$. For SGD we also tried various momentum parameters.

For clipped-SSTM and clipped-SGD we used clipping level of $1$ and $0.1$ respectively. Too small choice of the clipping level, e.g. $0.01$, slow downs the convergence significantly.

Another important parameter for clipped-SSTM here, was $a\_k\_ratio\_upper\_bound$ – we used it to upper bound the maximum ratio of $A_k/A_{k+1}$. Without this modification the method is to conservative. e.g., after $10^4$ steps $A_k/A_{k+1} \sim 0.9999$. Effectively, it can be seen as momentum parameter of SGD.

### D.1.2  Text classification

For BERT + CoLA the parameters of the methods were chosen as follows:

- Adam: $lr = 5e - 5$, $weight\_decay = 5e - 4$ and a batchsize of $32$

- SGD: $lr = 1e - 3$, $momentum = 0.9$ and a batchsize of $32$

- clipped-SSTM: $\nu = 1$, stepsize parameter $\alpha = 8e - 3$, norm clipping with clipping parameter $B = 1$ and a batchsize of $8 \times 32$

- clipped-SGD: $lr = 2e - 3$, $momentum = 0.9$, coordinate-wise clipping with clipping parameter $B = 0.1$ and a batchsize of $32$

There we used the same grid as in the previous task. The main difference here is that we didn't bound clipped-SSTM $A_k/A_{k+1}$ ratio – there are only $\sim 300$ steps of the method (because the batch size is $8 \cdot 32$), thus the the method is still not too conservative.

## D.2  On the relation between stepsize parameter $\alpha$ and batchsize

In our experiments, we noticed that clipped-SSTM show similar results when the ration $bs^2/\alpha$ is kept unchanged, where $bs$ is batchsize (see Fig. 3). We compare the performance of clipped-SSTM with $4$ different choices of $\alpha$ and the batchsize.
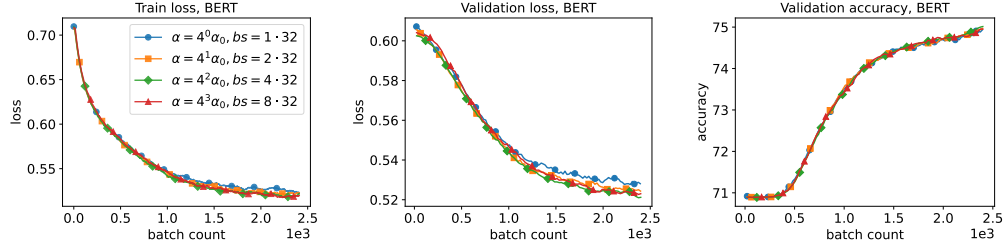
Figure 3: Train and validation loss + accuracy for clipped-SSTM with different parameters. Here $\alpha_0 = 0.000125$, $bs$ means batchsize. As we can see from the plots, increasing $\alpha$ 4 times and batchsize 2 times almost does not affect the method's behavior.

Thm. B.1 explains this phenomenon in the convex case. For the case of $\nu = 1$ we have (from (34) and (39)):

$$\alpha \sim \frac{1}{aM_1}, \quad \alpha_k \sim k\alpha, \quad m_k \sim \frac{Na\sigma^2\alpha_{k+1}^2}{C^2 R_0^2 \ln\frac{4N}{\beta}}, \quad N \sim \frac{a^{\frac{1}{2}}CR_0 M_1^{\frac{1}{2}}}{\varepsilon^{\frac{1}{2}}} \sim \frac{CR_0}{\alpha^{\frac{1}{2}}\varepsilon^{\frac{1}{2}}},$$

whence

$$m_k \sim \frac{CR_0 a\sigma^2\alpha^2(k+1)^2}{\alpha^{\frac{1}{2}}\varepsilon^{\frac{1}{2}}C^2 R_0^2 \ln\frac{4N}{\beta}} \sim \frac{\sigma^2\alpha^2(k+1)^2}{\alpha^{\frac{1}{2}}\alpha M_1 \varepsilon^{\frac{1}{2}}CR_0 \ln\frac{4N}{\beta}} \sim \alpha^{\frac{1}{2}},$$

where the dependencies on numerical constants and logarithmic factors are omitted. Therefore, the observed empirical relation between batchsize ($m_k$) and $\alpha$ correlates well with the established theoretical results for clipped-SSTM.

