# Escaping Saddle Points in Zeroth-order Optimization:
# The Power of Two-point Estimators

**Zhaolin Ren** [1]  **Yujie Tang** [2]  **Na Li** [1]

## Abstract

Two-point zeroth order methods are important in many applications of zeroth-order optimization, such as robotics, wind farms, power systems, online optimization, and adversarial robustness to black-box attacks in deep neural networks, where the problem may be high-dimensional and/or time-varying. Most problems in these applications are nonconvex and contain saddle points. While existing works have shown that zeroth-order methods utilizing $\Omega(d)$ function valuations per iteration (with $d$ denoting the problem dimension) can escape saddle points efficiently, it remains an open question if zeroth-order methods based on two-point estimators can escape saddle points. In this paper, we show that by adding an appropriate isotropic perturbation at each iteration, a zeroth-order algorithm based on $2m$ (for any $1 \leq m \leq d$) function evaluations per iteration can not only find $\epsilon$-second order stationary points polynomially fast, but do so using only $\tilde{O}(d/m\epsilon^2\bar{\psi})$ function evaluations, where $\bar{\psi} \geq \tilde{\Omega}(\sqrt{\epsilon})$ is a parameter capturing the extent to which the function of interest exhibits the strict saddle property.

## 1. Introduction

Two-point estimators, which approximate the gradient using two function evaluations per iteration, have been widely studied by researchers in the zeroth-order optimization literature, in convex (Nesterov & Spokoiny, 2017; Duchi et al., 2015; Shamir, 2017), nonconvex (Nesterov & Spokoiny, 2017), online (Shamir, 2017), as well as distributed settings (Tang et al., 2019). A key reason for doing so is that for applications of zeroth-order optimization arising in robotics

(Li et al., 2022), wind farms (Tang et al., 2020a), power systems (Chen et al., 2020), online (time-varying) optimization (Shamir, 2017), learning-based control (Malik et al., 2019; Li et al., 2021), and improving adversarial robustness to black-box attacks in deep neural networks (Chen et al., 2017), it may be costly or impractical to wait for $\Omega(d)$ (where $d$ denotes the problem dimension) function evaluations per iteration to make a step. This is especially true for high-dimensional and problems with time-varying noise. See Appendix A for more discussion.

However, despite the advantages of zeroth-order methods with two-point estimators, there has been a lack of existing work studying the ability of two-point estimators to escape saddle points in nonconvex optimization problems. Since nonconvex problems arise often in practice, it is crucial to know if two-point algorithms can efficiently escape saddle points of nonconvex functions and converge to second-order stationary points (see Definition 1 for a definition).

To motivate the challenges of escaping saddle points using two-point zeroth-order methods, we begin with a review of escaping saddle points using first-order methods. The problem of efficiently escaping saddle points in deterministic first-order optimization (with exact gradients) has been carefully studied in several earlier works (Jin et al., 2017; 2018b). A key idea in these works is the injection of an isotropic perturbation whenever the gradient is small, facilitating escape from a saddle if a negative curvature direction exists even without actively identifying the direction. However, the analysis of efficient saddle point escape for stochastic gradient methods is often more complicated. In general, the behavior of the stochastic gradient near the saddle point can be difficult to characterize. Hence, strong concentration assumptions are typically made on the stochastic gradients being used, such as subGaussianity, boundedness of the variance or a bounded gradient estimator (Ge et al., 2015; Daneshmand et al., 2018; Xu et al., 2018; Fang et al., 2019; Roy et al., 2020; Vlaski & Sayed, 2021b), creating an analytical issue when such idealized assumptions fail to hold.

Indeed, though zeroth-order methods can be viewed as stochastic gradient methods, common zeroth order estimators, such as two-point estimators (Nesterov & Spokoiny, 2017), are not subGaussian, and can have unbounded vari-

---

[1] John A. Paulson School of Engineering and Applied Sciences, Harvard University [2] Department of Industrial Engineering & Management at Peking University. Correspondence to: Zhaolin Ren <zhaolinren@g.harvard.edu>, Na Li <nali@seas.harvard.edu>.

ance. For instance, it can be shown that the variance of the two-point estimator is on the order of $\Omega(d\|\nabla f(x)\|^2)$ (Nesterov & Spokoiny, 2017), with both a dependence on the problem dimension $d$ as well as on the norm of the gradient, which can be unbounded. Due to non-subGaussianity and unboundedness, it is tricky to bound the effect of such zeroth-order estimators and establish tight concentration inequalities that facilitate its escape near saddle points. In addition, the large variance of the zeroth-order estimator is also an issue in non-saddle regions, i.e. when the gradient is large. While this is not an issue to show function improvement in expectation, as we discuss later, this becomes an issue when guaranteeing high probability bounds.

Due to these difficulties, previous works on escaping saddle points in zeroth-order optimization have exclusively focused on approaches requiring $\Omega(d)$ function evaluations per iteration to accurately estimate the gradient (Jin et al., 2018a; Bai et al., 2020; Vlatakis-Gkaragkounis et al., 2019), or in some cases negative curvature directions (Zhang et al., 2022; Lucchi et al., 2021) or the Hessian itself (Balasubramanian & Ghadimi, 2022), reducing in a sense the zeroth-order problem back to a first-order one. However, as explained earlier, two-point zeroth-order algorithms are important for high-dimensional and/or time-varying problems in many applications areas. This raises an important question:

**Can two-point zeroth-order methods escape saddle points and reach approximate second order stationary points efficiently?**

**Our Contribution.** In this work, we show that by adding an appropriate isotropic perturbation at each iteration, a zeroth-order algorithm based on *any* number $m$ of pairs ($1 \le m \le d$) of function evaluations per iteration can not only find ($\epsilon, \sqrt{\epsilon}$)-second order stationary points (cf. the definition later in Definition 1) polynomially fast, but do so using only $\tilde{O}(\mathrm{polylog}(\frac{1}{\cdot})d/\epsilon^{2.5})$ function evaluations, with a probability of at least $1 - \delta$. In particular, this proves that using a single two-point zeroth-order estimator at each iteration (with appropriate perturbation) suffices to efficiently escape saddle points in zeroth-order optimization, with high probability. Moreover, for functions that are ($\epsilon, \psi, O(\sqrt{\epsilon})$) strict-saddle (see Definition 3 for a definition of strict saddle functions), our results become $\tilde{O}(\mathrm{polylog}(\frac{1}{\cdot})d/\psi\epsilon^2)$, which is a significant improvement when $\psi \gg \epsilon$; strict saddle functions have been identified as an important class of functions in nonconvex optimization, with several well-known examples such as tensor decomposition (Ge et al., 2015), dictionary learning and phase retrieval (Sun et al., 2015). A comparison of our results with existing zeroth-order and first-order methods is shown in Table 1. We also provide numerical results in Section 4 showing that our proposed two-point algorithm requires fewer total function evaluations to converge than zeroth order methods that use $2d$ function evaluations per iteration, for a nonconvex test function proposed in (Du et al., 2017).

To overcome the theoretical challenges that were discussed earlier, we i) first show, via a careful analysis, that zeroth order methods can make function value improvement across iterates with large gradients with high probability, even when only a single two-point estimator (which can have significant variance at large gradients) is used per iteration. ii) Second, near saddle points, we overcome issues caused by the unbounded variance and non-subGaussinity of zeroth-order gradient estimators by developing new technical tools, including novel martingale concentration inequalities involving Gaussian vectors, to tightly bound such terms. In turn, this allows us to show that the noise emanating from the zeroth-order estimators will not overwhelm the effect of the additional isotropic perturbative noise, facilitating escape along negative curvature directions. To the best of our knowledge, both analyses are novel, and may be independent contributions on their own.

*Related Work.* Due to space considerations, we defer a full discussion of related work to Appendix A.

## 2. Problem Setup

We make the following assumptions on the class of functions $f : \mathbb{R}^d \to \mathbb{R}$ which we consider.

**Assumption 1** (Properties of $f$). *We suppose that $f : \mathbb{R}^d \to \mathbb{R}$ satisfies the following properties:*

1. *$f$ is twice-differentiable and lower bounded, i.e. $f := \min_x f(x) > -\infty$.*

2. *$f$ is $L$-gradient Lipschitz, i.e.*

$$\|\nabla f(x) - \nabla f(y)\| \le L\|x - y\| \ \ \forall x, y \in \mathbb{R}^d.$$

3. *$f$ is $\rho$-Hessian Lipschitz, i.e.*

$$\nabla^2 f(x) - \nabla^2 f(y) \ \ \le \rho\|x - y\| \ \ \forall x, y \in \mathbb{R}^d.$$

In our work, we focus on finding approximate second order stationary points, defined below.

**Definition 1.** A point $x \in \mathbb{R}^d$ is an ($\epsilon, \varphi$)-second order stationary point if

$$\|\nabla f(x)\| < \epsilon, \quad \text{and} \quad \lambda_{\min}(\nabla^2 f(x)) > -\varphi.$$

We define an ($\epsilon, \varphi$)-approximate saddle point as follows.

**Definition 2.** A point $x \in \mathbb{R}^d$ is an ($\epsilon, \varphi$)-approximate saddle point, if

$$\|\nabla f(x)\| < \epsilon, \quad \text{and} \quad \lambda_{\min}(\nabla^2 f(x)) \le -\varphi.$$

| | | Iteration Complexity | Fun. Evaluations. per iter |
|---|---|---|---|
| First-order | (Jin et al., 2017) (deterministic) | $\tilde{O}\left(\frac{1}{\epsilon^2}\right)$ | — |
| | (Fang et al., 2019) (SGD) | $\tilde{O}\left(\frac{1}{\epsilon^{3.5}}\right)$ | — |
| Zeroth-order | (Jin et al., 2018a) | $\tilde{O}\left(\frac{1}{\epsilon^2}\right)$ | $\tilde{O}\left(\frac{d^2}{\epsilon^3}\right)$ |
| | (Bai et al., 2020) | $\tilde{O}\left(\frac{1}{\epsilon^2}\right)$ | $\tilde{O}\left(\frac{d^2}{\epsilon^8}\right)$ |
| | (Vlatakis-Gkaragkounis et al., 2019) | $\tilde{O}\left(\frac{1}{\epsilon^2}\right)$ | $\tilde{O}(d)$ |
| | (Balasubramanian & Ghadimi, 2022) | $\tilde{O}\left(\frac{1}{\epsilon^{1.5}}\right)$ | $\tilde{O}\left(\frac{d}{\epsilon^2} + \frac{d^4}{\epsilon}\right)$ |
| | (Lucchi et al., 2021)$^{y}$ | $\tilde{O}\left(\frac{1}{\epsilon^2}\right)$ | $\tilde{O}\left(\frac{d}{\epsilon^{2.3}}\right)$ |
| | (Zhang et al., 2022) | $\tilde{O}\left(\frac{1}{\epsilon^2}\right)$ | $\tilde{O}(d)$ |
| | Algorithm 1 (this paper, $1 \le m \le d$)‡ | $\tilde{O}\left(\frac{d}{\epsilon^2 \bar{\psi} m}\right)$ | $2m$ |

Table 1: Selected comparison of convergence results to $(\epsilon, O(\sqrt{\epsilon})$-second order stationary points in smooth, nonconvex functions; for $^{y}$, the convergence is to $(\epsilon, \epsilon^{2/3})$-second order stationary points. For $^{z}$, the term $\bar{\psi}$ in the denominator is (i) $\psi$ when the function $f$ is $(\epsilon, \psi, O(\sqrt{\epsilon}))$-strict saddle for a $\psi > O(\sqrt{\epsilon})$ (see Definition 3 for a definition) and (ii) $O(\sqrt{\epsilon})$ if otherwise.

Following past convention (Jin et al., 2019a), we will focus in particular on escaping $(\epsilon, \sqrt{\rho\epsilon})$-saddle points. For notational simplicity, in following text, we refer to $(\epsilon, \sqrt{\rho\epsilon})$-saddle points simply as $\epsilon$-saddle points and $(\epsilon, \sqrt{\rho\epsilon})$-second order stationary points as $\epsilon$-second order stationary points. Beyond the definition of $\epsilon$-approximate saddle points above, it is known that many nonconvex functions with saddle points, such as orthogonal tensor decomposition (Ge et al., 2015), phase retrieval and dictionary learning (Sun et al., 2015), satisfy what is known as a *strict saddle* condition (Ge et al., 2015). For the Hessians of the saddle points of such functions, there is always a strict negative eigenvalue whose magnitude is bounded from below. We provide a precise definition below.

**Definition 3.** A twice-differential function $f(x)$ is $(\epsilon, \psi, \varrho)$-strict saddle for any $\psi > \varrho > 0$, if for any point $x$, either

1. $\|\nabla f(x)\| \ge \epsilon$ holds,

2. or when $\|\nabla f(x)\| < \epsilon$ holds, either

   (a) $\lambda_{\min}(\nabla^2 f(x)) \le -\psi$, or
   (b) $\lambda_{\min}(\nabla^2 f(x)) > -\varrho$.

In our work, we consider the following batch symmetric two-point zeroth-order estimator.

**Definition 4** ((Batch) two-point zeroth-order estimator with perturbation). We define a $m$-batch two-point zeroth order estimator as follows:

$$g_u^{(m)}(x) := \frac{1}{m} \sum_{i=1}^{m} \frac{f(x + uZ_i) - f(x - uZ_i)}{2u} Z_i, \quad (1)$$

where $Z_i \overset{i.i.d}{\sim} N(0, I)$, and $u > 0$ is a smoothing radius.

---

**Algorithm 1** Zeroth-order perturbed gradient descent (ZOPGD)

**input :** $x_0$, horizon $T$, step-size $\eta$, smoothing radius $u$, perturbation radius $r$, batch size $m$

**for** *step* $t = 0, \dots, T$ **do**

　　Sample $Z^{(m)} = \{Z_{t,i}\}_{i=1}^{m} \sim N(0, I)$ to compute $g_u^{(m)}(x_t)$, defined in Eq. (1).

　　Update $x_{t+1} = x_t - \eta\, g_u^{(m)}(x_t) + Y_t$, where $Y_t \sim N(0, \frac{r^2}{d} I)$

---

Such $2m$ zeroth-order gradient estimators have frequently been studied in zeroth-order optimization works (see e.g. (Nesterov & Spokoiny, 2017)). To facilitate efficient escape from saddle points, our proposed Algorithm 1 adds isotropic perturbation at each iteration.

We now state an informal version of our main result, and follow that with a few remarks.

**Theorem 1** (Main result, informal version of Theorem 2). *Consider running Algorithm 1. Let $\tilde{O}$ hide polylogarithmic terms in $\delta$ and other parameters. Suppose $\delta \in (0, 1/e]$. Suppose $\sqrt{\rho\epsilon} \le \min\{1, L\}^{1}$, such that $\bar{\psi} \le \min\{1, L\}$, where*

$$\bar{\psi} := \begin{cases} \min\{\psi, \sqrt{Lg}\} & \text{if } \varrho > \sqrt{\rho\epsilon} \text{ s.t. } f(\cdot) \text{ is } (\epsilon, \psi, \sqrt{\rho\epsilon})\text{-strict saddle} \\ \sqrt{\rho\epsilon} & \text{if otherwise.} \end{cases}$$

(2)

---

[1] In our paper, we focus on the case $\sqrt{\rho\epsilon} \le L$; otherwise, by the $L$-Lipschitz assumption, $\lambda_{\min}(\nabla^2 f(x)) \ge -L$ for all $x \in \mathbb{R}^d$, which implies $\epsilon$-first order stationary points are also $\epsilon$-second order stationary points.

*Suppose*

$$\eta = \mathcal{O}\left(\frac{\min\{\rho^{-}, \rho\bar{\psi}, \bar{r}g\}}{\rho^{-}d}\right); \quad r = \mathcal{O}(\cdot); \quad = \mathcal{O}\left(\frac{m}{d\max\{L, L^2 g\}}\right);$$

*Then, in*

$$
\begin{aligned}
T &= \tilde{\Omega}\left(\frac{(f(x_0) - f^*)}{\eta\epsilon^2} + \frac{\rho^2(f(x_0) - f^*)}{\eta\bar{\psi}^4}\right) \\
&= \tilde{\Omega}\left(\frac{d\max\{L, L^2\}\rho^2(f(x_0) - f^*)}{m\bar{\psi}\epsilon^2}\right)
\end{aligned}
$$

*iterations (with each iteration using $2m$ function evaluations), with probability at least $1 - \delta$, at least half the iterates are $(\epsilon, \sqrt{\rho\epsilon})$-second-order stationary points.*

*Remark* 1. As the choice of $\eta$ in Proposition 4 (Appendix D) and Theorem 2 (Appendix F) respectively imply, the $\tilde{\Omega}\left(\frac{f(x_0) - f^*}{\eta\epsilon^2}\right)$ term in the sample complexity comes from the large gradient iterations (Proposition 4), whereas the $\tilde{\Omega}\left(\frac{\rho^2(f(x_0) - f^*)}{\eta\bar{\psi}^4}\right)$ term comes from the escape saddle point phase.

*Remark* 2. As a corollary of Theorem 1, for functions $f$ which are $(\epsilon, \psi, \sqrt{\rho\epsilon})$ strict saddle, assuming that $\psi \geq \sqrt{\rho\epsilon}$, the sample complexity of our algorithm scales as $\tilde{\Omega}\left(\frac{d\max\{L^2, Lg(f(x_0) - f^*)\}}{m\epsilon^2\psi}\right)$, which scales as $\tilde{\Omega}\left(\frac{d}{m\epsilon^2}\right)$ when $\psi$ is of size $\Omega(1)$. Thus, in this setting, for two-point estimators, where $m = 1$, the dependence on $d$ and $\epsilon$ in our sample complexity (as measured by function evaluations) matches that achieved by the algorithms in (Vlatakis-Gkaragkounis et al., 2019; Zhang et al., 2022), which have to use $2d$ function evaluations per iteration to estimate the gradient.

**Comparison to gradient-based methods.** For first-order escape saddle point algorithms, standard perturbation-based methods (without acceleration) can find a $(\epsilon, O(\sqrt{\epsilon}))$-second-order stationary point using $\tilde{O}(1/\epsilon^2)$ iterations for deterministic GD (Jin et al., 2019a), while for standard SGD the best-known rates are slower at $\tilde{O}(1/\epsilon^{3.5})$ (Fang et al., 2019). In contrast, our sample complexity (as measured by the total number of function evaluations) is $\tilde{O}\left(\frac{d}{\epsilon^2\bar{\psi}}\right)$, where $\bar{\psi}$ is defined in Eq. (2). The extra (linear) dependence on $d$ is typical for zeroth-order algorithms (see e.g. (Nesterov & Spokoiny, 2017)); intuitively, gradient calculation for $d$-dimensional functions requires $O(d)$ calculations agnostically, so it makes sense that zeroth-order algorithms requires $d$ times more iterations. For general non strict-saddle functions, our dependence on $\epsilon$ sits between that of the deterministic methods and SGD methods, and suggests the benefit of a specialized treatment of zeroth-order methods over considering them simply as a subclass of SGD methods. Moreover, for $(\epsilon, \psi, \sqrt{\rho\epsilon})$- strict-saddle functions where $\psi = \Omega(1)$, our sample complexity becomes $\tilde{O}(\frac{d}{\epsilon^2})$, with an $\epsilon$ dependence that matches that of the best existing

sample complexity for non-accelerated first-order escape saddle point methods (Jin et al., 2017)

**Comparison to existing zeroth-order methods.** As Table 1 suggests, our sample complexity significantly outperforms that of (Jin et al., 2018a), (Bai et al., 2020), (Balasubramanian & Ghadimi, 2022), and also that in (Lucchi et al., 2021), which is a random search method. We note that the sample complexity in (Vlatakis-Gkaragkounis et al., 2019; Zhang et al., 2022) outperform our method, with a function evaluation complexity of $\tilde{O}\left(\frac{d}{\epsilon^2}\right)$. However, for for $(\epsilon, \psi, \sqrt{\rho\epsilon})$- strict-saddle functions where $\psi = \Omega(1)$, our sample complexity becomes $\tilde{O}(\frac{d}{\epsilon^2})$, which matches the sample complexity in (Vlatakis-Gkaragkounis et al., 2019; Zhang et al., 2022). Moreover, a key limitation of their methods is a requirement to use $\Omega(d)$ function evaluations to estimate the gradient at each iteration, which may not be practical in realistic applications when $d$ is large. In contrast, our method supports any number of function evaluations at each iteration between 1 to $d$. Moreover, numerically, we found that for a test nonconvex function proposed in (Du et al., 2017), our method (with two-point estimators) takes fewer function evaluations to escape saddle points and converge to the global minimum than the methods in (Vlatakis-Gkaragkounis et al., 2019; Zhang et al., 2022); see Section 4 for details.

# 3. Proof strategy and key challenges in the zeroth-order setting

Broadly speaking, our proof include two major parts, i) characterizing the progress made in iterations when the gradient is large (which we can define to be iterations $t$ where $\|\nabla f(x_t)\| \geq \epsilon$) (Section 3.1), ii) and iterations when we are at an $\epsilon$-approximate saddle point (where progress may be made along the negative eigendirection of the Hessian matrix) (Section 3.2). While the approach is similar to the first-order case (e.g. (Jin et al., 2019a)), the zeroth-order setting brings forth several unique challenges. In the rest of this section, we explain these challenges, sketch out our high-level proof outlines, and provide statements of the main technical results. Due to limited space, we defer the full proof to the Appendix.

## 3.1. Showing function decrease when gradients are large

**Challenge.** Due to the noise in two-point (or $2m$ where $m$ is a small constant) zeroth-order gradient, even when the gradient is large, it may not always be possible to make progress at each iteration, especially when $m < d$ is used in the gradient estimation equation in Eq. (1). While it is tempting to use an expectation-based argument to handle this issue, it is known that expectation-based function

decrease arguments are insufficient for the purpose of escaping saddle points (see e.g. Proposition 1 in (Ziyin et al., 2021)). We tackle this issue by using high-probability arguments instead; we note that achieving these high-probability bounds is highly nontrivial due to the large variance of the two-point zeroth-order estimator (scaling with $d$ times the squared norm of the gradient). Hence, any single iteration of the zeroth-order method may in fact lead to a function increase rather than decrease.

**High-level proof outline. (i)** We first characterize the function value change for our proposed algorithm (Lemma 1). **(ii)** Next, we tackle the issue of the possibility that the function value might increase for any given iteration. The key idea here is that across any small consecutive number of iterations, there will be one iteration where the zeroth-order estimator is sufficiently aligned with the gradient direction (Lemma 14 in Appendix D). **(iii)** Along with a series of other technical results in Appendix D, we then show that the function makes sufficient progress across the duration of the algorithm, with high probability (Proposition 1). To more concretely illustrate the key analytical challenge, we next introduce the following function decrease lemma, proved in Appendix D.

**Lemma 1** (Function decrease for batch zeroth-order optimization). *Suppose at each time $t$, the algorithm performs the update step (with batch-size parameter $1 \leq m \leq d$)*

$$x_{t+1} = x_t - \eta \left( g_u^{(m)}(x_t) + Y_t \right),$$

*where*

$$g_u^{(m)}(x_t) = \frac{1}{m} \sum_{i=1}^{m} \frac{f(x_t + u Z_{t,i}) - f(x_t - u Z_{t,i})}{2u} Z_{t,i},$$

*where each $Z_{t,i}$ is drawn i.i.d from $N(0,I)$, $u > 0$ is the smoothing radius, and $Y_t \sim N(0, \frac{r^2}{d} I)$ with $r > 0$ denoting the perturbation radius.*

*Then, there exist absolute constants $c_1 > 0, C_1 \geq 1$ such that, for any $T \in \mathbb{Z}^+$ and $T \geq \tau > 0$, $\alpha > 0$ and $\delta \in (0, 1/e]$, upon defining $\mathcal{H}_{0,\tau}(\delta)$ to be the event on which the inequality*

$$f(x_\tau) - f(x_0) \tag{3}$$

$$\leq -\frac{3}{4} \eta \sum_{t=0}^{\tau-1} \frac{1}{m} \sum_{i=1}^{m} \left( Z_{t,i}^\top \nabla f(x_t) \right)^2 \tag{4}$$

$$+ \left( -\eta + \frac{c_1 L \eta^2 \chi^3 d}{m} \right) \sum_{t=0}^{\tau-1} \|\nabla f(x_t)\|^2$$

$$+ \eta u^4 \chi^2 c_1 d^3 \left( \log \frac{T}{\eta} \right)^3 + \eta^3 L^2 u^4 \chi^2 c_1 d^4 \left( \log \frac{T}{\eta} \right)^4$$

$$+ \eta c_1 r^2 (\eta + L) \log \frac{T}{\eta} + \eta c_1 L \chi^2 r^2 \tag{5}$$

*is satisfied (where $\chi := \log(C_1 dmT/\delta)$), we have*

$$\mathbb{P}(\mathcal{H}_{0,\tau}(\delta)) \geq 1 - \frac{(\alpha + 4)}{T} \delta; \quad \mathbb{P}(\cap_{\tau=1}^{\tau^\theta} \mathcal{H}_{0,\tau}(\delta)) \geq 1 - \frac{5\delta^\theta}{T}$$

*for any $0 \leq \tau^\theta \leq T$.*

Our goal is to show that we can arrive at a contradiction $f(x_T) < \min_x f(x)$ when there is a large number of steps at which $\|\nabla f(x_t)\| \geq \epsilon$ (Proposition 1). As we can see from Eq. (5), this implies that we need to prove a lower bound of the form

$$\sum_{t=0}^{T-1} \frac{1}{m} \sum_{i=1}^{n} \left( Z_{t,i}^\top \nabla f(x_t) \right)^2 \geq \left( 1 + \frac{c_1 L \eta \chi^3 d}{m} \right) \sum_{t=0}^{T-1} \frac{1}{\alpha} \|\nabla f(x_t)\|^2 \tag{6}$$

for some $\alpha$ which is not too large (an example would be picking $\alpha$ such that it only scales logarithmically in the problem parameters). However, it is tricky to prove such a lower-bound in the zeroth-order setting. In particular, for small batch-sizes $m$, $\frac{1}{m} \sum_{i=1}^{m} \left( Z_{t,i}^\top \nabla f(x_t) \right)^2$ could be small even as $\|\nabla f(x_t)\|^2$ is large; this is because for each $i \in [m]$, $Z_{t,i}$ could have a negligible component in the $\nabla f(x_t)$ direction. This necessitates a more delicate analysis to prove a bound similar to Eq. (6). Due to space reasons, we defer our more detailed proof approach outline to Appendix D (see the discussion immediately following Lemma 1). The results in Appendix D culminates in the following result which limits the number of large-gradient.

**Proposition 1** (Bound on number of iterates with large gradients, informal version of Proposition 4). *Let $\delta \in (0, 1/e]$ be arbitrary. Letting $\tilde{O}$ hide polylogarithmic dependencies on $\delta$ (and other parameters), consider choosing $u, r, \eta$ and $T$ such that*

$$u = \tilde{O}\left( \frac{\rho_-}{\rho - d} \right); \quad r = O(\epsilon); \quad \eta = \tilde{\Theta}\left( \frac{m}{dL} \right);$$

$$T = \tilde{\Theta}\left( \frac{(f(x_0) - f^*) + \epsilon^2 =L}{\eta \epsilon^2} \right).$$

*Then, with probability at least $1 - O(\delta)$, there are at most $T/4$ iterations for which $\|\nabla f(x_t)\| \geq \epsilon$.*

### 3.2. Making progress near saddle points

**Challenge.** The noise in two-point zeroth-order estimators makes the analysis around $\epsilon-$approximate saddle points challenging, because the concentration properties of the (non-subGaussian) noise are hard to characterize. Intuitively, a noisier estimator might facilitate easier escape from saddle point. However, without an appropriate concentration bound, the noise may behave in unpredictable ways, preventing escape from saddle regions. Previous analysis of saddle point escape using stochastic estimators typically requires these estimators to satisfy subGaussian properties (Jin et al., 2019a; Fang et al., 2019), which zeroth-order estimators do not satisfy.

**High-level proof outline. (i)** We first prove a technical result showing that the travelling distance of the iterates can be bounded in terms of the function value decrease

(i.e., Improve or Localize, Lemma 2). **(ii)** Next, at any $\epsilon$-saddle point, we consider a coupling argument and define two sequences running near-identical zeroth-order dynamics, differing only in the sign of their perturbative term along the minimum eigendirection of $H$, which denotes the Hessian of the saddle (Lemma 3). Using Lemma 2 in point (i), if we assume for contradiction that the two sequences both "get stuck" and make little function value progress, the dynamics of the difference between the two sequences will remain small as both sequences remain close to the saddle point. **iii)** However, since the perturbation vectors of the two sequences differ in the (most) negative direction of $H$, the norm of the the difference of the two sequences will grow exponentially so long as *a).* the sequences remain close to the saddle point (and thus the Hessian has a negative curvature direction) and *b).* the effect of the zeroth-order stochastic noise can be controlled. This leads to a contradiction, implying that sufficient function decrease must have been made (Proposition 5 in Appendix E.3). **(iv)** To show that the zeroth-order stochastic noise can be controlled, we prove one technical result (Proposition 2), providing a concentration bound for the product of (possibly unbounded) subGaussian random vectors that scales linearly with the dimension $d$. This enables us to control the effect of the zeroth-order noise near saddle points, and is essential in showing that the eventual sample complexity scales linearly with $d$.

We provide a more detailed proof sketch below, where we elaborate more on our analytical challenges and ideas. We first introduce an informal statement of a key technical result that bounds, with high probability, the travelling distance of the iterates in terms of the function value decrease.

**Lemma 2** (Improve or Localize, informal version of Lemma 23). *Consider the perturbed zeroth-order update Algorithm 1. Let $\delta \in (0, 1/e)$ be arbitrary. Consider any $T_s = \tilde{\Omega}\left(\frac{1}{m}\log(1/\delta)\right)$, and any $t_0 \geq 0$. For any $F > 0$, suppose $f(x_{T_s+t_0}) - f(x_{t_0}) > -F$, i.e. $f(x_{t_0}) - f(x_{T_s+t_0}) < F$. Letting $\tilde{O}$ hide polylogarithmic terms involving $\delta$, suppose*

$$u = \mathcal{O}\left(\frac{\min\{\rho^-; \rho\bar{r}g\}}{\rho^{-d}}\right); \quad r = \mathcal{O}\left(\min\left\{; \frac{F}{T_s}\right\}\right);$$

$$= \mathcal{O}\left(\frac{m^{\rho^-}}{dL}\right):$$

*Then, with probability at least $1 - O\left(\frac{T_s\delta}{T}\right)$ (here $T \geq T_s$ denotes the total number of iterations), for each $\tau \in \{0, 1, \dots, T_s\}$, we have that*

$$\|x_{t_0+\tau} - x_{t_0}\|^2 \leq \phi_{T_s}(\delta, F), \quad where$$

$$\phi_{T_s}(\delta, F) = \tilde{O}\left(\max\left\{T_s, \frac{d}{m}\right\}\eta F + \tilde{O}(\eta^2\epsilon^2)\right).$$

Intuitively, the above result shows that if little function value improvement has been made, then the algorithm's iterates have not moved much, such that it remains approximately in a saddle region if it started out in a saddle region. Next, Lemma 3 formally introduces the coupling we have mentioned, setting the stage for the rest of our arguments. For notational convenience, in this section, unless otherwise specified, we will assume that the initial iterate $x_0$ is an $\epsilon$-saddle point.

**Lemma 3.** *Suppose $x_0$ is an $\epsilon$-approximate saddle point. Without loss of generality, suppose that the minimum eigendirection of $H := \nabla^2 f(x_0)$ is the $e_1$ direction (i.e. the first basis vector in $\mathbb{R}^d$), and let $\gamma$ to denote $-\lambda_{\min}(\nabla^2 f(x_0))$ (note $\gamma \geq \sqrt{\rho\epsilon}$). Consider the following coupling mechanism, where we run the zeroth-order gradient dynamics, starting with $x_0$, with two isotropic noise sequences, $Y_t$ and $Y_t^0$ respectively, where $(Y_t)_1 = -(Y_t)_1^0$, and $(Y_t)_j = (Y_t)_j^0$ for all other $j \neq 1$. Suppose that the sequence $\{Z_{t,i}\}_{t\geq T, i\in[m]}$ is the same for both sequences. Let $\{x_t\}$ denote the sequence with the $\{Y_t\}$ noise sequence, and let the $\{x_t^0\}$ denote the sequence with the $\{Y_t^0\}$ noise sequence, where $x_0^0 = x_0$, and*

$$x_{t+1}^0 = x_t^0 - \eta\left(\frac{\sum_{i=1}^m Z_{t,i}Z_{t,i}^\top \nabla f(x_t^0) + \frac{u}{2}Z_{t,i}Z_{t,i}^\top \tilde{H}_{t,i}^0 Z_{t,i}}{m}\right) + Y_t^0;$$

*and $\tilde{H}_{t,i}^0 := \frac{H_{t;i;+}^0 - H_{t;i;-}^0}{2}$, with $H_{t,i,+}^0 = \nabla^2 f(x_t^0 + \alpha_{t,i,+}^0 uZ_i^0)$ for some $\alpha_{t,i,+}^0 \in [0,1]$, and $H_{t,i,-}^0 = \nabla^2 f(x_t - \alpha_{t,i,-}^0 uZ_i^0)$ for some $\alpha_{t,i,-}^0 \in [0,1]$. Then, for any $t \geq 0$,*

$$\hat{x}_{t+1} := x_{t+1} - x_{t+1}^0$$

$$= \underbrace{\sum_{\tau=0}^t (I - \eta H)^{t-\tau}\eta\,\hat{g}_0(\tau)}_{W_{g_0}(t+1)} - \underbrace{\sum_{\tau=0}^t (I - \eta H)^{t-\tau}\eta(H_\tau - H)\hat{x}_\tau}_{W_H(t+1)}$$

$$- \underbrace{\sum_{\tau=0}^t (I - \eta H)^{t-\tau}\eta\,\hat{u}(\tau)}_{W_U(t+1)} + \underbrace{\sum_{\tau=0}^t (I - \eta H)^{t-\tau}\eta\,\hat{Y}_\tau}_{W_P(t+1)}$$

*where*

$$g_0(t) = \frac{1}{m}\sum_{i=1}^n (Z_{t,i}Z_{t,i}^\top - I)\nabla f(x_t);$$

$$g_0^0(t) = \frac{1}{m}\sum_{i=1}^n (Z_{t,i}(Z_{t,i})^\top - I)\nabla f(x_t^0);$$

$$\hat{g}_0(t) = g_0(t) - g_0^0(t); \quad u(t) = \frac{1}{m}\sum_{i=1}^n \frac{u}{2}Z_{t,i}Z_{t,i}\tilde{H}_{t,i}Z_{t,i};$$

$$u^0(t) = \frac{1}{m}\sum_{i=1}^n \frac{u}{2}Z_{t,i}Z_{t,i}\tilde{H}_{t,i}^0 Z_{t,i}; \quad \hat{u}(t) = u(t) - u^0(t);$$

$$\hat{Y}_t = Y_t - Y_t^0; \quad H_t = \int_0^1 \nabla^2 f(ax_t + (1-a)x_t^0)da:$$

Our goal is to show that the dominating term in the evolution of the difference dynamics comes from the $W_p$ term involving the additional perturbation. To this end, we need to bound the remaining terms, $W_{g_0}, W_H, W_u$. A key technical challenge is to find a precise concentration bound for the $W_{g_0}(t+1)$ term, where

$$W_{g_0}(t+1)$$
$$= \sum_{\tau=0}^{t} (I - \eta H)^{t-\tau} \left( \frac{\sum_{i=1}^{m} (Z_{\tau,i} Z_{\tau,i}^\top - I)(\nabla f(x_\tau) - \nabla f(x_\tau^\ell))}{m} \right).$$

For the simplicity of discussion, we assume for the time being that $m = 1$, and drop the $i$ index in the subscript of $Z_{\tau,i}$. Since $\mathbb{E}[Z_\tau Z_\tau^\top] = I$, heuristically, assuming that $Z_\tau Z_\tau^\top - I$ satisfies "nice" concentration properties, utilizing the independence of the $Z_\tau$'s across time and the fact that $(I - \eta H) \preceq (1 + \eta\gamma)I$, we would like to show that with high probability,

$$\begin{aligned}
&W_{g_0}(t) \\
&\lesssim \eta \sqrt{\sum_{\tau=0}^{t-1} (1+\eta\gamma)^{2(t-1-\tau)} \mathbb{E}\left[ \left\| (Z_\tau Z_\tau^\top - I)(\nabla f(x_\tau) - \nabla f(x_\tau^\ell)) \right\|^2 \mid \mathcal{F}_{\tau-1} \right]}
\end{aligned} \tag{7}$$

where $\mathcal{F}_{\tau-1}$ is a sigma-algebra containing all randomness up to and including iteration $\tau - 1$, such that $x_\tau$ and $x_\tau^\ell$ are both in $\mathcal{F}_{\tau-1}$, but $Z_\tau$ is not. Then, assuming that Eq. (7) holds, since

$$\mathbb{E}\left[ \left\| (Z_\tau Z_\tau^\top - I)(\nabla f(x_\tau) - \nabla f(x_\tau^\ell)) \right\|^2 \mid \mathcal{F}_{\tau-1} \right]$$
$$= O(d) \left\| \nabla f(x_\tau) - \nabla f(x_\tau^\ell) \right\|^2;$$

it follows that

$$\| W_{g_0}(t) \| \lesssim \eta \sqrt{O(d) \sum_{\tau=0}^{t-1} (1+\eta\gamma)^{2(t-1-\tau)} \| \nabla f(x_\tau) - \nabla f(x_\tau^\ell) \|^2}$$

With this bound on $\| W_{g_0}(t) \|$, we eventually prove in Proposition 5 in Appendix E.3 that our algorithm escapes any $\epsilon-$saddle point with constant probability and that the $O(d)$ term appearing in the square root term above will eventually lead to an $O(d)$ dependence in the sample complexity[2]. We note that the $O(d)$ dimension dependence matches that of the best-known existing upper bound for finding first-order stationary points in smooth nonconvex zeroth-order optimization (Nesterov & Spokoiny, 2017), and has been conjectured to be the best possible dimension dependence for general smooth nonconvex zeroth-order optimization (Balasubramanian & Ghadimi, 2022).

