
Gradient-Free Approaches is a Key to an Efficient Interaction with Markovian Stochasticity

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 This paper deals with stochastic optimization problems involving Markovian
2 noise with a zero-order oracle. We present and analyze a novel derivative-
3 free method for solving such problems in strongly convex smooth and
4 non-smooth settings with both one-point and two-point feedback oracles.
5 Using a randomized batching scheme, we show that when mixing time τ of
6 the underlying noise sequence is less than the dimension of the problem d , the
7 convergence estimates of our method do not depend on τ . This observation
8 provides an efficient way to interact with Markovian stochasticity: instead
9 of invoking the expensive first-order oracle, one should use the zero-order
10 oracle. Finally, we complement our upper bounds with the corresponding
11 lower bounds. This confirms the optimality of our results.

12 1 Introduction

13 Stochasticity is a fundamental aspect of many optimization problems, naturally arising
14 in the field of machine learning [48, 28]. Stochastic gradient descent (SGD) [45] and its
15 accelerated variants [38, 25] have become a de facto optimizers for modern large models
16 training. Theoretical properties of SGD have been extensively studied under various statistical
17 frameworks [36, 24, 10, 56], often relying on the assumption that noise is independent
18 and identically distributed (i.i.d.). However, in many real-world applications — including
19 reinforcement learning (RL) [6, 16], distributed optimization [35, 31], and bandit problems
20 [3] — noise is not i.i.d., instead exhibiting correlations or *Markovian structure*.

21 For instance, in the mentioned growing field of RL, sequential interactions with the environ-
22 ment induce state-dependent structure of the noise, creating a need for non-i.i.d. noise aware
23 algorithms. Although several gradient-based methods for Markovian stochastic oracles have
24 been studied in the past decade [14, 18], policy optimization in RL is based solely on reward
25 feedback, making traditional methods inapplicable, since there is no access to first-order
26 information [46, 9, 19]. *Zero-order optimization* (ZOO) methods are specifically developed
27 to address such problems, and are used in scenarios where gradients are unavailable or
28 prohibitively expensive to compute. Apart from RL, ZOO techniques are widely employed in
29 adversarial attack generation [8], hyperparameter tuning [47, 57], continuous bandits [7, 49]
30 and other applications [54, 33]. While the literature on ZOO is extensive, this work is, to
31 our knowledge, *the first study of optimization problem with both zero-order information and*
32 *Markovian noise*, aimed at developing an optimal algorithm for a large family of problems
33 from the intersection of these two areas.

34 1.1 Related works

35 \diamond **Zero-order** methods is one of the key and oldest areas of optimization. There are various
 36 zero-order approaches, here we can briefly highlight, e.g., one-dimensional methods [32, 42]
 37 or their high-dimensional analogues [41], ellipsoid algorithms [58] and searches along random
 38 directions [4]. Currently, the most popular and most studied mechanism behind ZOO
 39 methods is the finite-difference approximation of the gradient described in [43, 20, 40]. The
 40 idea is simple: querying two sufficiently close points is essentially equivalent to finding a
 41 value of the directional derivative of the function:

$$\langle \nabla f(x), e \rangle \approx \frac{f(x + te) - f(x)}{t} \approx \frac{f(x + te) - f(x - te)}{2t}, \quad (1)$$

42 where e is a random direction. It can be a random coordinate, a vector from the Euclidean
 43 sphere or a sample of the Gaussian distribution. The approximation (1) in turn leads back to
 44 the gradient methods or coordinate algorithms of Nesterov [39]. There are, however, several
 45 differences:

- 46 • First, to get full gradient information, the algorithm would need d queries instead of one
 47 gradient oracle call (here d is the dimension of x).
- 48 • Second, if the ZO oracle is inexact, i.e. only noisy values of function are available, then
 49 finite difference schemes can fail if noise components do not cancel out.

50 The setting of the second point, when function evaluations experience zero-mean additive
 51 perturbations, is called *Stochastic ZOO*. The stochasticity, as noted before, is abundant in
 52 the modern optimization world. To tackle this issue, additional assumptions about the noise
 53 structure are required. Here we briefly discuss two main ideas adopted in the literature, and
 54 refer the reader to Section 2 for precise definitions.

55 In the case of *two-point feedback*, we assume that for a fixed value of the noise variable one
 56 can call the stochastic zero-order oracle at least twice. It means that we can compute the
 57 finite difference approximation of the following form:

$$p(x, \xi, e) = \frac{f(x + te, \xi) - f(x - te, \xi)}{2t} \approx \langle \nabla_x f(x, \xi), e \rangle \quad (2)$$

58 Such approximation produces an estimate for the directional derivative of a noisy realization
 59 $f(\cdot, \xi)$ of the function f . As mentioned before, the approximation (2) can be used instead of
 60 the (stochastic) gradient in first-order methods. In the case of independent randomness, a
 61 large number of works are based on this idea. There are results for both non-smooth and
 62 smooth convex problems built on classical and accelerated gradient methods of Nesterov and
 63 Spokoiny [40]. In the scope of our paper, we are interested in the results for smooth strongly
 64 convex problems from [17], namely estimates on zero-order oracle calls to achieve ε -solution
 65 in terms of $\|x - x^*\|$: $\mathcal{O}(\frac{d\sigma_2^2}{\mu^2\varepsilon})$. Here σ_2 is introduced as the variance of the gradient, i.e. it is
 66 assumed that $\mathbb{E}_\xi \nabla f(x, \xi) = \nabla f(x)$ and $\mathbb{E}_\xi \|\nabla f(x, \xi) - \nabla f(x)\|^2 \leq \sigma_2^2$. The main limitation
 67 of two-point approach is that several evaluations with the same noise variable are required,
 68 which is well suited for problems like empirical risk optimization [34], but can be a major
 69 barrier for RL or online optimization.

70 In the *one-point feedback* setting, a more general stochasticity is assumed. In this case, each
 71 call to the zero-order oracle generates a new randomness. Now the approximation (1) looks
 72 as follows

$$p(x, \xi^\pm, e) = \frac{f(x + te, \xi^+) - f(x - te, \xi^-)}{2t} \quad (3)$$

73 Using different ξ^+ and ξ^- in (3) renders any conditions on the properties of $\nabla f(\cdot, \xi)$ useless.
 74 Instead, it is assumed that $\mathbb{E}_\xi f(x, \xi) = f(x)$ and $\mathbb{E}_\xi |f(x, \xi) - f(x)|^2 \leq \sigma_1^2$. With one-point
 75 feedback, the major problem is choosing the right shift t for the finite difference scheme.
 76 Picking it too small results in an amplification of the additive noise, and taking t too big
 77 leads to a poor gradient estimate. Because of this variance trade-off, the optimal rate for
 78 methods with one-point approximation is worse than for two-point feedback. In particular,
 79 for smooth strongly convex problems we have the following estimate on zero-order oracle
 80 calls [23]: $\mathcal{O}(\frac{d^2\sigma_1^2}{\mu^3\varepsilon^{\frac{1}{2}}})$.

with independent stochasticity. The key technique behind this acceleration is described in Section 2.1. The theory is also numerically validated in Section 3.

◊ **Non-smooth problems.** We also consider non-smooth problems with Markovian noise. Using the smoothing technique we come up with a corresponding upper bounds in this case, as shown in Figure 1. The details of these bounds are presented in Appendix B.2.

Figure 1: Summary of upper bounds. For notation, see Table 1

	Smooth		Non-smooth	
	IID	Markov.	IID	Markov.
FO	$\frac{\sigma_2^2}{\mu^2 \varepsilon}$ [45]	$\tau \frac{\sigma_2^2}{\mu^2 \varepsilon}$ [5]	$\frac{G^2}{\mu^2 \varepsilon}$ [50]	$\tau \frac{G^2}{\mu^2 \varepsilon}$ [14] ¹
ZO 2P	$d \frac{\sigma_2^2}{\mu^2 \varepsilon}$ [30]	$(d + \tau) \frac{\sigma_2^2}{\mu^2 \varepsilon}$	$d \frac{G^2}{\mu^2 \varepsilon}$ [22]	$(d + \tau) \frac{G^2}{\mu^2 \varepsilon}$
ZO 1P	$d^2 \frac{\sigma_1^2}{\mu^3 \varepsilon^2}$ [2] ²	$d(d + \tau) \frac{L \sigma_1^2}{\mu^3 \varepsilon^2}$	$d^2 \frac{\sigma_1^2 G^2}{\mu^4 \varepsilon^3}$ [23]	$d(d + \tau) \frac{\sigma_1^2 G^2}{\mu^4 \varepsilon^3}$

◊ **Computational efficiency.** First, as noted above, our method gives the same oracle complexity for any $\tau \leq d$. Moreover, if we assume that calling a zero-order oracle is d times cheaper than computing the corresponding gradient, then the gradient method with Markov noise will require resources proportionally to $d \cdot \tau$ — the cost of one oracle call is d and the complexity scales as τ for the first-order method from Figure 1. At the same time, the resource complexity of our zero-order method is proportional to $d + \tau$.

◊ **Lower bounds.** In Section 2.3 we establish the first information-theoretic lower bounds for solving Markovian optimization problems with one-point and two-point feedback. Our results match the convergence guarantee of our algorithm up to logarithmic factors, showing that the analysis is accurate and no further improvement is possible.

Table 1: Notations & Definitions

Sym.	Definition	Sym.	Definition
$\ \cdot\ , \langle \cdot, \cdot \rangle$	Norm, dot product, assumed Euclidean by default	ε	$\ x - x^*\ ^2$
\mathbf{Z}, \mathcal{Z}	Complete separable metric space, its Borel σ -algebra	d	Problem dimension
\mathbf{Q}	Markov kernel on $\mathbf{Z} \times \mathcal{Z}$	L	Gradient’s Lipschitz constant
$\mathbb{P}_\xi, \mathbb{E}_\xi$	Probability, Expectation under initial distribution ξ^3	μ	Strong convexity constant
$\{Z_k\}$	Canonical process with kernel \mathbf{Q}	G	Function’s Lipschitz constant
RB_2^d, RS_2^d	Uniform distribution on unit ℓ_2 -ball, -sphere	σ_1^2	$ F(x, Z) - f(x) ^2 \leq \sigma_1^2$
e	Random direction, $e \sim RS_2^d$	σ_2^2	$\ \nabla F(x, Z) - \nabla f(x)\ ^2 \leq \sigma_2^2$
$a_n \lesssim b_n$	$\exists c \in \mathbb{R}$ (problem-independent): $a_n \leq cb_n$ for all n	τ	Mixing time of Z
$a_n \simeq b_n$	$a_n \lesssim b_n$ and $b_n \lesssim a_n$	g, \hat{g}	Gradient estimators
$T = \tilde{\mathcal{O}}(S)$	$T \leq \text{poly}(\log S) \cdot S$ as $\varepsilon \rightarrow 0$	$f_t(x)$	$\mathbb{E}_r[f(x + tr)], r \sim RB_2^d$

2 Main results

We are now ready for a more formal presentation. In this paper, we study the minimization problem

$$\min_{x \in \mathbb{R}^d} f(x) := \mathbb{E}_{Z \sim \pi} [F(x, Z)], \quad (4)$$

where π is an unknown distribution and access to the function f (not to its gradient ∇f) is available through a stochastic one-point or two-point oracle $F(x, Z)$.

In our analysis, we will use a set of assumptions on the underlying function f and its oracle, starting with smoothness and convexity:

Assumption 1. *The function f is L -smooth on \mathbb{R}^d with $L > 0$, i.e., it is differentiable and there is a constant $L > 0$ such that the following inequality holds for all $x, y \in \mathbb{R}^d$:*

$$\|\nabla f(x) - \nabla f(y)\| \leq L\|x - y\|.$$

In the two-point feedback setting, we require the following generalization:

Assumption 1’. *For all $Z \in \mathcal{Z}$ the function $F(\cdot, Z)$ is L -smooth on \mathbb{R}^d .*

Note that the uniform 1’ implies 1.

¹The authors consider general convex case. Using standard restart technique, we get the corresponding bound in the strongly convex case.

²The noise is assumed to be point-independent.

³By construction, for any $A \in \mathcal{Z}$, we have $\mathbb{P}_\xi(Z_k \in A \mid Z_{k-1}) = \mathbf{Q}(Z_{k-1}, A)$, \mathbb{P}_ξ -a.s.

Assumption 2. The function f is μ -strongly convex on \mathbb{R}^d , i.e., it is continuously differentiable and there is a constant $\mu > 0$ such that the following inequality holds for all $x, y \in \mathbb{R}^d$:

$$\frac{\mu}{2} \|x - y\|^2 \leq f(x) - f(y) - \langle \nabla f(y), x - y \rangle. \quad (5)$$

We now turn to assumptions on the sequence of noise states $\{Z_i\}_{i=0}^\infty$. Specifically, we consider the case where $\{Z_i\}_{i=0}^\infty$ forms a time-homogeneous Markov chain. Let Q denote the corresponding Markov kernel. We impose the following assumption on Q to characterize its mixing properties:

Assumption 3. $\{Z_i\}_{i=0}^\infty$ is a stationary Markov chain on (Z, \mathcal{Z}) with Markov kernel Q and unique invariant distribution π . Moreover, Q is uniformly geometrically ergodic with mixing time $\tau \in \mathbb{N}$, i.e., for every $k \in \mathbb{N}$,

$$\Delta(Q^k) = \sup_{z, z' \in Z} (1/2) \|Q^k(z, \cdot) - Q^k(z', \cdot)\|_{\text{TV}} \leq (1/4)^{\lfloor k/\tau \rfloor}. \quad (6)$$

Assumption 3 is common in the literature on Markovian stochasticity [14, 12, 13, 5, 52]. It includes, for instance, irreducible aperiodic finite Markov chains [18]. The mixing time τ reflects how quickly the distribution of the chain approaches stationarity, providing a natural measure of the temporal dependence in the data.

Next, we specify our assumptions on the oracle. As discussed in Section 1.1, these assumptions differ based on the type of feedback.

Assumption 4 (for one-point). For all $x \in \mathbb{R}^d$ it holds that $\mathbb{E}_\pi[F(x, Z)] = f(x)$. Moreover, for all $Z \in Z$ and $x \in \mathbb{R}^d$ it holds that

$$|F(x, Z) - f(x)|^2 \leq \sigma_1^2,$$

Assumption 4' (for two-point). For all $x \in \mathbb{R}^d$ it holds that $\mathbb{E}_\pi[\nabla F(x, Z)] = \nabla f(x)$. Moreover, for all $Z \in Z$ and $x \in \mathbb{R}^d$ it holds that

$$\|\nabla F(x, Z) - \nabla f(x)\|^2 \leq \sigma_2^2.$$

Recent works on stochastic ZOO methods have considered milder assumptions, such as bounded variance (see Section 1.1). However, the uniform boundedness assumed in Assumptions 4 and 4', is standard in analyses under Markovian noise [14, 12, 13, 5, 52]. These assumptions can be relaxed under stronger conditions, e.g., uniform convexity and smoothness of $F(\cdot, Z)$ [18].

Assumptions 3 and 4 allow us to reduce the variance of the noise via batching, similarly to i.i.d. setting. This is captured in the following technical lemma:

Lemma 1. Let Assumptions 3 and 4(4') hold. Then for any $n \geq 1$ and $x \in \mathbb{R}^d$ and any initial distribution ξ on (Z, \mathcal{Z}) , we have

$$\mathbb{E}_\xi \left[\frac{1}{n} \sum_{i=1}^n F(x, Z_i) - f(x) \right]^2 \lesssim \frac{\tau}{n} \sigma_1^2, \quad \mathbb{E}_\xi \left\| \frac{1}{n} \sum_{i=1}^n \nabla F(x, Z_i) - \nabla f(x) \right\|^2 \lesssim \frac{\tau}{n} \sigma_2^2.$$

2.1 Batching technique

In this section, we describe the main tools used to establish the $(d + \tau)$ -type scaling of the error rate. We will focus on reducing the variance and bias of gradient estimators using a specialized batching approach.