## D.3 Evolution of the noise distribution

In this section, we provide our empirical study of the noise distribution evolution along the trajectories of different optimizers. As one can see from the plots, the noise distribution for ResNet-18 + ImageNet-100 task is always close to Gaussian distribution, whereas for BERT + CoLA task it is significantly heavy-tailed.
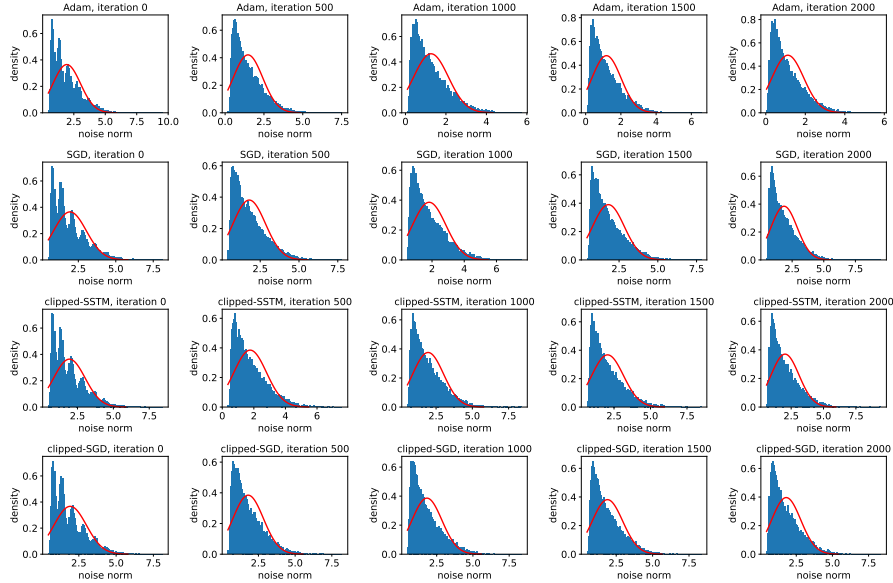


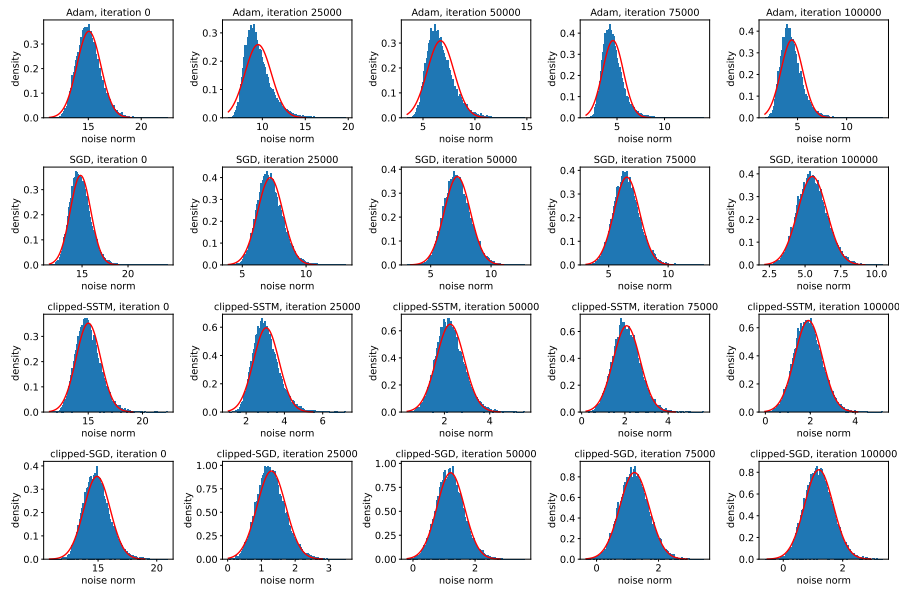Figure 4: Evolution of the noise distribution for BERT + CoLA task.

Figure 5: Evolution of the noise distribution for `ResNet-18` + `ImageNet-100` task.