**Key technical challenge** The key challenge in the above argument is to show that an equation in the form of Eq. (7) could in fact hold. At first glance, that an inequality such as Eq. (7) should hold is rather

[2]For general $1 \leq m \leq d$, there will also be an $O(1/m)$ dependence in the sample complexity.

non-obvious — this is because while the variable $(Z_\tau Z_\tau^\top - I)(\nabla f(x_\tau) - \nabla f(x_\tau^\ell)) \mid \mathcal{F}_{\tau-1}$ is mean-zero, it is subExponential rather than subGaussian. In fact, even in the subGaussian case, given a sequence of random vectors $\boldsymbol{x}_0, \ldots, \boldsymbol{x}_{t-1}$, such that each $\mathbb{E}[\boldsymbol{x}_\tau \mid \mathcal{F}_{\tau-1}] = 0$, and that each $\boldsymbol{x}_\tau \mid \mathcal{F}_{\tau-1}$ is norm-subGaussian with parameter $\sigma_\tau \in \mathcal{F}_{\tau-1}$ (which is an appropriate generalization of subGaussianity for vectors, proposed in (Jin et al., 2019b)), proving a concentration inequality of the form $\mathbb{P}\left[ \left\| \sum_{\tau=0}^{t-1} \boldsymbol{x}_\tau \right\| \leq \tilde{O}\left( \sqrt{\sum_{\tau=0}^{t-1} \sigma_\tau^2} \right) \right]$ is a very delicate matter.

In our case, the analogue of $\boldsymbol{x}_\tau$ is $(I - \eta H)^{t-1-\tau}(Z_\tau Z_\tau^\top - I)(\nabla f(x_\tau) - \nabla f(x_\tau^\ell))$, while the analogue of $\sigma_\tau^2$ is $(1 + \eta\gamma)^{2(t-1-\tau)} \mathbb{E}\left[ \| (Z_\tau Z_\tau^\top - I)(\nabla f(x_\tau) - \nabla f(x_\tau^\ell)) \|^2 \mid \mathcal{F}_{\tau-1} \right]$. Existing techniques (cf. (Tropp et al., 2015; Jin et al., 2019b)) rely crucially on subGaussian properties that allow for each $\tau$ the moment-generating function $\mathbb{E}[e^{\theta Y} \mid \mathcal{F}_{\tau-1}]$ to be defined for any fixed (and non-random) $\theta > 0$, where $Y_\tau$ takes the form

$$Y_\tau = \begin{pmatrix} 0 & \boldsymbol{x}_\tau^\top \\ \boldsymbol{x}_\tau & 0 \end{pmatrix};$$

such that $\mathbb{E}[Y_\tau \mid \mathcal{F}_{\tau-1}] = 0$ (since $\mathbb{E}[\boldsymbol{x}_\tau \mid \mathcal{F}_{\tau-1}] = 0$), and the eigenvalues of $Y_\tau$ are $\pm \| \boldsymbol{x}_\tau \|$. In the case when $\boldsymbol{x}_\tau$ is merely subExponential, the Moment Generating Function (MGF), $\mathbb{E}[e^{\theta Y} \mid \mathcal{F}_{\tau-1}]$, will no longer be well-defined at any fixed (and non-random) $\theta > 0$. This poses a challenge in our setting, since $\boldsymbol{x}_\tau$ takes the form $(I - \eta H)^{t-1-\tau}(Z_\tau Z_\tau^\top - I)(\nabla f(x_\tau) - \nabla f(x_\tau^\ell))$, which is subExponential rather than subGaussian. While it may be possible to force $(I - \eta H)^{t-1-\tau}(Z_\tau Z_\tau^\top - I)(\nabla f(x_\tau) - \nabla f(x_\tau^\ell))$ to be sub-Gaussian, say by normalizing $Z_\tau$ to have norm $\sqrt{d}$ (note any bounded random vector is also subGaussian), such that $\left\| (Z_\tau Z_\tau^\top - I)g \right\|^2 \leq d^2 \|g\|^2$ for any vector $g \in \mathbb{R}^d$, a careful examination of the argument in Proposition 5 would show that this results in a $O(d^2)$ rather than $O(d)$ dependence in the sample complexity, incurring a heavy price on the overall sample complexity (extra factor of $d$) if $d$ is large.

**Our solution** To overcome the issue, we build on the following observation: with high probability, for any vector $g \in \mathbb{R}^d$, $Z_\tau^\top g$ is bounded within some log factor of $\|g\|$. On the event $\{ Z_\tau^\top g = \tilde{O}(\|g\|) \}$, the variable

$$(Z_\tau Z_\tau^\top - I)g = Z_\tau(Z_\tau^\top g) - g \approx Z_\tau \|g\| - g$$

behaves approximately like a subGaussian random vector since $Z_\tau \sim N(0, I_d)$. Based on this intuition, after some careful analysis, we can show that $(Z_\tau Z_\tau^\top - I)(\nabla f(x_\tau) - \nabla f(x_\tau^\ell)) \mid \mathcal{F}_{\tau-1}$ is subGaussian on the event that $Z_\tau^\top \nabla f(x_\tau)$ is bounded within some log factor of $\| \nabla f(x_\tau) \|$, which happens with high probability. This then allows us to show that on this event, the corresponding MGF is well-defined for all fixed $\theta > 0$, enabling us

to prove a concentration inequality of the form Eq. (7). This intuition is crystallized in the following proposition, which proves a more general bound than what we strictly need. For notational simplicity, we introduce the function $\text{lr}(x) := \log(x \log(x))$.

**Proposition 2.** Let $F_t, t \geq 1$ be a filtration. Let $(Z_t)_{t \geq 0}$ be a sequence of random vectors following the distribution $N(0, I)$ such that $Z_t \in F_t$ and is independent of $F_{t-1}$, and let $(v_t)_{t \geq 0}$ be a sequence of random vectors such that $v_t \in F_{t-1}$. For each $\geq 0$, let

$$W = \sum_{t=0}^{\infty-1} M_t(Z_t Z_t^\top - I) v_t,$$

where each $M_t$ is a deterministic matrix of appropriate dimension. Then, there exist some absolute constants $c_0, C > 0$ such that for any $\in Z^+$ and $\in (0, 1/e]$, the following statements hold:

1. For any $> 0$, with probability at least $1 - $, we have

$$\|W\| \leq c_0 \sqrt{\sum_{t=0}^{\infty-1} \|M_t\|_2^2 d(\text{lr}(C/))^2 \|v_t\|^2} + \frac{1}{}\log(Cd/).$$

2. For any $B > b > 0$, with probability at least $1 - $,

either $\sum_{t=0}^{\infty-1} \|M_t\|_2^2 d(\text{lr}(C/))^2 \|v_t\|^2 \geq B$; or

$$\|W\| \leq \sqrt{\max\left(\sum_{t=0}^{\infty-1} \|M_t\|_2^2 d(\text{lr}(C/))^2 \|v_t\|^2, b\right)}$$
$$\cdot c_0 \sqrt{(\log(C\sqrt{d}/) + \log(\log(B/b) + 1))}$$

Moreover, as is clear from the bounds above, we may pick $C \geq 1$ such that $\log \frac{C}{} \geq 1, \forall \in (0, \frac{1}{e}]$.

With this result, along with a series of other technical results in Appendix E.3, we can show that the algorithm makes a function decrease of $\delta F$ with $\Omega(1)$ probability near an $$ saddle point (Proposition 5 in Appendix E.3). Armed with Proposition 5, as well as Proposition 1, the main result in Theorem 1 then follows. The complete detailed analysis can be found in Appendix E (escaping saddle point) and Appendix F (main result).

## 4. Simulations

We test the performance of our proposed algorithm with two-point estimators (ZOPGD-2pt) against existing zeroth-order benchmarks using the octopus function (proposed in (Du et al., 2017))[3]. It is known that the octopus function defined on $R^d$, which chains $d$ saddle points sequentially,

---

[3] Our code can be found at https://github.com/rafflesintown/escape-saddle-points-2pt

takes exponential (in $d$) time for exact gradient descent to escape; it has thus emerged as a popular benchmark to evaluate algorithms that seek to escape saddle points. In our experiments, we compare the performance of our two-point estimator algorithm (ZOPGD-2pt) with PAGD (Algorithm 1 in (Vlatakis-Gkaragkounis et al., 2019)) and ZO-GD-NCF (see (Zhang et al., 2022)), which are the only two existing zeroth-order algorithms that have (a) $\tilde{O}(d^{2})$ sample complexity for escaping saddle points (with the latter algorithm yielding the tightest bounds), and (b) performed the best empirically on escaping saddle points (see the simulation results in (Zhang et al., 2022)). Both PAGD and ZO-GD-NCF have to use $2d$ function evaluations per iteration to estimate the gradient while our algorithm only needs to use $2$ function evaluations. We plot the function value against the number of function evaluations.

We tested the algorithms for $d = 10$ and $d = 30$. To account for the stochasticity in the algorithms, for each algorithm, we computed the average and standard deviation over 30 trials, and plotted the mean trajectory with an additional band that represents $1.5$ times the standard deviation. For our algorithm's hyperparameters, we picked

$$ = \frac{1}{4dL}; u = 10^{-2}; r = 0.05; m = 1.$$

Note $m = 1$ corresponds to using a two-point estimator. For PAGD, we used the hyperparameters listed in their paper, and for ZO-GD-NCF, we used the code from their Neurips submission. For initialization, we chose a random $x_0$ near the saddle point at the origin, drawn from $N(0, 10^{-3} I_{d \times d})$

As we see in Fig. 2, our algorithm reaches the global minimum of the octopus function in significantly fewer function evaluations than PAGD and ZO-GD-NCF (approximately 2.5 times faster than ZO-GD-NCF, and approximately 3 times faster than PAGD), despite our algorithm only using $2$ function evaluations per iteration compared to $2d$ function evaluations per iteration for both PAGD and ZO-GD-NCF. This suggests that in addition to our theoretical convergence guarantees, there can also be empirical benefits to using two-point estimators versus existing $2d$-point estimators in the zeroth-order escaping saddle point literature.

## 5. Conclusion

In this paper, we proved that using two function evaluations per iteration suffices to escape saddle points and reach approximate second order stationary points efficiently in zeroth-order optimization. Along the way, we also gave the first analysis of high-probability function change using two (or more)-point zeroth-order gradient estimators, as well as a novel concentration bound for sums of subExponential (but not subGaussian) vectors which are each the products of Gaussian vectors. These technical contributions may be

Figure 1: Performance on toy octopus function, with $d = 30$

of independent interest to researchers working in zeroth-order optimization as well as general stochastic optimization. Finally, we provided numerical evidence supporting the theoretical convergence results.

## 6. Acknowledgements

## References

Adolphs, L., Daneshmand, H., Lucchi, A., and Hofmann, T. Local saddle point optimization: A curvature exploitation approach. In *The 22nd International Conference on Artificial Intelligence and Statistics*, pp. 486–495. PMLR, 2019.

Antonakopoulos, K., Mertikopoulos, P., Piliouras, G., and Wang, X. Adagrad avoids saddle points. In *International Conference on Machine Learning*, pp. 731–771. PMLR, 2022.

Avdiukhin, D. and Yaroslavtsev, G. Escaping saddle points with compressed sgd. *Advances in Neural Information Processing Systems*, 34:10273–10284, 2021.

Avdiukhin, D., Jin, C., and Yaroslavtsev, G. Escaping saddle points with inequality constraints via noisy sticky projected gradient descent. In *11th Annual Workshop on Optimization for Machine Learning*, 2019.

Bai, Q., Agarwal, M., and Aggarwal, V. Escaping saddle points for zeroth-order non-convex optimization using estimated gradient descent. In *2020 54th Annual Conference on Information Sciences and Systems (CISS)*, pp. 1–6. IEEE, 2020.

Balasubramanian, K. and Ghadimi, S. Zeroth-order nonconvex stochastic optimization: Handling constraints, high dimensionality, and saddle points. *Foundations of Computational Mathematics*, 22(1):35–76, 2022.

Chen, P.-Y., Zhang, H., Sharma, Y., Yi, J., and Hsieh, C.-J. Zoo: Zeroth order optimization based black-box attacks to deep neural networks without training substitute models. In *Proceedings of the 10th ACM workshop on artificial intelligence and security*, pp. 15–26, 2017.

Chen, Y., Bernstein, A., Devraj, A., and Meyn, S. Model-free primal-dual methods for network optimization with application to real-time optimal power flow. In *2020 American Control Conference (ACC)*, pp. 3140–3147. IEEE, 2020.

Criscitiello, C. and Boumal, N. Efficiently escaping saddle points on manifolds. *Advances in Neural Information Processing Systems*, 32, 2019.

Daneshmand, H., Kohler, J., Lucchi, A., and Hofmann, T. Escaping saddles with stochastic gradients. In *International Conference on Machine Learning*, pp. 1155–1164. PMLR, 2018.

Du, S. S., Jin, C., Lee, J. D., Jordan, M. I., Singh, A., and Poczos, B. Gradient descent can take exponential time to escape saddle points. In *Advances in neural information processing systems*, pp. 1067–1077, 2017.

9

Duchi, J. C., Jordan, M. I., Wainwright, M. J., and Wibisono, A. Optimal rates for zero-order convex optimization: The power of two function evaluations. *IEEE Transactions on Information Theory*, 61(5):2788–2806, 2015.

Fang, C., Lin, Z., and Zhang, T. Sharp analysis for nonconvex sgd escaping from saddle points. *Conference on Learning Theory*, pp. 1192–1234. PMLR, 2019.

Flaxman, A. D., Kalai, A. T., and McMahan, H. B. Online convex optimization in the bandit setting: gradient descent without a gradient. In *Proceedings of the Sixteenth Annual ACM-SIAM Symposium on Discrete Algorithms*, pp. 385–394, 2005.

Ge, R., Huang, F., Jin, C., and Yuan, Y. Escaping from saddle points—online stochastic gradient for tensor decomposition. In *Conference on Learning Theory*, pp. 797–842, 2015.

Ge, R., Li, Z., Wang, W., and Wang, X. Stabilized svrg: Simple variance reduction for nonconvex optimization. In *Conference on learning theory*, pp. 1394–1448. PMLR, 2019.

Han, A. and Gao, J. Escape saddle points faster on manifolds via perturbed riemannian stochastic recursive gradient. *arXiv preprint arXiv:2010.12191*, 2020.

Huang, M. Escaping saddle points for nonsmooth weakly convex functions via perturbed proximal algorithms. *arXiv preprint arXiv:2102.02837*, 2021.

Huang, M., Ji, K., Ma, S., and Lai, L. Ef ciently escaping saddle points in bilevel optimization. *arXiv preprint arXiv:2202.03684*, 2022.

Jin, C., Ge, R., Netrapalli, P., Kakade, S. M., and Jordan, M. I. How to escape saddle points ef ciently. *International Conference on Machine Learning*, pp. 1724–1732. PMLR, 2017.

Jin, C., Liu, L. T., Ge, R., and Jordan, M. I. On the local minima of the empirical risk. *Advances in neural information processing systems*, 31, 2018a.

Jin, C., Netrapalli, P., and Jordan, M. I. Accelerated gradient descent escapes saddle points faster than gradient descent. In *Conference On Learning Theory*, pp. 1042–1085. PMLR, 2018b.

Jin, C., Netrapalli, P., Ge, R., Kakade, S. M., and Jordan, M. I. On nonconvex optimization for machine learning: Gradients, stochasticity, and saddle points. *arXiv preprint arXiv:1902.04811*, 2019a.

Jin, C., Netrapalli, P., Ge, R., Kakade, S. M., and Jordan, M. I. A short note on concentration inequalities for random vectors with subgaussian norm. *arXiv preprint arXiv:1902.03736*, 2019b.

Larson, J., Menickelly, M., and Wild, S. M. Derivative-free optimization methods. *Acta Numerica*, 28:287–404, 2019.

Lee, J. D., Panageas, I., Piliouras, G., Simchowitz, M., Jordan, M. I., and Recht, B. First-order methods almost always avoid strict saddle points. *Mathematical programming*, 176(1):311–337, 2019.

Li, J., Balasubramanian, K., and Ma, S. Stochastic zeroth-order riemannian derivative estimation and optimization. *Mathematics of Operations Research*, 2022.

Li, Y., Tang, Y., Zhang, R., and Li, N. Distributed reinforcement learning for decentralized linear quadratic control: A derivative-free policy optimization approach. *IEEE Transactions on Automatic Control*, 2021.

Li, Z. Ssrgd: Simple stochastic recursive gradient descent for escaping saddle points. *Advances in Neural Information Processing Systems*, 32, 2019.

Sang, G., Tong, Q., Zhu, C., and Bi, J. Escaping saddle points with stochastically controlled stochastic gradient methods. *arXiv preprint arXiv:2103.04413*, 2021.

Lucchi, A., Orvieto, A., and Solomou, A. On the second-order convergence properties of random search methods. *Advances in Neural Information Processing Systems*, 34: 25633–25645, 2021.

Malik, D., Pananjady, A., Bhatia, K., Khamaru, K., Bartlett, P., and Wainwright, M. Derivative-free methods for policy optimization: Guarantees for linear quadratic systems. In The 22nd international conference on arti cial intelligence and statistics, pp. 2916–2925. PMLR, 2019.

Mokhtari, A., Ozdaglar, A., and Jadbabaie, A. Escaping saddle points in constrained optimization. *Advances in Neural Information Processing Systems*, 31, 2018.

Nesterov, Y. and Spokoiny, V. Random gradient-free minimization of convex functions. *Foundations of Computational Mathematics*, 17(2):527–566, 2017.

Panageas, I., Piliouras, G., and Wang, X. First-order methods almost always avoid saddle points: The case of vanishing step-sizes. *Advances in Neural Information Processing Systems*, 32, 2019.

Reddi, S., Zaheer, M., Sra, S., Poczos, B., Bach, F., Salakhutdinov, R., and Smola, A. A generic approach for escaping saddle points. In *International conference on arti cial intelligence and statistics*, pp. 1233–1242. PMLR, 2018.

Roy, A., Balasubramanian, K., Ghadimi, S., and Mohapatra, P. Escaping saddle-point faster under interpolation-like conditions. *Advances in Neural Information Processing Systems*, 33:12414–12425, 2020.

Shamir, O. An optimal algorithm for bandit and zero-order convex optimization with two-point feedback. *The Journal of Machine Learning Research*, 18(1):1703–1713, 2017.

Staib, M., Reddi, S., Kale, S., Kumar, S., and Sra, S. Escaping saddle points with adaptive gradient methods. In *International Conference on Machine Learning*, pp. 5956–5965. PMLR, 2019.

Sun, J., Qu, Q., and Wright, J. When are nonconvex problems not scary? *arXiv preprint arXiv:1510.06096*, 2015.

Sun, T., Li, D., Quan, Z., Jiang, H., Li, S., and Dou, Y. Heavy-ball algorithms always escape saddle points. *arXiv preprint arXiv:1907.09697*, 2019a.

Sun, Y., Flammarion, N., and Fazel, M. Escaping from saddle points on riemannian manifolds. *Advances in Neural Information Processing Systems*, 32, 2019b.

Tang, H., Lian, X., Qiu, S., Yuan, L., Zhang, C., Zhang, T., and Liu, J. Deepsqueeze: Parallel stochastic gradient descent with double-pass error-compensated compression. *arXiv preprint arXiv:1907.07346*, 2019.

Tang, Y., Ren, Z., and Li, N. Zeroth-order feedback optimization for cooperative multi-agent systems. In *2020 59th IEEE Conference on Decision and Control (CDC)*, pp. 3649–3656. IEEE, 2020a.

Tang, Y., Zhang, J., and Li, N. Distributed zero-order algorithms for nonconvex multiagent optimization. *IEEE Transactions on Control of Network Systems*, 8(1):269–281, 2020b.

Tropp, J. A. et al. An introduction to matrix concentration inequalities. *Foundations and Trends® in Machine Learning*, 8(1-2):1–230, 2015.

Vershynin, R. *High-dimensional probability: An introduction with applications in data science*, volume 47. Cambridge university press, 2018.

Vladimirova, M., Girard, S., Nguyen, H., and Arbel, J. Sub-weibull distributions: Generalizing sub-gaussian and sub-exponential properties to heavier tailed distributions. *Stat*, 9(1):e318, 2020.

Vlaski, S. and Sayed, A. H. Distributed learning in non-convex environments—part ii: Polynomial escape from saddle-points. *IEEE Transactions on Signal Processing*, 69:1257–1270, 2021a.

Vlaski, S. and Sayed, A. H. Second-order guarantees of stochastic gradient descent in non-convex optimization. *IEEE Transactions on Automatic Control*, 2021b.

Vlatakis-Gkaragkounis, E.-V., Flokas, L., and Piliouras, G. Efficiently avoiding saddle points with zero order methods: No gradients required. *Advances in Neural Information Processing Systems*, 32, 2019.

Wang, J.-K., Lin, C.-H., and Abernethy, J. Escaping saddle points faster with stochastic momentum. *arXiv preprint arXiv:2106.02985*, 2021.

Xu, Y., Jin, R., and Yang, T. First-order stochastic algorithms for escaping from saddle points in almost linear time. *Advances in neural information processing systems*, 31, 2018.

Zhang, C. and Li, T. Escape saddle points by a simple gradient-descent based algorithm. *Advances in Neural Information Processing Systems*, 34:8545–8556, 2021.

Zhang, H., Xiong, H., and Gu, B. Zeroth-order negative curvature finding: Escaping saddle points without gradients. *arXiv preprint arXiv:2210.01496*, 2022.

Ziyin, L., Li, B., Simon, J. B., and Ueda, M. Sgd with a constant large learning rate can converge to local maxima. *arXiv preprint arXiv:2107.11774*, 2021.

# A. Related Work

**Two-point methods in zeroth-order optimization.** Two-point (or in general $2m$-point, where $1 \leq m < d$ with d being the problem dimension) estimators, which approximate the gradient using two ($2m$) function evaluations per iteration, have been widely studied by researchers in the zeroth-order optimization literature, in convex (Nesterov & Spokoiny, 2017; Duchi et al., 2015; Shamir, 2017), nonconvex (Nesterov & Spokoiny, 2017), online (Shamir, 2017), as well as distributed settings (Tang et al., 2019). A key reason for doing so is that for applications of zeroth-order optimization arising in robotics (Li et al., 2022), wind farms (Tang et al., 2020a), power systems (Chen et al., 2020), online (time-varying) optimization (Shamir, 2017), learning-based control (Malik et al., 2019; Li et al., 2021), and improving adversarial robustness to black-box attacks in deep neural networks (Chen et al., 2017), it may be costly or impractical to wait for $\Theta(d)$ (where d denotes the problem dimension) function evaluations per iteration to make a step. This is especially true for high-dimensional and/or time-varying problems. Indeed, for high-dimensional problems, two-point estimators can make swift progress even in the initial stage compared to $2d$-point estimator, and can reach a higher-quality solution if computation is limited (Tang et al., 2020b; Chen et al., 2017). For instance, consider the work in (Chen et al., 2017), which studies the use of zeroth-order estimators to perform black-box attacks on deep neural networks, in order to identify (and then defend against) adversarial images that may lead to misclassiﬁcation. In the paper, the authors employed two-point zeroth-order estimators, due to the high computational cost of using $2d$ function evaluations per iteration for hundreds of iterations (here $d$ is the dimension of an image, which in this case is over 20000). The authors showed empirically that their two-point estimators worked well; however there were no accompanying theoretical results.

For online or time-varying environments, two-points estimators also often preferable. Since zeroth-order methods are often used in physical systems whose environment drifts or changes over time, this leads naturally to time-varying or online optimization. For these problems, $2d$-point estimators will not produce a good estimation because the underlying function can drift to a very different problem while waiting for the $2d$ function evaluations. Indeed, the fewer function evaluations an optimization procedure needs, the faster it can catch up with the time-varying environment. In fact, for online optimization, it has been shown that two points estimator is optimal for convex Lipschitz functions (Shamir, 2017). Thus, two-point estimators are a natural ﬁt for time-varying online optimization problems.

**Saddle point escape with access to deterministic gradient.** While standard gradient descent can escape saddle points asymptotically (Lee et al., 2019; Panageas et al., 2019), it is known that standard gradient descent may take exponential time to escape saddle points (Du et al., 2017). Hence, when access to deterministic gradient is available, research has centered on escaping saddle points with adding perturbation (Jin et al., 2017), momentum/acceleration based methods (Jin et al., 2018b; Sun et al., 2019a; Staib et al., 2019), or gradient-based robust Hessian power/curvature exploitation methods (Zhang & Li, 2021; Adolphs et al., 2019). In addition, there has also been work on escaping saddle points devoted to speciﬁc optimization settings, such as constrained optimization (Mokhtari et al., 2018; Avdiukhin et al., 2019), optimization of weakly convex functions (Huang, 2021), bilevel optimization (Huang et al., 2022), as well as on general manifolds (Sun et al., 2019b; Criscitiello & Boumal, 2019; Han & Gao, 2020).

**Saddle point escape in stochastic gradient descent (SGD).** In practice, only stochastic gradient estimators are available in many problems. While SGD may converge to local maxima in worst-case scenarios (Ziyin et al., 2021), under assumptions such as bounded variance or subGaussian noise, there have been many works that have studied the problem of saddle point escape in SGD (Ge et al., 2015; Daneshmand et al., 2018; Xu et al., 2018; Jin et al., 2019a; Vlaski & Sayed, 2021b). The best existing rate (without considering momentum/variance reduction techniques) appears to belong to that of (Fang et al., 2019), which converges to second order stationary points using $O(\epsilon^{-3.5})$ stochastic gradients. While zeroth-order gradient estimators may also be viewed as stochastic gradients, they typically do not satisfy the bounded/subGaussian noise assumptions that are assumed in these works, making a direct comparison inappropriate. Escaping saddle point via momentum methods in SGD has also been studied (Wang et al., 2021; Antonakopoulos et al., 2022); while we do not consider incorporating momentum in our works, this may be interesting future work. A number of papers has also considered the specialized setting of escaping saddle points in nonconvex ﬁnite-sum optimization (Reddi et al., 2018; Liang et al., 2021), with many considering the case where variance-reduction is used (Ge et al., 2019; Li, 2019). While the ﬁnite-sum problem is quite different from our problem, the variance reduction approach considered in these works may be a relevant future direction. The saddle point escape problem has also been studied in other speciﬁc settings such as compressed optimization (Avdiukhin & Yaroslavtsev, 2021), distributed optimization (Vlaski & Sayed, 2021a), or in the overparameterization case (Roy et al., 2020).

**Saddle point escape with zeroth-order information.** The problem of escaping saddle points in zeroth-order optimization

has been studied less often, and we have already listed all known works comparable to our work in the introduction (Bai et al., 2020; Vlatakis-Gkaragkounis et al., 2019; Balasubramanian & Ghadimi, 2022); a more detailed comparison of these works with our results has been provided in the discussion following the statement of our main result Theorem 1. We would like to mention that (Roy et al., 2020) also includes a convergence result $\tilde{O}(d^{1.5}_{4.5})$ for the case with noisy function evaluations, which is incomparable to our existing work which focuses on the case with exact function evaluation. In addition, (Roy et al., 2020) also makes a subGaussian assumption on the estimator noise, which zeroth-order estimators in our paper do not satisfy. Nonetheless, considering the extension to noisy function evaluations will make for important future work.

Zeroth-order optimization. Our work rests on a line of research in zeroth-order optimization which focuses on constructing gradient estimators using zeroth-order function values (Flaxman et al., 2005; Duchi et al., 2015; Nesterov & Spokoiny, 2017; Shamir, 2017; Larson et al., 2019). As we have discussed, for smooth nonconvex functions, it is known that two-point zeroth-order estimators suffice to find first-order stationary points using $\tilde{O}(d\epsilon^{-2})$ function evaluations (Nesterov & Spokoiny, 2017). Our work studies the more complicated problem of reaching second order stationary points, attaining a rate of $\tilde{O}(d\epsilon^{-2.5})$.

## B. Proof Roadmap

We begin by introducing several key concentration inequalities in Appendix C which we will frequently use in our proofs. We then describe in detail (and prove) the sequence of results that lead up to Proposition 4 in Appendix D, showing that there cannot be too many iterations with large gradients. Next, we describe the saddle point argument in detail, and prove Proposition 5 in Appendix E.3. Finally, we combine these results and prove our main result Theorem 2 (whose informal version is Theorem 1) in Appendix F

Throughout our proofs, absolute constants, as denoted by (e.g. $c, C$), may change from line to line. However, within the same proof, for clarity, we try to index different constants differently. We assume $\epsilon \le 1$ and $m \ge d$.

Notations. We shall denote the conditional expectation and conditional probability by $E_F[\cdot] = E[\cdot \mid F]$ and $P_F(\cdot) = P(\cdot \mid F)$ where $F$ is a sigma-algebra.

## C. Concentration inequalities

This section serves to introduce several probabilistic results which will be useful for our main proofs in subsequent sections. We first introduce subGaussian, subExponential and norm-subGaussian random vectors in Appendix C.1. Next, in Appendix C.2, we provide concentration bounds for norm-subGaussian and subExponential random vectors. We then prove a novel concentration inequality involving products of subGaussian random vectors in Appendix C.3. We conclude by stating some concentration bounds for Appendix C.4 random variables.

C.1. subGaussian, subExponential and norm-subGaussian random vectors

We first define subGaussian and subExponential random vectors. A detailed reference for these concepts can be found in (Vershynin, 2018).

Definition 5 (subGaussian and subExponential random vectors) A random vector $x \in R^d$ is $\sigma$-subGaussian (SG($\sigma$)), if there exists $\sigma > 0$ such that for any unit vector $g \in S^{d-1}$,

$$E[\exp(\lambda \langle g, x - E[x] \rangle)] \le \exp(\lambda^2 \sigma^2 / 2) \quad \forall \lambda \in R.$$

Meanwhile, a random vector $x \in R^d$ is $\sigma$-subExponential (SE($\sigma$)), if there exists $\sigma > 0$ such that for any unit vector $g \in S^{d-1}$,

$$E[\exp(\lambda \langle g, x - E[x] \rangle)] \le \exp(\lambda^2 \sigma^2 / 2) \quad \forall |\lambda| \le \frac{1}{\sigma}$$

An alternative concentration property for random vectors revolving around its norm, known as norm-subGaussianity (Jin et al., 2019b), is also relevant.

Definition 6 (norm-subGaussian random vectors) A random vector $x \in \mathbb{R}^d$ is $\sigma$-norm-subGaussian (nSG($\sigma$)), if there exists $\sigma > 0$ such that

$$P(\|x - \mathbb{E}x\| \geq s) \leq 2e^{-\frac{s^2}{2\sigma^2}} \quad \forall s \geq 0.$$

We recall the following result which provides several examples of nSG random vectors. In particular, it tells us a random vector $x \in \mathbb{R}^d$ that is ($\sigma = \sqrt{d}$) subGaussian is also $\sigma$-subGaussian.

Lemma 4 (Lemma 1 in (Jin et al., 2019b)) There exists absolute constant $c$ such that the following random vectors are all nSG($c\sigma$).

1. A bounded random vector $x \in \mathbb{R}^d$ so that $\|x\| \leq \sigma$.

2. A random vector $x \in \mathbb{R}^d$, where $x = \xi e_1$ and the random variable $\xi \in \mathbb{R}$ is $\sigma$-subGaussian.

3. A random vector $x \in \mathbb{R}^d$ that is ($\sigma = \sqrt{d}$) subGaussian

In addition, if $x \in \mathbb{R}^d$ is zero-mean nSG($\sigma$), its component along a single direction is also subGaussian.

Lemma 5. Suppose $x \in \mathbb{R}^d$ is zero-mean nSG($\sigma$). Then, for any fixed vector $v \in \mathbb{R}^d$, $\langle v, x \rangle$ is zero-mean $\|v\|\sigma$-subGaussian.

Proof. Without loss of generality, we assume that $v \in S^{d-1}$ is a unit vector. That $\langle v, x \rangle$ is zero-mean follows directly from $x$ being zero-mean and $v$ being fixed. Meanwhile, since $|\langle v, x \rangle| \leq \|v\| \|x\| = \|x\|$, for any $s \geq 0$, it follows that

$$P(|\langle v, x \rangle| \geq s) \leq P(\|x\| \geq s) \leq 2e^{-\frac{s^2}{2\sigma^2}};$$

where the last inequality follows from the fact that $x$ is zero-mean and also nSG($\sigma$). Hence $\langle v, x \rangle$ is zero-mean SG($\sigma$), as desired. $\square$

C.2. Concentration bounds for norm-subGaussian and subExponential random vectors

We begin by giving some concentration bounds for norm-subGaussian random vectors. To do so, we introduce the following condition.

Condition 1. Consider random vectors $x_1, \ldots, x_n \in \mathbb{R}^d$, and corresponding filtrations $\mathcal{F}_i$ generated by $(x_1, \ldots, x_i)$. We assume $x_i | \mathcal{F}_{i-1}$ is zero-mean, nSG($\sigma_i$), with $\sigma_i \in \mathcal{F}_{i-1}$, i.e,

$$\mathbb{E}[x_i | \mathcal{F}_{i-1}] = 0;$$

and

$$P(\|x_i\| \geq s | \mathcal{F}_{i-1}) \leq 2e^{-\frac{s^2}{2\sigma_i^2}} \quad \forall s \geq 0;$$

where $\sigma_i$ is a measurable function of $(x_1, \ldots, x_{i-1})$ for each $i$.