We begin by fixing a common building block of our gradient estimators at a point x for both one-point and two-point feedback, as introduced in Section 1.1:

$$\hat{g}(x, Z^{(\pm)}, e) = d \cdot p(x, Z^{(\pm)}, e) \cdot e = e \cdot \begin{cases} d \frac{F(x + te, Z^+) - F(x - te, Z^-)}{2t} & \text{(one-point),} \\ d \frac{F(x + te, Z) - F(x - te, Z)}{2t} & \text{(two-point).} \end{cases}$$

These estimators exhibit a twofold randomness that affects how rapidly they concentrate around the true gradient, as we will discuss below.

For clarity, we focus our discussion on the one-point case, although our conclusions extend to the two-point case as well.

A widely used variance reduction technique is *mini-batching*, where one computes $F(x, Z_i)$ over a batch of noise variables $\{Z_i\}_{i=1}^n$. The mini-batch gradient estimator is given by:

$$\hat{g}_{mb}(x) = \frac{1}{n} \sum_{i=1}^n \hat{g}(x, Z_i^\pm, e) = e \cdot d \left(\overbrace{\frac{1}{n} \sum_{i=1}^n p(x, Z_i^\pm, e)}^{p_{mb}} \right).$$

Let us estimate the scaling of its variance $\mathbb{E}_e \mathbb{E}_Z \|\hat{g}_{mb} - \nabla f\|^2$ with the noise level σ_1^2 . As $E_Z \hat{g}_{mb} \approx d \frac{f(x+te) - f(x-te)}{2t} \approx d \langle \nabla f, e \rangle$ we would like to estimate the following for any fixed direction e :

$$\mathbb{E}_Z [d \cdot p_{mb}(x) - d \langle \nabla f, e \rangle]^2 \approx \frac{d^2}{t^2} \mathbb{E}_Z \left[\frac{1}{n} \sum_{i=1}^n F(x+te, Z_i^+) - f(x+te) \right]^2 \stackrel{(1)}{\approx} \frac{d^2 \tau \sigma_1^2}{n t^2}. \quad (7)$$

With that, we bound the variance:

$$\mathbb{E}_e \mathbb{E}_Z \|\hat{g}_{mb} - \nabla f\|^2 \gtrsim \mathbb{E}_e \mathbb{E}_Z \|\hat{g}_{mb} - \mathbb{E}_Z \hat{g}_{mb}\|^2 \approx \mathbb{E}_e \mathbb{E}_Z \|\hat{g}_{mb} - d \langle \nabla f, e \rangle\|^2 \stackrel{(7)}{\approx} \frac{d^2 \tau \sigma_1^2}{n t^2}. \quad (8)$$

Can the mini-batching scheme be improved?

This subsection explores an unexpected source of improvement that contradicts our initial hypothesis. Specifically, we identify an inefficiency in the current use of samples Z_i , which becomes evident from two perspectives. Equation (8) shows the variance scales as $\frac{\tau}{n}$. If we could reduce τ by a factor of k , we would need k -times fewer samples to maintain the same variance. This leads us to the idea of sparsified sampling. We partition the Markov noise chain $\{Z_i\}$ into k subchains $\{Z_{k \cdot i + r}\}$ for $r = 0 \dots k-1$. This corresponds to a mixing time of $\lceil \frac{\tau}{k} \rceil$ for each subchain (see (3)), effectively reducing temporal correlation - a natural consequence of sampling every k -th element of the original chain. Thus, sampling from any single subchain could yield a $\min(k, \tau)$ -fold reduction in the number of samples needed (although such procedure would still require all intermediate oracle calls, yielding no computational speedup).

For a concrete illustration of that inefficiency, consider a lazy Markov chain that remains in the same state for (an average of) τ steps before transitioning uniformly at random. In such a case, all oracle queries $F(x, Z)$ for a fixed x return the same value for τ consecutive steps. Therefore, retaining only every τ -th estimate \hat{g} would yield a mini-batch of equivalent quality.

In summary, we observe that the mini-batching scheme could, in principle, operate just as effectively by retaining only every k -th sample and discarding the rest. This might suggest that better utilization of the samples is possible. First order methods, nevertheless, are unable to exploit this redundancy (as shown by [5]'s lower bound) and are effectively forced to wait out the τ -step mixing window. In contrast, we can exploit this structure by querying finite differences along different directions to estimate the gradient better. Specifically, we construct d subchains, and use the sample from the r -th subchain $Z_{d \cdot i + r}$ to estimate r -th partial derivative $\frac{F(x+te_r, Z) - F(x-te_r, Z)}{2t}$, effectively restoring the full gradient coordinate-wise.

Let us estimate the resulting variance reduction. First, we achieve a d -fold reduction by reconstructing all d gradient coordinates. Second, each coordinate now operates on a chain with mixing time $\lceil \frac{\tau}{d} \rceil$, yielding an additional factor of $\min(d, \tau)$. However, because batches are now split across d coordinates, each batch is d times smaller than before, introducing a factor of d loss. The net variance reduction is therefore $\min(d, \tau)$, and the final scaling becomes $d \cdot \frac{d\tau}{\min(d, \tau)} = d \cdot \max(d, \tau) \simeq d(d + \tau)$.

Random directions

This insight can be extended to a simpler yet equally effective method. Instead of assigning directions deterministically, we associate each sample with a random direction $e \in RS_2^d$, forming the estimator:

$$\hat{g}_{rd}[n](x, Z, e) = \frac{1}{n} \sum_{i=1}^n \hat{g}(x, Z_i, e_i).$$

While the above discussion was intuitive, we now outline a more formal approach (see Lemma 5 for details). As lazy Markov chain is effectively equivalent to stochastic i.i.d. τ -point feedback setting, we follow Corollary 2 of [15], who decompose the total variance into two terms:

$$\mathbb{E}\|\hat{g}_{rd} - \nabla f(x)\|^2 \leq 2\mathbb{E}\|\hat{g}_{rd} - \mathbb{E}_e \hat{g}_{rd}\|^2 + 2\mathbb{E}\|\mathbb{E}_e \hat{g}_{rd} - \nabla f(x)\|^2.$$

Each of the two terms individually eliminates one factor from the $d^2\tau$ dependence.

The first term:

$$\begin{aligned} \mathbb{E}\|\hat{g}_{rd} - \mathbb{E}_e \hat{g}_{rd}\|^2 &= \mathbb{E}_Z \mathbb{E}_e \left\| \frac{1}{n} \sum_{i=1}^n \underbrace{[\hat{g}(x, Z_i, e_i) - \mathbb{E}_{e_i} \hat{g}(x, Z_i, e_i)]}_{\mathbb{E}_e[\cdot]=0, \text{ independent w.r.t. } e} \right\|^2 \\ &= \frac{1}{n^2} \sum_{i=1}^n \mathbb{E}\|\hat{g}(x, Z_i, e_i) - \mathbb{E}_{e_i} \hat{g}(x, Z_i, e_i)\|^2 \end{aligned}$$

is independent of τ since Assumption 4 bounds each term directly.

For the second term, we observe that $\mathbb{E}_e \hat{g}_{rd} = \mathbb{E}_e \hat{g}_{mb}$, and thus the bound involves $\mathbb{E}\|\mathbb{E}_e \hat{g}_{mb} - \nabla f(x)\|^2$. This is crucially different from the $d^2\tau$ dependence that appeared in the mini-batch case, when we considered $\mathbb{E}\|\hat{g}_{mb} - \nabla f(x)\|^2$. Intuitively, the expectation over directions helps recover the full gradient rather than a directional component, thereby reducing variance with respect to d .

Multilevel Monte Carlo

The estimator \hat{g}_{rd} is not our final construction. While it controls variance, the temporal correlation in noise may introduce significant bias. A well-established approach to mitigating this is MLMC, widely used in the statistical literature [27, 26], and more recently in gradient optimization [13, 5]. Here is our interpretation.

With parameters J, l, M, B from Table 2, $\{Z_i\}$ - $2^J l$ samples from Z and $\{e_i\}$ - random directions we introduce MLMC estimator:

$$\hat{g}_{ml}(x) = \hat{g}_{rd}[l](x) + \begin{cases} 2^J [\hat{g}_{rd}[2^J l](x) - \hat{g}_{rd}[2^{J-1} l](x)], & \text{if } 2^J \leq M \\ 0, & \text{otherwise} \end{cases}$$

\hat{g}_{ml} is our final gradient estimator, with the following guarantees:

Lemma 2 (for one-point). *Let Assumptions 1, 3 and 4 hold. For any initial distribution¹ ξ on (Z, \mathcal{Z}) the gradient estimates \hat{g}_{ml} satisfy $\mathbb{E}[\hat{g}_{ml}] = \mathbb{E}[\hat{g}_{rd}[2^{\lceil \log_2 M \rceil} l]]$. Moreover,*

$$\begin{aligned} \mathbb{E}\|\nabla f_t(x) - \hat{g}_{ml}(x)\|^2 &\lesssim \frac{d\|\nabla f(x)\|^2}{B} + \frac{d^2 L^2 t^2}{B} + \frac{d(d+\tau)\sigma_1^2}{Bt^2}, \\ \|\nabla f_t(x) - \mathbb{E}[\hat{g}_{ml}(x)]\|^2 &\lesssim \frac{d\tau\sigma_1^2}{t^2 BM}. \end{aligned}$$

One can note that although \hat{g}_{ml} requires, on average, $\mathbb{E}[2^J l B] = \log_2^2 M \cdot B$ oracle calls, the variance is only reduced by a factor of B . In contrast, the bias is reduced significantly - by a factor of BM .

2.2 Algorithm

We now present the full version of Algorithm 1, which incorporates the gradient estimators discussed in the previous section and uses a slightly modified variant of Nesterov's Accelerated Gradient Descent at its core.

While technically we prove four separate upper bounds covering both one- and two-point feedback under smooth and non-smooth assumptions, they follow the same scheme which we will illustrate in the one-point smooth case.

¹Note that \hat{g}_{ml} (specifically Z_1) indirectly depends on the chain's initial distribution. As our algorithm is going to repeatedly call \hat{g}_{ml} , next iteration's initial distribution is current iteration's final distribution. This fact makes the estimates correlated. We sidestep this problem by assuming any initial distribution.

Table 2: Parameters of Algorithm 1

Hyperparameters		Momentums		Batch hidden parameters	
γ	Stepsize, $\in (0; \frac{3}{4L}]$	β	$\sqrt{\frac{4p^2\mu\gamma}{3}}$	$2^J l$	Batch size. If $2^J l > M$, then 0
t	Approximation step	η	$\frac{3\beta}{2p\mu\gamma} = \sqrt{\frac{3}{\mu\gamma}}$	J	Random, $J \sim \text{Geom}(1/2)$
B	Batch size multiplier	θ	$\frac{p\eta^{-1}-1}{\beta p\eta^{-1}-1}$	M	Batch size limit, $M = \frac{1}{p} + \frac{2}{\beta}$
N	Number of iterations	p	See Appendix	l	$(\lfloor \log_2 M \rfloor + 1) \cdot B$

Lemma 4 establishes key properties of the smoothed objective function. Lemma 5 provides bounds on the bias and variance of the baseline estimator \hat{g}_{rd} . Lemma 2 then quantifies how the MLMC scheme amplifies or reduces these statistics. Finally, in Appendix D.4, we combine the results of these lemmas to prove the first part of Theorem 1, bounding Algorithm 1’s error. By tuning the parameters appropriately, we obtain the following iteration complexity bound:

Theorem 1. *Let Assumptions 1 to 4 hold, and consider problem (4) solved by Algorithm 1. Then, for any target accuracy ε and batch size multiplier B (see Tables 1 and 2 for notation), and for a suitable choice of γ, t, p , the number of oracle calls required to ensure $\mathbb{E}\|x^N - x^*\|^2 \leq \varepsilon$ is bounded by*

$$B \cdot \tilde{\mathcal{O}} \left(\max \left[1, \frac{d}{B} \right] \sqrt{\frac{L}{\mu}} \log \frac{1}{\varepsilon} + \frac{Ld(d+\tau)\sigma_1^2}{B\mu^3\varepsilon^2} \right) \quad \text{one-point oracle calls.}$$

Theorem 1’. *Let Assumptions 1’ to 4’ hold, and consider problem (4) solved by Algorithm 1. Then, for any target accuracy ε and batch size multiplier B (see Tables 1 and 2 for notation), and for a suitable choice of γ, t, p , the number of oracle calls required to ensure $\mathbb{E}\|x^N - x^*\|^2 \leq \varepsilon$ is bounded by*

$$B \cdot \tilde{\mathcal{O}} \left(\max \left[1, \frac{d}{B} \right] \sqrt{\frac{L}{\mu}} \log \frac{1}{\varepsilon} + \frac{(d+\tau)\sigma_2^2}{B\mu^2\varepsilon} \right) \quad \text{two-point oracle calls.}$$

Remark. The *iteration complexity* of the algorithm, i.e., the number of iterates x^k generated (equal to the oracle complexity divided by B), is bound by $\tilde{\mathcal{O}} \left(\sqrt{\frac{L}{\mu}} \log \frac{1}{\varepsilon} \right)$ as the batch size multiplier B goes to infinity. This matches the optimal convergence rates for optimization with *exact* gradients [38].

2.3 Lower bounds

Here we present theorems demonstrating that no algorithm can asymptotically outperform Algorithm 1 in the smooth, strongly convex setting with either one- or two-point feedback.

Theorem 2. *(Lower bounds) For any (possibly randomized) algorithm that solves the problem (4), there exists a function f that satisfies Assumptions 1 to 4 (1’ to 4’), s.t. in order to achieve ε -approximate solution in expectation $\mathbb{E}\|x^N - x^*\|^2 \leq \varepsilon$, the algorithm needs at least*

$$\Omega \left(\frac{d(d+\tau)\sigma_1^2}{\mu^2\varepsilon^2} \right) \quad \text{one-point or} \quad \Omega \left(\frac{(d+\tau)\sigma_2^2}{\mu^2\varepsilon} \right) \quad \text{two-point oracle calls.}$$

Remark. These results assume bounded second moments rather than uniform noise bounds. We explain how to adapt them to our setting, incurring only logarithmic overheads, in Appendix F.2.

Discussion. We now compare our results to existing work. Akhavan et al. [2] analyze a special case of the one-point setting where the noise is independent of the query points. This

aligns with our one-point oracle model and allows i.i.d. sampling as a Markov chain with fixed mixing time $\tau = 1$. The only factor they do not consider is σ_1^2 , which, however, appears in their proof with additional μ^2 factor if used with scaled Gaussian noise. We discuss this further in Appendix F.

In the work of Beznosikov et al. [5], a first-order Markovian oracle is considered, but the hard instance problem is a one-dimensional quadratic function, which makes first-order and zero-order information equivalent. Their result therefore corresponds to the $d = 1$ case in the two-point regime. Duchi et al. [15] provide tight lower bounds for general convex functions under two-point feedback. Their techniques can be extended to the strongly convex case by incorporating a shared quadratic component across the hard instances, as detailed in Appendix F, Theorem 10, yielding the bound we state for the two-point oracle with $\tau = 1$.

Our novel contribution lies in establishing a lower bound that scales as $d\tau$ in the one-point regime for large τ ; see Theorem 8. While our analysis relies on classical tools such as multidimensional hypothesis testing, the Markovian structure requires new bound on distances between joint distributions and the use of clipping. Detailed proofs, discussions, and further remarks on clipping appear in Appendix F.

3 Experiments

This section empirically supports our theoretical convergence rates and lower bounds, with particular focus on the stochastic component where we claim linear scaling in $d + \tau$ instead of $d\tau$.