For norm subGaussian random vectors satisfying Condition 1, we first have the following bound.

Lemma 6. Suppose $(x_1, \ldots, x_n) \in \mathbb{R}^d$ satisfy Condition 1, i.e. each $x_i | \mathcal{F}_{i-1}$ is mean-zero, nSG($\sigma_i$) with $\sigma_i \in \mathcal{F}_{i-1}$. Let $\{u_i\}$ denote a sequence of random vectors such that $u_i \in \mathcal{F}_{i-1}$ for every $i \in [n]$. Then, there exists an absolute constant $c$ such that for any $\delta \in (0,1)$ and $\lambda > 0$, with probability at least $1 - \delta$,

$$\sum_{i=1}^n \langle u_i, x_i \rangle \leq c\lambda \sum_{i=1}^n \|u_i\|^2 \sigma_i^2 + \frac{1}{\lambda} \log(1/\delta).$$

Proof. We note that if $x_i$ is mean-zero and nSG($\sigma$), then by Lemma 5 $\langle u_i, x_i \rangle | \mathcal{F}_{i-1}$ is zero-mean and $\|u_i\|\sigma_i$-subGaussian. The rest of the proof follows from the proof of Lemma 39 in (Jin et al., 2019a) (key idea is exponentiate

and then apply Markov's inequality). For completeness, we restate the proof here. Observe that, for any $i$, since $\langle u_i, x_i \rangle$ is $\kappa \|u_i\| \sigma_i$-subGaussian, for any $\lambda > 0$, we have that

$$\mathbb{E}\left[\exp(\lambda \langle u_i, x_i \rangle) \mid F_{i-1}\right] \leq \exp(\lambda^2 \kappa \|u_i\|^2 \sigma_i^2 / 2)$$

For any $\lambda > 0$ and $s \geq 0$, observe that

$$P\left(\sum_{i=1}^n \langle u_i, x_i \rangle - \lambda \kappa \|u_i\|^2 \sigma_i^2 / 2 \geq s\right)$$

$$= P\left(\exp\left(\lambda \sum_{i=1}^n \langle u_i, x_i \rangle - \lambda^2 \kappa \|u_i\|^2 \sigma_i^2 / 2\right) \geq \exp(\lambda s)\right)$$

$$\leq \mathbb{E}\left[\exp\left(\lambda \sum_{i=1}^n \langle u_i, x_i \rangle - \lambda^2 \kappa \|u_i\|^2 \sigma_i^2 / 2\right)\right] \exp(-\lambda s)$$

$$= \mathbb{E}\left[\mathbb{E}\left[\exp\left(\lambda \sum_{i=1}^n \langle u_i, x_i \rangle - \lambda^2 \kappa \|u_i\|^2 \sigma_i^2 / 2\right) \mid F_{n-1}\right]\right] \exp(-\lambda s)$$

$$= \mathbb{E}\left[\exp\left(\lambda \sum_{i=1}^{n-1} \langle u_i, x_i \rangle - \lambda^2 \kappa \|u_i\|^2 \sigma_i^2 / 2\right) \mathbb{E}\left[\exp\left(\lambda \langle u_n, x_n \rangle - \lambda^2 \kappa \|u_n\|^2 \sigma_n^2 / 2\right) \mid F_{n-1}\right]\right] \exp(-\lambda s)$$

$$\overset{(i)}{\leq} \mathbb{E}\left[\exp\left(\lambda \sum_{i=1}^{n-1} \langle u_i, x_i \rangle - \lambda^2 \kappa \|u_i\|^2 \sigma_i^2 / 2\right)\right] \exp(-\lambda s) \leq \cdots \leq \exp(-\lambda s)$$

Above, (i) follows from the fact that $\langle u_i, x_i \rangle \mid F_{i-1}$ is zero-mean and $\kappa \|u_i\| \sigma_i$-subGaussian for each $i \in [n]$. The final result then follows by picking $\lambda = \frac{1}{2}$ and $s = \log(1/\delta)$. $\square$

Assuming Condition 1, the following concentration result also holds for a sequence of nSG random vectors.

Lemma 7 (Lemma 6, Corollary 7 and Corollary 8 in (Jin et al., 2019b) combined). Suppose $(x_1, \ldots, x_n) \in \mathbb{R}^d$ satisfy Condition 1. Then, there exists an absolute constant $c$ such that for any fixed $\delta \in (0, 1)$, $\lambda > 0$, with probability at least $1 - \delta$,

$$\left\| \sum_{i=1}^n x_i \right\| \leq c \lambda \sum_{i=1}^n \sigma_i^2 + \frac{1}{\lambda} \log(2d/\delta).$$

Moreover, there are two corollaries.

1. (Corollary 7 in (Jin et al., 2019b)) When $\{\sigma_i\}$ is deterministic, there exists an absolute constant $c$ such that for any fixed $\delta \in (0, 1)$, with probability at least $1 - \delta$.

$$\left\| \sum_{i=1}^n x_i \right\| \leq c \sqrt{\log(2d/\delta) \sum_{i=1}^n \sigma_i^2}$$

2. (Corollary 8 in (Jin et al., 2019b)) Suppose that the $\{\sigma_i\}$ sequence is random. Then, there exists an absolute constant $c$ such that for any fixed $\delta \in (0, 1)$ and $B > b > 0$, with probability at least $1 - \delta$:

$$\text{either} \sum_{i=1}^n \sigma_i^2 \geq B \quad \text{or} \quad \left\| \sum_{i=1}^n x_i \right\| \leq c \sqrt{\max\left(\sum_{i=1}^n \sigma_i^2, b\right) \left(\log(2d/\delta) + \log(\log(B/b))\right)}$$

We state here a Bernstein-type concentration inequality for sub-exponential random variables, which we also need.

15

Lemma 8 (Bernstein concentration inequality) Consider a sequence of independently distributed subexponential variables $x_1, \ldots, x_n \in \mathbb{R}$, with mean $E[x_i] \leq c^0$ for some $c^0 > 0$ and each $i \in [n]$. Then, there exists an absolute constant $C > 0$, such that for any $\delta \in (0, 1)$, with probability at least $1 - \delta$,

$$\sum_{i=1}^{n} x_i \leq C \cdot (n + \log(1/\delta)). \tag{8}$$

Proof. The result of Eq. (8) follows by applying Bernstein's inequality to $\sum_{i=1}^{n} x_i - E[x_i]$ (so each summand is mean-zero). Per Bernstein's inequality, (cf. Theorem 2.8.1 in (Vershynin, 2018)), there exists an absolute constant $c$ such that for any $s \geq 0$,

$$P\left(\sum_{i=1}^{n}(x_i - E[x_i]) \geq s\right) \leq \exp\left(-c \min\left\{\frac{s^2}{n\sigma^2}; \frac{s}{\sigma}\right\}\right).$$

Pick $s = \sigma\left(n + \frac{\log(1/\delta)}{c}\right)$. Then,

$$\min\left\{\frac{s^2}{n\sigma^2}; \frac{s}{\sigma}\right\} = \min\left\{n + 2\frac{\log(1/\delta)}{c} + \frac{(\log(1/\delta))^2}{c^2 n}; n + \frac{\log(1/\delta)}{c}\right\} = n + \frac{\log(1/\delta)}{c}.$$

Continuing, we have that

$$P\left(\sum_{i=1}^{n}(x_i - E[x_i]) \geq s\right) \leq \exp\left(-c \min\left\{\frac{s^2}{n\sigma^2}; \frac{s}{\sigma}\right\}\right) \leq \exp\left(-c\left(n + \frac{\log(1/\delta)}{c}\right)\right).$$

Thus, it follows that with probability at least $1 - \delta$,

$$\sum_{i=1}^{n}(x_i - E[x_i]) \leq \sigma\left(n + \frac{\log(1/\delta)}{c}\right) \implies \sum_{i=1}^{n} x_i \leq \sigma\left(n + \frac{\log(1/\delta)}{c}\right) + nc^0,$$

where implication holds since by assumption $E[x_i] \leq c^0$ for some $c^0 > 0$. Then, by setting $C = \max\{1 + c^0; 1/c\}$, the desired result follows. □

C.3. A novel concentration inequality for the zeroth-order setting

In the zeroth-order setting, we will frequently have to bound the norms of terms of the form

$$W = \sum_{t=0}^{\tau-1} M_t(Z_t Z_t^\top - I)v_t, \tag{9}$$

where $M_t$ is a known and fixed quantity, while $Z_t$ is random, and $v_t$ depends on $x_0$ and the history of previous $\{Z_j g_j\}_{j=0}^{t-1}$'s, and is hence $\mathcal{F}_{t-1}$-measurable. For our purposes, it suffices to consider $Z_t \sim N(0, I)$.

To see why such a bound will be useful, as mentioned in the main text and as we will see again later in the full proofs, in the analysis of escaping saddle points, we need to bound a term of the form

$$W_{g_0}(\tau) = \sum_{t=0}^{\tau-1}(I - \eta H)^{\tau-1-t}(Z_t Z_t^\top - I)(\nabla f(x_t) - \nabla f(x_t^0)),$$

where $H = \nabla^2 f(x_0)$ (assuming that $x_0$ is an $\epsilon$-saddle point), and $x_t$ and $x_t^0$ are two coupled sequences. Comparing with Eq. (9), we see that for the equation above, we can define $M_t = (I - \eta H)^{\tau-1-t}$ (a fixed and known quantity) and $v_t = \nabla f(x_t) - \nabla f(x_t^0)$ (clearly, $\nabla f(x_t) - \nabla f(x_t^0)$ is $\mathcal{F}_{t-1}$-measurable). This motivates why we wish to bound terms of the form Eq. (9).

Observe that each $(Z_t Z_t^\top - I)v_t \mid \mathcal{F}_{t-1}$ term is subExponential rather than subGaussian. While it is possible to define norm-subExponential vectors in analogous way to norm-subGaussian vectors, the corresponding moment generating function

(MGF) for subExponential random variables is not defined on the entirety of $\mathbb{R}$. When bounding a sum in the form of $\sum_{t=0}^{\tau-1}(Z_t Z_t^\top - I)v_t$, this creates a subtle but challenging technical issue.

Following the intuition outlined in the main text, we bypass this difficulty by proving the following result. For notational simplicity, we introduce the function

$$\mathrm{lr}(x) := \log(x\log(x)): \tag{10}$$

We now recall Proposition 2 which we first introduced in the main text.

**Proposition 2.** Let $F_t; t \geq 1$ be a filtration. Let $(Z_t)_{t \geq 0}$ be a sequence of random vectors following the distribution $N(0; I)$ such that $Z_t \in F_t$ and is independent of $F_{t-1}$, and let $(v_t)_{t \geq 0}$ be a sequence of random vectors such that $v_t \in F_{t-1}$. For each $\tau \geq 0$, let

$$W_\tau = \sum_{t=0}^{\tau-1} M_t(Z_t Z_t^\top - I)v_t;$$

where each $M_t$ is a deterministic matrix of appropriate dimension. Then, there exist some absolute constants $c, C > 0$ such that for any $\tau \in Z^+$ and $\delta \in (0; 1=e]$, the following statements hold:

1. For any $\zeta > 0$, with probability at least $1-\delta$, we have

$$\|W_\tau\| \leq c \cdot \sqrt{\sum_{t=0}^{\tau-1} \|M_t\|^2 d(\mathrm{lr}(C\tau=\delta))^2 \|v_t\|^2} + \frac{1}{\zeta}\log(C\tau d=\delta):$$

2. For any $B > b > 0$, with probability at least $1-\delta$,

$$\text{either} \quad \sum_{t=0}^{\tau-1} \|M_t\|^2 d(\mathrm{lr}(C\tau=\delta))^2 \|v_t\|^2 \geq B; \quad \text{or}$$

$$\|W_\tau\| \leq c \cdot \sqrt{\max\left(\sum_{t=0}^{\tau-1} \|M_t\|^2 d(\mathrm{lr}(C\tau=\delta))^2 \|v_t\|^2; b\right)} \cdot \sqrt{(\log(C\tau d=\delta) + \log(\log(B=b) + 1))}$$

Moreover, as is clear from the bounds above, we may pick $C \geq 1$ such that $\log\frac{C}{\delta} \geq 1; 8\delta \in (0; \frac{1}{e}]$.

Proof. We will focus on proving the first point, since the second follows as a natural corollary of our proof of the first part and the proof of Corollary 8 in (Jin et al., 2019b). For simplicity, we shall assume $v_t = 0$ in the intermediate steps; extension to the general case is straightforward.

First of all, for $0 \leq \lambda < 1$, let

$$g(\lambda; \delta) = \sqrt{\frac{\pi}{2}} Z^{\sqrt{2\,\mathrm{lr}(1=\delta)}} (x^2-1)e^{-x^2=2}\,dx = \sqrt{\frac{\pi}{2}} \cdot e^{-\lambda^2=2} \cdot \sqrt{\frac{2\,\mathrm{lr}(1=\delta)}{\log(1=\delta)}}:$$

It's not hard to see that for a fixed $\delta \in (0; 1=e]$, $g(\lambda; \delta)$ is continuous and strictly increasing over $\lambda \in [0; 1)$. Then, since $\frac{\log x}{x} + 1 \geq x$ for $x \geq 1$, by plugging in $x = \log(1=\delta)$, we get

$$\frac{\mathrm{lr}(1=\delta)}{(\log(1=\delta))^2} = \frac{\log\log(1=\delta) + \log(1=\delta)}{(\log(1=\delta))^2} = \frac{1}{\log(1=\delta)}\left(\frac{\log\log(1=\delta)}{\log(1=\delta)} + 1\right) \geq 1;$$

which leads to

$$g(\sqrt{2}; \delta) = \sqrt{\frac{\pi}{2}} \cdot 2e^{-2} \cdot \sqrt{\frac{2\,\mathrm{lr}(1=\delta)}{\log(1=\delta)}} \geq \sqrt{\frac{\pi}{2}} \cdot 2e^{-2=e^2} \cdot \sqrt{\frac{\pi}{2}} > 0$$

17

for $\lambda \in (0; 1=e]$. Furthermore, we obviously have $g(0; \lambda) < 0$. Therefore $g(\alpha; \lambda) = 0$ has a unique solution in $(0; 2)$, which we denote by $\alpha(\lambda)$.[4] These results imply that, for a random variable $Z$ following the standard normal distribution, we have

$$E\left[(Z^2 - 1)\mathbf{1}_{|\alpha(\lambda)| \ge |Z|}\right]_{|Z| \ge \sqrt{2\ln(1=\lambda)}} = \frac{1}{\sqrt{2\pi}}\int_{\sqrt{2\ln(1=\lambda)}}^{\alpha(\lambda)} (x^2 - 1)e^{-x^2=2}\,dx = g(h(\lambda); \lambda) = 0$$

and

$$P(\alpha(\lambda) \ge |Z| \ge \sqrt{2\ln(1=\lambda)}) = 1 - 2\left(\frac{1}{\sqrt{2\pi}}\int_{\sqrt{2\ln(1=\lambda)}}^{\infty} e^{-x^2=2}\,dx + \frac{1}{\sqrt{2\pi}}\int_0^{\alpha(\lambda)} e^{-x^2=2}\,dx\right)$$

$$\le 1 - 2\left(\frac{1}{2}\exp\left(-\frac{2\ln(1=\lambda)}{2}\right) + \frac{\alpha(\lambda)}{\sqrt{2\pi}}\right) = 1 - 2\left(\frac{\lambda}{2\log(1=\lambda)} + \frac{\alpha(\lambda)}{\sqrt{2\pi}}\right)$$

$$\le 1 - 2\lambda\left(\frac{1}{2} + \frac{\alpha}{\sqrt{2\pi}}\right) \le 1 - C\lambda$$

for any $\lambda \in (0; 1=e]$, where we define the absolute constant $C = 2(1=2 + 2=\sqrt{2\pi})$.

Now we let $A_t$ denote the event

$$A_t = \left\{\alpha(\lambda) \ge \frac{Z_t^\top v_t}{\|v_t\|} \ge \sqrt{2\ln(1=\lambda)}\right\}:$$

Since $Z_t^\top v_t = \|v_t\|$ conditioned on $F_{t-1}$ follows the standard normal distribution, we have

$$P_{F_{t-1}}(A_t) \ge 1 - C\lambda;\tag{11}$$

and

$$E_{F_{t-1}}\left[v_t^\top Z_t Z_t^\top I v_t \mathbf{1}_{A_t}\right] = 0:$$

Moreover, for any random vector $u \in F_{t-1}$ that is orthogonal to $v_t$, we have

$$E_{F_{t-1}}\left[u^\top Z_t Z_t^\top I v_t \mathbf{1}_{A_t}\right] = E_{F_{t-1}}\left[u^\top Z_t\right] \cdot E_{F_{t-1}}\left[Z_t^\top v_t \mathbf{1}_{A_t}\right] = 0;$$

where we used the fact that $Z_t^\top u$ is independent of $Z_t^\top v_t$ conditioned on $F_{t-1}$. Therefore

$$E_{F_{t-1}}\left[(Z_t Z_t^\top - I)v_t \mathbf{1}_{A_t}\right] = 0:$$

Consider defining then the random variable $Q_t$ by

$$Q_t := (Z_t Z_t^\top - I)v_t \cdot \mathbf{1}_{A_t}:$$

We now show that $Q_t \mid F_{t-1}$ is norm-subGaussian. Let $u \in R^d$ with $\|u\| = 1$ be arbitrary. We have

$$u^\top Q_t = u^\top (Z_t Z_t^\top - I)v_t \cdot \mathbf{1}_{A_t}$$

$$= u^\top\left[\frac{v_t v_t^\top}{\|v_t\|^2} + I - \frac{v_t v_t^\top}{\|v_t\|^2}\right](Z_t Z_t^\top - I)v_t \cdot \mathbf{1}_{A_t}$$

$$= u^\top v_t\left[\frac{|Z_t^\top v_t|^2}{\|v_t\|^2} - 1\right]\mathbf{1}_{A_t} + u^\top\left[I - \frac{v_t v_t^\top}{\|v_t\|^2}\right](Z_t Z_t^\top - I)v_t \cdot \mathbf{1}_{A_t}$$

$$= u^\top v_t\left[\frac{|Z_t^\top v_t|^2}{\|v_t\|^2} - 1\right]\mathbf{1}_{A_t} + u_\perp^\top Z_t Z_t^\top v_t \cdot \mathbf{1}_{A_t};$$

___

[4] By letting $W_0(x)$ denote the the principal branch of the Lambert $W$ function, it can be shown that

$$\alpha(\lambda) = \sqrt{-W_0\left(-\frac{\lambda^2 \cdot 2\ln(1=\lambda)}{(\log(1=\lambda))^2}\right)}:$$

where we denote $\tilde{u} = \left(I - \frac{v_t v_t^\top}{\|v_t\|^2}\right) u$. Since

$$u^\top v_t \left(\frac{|Z_t^\top v_t|^2}{\|v_t\|^2} - 1\right) \mathbf{1}_{A_t} \leq |u^\top v_t|(2\ln(1/\delta) - 1),$$

we see that $u^\top v_t \left(\frac{|Z_t^\top v_t|^2}{\|v_t\|^2} - 1\right) \mathbf{1}_{A_t}$ conditioned on $\mathcal{F}_{t-1}$ is $|u^\top v_t|(2\ln(1/\delta) - 1)$-subGaussian. Furthermore, since $|\tilde{u}^\top Z_t Z_t^\top v_t \mathbf{1}_{A_t}| \leq |Z_t^\top \tilde{u}| \cdot \sqrt{2\ln(1/\delta)}\|v_t\|$, we have

$$P_{\mathcal{F}_{t-1}}\left(|\tilde{u}^\top Z_t Z_t^\top v_t \mathbf{1}_{A_t}| \geq s\right) \leq P_{\mathcal{F}_{t-1}}\left(|Z_t^\top \tilde{u}| \geq \frac{s}{\sqrt{2\ln(1/\delta)}\|v_t\|}\right);$$

and since $Z_t^\top \tilde{u} = \|\tilde{u}\| \mid \mathcal{F}_{t-1}$ follows the standard normal distribution, we see that $\tilde{u}^\top Z_t Z_t^\top v_t \mathbf{1}_{A_t}$ is a $\sqrt{2\ln(1/\delta)}\|\tilde{u}\|\|v_t\|$-subGaussian variable. Note that $u^\top Q_t$ is just the sum of $u^\top v_t \left(\frac{|Z_t^\top v_t|^2}{\|v_t\|^2} - 1\right)\mathbf{1}_{A_t}$ and $\tilde{u}^\top Z_t Z_t^\top v_t \mathbf{1}_{A_t}$, we can conclude that $u^\top Q_t$ is subGaussian with parameter

$$(2\ln(1/\delta) - 1)|u^\top v_t| + \sqrt{2\ln(1/\delta)}\|\tilde{u}\|\|v_t\|$$
$$\leq 2\ln(1/\delta)(|u^\top v_t| + \|\tilde{u}\|\|v_t\|) \leq 2\sqrt{2}\ln(1/\delta)\sqrt{|u^\top v_t|^2 + \|\tilde{u}\|^2\|v_t\|^2}$$
$$= 2\sqrt{2}\ln(1/\delta)\|v_t\|;$$

whenever $\delta \in (0, 1/e]$. Consequently, by Lemma 1 in ([Jin et al., 2019b]), we see that $Q_t \mid \mathcal{F}_{t-1}$ is $8\ln(1/\delta)\sqrt{d}\|v_t\|$-norm-subGaussian.

It follows easily that $M_t Q_t \mid \mathcal{F}_{t-1}$ is mean-zero and $8\ln(1/\delta)\|M_t\|_2\|v_t\|\sqrt{d}$-norm-subGaussian. Hence, by Lemma 6 in ([Jin et al., 2019a]), we know that there exists an absolute constant $c$ such that for any $\iota > 0$ and $\delta > 0$, we have that with probability at least $1 - \delta$,

$$\left\|\sum_{t=0}^{\tau-1} M_t Q_t\right\| \leq c \sqrt{\sum_{t=0}^{\tau-1} d(\ln(1/\delta))^2\|M_t\|_2^2\|v_t\|^2 + \frac{1}{\iota}\log(2d/\delta)}.$$

Now, consider denoting the event

$$A := \bigcap_{t=0}^{\tau-1} A_t = \left\{Z_t^\top v_t \in [-2\sqrt{\ln(1/\delta)}\|v_t\|, \sqrt{2\ln(1/\delta)}\|v_t\|] ; \forall t \in \{0, \ldots, \tau-1\}\right\}$$

By the union bound and Eq. (11), we note that

$$P(A) \geq 1 - \tau\delta C:$$

Moreover, note that on the event $A$, $\sum_{t=0}^{\tau-1} M_t Q_t = \sum_{t=0}^{\tau-1} M_t (Z_t Z_t^\top - I) v_t$. Hence,

$$P\left(\left\|\sum_{t=0}^{\tau-1} M_t (Z_t Z_t^\top - I) v_t\right\| \geq c\sqrt{\sum_{t=0}^{\tau-1} d(\ln(1/\delta))^2\|M_t\|_2^2\|v_t\|^2 + \frac{1}{\iota}\log(2d/\delta)}\right)$$
$$\leq P\left(\left\|\sum_{t=0}^{\tau-1} M_t Y_t\right\| \geq c\sqrt{\sum_{t=0}^{\tau-1} d(\ln(1/\delta))^2\|M_t\|_2^2\|v_t\|^2 + \frac{1}{\iota}\log(2d/\delta)}; \text{ and } A \text{ happens}\right)$$
$$\leq 1 - P\left(\left\|\sum_{t=0}^{\tau-1} M_t Y_t\right\| \leq c\sqrt{\sum_{t=0}^{\tau-1} d(\ln(1/\delta))^2\|M_t\|_2^2\|v_t\|^2 + \frac{1}{\iota}\log(2d/\delta)}\right) + P(A^c)$$
$$\leq 1 - (\delta + \tau\delta C):$$

Now, by rescaling $\delta$ to $\delta = \delta/(C\tau + 1)$, we get the desired result. Note that this $C$ is different from the $C$ in the statement of the lemma by an absolute multiplicative factor. $\qquad\square$

19

## C.4. sub-Weibull random variables

In our work, we occasionally require bounding sums of heavy-tailed distribution, e.g. higher powers of $Z$, where $Z \sim N(0, I)$. To this end, we consider the following definition of sub-Weibull random variables.

**Definition 7.** We say that a random variable $X \in \mathbb{R}$ is sub-Weibull $(K, \theta)$ for some $K, \theta > 0$,

$$P(|X| \geq s) \leq 2\exp(-(s/K)^{1/\theta}) \quad \forall s \geq 0.$$

For instance, the standard normal distribution is sub-Weibull $(1, \frac{1}{2})$. From the way we define the tail parameter $\theta$, the larger the $\theta$, the heavier the tail of the distribution.

In our work, we need to show that the sum of sub-Weibull random variables is again sub-Weibull, which is ensured by the following result

**Lemma 9.** Suppose $X$ and $Y$ are sub-Weibull $(K_X, \theta)$ and sub-Weibull $(K_Y, \theta)$ respectively. Then $XY$ is sub-Weibull $(C(K_X \cdot K_Y), 2\theta)$ and $X + Y$ is sub-Weibull $(C(K_X + K_Y), \theta)$ for some absolute constant $C > 0$.

A helpful result is the following, which bounds the sum of identically distributed sub-Weibull random variables.

**Lemma 10** (Corollary 3.1 in (Vladimirova et al., 2020)). Suppose $X_1, \ldots, X_n$ are identically distributed $(K^0, \theta)$ sub-Weibull random variables. Then, for some absolute constant $c > 0$, for all $s \geq ncK^0$, we have

$$P\left(\left|\sum_{i=1}^{n} X_i\right| \geq s\right) \leq \exp\left(-\left(\frac{s}{ncK^0}\right)^{1/\theta}\right)$$

In our work, we frequently need to bound sums of the $k$th power of the norm of a standard $d$-dimensional Gaussian. We do so using Lemma 10.

**Lemma 11.** Suppose $X_i \overset{i.i.d.}{\sim} N(0, I_d)$ for $i \in [n]$. Then, for any $k \in \mathbb{Z}^+$, there exists absolute constants $c, C > 0$ such that for any $\delta \in (0, 1)$, with probability at least $1 - \delta$,

$$\sum_{i=1}^{n} \|X_i\|^{2k} \leq nCc^k d^k (1 + (\log(1/\delta))^k).$$

In particular, for any $\delta \in (0, 1/e)$ such that $\log(1/\delta) \geq 1$, it follows that

$$\sum_{i=1}^{n} \|X_i\|^{2k} \leq 2nCc^k d^k (\log(1/\delta))^k.$$

**Proof.** First, observe that for any $j \in [d]$, $(X_i)_j^2$, being subExponential, is $(1, 1)$-subWeibull. Then, by Lemma 9, $\|X_i\|^2 = \sum_{j=1}^{d} (X_i)_j^2$ is $(cd, 1)$ for some absolute constant $c$. Now, it follows from definition of sub-Weibullness in Definition 7 that $\|X_i\|^{2k}$ is $(c^k d^k, k)$-subWeibull. Hence, applying Lemma 10, we have that there exists absolute constant $C > 0$ such that for any $s \geq nCc^k d^k$,

$$P\left(\left|\sum_{i=1}^{n} \|X_i\|^{2k}\right| \geq s\right) \leq \exp\left(-\left(\frac{s}{nCc^k d^k}\right)^{1/k}\right)$$

Choosing $s = (1 + (\log(1/\delta))^k)nCc^k d^k$, we arrive then at the desired result. □

## C.5. Supermartingale concentration inequalities

We first state and prove a supermartingale-type concentration inequality of the form we later require.

**Lemma 12.** Consider a filtration of sigma-algebras $F_0 \subseteq F_1 \subseteq \cdots \subseteq F_{n-1} \subseteq F_n$ and a sequence of random variables $X_1, \ldots, X_n$ such that $X_i \in F_i$. Suppose that

$$P_{F_{i-1}}(X_i \leq a) = 1 \quad \text{and} \quad P_{F_{i-1}}(X_i \geq b) \leq p \tag{12}$$

20

for some $a, b > 0$ and $0 < p \leq \frac{1}{2}$. Then, for any $0 < \lambda \leq b$ such that $\mu - b + \eta \leq j - \mu \leq \frac{1-p}{p}(a + \eta)$, we have

$$P\left(\sum_{i=1}^{n} X_i \geq \mu n + s\right) \leq \exp\left(-\frac{s^2}{4n(b-\eta)^2}\right); \quad \forall s > 0.$$

**Proof.** Observe that by Markov's inequality, for any $\lambda > 0$,

$$P\left(\sum_{i=1}^{n} X_i \geq \mu n + s\right) = P\left(\exp\left(\lambda \sum_{i=1}^{n}(X_i + \eta)\right) \geq \exp(\lambda s)\right) \leq \frac{E[\exp(\lambda \sum_{i=1}^{n}(X_i + \eta))]}{\exp(\lambda s)}.$$

Now, observe that

$$E\left[\exp\left(\lambda \sum_{i=1}^{n}(X_i + \eta)\right)\right] = E\left[E_{F_{n-1}}\left[\exp\left(\lambda \sum_{i=1}^{n}(X_i + \eta)\right)\right]\right]$$

$$= E\left[\exp\left(\lambda \sum_{i=1}^{n-1}(X_i + \eta)\right) E_{F_{n-1}}[\exp(\lambda(X_n + \eta))]\right]. \qquad (13)$$

Let us now compute $E_{F_{n-1}}[\exp(\lambda(X_n + \eta))]$:

$$E_{F_{n-1}}[\exp(\lambda(X_n + \eta))]$$

$$= \int_{(-1; -b]} \exp(\lambda(x + \eta)) P_{F_{n-1}}(X_n \in dx) + \int_{(-b; a]} \exp(\lambda(x + \eta)) P_{F_{n-1}}(X_n \in dx)$$

$$\leq P_{F_{n-1}}(X_n \leq -b) \exp(\lambda(-b + \eta)) + P_{F_{n-1}}(-b < X_n \leq a) \exp(\lambda(a + \eta))$$

$$\leq p \exp(\lambda(-b + \eta)) + (1 - p) \exp(\lambda(a + \eta)).$$

Then observe that by our choice of $\eta$, $-b + \eta < 0$, and that $j - \mu \leq (a + \eta)\frac{1-p}{p}$. Since we assumed $p \leq \frac{1}{2}$, this means that $\frac{1-p}{p} \geq 1$ and so for any $k \geq 1$,

$$j - b + \eta \leq j \leq (a + \eta)\frac{1-p}{p} =) j - b + \eta \leq j \leq (a + \eta)\left(\frac{1-p}{p}\right)^{1-k} =) p j - b + \eta^k \leq (1-p)(a + \eta)^k.$$

Consequently, by Taylor expansion,

$$p \exp(\lambda(-b + \eta)) + (1 - p) \exp(\lambda(a + \eta))$$

$$= 1 + \sum_{k=1}^{\infty} \frac{\lambda^k (p(-b + \eta)^k + (1-p)(a + \eta)^k)}{k!} \leq 1 + \sum_{k=1}^{\infty} \frac{\lambda^k (p(-b + \eta)^k + p j - b + \eta^k)}{k!}$$

$$= 1 + \sum_{k=1}^{\infty} \frac{\lambda^{2k} 2 p j - b + \eta^{2k}}{(2k)!} \leq 1 + \sum_{k=1}^{\infty} \frac{\lambda^{2k} j - b + \eta^{2k}}{(k)!}$$

$$= \exp(\lambda^2 (-b + \eta)^2);$$

which leads to

$$E_{F_{n-1}}[\exp(\lambda(X_n + \eta))] \leq \exp(\lambda^2 (-b + \eta)^2).$$

Now, continuing from Eq. (13), we have that

$$E\left[\exp\left(\lambda \sum_{i=1}^{n}(X_i + \eta)\right)\right] \leq E\left[\exp\left(\lambda \sum_{i=1}^{n-1}(X_i + \eta)\right) E_{F_{n-1}}[\exp(\lambda(X_n + \eta))]\right]$$

$$\leq E\left[\exp\left(\lambda \sum_{i=1}^{n-1}(X_i + \eta)\right)\right] \exp(\lambda^2 (b - \eta)^2)$$

21

$$\vdots$$

$$\exp(n\lambda^2(b-\beta)^2).$$

Thus, for any $\lambda > 0$ and $s \geq 0$,

$$P\left(\sum_{i=1}^{n} X_i \geq n\beta + s\right) \leq \frac{E[\exp(\lambda(\sum_{i=1}^{n}(X_i - \beta + \beta)))]}{\exp(\lambda s)}$$

$$\leq \exp(n\lambda^2(b-\beta)^2 - \lambda s).$$

By finding the minimizing $\lambda$, we find that

$$P\left(\sum_{i=1}^{n} X_i \geq n\beta + s\right) \leq \exp\left(-\frac{s^2}{4n(b-\beta)^2}\right);$$

which completes the proof. □

We will later require a weakened form of a supermartingale concentration inequality, as stated and proven below.

**Proposition 3** (Weakened supermartingale concentration inequality). Consider a filtration of sigma-algebras $F_0 \subseteq F_1 \subseteq \cdots \subseteq F_n$ and a sequence of random variables $X_1, \ldots, X_n$ such that $X_i \in F_i$. Consider for each $i \in \{1, \ldots, n\}$ a bad set $B_i$ where $1_{B_i} \in F_{i-1}$, and suppose

$$P_{F_{i-1}}(X_i 1_{B_i^c} \leq a) = 1 \quad \text{and} \quad P_{F_{i-1}}(X_i 1_{B_i^c} \geq b) \leq \beta$$

for some $a, b > 0$ and $0 \leq p \leq 1/2$. Then, for any $0 < \beta \leq b$ such that $j \leq b\beta + j \leq \frac{1-p}{p}(a + \beta)$, we have

$$P\left(\sum_{i=1}^{n} X_i \geq n\beta + s\right) \leq \exp\left(-\frac{s^2}{4n(b-\beta)^2}\right) + \sum_{i=1}^{n} P(X_i \in B_i); \quad \forall s > 0.$$

**Proof.** We define $Q_i := X_i 1_{B_i^c}$. We can then apply Lemma 12 and get

$$P\left(\sum_{i=1}^{n} Q_i \geq n\beta + s\right) \leq \exp\left(-\frac{s^2}{4n(b-\beta)^2}\right).$$

Since $P(X_i \neq Q_i \text{ for some } i \in [n]) \leq \sum_i P(X_i \in B_i)$, it follows that

$$P\left(\sum_{i=1}^{n} X_i \geq n\beta + s\right) \leq \exp\left(-\frac{s^2}{4n(b-\beta)^2}\right) + \sum_{i=1}^{n} P(X_i \in B_i);$$

which completes the proof. □

## D. Function decrease in large gradient regime

In this section, we show that sufficient function decrease can be made across the iterations with large gradients. We first restate and prove the function decrease lemma (Lemma 1), first introduced in the main text. We then provide a detailed roadmap of our proof in the subsequent discussion following the proof of Lemma 1.

**Lemma 1** (Function decrease for batch zeroth-order optimization). Suppose at each time $t$, the algorithm performs the update step (with batch-size parameter $m \leq d$)

$$x_{t+1} = x_t - \eta g_u^{(m)}(x_t) + Y_t;$$

where

$$g_u^{(m)}(x_t) = \frac{1}{m}\sum_{i=1}^{m}\frac{f(x_t + uZ_{t;i}) - f(x_t - uZ_{t;i})}{2u}Z_{t;i};$$

where each $Z_{t;i}$ is drawn i.i.d from $N(0; I)$, $u > 0$ is the smoothing radius, and $Y_t \sim N(0; \frac{r^2}{d}I)$ with $r > 0$ denoting the perturbation radius.