Setup. Our setup repeats the problem we used to prove the lower bounds (see Appendix F and [51]). We consider a quadratic objective $f(x) = \frac{1}{2}\|x\|^2$ and a two-point Markovian oracle $F(x, Z) = f(x) + \langle x, Z \rangle$. The noise sequence $\{Z_i\}$ is a lazily updated standard Gaussian vector with variance σ_2^2 . Figure 2 illustrates how the optimization error of Algorithm 1 scales with mixing time, problem dimension, and different values of σ_2^2 .

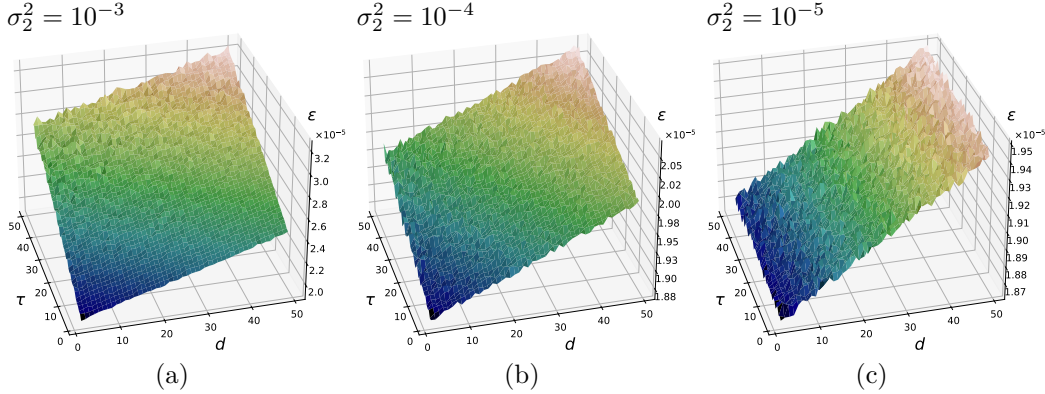


Figure 2: Optimization error $\varepsilon = \|x^N - x^*\|^2$ after $N = 10^3$ iterations. Starting point error $\|x_0 - x^*\|^2 = 10^{-2}$. Stepsize $\gamma = 10^{-3}$, $t = 10^{-5}$. The results are averaged over 10^4 runs.

Discussion. The results confirm the linear dependence of the error on both the problem dimension d and the mixing time τ . The noise parameter σ^2 controls the influence of the stochastic part. In Fig. (a), where $\sigma_2^2 = 10^{-3}$, the stochastic component dominates, while in Fig. (c), with $\sigma_2^2 = 10^{-5}$, it is negligible. Fig. (b) shows an intermediate regime that smoothly interpolates between the two, yet maintains the linear scaling. The deterministic part (c) shows no dependence on mixing time, but grows linearly with d , which aligns with our theory (Theorem 1'). The stochastic part (a) scales as $(d + \tau)$, also matching the bound from the Theorem 1'.

References

- [1] Arya Akhavan, Massimiliano Pontil, and Alexandre Tsybakov. Exploiting higher order smoothness in derivative-free optimization and continuous bandits. *Advances in Neural Information Processing Systems*, 33:9017–9027, 2020.

- [2] Arya Akhavan, Evgenii Chzhen, Massimiliano Pontil, and Alexandre B Tsybakov. Gradient-free optimization of highly smooth functions: improved analysis and a new algorithm. *Journal of Machine Learning Research*, 25(370):1–50, 2024.
- [3] Peter Auer. Finite-time analysis of the multiarmed bandit problem. *Machine Learning*, 47:235–256, 2002.
- [4] El Houcine Bergou, Eduard Gorbunov, and Peter Richtárik. Stochastic three points method for unconstrained smooth minimization. *SIAM Journal on Optimization*, 30(4): 2726–2749, 2020.
- [5] Aleksandr Beznosikov, Sergey Samsonov, Marina Sheshukova, Alexander Gasnikov, Alexey Naumov, and Eric Moulines. First order methods with markovian noise: from acceleration to variational inequalities. *Advances in Neural Information Processing Systems*, 36, 2024.
- [6] Jalaj Bhandari, Daniel Russo, and Raghav Singal. A finite time analysis of temporal difference learning with linear function approximation. In *Conference on learning theory*, pages 1691–1692. PMLR, 2018.
- [7] Sébastien Bubeck, Nicolo Cesa-Bianchi, et al. Regret analysis of stochastic and non-stochastic multi-armed bandit problems. *Foundations and Trends® in Machine Learning*, 5(1):1–122, 2012.
- [8] Pin-Yu Chen, Huan Zhang, Yash Sharma, Jinfeng Yi, and Cho-Jui Hsieh. 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*, pages 15–26, 2017.
- [9] Krzysztof Choromanski, Mark Rowland, Vikas Sindhwani, Richard Turner, and Adrian Weller. Structured evolution with compact architectures for scalable policy optimization. In *International Conference on Machine Learning*, pages 970–978. PMLR, 2018.
- [10] Aymeric Dieuleveut, Nicolas Flammarion, and Francis Bach. Harder, better, faster, stronger convergence rates for least-squares regression. *Journal of Machine Learning Research*, 18(101):1–51, 2017.
- [11] Thinh T Doan. Finite-time analysis of markov gradient descent. *IEEE Transactions on Automatic Control*, 68(4):2140–2153, 2022.
- [12] Thinh T Doan, Lam M Nguyen, Nhan H Pham, and Justin Romberg. Convergence rates of accelerated markov gradient descent with applications in reinforcement learning. *arXiv preprint arXiv:2002.02873*, 2020.
- [13] Ron Dorfman and Kfir Yehuda Levy. Adapting to mixing time in stochastic optimization with markovian data. In *International Conference on Machine Learning*, pages 5429–5446. PMLR, 2022.
- [14] John C Duchi, Alekh Agarwal, Mikael Johansson, and Michael I Jordan. Ergodic mirror descent. *SIAM Journal on Optimization*, 22(4):1549–1578, 2012.
- [15] John C Duchi, Michael I Jordan, Martin J Wainwright, and Andre Wibisono. Optimal rates for zero-order convex optimization: The power of two function evaluations. *IEEE Transactions on Information Theory*, 61(5):2788–2806, 2015.
- [16] Alain Durmus, Eric Moulines, Alexey Naumov, Sergey Samsonov, and Hoi-To Wai. On the stability of random matrix product with markovian noise: Application to linear stochastic approximation and td learning. In *Conference on Learning Theory*, pages 1711–1752. PMLR, 2021.
- [17] Pavel Dvurechensky, Eduard Gorbunov, and Alexander Gasnikov. An accelerated directional derivative method for smooth stochastic convex optimization. *European Journal of Operational Research*, 290(2):601–621, 2021.

- [18] Mathieu Even. Stochastic gradient descent under markovian sampling schemes. In *International Conference on Machine Learning*, pages 9412–9439. PMLR, 2023.
- [19] Maryam Fazel, Rong Ge, Sham Kakade, and Mehran Mesbahi. Global convergence of policy gradient methods for the linear quadratic regulator. In *International conference on machine learning*, pages 1467–1476. PMLR, 2018.
- [20] Abraham D. Flaxman, Adam Tauman Kalai, and H. Brendan McMahan. Online convex optimization in the bandit setting: gradient descent without a gradient. In *Proceedings of the Sixteenth Annual ACM-SIAM Symposium on Discrete Algorithms*, SODA '05, page 385–394, USA, 2005. Society for Industrial and Applied Mathematics. ISBN 0898715857.
- [21] Alexander Gasnikov, Darina Dvinskikh, Pavel Dvurechensky, Eduard Gorbunov, Aleksandr Beznosikov, and Alexander Lobanov. *Randomized Gradient-Free Methods in Convex Optimization*, pages 1–15. Springer International Publishing, Cham, 2020. ISBN 978-3-030-54621-2. doi: 10.1007/978-3-030-54621-2_859-1. URL https://doi.org/10.1007/978-3-030-54621-2_859-1.
- [22] Alexander Gasnikov, Anton Novitskii, Vasilii Novitskii, Farshed Abdukhakimov, Dmitry Kamzolov, Aleksandr Beznosikov, Martin Takac, Pavel Dvurechensky, and Bin Gu. The power of first-order smooth optimization for black-box non-smooth problems. In Kamalika Chaudhuri, Stefanie Jegelka, Le Song, Csaba Szepesvari, Gang Niu, and Sivan Sabato, editors, *Proceedings of the 39th International Conference on Machine Learning*, volume 162 of *Proceedings of Machine Learning Research*, pages 7241–7265. PMLR, 17–23 Jul 2022. URL <https://proceedings.mlr.press/v162/gasnikov22a.html>.
- [23] Alexander V Gasnikov, Ekaterina A Krymova, Anastasia A Lagunovskaya, Ilnura N Usmanova, and Fedor A Fedorenko. Stochastic online optimization. single-point and multi-point non-linear multi-armed bandits. convex and strongly-convex case. *Automation and remote control*, 78:224–234, 2017.
- [24] Saeed Ghadimi and Guanghui Lan. Stochastic first-and zeroth-order methods for nonconvex stochastic programming. *SIAM journal on optimization*, 23(4):2341–2368, 2013.
- [25] Saeed Ghadimi and Guanghui Lan. Accelerated gradient methods for nonconvex nonlinear and stochastic programming. *Mathematical Programming*, 156(1):59–99, 2016.
- [26] Michael B. Giles. Multilevel monte carlo path simulation. *Operations Research*, 56(3): 607–617, 2008. doi: 10.1287/opre.1070.0496. URL <https://doi.org/10.1287/opre.1070.0496>.
- [27] Peter W. Glynn and Chang-Han Rhee. Exact estimation for markov chain equilibrium expectations. *Journal of Applied Probability*, 51A:377–389, 2014. ISSN 00219002. URL <http://www.jstor.org/stable/43284129>.
- [28] Ian Goodfellow, Yoshua Bengio, and Aaron Courville. *Deep Learning*. MIT Press, 2016. <http://www.deeplearningbook.org>.
- [29] Eduard Gorbunov, Pavel Dvurechensky, and Alexander Gasnikov. An accelerated method for derivative-free smooth stochastic convex optimization. *SIAM Journal on Optimization*, 32(2):1210–1238, 2022. doi: 10.1137/19M1259225. URL <https://doi.org/10.1137/19M1259225>.
- [30] Elad Hazan and Satyen Kale. Beyond the regret minimization barrier: Optimal algorithms for stochastic strongly-convex optimization. *Journal of Machine Learning Research*, 15(71):2489–2512, 2014. URL <http://jmlr.org/papers/v15/hazan14a.html>.
- [31] Bjorn Johansson, Maben Rabi, and Mikael Johansson. A simple peer-to-peer algorithm for distributed optimization in sensor networks. In *2007 46th IEEE Conference on Decision and Control*, pages 4705–4710, 2007. doi: 10.1109/CDC.2007.4434888.

- [32] J. Kiefer. Sequential minimax search for a maximum. *Proceedings of the American Mathematical Society*, 4(3):502–506, 1953. ISSN 00029939, 10886826. URL <http://www.jstor.org/stable/2032161>.
- [33] Xiangru Lian, Yijun Huang, Yuncheng Li, and Ji Liu. Asynchronous parallel stochastic gradient for nonconvex optimization. *Advances in neural information processing systems*, 28, 2015.
- [34] Sijia Liu, Bhavya Kailkhura, Pin-Yu Chen, Paishun Ting, Shiyu Chang, and Lisa Amini. Zeroth-order stochastic variance reduction for nonconvex optimization. *Advances in Neural Information Processing Systems*, 31, 2018.
- [35] Cassio G. Lopes and Ali H. Sayed. Incremental adaptive strategies over distributed networks. *IEEE Transactions on Signal Processing*, 55(8):4064–4077, 2007. doi: 10.1109/TSP.2007.896034.
- [36] Eric Moulines and Francis Bach. Non-asymptotic analysis of stochastic approximation algorithms for machine learning. In J. Shawe-Taylor, R. Zemel, P. Bartlett, F. Pereira, and K.Q. Weinberger, editors, *Advances in Neural Information Processing Systems*, volume 24. Curran Associates, Inc., 2011. URL https://proceedings.neurips.cc/paper_files/paper/2011/file/40008b9a5380fcacce3976bf7c08af5b-Paper.pdf.
- [37] Deanna Needell, Rachel Ward, and Nati Srebro. Stochastic gradient descent, weighted sampling, and the randomized kaczmarz algorithm. *Advances in neural information processing systems*, 27, 2014.
- [38] Yurii Nesterov. A method for solving the convex programming problem with convergence rate $o(1/k^2)$. In *Doklad nauk Sssr*, volume 269, page 543, 1983.
- [39] Yurii Nesterov. Efficiency of coordinate descent methods on huge-scale optimization problems. *SIAM Journal on Optimization*, 22(2):341–362, 2012.
- [40] Yurii Nesterov and Vladimir Spokoiny. Random gradient-free minimization of convex functions. *Foundations of Computational Mathematics*, 17(2):527–566, 2017.
- [41] Donald J Newman. Location of the maximum on unimodal surfaces. *Journal of the ACM (JACM)*, 12(3):395–398, 1965.
- [42] J. Nocedal and S. Wright. *Numerical Optimization*. Springer Series in Operations Research and Financial Engineering. Springer New York, 2006. ISBN 9780387227429. URL <https://books.google.ru/books?id=7wDpBwAAQBAJ>.
- [43] Boris Polyak. *Introduction to Optimization*. Optimization Software - Inc., Publications Division, 1987.
- [44] Yuyang Qiu, Uday Shanbhag, and Farzad Yousefian. Zeroth-order methods for non-differentiable, nonconvex, and hierarchical federated optimization. In A. Oh, T. Naumann, A. Globerson, K. Saenko, M. Hardt, and S. Levine, editors, *Advances in Neural Information Processing Systems*, volume 36, pages 3425–3438. Curran Associates, Inc., 2023. URL https://proceedings.neurips.cc/paper_files/paper/2023/file/0a70c9cd8179fe6f8f6135fafa2a8798-Paper-Conference.pdf.
- [45] Herbert Robbins and Sutton Monro. A stochastic approximation method. *The annals of mathematical statistics*, pages 400–407, 1951.
- [46] Tim Salimans, Jonathan Ho, Xi Chen, Szymon Sidor, and Ilya Sutskever. Evolution strategies as a scalable alternative to reinforcement learning. *arXiv preprint arXiv:1703.03864*, 2017.
- [47] Bobak Shahriari, Kevin Swersky, Ziyu Wang, Ryan P. Adams, and Nando de Freitas. Taking the human out of the loop: A review of bayesian optimization. *Proceedings of the IEEE*, 104(1):148–175, 2016. doi: 10.1109/JPROC.2015.2494218.

- [48] Shai Shalev-Shwartz and Shai Ben-David. *Understanding machine learning: From theory to algorithms*. Cambridge university press, 2014.
- [49] Ohad Shamir. An optimal algorithm for bandit and zero-order convex optimization with two-point feedback. *Journal of Machine Learning Research*, 18(52):1–11, 2017.
- [50] Ohad Shamir and Tong Zhang. Stochastic gradient descent for non-smooth optimization: Convergence results and optimal averaging schemes. In Sanjoy Dasgupta and David McAllester, editors, *Proceedings of the 30th International Conference on Machine Learning*, volume 28 of *Proceedings of Machine Learning Research*, pages 71–79, Atlanta, Georgia, USA, 17–19 Jun 2013. PMLR. URL <https://proceedings.mlr.press/v28/shamir13.html>.
- [51] Alexander Shapiro, Darinka Dentcheva, and Andrzej Ruszczyński. *Lectures on Stochastic Programming*. Society for Industrial and Applied Mathematics, 2009. doi: 10.1137/1.9780898718751. URL <https://epubs.siam.org/doi/abs/10.1137/1.9780898718751>.
- [52] Vladimir Solodkin, Andrew Veprikov, and Aleksandr Beznosikov. Methods for optimization problems with markovian stochasticity and non-euclidean geometry. *arXiv preprint arXiv:2408.01848*, 2024.
- [53] Sebastian U. Stich. Unified optimal analysis of the (stochastic) gradient method, 2019. URL <https://arxiv.org/abs/1907.04232>.
- [54] Ben Taskar, Vassil Chatalbashev, Daphne Koller, and Carlos Guestrin. Learning structured prediction models: A large margin approach. In *Proceedings of the 22nd international conference on Machine learning*, pages 896–903, 2005.
- [55] Alexandre B. Tsybakov. *Lower bounds on the minimax risk*, pages 77–135. Springer New York, New York, NY, 2009. ISBN 978-0-387-79052-7. doi: 10.1007/978-0-387-79052-7_2. URL https://doi.org/10.1007/978-0-387-79052-7_2.
- [56] Sharan Vaswani, Francis Bach, and Mark Schmidt. Fast and faster convergence of sgd for over-parameterized models and an accelerated perceptron. In *The 22nd international conference on artificial intelligence and statistics*, pages 1195–1204. PMLR, 2019.
- [57] Jian Wu, Saul Toscano-Palmerin, Peter I Frazier, and Andrew Gordon Wilson. Practical multi-fidelity bayesian optimization for hyperparameter tuning. In *Uncertainty in Artificial Intelligence*, pages 788–798. PMLR, 2020.
- [58] David B Yudin and Arkadi S Nemirovskii. Informational complexity and efficient methods for the solution of convex extremal problems. *Matekon*, 13(2):22–45, 1976.
- [59] Yawei Zhao. Markov chain mirror descent on data federation. *arXiv preprint arXiv:2309.14775*, 2023.