Then, there exist absolute constants $c > 0$, $C_1 \geq 1$ such that, for any $T \in Z^+$ and $T > 0$, $\eta > 0$ and $\epsilon \in (0; 1 = e]$, upon defining $H_{0;\tau}(\epsilon)$ to be the event on which the inequality

$$f(x_\tau) \leq f(x_0) \tag{3}$$

$$- \frac{3}{4} \sum_{t=0}^{\tau - 1} \eta \frac{1}{m} \sum_{i=1}^{m} \|Z_{t;i}^> r f(x_t)\|^2 \tag{4}$$

$$+ \left( -\eta + \frac{c_1 L \eta^2 \iota^3 d}{m} \right) \sum_{t=0}^{\tau - 1} \|r f(x_t)\|^2$$

$$+ \eta u^4 \rho^2 c_1 d^3 \log \frac{T}{\epsilon}^3 + \eta L^2 u^4 \rho^2 c_1 d^4 \log \frac{T}{\epsilon}^4$$

$$+ \eta c_1 r^2 (\eta + L\eta^2) \log \frac{T}{\epsilon} + \eta c_1 L^2 \eta^2 r^2 \tag{5}$$

is satisfied (where $\iota := \log(C_1 dmT = \epsilon)$), we have

$$P(H_{0;\tau}(\epsilon)) \geq 1 - \frac{(\tau + 4)\epsilon^0}{T}; \qquad P(\setminus_{\tau=1}^0 H_{0;\tau}(\epsilon)) \geq 1 - \frac{5\epsilon^0}{T}$$

for any $0 \leq \tau^0 \leq T$.

Proof. First, for each $t \in \{-1; \ldots; \}$, we define $F_t$ to be the sigma-algebra generated by

$$x_0; \quad (f Z_{0;i} g_{i=1}^m; \ldots; f Z_{t;i} g_{i=1}^m); \quad (Y_0; \ldots; Y_t):$$

Note that $F_{-1}$ is the sigma-algebra generated only by $x_0$.

By Taylor expansion, for any $x, y \in R^d$, there exists $\theta \in [0; 1]$ such that $f(x + y) = f(x) + h r f(x); y i + \frac{1}{2} y^> r^2 f(x + \theta y) y$. Therefore

$$\frac{f(x_t + uZ_{t;i}) - f(x_t - uZ_{t;i})}{2u} = h r f(x); Z_{t;i} i + \frac{u}{2} Z_{t;i}^> H_{t;i} Z_{t;i}$$

with

$$H_{t;i} = \frac{r^2 f(x + \theta_{i;+} uZ_{t;i}) - r^2 f(x - \theta_{i;-} uZ_{t;i})}{2}$$

for some $\theta_{i;\pm} \in [0; 1]$, and

$$x_{t+1} = x_t - \eta \left( \frac{1}{m} \sum_{i=1}^{m} Z_{t;i} Z_{t;i}^> r f(x_t) + \frac{u}{2} Z_{t;i} Z_{t;i}^> H_{t;i} Z_{t;i} + Y_t \right) \tag{14}$$

By the $\rho$-Hessian Lipschitz property of $f$, it follows that $\|H_{t;i}\| \leq \rho u \|Z_{t;i}\|$

Observe that

$$f(x_{t+1}) \overset{(i)}{\leq} f(x_t) + h x_{t+1} - x_t; r f(x_t) i + \frac{L}{2} \|x_{t+1} - x_t\|^2$$

$$\overset{(ii)}{=} f(x_t) - \eta \frac{1}{m} \sum_{i=1}^{m} \|Z_{t;i}^> r f(x_t)\|^2 - \eta \frac{1}{m} \sum_{i=1}^{m} \frac{u}{2} Z_{t;i}^> r f(x_t) Z_{t;i}^> H_{t;i} Z_{t;i} - \eta h r f(x_t); Y_t i$$

$$+ \frac{L\eta^2}{2} \left\| \frac{1}{m} \sum_{i=1}^{m} Z_{t;i} Z_{t;i}^> r f(x_t) + \frac{u}{2} Z_{t;i} Z_{t;i}^> H_{t;i} Z_{t;i} + Y_t \right\|^2$$

$$\overset{(iii)}{\leq} f(x_t) - \eta \frac{1}{m} \sum_{i=1}^{m} \|Z_{t;i}^> r f(x_t)\|^2 + \frac{\eta}{m} \sum_{i=1}^{m} \left( \frac{\|Z_{t;i}^> r f(x_t)\|^2}{4} + \frac{u^2 \|Z_{t;i}^> H_{t;i} Z_{t;i}\|^2}{4} \right) - \eta h r f(x_t); Y_t i$$

23

$$+ \frac{L^2}{2} @2 \left[\frac{1}{m}\sum_{i=1}^{m} Z_{t;i} Z_{t;i}^\top \nabla f(x_t)\right]^2 + u^2 \left[\frac{1}{m}\sum_{i=1}^{m} Z_{t;i} Z_{t;i}^\top H_{t;i} Z_{t;i}\right]^2 + 4kY_t k^2 A$$

$$\overset{(iv)}{\leq} f(x_t) - \frac{3}{4m}\sum_{i=1}^{m} \left[Z_{t;i}^\top \nabla f(x_t)\right]^2 + \frac{u^2}{m}\sum_{i=1}^{m} \frac{u^2 L^2 kZ_{t;i} k^6}{4} - \langle hr f(x_t); Y_t i\rangle$$

$$+ \frac{L^2}{2} @2 \left[\frac{1}{m}\sum_{i=1}^{m} Z_{t;i} Z_{t;i}^\top \nabla f(x_t)\right]^2 + \frac{u^2}{m}\sum_{i=1}^{m} u^2 L^2 kZ_{t;i} k^8 + 4kY_t k^2 A$$

$$\leq f(x_t) - \frac{3}{4m}\sum_{i=1}^{m} \left[Z_{t;i}^\top \nabla f(x_t)\right]^2 + \frac{u^4 L^2}{4m}\sum_{i=1}^{m} kZ_{t;i} k^6 + \frac{L^2 u^4 L^2}{2m}\sum_{i=1}^{m} kZ_{t;i} k^8 - \langle hr f(x_t); Y_t i\rangle$$

$$+ \frac{L^2}{2} @2 \left[\frac{1}{m}\sum_{i=1}^{m} Z_{t;i} Z_{t;i}^\top \nabla f(x_t)\right]^2 + 4kY_t k^2 A \qquad (15)$$

Above, to derive (i), we used the $L$-smoothness of $f$. To derive (ii), we used the expression for $(x_{t+1} - x_t)$ shown in Eq. (14). To derive (iii), we used the fact that $u(a^2 + b^2) \geq 2$ for any $a; b \in \mathbb{R}_{\geq 0}$, as well as two applications of the fact that $ka + bk^2 \leq 2(kak^2 + kbk^2)$ for any two vectors $a; b \in \mathbb{R}^d$. To derive (iv), we used the fact that $H_{t;i} \leq u L kZ_{t;i} k$.

To continue from Eq. (15), we first observe that we can rewrite

$$Z_{t;i} Z_{t;i}^\top \nabla f(x_t) = (Z_{t;i} Z_{t;i}^\top - I)\nabla f(x_t) + \nabla f(x_t);$$

so that

$$\left[\frac{1}{m}\sum_{i=1}^{m} Z_{t;i} Z_{t;i}^\top \nabla f(x_t)\right]^2 \leq 2\left[\frac{1}{m}\sum_{i=1}^{m}(Z_{t;i} Z_{t;i}^\top - I)\nabla f(x_t)\right]^2 + 2kr f(x_t)k^2:$$

Observe that we can apply the bound in Proposition 2 to $\sum_{i=1}^{m}(Z_{t;i} Z_{t;i}^\top - I)\nabla f(x_t)$, and since $Z_{t;i}$ is independent of $F_{t-1}$ for all $i$, we know there exist absolute constants $c_1 > 0; C_1 \geq 1$ such that for any $\delta \in (0; 1=e]$ and $\epsilon > 0$, with probability at least $1 - \delta$ conditioned on $F_{t-1}$,

$$\left\|\sum_{i=1}^{m}(Z_{t;i} Z_{t;i}^\top - I)\nabla f(x_t)\right\| \leq c_1 \sum_{i=1}^{m} d(\ln(C_1 m=\delta))^2 kr f(x_t)k^2 + \frac{1}{\epsilon}\log(C_1 dm=\delta)$$

$$= c_1 md(\ln(C_1 m=\delta))^2 kr f(x_t)k^2 + \frac{1}{\epsilon}\log(C_1 dm=\delta): \qquad (16)$$

Moreover, since $C_1 \geq 1$, $\log(C_1 dm=\delta)$ and $\ln(C_1 m=\delta)$ both are at least $1$ as long as $\delta \in (0; 1=e]$. Observe that conditioned on $F_{t-1}$, $\nabla f(x_t)$ is fixed. Hence, we can pick

$$\epsilon = \sqrt{\frac{1}{c_1 md \ln(C_1 dm=\delta) kr f(x_t)k}}$$

which is $F_{t-1}$-measurable, and plug it into Eq. (16) to find that the probability conditioned on $F_{t-1}$ of the following event

$$\left\|\sum_{i=1}^{m}(Z_{t;i} Z_{t;i}^\top - I)\nabla f(x_t)\right\| \leq 2\sqrt{c_1}(\ln(C_1 dm=\delta))^{3=2}\sqrt{mdkr f(x_t)k} \qquad (17)$$

is at least $1 - \delta$. By taking the total expectation, it follows that the event has a total probability at least $1 - \delta$. Thus, with probability at least $1 - \delta$,

$$\left[\frac{1}{m}\sum_{i=1}^{m} Z_{t;i} Z_{t;i}^\top \nabla f(x_t)\right]^2 \leq 2\left[\frac{1}{m}\sum_{i=1}^{m}(Z_{t;i} Z_{t;i}^\top - I)\nabla f(x_t)\right]^2 + 2kr f(x_t)k^2$$

$$4c_1(\ln(C_1 dm=))^3 \frac{d}{m} kr f(x_t)k^2 + 2 kr f(x_t)k^2$$

$$c_2(\ln(C_1 dm=))^3 \frac{d}{m} kr f(x_t)k^2; \tag{18}$$

where the last inequality comes from the fact that $\ln(C_1 dm=) \geq 1$, our assumption at the outset of the appendix that $d \leq m$, and denoting $c_2 := 4c_1 + 2$.

Denote the event $H_{0;}(\ )$ as the event that

$$f(x) \quad f(x_0) \leq \sum_{t=0}^{X-1} \frac{3}{4m} \sum_{i=1}^{X^m} Z_{t;i}^{\geq} r f(x_t)^2 + L^2 \frac{c_2 d(\ln(C_1 dm=))^3}{m} \sum_{t=0}^{X-1} kr f(x_t)k^2$$

$$+ \frac{u^4 {}^2}{4m} \sum_{t=0}^{X-1}\sum_{i=1}^{X^m} kZ_{t;i} k^6 + \frac{L^2 u^4 {}^2}{2m} \sum_{t=0}^{X-1}\sum_{i=1}^{X^m} kZ_{t;i} k^8$$

$$\sum_{t=0}^{X-1} hr f(x_t); Y_t i + 2L^2 \sum_{t=0}^{X-1} kY_t k^2 \tag{19}$$

holds.

Now, continuing from Eq. (15), and using the bound in Eq. (18), summing over the iterations $t$ from $0$ to $-1$, we find using the union bound that $P(\bigcap_{=1}^{0} H_{0;}(\ )) \geq 1 - {}^0$, $P(H_{0;}(\ )) \geq 1 - :$

Now, by Lemma 6, for any $2 (0; 1); \ > 0$, with probability at least $1 - $, there exists an absolute constant $c_3 > 0$ such that

$$\sum_{t=0}^{X-1} hr f(x_t); Y_t i \leq \frac{1}{} \sum_{t=0}^{X-1} kr f(x_t)k^2 + c_3 r^2 \log(1=) : \tag{20}$$

Meanwhile, since $Y_t \sim N(0; (r^2=d)I)$, $kY_t k^2$ is sub-exponential with sub-exponential norm $cr^2$ for some absolute constant $c > 0$, and by Bernstein's inequality (Lemma 8), there exists some absolute constant $c_4$ such that

$$\sum_{t=0}^{X-1} kY_t k^2 \leq c_4 r^2 ( + \log(1=)) \tag{21}$$

with probability at least $1 - $.

To bound $\sum_{t=0}^{P-1} \frac{1}{m} \sum_{i=1}^{P^m} kZ_{t;i} k^6$ and $\sum_{t=0}^{P-1} \frac{1}{m} \sum_{i=1}^{P^m} kZ_{t;i} k^8$, both sums of heavy tailed Gaussian moments, we use Lemma 11, which states that for any $k \in 2 \mathbb{Z}^+$ and $2 (0; 1)$, with probability at least $1 - $,

$$\frac{1}{m} \sum_{t=0}^{X-1}\sum_{i=1}^{X^m} kZ_{t;i} k^{2k} \leq c_5 (c_6)^k d^k (1 + (\log(1=))^k) \tag{22}$$

for some absolute constants $c_5, c_6 > 0$. As in the statement of the proof, using $= \ln(C_1 dm=)$ to ease the notation, denote the event that

$$f(x) \quad f(x_0) \leq \frac{3}{4} \sum_{t=0}^{X-1} \frac{1}{m} \sum_{i=1}^{X^m} Z_{t;i}^{\geq} r f(x_t)^2 + \left( - + \frac{c_2 L^2 {}^3 d}{m} \right) \sum_{t=0}^{X-1} kr f(x_t)k^2$$

$$+ \frac{u^4 {}^2}{2} c_5 c_6^3 d^3 \left(\log \frac{1}{}\right)^3 + L^2 u^4 {}^2 c_5 c_6^4 d^4 \left(\log \frac{1}{}\right)^4$$

$$+ (c_3 r^2 + 2c_4 Lr^2) \log \frac{1}{} + 2c_4 L^2 r^2$$

holds as $H_{0;}(\ )$.

25

Plugging Eq. (20), Eq. (21), and Eq. (22) into Eq. (19), by union bound, we see that

$$P(\setminus_{\ell=1}^{\delta_0} H_{0;\ell}(\epsilon)) \geq 1 - (\delta_0 + 4\delta_0) = 1 - 5\delta_0; \qquad P(H_{0;\ell}) \geq 1 - (\delta + 4\delta):$$

The final result then follows by rescaling $\delta$ to $\frac{\delta}{T}$ and denoting $c_1 := \max\{c_2; c_3; 2c_4; c_5 c_6^3 = 2; c_5 c_6^4\}$. □

Outline of proof approach. Similar to the first-order setting, our goal is to show that we can arrive at a contradiction $f(x_T) < \min_x f(x)$ when there is a large number of steps at which $\|\nabla f(x_t)\| \geq \epsilon$. Roughly speaking, as Eq. (5) shows, we need to prove a lower bound of the form

$$\sum_{t=0}^{T-1} \frac{1}{m} \sum_{i=1}^{m} \langle Z_{t;i}, \nabla f(x_t)\rangle^2 \geq \left(\frac{1}{\eta} + \frac{c_1 L^3 d}{m}\right) \sum_{t=0}^{T-1} \|\nabla f(x_t)\|^2 \tag{23}$$

for some $\eta$ which is not too large (an example would be picking $\eta$ such that it only scales logarithmically in the problem parameters). However, it is tricky to prove such a lower-bound in the zeroth-order setting. In particular, for small batch-sizes $m$, $\frac{1}{m} \sum_{i=1}^{m} \langle Z_{t;i}, \nabla f(x_t)\rangle^2$ could be small even when $\|\nabla f(x_t)\|^2$ is large; this is because for each $i \in [m]$, $Z_{t;i}$ could have a negligible component in the $\nabla f(x_t)$ direction. This necessitates a more careful analysis to prove a bound similar to Eq. (23). We do so using the following approach.

1. Intuitively, whilst for each individual iteration $t$, $\frac{1}{m} \sum_{i=1}^{m} \langle Z_{t;i}, \nabla f(x_t)\rangle^2$ could be small even when $\|\nabla f(x_t)\|^2$ is large, in a small number of (consecutive) iterations $\{t_0; \ldots; t_0 + t_f\}$, with high probability, there will be at least one iteration $t$ within $\{t_0; \ldots; t_0 + t_f - 1\}$, such that $\frac{1}{m} \sum_{i=1}^{m} \langle Z_{t;i}, \nabla f(x_t)\rangle^2 = \Theta(\|\nabla f(x_t)\|^2)$. We formalize this intuition in Lemma 14. Thus, we consider breaking the time-steps into chunks where each chunk has $t_f$ consecutive iterations.

2. Consider any such interval $\{t_0; \ldots; t_0 + t_f - 1\}$. There are two cases to consider.

   (a) The first case is when the gradient throughout all $t_f$ iterations is large enough to dominate the perturbation terms. Intuitively, in this case, it is not hard to see that given appropriate parameter choices, the gradient will change little throughout the $t_f$ iterations. In fact, as we formalize in Lemma 16, for an appropriate choice of $\eta$ and $\mu$, we can show that

   $$\frac{1}{2}\|\nabla f(x_{t_0})\| \leq \|\nabla f(x_t)\| \leq 2\|\nabla f(x_{t_0})\| \qquad \forall t \in \{t_0; \ldots; t_0 + t_f - 1\}:$$

   As a result, combined with point 1, we see that

   $$\sum_{t=t_0}^{t_0 + t_f - 1} \frac{1}{m} \sum_{i=1}^{m} \langle Z_{t;i}, \nabla f(x_t)\rangle^2 \geq \Omega(\|\nabla f(x_{t_0})\|^2):$$

   Thus, by choosing $\eta$ and $\mu$ judiciously, for such intervals, it is possible to show that

   $$\sum_{t=t_0}^{t_0 + t_f - 1} \frac{1}{m} \sum_{i=1}^{m} \langle Z_{t;i}, \nabla f(x_t)\rangle^2 \geq \Omega(\|\nabla f(x_{t_0})\|^2) \geq \left(\frac{1}{\eta} + \frac{c_1 L^3 d}{m}\right) \sum_{t=t_0}^{t_0 + t_f - 1} \|\nabla f(x_t)\|^2$$
   $$= \left(\frac{1}{\eta} + \frac{c_1 L^3 d}{m}\right) t_f \|\nabla f(x_{t_0})\|^2$$

   Thus, in these intervals, it is possible to obtain function improvement on the order of $\Omega(\|\nabla f(x_{t_0})\|^2)$.

   (b) The remaining case is when the gradient is small and dominated by the perturbation terms in any one of the $t_f$ iterations. In this case, as we show in Lemma 17, for each of the $t_f$ iterations, the gradient will be small and on the same scale as the perturbation terms. In turn, by choosing $\eta$ and $\mu$ appropriately, we can make the perturbation terms small. Thus, whilst these intervals may not contribute to function decrease, they also contribute little in the way of function increase.

26

3. When there are at least $\bar{T} = 4$ iterations with large gradient (i.e. $\|\nabla f(x_t)\| \geq \epsilon$), assuming $t_f$ divides $T$, it follows that there are at least $\bar{T} = (4t_f)$ intervals of length $t_f$ where one iteration in the interval contains a large gradient. By choosing $u$, $r$ and appropriately such they are dominated by it is possible to show that with high probability, such an interval cannot belong to the second case above, and must instead be from the first case. Since $\|\nabla f(x_{t_0})\| \geq \|\nabla f(x_{t_0})\|$ for each $t \in \{t_0, \ldots, t_0 + t_f - 1\}$ in this case, and we know that one of the iterations has a gradient with size at least $\epsilon$, it follows that we make function decrease progress of at least $(\epsilon^2)$ for such intervals. By appropriately choosing $u$ and $r$ to limit the effects of the intervals of the second form, we can then show a contradiction of the form $f(x_0) - f(x_T) < f$. We demonstrate this formally in Proposition 4.

We formalize our approach in the following series of results. First, for analytical convenience, we prove the following result showing that for any $t$, the perturbation terms $\|Y_t\|$ and $\frac{1}{m} \sum_{i=1}^{m} \|Z_{t;i}\|^4$ are bounded with high probability.

**Lemma 13.** There exists an absolute constant $c_3 > 0$ such that, for any $t \in \mathbb{N}$, the event

$$G_t(\delta) := \left\{ \|Y_t\|^2 \leq c_3^2 r^2 \left(1 + \frac{\log(T/\delta)}{d}\right) \quad \text{and} \quad \frac{1}{m} \sum_{i=1}^{m} \|Z_{t;i}\|^4 \leq 2c_3 d^2 \left(\log \frac{T}{\delta}\right)^2 \right\}$$

has probability at least $1 - 2\delta/T$ for any $\delta \in (0, 1/e]$.

*Proof.* Noting that $Y_t \sim N(0, (r^2/d) I)$, by applying Bernstein's inequality (Lemma 8), it can be shown that with probability at least $\delta/T$,

$$\|Y_t\|^2 \leq c_3^2 r^2 \left(1 + \frac{\log(T/\delta)}{d}\right);$$

where $c_3 > 0$ is some absolute constant. Then by using Lemma 11, applying the union bound, and redefining the constant $c_3$, we complete the proof. $\square$

Next, in Lemma 14, we show that in a small number of iterations, with high probability, there exists some iteration $t$ such that $\frac{1}{m} \sum_{i=1}^{m} \langle Z_{t;i}, \nabla f(x_t) \rangle^2 \geq \frac{1}{2} k r \|\nabla f(x_t)\|^2$.

**Lemma 14.** There exists an absolute constant $c_2 \geq 1$ such that, upon defining

$$t_f(\delta) = \left\lceil \frac{c_2}{m} \log \frac{T}{\delta} \right\rceil; \qquad \delta > 0;$$

and defining the event

$$B_{t_0}(\delta; k) := \bigcup_{t=t_0}^{t_0+t_f k - 1} \left\{ \frac{1}{m} \sum_{i=1}^{m} \langle Z_{t;i}, \nabla f(x_t) \rangle^2 \geq \frac{1}{2} k r \|\nabla f(x_t)\|^2 \right\};$$

we have

$$P(B_{t_0}(\delta; k)) \geq 1 - \frac{\delta}{T};$$

for any $\delta \in (0, 1)$, $t_0 \in \mathbb{N}$ and $k \geq t_f(\delta)$.

*Proof.* Denote the event

$$E_t = \left\{ \frac{1}{m} \sum_{i=1}^{m} |\langle Z_{t;i}, \nabla f(x_t) \rangle|^2 < \frac{1}{2} k r \|\nabla f(x_t)\|^2 \right\}:$$

Observe that, conditioned on $F_{t-1}$, the set of random variables $\left\{ k r \|\nabla f(x_t)\|^2 - \langle Z_{t;i}, \nabla f(x_t) \rangle^2 \right\}_{i=1}^{m}$ are independent, mean-zero, and subexponential with subexponential norm $c k r \|\nabla f(x_t)\|^2$ for some absolute constant $c > 0$. Hence

$$P_{F_{t-1}}(E_t) = P_{F_{t-1}}\left( \frac{1}{m} \sum_{i=1}^{m} |\langle Z_{t;i}, \nabla f(x_t) \rangle|^2 < \frac{1}{2} k r \|\nabla f(x_t)\|^2 \right)$$

$$= P_{F_{t-1}}\left( \sum_{i=1}^{m} \left( k r \|\nabla f(x_t)\|^2 - \langle Z_{t;i}, \nabla f(x_t) \rangle^2 \right) > \frac{m}{2} k r \|\nabla f(x_t)\|^2 \right)$$

27

$$\exp(-c^0 m);$$

where $c^0$ is some positive absolute constant. Then, for any $t_0, k \in \mathbb{N}$,

$$P\left[\frac{1}{m}\sum_{i=1}^{m} \left(Z_{t;i}^{\top} \nabla f(x_t)\right)^2 < \frac{1}{2}\|\nabla f(x_t)\|^2 \text{ for every } t \in [t_0; t_0 + k)\right]$$

$$= E\left[\prod_{t=t_0}^{t_0+k-1} \mathbb{1}_{E_t}\right] = E\left[\prod_{t=t_0}^{t_0+k-2} \mathbb{1}_{E_t} \cdot E_{F_{t_0+k-2}}\left[\mathbb{1}_{E_{t_0+k-1}}\right]\right]$$

$$\leq \exp(-c^0 m) \cdot E\left[\prod_{t=t_0}^{t_0+k-2} \mathbb{1}_{E_t}\right] \leq \cdots \leq \exp(-c^0 mk):$$

Therefore, by letting $c_2 = \max\{1; 1=c^0\}$ and

$$k \geq t_f(\eta) = \frac{c_2}{m}\log\frac{T}{\eta};$$

we get

$$P\left[\frac{1}{m}\sum_{i=1}^{m} \left(Z_{t;i}^{\top} \nabla f(x_t)\right)^2 < \frac{1}{2}\|\nabla f(x_t)\|^2 \text{ for every } t \in [t_0; t_0 + k)\right] \leq \frac{\eta}{T};$$

which completes the proof. □

The term $t_f(\eta)$ will frequently appear in the proofs to come; in the sequel we denote

$$t_f(\eta) := \frac{c_2}{m}\log\frac{T}{\eta}; \qquad \eta \in (0; 1=e]; \tag{24}$$

where $c_2 \geq 1$ is the absolute constant defined in Lemma 14.

We next show that with high probability, the norm difference term $\|\nabla f(x_{t+1}) - \nabla f(x_t)\|$ can be bounded in terms of $\|\nabla f(x_t)\|$ and the perturbation terms $\frac{u}{2m}\sum_{i=1}^{m} Z_{t;i} Z_{t;i}^{\top} H_{t;i} Z_{t;i}$ as well as $\|Y_t\|$.

Lemma 15. Define

$$A_t(\eta) := \left\{\|\nabla f(x_{t+1}) - \nabla f(x_t)\| \leq \frac{\|\nabla f(x_t)\|}{8t_f(\eta)} + L\left(\frac{u}{2m}\sum_{i=1}^{m} Z_{t;i} Z_{t;i}^{\top} H_{t;i} Z_{t;i} + \|Y_t\|\right)\right\} \tag{25}$$

where $t_f(\eta)$ is defined in Eq.(24), and let $C_1 \geq 1$ be the corresponding absolute constants defined in Lemma 1. Then there exists an absolute constant $c_4 > 0$ such that, whenever $\eta$ satisfies

$$L \leq \frac{c_4(\ln(C_1 dmT=\eta))^{3=2}\sqrt{d}}{\sqrt{m}} \cdot \frac{1}{8t_f(\eta)}; \tag{26}$$

we have

$$P(A_t(\eta)) \geq 1 - \frac{\eta}{T}$$

for any $\eta \in (0; 1=e]$ and $t \in \mathbb{Z}^+$.

Proof. Since $\nabla f$ is $L$-Lipschitz, following the zeroth-order update step, we see that

$$\|\nabla f(x_{t+1}) - \nabla f(x_t)\| \leq L\|x_{t+1} - x_t\| \tag{27}$$

$$= L\left\|\frac{1}{m}\sum_{i=1}^{m} Z_{t;i} Z_{t;i}^{\top} \nabla f(x_t) + \frac{u}{2m}\sum_{i=1}^{m} Z_{t;i} Z_{t;i}^{\top} H_{t;i} Z_{t;i} + Y_t\right\|: \tag{28}$$

Now, it follows from Eq. (18) (with a slight modification in the absolute constant terms since here the norm is not squared) that there exists some absolute constant $c_4 > 0$ such that for any $\delta \in (0, 1/e]$, we have that with probability at least $1 - \delta/T$, the event

$$\frac{1}{m}\sum_{i=1}^{m} z_{t;i}\, z_{t;i}^{\top}\, \nabla f(x_t) - \nabla f(x_t) \ge -c_4 (\ln(C_1 dmT/\delta))^{3/2}\sqrt{\frac{d}{m}}\|\nabla f(x_t)\|;$$

Hence, continuing from Eq. (28), it follows that with probability at least $1 - \delta/T$,

$$\|\nabla f(x_{t+1}) - \nabla f(x_t)\|$$
$$\le \eta\left( c_4 (\ln(C_1 dmT/\delta))^{3/2}\sqrt{\frac{d}{m}}\|\nabla f(x_t)\| + \left\| \frac{\mu}{2m}\sum_{i=1}^{m} z_{t;i}\, z_{t;i}^{\top}\, H_{t;i}\, z_{t;i} \right\| + \|Y_t\| \right);$$

and by plugging in the condition Eq. (26), we see that the event

$$A_t(\delta) = \left\{ \|\nabla f(x_{t+1}) - \nabla f(x_t)\| \le \frac{\|\nabla f(x_t)\|}{8 t_f(\delta)} + \eta\left( \left\| \frac{\mu}{2m}\sum_{i=1}^{m} z_{t;i}\, z_{t;i}^{\top}\, H_{t;i}\, z_{t;i} \right\| + \|Y_t\| \right) \right\}$$

has probability at least $1 - \delta/T$.                                                    □

We show now that if the norm of the gradient dominates the norm of the perturbation terms, and we choose the step-size sufficiently small, then in a small number of iterations, the norm of the gradient does not change very much. For notational simplicity, we denote the event

$$E(t_1; t_2; \delta) := \bigcap_{t=t_1}^{t_1+t_2-1}\left\{ \|\nabla f(x_t)\| > 8 t_f(\delta)\, \eta\left( \left\| \frac{\mu}{2}\frac{1}{m}\sum_{i=1}^{m} z_{t;i}\, z_{t;i}^{\top}\, H_{t;i}\, z_{t;i} \right\| + \|Y_t\| \right) \right\}:$$

**Lemma 16.** Let $\delta \in (0, 1/e]$ and $T \in \mathbb{Z}^+$ be such that $T > 2 t_f(\delta) + 1$. Consider any positive integer $t_f^0 \le 2 t_f(\delta)$, and any $t_0 \in \{0, \ldots, T - 1 - t_f^0\}$. Suppose $\eta$ satisfies the condition Eq. (26). Then, on the event

$$E(t_0; t_f^0; \delta) \setminus \left(\bigcap_{t=t_0}^{t_0+t_f^0-1} A_t(\delta)\right);$$

we have

$$\frac{1}{2}\|\nabla f(x_0)\| \le \|\nabla f(x_t)\| \le 2\|\nabla f(x_0)\|$$

for all $t \in \{t_0, \ldots, t_0 + t_f^0 - 1\}$.

Proof. By plugging

$$\|\nabla f(x_t)\| > 8 t_f(\delta)\, \eta\left( \left\| \frac{\mu}{2}\frac{1}{m}\sum_{i=1}^{m} z_{t;i}\, z_{t;i}^{\top}\, H_{t;i}\, z_{t;i} \right\| + \|Y_t\| \right)$$

into the definition of $A_t(\delta)$, we see that, on the event $E(t_0; t_f^0; \delta) \setminus \left(\bigcap_{t=t_0}^{t_0+t_f^0-1} A_t(\delta)\right)$, we have

$$\|\nabla f(x_{t+1}) - \nabla f(x_t)\| \le \frac{\|\nabla f(x_t)\|}{4 t_f(\delta)};$$

and consequently,

$$\left(1 - \frac{1}{4 t_f(\delta)}\right)\|\nabla f(x_t)\| \le \|\nabla f(x_{t+1})\| \le \left(1 + \frac{1}{4 t_f(\delta)}\right)\|\nabla f(x_t)\|;$$

which leads to

$$\left(1 - \frac{1}{4 t_f(\delta)}\right)^{t-t_0}\|\nabla f(x_0)\| \le \|\nabla f(x_t)\| \le \left(1 + \frac{1}{4 t_f(\delta)}\right)^{t-t_0}\|\nabla f(x_0)\|$$

for all $t \in \{t_0, \ldots, t_0 + t_f^0\}$. Then, since $(1 + 1/(4x))^{2x} \le 2$ and $(1 - 1/(4x))^{2x} \ge 1/2$ for any $x \ge 1$, noting that $t_f^0 \le 2 t_f(\delta)$, we get the desired result.                                                    □

Conversely, in the following result, we show that in a small number of consecutive iterations, if the gradient is smaller than the perturbation terms in any one of the iterations, then for each of the iterations in this range, the gradient will be small and be on the same scale as the size of the perturbation terms.

Lemma 17. Let $\epsilon \in (0; 1/e]$ and $T \in Z^+$ be such that $T > 2t_f(\epsilon) + 1$. Consider any positive integer $t_f^0 \leq 2t_f(\epsilon)$, and any $t_0 \in \{0; \ldots; T - 1 - t_f^0\}$. Suppose satisfies the condition Eq. (26). Then, on the event

$$
E^c(t_0; t_f^0; \epsilon) \setminus \left( \bigcap_{t=t_0}^{t_0+t_f^0-1} A_t(\epsilon) \right) \setminus \left( \bigcap_{t=t_0}^{t_0+t_f^0-1} G_t(\epsilon) \right);
$$

we have

$$
\| \nabla f(x_t) \| \leq c_5 t_f(\epsilon) L \left( u^2 d^2 \left( \log \frac{T}{\epsilon} \right)^2 + \sqrt{1 + \frac{\log(T=\epsilon)}{d}} r \right) \quad 8t \in \{t_0; t_0 + 1; \ldots; t_0 + t_f^0 - 1\};
$$

where $c_5$ is some absolute constant.

Proof. Let $t^0$ be the first iteration in $\{t_0; t_0 + 1; \ldots; t_0 + t_f^0 - 1\}$ such that

$$
\| \nabla f(x_{t^0}) \| \leq 8t_f(\epsilon) L \left( \frac{u}{2} \frac{1}{m} \sum_{i=1}^{m} Z_{t^0;i} Z_{t^0;i}^> H_{t^0;i} Z_{t^0;i} + \| Y_{t^0} \| \right): \tag{29}
$$

Since we are working on an event which is a subset of $E^c(t_0; t_f^0; \epsilon)$, $t^0$ is well-defined. By $\| H_{t^0;i} \| \leq u \| Z_{t^0;i} \|$, we see that

$$
\| \nabla f(x_{t^0}) \| \leq 8t_f(\epsilon) L \left( \frac{u^2}{2m} \sum_{i=1}^{m} \| Z_{t^0;i} \|^4 + \| Y_{t^0} \| \right)
$$

$$
\leq 8t_f(\epsilon) L \left( c_3 u^2 d^2 \left( \log \frac{T}{\epsilon} \right)^2 + c_3 \sqrt{1 + \frac{\log(T=\epsilon)}{d}} r \right);
$$

where we used the definition of $G_t(\epsilon)$.

Recall that $t^0$ is the first time step such that Eq. (29) holds. By deriving similarly as in the proof of Lemma 16, we can show that for any $j \in \{t_0; t_0 + 1; \ldots; t^0 - 1\}$,

$$
\| \nabla f(x_j) \| \leq 2 \| \nabla f(x_{t^0}) \| \leq 16 t_f(\epsilon) L c_3 \left( u^2 d^2 \left( \log \frac{T}{\epsilon} \right)^2 + \sqrt{1 + \frac{\log(T=\epsilon)}{d}} r \right):
$$

Meanwhile, for iterations $t \in [t^0; t_0 + t_f^0)$, by using the definitions of $A_t(\epsilon)$ and $G_t(\epsilon)$, we have

$$
\| \nabla f(x_{t+1}) \| \leq \left( 1 + \frac{1}{8t_f(\epsilon)} \right) \| \nabla f(x_t) \| + L c_3 \left( u^2 d^2 \left( \log \frac{T}{\epsilon} \right)^2 + \sqrt{1 + \frac{\log(T=\epsilon)}{d}} r \right)
$$

$$
= \left( 1 + \frac{1}{8t_f(\epsilon)} \right)^{t+1-t^0} \| \nabla f(x_{t^0}) \|
$$

$$
+ \sum_{i=0}^{t-t^0} \left( 1 + \frac{1}{8t_f(\epsilon)} \right)^{t-t^0-i} L c_3 \left( u^2 d^2 \left( \log \frac{T}{\epsilon} \right)^2 + \sqrt{1 + \frac{\log(T=\epsilon)}{d}} r \right)
$$

$$
\leq \left( 1 + \frac{1}{8t_f(\epsilon)} \right)^{t_f^0} \| \nabla f(x_{t^0}) \|
$$

$$
+ 8t_f(\epsilon) \left( \left( 1 + \frac{1}{8t_f(\epsilon)} \right)^{t_f^0} - 1 \right) L c_3 \left( u^2 d^2 \left( \log \frac{T}{\epsilon} \right)^2 + \sqrt{1 + \frac{\log(T=\epsilon)}{d}} r \right)
$$

$$e^{1/4} \cdot 8t_f(\eta) L c_3 \eta^2 u^2 d^2 \left( \log \frac{T}{\eta} \right)^2 + \left( 1 + \frac{\log(T/\delta)}{d} \right)^r$$

$$+ 8t_f(\eta)(e^{1/4} - 1) L c_3 \eta^2 u^2 d^2 \left( \log \frac{T}{\eta} \right)^2 + \left( 1 + \frac{\log(T/\delta)}{d} \right)^r$$

$$\leq 16 t_f(\eta) L c_3 \eta^2 u^2 d^2 \left( \log \frac{T}{\eta} \right)^2 + \left( 1 + \frac{\log(T/\delta)}{d} \right)^r ;$$

where we used $\eta^0 \leq 2 t_f(\eta)$ and the fact that $(1 + 1/(8x))^{2x} \leq e^{1/4}$ for all $x > 0$. By defining $c_5 := 16 c_3$, we complete the proof. $\square$

We next derive a useful result showing that the function change $f(x_\tau) - f(x_0)$ can be decomposed into one component arising from intervals when the gradient dominates noise (which improves function value) and another component arising from intervals with small gradient which may add to function value but whose contributions are bounded in terms of $\eta$ and $r$. For now, we focus on the case $\tau \geq t_f(\eta)$, since it will be useful to us in proving that there cannot be more than $T/4$ iterations with large gradient.

**Lemma 18** (Function change for large $\tau$). *Let $c_1 > 0; c_4 > 0; c_5 > 0; C_1 \geq 1$ be the absolute constants defined in the statements of the previous lemmas. Let $\eta \in (0; 1/e]$, and let $\tau \geq t_f(\eta))$ be arbitrary. Consider splitting $\{0; 1; \ldots; \tau - 1\}$ into $K := \lfloor \tau = t_f(\eta) \rfloor$ intervals:*

$$J_k = \{ k t_f(\eta); \ldots; (k + 1) t_f(\eta) - 1 \}; \quad 0 \leq k < K - 1;$$
$$J_{K-1} = \{ (K-1) t_f(\eta); \ldots; \tau - 1 \}:$$

*Let $I_1$ denote the set of indices $k$ such that for every time-step $t$ in the interval $J_k$, the gradient dominates the noise terms as*

$$\| r f(x_t) \| > 8 t_f(\eta) L \left( \frac{u}{2} + \frac{1}{m} \sum_{i=1}^m Z_{t;i} Z_{t;i}^\top + \| H_{t;i} Z_{t;i} \| + \| Y_t \| \right) : \tag{30}$$

*Suppose we choose $\eta$ such that*

$$\frac{1}{L t_f(\eta)} \geq \min \left( \frac{p\bar{m}}{8 c_4 (\log(C_1 d m T/\delta))^{3=2} p\bar{d}} ; \frac{m}{128 c_1 (\log(C_1 d m T/\delta))^3 d} \right) : \tag{31}$$

*Then, on the event*

$$E(\eta) := H(\eta) \setminus \left( \bigcup_{t=0}^{\tau-1} A_t(\eta) \right) \setminus \left( \bigcup_{t=0}^{\tau-1} G_t(\eta) \right) \setminus \left( \bigcup_{k=0}^{K-2} B_{k t_f(\eta)}(\eta; t_f(\eta)) \right) \setminus B_{(K-1) t_f(\eta)}(\eta; \tau - (K-1) t_f(\eta));$$

*we have the following upper bound on function value change:*

$$f(x_\tau) - f(x_0) \leq -\frac{\eta}{2} \sum_{k \in I_1} \min_{t \in J_k} \| r f(x_t) \|^2 + \frac{c_5^2}{64} \eta^3 t_f(\eta)^2 L^2 u^2 d^2 \left( \left( \log \frac{T}{\eta} \right)^2 + \left( \frac{p}{2 \log(T/\delta)} r \right)^{!2} \right)$$

$$+ \eta u^4 \eta^2 c_1 d^3 \left( \log \frac{T}{\eta} \right)^3 + L^2 u^4 \eta^2 c_1 d^4 \left( \log \frac{T}{\eta} \right)^4$$

$$+ c_1 r^2 (128 t_f(\eta) + \eta L) \log \frac{T}{\eta} + c_1 L \eta^2 r^2 : \tag{32}$$

*Moreover, $P(E(\eta)) \geq 1 - \frac{(5\tau + 4)\delta}{T}$.*

**Proof.** Without loss of generality, we may assume that $\tau$ is a multiple of $t_f(\eta)$.[5] Then, any interval $J_k = \{ t_0; \ldots; t_0 + t_f(\eta) - 1 \}$ belongs to one of the following two cases:

---
[5]To accommodate the last interval which has length at most $2 t_f(\eta) - 1$, we note that the results we require for the proof, namely Lemma 14, Lemma 16 and Lemma 17, all hold for any interval length $\eta^0 \leq 2 t_f(\eta)$.

Case 1) (Gradient dominates noise): Recall that this means that for every $t \ge 2 \in J_k$, we have

$$\|\nabla f(x_t)\| > 8t_f(\epsilon)L\left(\frac{u}{2} \cdot \frac{1}{m}\sum_{i=1}^{m} Z_{t;i}Z_{t;i}^\top H_{t;i}Z_{t;i} + \|Y_t\|\right):$$

By our choice of $\gamma$ in Eq. (31), we can apply Lemma 16 to get

$$\min_{t \ge 2 \in J_k}\|\nabla f(x_t)\| \ge \frac{1}{4}\max_{t \ge 2 \in J_k}\|\nabla f(x_t)\|:$$

We now consider the two cases when $J_k$ has fewer than $t_f(\epsilon)$ iterations and when $d = J_k$. f Note also that on the event $B_{kt_f(\epsilon)}(\epsilon; t_f(\epsilon))$, there exists some $t \ge 2 \in J_k$ such that

$$\frac{1}{m}\sum_{i=1}^{m} Z_{t;i}^\top \nabla f(x_t)^2 \ge \frac{1}{2}\|\nabla f(x_t)\|^2:$$

This implies then that

$$\frac{1}{4}\sum_{t \ge 2 \in J_k}\frac{1}{m}\sum_{i=1}^{m} Z_{t;i}^\top \nabla f(x_t)^2 \ge \frac{1}{4}\min_{t \ge 2 \in J_k}\|\nabla f(x_t)\|^2 \ge \frac{1}{64}\max_{t \ge 2 \in J_k}\|\nabla f(x_t)\|^2$$

$$\ge \frac{1}{64t_f(\epsilon)}\sum_{t \ge 2 \in J_k}\|\nabla f(x_t)\|^2: \tag{33}$$

Thus by setting $\lambda = 128t_f(\epsilon)$ in Eq. (5) and by choosing $m$ such that

$$\frac{c_1 L^2 \epsilon^3 d}{m} \le \lambda - \gamma = \frac{1}{128t_f(\epsilon)} \implies m \ge \frac{m}{128c_1 L t_f(\epsilon)d^3};$$

it follows that

$$\frac{3}{4}\sum_{t \ge 2 \in J_k}\frac{1}{m}\sum_{i=1}^{m} Z_{t;i}^\top \nabla f(x_t)^2 + \left(\frac{1}{128t_f(\epsilon)} + \frac{c_1 L^2 \epsilon^3 d}{m}\right)\sum_{t \ge 2 \in J_k}\|\nabla f(x_t)\|^2$$

$$= \frac{3}{4}\sum_{t \ge 2 \in J_k}\frac{1}{m}\sum_{i=1}^{m} Z_{t;i}^\top \nabla f(x_t)^2 + \frac{1}{64t_f(\epsilon)}\sum_{t \ge 2 \in J_k}\|\nabla f(x_t)\|^2$$

$$\ge \frac{1}{2}\sum_{t \ge 2 \in J_k}\frac{1}{m}\sum_{i=1}^{m} Z_{t;i}^\top \nabla f(x_t)^2$$

$$\ge \frac{1}{2}\min_{t \ge 2 \in J_k}\|\nabla f(x_t)\|^2 \tag{34}$$

Case 2) (Gradient does not dominate noise): there exists some $t \ge 2 \in J_k$ such that

$$\|\nabla f(x_t)\| \le 8t_f(\epsilon)L\left(\frac{u}{2} \cdot \frac{1}{m}\sum_{i=1}^{m} Z_{t;i}Z_{t;i}^\top H_{t;i}Z_{t;i} + \|Y_t\|\right):$$

By our choice of $\gamma$ in Eq. (31), we can apply Lemma 17 to get

$$\|\nabla f(x_t)\| \le c_5 t_f(\epsilon)L\left(u^2 d^2\left(\log\frac{T}{\epsilon}\right)^2 + \left(1 + \frac{\log(T=\epsilon)}{d}\right)r^r\right) \qquad 8t \ge 2 \in J_k:$$

Hence, by setting $\lambda = 128t_f(\epsilon)$ in Eq. (5) and choosing $m$ such that

$$\frac{c_1 L^2 \epsilon^3 d}{m} \le \lambda - \gamma = \frac{1}{128t_f(\epsilon)};$$

it follows that

$$\frac{}{128 t_f(\ )} + \frac{c_1 L^2 \ ^3 d}{m} \sum_{t 2 J_k} kr\ f\ (x_t)k^2$$

$$\frac{}{64 t_f(\ )} \sum_{t 2 J_k} c_5 t_f(\ )\ L\ u^2 d^2\ \left(\log \frac{T}{} \right)^2 + \left(1 + \frac{\log(T =\ )}{d}\right)^r r^2$$

$$\frac{c_5^2}{64} t_f(\ )^2\ ^3 L^2\ u^2 d^2\ \left(\log \frac{T}{}\right)^2 + \left(1 + \frac{\log(T =\ )}{d}\right)^r r^2 \tag{35}$$

Without loss of generality, we may assume that is a multiple of $t_f(\ )$.[6] Then, any interval $J_k = f t_0; \dots; t_0 + t_f(\ )\ 1g$ belongs to one of the following two cases:

Having studied the two cases, we may now proceed to use them to complete the proof. Let $I_1^c$ denote the complement of $I_1$ in $f 0; 1; \dots; K\ 1g$. Then,

$$\frac{3}{4} \sum_{t=0}^{1} \frac{1}{m} \sum_{i=1}^{m} Z_{t;i}^{\geq} r\ f\ (x_t)^2 + \ + \frac{c_1 L^2\ ^3 d}{m} \sum_{t=0}^{1} kr\ f\ (x_t)k^2$$

$$= \sum_{k 2 I_1} \left( \frac{3}{4} \sum_{t=2 J_k} \frac{1}{m} \sum_{i=1}^{m} Z_{t;i}^{\geq} r\ f\ (x_t)^2 + \frac{}{128 t_f(\ )} + \frac{c_1 L^2\ ^3 d}{m} \sum_{t 2 J_k} kr\ f\ (x_t)k^2 \right)$$

$$+ \sum_{k 2 I_1^c} \left( \frac{3}{4} \sum_{t 2 J_k} \frac{1}{m} \sum_{i=1}^{m} Z_{t;i}^{\geq} r\ f\ (x_t)^2 + \frac{}{128 t_f(\ )} + \frac{c_1 L^2\ ^3 d}{m} \sum_{t 2 J_k} kr\ f\ (x_t)k^2 \right)$$

$$\sum_{k 2 I_1} \frac{}{2} \min_{t 2 J_k} kr\ f\ (x_t)k^2 + \sum_{k 2 I_1^c} t_f(\ ) @ \frac{c_5^2}{64} t_f(\ )^2\ ^3 L^2\ u^2 d^2\ \left(\log \frac{T}{}\right)^2 + \left(1 + \frac{\log(T =\ )}{d}\right)^r r^2 \Big]^1 A$$

$$\sum_{k 2 I_1} \frac{}{2} \min_{t 2 J_k} kr\ f\ (x_t)k^2 + \frac{c_5^2}{64} t_f(\ )^2\ ^3 L^2\ u^2 d^2\ \left(\log \frac{T}{}\right)^2 + \left(1 + \frac{\log(T =\ )}{d}\right)^r r^2 : \tag{36}$$

and so by Eq. (5),

$$f\ (x)\ \ f\ (x_0) \sum_{k 2 I_1} \frac{}{2} \min_{t 2 J_k} kr\ f\ (x_t)k^2 + \frac{c_5^2}{64} t_f(\ )^2\ ^3 L^2\ u^2 d^2\ \left(\log \frac{T}{}\right)^2 + \left(1 + \frac{\log(T =\ )}{d}\right)^r r^2$$

$$+\ u^4\ ^2\ c_1 d^3\ \left(\log \frac{T}{}\right)^3 + L^2 u^4\ ^2\ c_1 d^4\ \left(\log \frac{T}{}\right)^4$$

$$+\ c_1 r^2(\ + L)\ \log \frac{T}{} + c_1 L^2 r^2:$$

Note that we choose $= 128 t_f(\ )$. In addition, observe that by our choice of (such that $\frac{1}{e}$), it follows that $q \frac{}{1 + \frac{\log(T =\ )}{d}} p \frac{}{2 \log(T =\ )}$.

We can now complete our proof by using the union bound (suppressing the dependence of some of the events on for notational simplicity) to derive

$$P(E^c)\ \ P(H^c) + \sum_{t=0}^{1} P(A_t^c) + \sum_{t=0}^{1} P(G_t^c) + \sum_{k=0}^{K\ 1} P(B_{k t_f(\ )}^c(\ ; t_f(\ )))$$

$$\frac{(\ +4)}{T} + \frac{}{T} + 2\frac{}{T} + \frac{K}{T}\ \frac{(5\ +4)}{T} : \qquad \square$$

---

[6]To accommodate the last interval which has length at most $2 t_f(\ )\ 1$, we note that the results we require for the proof, namely Lemma 14, Lemma 16 and Lemma 17, all hold for any interval length $2 t_f(\ )$.

We are now ready to show that if sufficiently many iterations have a large gradient, then with high probability, the function value of the last iterate $f(x_T)$, will be less than $\min_x f(x)$, a contradiction. Hence this limits the number of iterations that can have a large gradient.

**Proposition 4.** Let $c_1 > 0; c_2 \geq 1; c_4 > 0; c_5 > 0; C_1 \geq 1$ be the absolute constants defined in the statements of the previous lemmas, and let $\epsilon \in (0; 1=e]$ be arbitrary. Suppose we choose $u, r, \sigma$ and $T$ such that

$$u \leq \frac{\epsilon^{p-\delta}}{d^{p}\sqrt{\log(T=\delta)}} \min\left\{\frac{1}{64c_5^2 c_2}; \frac{1}{2048c_1 c_2}\right\}^{1=4}; \quad r \leq \min\left\{\frac{1}{8c_5}\epsilon^{p}\frac{\delta}{2c_2}; \frac{1}{32}\epsilon^{p}\frac{\delta}{c_1}\right\};$$

$$\frac{1}{Lt_f(\delta)} \min\left\{\frac{1}{\log(T=\delta)}; \frac{1}{8c_4(\ln(C_1 dmT=\delta))^{3=2}}\epsilon^{p}\frac{\delta}{d}; \frac{m}{128c_1(\ln(C_1 dmT=\delta))^3 d}\right\};$$

$$T \geq \max\left\{\frac{256 t_f(\delta)(f(x_0) - f^* + \sigma^2=L)}{\epsilon^2}; 4\right\}:$$

Then, with probability at least $1 - 6\delta$, there are at most $T=4$ iterations for which $\|\nabla f(x_t)\| \geq \epsilon$.

**Proof.** Without loss of generality, we assume that $T$ is a multiple of $t_f(\delta)$, and we similarly split $\{0; 1; \ldots; T\}$ into $K = \lfloor T=t_f(\delta)\rfloor$ intervals $J_0; \ldots; J_{K-1}$. Let $I_1$ denote the set of indices $k$ such that for every $t \in J_k$,

$$\|\nabla f(x_t)\| > 8t_f(\delta) L\left[\frac{u}{2} + \frac{1}{m}\sum_{i=1}^{m}\left\|Z_{t;i} Z_{t;i}^{\geq} H_{t;i} Z_{t;i}\right\| + \|Y_t\|\right]: \tag{37}$$

We let $I_1^c$ denote the complement of $I_1$ in $\{0; 1; \ldots; K-1\}$. We denote

$$E_T(\delta) := H_T(\delta) \setminus \left(\bigcap_{t=0}^{T-1} A_t(\delta)\right) \setminus \left(\bigcap_{t=0}^{T-1} G_t(\delta)\right) \setminus \left(\bigcap_{k=0}^{K-1} B_{kt_f(\delta)}(\delta; t_f(\delta))\right):$$

In the remaining part of the proof, unless otherwise stated, we shall always assume that we are working on the event $E_T(\delta)$.

By Lemma 18 with $\tau = T$ and our choices of $\sigma$ and $\delta$ in the statement of the lemma, we have

$$f(x_T) - f(x_0) \leq -\sum_{k\in I_1}\frac{\tau}{2}\min_{t\in J_k}\|\nabla f(x_t)\|^2 + T\frac{c_5^2}{64}t_f(\delta)^2 \epsilon^3 L^2 u^2 d^2\left(\log\frac{T}{\delta}\right)^2 + \epsilon^{p}\frac{\delta}{2\log(T=\delta)}r^{!2}$$

$$+ T\epsilon u^4 \sigma^2 c_1 d^3\left(\log\frac{T}{\delta}\right)^3 + TL^2 u^4 \sigma^2 c_1 d^4\left(\log\frac{T}{\delta}\right)^4$$

$$+ \epsilon c_1 r^2(128 t_f(\delta) + \tau L)\log\frac{T}{\delta} + Tc_1 L^2 \sigma^2 r^2: \tag{38}$$

Suppose that there are at least $T=4$ iterations where $\|\nabla f(x_t)\| \geq \epsilon$. Let $I$ denote the set of indices $k$ for which there exists some $t \in J_k$ with $\|\nabla f(x_t)\| \geq \epsilon$. Then, by the pigeonhole principle, the set $I$ has at least $\lceil T=(4t_f(\delta))\rceil$ members. Note that, by our choices of the parameters $u; \sigma; r$, it can be shown that

$$c_5 t_f(\delta) L\left[u^2 d^2\left(\log\frac{T}{\delta}\right)^2 + \left(1 + \sqrt{\frac{\log(T=\delta)}{d}}\right)r^{!}\right] < \epsilon; \tag{39}$$

while by Lemma 17, if $k$ is in $I_1^c$, we have

$$\|\nabla f(x_t)\| \leq c_5 t_f(\delta) L\left[u^2 d^2 \log(T=\delta) + \left(1 + \sqrt{\frac{\log(T=\delta)}{d}}\right)r^{!}\right]; \quad \forall t \in J_k:$$

This implies that $I \subseteq I_1$.

Observe that by Lemma 16, for any $k \in I_1$, we have

$$\frac{1}{2}\|\nabla f(x_{kt_f(\delta)})\| \leq \|\nabla f(x_t)\| \leq 2\|\nabla f(x_{kt_f(\delta)})\|; \quad \forall t \in J_k:$$

34

This implies in particular that for any $k \geq 1$, we have $\min_{t \in \mathcal{J}_k} \|\nabla f(x_t)\|^2 \geq \frac{1}{16}\epsilon^2$, and consequently

$$\sum_{k \geq I_1} \frac{\tau}{2} \min_{t \in \mathcal{J}_k} \|\nabla f(x_t)\|^2 \geq \sum_{k \geq I} \frac{\tau}{2} \cdot \frac{\epsilon^2}{16} \geq \frac{T}{4t_f(\epsilon)} \cdot \frac{\tau}{2} \cdot \frac{\epsilon^2}{16} = \frac{T\tau\epsilon^2}{128 t_f(\epsilon)}.$$

Hence, by Eq. (38),

$$f(x_T) - f(x_0) \leq -\frac{T\tau\epsilon^2}{128 t_f(\epsilon)} + T\frac{c_5^2}{64}t_f(\epsilon)^2\mu^3 L^2\tau u^2 d^2\left(\log\frac{T}{\delta}\right)^2 + \frac{p}{2\log(T/\delta)}r\right)^2$$

$$+ T\mu u^4\tau^2 c_1 d^3\left(\log\frac{T}{\delta}\right)^3 + T(\mu L)u^4\tau^2 c_1 d^4\left(\log\frac{T}{\delta}\right)^4$$

$$+ c_1 r^2(128 t_f(\epsilon) + \mu L)\tau\log\frac{T}{\delta} + T\mu c_1 L\tau^2 r^2. \tag{40}$$

Now, by our choices of $u$, $r$ and $\mu$, we have

$$T\frac{c_5^2}{64}t_f(\epsilon)^2\mu^3 L^2\tau u^2 d^2\left(\left(\log\frac{T}{\delta}\right)^2 + \frac{p}{2\log(T/\delta)}r\right)^2$$

$$\leq T\left(\frac{c_5^2}{32}t_f(\epsilon)^2(\mu L)^2\tau u^4 d^4\mu^2\left(\log\frac{T}{\delta}\right)^4 + 2\log(T/\delta)r^2\right)$$

$$\leq T\left(\frac{\epsilon^2}{2048 c_2\tau\left(\log\frac{T}{\delta}\right)^2} + \frac{\epsilon^2}{2048 c_2\log(T/\delta)}\right) \leq \frac{T\tau\epsilon^2}{512 t_f(\epsilon)};$$

where we used $\log(T/\delta) \geq 1$ and $2c_2\log(T/\delta) \leq t_f(\epsilon)$. We also have

$$T\mu u^4\tau^2 c_1 d^3\left(\log\frac{T}{\delta}\right)^3 + T(\mu L)u^4\tau^2 c_1 d^4\left(\log\frac{T}{\delta}\right)^4 + T\mu c_1 L\tau^2 r^2$$

$$\leq T\frac{\epsilon^2}{2048 c_2 d\log(T/\delta)} + T\frac{\epsilon^2}{2048 c_2 t_f(\epsilon)\log(T/\delta)} + T\frac{\epsilon^2}{1024 t_f(\epsilon)\log(T/\delta)}$$

$$\leq \frac{T\tau\epsilon^2}{512 t_f(\epsilon)};$$

where we used $c_2 d\log(T/\delta) \geq t_f(\epsilon)$, $c_2 \geq 1$ and $\log(T/\delta) \geq 1$. Finally,

$$c_1 r^2(128 t_f(\epsilon) + \mu L)\tau\log\frac{T}{\delta} \leq \frac{(128 t_f(\epsilon) + 1)\tau\epsilon^2}{1024 L t_f(\epsilon)} < \frac{\epsilon^2}{L}.$$

By plugging these bounds into Eq. (40), we get

$$f(x_T) - f(x_0) < -\frac{T\tau\epsilon^2}{128 t_f(\epsilon)} + \frac{T\tau\epsilon^2}{512 t_f(\epsilon)} + \frac{T\tau\epsilon^2}{512 t_f(\epsilon)} + \frac{\epsilon^2}{L} \leq -\frac{T\tau\epsilon^2}{256 t_f(\epsilon)} + \frac{\epsilon^2}{L}.$$

Therefore, as long as

$$T \geq \frac{256 t_f(\epsilon)\left((f(x_0) - f^\star) + \epsilon^2/L\right)}{\tau\epsilon^2};$$

we will get $f(x_T) < f^\star$, which is a contradiction. Thus, we can conclude that on the event $E(\eta)$, there are at most $T/4$ iterations for which $\|\nabla f(x_t)\| \geq \epsilon$.

We can now complete our proof by using the union bound (suppressing the dependence of some of the events on $\eta$ for notational simplicity) to derive

$$P(E_T^c) \leq P(H_T^c) + \sum_{t=0}^{T-1}P(A_t^c) + \sum_{t=0}^{T-1}P(G_t^c) + \sum_{k=0}^{K-1}P(B_{kt_f(\epsilon)}^c(\eta; t_f(\epsilon)))$$

$$\leq \frac{(T+4)\delta}{T} + \delta + 2\delta + \frac{K\delta}{T} \leq 6\delta. \qquad \square$$

35

# E. Escaping saddle point

In this section, we rst show that the travelling distance of the iterates can be bounded in terms of the function value improvement (Appendix E.2). Utilizing this result, as well as Proposition 2 in Appendix C.3 which provides a concentration bound on the the zeroth-order noise, we then prove that suf cient function value decrease can be made near a saddle point in Appendix E.3.

## E.1. Key quantities and notation

We will use $\gamma$ to denote $\lambda_{\min}(r^2 f(x_0))$, where we know that $\gamma \geq \sqrt{\rho \epsilon}$.

## E.2. Improve or Localize

In this subsection, we aim to bound the movement of the iterates across a number of steps in terms of the function value improvement made during these number of steps.

We rst state a simple result separating the norm of the difference between $x_{t_0+\tau}$ and $x_{t_0}$ into a few different terms.

Lemma 19. Consider the perturbed zeroth-order update Algorithm 1. Then, for any $t_0 \in N$ and $\tau \in N$,

$$\|x_{t_0+\tau} - x_{t_0}\|^2 \leq V_1(t_0; \tau) + V_2(t_0; \tau) + V_3(t_0; \tau) + V_4(t_0; \tau); \tag{41}$$

where

$$V_1(t_0; \tau) := 8\eta^2 \sum_{t=t_0}^{t_0+\tau-1} \tau \|r f(x_t)\|^2; \quad V_2(t_0; \tau) := 8\eta^2 \sum_{t=t_0}^{t_0+\tau-1} \tau \left\| \frac{1}{m} \sum_{i=1}^m (Z_{t;i} Z_{t;i}^\top - I) r f(x_t) \right\|^2$$
$$V_3(t_0; \tau) := 4\eta^2 \sum_{t=t_0}^{t_0+\tau-1} \tau \|Y_t\|^2; \quad V_4(t_0; \tau) := 4\eta^2 \sum_{t=t_0}^{t_0+\tau-1} \tau \left\| \frac{1}{m} \sum_{i=1}^m u Z_{t;i} Z_{t;i}^\top H_{t;i} Z_{t;i} \right\|^2: \tag{42}$$

Proof. For notational convenience, let $t_0 := 0$. Then, applying the form of the perturbed zeroth-order update in Algorithm 1, we get

$$\|x_\tau - x_0\|^2$$
$$= \left\| \sum_{t=0}^{\tau-1} x_{t+1} - x_t \right\|^2$$
$$= \eta^2 \left\| \sum_{t=0}^{\tau-1} \frac{1}{m} \sum_{i=1}^m Z_{t;i} Z_{t;i}^\top r f(x_t) + \frac{1}{m} \sum_{i=1}^m u Z_{t;i} Z_{t;i}^\top H_{t;i} Z_{t;i} + Y_t \right\|^2$$
$$\leq 4\eta^2 \left\| \sum_{t=0}^{\tau-1} \frac{1}{m} \sum_{i=1}^m Z_{t;i} Z_{t;i}^\top r f(x_t) \right\|^2 + 4\eta^2 \left\| \sum_{t=0}^{\tau-1} \frac{1}{m} \sum_{i=1}^m u Z_{t;i} Z_{t;i}^\top H_{t;i} Z_{t;i} \right\|^2 + 4\eta^2 \left\| \sum_{t=0}^{\tau-1} Y_t \right\|^2$$
$$\leq 4\eta^2 \left\| \sum_{t=0}^{\tau-1} \frac{1}{m} \sum_{i=1}^m (Z_{t;i} Z_{t;i}^\top - I) r f(x_t) + \sum_{t=0}^{\tau-1} r f(x_t) \right\|^2 + 4\eta^2 \left\| \sum_{t=0}^{\tau-1} \frac{1}{m} \sum_{i=1}^m u Z_{t;i} Z_{t;i}^\top H_{t;i} Z_{t;i} \right\|^2 + 4\eta^2 \left\| \sum_{t=0}^{\tau-1} Y_t \right\|^2$$
$$\leq \underbrace{8\eta^2 \sum_{t=0}^{\tau-1} \tau \|r f(x_t)\|^2}_{V_1(0; \tau)} + \underbrace{8\eta^2 \sum_{t=0}^{\tau-1} \tau \left\| \frac{1}{m} \sum_{i=1}^m (Z_{t;i} Z_{t;i}^\top - I) r f(x_t) \right\|^2}_{V_2(0; \tau)} + \underbrace{4\eta^2 \left\| \sum_{t=0}^{\tau-1} Y_t \right\|^2}_{V_3(0; \tau)} + \underbrace{4\eta^2 \sum_{t=0}^{\tau-1} \tau \left\| \frac{1}{m} \sum_{i=1}^m u Z_{t;i} Z_{t;i}^\top H_{t;i} Z_{t;i} \right\|^2}_{V_4(0; \tau)}:$$

□

We now proceed to bound the terms $V_1(t_0; \tau); V_2(t_0; \tau); V_3(t_0; \tau)$ and $V_4(t_0; \tau)$.

First, we have the following result bounding $V_1(t_0; \tau)$.

Lemma 20. Let $c_1 > 0; c_2 \geq 1; c_4 > 0; c_5 > 0; C_1 \geq 1$ be the absolute constants de ned in the statements of the previous lemmas, and let $\iota \in (0; 1=e]$ be arbitrary.

Suppose we choose $\epsilon$ such that

$$\frac{1}{Lt_f(\epsilon)} \le \min\left\{\frac{p^{\frac{m}{d}}}{8c_4(lr(C_1dmT\epsilon))^{3-2p}}, \frac{m}{128c_1(lr(C_1dmT\epsilon))^3d}\right\}:$$

There are two cases to consider.

1. The first is when $\tau \ge t_f(\epsilon)$. In this case, split $\{t_0; t_0 + 1; \ldots; t_0 + \tau - 1\}$ into $K := \lfloor \tau = t_f(\epsilon)\rfloor$ intervals:

$$J_k = \{t_0 + kt_f(\epsilon); \ldots; t_0 + (k+1)t_f(\epsilon) - 1\}; \quad 0 \le k < K - 1;$$
$$J_{K-1} = \{t_0 + (K-1)t_f(\epsilon); \ldots; t_0 + \tau - 1\}:$$

Then, on the event

$$E_{t_0;\tau}(\epsilon) := H_{t_0;\tau}(\epsilon) \setminus \left(\bigcup_{t=t_0}^{t_0+\tau-1} A_t(\epsilon)\right) \setminus \left(\bigcup_{t=t_0}^{t_0+\tau-1} G_t(\epsilon)\right) \setminus \left(\bigcup_{k=0}^{K-2} B_{t_0+kt_f(\epsilon)}(\epsilon; t_f(\epsilon))\right) \setminus B_{t_0+(K-1)t_f(\epsilon)}(\epsilon; \tau - (K-1)t_f(\epsilon));$$

we have that

$$V_1(t_0; \tau) \ge 8\epsilon^2 \sum_{t=t_0}^{t_0+\tau-1} kr f(x_t)k^2$$
$$- 64\epsilon t_f(\epsilon)((f(x_0) - f(x_\tau)) + N_{u;r}(\epsilon; \tau));$$

where

$$N_{u;r}(\epsilon; \tau) := \frac{c_5^2}{64}\epsilon^3 t_f(\epsilon)^2 L^2 \tau u^2 d^2 \left(\log\frac{T}{\epsilon}\right)^2 + p\frac{1}{2\log(T=\epsilon)r}^{!2}$$
$$+ \sigma u^{4}\epsilon^2 c_1 d^3\left(\log\frac{T}{\epsilon}\right)^3 + L^2 u^4 \epsilon^2 c_1 d^4\left(\log\frac{T}{\epsilon}\right)^4$$
$$+ \sigma c_1 r^2(128 t_f(\epsilon) + L\epsilon)\log\frac{T}{\epsilon} + c_1 L^2 r^2$$
$$+ c_5^2 t_f^3(\epsilon)\epsilon^3 L^2 \tau u^2 d^2 \left(\log(T=\epsilon) + p\frac{1}{2\log(T=\epsilon)r}\right)^2: \tag{43}$$

2. The second is when $\tau < t_f(\epsilon)$. Suppose we choose $u$ and $r$ such that

$$u \le \frac{p\bar\epsilon}{d^p\overline{\log(T=\epsilon)}}\min\left\{\frac{1}{64c_5^2c_2}; \frac{1}{2048c_1c_2}\right\}^{1=4}; \quad r \le \min\left\{\frac{1}{8c_5}p\frac{\bar\epsilon}{2c_2}; \frac{1}{32}p\frac{\bar\epsilon}{c_1}\right\}:$$

Suppose the event $\bigcup_{t=t_0}^{t_0+\tau-1}(A_t(\epsilon) \setminus G_t(\epsilon))$ holds. Suppose also that $kr f(x_{t_0})k \ge \epsilon$: Then,

$$V_1(t_0; \tau) \ge 32\sigma^2\epsilon^2\tau^2 - 32\sigma^2(t_f(\epsilon))^2\epsilon^2$$

Proof. 1. We first consider the case where $\tau \ge t_f(\epsilon)$. Let $I_1$ denote the set of indices $k$ such that for every time-step $t$ in the interval $J_k$, the gradient dominates the noise terms as

$$kr f(x_t)k > 8t_f(\epsilon)L\left(\frac{u}{2} + \frac{1}{m}\sum_{i=1}^{m}Z_{t;i}Z_{t;i}^{\ge} H_{t;i}Z_{t;i}\right) + kY_tk: \tag{44}$$

WLOG, we may assume that $t_0 := 0$, and denote $V_1(\tau) := V_1(0; \tau)$. WLOG, we also assume that $\tau$ is a multiple of $t_f(\epsilon)$. From Lemma 18, on the event that $E_\tau(\epsilon)$ holds and by our choice of $\epsilon$, we have

$$f(x_\tau) - f(x_0) \le -\sum_{k\ge I_1}\frac{1}{2}\min_{t\ge J_k}kr f(x_t)k^2\epsilon + \frac{c_5^2}{64}\epsilon^3 t_f(\epsilon)^2 L^2 \tau u^2 d^2\left(\log\frac{T}{\epsilon}\right)^2 + p\frac{1}{2\log(T=\epsilon)r}^{!2}$$

37

$$+ \quad u^4 \eta^2 \ c_1 d^3 \ \left(\log \frac{T}{\delta}\right)^3 + \ L^2 u^4 \eta^2 \ c_1 d^4 \ \left(\log \frac{T}{\delta}\right)^4$$

$$+ \ c_1 r^2 (128 t_f(\epsilon) + \ L) \log \frac{T}{\delta} + \ c_1 L^2 r^2:$$

By Lemma 16 (and our choice of $\eta$), it follows that for any $k \in I_1$, on the event $\bigcap_{t \in J_k} A_t(\epsilon)$, we have

$$\sum_{t \in J_k} \|\nabla f(x_t)\|^2 \ \geq \ 4 t_f \min_{t \in J_k} \|\nabla f(x_t)\|^2:$$

Thus, on the event that $E(\epsilon)$ holds, for our choice of $\eta$, we have

$$\sum_{k \in I_1} \sum_{t \in J_k} \|\nabla f(x_t)\|^2 \ \geq \ 4 t_f(\epsilon) \sum_{k \in I_1} \min_{t \in J_k} \|\nabla f(x_t)\|^2$$

$$\geq 8 t_f(\epsilon) \sum_{k \in I_1} \frac{1}{2} \min_{t \in J_k} \|\nabla f(x_t)\|^2$$

$$\geq 8 t_f(\epsilon) @ (f(x_0) - f(x_*)) + \frac{c_5^2}{64} \eta^3 t_f(\epsilon)^2 L^2 \ u^2 d^2 \ \left(\log \frac{T}{\delta}\right)^2 + \ \left(p \frac{1}{2 \log(T/\delta)} r\right)^2 A$$

$$+ \ 8 t_f(\epsilon) \left( u^4 \eta^2 \ c_1 d^3 \ \left(\log \frac{T}{\delta}\right)^3 + \ L^2 u^4 \eta^2 \ c_1 d^4 \ \left(\log \frac{T}{\delta}\right)^4 \right)$$

$$+ \ 8 t_f(\epsilon) \left( c_1 r^2 (128 t_f(\epsilon) + \ L) \log \frac{T}{\delta} + \ c_1 L^2 r^2 \right):$$

Similarly, for any $k \in I_1^c$ (where $I_1^c$ denotes the complement of $I_1$ in $\{0, 1, \ldots, K - 1\}$, i.e. intervals where the gradient is smaller than than the perturbation terms in some iteration), on the event $(\bigcap_{t \in J_k} A_t(\epsilon)) \setminus (\bigcap_{t \in J_k} G_t(\epsilon))$, by Lemma 17 (and our choice of $\eta$), we have

$$\|\nabla f(x_t)\| \ \leq \ c_5 t_f(\epsilon) L \ u^2 d^2 \ \log(T/\delta) + \ p \frac{1}{2 \log(T/\delta)} r \ ; \qquad 8 t \in J_k:$$

On the event that $E(\epsilon)$ holds, this gives us then

$$\sum_{k \in I_1^c} \sum_{t \in J_k} \|\nabla f(x_t)\|^2 \ \leq \ c_5^2 t_f^2(\epsilon) \eta^2 L^2 \ u^2 d^2 \ \log(T/\delta) + \ \left(p \frac{1}{2 \log(T/\delta)} r\right)^2 :$$

Hence, on the event that $E(\epsilon)$ holds, we have that

$$\sum_{t=0}^{T-1} \|\nabla f(x_t)\|^2 = \ \sum_{k \in I_1} \sum_{t \in J_k} \|\nabla f(x_t)\|^2 + \ \sum_{k \in I_1^c} \sum_{t \in J_k} \|\nabla f(x_t)\|^2$$

$$\geq 8 t_f(\epsilon) @ (f(x_0) - f(x_*)) + \frac{c_5^2}{64} \eta^3 t_f(\epsilon)^2 L^2 \ u^2 d^2 \ \left(\log \frac{T}{\delta}\right)^2 + \ \left(p \frac{1}{2 \log(T/\delta)} r\right)^2 A$$

$$+ \ 8 t_f(\epsilon) \left( u^4 \eta^2 \ c_1 d^3 \ \left(\log \frac{T}{\delta}\right)^3 + \ L^2 u^4 \eta^2 \ c_1 d^4 \ \left(\log \frac{T}{\delta}\right)^4 \right)$$

$$+ \ 8 t_f(\epsilon) \left( c_1 r^2 (128 t_f(\epsilon) + \ L) \log \frac{T}{\delta} + \ c_1 L^2 r^2 \right)$$

$$+ \ 8 t_f(\epsilon) \left( c_5^2 t_f^2(\epsilon) \eta^2 L^2 \ u^2 d^2 \ \log(T/\delta) + \ \left(p \frac{1}{2 \log(T/\delta)} r\right)^2 \right):$$

This yields the final result for the case $\tau \leq t_f(\epsilon)$.