A Appendix overview

In this section, the overall structure of the technical appendices is presented.

In Appendix B, we introduce the additional adversarial robustness of the Algorithm 1 and present a formal statements for our results in the non-smooth case.

In Appendix C, we define the shorthanded notation used in the proof of upper bounds.

In Appendices D and E, we gradually introduce all lemmas and proofs of our theorems in one-point and two-point setting respectively, for both smooth and non-smooth problems.

In Appendix F we present our lower bounds and provide a more detailed overview of the related results.

Finally, in Appendix G, we formally state the common-knowledge facts that we use.

B Additional results

B.1 Adversarial noise

In addition to the main results that show optimal scaling with the stochastic noise, we also prove a *robustness* of our algorithm. Precisely, the oracle F considered in this paper may return its values with an additive, non-random, potentially adversarial error $\Delta(x) \leq \Delta$.

$$\hat{F}(x, Z) = F(x, Z) + \Delta(x). \quad (9)$$

We will prove that this have no effect of the convergence guarantees of our algorithm for any Δ within a tolerable threshold. This threshold varies between smooth and non-smooth case, but not between one-point and two-point settings. The precise bounds are presented in the theorems in Appendices D and E.

B.2 Non-smooth

In the non-smooth case, we consider a similar set of assumptions, however f is no longer necessarily smooth or even differentiable.

Assumption 5. *The function f is μ -strongly convex on \mathbb{R}^d , i.e., there is a constant $\mu > 0$ such that the following inequality holds for all $x, y \in \mathbb{R}^d$ and $\lambda \in [0; 1]$:*

$$f(\lambda x + (1 - \lambda)y) \leq \lambda f(x) + (1 - \lambda)f(y) - \lambda(1 - \lambda)\frac{\mu}{2}\|x - y\|^2$$

Assumption 6. *The function f is G -Lipschitz on \mathbb{R}^d , i.e., there is a constant $G > 0$ such that the following inequality holds for all $x, y \in \mathbb{R}^d$:*

$$|f(x) - f(y)| \leq G\|x - y\|.$$

Again, for the two-point case, we need the generalization:

Assumption 6'. *For all $Z \in \mathcal{Z}$ the function $F(\cdot, Z)$ is G -Lipschitz on \mathbb{R}^d .*

Regarding the noise levels, we keep Assumption 4 for the one-point case.

For the two-point case, however, we cannot keep Assumption 4', as f is no longer differentiable. Instead, we will also use function unbiasedness. In that case, we will not use any additional assumptions on noise variance, as gradient of the smoothed function is already bounded by G as it is Lipschitz and differentiable.

Assumption 7. *For all $x \in \mathbb{R}^d$ it holds that $\mathbb{E}_\pi F(x, Z) = f(x)$.*

Theorem 3. *Let Assumptions 3, 4, 5 and 6 hold, and consider problem (4) solved by Algorithm 1. Then, for any target accuracy ε and batch size multiplier B (see Tables 1 and 2 for notation), and for a suitable choice of γ, t, p , the number of oracle calls required to ensure $\mathbb{E}\|x^N - x^*\|^2 \leq \varepsilon$ is bounded by*

$$B \cdot \tilde{\mathcal{O}} \left(\sqrt{\frac{dG^2}{\mu^2\varepsilon}} \log \frac{1}{\varepsilon} + \frac{d(d + \tau)\sigma_1^2 G^2}{B\mu^4 \varepsilon^3} + \frac{dG^2}{B\mu^2 \varepsilon} \right) \quad \text{one-point oracle calls.}$$

564 We present the following theorems.

565 **Theorem 3'.** Assume Assumption 5, 6', 3 and 7 hold, and consider problem (4) solved by
 566 Algorithm 1. Then, for any target accuracy ε and batch size multiplier B (see Tables 1 and 2
 567 for notation), and for a suitable choice of γ, t, p , the number of oracle calls required to ensure
 568 $\mathbb{E}\|x^N - x^*\|^2 \leq \varepsilon$ is bounded by

$$B \cdot \tilde{O} \left(\sqrt{\frac{\sqrt{d}G^2}{\mu^2\varepsilon}} \log \frac{1}{\varepsilon} + \frac{(d + \tau)G^2}{B\mu^2\varepsilon} \right) \quad \text{two-point oracle calls.}$$

569 As we can see, there is no dependence on the mixing time as long as it is less then the
 570 dimension of the problem. Our results coincide with previous work under i.i.d. noise when
 571 applied with $\tau = 1$, as previously claimed in Figure 1.

572 C Notations and definitions.

573 In this section we define the shorthand notation used in the proof of upper bounds. For
 574 general notations and definitions, see Tables 1 and 2.

575 Markovian error:

$$h(x, Z) := F(x, Z) - f(x) \quad (10)$$

576 Single sample gradient estimators:

$$\hat{g}_i := d \frac{\hat{F}(x + te_i, Z_i^{(+)}) - \hat{F}(x - te_i, Z_i^{(-)})}{2t} e_i \quad (11)$$

$$\tilde{g}_i := d \frac{F(x + te_i, Z_i^{(+)}) - F(x - te_i, Z_i^{(-)})}{2t} e_i \quad (12)$$

$$\stackrel{(10)}{=} d \frac{f(x + te_i) + h(x + te_i, Z_i^{(+)}) - f(x - te_i) - h(x - te_i, Z_i^{(-)})}{2t} e_i$$

$$g_i := d \frac{f(x + te_i) - f(x - te_i)}{2t} e_i \quad (13)$$

577 Batched gradient estimators:

$$\hat{g}^j := \hat{g}_{rd}[2^j l] = \frac{1}{2^j l} \sum_{i=1}^{2^j l} \hat{g}_i \quad (14)$$

(Not to be confused with \hat{g}^k , which is \hat{g}_{ml} calculated on k -th iteration)

$$\tilde{g}^j := \frac{1}{2^j l} \sum_{i=1}^{2^j l} \tilde{g}_i \quad (15)$$

$$g^j := \frac{1}{2^j l} \sum_{i=1}^{2^j l} g_i \quad (16)$$

578 Directional gradients:

$$\nabla_{e_i} f(x_0) := d \langle \nabla f(x_0), e_i \rangle e_i \quad (17)$$

$$\nabla_{e_i} F_i := d \langle \nabla F(x, Z_i), e_i \rangle e_i \quad (18)$$

579 Misc:

$$\mathbb{E}_e := \mathbb{E}_{e_1, e_2, \dots, e_{2^j l}} \quad (19)$$

$$\mathbb{E}_Z := \mathbb{E}_{Z_1, Z_2, \dots, Z_{2^j l}}, \text{ where } Z_1 \sim \xi - \text{arbitrary initial distribution on } (Z, \mathcal{Z})$$

$$\mathbb{E} := \mathbb{E}_Z \mathbb{E}_e$$

$$\mathcal{F}_k := \sigma(x^1, x^2, \dots, x^k) - \text{sigma algebra of first } k \text{ iterations}$$

$$\mathbb{E}_k[\cdot] := \mathbb{E}[\cdot | \mathcal{F}_k]$$

580

$$r^N := \frac{1}{\mu} (f(x_f^N) - f(x^*)) + \|x^N - x^*\|^2 \quad (20)$$

581 D Proofs of one-point results

582 D.1 Markov variance reduction

583 **Lemma 3** (Extended version of Lemma 1). *Let Assumptions 3 and 4(4') hold. Then for*
 584 *any $n \geq 1$ and $x \in \mathbb{R}^d$ and any initial distribution ξ on (Z, \mathcal{Z}) , we have*

$$\mathbb{E}_Z \left[\left(\frac{1}{n} \sum_{i=1}^n \mathbb{E}_{e_i} [h(x + te_i, Z_i) e_i] \right)^2 \right] \lesssim \frac{\tau}{dn} \sigma_1^2, \quad (21)$$

$$\mathbb{E}_Z \left[\left\| \frac{1}{n} \sum_{i=1}^n \nabla F(x, Z_i) - \nabla f(x) \right\|^2 \right] \lesssim \frac{\tau}{n} \sigma_2^2, \quad (22)$$

585 *Proof.* The proof of (22) can be found in Lemma 1 of Beznosikov et al. [5].

586 The proof under Assumption 4 relies on the fact that aforementioned Lemma 1 requires just
 587 the two following conditions from the stochastic realizations $\nabla F(x, Z_i)$:

$$\begin{cases} \mathbb{E}_\pi \nabla F(x, Z_i) = \nabla f(x) \\ \|\nabla F(x, Z_i) - \nabla f(x)\|^2 \leq \sigma_2^2 \end{cases}$$

588 Denote $h_t(x, Z_i) := \mathbb{E}_e [h(x + te, Z_i) e], e \sim RS_2^d(1)$.

589 Thus (21) $\Leftrightarrow \mathbb{E}_Z \left[\left(\frac{1}{n} \sum_{i=1}^n h_t(x, Z_i) \right)^2 \right] \lesssim \frac{\tau}{n} \frac{\sigma_1^2}{d} \Leftrightarrow \begin{cases} \mathbb{E}_\pi h_t(x, Z_i) = 0 \\ \|h_t(x, Z_i)\|^2 \lesssim \frac{\sigma_1^2}{d} \end{cases}$

590 Let's prove both of these equations, starting with unbiasedness:

$$\mathbb{E}_\pi h_t(x, Z_i) = \mathbb{E}_\pi \mathbb{E}_e [h(x + te, Z_i) e] = \mathbb{E}_e \mathbb{E}_\pi [h(x + te, Z_i) e] \stackrel{(4)}{=} \mathbb{E}_e 0 = 0$$

591

$$\begin{aligned} \|h_t(x, Z_i)\|^2 &= \|\mathbb{E}_e [h(x + te, Z_i) e]\|^2 \\ &= \langle \mathbb{E}_e [h(x + te, Z_i) e], h_t(x, Z_i) \rangle \\ &\stackrel{\textcircled{1}}{=} \mathbb{E}_e [h(x + te, Z_i) \cdot \langle e, h_t(x, Z_i) \rangle] \\ &\stackrel{\textcircled{2}}{\leq} \sqrt{\mathbb{E}_e h(x + te, Z_i)^2} \cdot \sqrt{\mathbb{E}_e \langle e, h_t(x, Z_i) \rangle^2} \\ &\stackrel{(81)}{=} \sqrt{\mathbb{E}_e h(x + te, Z_i)^2} \cdot \sqrt{\frac{1}{d} \|h_t(x, Z_i)\|^2} \\ &\stackrel{(4)}{\leq} \sqrt{\sigma_1^2} \cdot \sqrt{\frac{1}{d} \|h_t(x, Z_i)\|^2}, \end{aligned}$$

592 where $\textcircled{1}$ holds since $h_t(x, Z_i)$ does not depend on e , and $\textcircled{2}$ is a Cauchy-Schwartz inequality
 593 for the following dot product: $\langle x(e), y(e) \rangle := \mathbb{E}_e [x \cdot y]$.

594 To conclude the proof we square the inequality we got:

595

$$\|h_t(x, Z_i)\|^2 \leq \frac{\sqrt{\sigma_1^2}}{\sqrt{d}} \cdot \sqrt{\|h_t(x, Z_i)\|^2} \Rightarrow \|h_t(x, Z_i)\|^2 \leq \frac{\sigma_1^2}{d}.$$

596

□

597 D.2 Properties of smoothed function

598 The following lemma establishes key properties of the l_2 -ball smoothed function

599 **Lemma 4.** *Assume f is convex. Then the following holds for all $x \in \mathbb{R}^d$*

$$\begin{aligned} &\text{If } f \text{ is } L\text{-smooth} / G\text{-Lipschitz} / \mu\text{-strongly convex [Assumptions 1, 2 and 6],} \\ &\text{then } f_t \text{ from (1) is also } L\text{-smooth} / G\text{-Lipschitz} / \mu\text{-strongly convex.} \end{aligned} \quad (23)$$

$$\nabla f_t(x) = \mathbb{E}_e [g(x)], \quad (24)$$

$$f_t(x) \geq f(x), \quad (25)$$

If f is additionally G -Lipschitz:

$$f_t(x) \leq f(x) + Gt, \quad (26)$$

$$f_t \text{ is } L\text{-smooth with } L = \frac{\sqrt{d}G}{t}, \quad (27)$$

If f is additionally L -smooth:

$$f_t(x) \leq f(x) + Lt^2, \quad (28)$$

$$\|\nabla f(x) - \nabla f_t(x)\|^2 \leq L^2 t^2, \quad (29)$$

$$\|\nabla f_t(x)\|^2 \geq \frac{1}{2} \|\nabla f(x)\|^2 - L^2 t^2. \quad (30)$$

Proof. Proving (23), we start with G -Lipschitzness:

$$\begin{aligned} |f_t(x) - f_t(y)| &= |\mathbb{E}_r [f(x + tr) - f(y + tr)]| \\ &\stackrel{(80)}{\leq} \mathbb{E}_r |f(x + tr) - f(y + tr)| \\ &\stackrel{(6)}{\leq} \mathbb{E}_r G \|x - y\| = G \|x - y\|. \end{aligned}$$

Next, L -smoothness is analogous. Finally, μ -strong convexity of f_t , (24), (25) and (28) are proven in Lemmas A2-A3 of [1].

(26) and (27) can be seen in section 4.1 of Gasnikov et al. [21].

We prove the rest of inequalities in order.