38

2. We next consider the case where $\tau < t_f(\epsilon)$. Recall the notation that

$$E(t_0; t_0 + \tau; \epsilon) := \left\{ \bigcup_{t=t_0}^{t_0+\tau-1} \left( \|r f(x_t)\| > 8t_f(\epsilon) L \left( \frac{u}{2} \frac{1}{m} \sum_{i=1}^m Z_{t;i} Z_{t;i}^\top H_{t;i} Z_{t;i} \right) + \|Y_t\| \right) \right\}$$

There are two cases to consider.

(a) On the event $E(t_0; t_0 + \tau; \epsilon) \setminus \bigcup_{t=t_0}^{t_0+\tau-1} A_t(\epsilon)$; we have by Lemma 16 that $\|r f(x_t)\| \le 2\|r f(x_0)\|$ for each $t \in \{0, 1, \dots, \tau - 1\}$. Then,

$$V_1(t_0; \tau) = 8\tau^2 \sum_{t=t_0}^{t_0+\tau-1} \|r f(x_t)\|^2 \le 8\tau^2 \cdot \tau \cdot 4\|r f(x_0)\|^2 \le 32\tau^2 \cdot \tau \cdot \epsilon^2;$$

where the final inequality uses the assumption that $\|r f(x_0)\| \le \epsilon$:

(b) Suppose the event $E^c(t_0; t_0 + \tau; \epsilon) \setminus \bigcup_{t=t_0}^{t_0+\tau-1} A_t(\epsilon) \setminus \bigcup_{t=t_0}^{t_0+\tau-1} G_t(\epsilon)$ holds. In this case, by Lemma 17, we have that for each $t \in \{t_0; t_0 + 1; \dots; t_0 + \tau - 1\}$

$$\|r f(x_t)\| \le c_5 t_f(\epsilon) L \left( u^2 d^2 \left( \log \frac{T}{\delta} \right)^2 + \sqrt{r} \left( 1 + \frac{\log(T=\delta)}{d} r \right) \right)!$$

;

where the final inequality follows by our choice of $u$ and $r$ (cf. Eq. (39)). Hence,

$$V_1(t_0; \tau) = 8\tau^2 \sum_{t=t_0}^{t_0+\tau-1} \|r f(x_t)\|^2$$

$$\le 8\tau^2 \cdot \tau \left( c_5 t_f(\epsilon) L \left( u^2 d^2 \left( \log \frac{T}{\delta} \right)^2 + \sqrt{r} \left( 1 + \frac{\log(T=\delta)}{d} r \right) \right)!! \right)^2$$

$$\le 8\tau^2 \cdot \tau \cdot \epsilon^2 < 32\tau^2 \cdot \tau \cdot \epsilon^2:$$

The final result for the case $\tau < t_f(\epsilon)$ then follows.

$\square$

We proceed to bound $V_2(t_0; \tau)$.

**Lemma 21.** Let $c_1 > 0; c_2 \ge 1; c_4 > 0; c_5 > 0; C_1 \ge 1$ be the absolute constants defined in the statements of the previous lemmas, and let $\rho \in (0; 1=e]$ be arbitrary and $\epsilon > 0$ be arbitrary. Suppose we choose $\eta$ such that

$$\frac{1}{L t_f(\epsilon)} \ge \eta \ge \min \left\{ \frac{\rho \sqrt{m}}{8 c_4 (\ln(C_1 d m T = \rho))^{3=2}} \frac{\epsilon}{d}; \frac{m}{128 c_1 (\ln(C_1 d m T = \rho))^3 d} \right\}:$$

Let $T_s$ denote an integer such that $T_s \ge \max\{\tau; t_f(\epsilon)\}$, and for any $F > 0$, define

$$B(\tau; F) := \frac{8 t_f(\epsilon)(F + N_{u;r}(T_s; \delta))}{\eta} \left( T_s + \frac{d}{m} (\ln(C T^2 = \delta))^2 \right); \qquad b(\tau; F) := \frac{t_f(\epsilon) F}{\eta}:$$

Let $c^0; C > 0$ denote the same constants as in the statement of Proposition 2. Denote the event that

either $\sum_{t=t_0}^{t_0+\tau-1} \frac{d}{m} (\ln(C T^2 = \delta))^2 \|r f(x_t)\|^2 \le B(\tau; F)$

or $\frac{V_2(t_0; \tau)}{8\tau^2} \le c^0 \sqrt{\max \left( \sum_{t=t_0}^{t_0+\tau-1} \frac{d}{m} (\ln(C T^2 = \delta))^2 \|r f(x_t)\|^2; b(\tau; F) \right) \left( \log \frac{C T^2}{\delta} + \log \left( \log \frac{B(\tau; F)}{b(\tau; F)} + 1 \right) \right)}$

holds as $L_{t_0;\delta}(\epsilon;F)$[7]. We show that $P(L_{t_0;\delta}(\epsilon;F)) \geq 1 - \frac{\delta}{T}$: Finally, denote the event $M_{t_0;T_s}(F)$ as the event that $f(x_{t_0}) - f(x_{t_0+T_s}) < F$.

Then, on the event $L_{t_0;\delta}(\epsilon) \setminus E_{t_0;T_s}(\epsilon) \setminus M_{t_0;T_s}(F)$ (where $E_{0;T_s}(\epsilon)$ is as defined in Lemma 20),

$$V_2(t_0;\epsilon) \geq 8c^{\alpha 2}\gamma_1(\epsilon;F)t_f(\epsilon)\max\left\{\frac{8d}{m}(\text{lr}(CT^2=\epsilon))^2(F + N_{u;r}(T_s;\delta)); F\right\};\qquad(45)$$

where

$$\gamma_1(\epsilon;F) := \log\left(\frac{CT^2}{\epsilon}\right) + \log\left(\log\left(\frac{B(\epsilon;F)}{b_1(\epsilon;F)}\right) + 1\right):$$

Proof. We note that $P(L_{t_0;\delta}(\epsilon;F)) \geq 1 - \frac{\delta}{T}$: is a direct consequence of Proposition 2. In the rest of the proof, without loss of generality, we assume that $t_0 = 0$ for notational simplicity. On the event $L_{t_0;\delta}(\epsilon;F) \setminus E_{0;T_s}(\epsilon) \setminus M_{t_0;T_s}(F)$, suppose that

$$\sum_{t=0}^{X-1}\frac{d}{m}(\text{lr}(CT^2=\epsilon))^2 kr\|\nabla f(x_t)\|^2 \geq B(\epsilon;F) = \frac{8t_f(\epsilon)(F + N_{u;r}(T_s;\delta))}{\gamma T_s} + \frac{d}{m}\gamma(\text{lr}(CT^2=\epsilon))^2$$

$$=)\quad \sum_{t=0}^{X-1} kr\|\nabla f(x_t)\|^2 \geq 8t_f(\epsilon)(F + N_{u;r}(T_s;\delta))$$

$$=)\quad \sum_{t=0}^{TX-1} kr\|\nabla f(x_t)\|^2 \geq 8t_f(\epsilon)(F + N_{u;r}(T_s;\delta))$$

$$=)\quad 8\gamma^2 T_s\sum_{t=0}^{TX-1} kr\|\nabla f(x_t)\|^2 \geq 64\gamma T_s t_f(\epsilon)(F + N_{u;r}(T_s;\delta))$$

$$=)\quad 8\gamma^2 T_s\sum_{t=0}^{TX-1} kr\|\nabla f(x_t)\|^2 \geq 64\gamma T_s t_f(\epsilon)(f(x_0) - f(x_{T_s}) + N_{u;r}(T_s;\delta));\quad \text{since} f(x_0) - f(x_{T_s}) \geq F$$

$$()\quad V_1(0;T_s) \geq 64\gamma T_s t_f(\epsilon)(f(x_0) - f(x_{T_s}) + N_{u;r}(T_s;\delta));$$

where we note the last equation contradicts Lemma 20. For notational simplicity, denote

$$\gamma(\epsilon;F) := \log\left(\frac{CT^2}{\epsilon}\right) + \log\left(\log\left(\frac{B(\epsilon;F)}{b(\epsilon;F)}\right) + 1\right):$$

Observe that $\gamma_1$ is larger than $\gamma$ for every $\epsilon \geq 1$. Since $L_{t_0;\delta}(\epsilon;F)$ holds, we must have then that

$$s\frac{V_2(0;\epsilon)}{8\gamma^2} \geq c^{\alpha}\max\left(\sqrt{\sum_{t=0}^{X-1}\frac{d}{m}(\text{lr}(CT^2=\epsilon))^2 kr\|\nabla f(x_t)\|^2; b(\epsilon;F)}\right)\gamma_1(\epsilon;F):$$

Now, continuing, recalling the definition of $V_1(0;T_s) = 8\gamma^2 T_s\sum_{t=0}^{T_s-1}kr\|\nabla f(x_t)\|^2$

$$V_2(0;\epsilon) \geq c^{\alpha 2}\gamma_1(\epsilon;F)\max\left(8\gamma^2\sum_{t=0}^{X-1}\frac{d}{m}(\text{lr}(CT^2=\epsilon))^2 kr\|\nabla f(x_t)\|^2; 8\gamma^2 b(\epsilon;F)\right)$$

$$\geq c^{\alpha 2}\gamma_1(\epsilon;F)\max\left(8\gamma^2\sum_{t=0}^{TX-1}\frac{d}{m}(\text{lr}(CT^2=\epsilon))^2 kr\|\nabla f(x_t)\|^2; 8\gamma^2 b(\epsilon;F)\right)$$

$$\geq c^{\alpha 2}\gamma_1(\epsilon;F)\max\left(\frac{d}{m}(\text{lr}(CT^2=\epsilon))^2\frac{V_1(0;T_s)}{T_s}; 8t_f(\epsilon)F\right)$$

---

[7]We note that by construction $B(\epsilon;F) \geq b(\epsilon;F)$

$$\overset{(i)}{\leq} c^{02} \, \beta_1(\delta; F) \max_\gamma \left[ \frac{d}{m} (\ln(CT^2/\delta))^2 (64 \, t_f(\gamma)(f(x_0) - f(x_{T_s}) + N_{u;r}(T_s; \delta))); 8 \, t_f(\gamma) F \right]$$

$$\overset{(ii)}{\leq} c^{02} \, \beta_1(\delta; F) \max_\gamma \left[ \frac{d}{m} (\ln(CT^2/\delta))^2 (64 \, t_f(\gamma)(F + N_{u;r}(T_s; \delta))); 8 \, t_f(\gamma) F \right]$$

$$= c^{02} \, \beta_1(\delta; F)(8 \, t_f(\gamma)) \max_\gamma \left[ \frac{d}{m} (\ln(CT^2/\delta))^2 (8(F + N_{u;r}(T_s; \delta))); F \right].$$

We note that (i) is a consequence of Lemma 20, while (ii) comes from our assumption that the event $M_{ey;r_s}(F)$ holds, i.e. $f(x_{t_0}) - f(x_{t_0 + T_s}) \leq F$.

□

We next bound $V_3(t_0; \delta)$ and $V_4(t_0; \delta)$.

Lemma 22. Let $c > 0$ denote the same constant in Lemma 7. Consider any arbitrary $0 \leq \delta \leq 1 = e$; and let $t_f(\gamma)$ be arbitrary. Let $N_{t_0; \gamma}(\delta)$ denote the event that

$$V_3(t_0; \delta) := 4\gamma^2 \sum_{t=t_0}^{t_0 \vee \tau - 1} \|Y_t\|^2 \leq 4c_6 \gamma^2 \log(2dT/\delta) r^2;$$

where $c_6 > 0$ is an absolute constant. Then, by Lemma 7, $P(N_{t_0; \gamma}(\delta)) \geq 1 - \frac{\delta}{T}$. Denote the event

$$O_t(\delta) := \left( \frac{1}{m} \sum_{i=1}^{m} \|Z_{t;i}\|^8 \leq c_7 d^4 \log^4 \left( \frac{T}{\delta} \right) \right);$$

where $c_7 > 0$ is an absolute constant. Then, on the event $\cap_{t=t_0}^{t_0 \vee \tau - 1} O_t(\delta)$, we have

$$V_4(t_0; \delta) \leq 4c_7 \gamma^2 \beta^2 L^2 u^4 d^4 \log^4 \left( \frac{T}{\delta} \right).$$

Moreover, for each $t$, $P(O_t(\delta)) \geq 1 - \frac{\delta}{T}$.

Proof. The proof for $V_3(t_0; \delta)$ follows directly from Lemma 7, by picking $c_6$ to be the $c$ that appears in the statement of Lemma 7. Meanwhile, observe that

$$V_4(t_0; \delta) = 4\gamma^2 \sum_{t=t_0}^{t_0 \vee \tau - 1} \left\| \frac{1}{m} \sum_{i=1}^{m} u Z_{t;i} Z_{t;i}^\geq H_{t;i} Z_{t;i} \right\|^2$$

$$\leq 4\gamma^2 @ \sum_{t=t_0}^{t_0 \vee \tau - 1} \left( \frac{1}{m} \sum_{i=1}^{m} \|u Z_{t;i} Z_{t;i}^\geq H_{t;i} Z_{t;i}\| \right)^2 A$$

$$\overset{(iii)}{\leq} 4\gamma^2 \sum_{t=t_0}^{t_0 \vee \tau - 1} \frac{1}{m} \sum_{i=1}^{m} \beta^2 u^4 L \|Z_{t;i}\|^8$$

$$\leq 4c_7 \gamma^2 \beta^2 L^2 u^4 d^4 \log^4 \left( \frac{T}{\delta} \right).$$

Above, to derive (iii), we used the bound that $\|H_{t;i}\| \leq u L \|Z_{t;i}\|$. The final inequality is a consequence of our assumption that $\cap_{t=t_0}^{t_0+\tau-1} O_t(\delta)$ holds. Finally, the result that $P(O_t(\delta)) \geq 1 - \frac{\delta}{T}$ holds due to Lemma 11, where we note that we may pick the absolute constant $c_7$ to be equal to $2Cc^4$, where $c; C > 0$ are the absolute constants that appear in the statement of Lemma 11.

□

Finally, combining the earlier results, we have the following technical result, which bounds the travelling distance of the iterates in terms of the decrease in function value decrease.

41

**Lemma 23** (Improve or Localize) Consider the perturbed zeroth-order update Algorithm 1. Let $c_0 > 0; c_1 > 0; c_2 \geq 1; c_4 > 0; c_5 > 0; c_6 > 0; c_7 > 0; C_1 \geq 1$ be the absolute constants defined in the statements of the previous lemmas, and let $\delta \in (0, 1/e]$ be arbitrary. Consider any $T_s \leq t_f(\delta)$. For any $\mathscr{F} > 0$, suppose $f(x_{T_s}) - f(x_0) > -\mathscr{F}$; i.e. $f(x_0) - f(x_{T_s}) < \mathscr{F}$. Suppose that the event

$$P_{t_0;T_s}(\delta;\mathscr{F}) := \bigcap_{\tau=1}^{T_s} \left( L_{t_0;\tau}(\delta;\mathscr{F}) \cap N_{t_0;\tau}(\delta) \right) \cap \bigcap_{t=t_0}^{t_0+T_s-1} O_t(\delta) \cap A_t(\delta) \cap G_t(\delta) \cap \bigcap_{\tau=t_f(\delta)}^{T_s-1} E_{t_0;\tau}(\delta)$$

holds, where the events $E_{t_0;\tau}(\delta); L_{t_0;\tau}(\delta); N_{t_0;\tau}(\delta); O_t(\delta)$ are as defined in Lemma 20, Lemma 21 and Lemma 22, and $G_t(\delta)$ and $A_t(\delta)$ are as defined in Lemma 13 and Lemma 15. Suppose we choose $u$ and $r$ such that

$$u \leq \frac{\delta}{d} \sqrt{\frac{\rho}{\rho \log(T/\delta)}} \min\left\{ \frac{1}{64 c_5^2 c_2}; \frac{1}{2048 c_1 c_2} \right\}^{1/4}; \quad r \leq \min\left\{ \frac{1}{8 c_5} \sqrt{\frac{\rho}{2 c_2}}; \frac{1}{32} \sqrt{\frac{\rho}{c_1}} \right\};$$

$$\frac{1}{L t_f(\delta)} \min\left\{ \frac{1}{\log(T/\delta)}; \frac{\sqrt{\frac{\rho}{m}}}{8 c_4 (\mathrm{lr}(C_1 dmT/\delta))^{3/2}} \sqrt{\frac{\rho}{d}}; \frac{m}{128 c_1 (\mathrm{lr}(C_1 dmT/\delta))^3 d} \right\};$$

Suppose $\delta \leq \min\left\{ 1; \frac{1}{t_f(\delta)}; \frac{1}{t_f L} \right\}$. Suppose also we pick $u$ and $r$ small enough such that

$$u \leq \frac{r^{1/2}}{d \log(T/\delta)^{1/2}}; \quad r^2 \leq \min\left\{ \frac{\mathscr{F}}{T_s \log(T/\delta) \left( \frac{65 c_5^2}{8} + 132 c_1 + 1 \right)}; \frac{\mathscr{F}}{4 c_6 \log(2dT/\delta) + 4 c_7 T_s} \right\}^{9/8};$$

Then, for each $\tau \in \{0, 1, \ldots, T_s\}$, we have that

$$\|x_{t_0+\tau} - x_{t_0}\|^2 \leq \mathscr{S}_{T_s}(\delta;\mathscr{F});$$

where

$$\mathscr{S}_{T_s}(\delta;\mathscr{F}) \leq \max\left\{ 128 T_s t_f(\delta)\mathscr{F}; 32\sigma^2 (t_f(\delta))^2\eta^2 + 8 c_0 \sigma^2 \Delta_1(\delta;\mathscr{F}) t_f(\delta) \max\left\{ \frac{16d}{m}(\mathrm{lr}(CT^2/\delta))^2\mathscr{F}; T_s\mathscr{F} \right\} \right\} + T_s t_f(\delta)\mathscr{F};$$

where $\Delta_1(\delta;\mathscr{F})$ is defined as in Lemma 21. Moreover $\mathbb{P}(P_{t_0;T_s}(\delta;\mathscr{F})) \geq 1 - \frac{12 T_s}{T}$.

Proof. We recall that

$$\|x_{t_0+\tau} - x_{t_0}\|^2 \leq 8\eta^2 \underbrace{\left\| \sum_{t=t_0}^{t_0+\tau-1} \eta r \nabla f(x_t) \right\|^2}_{V_1(t_0;\tau)} + 8\eta^2 \underbrace{\left\| \sum_{t=t_0}^{t_0+\tau-1} \frac{1}{m}\sum_{i=1}^{m}(Z_{t;i} Z_{t;i}^\top - I)r \nabla f(x_t) \right\|^2}_{V_2(t_0;\tau)}$$

$$+ 4\eta^2 \underbrace{\left\| \sum_{t=t_0}^{t_0+\tau-1} Y_t \right\|^2}_{V_3(t_0;\tau)} + 4\eta^2 \underbrace{\left\| \sum_{t=t_0}^{t_0+\tau-1} \frac{1}{m}\sum_{i=1}^{m} u Z_{t;i} Z_{t;i}^\top H_{t;i} Z_{t;i} \right\|^2}_{V_4(t_0;\tau)}$$

By Lemma 20, Lemma 21, and Lemma 22, which bound $V_1(t_0;\tau); V_2(t_0;\tau);$ and $V_3(t_0;\tau); V_4(t_0;\tau)$ respectively, on the event $P_{t_0;T_s}(\delta;\mathscr{F})$, we have, for any $0 \leq \tau \leq T_s$,

$$\|x_\tau - x_0\|^2 \leq V_1(0;\tau) + V_2(0;\tau) + V_3(0;\tau) + V_4(0;\tau)$$
$$\leq \max\left\{ 64\eta t_f(\delta)(\mathscr{F} + N_{u;r}(\delta;\tau)); 32\sigma^2(t_f(\delta))^2\eta^2 \right.$$
$$+ 8 c_0 \sigma^2 \Delta_1(\delta;\mathscr{F}) t_f(\delta) \max\left\{ \frac{8d}{m}(\mathrm{lr}(CT^2/\delta))^2 (\mathscr{F} + N_{u;r}(T_s;\delta)); \mathscr{F} \right\}$$
$$\left. + 4 c_6 \eta^2 \log(2dT/\delta)r^2 + 4 c_7 \eta^2 \sigma^2\eta^2 u^4 d^4 (\log(T/\delta))^4; \right.$$

42

where $N_{u;r}(\cdot;\cdot)$ is defined as in Lemma 20.

For the simplified bound (which does not contain $N_{u;r}(\cdot;\cdot)$), it remains for us to show that our choice of $u$ and $r$ ensures that $N_{u;r}(T_s;\cdot) \leq F$ and

$$4c_6 \sigma^2 T_s \log(2dT/\delta)r^2 + 4c_7 \sigma^2 T_s^2 \sigma^2 u^4 d^4 (\log(T/\delta))^4 \leq T_s t_f(\epsilon)F:$$

First, our choice of $u$ ensures that

$$u^4 d^4 \sigma^2 (\log(T/\delta))^4 \leq r^2:$$

Next, recall that

$$N_{u;r}(\cdot;\cdot) := \frac{c_5^2}{64}\sigma^3 t_f(\epsilon)^2 L^2 \sigma u^2 d^2 \left(\log\frac{T}{\delta}\right)^2 + \left(p\frac{\sigma}{2\log(T/\delta)r}\right)^2$$

$$+ \sigma u^{4}\sigma^2 c_1 d^3 \left(\log\frac{T}{\delta}\right)^3 + L\sigma^2 u^4 \sigma^2 c_1 d^4 \left(\log\frac{T}{\delta}\right)^4$$

$$+ c_1 r^2 (128 t_f(\epsilon) + L)\log\frac{T}{\delta} + c_1 L\sigma^2 r^2$$

$$+ c_5^2 t_f^3(\epsilon)\sigma^3 L^2 \sigma u^2 d^2 \log(T/\delta) + \left(p\frac{\sigma}{2\log(T/\delta)r}\right)^2:$$

Recalling our choice of $\eta$ such that

$$\eta \leq \min\left\{1;\frac{1}{t_f(\epsilon)};\frac{1}{t_f(\epsilon)L}\right\}g;$$

it follows that

$$N_{u;r}(T_s;\cdot) \leq T_s r^2 \left[\frac{8c_5^2}{64}\log(T/\delta) + 2c_1 + 2c_1 + (128c_1 + 1)\log(T/\delta) + c_1 + 8c_5^2\log(T/\delta)\right]$$

$$\leq T_s r^2 \log(T/\delta)\left[\frac{65c_5^2}{8} + 132c_1 + 1\right] \leq F;$$

where the last inequality follows choosing $r$ such that $r^2 \leq \dfrac{F}{T_s \log(T/\delta)\left[\frac{65c_5^2}{8}+132c_1+1\right]}$. Similarly, we have

$$4c_6 \sigma^2 T_s \log(2dT/\delta)r^2 + 4c_7 \sigma^2 T_s^2 \sigma^2 u^4 d^4 (\log(T/\delta))^4$$
$$\leq T_s t_f(\epsilon)\left[4c_6 \sigma \log(2dT/\delta)r^2 + 4c_7 T_s \sigma^2 u^4 d^4(\log(T/\delta))^4\right]$$
$$\leq T_s t_f(\epsilon)\left[4c_6 \sigma \log(2dT/\delta)r^2 + 4c_7 T_s r^2\right]$$

By choosing $r$ such that

$$r^2 \leq \frac{F}{4c_6 \log(2dT/\delta) + 4c_7 T_s};$$

it follows that

$$4c_6 \sigma^2 T_s \log(2dT/\delta)r^2 + 4c_7 \sigma^2 T_s^2 \sigma^2 u^4 d^4 (\log(T/\delta))^4 \leq T_s t_f(\epsilon)F;$$

as desired.

We next lower bound the probability of

$$P_{t_0;T_s}(\cdot;F) := \bigcap_{\tau=1}^{T_s}(L_{t_0;\tau}(\cdot;F))\cap N_{t_0;\tau}(\cdot)) \cap \bigcap_{t=t_0}^{t_0+T_s-1}O_t(\cdot)\cap A_t(\cdot)\cap G_t(\cdot) \cap \bigcap_{\tau=t_f(\epsilon)}^{T_s}E_{t_0;\tau}(\cdot):$$

Observe that

$$\bigcap_{\tau=t_f(\epsilon)}^{T_s}E_{t_0;\tau}(\cdot)$$

$$= \left\| \sum_{t=t_0}^{T_s} \nabla_{t_f(\cdot)} H_{t_0;}(\cdot) \right\| \left( \prod_{t=t_0}^{t_0 + 1} A_t(\cdot) \backslash G_t(\cdot) \right) \left( \prod_{k=0}^{K-2} B_{t_0 + k t_f(\cdot)}(\cdot; t_f(\cdot)) \backslash B_{t_0 + (K-1)t_f(\cdot)}(\cdot; (K-1)t_f(\cdot)) \right)$$

$$= \left\| \sum_{t=t_f(\cdot)}^{T_s} H_{t_0;}(\cdot) \right\| \left( \prod_{k=0}^{K-2} B_{t_0 + k t_f(\cdot)}(\cdot; t_f(\cdot)) \backslash B_{t_0 + (K-1)t_f(\cdot)}(\cdot; (K-1)t_f(\cdot)) \right) \left( \prod_{t=t_0}^{T_f - 1} A_t(\cdot) \backslash G_t(\cdot) \right) :$$

Note this implies that $\left\| \sum_{t=t_f(\cdot)}^{T_s} E_{t_0;}(\cdot) \right\| \left\| \sum_{t=t_0}^{T_s-1} A_t(\cdot) \backslash G_t(\cdot) \right\| = \left\| \sum_{t=t_f(\cdot)}^{T_s} E_{t_0;}(\cdot) \right\|$ We note that by Lemma 1,

$$P\left( \left\| \sum_{t=t_f(\cdot)}^{T_s} H_{t_0;}(\cdot) \right\|^c \right) \quad \frac{5T_s}{T}:$$

Meanwhile, we note that

$$\left\| \sum_{t=t_0}^{T_s-1} B_t(\cdot; t_f(\cdot)) \right\| \quad \left\| \sum_{t=t_f(\cdot)}^{T_s} \prod_{k=0}^{K-2} B_{t_0 + k t_f(\cdot)}(\cdot; t_f(\cdot)) \backslash B_{t_0 + (K-1)t_f(\cdot)}(\cdot; (K-1)t_f(\cdot)) \right\| :$$

Hence, by Lemma 14, we have that

$$P\left( \left\| \sum_{t_f(\cdot)}^{T_s} \prod_{k=0}^{K-2} B_{t_0 + k t_f(\cdot)}(\cdot; t_f(\cdot)) \backslash B_{t_0 + (K-1)t_f(\cdot)}(\cdot; (K-1)t_f(\cdot)) \right\|^c \right)$$

$$P\left( \left\| \sum_{t=t_0}^{T_s-1} B_t(\cdot; t_f(\cdot)) \right\|^c \right) \quad \frac{T_s}{T}:$$

Meanwhile, by Lemma 13 and Lemma 15, we may bound

$$P\left( \left\| \sum_{t=t_0}^{T_f-1} A_t(\cdot) \backslash G_t(\cdot) \right\|^c \right) \quad \frac{T_s}{T} + \frac{2T_s}{T} = \frac{3T_s}{T}:$$

Hence, it follows that

$$P\left( \left\| \sum_{t=t_f(\cdot)}^{T_s} E_{t_0;}(\cdot) \right\| \left\| \sum_{t=t_0}^{T_s-1} A_t(\cdot) \backslash G_t(\cdot) \right\|^c \right) \quad \frac{5T_s}{T} + \frac{T_s}{T} + \frac{3T_s}{T} = \frac{9T_s}{T}:$$

Meanwhile, it follows from our results in the preceding lemmas that

$$P\left( \left\| \sum_{=1}^{T_s} (L_{t_0;}(\cdot; F) \backslash N_{t_0;}(\cdot)) \right\| \left\| \sum_{t=t_0}^{T_s-1} O_t(\cdot) \right\|^c \right) \quad \frac{3T_s}{T}:$$

Hence, it follows that $P(P_{t_0;T_s}(\cdot; F)) \quad 1 \quad \frac{12T_s}{T}.$

$\square$

## E.3. Proving function value decrease near saddle point

We next build on the technical result earlier to prove that each time we are near the saddle point, there is a constant probability of making significant function value decrease. We briefly provide a high-level proof outline below. In our proof, we introduce a coupling argument connecting two closely-related sequences both starting from the saddle, differing only in the sign of their perturbative term along the minimum eigendirection of the Hessian at the saddle. Specifically, when function decrease from a saddle is not sufficiently large, due to the earlier technical result, we know that the coupled sequences will remain within a radius of the original saddle for a large number (which we will denote $T_s$) of iterations. We then utilize this fact to show that the difference of the coupled sequence will (with some constant probability) grow exponentially large, eventually moving out of their specified radius within $T_s$ iterations, leading to a contradiction.

Our first result formally introduces the coupling, setting the stage for the rest of our arguments. For notational convenience, in this section, unless otherwise specified, we will often assume that the initial iterate is an -saddle point.