Proof of (29):

$$\begin{aligned} \|\nabla f(x) - \nabla f_t(x)\|^2 &= \|\nabla f(x) - \mathbb{E}_r \nabla f(x + tr)\|^2 \\ &= \|\mathbb{E}_r [\nabla f(x) - \nabla f(x + tr)]\|^2 \\ &\stackrel{(80)}{\leq} \mathbb{E}_r \|\nabla f(x) - \nabla f(x + tr)\|^2 \\ &\stackrel{(1)}{\leq} \mathbb{E}_r L^2 t^2 = L^2 t^2. \end{aligned}$$

613

614 Proof of (30):

$$\begin{aligned} \|\nabla f_t(x)\|^2 &= \|\nabla f(x) + [\nabla f_t(x) - \nabla f(x)]\|^2 \\ &\stackrel{\textcircled{1}}{\geq} \frac{1}{2} \|\nabla f(x)\|^2 - \|\nabla f_t(x) - \nabla f(x)\|^2 \\ &\stackrel{(29)}{\geq} \frac{1}{2} \|\nabla f(x)\|^2 - L^2 t^2, \end{aligned}$$

615 where $\textcircled{1}$ uses that $\|a + b\|^2 \geq 1/2 \|a\|^2 - \|b\|^2$. \square

616 D.3 Inequalities for gradient approximation

617 **Lemma 5.** Assume Assumption 1, Assumption 3 and Assumption 4. Then the following
618 inequalities hold for any initial distribution ξ on $(\mathcal{Z}, \mathcal{Z})$ and for all $x \in \mathbb{R}^d$:

$$\|\hat{g}^j - \tilde{g}^j\|^2 \leq \frac{d^2 \Delta^2}{t^2}, \quad (31)$$

$$\mathbb{E} \|\tilde{g}_i - g_i\|^2 \leq \frac{d^2 \sigma_1^2}{t^2}, \quad (32)$$

$$\mathbb{E} \|\mathbb{E}_e [\tilde{g}^j - g^j]\|^2 \leq \frac{d C_1 \tau \sigma_1^2}{t^2 2^j l}, \quad (33)$$

$$\mathbb{E}\|g_i - \nabla_{e_i} f\|^2 \leq \frac{d^2 L^2 t^2}{4}, \quad (34)$$

$$\mathbb{E}\|\tilde{g}^j - \mathbb{E}_e \tilde{g}^j\|^2 \leq \frac{3}{2^j l} \left[\frac{d^2 \sigma_1^2}{t^2} + \frac{d^2 L^2 t^2}{4} + d\|\nabla f\|^2 \right], \quad (35)$$

$$\mathbb{E}\|\tilde{g}^j - \mathbb{E}_e g^j\|^2 \lesssim \frac{d(d+\tau)\sigma_1^2}{t^2 2^j l} + \frac{d^2 L^2 t^2}{2^j l} + \frac{d\|\nabla f\|^2}{2^j l}, \quad (36)$$

$$\mathbb{E}\|\hat{g}^j - \nabla f_t\|^2 \lesssim \frac{d^2 \Delta^2}{t^2} + \frac{d(d+\tau)\sigma_1^2}{t^2 2^j l} + \frac{d^2 L^2 t^2}{2^j l} + \frac{d\|\nabla f\|^2}{2^j l}, \quad (37)$$

$$\|\mathbb{E} \hat{g}^j - \nabla f_t\|^2 \leq \frac{2d^2 \Delta^2}{t^2} + \frac{2dC_1 \tau \sigma_1^2}{t^2 2^j l}. \quad (38)$$

$$(39)$$

619 *Proof.* We prove all estimates one by one, starting with (31):

$$\begin{aligned} \|\hat{g}^j - \tilde{g}^j\|^2 &\stackrel{(14),(15)}{=} \left\| \frac{1}{2^j l} \sum_{i=1}^{2^j l} [\hat{g}_i - \tilde{g}_i] \right\|^2 \\ &\stackrel{(11),(12)}{=} \frac{d^2}{4t^2} \left\| \frac{1}{2^j l} \sum_{i=1}^{2^j l} [\hat{F}(x + te_i, Z_i^+) - \hat{F}(x - te_i, Z_i^-) \right. \\ &\quad \left. - F(x + te_i, Z_i^+) + F(x - te_i, Z_i^-)] e_i \right\|^2 \\ &\stackrel{(9)}{=} \frac{d^2}{4t^2} \left\| \frac{1}{2^j l} \sum_{i=1}^{2^j l} [\Delta(x + te_i) - \Delta(x - te_i)] e_i \right\|^2 \\ &\stackrel{(77)}{\leq} \frac{d^2}{4t^2 2^j l} \sum_{i=1}^{2^j l} \|\Delta(x + te_i) - \Delta(x - te_i)\|^2 \|e_i\|^2 \\ &\stackrel{\|e_i\|=1}{=} \frac{d^2}{4t^2 2^j l} \sum_{i=1}^{2^j l} |\Delta(x + te_i) - \Delta(x - te_i)|^2 \\ &\stackrel{(9)}{\leq} \frac{d^2}{4t^2} 4\Delta^2 \\ &= \frac{d^2 \Delta^2}{t^2}. \end{aligned}$$

620 Proof of (32):

$$\begin{aligned} \mathbb{E}\|\tilde{g}_i - g_i\|^2 &\stackrel{(12),(13)}{=} \mathbb{E} \left\| d \frac{h(x + te_i, Z_i^+) - h(x - te_i, Z_i^-)}{2t} e_i \right\|^2 \\ &\stackrel{\|e_i\|=1}{=} \frac{d^2}{4t^2} \mathbb{E} [h(x + te_i, Z_i^+) - h(x - te_i, Z_i^-)]^2 \\ &\stackrel{(77),(4)}{\leq} \frac{d^2 \sigma_1^2}{t^2}. \end{aligned}$$

621 Proof of (33):

$$\begin{aligned} &\mathbb{E}\|\mathbb{E}_e [\tilde{g}^j - g^j]\|^2 \\ &\stackrel{(12),(16)}{=} \mathbb{E} \left\| \frac{1}{2^j l} \sum_{i=1}^{2^j l} \mathbb{E}_e \left[d \frac{h(x + te_i, Z_i^+) - h(x - te_i, Z_i^-)}{2t} e_i \right] \right\|^2 \\ &= \frac{d^2}{t^2} \mathbb{E} \left\| \frac{1}{2^j l} \sum_{i=1}^{2^j l} \mathbb{E}_e \left[\frac{h(x + te_i, Z_i^+) e_i - h(x - te_i, Z_i^-) e_i}{2} \right] \right\|^2 \end{aligned}$$

$$\begin{aligned}
& \stackrel{(77)}{\leq} \frac{d^2}{t^2} \frac{1}{2} \left[\mathbb{E} \left\| \frac{1}{2^j l} \sum_{i=1}^{2^j l} \mathbb{E}_e [h(x + te_i, Z_i^+) e_i] \right\|^2 + \mathbb{E} \left\| \frac{1}{2^j l} \sum_{i=1}^{2^j l} \mathbb{E}_e [h(x - te_i, Z_i^-) e_i] \right\|^2 \right] \\
& \stackrel{(21)}{\leq} \frac{dC_1 \tau \sigma_1^2}{t^2 2^j l}.
\end{aligned}$$

622 Proof of (34):

$$\begin{aligned}
& \mathbb{E} \|g_i - \nabla_{e_i} f\|^2 \\
(13), (17) \quad & \mathbb{E} \left\| d \frac{f(x + te_i) - f(x - te_i)}{2t} e_i - d \langle \nabla f(x), e_i \rangle e_i \right\|^2 \\
& = d^2 \mathbb{E} \left| \frac{f(x + te_i) - f(x) + f(x) - f(x - te_i) - 2t \langle \nabla f(x), e_i \rangle}{2t} \right|^2 \\
& = d^2 \mathbb{E} \left| \frac{f(x + te_i) - f(x) - \langle \nabla f(x), te_i \rangle}{2t} + \frac{f(x) - f(x - te_i) + \langle \nabla f(x), -te_i \rangle}{2t} \right|^2 \\
& \stackrel{\textcircled{1}}{\leq} \frac{2d^2}{4t^2} \left(\frac{L^2 t^4}{4} + \frac{L^2 t^4}{4} \right) \\
& = \frac{d^2 L^2 t^2}{4},
\end{aligned}$$

623 where ① uses Assumption 1, (74) and (77).

624 Proof of (35):

$$\begin{aligned}
\mathbb{E} \|\tilde{g}^j - \mathbb{E}_e \tilde{g}^j\|^2 & \stackrel{(15)}{=} \mathbb{E}_Z \mathbb{E}_e \left\| \frac{1}{2^j l} \sum_{i=1}^{2^j l} [\tilde{g}_i - \mathbb{E}_{e_i} \tilde{g}_i] \right\|^2 \\
& \stackrel{\textcircled{1}}{=} \mathbb{E}_Z \mathbb{E}_e \frac{1}{2^{2j} l^2} \sum_{i=1}^{2^j l} \|\tilde{g}_i - \mathbb{E}_{e_i} \tilde{g}_i\|^2 \\
& \stackrel{(78)}{\leq} \frac{1}{2^{2j} l^2} \sum_{i=1}^{2^j l} \mathbb{E}_Z \mathbb{E}_e \|\tilde{g}_i\|^2 \\
& \stackrel{(77)}{\leq} \frac{3}{2^{2j} l^2} \sum_{i=1}^{2^j l} \mathbb{E} \left[\|\tilde{g}_i - g_i\|^2 + \|g_i - \nabla_{e_i} f\|^2 + \|\nabla_{e_i} f\|^2 \right] \\
& \stackrel{(32), (34), (81)}{\leq} \frac{3}{2^j l} \left[\frac{d^2 \sigma_1^2}{t^2} + \frac{d^2 L^2 t^2}{4} + d \|\nabla f\|^2 \right],
\end{aligned}$$

625 where ① holds, since \tilde{g}_i are independent w.r.t. e_i and $\mathbb{E}_e [\tilde{g}_i - \mathbb{E}_{e_i} [\tilde{g}_i]] = 0$.

626 Proof of (36):

$$\begin{aligned}
\mathbb{E} \|\tilde{g}^j - \mathbb{E}_e g^j\|^2 & \stackrel{(77)}{\leq} 2 \mathbb{E} \left[\|\tilde{g}^j - \mathbb{E}_e \tilde{g}^j\|^2 + \|\mathbb{E}_e \tilde{g}^j - \mathbb{E}_e g^j\|^2 \right] \\
& \stackrel{(35), (33)}{\leq} 2 \left[\frac{3}{2^j l} \left[\frac{d^2 \sigma_1^2}{t^2} + \frac{d^2 L^2 t^2}{4} + d \|\nabla f\|^2 \right] + \frac{dC_1 \tau \sigma_1^2}{t^2 2^j l} \right] \\
& \lesssim \frac{d(d + \tau) \sigma_1^2}{t^2 2^j l} + \frac{d^2 L^2 t^2}{2^j l} + \frac{d \|\nabla f\|^2}{2^j l}.
\end{aligned}$$

627 Proof of (37):

$$\begin{aligned}
\mathbb{E} \|\hat{g}^j - \nabla f_t\|^2 & \stackrel{(77)}{\leq} 2 \mathbb{E} \left[\|\hat{g}^j - \tilde{g}^j\|^2 + \|\tilde{g}^j - \mathbb{E}_e g^j\|^2 \right] \\
& \stackrel{(31), (36)}{\lesssim} \frac{d^2 \Delta^2}{t^2} + \frac{d(d + \tau) \sigma_1^2}{t^2 2^j l} + \frac{d^2 L^2 t^2}{2^j l} + \frac{d \|\nabla f\|^2}{2^j l}.
\end{aligned}$$

628 Proof of (38):

$$\begin{aligned}
\|\mathbb{E}\hat{g}^j - \nabla f_t\|^2 &\stackrel{(77)}{\leq} 2\|\mathbb{E}\hat{g}^j - \mathbb{E}\tilde{g}^j\|^2 + 2\|\mathbb{E}\tilde{g}^j - \nabla f_t\|^2 \\
&\stackrel{(24)}{=} 2\|\mathbb{E}\hat{g}^j - \mathbb{E}\tilde{g}^j\|^2 + 2\|\mathbb{E}_Z\mathbb{E}_e\tilde{g}^j - \mathbb{E}_e g^j\|^2 \\
&\stackrel{(80)}{\leq} 2\|\mathbb{E}\hat{g}^j - \mathbb{E}\tilde{g}^j\|^2 + 2\mathbb{E}_Z\|\mathbb{E}_e\tilde{g}^j - \mathbb{E}_e g^j\|^2 \\
&\stackrel{(31),(33)}{\leq} \frac{2d^2\Delta^2}{t^2} + \frac{2dC_1\tau\sigma_1^2}{t^2 2^j l}.
\end{aligned}$$

629

□

630 **Lemma 6.** Assume Assumption 6, Assumption 2, Assumption 4. Then the following
631 inequalities hold for any initial distribution ξ on (Z, \mathcal{Z}) and for all $x \in \mathbb{R}^d$:

$$\mathbb{E}\|g_i\|^2 \leq dG^2, \quad (40)$$

$$\mathbb{E}\|\tilde{g} - \mathbb{E}_e \tilde{g}\|^2 \leq \frac{2}{2^j l} \left[\frac{d^2\sigma_1^2}{t^2} + dG^2 \right], \quad (41)$$

$$\mathbb{E}\|\tilde{g} - \mathbb{E}_e g\|^2 \lesssim \frac{dC_1(d+\tau)\sigma_1^2}{t^2 2^j l} + \frac{dG^2}{2^j l}. \quad (42)$$

$$\mathbb{E}\|\hat{g} - \nabla f_t\|^2 \lesssim \frac{d^2\Delta^2}{t^2} + \frac{dC_1(d+\tau)\sigma_1^2}{t^2 2^j l} + \frac{dG^2}{2^j l}. \quad (43)$$

632 *Proof.*

633 Proof of (40):

$$\begin{aligned}
\mathbb{E}\|g_i\|^2 &\stackrel{(11)}{=} \frac{d^2}{4t^2} \mathbb{E}|f(x + te_i) - f(x - te_i)|^2 \\
&\stackrel{(77)}{\leq} \frac{d^2}{2t^2} \mathbb{E}\left[|f(x + te_i) - \mathbb{E}_{e_i} f(x + te_i)|^2 + |\mathbb{E}_{e_i} f(x + te_i) - f(x - te_i)|^2\right] \\
&\stackrel{\textcircled{1}}{\leq} \frac{d^2}{t^2} \mathbb{E}|f(x + te_i) - \mathbb{E}_{e_i} f(x + te_i)|^2 \\
&\stackrel{\textcircled{2}}{\lesssim} dG^2,
\end{aligned}$$

634 where $\textcircled{1}$ uses that the distribution of e_i is symmetric, and

635 $\textcircled{2}$ uses the fact that for f which is G -Lipshitz and $e \in RS_2^d(1)$ it holds that

636 $\mathbb{E}[f(e) - \mathbb{E}_e f(e)]^2 \lesssim \frac{G^2}{d}$ [same reasoning as [49], Lemma 9].

637 Proof of (41):

$$\begin{aligned}
\mathbb{E}\|\tilde{g}^j - \mathbb{E}_e \tilde{g}^j\|^2 &\stackrel{\textcircled{1}}{\leq} \frac{1}{2^{2j} l^2} \sum_{i=1}^{2^j l} \mathbb{E}_Z \mathbb{E}_e \|\tilde{g}_i\|^2 \\
&\stackrel{(77)}{\leq} \frac{2}{2^{2j} l^2} \sum_{i=1}^{2^j l} \mathbb{E} \left[\|\tilde{g}_i - g_i\|^2 + \|g_i\|^2 \right] \\
&\stackrel{(32),(40)}{\leq} \frac{2}{2^j l} \left[\frac{d^2\sigma_1^2}{t^2} + dG^2 \right],
\end{aligned}$$

638 where $\textcircled{1}$ is analogous to (35).

639 Proof of (42):

$$\begin{aligned}
\mathbb{E}\|\tilde{g}^j - \mathbb{E}_e g^j\|^2 &\stackrel{(77)}{\leq} 2\mathbb{E} \left[\|\tilde{g}^j - \mathbb{E}_e \tilde{g}^j\|^2 + \|\mathbb{E}_e \tilde{g}^j - \mathbb{E}_e g^j\|^2 \right] \\
&\stackrel{(41),(33)}{\leq} 2 \left[\frac{2}{2^j l} \left[\frac{d^2\sigma_1^2}{t^2} + dG^2 \right] + \frac{dC_1\tau\sigma_1^2}{t^2 2^j l} \right] \\
&\lesssim \frac{d(d+\tau)\sigma_1^2}{t^2 2^j l} + \frac{dG^2}{2^j l}.
\end{aligned}$$

640 Proof of (43):

$$\begin{aligned} \mathbb{E} \|\hat{g}^j - \nabla f_t\|^2 &\stackrel{(77)}{\leq} 2\mathbb{E} \left[\|\hat{g}^j - \tilde{g}^j\|^2 + \|\mathbb{E}_e \tilde{g}^j - \nabla f_t\|^2 \right] \\ &\stackrel{(31),(42)}{\lesssim} \frac{d^2 \Delta^2}{t^2} + \frac{d(d+\tau)\sigma_1^2}{t^2 2^j l} + \frac{dG^2}{2^j l}. \end{aligned}$$

641

□

642 **Lemma 7** (Lemma 2). *Let Assumptions 3 and 4 hold. For any initial distribution ξ on*
 643 *(Z, \mathcal{Z}) the gradient estimates \hat{g}_{ml} satisfy $\mathbb{E}[\hat{g}_{ml}] = \mathbb{E}[\hat{g}_{rd}[2^{\lfloor \log_2 M \rfloor} l]]$. Moreover,*

$$\|\nabla f_t(x) - \mathbb{E}[\hat{g}_{ml}]\|^2 \lesssim \frac{d^2 \Delta^2}{t^2} + \frac{d\tau\sigma_1^2}{t^2 MB}. \quad (44)$$

644 Moreover, under assumption Assumption 1

$$\mathbb{E}[\|\nabla f_t(x) - \hat{g}_{ml}\|^2] \lesssim \frac{d^2 \Delta^2}{t^2} + \frac{d(d+\tau)\sigma_1^2}{t^2 B} + \frac{d^2 L^2 t^2}{B} + \frac{d}{B} \|\nabla f\|^2. \quad (45)$$

645 While under assumption Assumption 6

$$\mathbb{E}[\|\nabla f_t(x) - \hat{g}_{ml}\|^2] \lesssim \frac{d^2 \Delta^2}{t^2} + \frac{d(d+\tau)\sigma_1^2}{t^2 B} + \frac{dG^2}{B}. \quad (46)$$

646 *Proof.* Recall that \hat{g}_{ml} is a sum of a baseline estimate $\hat{g}_{rd}[l] \stackrel{(14)}{=} \hat{g}^0$ and a refining term
 647 $2^J[\hat{g}^J - \hat{g}^{J-1}]$. To show that $\mathbb{E}[\hat{g}_{ml}] = \mathbb{E}\hat{g}^{\lfloor \log_2 M \rfloor}$, then, we use the law of total expectation:

$$\begin{aligned} \mathbb{E}[\hat{g}_{ml}] &= \mathbb{E}[\mathbb{E}_J[\hat{g}_{ml}]] = \mathbb{E}[\hat{g}^0] + \sum_{j=1}^{\lfloor \log_2 M \rfloor} \mathbb{P}\{J=j\} \cdot 2^j \mathbb{E}[\hat{g}^j - \hat{g}^{j-1}] \\ &= \mathbb{E}[\hat{g}^0] + \sum_{j=1}^{\lfloor \log_2 M \rfloor} \mathbb{E}[\hat{g}^j - \hat{g}^{j-1}] = \mathbb{E}\hat{g}^{\lfloor \log_2 M \rfloor}. \end{aligned} \quad (47)$$

648 This immediately helps us prove the statement (44):

$$\|\nabla f_t(x) - \mathbb{E}\hat{g}_{ml}\|^2 = \left\| \nabla f_t(x) - \mathbb{E}[\hat{g}^{\lfloor \log_2 M \rfloor}] \right\|^2 \stackrel{(38)}{\leq} \frac{2d^2 \Delta^2}{t^2} + \frac{2dC_1\tau\sigma_1^2}{t^2 2^{\lfloor \log_2 M \rfloor} l} \stackrel{l \geq B}{\lesssim} \frac{d^2 \Delta^2}{t^2} + \frac{d\tau\sigma_1^2}{t^2 MB}.$$

649 Proving the statement of (45) we also start with total expectation:

$$\begin{aligned} &\mathbb{E}[\|\nabla f(x) - \hat{g}_{ml}\|^2] \\ &\stackrel{(77)}{\leq} 2\mathbb{E}[\|\nabla f(x) - \hat{g}^0\|^2] + 2\mathbb{E}[\|\hat{g}_{ml} - \hat{g}^0\|^2] \\ &= 2\mathbb{E}[\|\nabla f(x) - \hat{g}^0\|^2] + 2 \sum_{j=1}^{\lfloor \log_2 M \rfloor} \mathbb{P}\{J=j\} \cdot 4^j \mathbb{E}[\|\hat{g}^j - \hat{g}^{j-1}\|^2] \\ &= 2\mathbb{E}[\|\nabla f(x) - \hat{g}^0\|^2] + 2 \sum_{j=1}^{\lfloor \log_2 M \rfloor} 2^j \mathbb{E}[\|\hat{g}^j - \hat{g}^{j-1}\|^2] \\ &\stackrel{\textcircled{1}}{=} 2\mathbb{E}[\|\nabla f(x) - \hat{g}^0\|^2] + 2 \sum_{j=1}^{\lfloor \log_2 M \rfloor} 2^j \mathbb{E}[\|\tilde{g}^j - \tilde{g}^{j-1}\|^2] \\ &\stackrel{(77)}{\leq} 2\mathbb{E}[\|\nabla f(x) - \hat{g}^0\|^2] + 4 \sum_{j=1}^{\lfloor \log_2 M \rfloor} 2^j (\mathbb{E}\|\tilde{g}^j - \mathbb{E}_e \tilde{g}^j\|^2 + \mathbb{E}\|\mathbb{E}_e \tilde{g}^{j-1} - \tilde{g}^{j-1}\|^2) \\ &\leq 2\mathbb{E}[\|\nabla f(x) - \hat{g}^0\|^2] + 16 \sum_{j=0}^{\lfloor \log_2 M \rfloor} 2^j \mathbb{E}[\|\mathbb{E}_e \tilde{g}^j - \tilde{g}^j\|^2] \\ &\stackrel{(37),(36)}{\lesssim} 2 \left[\frac{d^2 \Delta^2}{t^2} + \frac{d(d+\tau)\sigma_1^2}{t^2 l} + \frac{d^2 L^2 t^2}{l} + \frac{d}{l} \cdot \|\nabla f\|^2 \right] \end{aligned}$$

$$\begin{aligned}
& 16 \sum_{j=0}^{\lfloor \log_2 M \rfloor} 2^j \left[\frac{d(d+\tau)\sigma_1^2}{t^2 2^{jl}} + \frac{d^2 L^2 t^2}{2^{jl}} + \frac{d \|\nabla f\|^2}{2^{jl}} \right] \\
& \stackrel{l \geq \log_2 M \cdot B}{\lesssim} 2 \left[\frac{d^2 \Delta^2}{t^2} + \frac{d(d+\tau)\sigma_1^2}{t^2 B} + \frac{d^2 L^2 t^2}{B} + \frac{d}{B} \cdot \|\nabla f\|^2 \right] + \\
& 16 \left[\frac{d(d+\tau)\sigma_1^2}{t^2 B} + \frac{d^2 L^2 t^2}{B} + \frac{d \|\nabla f\|^2}{B} \right] \\
& \lesssim \frac{d^2 \Delta^2}{t^2} + \frac{d(d+\tau)\sigma_1^2}{t^2 B} + \frac{d^2 L^2 t^2}{B} + \frac{d}{B} \|\nabla f\|^2,
\end{aligned}$$

650 where ① uses that $\hat{g}^j - \hat{g}^{j-1} = \tilde{g}^j - \tilde{g}^{j-1}$, since $\tilde{g}^j - \hat{g}^j \stackrel{(31)}{=} \tilde{g}^{j-1} - \hat{g}^{j-1}$.

651 The proof of (46) is exactly the same, replacing (37) and (36) with (43) and (42).

$$\mathbb{E}[\|\nabla f(x) - \hat{g}_{ml}\|^2] \lesssim \frac{d^2 \Delta^2}{t^2} + \frac{d(d+\tau)\sigma_1^2}{t^2 B} + \frac{dG^2}{B}.$$

652

□

653 D.4 Proof of Theorem 1

654 The proof of Theorem 1 requires two technical Lemmas.

655 **Lemma 8.** Assume Assumptions 1 and 2. Then for the iterates of Algorithm 1 with
 656 $\theta = (p\eta^{-1} - 1)/(\beta p\eta^{-1} - 1)$, $\theta > 0$, $\eta \geq 1$, $p > 0$ and arbitrary $\alpha > 0$ it holds that

$$\begin{aligned}
\mathbb{E}_k[\|x^{k+1} - x^*\|^2] & \leq (1 + \alpha p \gamma \eta)(1 - \beta) \|x^k - x^*\|^2 + (1 + \alpha p \gamma \eta) \beta \|x_g^k - x^*\|^2 \\
& + (1 + \alpha p \gamma \eta)(\beta^2 - \beta) \|x^k - x_g^k\|^2 + p^2 \eta^2 \gamma^2 \mathbb{E}_k[\|\hat{g}^k\|^2] \\
& - 2\eta^2 \gamma \langle \nabla f(x_g^k), x_g^k + (p\eta^{-1} - 1)x_f^k - \eta^{-1} p x^* \rangle \\
& + \frac{p\eta\gamma}{\alpha} \|\mathbb{E}_k[\hat{g}^k] - \nabla f(x_g^k)\|^2.
\end{aligned} \tag{48}$$

657 **Lemma 9.** Assume Assumptions 1 and 2. Let problem (4) be solved by Algorithm 1. Then
 658 for any $u \in \mathbb{R}^d$, we get

$$\begin{aligned}
\mathbb{E}_k[f(x_f^{k+1})] & \leq f(u) - \langle \nabla f(x_g^k), u - x_g^k \rangle - \frac{\mu}{2} \|u - x_g^k\|^2 - \frac{\gamma}{2} \|\nabla f(x_g^k)\|^2 \\
& + \frac{\gamma}{2} \|\mathbb{E}_k[\hat{g}^k] - \nabla f(x_g^k)\|^2 + \frac{L\gamma^2}{2} \mathbb{E}_k[\|\hat{g}^k\|^2].
\end{aligned}$$

659 These are proven in Beznosikov et al. [5] as Lemmas 5 and 6, with a slightly different notation:
 660 \hat{f} corresponds to f and \hat{g} to g .

661 **Lemma 10** (stepsize tuning). Given an optimization error after N iterations bounded by

$$r^N \leq \exp(-N\Gamma a)r^0 + \Gamma b$$

662 and an upper bound on stepsize $\Gamma \leq \frac{1}{u}$ there exists a constant stepsize $\Gamma_0 \leq \frac{1}{u}$, such that

$$r^N = \tilde{\mathcal{O}} \left(\exp \left(-\frac{Na}{u} \right) r^0 + \frac{b}{aN} \right)$$

663 Equivalently, the number of iterations to get $r^N \lesssim \varepsilon$:

$$N = \tilde{\mathcal{O}} \left(\frac{u}{a} \ln \varepsilon^{-1} + \frac{b}{a\varepsilon} \right) \tag{49}$$

664 *Proof.* This setup is a simpler version of the one considered in Section 4 of Stich [53] and so
 665 we will tune Γ similarly to their Lemma 2:

$$\Gamma := \min \left(\frac{\ln \max(2, ar^0 N/b)}{aN}, \frac{1}{u} \right)$$

666 If $\frac{1}{u} < \frac{\ln \max(2, ar^0 N/b)}{aN}$, then $\Gamma := \frac{1}{u}$.

$$r^N \leq \exp\left(\frac{-Na}{u}\right) r^0 + \frac{b}{u} \leq \exp\left(\frac{-Na}{u}\right) r^0 + \frac{b \ln(\dots)}{aN} = \tilde{O}\left(\exp\left(-\frac{Na}{u}\right) r^0 + \frac{b}{aN}\right)$$

667 Otherwise $\frac{\ln \max(2, ar^0 N/b)}{aN} \leq \frac{1}{u}$ and $\Gamma := \frac{\ln \max(2, ar^0 N/b)}{aN}$, with $\Gamma b = \tilde{O}(\frac{b}{aN})$ immediately.

$$\exp(-N\Gamma a) r^0 = \exp(-\ln \max(2, ar^0 N/b)) = \frac{1}{\max(2, ar^0 N/b)}$$

668 If $ar^0 N/b > 2$, we also get $\tilde{O}(\frac{b}{aN})$, else $\frac{1}{2} \leq \frac{b}{aNr^0}$ and we get $\tilde{O}(\frac{b}{aN})$ as well.

669 To conclude the proof we should mitigate the fact that the stepsize currently depends on the
670 number of iterations. This can easily be done via a restart procedure which would run the
671 algorithm for $N = 1, 2, 4, \dots$ iterations with a stepsize $\Gamma(N)$. \square

672 **Theorem 4** (Theorem 1). *Let Assumptions 1 to 4 hold, and consider problem (4) solved by*
673 *Algorithm 1. Then, for a suitable choice of hidden parameters (with $p \simeq \frac{B}{B+d}$) and arbitrary*
674 *choice of free parameters (see Table 2), it holds that:*

$$\mathbb{E}r^N \lesssim \exp\left(-\sqrt{\frac{p^2 \mu \gamma N^2}{3}}\right) r^0 + \frac{p\sqrt{\gamma}}{\mu^{3/2}} \cdot \left[\sigma_1^2 \frac{d(d+\tau)}{t^2 B} + t^2 \frac{L^2 d^2}{B}\right] + \frac{\Delta^2 d^2}{\mu^2 t^2} + \frac{Lt^2}{\mu}$$

675 Moreover, for arbitrary $\varepsilon \gtrsim \frac{d\Delta\sqrt{L}}{\mu^{3/2}}$ and an appropriate choice of t and γ , the number of oracle
676 calls required to ensure $r^N \lesssim \varepsilon$ is bounded by

$$B \cdot \tilde{O}\left(\max\left[1, \frac{d}{B}\right] \sqrt{\frac{L}{\mu}} \log \frac{1}{\varepsilon} + \frac{Ld(d+\tau)\sigma_1^2}{B\mu^3 \varepsilon^2}\right) \quad \text{one-point oracle calls.}$$

677 *Proof.* Applying Lemma 9 with $u = x^*$ (for arbitrary x^*) and $u = x_f^k$ to f_t , we get:

$$\begin{aligned} \mathbb{E}_k[f_t(x_f^{k+1})] &\leq f_t(x^*) - \langle \nabla f_t(x_g^k), x^* - x_g^k \rangle - \frac{\mu}{2} \|x^* - x_g^k\|^2 - \frac{p\gamma}{2} \|\nabla f_t(x_g^k)\|^2 \\ &\quad + \frac{p\gamma}{2} \|\mathbb{E}_k[\hat{g}^k] - \nabla f_t(x_g^k)\|^2 + \frac{Lp^2\gamma^2}{2} \mathbb{E}_k[\|\hat{g}^k\|^2], \end{aligned} \quad (50)$$

678

$$\begin{aligned} \mathbb{E}_k[f_t(x_f^{k+1})] &\leq f_t(x_f^k) - \langle \nabla f_t(x_g^k), x_f^k - x_g^k \rangle - \frac{\mu}{2} \|x_f^k - x_g^k\|^2 - \frac{p\gamma}{2} \|\nabla f_t(x_g^k)\|^2 \\ &\quad + \frac{p\gamma}{2} \|\mathbb{E}_k[\hat{g}^k] - \nabla f_t(x_g^k)\|^2 + \frac{Lp^2\gamma^2}{2} \mathbb{E}_k[\|\hat{g}^k\|^2]. \end{aligned} \quad (51)$$