**Lemma 3.** Suppose $x_0$ is an $\epsilon$-approximate saddle point. Without loss of generality, suppose that the minimum eigendirection of $H := \nabla^2 f(x_0)$ is the $e_1$ direction (i.e. the first basis vector in $\mathbb{R}^d$), and let $\gamma$ to denote $\lambda_{min}(\nabla^2 f(x_0))$ (note $\gamma \leq -\sqrt{\rho\epsilon}$). Consider the following coupling mechanism, where we run the zeroth-order gradient dynamics, starting at $x_0$, with two isotropic noise sequences $Y_t$ and $Y_t^0$ respectively, where $(Y_t)_1 = -(Y_t)_1^0$, and $(Y_t)_j = (Y_t)_j^0$ for all other $j \neq 1$. Suppose that the sequence $\{Z_{t;i}\}_{t \geq 2T; i \in 2[m]}$ is the same for both sequences. Let $\{x_t\}$ denote the sequence with the $Y_t$ noise sequence, and let the $\{x_t^0\}$ denote the sequence with the $Y_t^0$ noise sequence, where $x_0^0 = x_0$; and

$$x_{t+1}^0 = x_t^0 - \left( \frac{\sum_{i=1}^m Z_{t;i} Z_{t;i}^\top \nabla f(x_t^0) + \frac{u}{2} Z_{t;i} Z_{t;i}^\top H_{t;i}^0 Z_{t;i}}{m} \right) + Y_t^0;$$

and $H_{t;i}^0 := \frac{H_{t;i;+}^0 + H_{t;i;-}^0}{2}$, with $H_{t;i;+}^0 = \nabla^2 f(x_t^0 + \theta_{t;i;+}^0 uZ_i^0)$ for some $\theta_{t;i;+}^0 \in [0;1]$, and $H_{t;i;-}^0 = \nabla^2 f(x_t^0 - \theta_{t;i;-}^0 uZ_i^0)$ for some $\theta_{t;i;-}^0 \in [0;1]$. Then, for any $t \geq 0$,

$$\hat{x}_{t+1} := x_{t+1} - x_{t+1}^0$$

$$= \sum_{\tau=0}^t \underbrace{(I - \eta H)^{t-\tau} \hat{\psi}_{g_0}(\tau)}_{W_{g_0}(t+1)} - \sum_{\tau=0}^t \underbrace{(I - \eta H)^{t-\tau} (H_\tau - H)\hat{x}_\tau}_{W_H(t+1)}$$

$$- \sum_{\tau=0}^t \underbrace{(I - \eta H)^{t-\tau} \hat{\psi}_u(\tau)}_{W_u(t+1)} + \sum_{\tau=0}^t \underbrace{(I - \eta H)^{t-\tau} \hat{Y}_\tau}_{W_p(t+1)}$$

where

$$\psi_{g_0}(t) = \frac{1}{m} \sum_{i=1}^m (Z_{t;i} Z_{t;i}^\top - I)\nabla f(x_t);$$

$$\psi_{g_0}^0(t) = \frac{1}{m} \sum_{i=1}^m (Z_{t;i} (Z_{t;i})^\top - I)\nabla f(x_t^0);$$

$$\hat{\psi}_{g_0}(t) = \psi_{g_0}(t) - \psi_{g_0}^0(t); \quad \psi_u(t) = \frac{1}{m} \sum_{i=1}^m \frac{u}{2} Z_{t;i} Z_{t;i} H_{t;i} Z_{t;i};$$

$$\psi_u^0(t) = \frac{1}{m} \sum_{i=1}^m \frac{u}{2} Z_{t;i} Z_{t;i} H_{t;i}^0 Z_{t;i}; \quad \hat{\psi}_u(t) = \psi_u(t) - \psi_u^0(t);$$

$$\hat{Y}_t = Y_t - Y_t^0; \quad H_t = \int_0^1 \nabla^2 f(ax_t + (1-a)x_t^0)da.$$

**Proof.** Observe that

$$\hat{x}_{t+1} := x_{t+1} - x_{t+1}^0$$

$$= x_t - \eta(\nabla f(x_t) + \psi_{g_0}(t) + \psi_u(t)Y_t) - x_t^0 - \eta(\nabla f(x_t^0) + \psi_{g_0}^0(t) + \psi_u^0(t) + Y_t^0)$$

$$= \hat{x}_t - \eta(\nabla f(x_t) - \nabla f(x_t^0)) + \eta(\psi_{g_0}(t) - \psi_{g_0}^0(t)) + \eta(\psi_u(t) - \psi_u^0(t)) + (Y_t - Y_t^0)$$

$$= \hat{x}_t - \eta H_t \hat{x}_t - \eta(H_t - H)\hat{x}_t - \eta\hat{\psi}_{g_0}(t) - \eta\hat{\psi}_u(t) - \hat{Y}_t$$

$$= \sum_{\tau=0}^t \underbrace{(I - \eta H)^{t-\tau} \hat{\psi}_{g_0}(\tau)}_{W_{g_0}(t+1)} - \sum_{\tau=0}^t \underbrace{(I - \eta H)^{t-\tau} (H_\tau - H)\hat{x}_\tau}_{W_H(t+1)} - \sum_{\tau=0}^t \underbrace{(I - \eta H)^{t-\tau} \hat{\psi}_u(\tau)}_{W_u(t+1)} + \sum_{\tau=0}^t \underbrace{(I - \eta H)^{t-\tau} \hat{Y}_\tau}_{W_p(t+1)}$$

where

$$\psi_{g_0}(t) = \frac{1}{m} \sum_{i=1}^m (Z_{t;i} Z_{t;i}^\top - I)\nabla f(x_t); \quad \psi_{g_0}^0(t) = \frac{1}{m} \sum_{i=1}^m (Z_{t;i} (Z_{t;i})^\top - I)\nabla f(x_t^0); \quad \hat{\psi}_{g_0}(t) = \psi_{g_0}(t) - \psi_{g_0}^0(t);$$

45

$$\nabla_u(t) = \frac{1}{m} \sum_{i=1}^{X^m} \frac{u}{2} Z_{t;i} Z_{t;i} H_{t;i} Z_{t;i} ; \quad \nabla_u^0(t) = \frac{1}{m} \sum_{i=1}^{X^m} \frac{u}{2} Z_{t;i} Z_{t;i} H_{t;i}^0 Z_{t;i} ; \quad \hat\nabla_u(t) = \nabla_u(t) \quad \nabla_u^0(t);$$

$$\hat Y_t = Y_t \quad Y_t^0; \quad H_t = \int_0^1 r^2 f(a x_t + (1 \quad a) x_t^0) da:$$

To derive the final equality, we utilized the fact that $x_0^0 = x_0$. This completes our proof. $\qquad\square$

Suppose $x_0$ is an $\epsilon$-saddle point. Recall that $\gamma > 0$ denotes $\lambda_{min}(r^2 f(x_0))$, where we know that $\gamma \quad \frac{p}{}$.

$$\alpha := \begin{cases} \min f ; 1; L g & \text{if } f() \text{ is } (;; \frac{p}{})\text{-strict saddle for any} > \frac{p}{} \\ \frac{p}{} & \text{otherwise} \end{cases}$$

In the sequel, for any $\nu \quad 0$, it is helpful to define the quantities

$$\phi(t)^2 := \frac{(1 + )^{2t}}{( )^2 + 2}; \qquad \psi(t)^2 := \frac{(1 + )^{2t} \quad 1}{( )^2 + 2}: \tag{46}$$

We next introduce some probabilistic events (and their implications) which, if true, can be used to bound the sizes of $kW_{g_0}(t + 1) k; kW_u(t + 1) k, kW_u(t + 1) k$ (and as we will see in the next result, indirectly bound $kW_H(t + 1) k$. These bounds will be useful in the final proof of making function value progress near a saddle point.

**Lemma 24.** We assume $2 (0; 1=e]$ throughout the lemma. Suppose that we pick $\mu, \nu, k$ and $\gamma$ as specified in Lemma 23. Suppose $T_s \quad t_f()$. Suppose also that

$$f(x_{T_s}) \quad f(x_0) > \quad F; \quad f(x_{T_s}^0) \quad f(x_0) > \quad F:$$

Then, we have the following results.

1. Let $S()$ denote the event

$$S() := \left\{ \max f k x_t \quad x_0 k^2; k x_t^0 \quad x_0 k^2 g \quad T_s(;F); \quad 80 \quad t \quad T_s \right\}:$$

   In addition, let $S_u()$ denote the event

$$S_u() := \left\{ kW_u(t + 1) k \quad (t + 1) \overset{p\frac{3}{}}{=} 2c_3 d^2 (\log(T=))^2 u^2; \quad 80 \quad t \quad T_s \quad 1 \right\};$$

   where $c_3$ is the same absolute constant as $c_3$ in the preceding lemmas. Then,

$$P(S() \setminus S_u()) \quad 1 \quad \frac{24T_s}{T}:$$

2. Consider defining the event $R_t()$, which is the event where

   either $\sum_{=0}^{X^t} (1 + )^{2(t )} \frac{dL^2}{m} x \quad x^0{}^2 (lr(CT^2=))^2 \quad G_{T_s}(;F);$ or

$$\frac{kW_{g_0}(t + 1) k}{c^0 \sum_{=u}^{v} \max \left( lr \frac{CT^2}{}{}^2 \sum_{=0}^{X^t} \frac{dL^2}{m}(1 + )^{2(t )} k x \quad x^0 k^2; g(t + 1) \right) \left( \log \frac{CdT^2}{} + \log \log \frac{G_{T_s}(;F)}{g(t + 1)} + 1 \right)}$$

   normalsize holds. Above, $c^0; C$ refer to the same constants as in Proposition 2, and

$$G_{T_s}(;F) := 8 \sum_{=0}^{T_X^{s} 1} (1 + )^2 \frac{dL^2}{m}(lr(CT^2=))^2 \quad T_s(;F) + \frac{(T_s) r}{60} \frac{}{p\bar d}{}^2; \qquad g(t + 1) := \frac{(t + 1) r}{60} \frac{}{p\bar d}{}^2:$$

46

Then, $P(R_t(\ )) \geq 1 - \frac{\delta}{T}$. Suppose the event

$$\setminus_{t=0}^{T_s - 1} R_t(\ ) \setminus S(\ )$$

holds. Then, the event $S_{g_0}(\ )$ holds, where

$$S_{g_0}(\ ) := \setminus_{t=0}^{T_s - 1} S_{g_0;t}(\ );$$

and $S_{g_0;t}(\ )$ is defined as

$$S_{g_0;t}(\ ) := \left\{ \ kW_{g_0}(t+1)k \geq \tau_1(\ ;F)c^0 t \max \left( lr \ \frac{CT^2}{\ } , \ 2X^t \frac{dL^2}{m}(1+\ )^{2(t-\ )} kx - x^0 k^2; g(t+1) \right) \right\};$$

where

$$\tau_1(\ ;F) := \left[ \log \frac{CdT^2}{\ } + \log \left( \log \frac{G_{T_s}(\ ;F)}{g(1)} \right) + 1 \right]:$$

3. In addition, let $S_p(\ )$ denote the event

$$S_p(\ ) := \left\{ kW_p(t+1)k \geq \frac{2^p \ 2\log(T=\ ) \ (t+1) \ r}{\overline{p} \ d} \ 80 \ t \ T_s \ 1 \right\}:$$

Then, $P(S_p(\ )) \geq 1 - \frac{T_s}{T}$.

Proof. We consider the three claims separately.

1. Note that our assumptions satisfy the conditions required in Lemma 23. Hence, by Lemma 23, on the event $R_{0;T_s}(\ ;F)$, we have that $kx - x^0 k^2 \leq \tau_{T_s}(\ ;F)$. Simultaneously, on the event $R_{0;T_s}(\ ;F)$, we know that $\setminus_{t=0}^{T_s - 1} G_t(\ )$ holds, i.e.

$$\frac{1}{m} \sum_{i=1}^m kZ_{t;i} k^4 \leq 2c_3 d^2 (\log(T=\ ))^2; \quad 80 \ t \ T_s \ 1: \tag{47}$$

Thus, for $W_u(t+1)$, we have that

$$kW_u(t+1)k = \sum_{\ =0}^{X^t} (I - H_\ )^t \hat{\ }_u(\ )$$

$$\leq \sum_{\ =0}^{X^t} (I - H_\ )^t \ _u(\ ) + \sum_{\ =0}^{X^t} (I - H_\ )^t \hat{\ }_u^0(\ )$$

$$\leq \sum_{\ =0}^{X^t} (1+\ )^t \left( \frac{1}{m}\sum_{i=1}^m \frac{u}{2} Z_{t;i} Z_{t;i} H_{t;i} Z_{t;i} + \frac{1}{m}\sum_{i=1}^m \frac{u}{2} Z_{t;i} Z_{t;i} H_{t;i}^0 Z_{t;i} \right)$$

$$\leq \sum_{\ =0}^{X^t} (1+\ )^t \ \frac{\ }{m}\sum_{i=1}^m kZ_{t;i} k^4 u^2$$

$$\overset{(iv)}{\leq} \sum_{\ =0}^{X^t} (1+\ )^t \ (2c_3) d^2 (\log(T 2=\ ))^2 u^2$$

$$\leq \frac{(1+\ )^{t+1}}{\ } \ 2c_3 \ Cd^2 (\log(T=\ ))^2 \ u^2$$

$$\overset{(v)}{=} (t+1) \frac{p \ (\ )^2 + 2}{p \ } \ 2c_3 \ d^2 (\log(T=\ ))^2 \ u^2$$

$$(t+1) p = \frac{3}{\ } \ 2c_3 \ d^2 (\log(T=\ ))^2 \ u^2$$

(vi) $(t+1)\, p \overset{p}{=} \frac{3}{} \; 2c_3\, d^2 (\log(T=))^2\, u^2$

where the inequality in (iv) holds due to Eq. (47), the equality in (v) holds due to the definition of $(t+1)$, and the inequality in (vi) used the fact that .

Hence the event

$$\Big\{ \sum_{t=0}^{T_s} \|kx_t \; x_0k^2 \quad T_s(\;;F)\Big\} \text{ and } \Big\{ \backslash S_u() \Big\}$$

holds with probability at least $1 \quad \frac{12T_s}{T}$.

Note that by the coupling, the distribution of $x^0$ is the same as that of . Thus, by the assumption $n(x_{T_s}^0) \quad f(x_0) >$ F, it follows by a similar argument that the bound $kx^0 \quad x_0k^2 \quad T_s(\;;F)$ also holds with probability at least $1 \quad \frac{12T_s}{T}$. The claim then follows by an application of the union bound.

2. For the second claim, observe first that the claim $\Pr(R_t()) \quad 1 \quad \frac{}{T}$ is a consequence of Proposition 2. Suppose next that $f(x_{T_s}) \quad f(x_0) > \quad F$. Then, by definition of the event $S()$, we know that

$$kx \quad x_0k^2 \quad T_s(\;;F); \qquad kx^0 \quad x_0k^2 \quad T_s(\;;F)$$

where $T_s(\;;F)$ is as defined in Lemma 23.

Suppose now that $R_t()$ holds true, and suppose for contradiction that

$$\sum_{=0}^{X^t} (1+)^{2(t)}\frac{dL^2}{m}\|kx \quad x^0k^2(\ln(CT^2=))^2$$

$$G_{T_s}(\;;F)$$

$$= 8 \sum_{=0}^{T-1X} (1+)^2 \frac{dL^2}{m}(\ln(CT^2=))^2 \; T_s(\;;F) + \left(\frac{(T_s)\,r}{60\,d}\right)^2:$$

This implies that there exists some $0 \quad t \quad T_s$ such that $kx \quad x^0k^2 \quad 8\,T_s(\;;F)$. However, we also know that on the event $S()$,

$$kx \quad x^0k^2 \quad 2kx \quad x_0k^2 + 2kx^0 \quad x_0k^2 \quad 4\,T_s(\;;F):$$

This leads to a contradiction. We must then have that

$$kW_{g_0}(t+1)k \quad 1(\;;F)c^0\sqrt[u]{\max\left\{\ln\left(\frac{CT^2}{} \quad ^2\sum_{=0}^{X^t}\frac{dL^2}{m}(1+)^{2(t)}\right)kx \quad x^0k^2;\,g(t+1)\right\}};$$

where

$$1(\;;F) := \quad \log\left(\sqrt{\frac{CdT^2}{}}\right) + \log\left(\log\left(\frac{G(\;;F)}{g(1)}\right) + 1\right)$$

3. Observe that

$$W_p(t+1) = \sum_{=0}^{X^t}(I \quad H)^t \hat{Y} = \sum_{=0}^{X^t}(1+)^t (2(Y)_1);$$

which means that $W_p(t+1)$ is a 1-dimensional Gaussian with variance

$$^2\sum_{=0}^{X^t}(1+)^{2(t)}\frac{4r^2}{d} = \frac{4\,^2r^2}{d}\frac{(1+)^{2(t+1)} \quad 1}{2 \quad +(\ )^2} = \frac{4\,^2r^2 \quad (t+1)^2}{d}: \tag{48}$$

Since $(t+1)$ $(t+1)$, using the subGaussianity of a Gaussian distribution, it follows that for any $t$, with probability at least $1$ $=T$,

$$\|W_p(t+1)\| \quad \frac{2^p \; 2\log(T=\;) \;(t+1)\;r}{\frac{p}{d}}:$$

$\square$

For any $F > 0$, we are now ready to show that the algorithm makes a function decrease with $F$ with $(1)$ probability near an $\epsilon$-saddle point.

**Proposition 5.** Suppose that $x_{t_0}$ is an $\epsilon$-approximate saddle point. Let $c_0 > 0; c_1 > 0; c_2$ $1; c_4 > 0; c_5 > 0; c_6 > 0; c_7 > 0; C_1$ $1$ be the absolute constants defined in the statements of the previous lemmas, and let $(d+1=e)$ be arbitrary. Consider any $F > 0$. As in the statement of Lemma 23, suppose we choose $u$ and $r$ such that

$$u \quad \frac{p^{-}}{d^p \; \log(T=\;)} \quad \min \quad \frac{1}{64c_5^2 c_2}; \frac{1}{2048c_1 c_2} \quad ^{1=4}; \quad r \quad \min \quad \frac{1}{8c_5} \; p^{-}_{2c_2}; \frac{1}{32} \; p^{-}_{c_1} \; ;$$

$$\frac{1}{Lt_f(\;)} \min \quad \frac{1}{\log(T=\;)}; \frac{1}{8c_4(lr(\;C_1 dmT=\;))^{3=2}} \; p^{-}_{d}; \frac{m}{128c_1(lr(\;C_1 dmT=\;))^3 d} \; :$$

Suppose we pick

$$T_s = \max \quad d - e; t_f(\;); 4 \quad ; \tag{49}$$

where

$$= \max \quad \log \quad 2^p \quad \frac{}{T_s(\;;F\;)} \frac{20^p \; d^{-p} \; ^{2\;2+2} \; !}{r} \quad ;1 \quad ;$$

$$:= \quad \frac{\min f\;;\;1; Lg}{p\,-} \quad \begin{array}{l} \text{if } f(\;) \text{ is } (\;;\;;\; \frac{p}{-})\text{-strict saddle for any} > \frac{p}{-} \\ \text{otherwise} \end{array}$$

Suppose in addition that $u;$ also satisfy the conditions

$$u \quad \frac{s \; \frac{r}{p^{-}}}{120 \; \frac{p}{3c_3} \; p^{-}_{d} \; d^2(\log(T=\;))^2}; \quad \max \quad \frac{1}{c^0 c_9 \; _1(\;;F\;)}; \frac{m}{360\,(c^0)^2 c_9^2 dL^2 \; lr \; ^{CT^2} \; _1(\;;F\;)^2}; \frac{1}{2} \; ;$$

where $_1(\;;F\;)$ is as defined in Lemma 23; $c^0; c_3; C > 0$ are the same constants as in the previous results, and $c_9 = 2^p \; 2 + \frac{1}{20}$. Suppose also that $_{T_s}(\;;F\;)$ satisfies the bound

$$_{T_s}(\;;F\;) \quad \frac{}{60c_9 \; \log(T=\;)} \quad ^2: \tag{50}$$

Then, with probability at least $\frac{1}{3}$ $\frac{13T_s}{T}$, $f(x_{t_0+T_s})$ $f(x_{t_0})$ $F$.

**Proof of Proposition 5.** Without loss of generality, we assume that $t_0 = 0$. By Lemma 3, we have

$$\hat{x}_{t+1}$$
$$:= x_{t+1} \quad x^0_{t+1}$$
$$= \quad \underbrace{\sum_{=t_0}^{X^t} (I \quad H\;)^t \; \hat{g}_0(\;)}_{W_{g_0}(t+1)} \quad \underbrace{\sum_{=t_0}^{X^t} (I \quad H\;)^t \;(H \quad H)\hat{x}}_{W_H(t+1)} \quad \underbrace{\sum_{=t_0}^{X^t} (I \quad H\;)^t \; \hat{}_u(\;)}_{W_u(t+1)} \quad \underbrace{\sum_{=t_0}^{X^t} (I \quad H\;)^t \; \hat{Y}}_{W_p(t+1)}$$

49

where

$$\Lambda_{g_0}(t) = \frac{1}{m}\sum_{i=1}^{n}(Z_{t;i}Z_{t;i}^{\top} - I)r f(x_t); \quad \Lambda_{g_0}^{0}(t) = \frac{1}{m}\sum_{i=1}^{n}(Z_{t;i}(Z_{t;i})^{\top} - I)r f(x_t^0); \quad \hat{\Lambda}_{g_0}(t) = \Lambda_{g_0}(t) - \Lambda_{g_0}^{0}(t);$$

$$\Lambda_u(t) = \frac{1}{m}\sum_{i=1}^{n}\frac{u}{2}Z_{t;i}Z_{t;i}H_{t;i}Z_{t;i}; \quad \Lambda_u^{0}(t) = \frac{1}{m}\sum_{i=1}^{n}\frac{u}{2}Z_{t;i}Z_{t;i}H_{t;i}^{0}Z_{t;i}; \quad \hat{\Lambda}_u(t) = \Lambda_u(t) - \Lambda_u^{0}(t);$$

$$\hat{Y}_t = Y_t - Y_t^0; \quad H_t = \int_0^1 r^2 f(ax_t + (1-a)x_t^0)da:$$

Recall that we define for $t \geq 0$,

$$\sigma(t)^2 := \frac{(1+\eta)^{2t}}{(\eta)^2+2}; \qquad \rho(t)^2 := \frac{(1+\eta)^{2t}-1}{(\eta)^2+2}:$$

Throughout the proof, we suppose for contradiction that

$$f(x_{T_s}) - f(x_0) > -F; \quad f(x_{T_s}^0) - f(x_0) > -F;$$

and assume the event

$$\bigcap_{t=0}^{T_s-1}R_t(\cdot) \cap S(\cdot) \cap S_u(\cdot) \cap S_p(\cdot)$$

holds, where the events intersected are defined in Lemma 24. Then, by Lemma 24, the event $S_{g_0}(\cdot)$ (also defined in Lemma 24) holds.[8]

Consider the following induction argument, where we seek to show that there exists an absolute constant $c_9$ such that for every $t \in \{0, 1, \ldots, T_s\}$,

$$\|x_t - x_t^0\| \leq c_9 \log(T/\delta)\frac{\sigma(t)\,r}{p\sqrt{d}}; \quad \text{and} \quad \max\{\|W_{g_0}(t)\|, \|W_H(t)\|, \|W_u(t)\|\} \leq \frac{\sigma(t+1)\,r}{p\sqrt{d}} \tag{51}$$

Combined with a lower bound on $\|W_p(t+1)\|$ (which makes use of the property that $W_p(t+1)$ is a 1-dimensional Gaussian), we will then use the inductive claim in Eq. (51) to show that

$$\|W_p(T_s)\| \geq 2\left(\|W_{g_0(T_s)}\| + \|W_H(T_s)\| + \|W_u(T_s)\|\right):$$

Since $W_p(t+1)$ is a 1-dimensional Gaussian random variable with a standard deviation that grows exponentially with $t$, by our choice of $T_s$, we will see that $x_{T_s} - x_{T_s}^0$ is larger than what we expect (since our assumptions imply that $\max\{\|x_{T_s} - x_0\|^2; \|x_{T_s}^0 - x_0\|^2\} \leq \sigma_{T_s}(\cdot; F)$, i.e. $x_{T_s}$ and $x_{T_s}^0$ both remain close to $x_0$ and hence close to each other). This yields a contradiction, implying that on the event we assumed to hold, i.e.

$$\bigcap_{t=0}^{T_s-1}R_t(\cdot) \cap S(\cdot) \cap S_p(\cdot)$$

the assumption

$$f(x_{T_s}) - f(x_0) > -F; \quad \text{and} f(x_{T_s}^0) - f(x_0) > -F$$

is not true, i.e. one of the sequences must have made function value progress of $F$ at least.

We proceed to prove Eq. (51). Observe that the claim holds for the base case $t = 0$, this is true since $x_0 = x_0^0$. Now suppose that this holds for all $\leq t$. We will seek to show that Eq. (51) holds for $t+1$ as well. We do so by bounding the norms of $W_{g_0}(t+1)$; $W_H(t+1)$; $W_u(t+1)$ and $W_p(t+1)$ respectively.

---

[8] We may also directly assume that $S_{g_0}(\cdot)$ also holds, but our way of reasoning prevents double counting of probabilities.

1. (Bounding $\|W_{g_0}(t+1)\|$) Since the event $\mathcal{E}_{g_0}(\cdot)$ holds, it follows that for each $0 \le t \le T_s - 1$, we have that

$$\|W_{g_0}(t+1)\| \le \psi_1(\cdot;F)c^0 \bigvee_{\tau=0}^{t} \max\left\{ l r \frac{CT^2}{\epsilon}, \frac{2X^t}{\epsilon} \frac{dL^2}{m}(1+\gamma)^{2(t-\tau)} \|x_\tau - x^0\|^2; g(t+1)\right\}$$

where

$$\psi_1(\cdot;F) := \left\lceil \log\left(\frac{CdT^2}{\epsilon}\right) + \log\left(\log\left(\frac{G_{T_s}(\cdot;F)}{g(1)}\right)\right) + 1\right\rceil;$$

and the terms $G_{T_s}(\cdot;F)$ and $g(1)$ are defined as in Lemma 24. Recall by the inductive claim in Eq. (51) that there exists $c_9 > 0$ such that

$$\|x_\tau - x^0\| \le c_9 \log(T=\epsilon)^{-p}\frac{(\tau)\,r}{d} \quad \forall\, 8 \le 0 \le \tau \le t:$$

Hence, it follows that

$$\|W_{g_0}(t+1)\| \le c^0 \psi_1(\cdot;F) \max\left\{ \sqrt[p]{t+1}\, l r \frac{CT^2}{\epsilon}, c_9^p \frac{dL}{\sqrt[p]{m}} \frac{(\tau)\,r}{\sqrt[p]{d}}; \frac{(t+1)\,r}{60\sqrt[p]{d}}\right\}:$$

Hence, noting the choice of $T_s$ in Eq. (49), by choosing $\epsilon$ such that

$$c^0 c_9 \psi_1(\cdot;F)\sqrt[p]{T_s}\, l r \frac{CT^2}{\epsilon} \sqrt[p]{c_9^p \frac{dL}{m}} \le \frac{1}{60} \Longleftrightarrow \epsilon \le \frac{m}{360\,(c^0)^2 c_9^2 dL^2\, l r \left(\frac{CT^2}{\epsilon}\right)^2 \psi_1(\cdot;F)^2}; \text{ and } \tag{52}$$

$$c^0 c_9 \psi_1(\cdot;F) \le 1:$$

it follows that

$$\|W_{g_0}(t+1)\| \le \frac{(t+1)\,r}{60\sqrt[p]{d}}:$$

2. Meanwhile, the term $\|W_H(t+1)\|$ can be bounded as follows. By the inductive assumption in Eq. (51), we have that

$$\|\hat{x}_\tau\| = \|x_\tau - x^0\| \le c_9 \log(T=\epsilon)^{-p}\frac{(\tau)\,r}{d} \quad \forall\, 8 \le 0 \le \tau \le t:$$

Moreover, on the event our proof assumes, we know that

$$\max\left\{\|x_\tau - x_0\|^2; \|x_0 - x^0\|^2 \le x_0 - x^0\|^2\right\} \le \mathcal{B}_{T_s}(\cdot;F):$$

Thus, using the $\rho$-Hessian Lipschitz property, we have

$$\|W_H(t+1)\| = \left\| \sum_{\tau=0}^{X^t} (I - \eta H)^{t-\tau}(H - \hat{H})\hat{x}_\tau \right\|$$

$$\le \sum_{\tau=0}^{X^t} (1+\gamma)^{t-\tau}\sqrt[p]{\mathcal{B}_{T_s}(\cdot;F)}\, c_9 \log(T=\epsilon)^{-p}\frac{(\tau)\,r}{d}$$

$$\le c_9(t+1)\log(T=\epsilon)^{-p}\sqrt[p]{\mathcal{B}_{T_s}(\cdot;F)}\frac{(t)\,r}{d}$$

$$\le c_9 T_s \log(T=\epsilon)^{-p}\sqrt[p]{\mathcal{B}_{T_s}(\cdot;F)}\frac{(t)\,r}{d}:$$

Given our choice of $T_s$ in Eq. (49), if

$$c_9 T_s \log(T=\epsilon)^{-p}\sqrt[p]{\mathcal{B}_{T_s}(\cdot;F)} \le \frac{1}{60} \Longleftrightarrow \mathcal{B}_{T_s}(\cdot;F) \le \left(\frac{1}{60 c_9 \log(T=\epsilon)}\right)^2$$

it follows that

$$\|W_H(t+1)\| \le \frac{(t+1)\,r}{60\sqrt[p]{d}}:$$

3. Meanwhile, for $W_u(t+1)$, since the event $\mathcal{S}_u(\epsilon)$ holds, we have that

$$\|W_u(t+1)\| \leq (t+1)^{p-\frac{3}{2}} \cdot 2c_3 \, d^2 (\log(T/\delta))^2 \cdot u^2.$$

Now, by picking

$$(t+1)^{p-\frac{3}{2}} \cdot 2c_3 \, d^2 (\log(T/\delta))^2 \cdot u^2 \leq \frac{(t+1)\, r}{60}\rho^{-\frac{1}{2}}d \quad \Leftrightarrow \quad u \leq s\frac{\rho^{-\frac{1}{2}}r}{120^{\frac{2}{p}}3c_3^{\frac{2}{p}}d\,d^2(\log(T/\delta))^2};$$

it follows that with probability $1-\delta/T$, $\|W_u(t+1)\| \leq \frac{(t+1)}{60}\frac{r}{\rho^{-\frac{1}{2}}d}$.

4. Meanwhile, observe that since $\mathcal{S}_p(\epsilon)$ holds, it follows that

$$W_p(t+1) \geq 2^{\frac{p}{2}-2\log(T/\delta)}\frac{(t+1)\, r}{\rho^{-\frac{1}{2}}d}.$$

Combining the bounds for $W_{g_0}; W_p; W_H$ and $W_u$, it follows that

$$\|\mathcal{X}_{t+1}\| \leq \|W_{g_0}(t+1)\| + \|W_p(t+1)\| + \|W_H(t+1)\| + \|W_u(t+1)\|$$
$$\leq \frac{(t+1)\, r}{\rho^{-\frac{1}{2}}d}\left(\frac{1}{60} + \frac{1}{60} + \frac{1}{60} + 2^{\frac{p}{2}-2\log(T/\delta)}\right)$$
$$\leq \frac{(t+1)\, r}{\rho^{-\frac{1}{2}}d}\left(\frac{1}{20} + 2^{\frac{p}{2}-2}\right)\log(T/\delta);$$

where the final inequality uses the fact that $T/\delta \geq 1 \geq e$ (which implies $\log(T/\delta) \geq 1$). Hence, we see that the first part of the inductive claim of Eq. (51) holds with the constant $c_0 := \frac{1}{20} + 2^{\frac{p}{2}-2}$, and the second part follows naturally as a consequence of our argument above.

Meanwhile, observe that for any $\beta$ such that $\beta \leq \frac{1}{2}$, we have that $(1+\beta)^{\frac{1}{\beta}-1} \leq 2$. Thus, by choosing $\beta$ such that $\beta \leq \frac{1}{2}$, we have that for any $y \leq \frac{1}{\beta}$,

$$(t+1)^2 \geq \frac{1}{2}(t+1)^{2-\beta}.$$

Hence, following Eq. (48), by choosing $T_s \leq \frac{1}{\beta}$, $W_p(T_s)$ is a 1-dimensional Gaussian with variance at least $\frac{2^{\frac{2}{p}}r^2}{d}\rho^{-(T_s)}$, such that with probability at least 2/3,

$$\|W_p(T_s)\| \geq \frac{(T_s)\, r}{10\rho^{-\frac{1}{2}}d}.$$

Simultaneously, we know that on the event

$$\bigcap_{t=0}^{T_s-1} \mathcal{R}_t(\epsilon) \cap \mathcal{S}(\epsilon) \cap \mathcal{S}_u(\epsilon) \cap \mathcal{S}_p(\epsilon);$$

we have

$$\|W_{g_0}(T_s)\| + \|W_H(T_s)\| + \|W_u(T_s)\| \leq \frac{3(T_s)\, r}{60\rho^{-\frac{1}{2}}d} = \frac{(T_s)\, r}{20\rho^{-\frac{1}{2}}d}.$$

We note that by Lemma 24, we have

$$\mathbb{P}\left(\bigcap_{t=0}^{T_s-1} \mathcal{R}_t(\epsilon) \cap \mathcal{S}(\epsilon) \cap \mathcal{S}_u(\epsilon) \cap \mathcal{S}_p(\epsilon)\right) \geq 1 - \left(\frac{24T_s}{T} + \frac{T_s}{T} + \frac{T_s}{T}\right) = 1 - \frac{26T_s}{T}.$$

52

Thus, with probability at least $2 - 3e^{-\frac{26T_s}{T}}$, we have

$$k^\top A_{T_s} k - \frac{1}{2} kW_p(T_s)k \geq \frac{(T_s)}{20} \rho \frac{r}{d}$$

Thus, choosing $T_s = \tau - \nu$, where

$$\nu = \max\left\{ \log_2\left( \frac{p}{T_s(\cdot;F)} 20^p \bar{d}^p \frac{\gamma^{2\alpha^2+2}}{r}! \right); 1 \right\};$$

noting that if $\alpha_1 = 2$, then $(1+\gamma)^{-\frac{1}{\alpha}} \geq (1+\gamma)^{-\frac{1}{\alpha}} \geq \frac{\gamma}{2}$, we have that with probability at least $2 - 3e^{-\frac{26T_s}{T}}$,

$$k^\top A_{T_s} k \geq \frac{(T_s)}{20} \rho \frac{r}{d} = \frac{r}{20} \rho p \frac{(1+\gamma)^{T_s}}{2\gamma + (\gamma)^2}$$

$$\geq \frac{r}{20} \rho \frac{(1+\gamma)^{\log_2\left(\frac{p}{T_s(\cdot;F)}20^p \bar{d}^p \frac{\gamma^{2\alpha^2+2}}{r}!\right)}}{2\gamma + (\gamma)^2}$$

$$\geq \frac{r}{20} \rho p \frac{r}{2\gamma + (\gamma)^2} 2^{\log_2\left(\frac{p}{T_s(\cdot;F)}20^p \bar{d}^p \frac{\gamma^{2\alpha^2+2}}{r}\right)} > 2 \frac{p}{T_s(\cdot;F)} > 2 \frac{p}{(T_s;\cdot)}.$$

Thus, at least one of $\|x_{T_s} - x_0\|$ and $\|x^0_{T_s} - x_0\|$ is larger than $2\frac{p}{(T_s;\cdot)}$, a contradiction. Since the two sequences have the same distribution, it follows that with probability at least $1 - 8e^{-\frac{13T_s}{T}}$, $f(x_{T_s}) - f(x_0) \leq -F$. $\qquad\square$

In the result above, we require an upper bound on the norm of $q_{T_s}(\cdot;F)$ to hold (i.e. equation 50), which in turn necessitates an upper bound on $F$, the function value improvement we can expect to make. Below, we show how to choose $F$ to be as large as possible (up to constants and logarithmic factors) whilst still satisfying equation 50, assuming that $u$ and $r$ are chosen appropriately small such that the dominant term of $\|q_{T_s}(\cdot;F)\|$ scales with $F$.

Lemma 25. Consider choosing $F$ such that

$$F = \frac{1}{2}\left(\frac{1}{60c_9\kappa\log(T=\delta)}\right)^2 \frac{1}{T_s t_f(\alpha)(129 + 8c^{02}_1(\cdot;F))(16(\text{lr}(CT^2=\delta))^2 + 1))}.$$

Suppose $\eta \leq \min\left\{1; \frac{1}{t_f(\alpha)}; \frac{1}{t_f L}\right\}$. Suppose we pick $u$ and $r$ small enough such that

$$u \leq \frac{r^{1=2}}{d\log(T=\delta)^{1=2}}; \quad r^2 \leq \min\left\{\frac{F}{2\kappa\log(T=\delta)\left(\frac{65c_5^2}{8} + 6c_1 + 1\right)}; \left(\frac{F}{4c_6\kappa\log(2dT=\delta) + \frac{8c_7}{\kappa}}\right)^{\frac{9}{=}}\right\}.$$

Then, $\|N_{u;r}(T_s;\cdot)\| \leq F$, and that

$$4c_6\gamma^2 T_s \log(2dT=\delta)r^2 + 4c_7\gamma^2 T_s^2\gamma^2 u^4 d^4 (\log(T=\delta))^4 \leq T_s t_f(\alpha)F.$$

Suppose in addition $\eta$ is small enough so that

$$32\gamma^2(t_f(\alpha))^2\eta^2 \leq \frac{1}{2}\left(\frac{1}{60c_9\kappa\log(T=\delta)}\right)^2.$$

Suppose also that $\alpha_1 \geq 9$ and $\kappa \geq \frac{m}{d}$, so that $T_s = \tau - \nu \geq \frac{d}{m}$. Then, the condition in Eq.(50) will be satisfied.

_____

[9]Without loss of generality, we may set $\alpha = 1$ if $f(\cdot)$ is $(\cdot;\cdot;\alpha\frac{p}{\cdot})$-strict saddle for any $\alpha > 1$.