679 Combining $2p\gamma\eta \cdot (50) + 2\gamma\eta(\eta - p) \cdot (51) + (48)$ we get:

$$\begin{aligned} &\mathbb{E}_k[\|x^{k+1} - x^*\|^2 + 2\gamma\eta^2 f_t(x_f^{k+1})] \\ &\leq (1 + \alpha p\gamma\eta)(1 - \beta) \|x^k - x^*\|^2 + (1 + \alpha p\gamma\eta)\beta \|x_g^k - x^*\|^2 \\ &\quad + (1 + \alpha p\gamma\eta)(\beta^2 - \beta) \|x^k - x_g^k\|^2 - 2\eta^2 \gamma \langle \nabla f_t(x_g^k), x_g^k \rangle + (p\eta^{-1} - 1)x_f^k - \eta^{-1} p x^* \\ &\quad + p^2 \eta^2 \gamma^2 \mathbb{E}_k[\|\hat{g}^k\|^2] + \frac{p\eta\gamma}{\alpha} \|\mathbb{E}_k[\hat{g}^k] - \nabla f_t(x_g^k)\|^2 \\ &\quad + 2p\gamma\eta \left(f_t(x^*) - \langle \nabla f_t(x_g^k), x^* - x_g^k \rangle - \frac{\mu}{2} \|x^* - x_g^k\|^2 - \frac{p\gamma}{2} \|\nabla f_t(x_g^k)\|^2 \right. \\ &\quad \left. + \frac{p\gamma}{2} \|\mathbb{E}_k[\hat{g}^k] - \nabla f_t(x_g^k)\|^2 + \frac{Lp^2\gamma^2}{2} \mathbb{E}_k[\|\hat{g}^k\|^2] \right) \\ &\quad + 2\gamma\eta(\eta - p) \left(f_t(x_f^k) - \langle \nabla f_t(x_g^k), x_f^k - x_g^k \rangle - \frac{\mu}{2} \|x_f^k - x_g^k\|^2 - \frac{p\gamma}{2} \|\nabla f_t(x_g^k)\|^2 \right. \\ &\quad \left. + \frac{p\gamma}{2} \|\mathbb{E}_k[\hat{g}^k] - \nabla f_t(x_g^k)\|^2 + \frac{Lp^2\gamma^2}{2} \mathbb{E}_k[\|\hat{g}^k\|^2] \right) \end{aligned}$$

$$\begin{aligned}
&= (1 + \alpha p \gamma \eta)(1 - \beta) \|x^k - x^*\|^2 + 2\gamma \eta (\eta - p) f_t(x_f^k) + 2p\gamma \eta f_t(x^*) \\
&\quad + ((1 + \alpha p \gamma \eta)\beta - p\gamma \eta \mu) \|x_g^k - x^*\|^2 \\
&\quad + (1 + \alpha p \gamma \eta)(\beta^2 - \beta) \|x^k - x_g^k\|^2 - p\gamma^2 \eta^2 \|\nabla f_t(x_g^k)\|^2 \\
&\quad + \left(\frac{p\eta\gamma}{\alpha} + p\gamma^2 \eta^2 \right) \|\mathbb{E}_k[\hat{g}^k] - \nabla f_t(x_g^k)\|^2 + (p^2 \eta^2 \gamma^2 + p^2 \gamma^3 \eta^2 L) \mathbb{E}_k[\|\hat{g}^k\|^2] \\
&\stackrel{(76)}{\leq} (1 + \alpha p \gamma \eta)(1 - \beta) \|x^k - x^*\|^2 + 2\gamma \eta (\eta - p) f_t(x_f^k) + 2p\gamma \eta f_t(x^*) \\
&\quad + ((1 + \alpha p \gamma \eta)\beta - p\gamma \eta \mu) \|x_g^k - x^*\|^2 \\
&\quad + (1 + \alpha p \gamma \eta)(\beta^2 - \beta) \|x^k - x_g^k\|^2 - p\gamma^2 \eta^2 \|\nabla f_t(x_g^k)\|^2 \\
&\quad + p\eta\gamma \left(\frac{1}{\alpha} + \gamma \eta \right) \|\mathbb{E}_k[\hat{g}^k] - \nabla f_t(x_g^k)\|^2 + 2p^2 \eta^2 \gamma^2 (1 + \gamma L) \mathbb{E}_k[\|\hat{g}^k - \nabla f_t(x_g^k)\|^2] \\
&\quad + 2p^2 \eta^2 \gamma^2 (1 + \gamma L) \mathbb{E}_k \left[\underbrace{\|\nabla f_t(x_g^k)\|^2}_{x_g^k \in \mathcal{F}_k} \right].
\end{aligned}$$

680 Choosing $\alpha = \frac{\beta}{2p\eta\gamma}$ gives:

$$\begin{aligned}
\beta &= \sqrt{4p^2 \mu \gamma / 3} \stackrel{\gamma \leq \frac{3}{4L}}{\leq} \sqrt{p^2 \mu / L} < 1, \\
(1 + \alpha p \eta \gamma)(1 - \beta) &= \left(1 + \frac{\beta}{2}\right)(1 - \beta) \leq \left(1 - \frac{\beta}{2}\right), \\
((1 + \alpha p \eta \gamma)\beta - p\mu \gamma \eta) &= \left(\beta + \frac{\beta^2}{2} - p\mu \gamma \eta\right) \stackrel{\beta < 1}{\leq} \left(\frac{3\beta}{2} - p\mu \gamma \eta\right) \stackrel{p\mu \gamma \eta = 3\beta/2}{\leq} 0.
\end{aligned}$$

681 Thus:

$$\begin{aligned}
&\mathbb{E}_k[\|x^{k+1} - x^*\|^2 + 2\gamma \eta^2 f_t(x_f^{k+1})] \\
&\leq (1 - \beta/2) \|x^k - x^*\|^2 + 2\gamma \eta (\eta - p) f_t(x_f^k) + 2p\gamma \eta f_t(x^*) \\
&\quad + p\eta^2 \gamma^2 (1 + 2p/\beta) \|\mathbb{E}_k[\hat{g}^k] - \nabla f_t(x_g^k)\|^2 \\
&\quad + 2p^2 \eta^2 \gamma^2 (1 + \gamma L) \mathbb{E}_k[\|\hat{g}^k - \nabla f_t(x_g^k)\|^2] \\
&\quad - p\gamma^2 \eta^2 (1 - 2p(1 + \gamma L)) \|\nabla f_t(x_g^k)\|^2.
\end{aligned}$$

682 Subtracting $2\gamma \eta^2 f_t(x^*)$ from both sides, we get:

$$\begin{aligned}
&\mathbb{E}_k[\|x^{k+1} - x^*\|^2 + 2\gamma \eta^2 (f_t(x_f^{k+1}) - f_t(x^*))] \\
&\leq (1 - \beta/2) \|x^k - x^*\|^2 + (1 - p/\eta) \cdot 2\gamma \eta^2 (f_t(x_f^k) - f_t(x^*)) \\
&\quad + p\eta^2 \gamma^2 (1 + 2p/\beta) \|\mathbb{E}_k[\hat{g}^k] - \nabla f_t(x_g^k)\|^2 \\
&\quad + 2p^2 \eta^2 \gamma^2 (1 + \gamma L) \mathbb{E}_k[\|\hat{g}^k - \nabla f_t(x_g^k)\|^2] \\
&\quad - p\gamma^2 \eta^2 (1 - 2p(1 + \gamma L)) \|\nabla f_t(x_g^k)\|^2 \\
&\stackrel{\beta/2 = p/\eta}{=} (1 - \beta/2) [\|x^k - x^*\|^2 + 2\gamma \eta^2 (f_t(x_f^k) - f_t(x^*))] \\
&\quad + p\eta^2 \gamma^2 (1 + 2p/\beta) \|\mathbb{E}_k[\hat{g}^k] - \nabla f_t(x_g^k)\|^2 \\
&\quad + 2p^2 \eta^2 \gamma^2 (1 + \gamma L) \mathbb{E}_k[\|\hat{g}^k - \nabla f_t(x_g^k)\|^2] \\
&\quad - p\gamma^2 \eta^2 (1 - 2p(1 + \gamma L)) \|\nabla f_t(x_g^k)\|^2.
\end{aligned}$$

683 Applying Lemma 7, one can obtain:

$$\begin{aligned}
& \mathbb{E}_k [\|x^{k+1} - x^*\|^2 + 2\gamma\eta^2(f_t(x_f^{k+1}) - f_t(x^*))] \\
& \lesssim (1 - \beta/2) \left[\|x^k - x^*\|^2 + 2\gamma\eta^2(f_t(x_f^k) - f_t(x^*)) \right] \\
& \quad + p\eta^2\gamma^2(1 + 2p/\beta) \cdot \left[\frac{d^2\Delta^2}{t^2} + \frac{d\tau\sigma_1^2}{t^2MB} \right] \\
& \quad + 2p^2\eta^2\gamma^2(1 + \gamma L) \cdot \left[\frac{d^2\Delta^2}{t^2} + \frac{d(d+\tau)\sigma_1^2}{t^2B} + \frac{d^2L^2t^2}{B} + \frac{d}{B} \|\nabla f(x_g^k)\|^2 \right] \\
& \quad - p\gamma^2\eta^2(1 - 2p(1 + \gamma L)) \|\nabla f_t(x_g^k)\|^2 \\
& = \frac{1}{M} = p(1 + 2p/\beta)^{-1} \left[\|x^k - x^*\|^2 + 2\gamma\eta^2(f_t(x_f^k) - f_t(x^*)) \right] \\
& \quad + p^2\eta^2\gamma^2 \cdot \left[\frac{d^2\Delta^2M}{t^2} + \frac{d\tau\sigma_1^2}{t^2B} \right] \\
& \quad + 2p^2\eta^2\gamma^2(1 + \gamma L) \cdot \left[\frac{d^2\Delta^2}{t^2} + \frac{d(d+\tau)\sigma_1^2}{t^2B} + \frac{d^2L^2t^2}{B} + \frac{d}{B} \|\nabla f(x_g^k)\|^2 \right] \\
& \quad - p\gamma^2\eta^2(1 - 2p(1 + \gamma L)) \|\nabla f_t(x_g^k)\|^2 \\
& \stackrel{(30)}{\lesssim} (1 - \beta/2) \left[\|x^k - x^*\|^2 + 2\gamma\eta^2(f_t(x_f^k) - f_t(x^*)) \right] \\
& \quad + \Delta^2 \cdot \left[\frac{p^2\eta^2\gamma^2d^2M + p^2\eta^2\gamma^2(1 + \gamma L)d^2}{t^2} \right] \\
& \quad + \|\nabla f_t(x_g^k)\|^2 \cdot \left[p^2\eta^2\gamma^2(1 + \gamma L)\frac{d}{B} - p\gamma^2\eta^2(1 - 2p(1 + \gamma L)) \right] \\
& \quad + \sigma_1^2 \cdot \left[\frac{p^2\eta^2\gamma^2d\tau + p^2\eta^2\gamma^2(1 + \gamma L)d(d+\tau)}{t^2B} \right] \\
& \quad + \frac{t^2}{B} \cdot p^2\eta^2\gamma^2(1 + \gamma L)L^2(d^2 + d) \\
& \stackrel{\gamma L < 1}{\lesssim} (1 - \beta/2) \left[\|x^k - x^*\|^2 + 2\gamma\eta^2(f_t(x_f^k) - f_t(x^*)) \right] \\
& \quad + p^2\eta^2\gamma^2 \cdot \left[\sigma_1^2 \frac{d(d+\tau)}{t^2B} + t^2 \frac{L^2d^2}{B} + \Delta^2 \frac{d^2M}{t^2} \right] \\
& \quad + \|\nabla f(x_g^k)\|^2 \cdot p\gamma^2\eta^2 \underbrace{\left[-1 + p(1 + \gamma L) \left(1 + \frac{d}{B} \right) \right]}_{=0 \text{ for } p \simeq \frac{B}{B+d}} \\
& \stackrel{p\eta\gamma=3\beta/(2\mu)}{\lesssim} (1 - \beta/2) \left[\|x^k - x^*\|^2 + 2\gamma\eta^2(f_t(x_f^k) - f_t(x^*)) \right] \\
& \quad + \frac{\beta^2}{\mu^2} \cdot \left[\sigma_1^2 \frac{d(d+\tau)}{t^2B} + t^2 \frac{L^2d^2}{B} + \Delta^2 \frac{d^2M}{t^2} \right].
\end{aligned}$$

684 Finally, we perform the recursion and substitute $\beta = \sqrt{4p^2\mu\gamma/3}$, $\eta = \sqrt{\frac{3}{\mu\gamma}}$,

685 $r_t^N = \|x^N - x^*\|^2 + \frac{1}{\mu}(f_t(x_f^N) - f_t(x^*))$:

$$\begin{aligned}
\mathbb{E}r_t^N & \lesssim \left(1 - \sqrt{\frac{p^2\mu\gamma}{3}} \right)^N r_t^0 \\
& \quad + \frac{\beta}{\mu^2} \cdot \left[\sigma_1^2 \frac{d(d+\tau)}{t^2B} + t^2 \frac{L^2d^2}{B} + \Delta^2 \frac{d^2M}{t^2} \right]
\end{aligned}$$

$$\begin{aligned}
&\lesssim \exp\left(-\sqrt{\frac{p^2\mu\gamma N^2}{3}}\right) r_t^0 \\
&\quad + \frac{p\sqrt{\gamma}}{\mu^{3/2}} \cdot \left[\sigma_1^2 \frac{d(d+\tau)}{t^2 B} + t^2 \frac{L^2 d^2}{B} + \Delta^2 \frac{d^2 M}{t^2}\right] \\
&\stackrel{\textcircled{1}}{\lesssim} \exp\left(-\sqrt{\frac{p^2\mu\gamma N^2}{3}}\right) r_t^0 \\
&\quad + \frac{p\sqrt{\gamma}}{\mu^{3/2}} \cdot \left[\sigma_1^2 \frac{d(d+\tau)}{t^2 B} + t^2 \frac{L^2 d^2}{B}\right] \\
&\quad + \frac{\Delta^2 d^2}{\mu^2 t^2},
\end{aligned}$$

686 where ① uses that $M \simeq \frac{1}{p} \left(1 + \frac{1}{\sqrt{\mu\gamma}}\right) \Rightarrow Mp\sqrt{\gamma} \simeq \sqrt{\gamma} + \frac{1}{\sqrt{\mu}} \leq \frac{1}{\sqrt{L}} + \frac{1}{\sqrt{\mu}} \lesssim \frac{1}{\sqrt{\mu}}$

687 Recall that x^* is arbitrary. Therefore by setting $x^* = \arg \min f(x)$, we may bound the error
688 for non-smoothed f :

$$\begin{aligned}
r^N &= \|x^N - x^*\|^2 + \frac{6}{\mu} (f(x_f^N) - f(x^*)) \\
&= \|x^N - x^*\|^2 + \frac{6}{\mu} \underbrace{(f(x_f^N) - f_t(x_f^N))}_{\leq 0 \text{ (28)}} \underbrace{(-f(x^*) + f_t(x^*))}_{\leq Lt^2 \text{ (28)}} + \frac{6}{\mu} (f_t(x_f^N) - f_t(x^*)) \\
&\leq r_t^N + 6 \frac{Lt^2}{\mu}
\end{aligned}$$

689 Thus we get

$$\begin{aligned}
\mathbb{E}r^N &\lesssim \exp\left(-\sqrt{\frac{p^2\mu\gamma N^2}{3}}\right) r^0 \\
&\quad + \frac{p\sqrt{\gamma}}{\mu^{3/2}} \cdot \left[\sigma_1^2 \frac{d(d+\tau)}{t^2 B} + t^2 \frac{L^2 d^2}{B}\right] \\
&\quad + \frac{\Delta^2 d^2}{\mu^2 t^2} + \frac{Lt^2}{\mu}
\end{aligned}$$

690 To finish the analysis we need to define t and γ , as well as the tolerable level of noise Δ .
691 Currently we are left with an expression of form:

$$\mathbb{E}r^N \lesssim \exp(-N\Gamma a)r^0 + \Gamma b + c, \Gamma \leq \frac{1}{u}$$

692 with

$$\begin{aligned}
\Gamma &= \sqrt{\gamma} \\
u &\simeq \sqrt{L} \\
a &\simeq p\sqrt{\mu} \\
b &\simeq \frac{p}{\mu^{3/2}} \cdot \left[\sigma_1^2 \frac{d(d+\tau)}{t^2 B} + t^2 \frac{L^2 d^2}{B}\right] \\
c &= \frac{\Delta^2 d^2}{\mu^2 t^2} + \frac{Lt^2}{\mu}
\end{aligned}$$

693 To get $c \lesssim \varepsilon$ we have to bound t :

$$\frac{d\Delta}{\mu\sqrt{\varepsilon}} \lesssim t \lesssim \frac{\sqrt{\mu\varepsilon}}{\sqrt{L}}$$

694 Thus we bound the adversarial noise $\varepsilon \gtrsim \frac{d\Delta\sqrt{L}}{\mu^{3/2}} \Leftrightarrow \Delta \lesssim \frac{\varepsilon\mu^{3/2}}{d\sqrt{L}}$.