*Proof.* We note that since $\frac{\iota}{\bar\psi} \le T_s \le \frac{2\iota}{\bar\psi}$, it follows by our choice of $r$ that $r$ also satisfies the condition

$$r^2 \le \min\left\{ \frac{F}{\eta T_s \log(T/\delta)\left(\frac{65c_5^2}{8} + 6c_1 + 1\right)}, \frac{F}{4c_6 \log(2dT/\delta) + 4c_7\eta T_s} \right\}.$$

Hence, our choice of $\eta, u$ and $r$ satisfies the conditions in Lemma 23, and it follows then that

$$\phi_{T_s}(\delta, F) \le \max\left\{ 128\eta T_s t_f(\delta)F, 32\eta^2(t_f(\delta))^2\epsilon^2 + 8c^{\prime 2}\beta_1(\delta; F)\eta t_f(\delta)\max\left\{\frac{16d}{m}(\mathrm{lr}(CT^2/\delta))^2F, T_sF\right\} + T_s\eta t_f(\delta)F \right\},$$

where $\beta_1(\delta; F)$ is as defined in Lemma 21.

The condition in Eq. (50) requires that

$$\phi_{T_s}(\delta, F) \le \left(\frac{\bar\psi}{60c_9\iota\rho\log(T/\delta)}\right)^2.$$

By our choice of $\eta$ such that

$$32\eta^2(t_f(\delta))^2\epsilon^2 \le \frac{1}{2}\left(\frac{\bar\psi}{60c_9\iota\rho\log(T/\delta)}\right)^2,$$

it suffices for us to show that

$$\frac{1}{2}\left(\frac{\bar\psi}{60c_9\iota\rho\log(T/\delta)}\right)^2 \ge 128\eta T_s t_f(\delta)F + 8c^{\prime 2}\beta_1(\delta; F)\eta t_f(\delta)\max\left\{\frac{16d}{m}(\mathrm{lr}(CT^2/\delta))^2F, T_sF\right\} + \eta T_s t_f(\delta)F$$

$$= 129\eta T_s t_f(\delta)F + 8c^{\prime 2}\beta_1(\delta; F)\eta t_f(\delta)\max\left\{\frac{16d}{m}(\mathrm{lr}(CT^2/\delta))^2F, T_sF\right\}.$$

By our assumption, we know that $T_s \ge \frac{d}{m}$. Thus, further simplifying indicates that it suffices for us to show

$$\frac{1}{2}\left(\frac{\bar\psi}{60c_9\iota\rho\log(T/\delta)}\right)^2 \ge 129\eta T_s t_f(\delta)F + 8c^{\prime 2}\beta_1(\delta; F)\eta t_f(\delta)\max\left\{16T_s(\mathrm{lr}(CT^2/\delta))^2F, T_sF\right\}. \tag{53}$$

By choosing $F$ such that

$$F \le \frac{1}{2}\left(\frac{\bar\psi}{60c_9\iota\rho\log(T/\delta)}\right)^2 \frac{1}{\eta T_s t_f(\delta)\left(129 + 8c^{\prime 2}\beta_1(\delta; F)\left(16(\mathrm{lr}(CT^2/\delta))^2 + 1\right)\right)},$$

we see that Eq. (53) is satisfied.

$\square$

*Remark* 3. Suppose without loss of generality that $T_s = \frac{\iota}{\eta\bar\psi}$. Then, as a consequence of Lemma 25, we note that the amortized function value progress of decreasing function value by $F$ over $T_s$ iterations is

$$\frac{F}{T_s} = \frac{1}{2}\left(\frac{\bar\psi}{60c_9\iota\rho\log(T/\delta)}\right)^2 \frac{1}{\eta T_s^2 t_f(\delta)\left(129 + 8c^{\prime 2}\beta_1(\delta; F)\left(16(\mathrm{lr}(CT^2/\delta))^2 + 1\right)\right)}$$

$$= \eta\frac{\bar\psi^4}{\rho^2}\frac{1}{2\iota^2}\frac{1}{(60c_9\iota\log(T/\delta))^2\,(t_f(\delta))\left(129 + 8c^{\prime 2}\beta_1(\delta; F)\left(16(\mathrm{lr}(CT^2/\delta))^2 + 1\right)\right)}$$

## F. Proving the main result (informal statement in Theorem 1, full statement in Theorem 2)

In this section, we prove our main result. First, we need an additional result (Lemma 26) showing that with high probability, we can bound the function value increase if a saddle appears within $t_f(\delta)$ iterations immediately after we have had $T_s$ iterations after the previous saddle. We note that such a bound is necessary because our earlier result upper bounding function increase in $\tau$ iterations (see Lemma 18) focused on the case where $\tau \ge t_f(\delta)$. Next, we state and prove Theorem 2, which is the precise version of Theorem 1 in the main text.

**Lemma 26** (Function change for small $\tau$). *Let $c_1 > 0, c_4 > 0, c_5 > 0, C_1 \geq 1$ be the absolute constants defined in the statements of the previous lemmas. Let $\delta \in (0, 1/e]$, and suppose $\tau < t_f(\delta)$.*

*Let $J$ denote the interval $\{0, 1 \ldots, \tau - 1\}$ where $\tau < t_f(\delta)$.*

*Suppose we choose $\eta$ such that*

$$\eta \leq \frac{1}{Lt_f(\delta)} \cdot \min\left\{\frac{\sqrt{m}}{8c_4(\mathrm{lr}(C_1 dmT/\delta))^{3/2}\sqrt{d}}, \frac{m}{128c_1(\mathrm{lr}(C_1 dmT/\delta))^3 d}\right\}. \tag{54}$$

*Suppose also we pick $u, r$ and $\eta$ as prescribed in the statement of Proposition 4.*

*Suppose that $\min_{t2J}\|\nabla f(x_t)\| \leq \epsilon$. Then, on the event*

$$\mathcal{D}_\tau(\delta) := \mathcal{H}_{0,\tau}(\delta) \cap \left(\bigcap_{t=0}^{\tau-1} \mathcal{A}_t(\delta)\right) \cap \left(\bigcap_{t=0}^{\tau-1} \mathcal{G}_t(\delta)\right),$$

*we have the following upper bound on function value change:*

$$f(x_\tau) - f(x_0) \leq \frac{\eta}{4}\epsilon^2 + t_f(\delta)\eta u^4 \rho^2 \cdot c_1 d^3 \left(\log\frac{T}{\delta}\right)^3 + t_f(\delta)L\eta^2 u^4 \rho^2 \cdot c_1 d^4 \left(\log\frac{T}{\delta}\right)^4$$

$$+ \eta c_1 r^2 (128 t_f(\delta) + \eta L)\log\frac{T}{\delta} + t_f(\delta)c_1 L\eta^2 r^2.$$

*Moreover, $\mathsf{P}(\mathcal{D}_\tau(\delta)) \geq 1 - \frac{(4t_f(\delta)+4)\delta}{T}$.*

*Proof.* Throughout the proof, we assume that the event $\mathcal{D}_\tau(\delta)$ holds.

Let $J$ denote $\{0, 1 \ldots, \tau - 1\}$ where $\tau < t_f(\delta)$. Then, $J$ belongs to one of the two following cases.

**Case 1**) (Gradient dominates noise): Recall that this means that for every $t \in J$, we have

$$\|\nabla f(x_t)\| > 8t_f(\delta)\eta L\left(\frac{u}{2}\frac{1}{m}\sum_{i=1}^{n} Z_{t,i}Z_{t,i}^>\tilde{H}_{t,i}Z_{t,i} + \|Y_t\|\right).$$

By our choice of $\eta$ in Eq. (31), we can apply Lemma 16 to get

$$\min_{t2J}\|\nabla f(x_t)\| \geq \frac{1}{4}\max_{t2J}\|\nabla f(x_t)\|.$$

Thus by setting $\alpha = 128t_f(\delta)$ in Eq. (5) and by choosing $\eta$ such that

$$\frac{c_1 L\eta^2\chi^3 d}{m} \leq \frac{\eta}{\alpha} = \frac{\eta}{128t_f(\delta)} \iff \eta \leq \frac{m}{128c_1 Lt_f(\delta)d\chi^3},$$

it follows that

$$-\frac{3\eta}{4}\sum_{t2J}\left(\frac{1}{m}\sum_{i=1}^{n} Z_{t,i}^>\nabla f(x_t)\right)^2 + \left(\frac{\eta}{128t_f(\delta)} + \frac{c_1 L\eta^2\chi^3 d}{m}\right)\sum_{t2J}\|\nabla f(x_t)\|^2$$

$$= -\frac{3\eta}{4}\sum_{t2J}\left(\frac{1}{m}\sum_{i=1}^{n} Z_{t,i}^>\nabla f(x_t)\right)^2 + \frac{\eta}{64t_f(\delta)}\sum_{t2J}\|\nabla f(x_t)\|^2$$

$$\leq \frac{\eta}{64t_f(\delta)}\sum_{t2J}\|\nabla f(x_t)\|^2$$

$$\leq \frac{\eta}{64t_f(\delta)}\sum_{t2J}\max_{t2J}\|\nabla f(x_t)\|^2$$

$$\leq \frac{16\eta}{64t_f(\delta)}\sum_{t2J}\min_{t2J}\|\nabla f(x_t)\|^2 \leq \frac{\eta}{4}\min_{t2J}\|\nabla f(x_t)\|^2 \leq \frac{\eta}{4}\epsilon^2, \tag{55}$$

where the final bound holds since we assumed $\min_{t2J}\|\nabla f(x_t)\| \leq \epsilon$.

**Case 2)** (Gradient does not dominate noise): there exists some $t \in J$ such that

$$\|\nabla f(x_t)\| \le 8t_f(\delta)\eta L \left( \frac{u}{2} \left\| \frac{1}{m}\sum_{i=1}^n Z_{t,i}Z_{t,i}^\top \tilde{H}_{t,i}Z_{t,i} \right\| + \|Y_t\| \right).$$

By our choice of $\eta$ in Eq. (31), we can apply Lemma 17 to get

$$\|\nabla f(x_t)\| \le c_5 t_f(\delta)\eta L \left( u^2 d^2 \rho \left(\log\frac{T}{\delta}\right)^2 + \left(1 + \sqrt{\frac{\log(T/\delta)}{d}}\right) r \right) \qquad \forall t \in J.$$

Note that, by our choices of the parameters $\eta, u, r$, it can be shown that

$$c_5 t_f(\delta)\eta L \left( u^2 d^2 \rho \left(\log\frac{T}{\delta}\right)^2 + \left(1 + \sqrt{\frac{\log(T/\delta)}{d}}\right) r \right) < \epsilon,$$

Hence, by setting $\alpha = 128 t_f(\delta)$ in Eq. (5) and choosing $\eta$ such that

$$\frac{c_1 L \eta^2 \chi^3 d}{m} \le \frac{\eta}{\alpha} = \frac{\eta}{128 t_f(\delta)},$$

it follows that

$$
\begin{aligned}
&\left(\frac{\eta}{128 t_f(\delta)} + \frac{c_1 L \eta^2 \chi^3 d}{m}\right) \sum_{t\in J} \|\nabla f(x_t)\|^2 \\
&\le \frac{\eta}{64 t_f(\delta)} \sum_{t\in J} \left( c_5 t_f(\delta)\eta L \left( u^2 d^2 \rho \left(\log\frac{T}{\delta}\right)^2 + \left(1 + \sqrt{\frac{\log(T/\delta)}{d}}\right) r \right) \right)^2 \\
&\le \frac{\eta}{64 t_f(\delta)} \sum_{t\in J} \epsilon^2 \\
&\le \frac{\eta}{64}\epsilon^2
\end{aligned}
\tag{56}
$$

Combining both cases above (Eq. (55) and Eq. (56)), we see that for the choice $\alpha = 128 t_f(\delta)$, the bound

$$-\frac{3\eta}{4}\sum_{t\in J}\left(\frac{1}{m}\sum_{i=1}^n Z_{t,i}^\top \nabla f(x_t)\right)^2 + \left(\frac{\eta}{128 t_f(\delta)} + \frac{c_1 L \eta^2 \chi^3 d}{m}\right)\sum_{t\in J}\|\nabla f(x_t)\|^2 \le \frac{\eta}{4}\epsilon^2 \tag{57}$$

always holds.

Recall by Eq. (5) that we have

$$
\begin{aligned}
f(x_\tau) - f(x_0) \le{} & -\frac{3\eta}{4}\sum_{t=0}^{\tau-1}\left(\frac{1}{m}\sum_{i=1}^n Z_{t,i}^\top \nabla f(x_t)\right)^2 + \left(\frac{\eta}{\alpha} + \frac{c_1 L \eta^2 \chi^3 d}{m}\right)\sum_{t=0}^{\tau-1}\|\nabla f(x_t)\|^2 \\
& + \tau \eta u^4 \rho^2 \cdot c_1 d^3 \left(\log\frac{T}{\delta}\right)^3 + \tau L \eta^2 u^4 \rho^2 \cdot c_1 d^4 \left(\log\frac{T}{\delta}\right)^4 \\
& + \eta c_1 r^2 (\alpha + \eta L)\log\frac{T}{\delta} + \tau c_1 L \eta^2 r^2.
\end{aligned}
$$

By plugging in Eq. (57) above, as well as the choice $\alpha = 128 t_f(\delta)$, we see that

$$
\begin{aligned}
f(x_\tau) - f(x_0) \le{} & \frac{\eta}{4}\epsilon^2 + t_f(\delta)\eta u^4 \rho^2 \cdot c_1 d^3 \left(\log\frac{T}{\delta}\right)^3 + t_f(\delta)L \eta^2 u^4 \rho^2 \cdot c_1 d^4 \left(\log\frac{T}{\delta}\right)^4 \\
& + \eta c_1 r^2 (128 t_f(\delta) + \eta L)\log\frac{T}{\delta} + t_f(\delta)c_1 L \eta^2 r^2.
\end{aligned}
$$

56

We can now complete our proof by using the union bound (suppressing the dependence of some of the events on $\delta$ for notational simplicity) to derive

$$\mathsf{P}(\mathcal{D}_\tau^c) \le \mathsf{P}(\mathcal{H}_\tau^c) + \sum_{t=0}^{\tau-1} \mathsf{P}(\mathcal{A}_t^c) + \sum_{t=0}^{\tau-1} \mathsf{P}(\mathcal{G}_t^c)$$

$$\le \frac{(\tau+4)\delta}{T} + \frac{\tau}{T}\delta + 2\frac{\tau}{T}\delta \le \frac{(4t_f(\delta)+4)}{T}\delta \qquad \square$$

Armed with Proposition 5 and Lemma 25, we are now ready to show for $T$ sufficiently large, with high probability, there can be no more than $T/4$ $\epsilon$-saddle points. Combined with Proposition 4, this yields the following result.

**Theorem 2.** *Suppose we pick $u, r, \eta$ such that they satisfy the conditions in Proposition 5 and Lemma 25. Suppose $F$ is chosen as prescribed in Lemma 25. Suppose that $\bar{\psi} \le 1$, so that $T_s \ge \frac{\iota}{\eta\bar{\psi}} \ge \frac{d}{mL}$[10]. Suppose we pick $T_s$ as prescribed in Proposition 5. Suppose in addition we pick $r$ such that*

$$r^2 \le \min\left\{\frac{\epsilon^2}{4(130c_1 t_f(\delta) + c_1\log(T/\delta) + c_1)}, \frac{F\bar{\psi}}{80\iota\log(T/\delta)\left(\frac{65c_5^2}{8} + 132c_1 + 1\right)}\right\}.$$

*Suppose also that we choose $\eta$ such that*

$$\eta \le \frac{0.1}{2\epsilon^2}\frac{\bar{\psi}}{2\iota}\frac{1}{2}\left(\frac{\bar{\psi}}{60c_9\iota\rho\log(T/\delta)}\right)^2 \frac{1}{t_f(\delta)\left(129 + 8c'^2\beta_1(\delta; F)\right)\left(16(\mathrm{lr}(CT^2/\delta))^2 + 1\right)}$$

*Suppose*

$$T \ge \left(\frac{256t_f(\delta)\left(f(x_0) - f\right) + \epsilon^2/L)}{\eta\epsilon^2}, \frac{\varphi\rho^2\left(f(x_0) - f\right)}{\eta\bar{\psi}^4}, 256\lceil\frac{\iota}{\eta\bar{\psi}}\rceil, 256t_f(\delta), 1024\right), \qquad (58)$$

*where*

$$\varphi := 20\left(2\iota^2(60c_9\iota\log(T/\delta))^2\left(t_f(\delta)\right)\left(129 + 8c'^2\beta_1(\delta; F)\right)\left(16(\mathrm{lr}(CT^2/\delta))^2 + 1\right)\right).$$

*Then, with probability at least $1 - 22\delta$, there are at least $T/2$ $\epsilon$-approximate second order stationary points.*

*Proof.* Consider defining the following sequence of stopping times:

$$\tau_1 = \inf_t\{t \le T : \|\nabla f(x_t)\| < \epsilon, \lambda_{\min}(\nabla^2 f(x_t)) \le -\sqrt{\rho\epsilon}\},$$

$$\tau_{i+1} = \inf_t\{t \le T : t > \tau_i + T_s, \|\nabla f(x_t)\| < \epsilon, \lambda_{\min}(\nabla^2 f(x_t)) \le -\sqrt{\rho\epsilon}\}, \quad \forall 1 \le i \le \lfloor T/T_s\rfloor.$$

$\square$

We note that if $\tau_i = T$, then $\tau_j = T$ for any $j > i$. Let $N_s$ denote the (random) number of saddle points encountered in $T$ iterations.

We observe that we can decompose the function change as

$$f(x_T) - f(x_0)$$
$$= (f(x_{\tau_{N_s}}) - f(x_0)) + (f(x_T) - f(x_{\tau_{N_s}}))$$

---

[10]Recall we focus on the case $\ge L$, since otherwise, by the $L$-Lipschitz assumption, $\lambda_{\min}(\nabla^2 f(x)) \ge -L$ for all $x \in \mathbb{R}^d$, i.e. $\epsilon$-first order stationary points are also $\epsilon$-second order stationary points.

$$= (f(x_{\tau_1}) - f(x_0)) + \sum_{i=1}^{N_s} (f(x_{\tau_i+T_s}) - f(x_{\tau_i})) + \sum_{i=1}^{N_s-1} \left( f(x_{\tau_{i+1}}) - f(x_{\tau_i+T_s}) \right) + (f(x_T) - f(x_{\tau_{N_s}}))$$

$$= \underbrace{\sum_{i=1}^{N_s} (f(x_{\tau_i+T_s}) - f(x_{\tau_i})) + (f(x_{\tau_1}) - f(x_0))}_{U_1} + \underbrace{\sum_{i=1}^{N_s-1} \left( f(x_{\tau_{i+1}}) - f(x_{\tau_i+T_s}) \right) + (f(x_T) - f(x_{\tau_{N_s}}))}_{U_2}.$$

We first consider $U_1$. Letting $x_j := x_T$ for any $j \geq T$, we have that

$$\sum_{i=1}^{N_s} f(x_{\tau_i+T_s}) - f(x_{\tau_i}) = \sum_{i=1}^{\lceil T/T_s \rceil} (f(x_{\tau_i+T_s}) - f(x_{\tau_i})) \, 1_{\tau_i < T}$$

Now, by Eq. (32), observe that with probability at least $1 - \frac{(5T_s+4)\delta}{T} \geq 1 - \frac{6T_s\delta}{T}$ (note $T_s \geq 4$), for any $1 \leq i \leq T/T_s$, we have that

$$(f(x_{\tau_i+T_s}) - f(x_{\tau_i})) \, 1_{\tau_i < T} \leq \tau \frac{c_5^2}{64} \eta^3 t_f(\delta)^2 L^2 \left( u^2 d^2 \rho \left( \log \frac{T}{\delta} \right)^2 + \left( \frac{2 \log(T/\delta) r}{} \right)^2 \right)$$

$$+ \tau \eta u^4 \rho^2 \cdot c_1 d^3 \left( \log \frac{T}{\delta} \right)^3 + \tau L \eta^2 u^4 \rho^2 \cdot c_1 d^4 \left( \log \frac{T}{\delta} \right)^4$$

$$+ \eta c_1 r^2 (128 t_f(\delta) + \eta L) \log \frac{T}{\delta} + \tau c_1 L \eta^2 r^2$$

$$:= M_{u,r,T_s}.$$

Suppose we pick $u, r$ such that $M_{u,r,T_s} \leq 0.1F$. Recall from Proposition 5 that with probability at least $1/3 - \frac{13T_s\delta}{T}$, $(f(x_{\tau_i+T_s}) - f(x_{\tau_i})) \, 1_{\tau_i < T} \leq -F$. Choosing $\delta$ such that $1/3 - \frac{13T_s\delta}{T} \geq 0.3$, and letting $\mu = 0.1F$, we note that $|-F + \mu| = 0.9F \geq \frac{0.7}{0.3} 0.2F \geq \frac{0.7}{0.3}(M_{u,r,T_s} + \mu)$.

Now, let $\mathcal{E}_{\tau_i}$ denote the bad event on which

$$\text{neither } (f(x_{\tau_i+T_s}) - f(x_{\tau_i})) \, 1_{\tau_i < T} \leq -F, \quad \text{nor } (f(x_{\tau_i+T_s}) - f(x_{\tau_i})) \, 1_{\tau_i < T} \leq M_{u,r,T_s} \leq 0.1F.$$

We know that $\mathcal{E}_{\tau_i}$ has probability at most $\frac{6T_s\delta}{T}$. Let $\mathcal{E}_\tau := \cup_{i=1}^{\lceil T/T_s \rceil} \mathcal{E}_{\tau_i}$, such that $\mathsf{P}(\mathcal{E}_\tau) \leq 6\delta$. Then, by applying the weakened supermartingale inequality in Proposition 3, we have

$$\mathsf{P}\left( \sum_{i=1}^{T/T_s} (f(x_{\tau_i+T_s}) - f(x_{\tau_i})) \, 1_{\tau_i < T} \geq -N_s 0.9F + s \right) \leq \mathsf{E}\left[ \exp\left( -\frac{s^2}{4N_s F^2} \right) \right] + \mathsf{P}(\mathcal{E}_\tau) \leq \exp\left( -\frac{s^2}{4(T/T_s)F^2} \right) + 6\delta.$$

Now, pick $s = 2F\sqrt{\log(1/\delta)T/T_s}$, then

$$\mathsf{P}\left( \sum_{i=1}^{T/T_s} (f(x_{\tau_i+T_s}) - f(x_{\tau_i})) \, 1_{\tau_i < T} \geq -N_s 0.9F + 2F\sqrt{\log(1/\delta)T/T_s} \right) \leq 7\delta.$$

Note that supposing for contradiction that there are at least $T/4$ saddles, we must then have that $N_s \geq T/(4T_s)$, such that

$$-N_s 0.9F + 2F\sqrt{\log(1/\delta)T/T_s} \leq F(-0.9T/(4T_s) + (2\sqrt{\log(1/\delta)T/T_s})) \leq F(-0.1T/T_s),$$

where we may ensure the last inequality by picking $T/T_s$ such that

$$T/T_s \geq \left( \frac{2}{0.125} \right)^2 \sqrt{\log(1/\delta)} = 256\sqrt{\log(1/\delta)}.$$

58

Note that our choice of $T$ ensures this.

Thus, with probability at least $1 - 7\delta$,

$$U_1 = \sum_{i=1}^{T/T_s} (f(x_{\tau_i + T_s}) - f(x_{\tau_i})) 1_{\tau_i < T} \le -(0.1T/T_s)F.$$

Next, we bound the summand $U_2$. Recall that

$$U_2 = (f(x_{\tau_1}) - f(x_0)) + \sum_{i=1}^{N-1} f(x_{\tau_{i+1}}) - f(x_{\tau_i + T_s}).$$

Without loss of generality, we may analyze each of the summands $f(x_{\tau_{i+1}}) - f(x_{\tau_i + T_s})$ in the same way as we treat $(f(x_{\tau_1}) - f(x_0))$. Let us then consider the summand $f(x_{\tau_1}) - f(x_0)$. There are two cases to consider.

1. The first is when $\tau_1 < t_f(\delta)$. In this case, since we know that $\|\nabla f(x_{\tau_1})\| \le \epsilon$ (as $x_{\tau_1}$ is an $\epsilon$-saddle point), it follows by Lemma 26 that

$$f(x_{\tau_1}) - f(x_0) \le \frac{\eta}{4}\epsilon^2 + t_f(\delta)\eta u^4 \rho^2 \cdot c_1 d^3 \left( \log \frac{T}{\delta} \right)^3 + t_f(\delta)L\eta^2 u^4 \rho^2 \cdot c_1 d^4 \left( \log \frac{T}{\delta} \right)^4$$

$$+ \eta c_1 r^2 (128 t_f(\delta) + \eta L) \log \frac{T}{\delta} + t_f(\delta)c_1 L\eta^2 r^2$$

with probability at least $1 - \frac{(4t_f(\delta)+4)\delta}{T}$.

2. The second case is when $\tau_1 \ge t_f(\delta)$. In this case, by Lemma 18, we have that

$$f(x_{\tau_1}) - f(x_0) \le \tau_1 \frac{c_5^2}{64}\eta^3 t_f(\delta)^2 L^2 \left( u^2 d^2 \rho \left( \log \frac{T}{\delta} \right)^2 + \rho \frac{p}{2\log(T/\delta)r} \right)^2$$

$$+ \tau_1 \eta u^4 \rho^2 \cdot c_1 d^3 \left( \log \frac{T}{\delta} \right)^3 + \tau_1 L\eta^2 u^4 \rho^2 \cdot c_1 d^4 \left( \log \frac{T}{\delta} \right)^4$$

$$+ \eta c_1 r^2 (128 t_f(\delta) + \eta L) \log \frac{T}{\delta} + \tau_1 c_1 L\eta^2 r^2.$$

with probability at least $1 - \frac{(5\tau_1 + 4)\delta}{T}$.

By our choice of $u$, we know that

$$t_f(\delta)\eta u^4 \rho^2 \cdot c_1 d^3 \left( \log \frac{T}{\delta} \right)^3 + t_f(\delta)L\eta^2 u^4 \rho^2 \cdot c_1 d^4 \left( \log \frac{T}{\delta} \right)^4 + \eta c_1 r^2 (128 t_f(\delta) + \eta L) \log \frac{T}{\delta} + t_f(\delta)c_1 L\eta^2 r^2$$

$$\le t_f(\delta)r^2 c_1 + t_f(\delta)r^2 c_1 + c_1 r^2 (128 t_f(\delta) + 1) \log(T/\delta) + c_1 r^2$$

$$= r^2 (130 c_1 t_f(\delta) + c_1 \log(T/\delta) + c_1).$$

Hence, by picking $r$ such that

$$r \le \frac{\epsilon^2}{4(130 c_1 t_f(\delta) + c_1 \log(T/\delta) + c_1)},$$

it follows that

$$\frac{\eta \epsilon^2}{4} \ge t_f(\delta)\eta u^4 \rho^2 \cdot c_1 d^3 \left( \log \frac{T}{\delta} \right)^3 + t_f(\delta)L\eta^2 u^4 \rho^2 \cdot c_1 d^4 \left( \log \frac{T}{\delta} \right)^4$$

$$+ \eta c_1 r^2 (128 t_f(\delta) + \eta L) \log \frac{T}{\delta} + t_f(\delta)c_1 L\eta^2 r^2.$$

59

Then, if $\tau_1 < t_f(\delta)$, with probability at least $1 - \frac{(5t_f(\delta)+4)}{\delta}$,

$$f(x_{\tau_1}) - f(x_0) \leq \frac{\eta\epsilon^2}{2}.$$

Suppose also that we pick $r$ such that

$$r^2 \leq \frac{F\sqrt{\rho\epsilon}}{80\iota \log(T/\delta)\left(\frac{65c_5^2}{8} + 132c_1 + 1\right)} \leq \frac{F}{40\eta T_s \log(T/\delta)\left(\frac{65c_5^2}{8} + 132c_1 + 1\right)}.$$

Then, it can be verified that

$$\frac{F}{40}\frac{T}{T_s} \geq T\frac{c_5^2}{64}\eta^3 t_f(\delta)^2 L^2 \left(u^2 d^2 \rho \left(\log\frac{T}{\delta}\right)^2 + \rho\frac{1}{2\log(T/\delta)r}\right)^2$$

$$+ T\eta u^4 \rho^2 \cdot c_1 d^3 \left(\log\frac{T}{\delta}\right)^3 + TL\eta^2 u^4 \rho^2 \cdot c_1 d^4 \left(\log\frac{T}{\delta}\right)^4$$

$$+ \frac{T}{T_s}\eta c_1 r^2 (128 t_f(\delta) + \eta L) \log\frac{T}{\delta} + Tc_1 L\eta^2 r^2.$$

Then, by a union bound, it follows that with probability at least $1 - 9\delta$,

$$U_2 = (f(x_{\tau_1}) - f(x_0)) + \sum_{i=1}^{N_s - 1} f(x_{\tau_{i+1}}) - f(x_{\tau_i + T_s})$$

$$\leq \frac{T}{T_s}\frac{\eta\epsilon^2}{2} + \frac{F}{40}\frac{T}{T_s}$$

Therefore, by the union bound, with probability at least $1 - 16\delta$,

$$f(x_{\tau_{N_s}}) - f(x_0) = U_1 + U_2 \leq \frac{T}{T_s}\left(-0.1F + \eta\epsilon^2/2 + \frac{F}{40}\right)$$

By recalling our choice of $F$ in Lemma 25, by choosing $\eta$ such that

$$\eta \leq \frac{0.1}{2\epsilon^2}\frac{\bar{\psi}}{2\iota}\frac{1}{2}\left(\frac{\bar{\psi}}{60c_9\iota\rho\log(T/\delta)}\right)^2\frac{1}{t_f(\delta)\left(129 + 8c'^2\beta_1(\delta; F)\right)\left(16(\mathrm{lr}(CT^2/\delta))^2 + 1\right)}$$

$$\leq \frac{0.1}{2\epsilon^2}\frac{1}{2}\left(\frac{\sqrt{\rho\epsilon}}{60c_9\iota\rho\log(T/\delta)}\right)^2\frac{1}{\eta T_s t_f(\delta)\left(129 + 8c'^2\beta_1(\delta; F)\right)\left(16(\mathrm{lr}(CT^2/\delta))^2 + 1\right)} = \frac{0.1F}{2\epsilon^2},$$

it follows that with probability at least $1 - 16\delta$,

$$f(x_{\tau_{N_s}}) - f(x_0) = U_1 + U_2$$

$$\leq \frac{T}{T_s}\left(-0.1F + \eta\epsilon^2/2 + \frac{F}{40}\right)$$

$$\leq \frac{T}{T_s}(-0.1F + 0.1F/4 + 0.1F/4) = \frac{T}{T_s}(-0.05F).$$

Choose $T$ such that

$$-(0.05T/T_s)F \leq -(f(x_0) - f^*) \iff T \geq \frac{20T_s(f(x_0) - f^*)}{F} \geq \frac{\varphi\rho^2(f(x_0) - f^*)}{\eta\bar{\psi}^4}$$

yields a contradiction, where

$$\varphi := 20\left(2\iota^2(60c_9\iota\log(T/\delta))^2(t_f(\delta))\left(129 + 8c'^2\beta_1(\delta; F)\right)\left(16(\mathrm{lr}(CT^2/\delta))^2 + 1\right)\right)$$

Hence, with probability at least $1 - 16\delta$, there cannot be more than $T/4$ saddle points. In addition, with probability at least $1 - 6\delta$, by Proposition 4, there cannot be more than $T/4$ iterates with $\|\nabla f(x_t)\| \geq \epsilon$. Hence, with probability at least $1 - 22\delta$, there are at least $T/2$ $\epsilon$-approximate second order stationary points.

# G. More complete discussion of simulations

We test the performance of our proposed algorithm with two-point estimators (ZOPGD-2pt) against existing zeroth-order benchmarks using the *octopus function* (proposed in (Du et al., 2017)) of varying dimensions.[11] It is known that the octopus function defined on $\mathbb{R}^d$, which chains $d$ saddle points sequentially, takes exponential (in $d$) time for exact gradient descent to escape; it has thus emerged as a popular benchmark to evaluate and compare the performance of algorithms that seek to escape saddle points. In our experiments, we compare the performance of our two-point estimator algorithm (ZOPGD-2pt) with PAGD (Algorithm 1 in (Vlatakis-Gkaragkounis et al., 2019)) and ZO-GD-NCF (see (Zhang et al., 2022)), which are the only two existing zeroth-order algorithms that have (a) a $\tilde{O}(d/\epsilon^2)$ sample complexity for escaping saddle points (with the latter algorithm yielding the tightest bounds), and (b) performed the best empirically on escaping saddle points (see the simulation results in (Zhang et al., 2022)). We note that both PAGD and ZO-GD-NCF have to use $2d$ function evaluations per iteration to estimate the gradient while our algorithm only needs to use 2 function evaluations. In our plots, we plot the function value against the number of function evaluations. For completeness, we also plot the performance of exact gradient descent (normalized such that its $x$-axis is also the number of function queries).

We tested the algorithms for $d = 10$ and $d = 30$. To account for the stochasticity in the algorithms, for each algorithm, we computed the average and standard deviation over 30 trials, and plotted the mean trajectory with an additional band that represents 1.5 times the standard deviation. For our algorithm's hyperparameters, we picked

$$\eta = \frac{1}{4dL}, u = 10^{-2}, r = 0.05, m = 1( \text{ i.e. two-point estimator}) \tag{59}$$
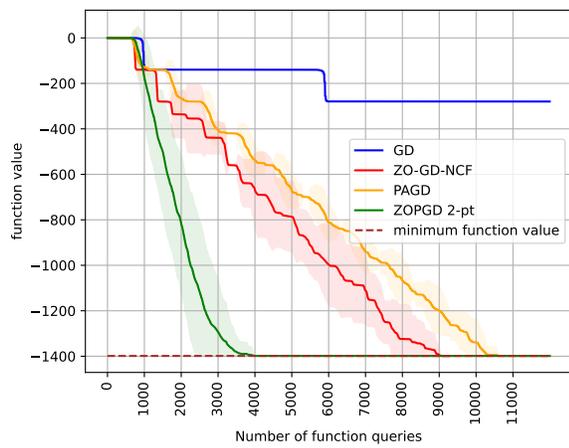
For PAGD, we used the hyperparameters listed in their paper, and for ZO-GD-NCF, we used the code from their Neurips submission. We note in particular that both methods used the step-size $\frac{1}{4L}$. For initialization, we chose a random $x_0$ near the saddle point at the origin, drawn from $N(0, 10^{-3}I_{d\times d})$[12] (fixed for all trials and all algorithms).

As we can see in Fig. 2, in both cases, our algorithm reaches the global minimum of the octopus function in significantly fewer function evaluations than PAGD and ZO-GD-NCF (approximately 2.5 times faster than ZO-GD-NCF, and approximately 3 times faster than PAGD), despite our algorithm only using 2 function evaluations per iteration compared to $2d$ function evaluations per iteration for both PAGD and ZO-GD-NCF. As a sanity check, we note that the number of function evaluations required for PAGD and ZO-GD-NCF to reach the global minimum approximately matches that in Figure 1 of (Zhang et al., 2022); here the correspondence is only approximate since (Zhang et al., 2022) only plots one trial while we compute the mean and standard deviation of 30 trials.

This result suggests that in addition to the theoretical convergence guarantees, there might also be empirical benefits to using two-point estimators versus existing $2d$-point estimators in the zeroth-order escaping saddle point literature.

---

[11]Our code can be found at https://github.com/rafflesintown/escape-saddle-points-2pt

[12]Using the random seed in our code, we note that $\|\nabla f(x_0)\| = 0.011$ for $d = 10$ and $\|\nabla f(x_0)\| = 0.030$ for $d = 30$.

(a) $d = 10$

(b) $d = 30$

Figure 2: Performance on toy octopus function, with $\tau = e, L = e, \gamma = 1$ (Here, $\tau, L, \gamma$ are parameters determining the properties of $f$. Our parameter choice is consistent with that in (Zhang et al., 2022). See (Du et al., 2017) for details about the definitions of $\tau, L$ and $\gamma$.).