695 Applying Lemma 10, to get $r^N \lesssim \varepsilon$ one would need N iterations:

$$N = \tilde{O} \left(\frac{1}{p} \sqrt{\frac{L}{\mu}} \log \frac{1}{\varepsilon} + \frac{d}{B\mu^2\varepsilon} \left[\frac{(d+\tau)\sigma_1^2}{t^2} + L^2 t^2 d \right] \right) \quad (52)$$

696 Recalling $p \simeq \frac{B}{B+d}$, as well as setting t to its upper bound, we get the total number of
697 iterations:

$$N = \tilde{O} \left(\left[1 + \frac{d}{B} \right] \sqrt{\frac{L}{\mu}} \log \frac{1}{\varepsilon} + \frac{Ld(d+\tau)\sigma_1^2}{\mu^3\varepsilon^2 B} \right)$$

698 Finally, as noted in Section 2.1, each \hat{g}_{ml} uses $\tilde{O}(B)$ oracle calls, thus the oracle complexity
699 is:

$$B \cdot \tilde{O} \left(\max \left[1, \frac{d}{B} \right] \sqrt{\frac{L}{\mu}} \log \frac{1}{\varepsilon} + \frac{Ld(d+\tau)\sigma_1^2}{B\mu^3\varepsilon^2} \right) \quad \text{one-point oracle calls.}$$

700

□

701 D.5 Proof of Theorem 3

702 **Theorem 5** (Theorem 3). *Let Assumptions 2 to 4 and 6 hold, and consider problem (4)
703 solved by Algorithm 1. Then, for a suitable choice of hidden parameters (with $p \simeq 1$) and
704 arbitrary choice of free parameters (see Table 2), it holds that:*

$$\mathbb{E}r^N \lesssim \exp \left(-\sqrt{\frac{\mu\gamma N^2}{3}} \right) r^0 + \frac{\sqrt{\gamma}}{\mu^{3/2}} \cdot \left[\sigma_1^2 \frac{dC_1(d+\tau)}{t^2 B} + \frac{G^2 d}{B} \right] + \frac{\Delta^2 d^2}{\mu^2 t^2} + \frac{Gt}{\mu}$$

705 Moreover, for arbitrary $\varepsilon \gtrsim \left[\frac{d\Delta G}{\mu^2} \right]^{2/3}$ and an appropriate choice of t and γ , the number of
706 oracle calls required to ensure $r^N \lesssim \varepsilon$ is bounded by

$$B \cdot \tilde{O} \left[\sqrt{\frac{\sqrt{d}G^2}{\mu^2\varepsilon}} \log \frac{1}{\varepsilon} + \frac{d(d+\tau)\sigma_1^2 G^2}{\mu^4 \varepsilon^3 B} + \frac{G^2 d}{B\mu^2\varepsilon} \right] \quad \text{one-point oracle calls.}$$

707 *Proof.* The proof is almost identical to the smooth case. The difference is we use (46) instead
708 of (45). With that $p \simeq 1$ is enough, as the term with $\|\nabla f(x_g^k)\|$ no longer exists. Additionally,
709 $\frac{d^2 L^2 t^2}{B} \rightarrow \frac{dG^2}{B}$. Finally, we may use Lemma 9 as smoothed function is indeed smooth (27).

$$\begin{aligned} \mathbb{E}r^N &\lesssim \exp \left(-\sqrt{\frac{p^2 \mu \gamma N^2}{3}} \right) r_t^0 \\ &\quad + \frac{p\sqrt{\gamma}}{\mu^{3/2}} \cdot \left[\sigma_1^2 \frac{dC_1(d+\tau)}{t^2 B} + \frac{G^2 d}{B} \right] \\ &\quad + \frac{\Delta^2 d^2}{\mu^2 t^2} + \underbrace{\frac{Gt}{\mu}}_{(26)} \\ &\stackrel{p \simeq 1}{\simeq} \exp \left(-\sqrt{\frac{\mu \gamma N^2}{3}} \right) r_t^0 \\ &\quad + \frac{\sqrt{\gamma}}{\mu^{3/2}} \cdot \left[\sigma_1^2 \frac{dC_1(d+\tau)}{t^2 B} + \frac{G^2 d}{B} \right] \\ &\quad + \frac{\Delta^2 d^2}{\mu^2 t^2} + \frac{Gt}{\mu} \end{aligned}$$

710 Applying Lemma 10 with:

$$\begin{aligned}\Gamma &= \sqrt{\gamma} \\ u &\simeq \sqrt{L} \stackrel{(27)}{\simeq} \sqrt{\frac{\sqrt{d}G}{t}} \\ a &= \sqrt{\mu} \\ b &= \frac{1}{\mu^{3/2}} \cdot \left[\sigma_1^2 \frac{dC_1(d+\tau)}{t^2 B} + \frac{G^2 d}{B} \right]\end{aligned}$$

711 We get that $r^N \lesssim \varepsilon$ takes N iterations:

$$N = \tilde{O} \left(\sqrt{\frac{\sqrt{d}G}{t\mu}} \log \frac{1}{\varepsilon} + \frac{d}{B\mu^2\varepsilon} \left[\frac{(d+\tau)\sigma_1^2}{t^2} + G^2 \right] \right).$$

712 To get $c \lesssim \varepsilon$ we have to bound t :

$$\frac{d\Delta}{\mu\sqrt{\varepsilon}} \lesssim t \lesssim \frac{\mu\varepsilon}{G}$$

713 Thus we bound the adversarial noise $\varepsilon \gtrsim \left[\frac{d\Delta G}{\mu^2} \right]^{2/3} \Leftrightarrow \Delta \lesssim \frac{\varepsilon^{3/2}\mu^2}{dG}$.

714 Substituting $L = \frac{\sqrt{d}G}{t}$, as well as setting t to its upper bound, we get the total number of
715 iterations:

$$N = \tilde{O} \left(\sqrt{\frac{\sqrt{d}G^2}{\mu^2\varepsilon}} \log \frac{1}{\varepsilon} + \frac{d(d+\tau)\sigma_1^2 G^2}{\mu^4 \varepsilon^3 B} + \frac{G^2 d}{B\mu^2\varepsilon} \right).$$

716 And the oracle complexity:

$$B \cdot \tilde{O} \left(\sqrt{\frac{\sqrt{d}G^2}{\mu^2\varepsilon}} \log \frac{1}{\varepsilon} + \frac{d(d+\tau)\sigma_1^2 G^2}{\mu^4 \varepsilon^3 B} + \frac{G^2 d}{B\mu^2\varepsilon} \right) \quad \text{one-point oracle calls.}$$

717

□

718 E Proofs of two-point results

719 The proofs for one- and two- point feedback will functionally differ only in Lemma 5 and
720 Lemma 7, while the rest of the machinery will be reused.

721 E.1 Inequalities for gradient approximation

722 **Lemma 5'.** Assume Assumption 1, Assumption 3 and Assumption 4. Then the following
723 inequalities hold for any initial distribution ξ on (Z, \mathcal{Z}) and for all $x \in \mathbb{R}^d$:

$$\mathbb{E}[\|\tilde{g}_i - \nabla_{e_i} F_i\|^2] \leq \frac{L^2 d^2 t^2}{4}, \quad (53)$$

$$\mathbb{E}\|\nabla_{e_i} F_i - \nabla_{e_i} f\|^2 \leq d\sigma_2^2, \quad (54)$$

$$\mathbb{E}\|\tilde{g}^j - \mathbb{E}_e \tilde{g}^j\|^2 \leq \frac{3}{2^j l} \left[d^2 t^2 L^2 / 4 + d\sigma_2^2 + d\|\nabla f\|^2 \right], \quad (55)$$

$$\mathbb{E}_Z \|\mathbb{E}_e \tilde{g}^j - \nabla f_t\|^2 \lesssim \frac{\tau}{2^j l} \sigma_2^2, \quad (56)$$

$$\|\mathbb{E} \hat{g}^j - \nabla f_t\|^2 \lesssim \frac{d^2 \Delta^2}{t^2} + \frac{\tau}{2^j l} \sigma_2^2, \quad (57)$$

$$\mathbb{E}\|\tilde{g}^j - \nabla f_t\|^2 \lesssim \frac{d^2 t^2 L^2}{2^j l} + \frac{d+\tau}{2^j l} \sigma_2^2 + \frac{d}{2^j l} \|\nabla f\|^2. \quad (58)$$

$$\mathbb{E}\|\hat{g}^j - \nabla f_t\|^2 \lesssim \frac{d^2 \Delta^2}{t^2} + \frac{d^2 t^2 L^2}{2^j l} + \frac{d+\tau}{2^j l} \sigma_2^2 + \frac{d}{2^j l} \|\nabla f\|^2. \quad (59)$$

724 *Proof.*

725 We prove all estimates one by one, starting with (53):

$$\begin{aligned}
& \mathbb{E}[\|\tilde{g}_i - \nabla_{e_i} F_i\|^2] \\
& \stackrel{(12),(18)}{=} d^2 \mathbb{E} \left[\left\| \frac{F(x + te_i, Z_i) - F(x - te_i, Z_i)}{2t} e_i - \langle \nabla F(x, Z_i), e_i \rangle e_i \right\|^2 \right] \\
& \stackrel{(1')+(74)}{\leq} \frac{L^2 d^2 t^2}{4}.
\end{aligned}$$

726 Proof of (54):

$$\mathbb{E} \|\nabla_{e_i} F_i - \nabla_{e_i} f\|^2 = \mathbb{E}_Z \mathbb{E}_{e_i} \|\nabla_{e_i} F_i - \nabla_{e_i} f\|^2 \stackrel{(17),(81)}{=} d \mathbb{E}_Z \|\nabla F_i - \nabla f\|^2 \stackrel{(4')}{\leq} d \sigma_2^2.$$

727 Proof of (55):

$$\begin{aligned}
\mathbb{E} \|\tilde{g}^j - \mathbb{E}_e \tilde{g}^j\|^2 & \stackrel{\text{like (35)}}{\leq} \frac{1}{2^{2j} l^2} \sum_{i=1}^{2^j l} \mathbb{E} \|\tilde{g}_i\|^2 \\
& \stackrel{(77)}{\leq} \frac{3}{2^{2j} l^2} \sum_{i=1}^{2^j l} \left[\mathbb{E} [\|\tilde{g}_i - \nabla_{e_i} F_i\|^2] + \right. \\
& \quad \left. \mathbb{E} [\|\nabla_{e_i} F_i - \nabla_{e_i} f\|^2] + \mathbb{E} [\|\nabla_{e_i} f\|^2] \right] \\
& \stackrel{(53),(54),(81)}{=} \frac{3}{2^j l} \left[d^2 t^2 L^2 / 4 + d \sigma_2^2 + d \|\nabla f\|^2 \right].
\end{aligned}$$

728 Proof of (56):

$$\begin{aligned}
\mathbb{E}_Z \|\mathbb{E}_e \tilde{g}^j - \nabla f_t\|^2 & \stackrel{(15)}{=} \mathbb{E}_Z \left\| \mathbb{E}_e \left[\frac{1}{2^j l} \sum_{i=1}^{2^j l} \tilde{g}_i \right] - \nabla f_t \right\|^2 \\
& \stackrel{(24)}{\leq} \mathbb{E}_Z \left\| \frac{1}{2^j l} \sum_{i=1}^{2^j l} \nabla F_t(x, Z_i) - \nabla f_t(x) \right\|^2 \\
& \stackrel{\textcircled{1}}{\lesssim} \frac{\tau}{2^j l} \sigma_2^2,
\end{aligned}$$

729 where ① uses (22) for $\nabla F_t, \nabla f_t$. Let us verify that Assumption 4' holds:

730 Unbiasedness:

$$\begin{aligned}
\mathbb{E}_\pi \nabla F_t(x, Z) &= \mathbb{E}_\pi \nabla [\mathbb{E}_r F(x + tr, Z)] = \\
&= \mathbb{E}_\pi \mathbb{E}_r \nabla F(x + tr, Z) = \mathbb{E}_r \mathbb{E}_\pi \nabla F(x + tr, Z) = \mathbb{E}_r \nabla f(x + tr) = \nabla f_t(x).
\end{aligned}$$

731 Variance:

$$\begin{aligned}
\|\nabla F_t(x, Z) - \nabla f_t(x)\|^2 &= \|\mathbb{E}_r \nabla F(x + tr, Z) - \nabla f(x + tr)\|^2 \stackrel{(80)}{\leq} \\
&= \mathbb{E}_r \|\nabla F(x + tr, Z) - \nabla f(x + tr)\|^2 \stackrel{(4')}{\leq} \mathbb{E}_r \sigma_2^2 = \sigma_2^2.
\end{aligned}$$

732 Proof of (57):

$$\begin{aligned}
\|\mathbb{E} \hat{g}^j - \nabla f_t\|^2 & \stackrel{(77)}{\leq} 2 \|\mathbb{E} \hat{g}^j - \mathbb{E} \tilde{g}^j\|^2 + 2 \|\mathbb{E} \tilde{g}^j - \nabla f_t\|^2 \\
& \stackrel{(80)}{\leq} 2 \|\mathbb{E} \hat{g}^j - \mathbb{E} \tilde{g}^j\|^2 + 2 \mathbb{E}_Z \|\mathbb{E}_e \tilde{g}^j - \nabla f_t\|^2
\end{aligned}$$

$$\stackrel{(31),(56)}{\lesssim} \frac{d^2 \Delta^2}{t^2} + \frac{\tau}{2^j l} \sigma_2^2.$$

733 Proof of (58):

$$\begin{aligned} \mathbb{E} \|\tilde{g}^j - \nabla f_t\|^2 &\stackrel{(77)}{\leq} 2\mathbb{E} \|\tilde{g}^j - \mathbb{E}_e \tilde{g}^j\|^2 + 2\mathbb{E} \|\mathbb{E}_e \tilde{g}^j - \nabla f_t\|^2 \\ &\stackrel{(55),(56)}{\lesssim} \frac{1}{2^j l} \left[d^2 t^2 L^2 + d\sigma_2^2 + d\|\nabla f\|^2 \right] + \frac{\tau}{2^j l} \sigma_2^2 \\ &\lesssim \frac{d^2 t^2 L^2}{2^j l} + \frac{d+\tau}{2^j l} \sigma_2^2 + \frac{d}{2^j l} \|\nabla f\|^2. \end{aligned}$$

734 Proof of (59):

$$\begin{aligned} \mathbb{E} \|\hat{g}^j - \nabla f_t\|^2 &\stackrel{(77)}{\leq} 2\mathbb{E} \|\hat{g}^j - \tilde{g}^j\|^2 + 2\mathbb{E} \|\tilde{g}^j - \nabla f_t\|^2 \\ &\stackrel{(31),(58)}{\lesssim} \frac{d^2 \Delta^2}{t^2} + \frac{d^2 t^2 L^2}{2^j l} + \frac{d+\tau}{2^j l} \sigma_2^2 + \frac{d}{2^j l} \|\nabla f\|^2. \end{aligned}$$

735

□

736 **Lemma 11.** Assume Assumption 6', Assumption 2, Assumption 7. Then the following
737 inequalities hold for any initial distribution ξ on (Z, \mathcal{Z}) and for all $x \in \mathbb{R}^d$:

$$\mathbb{E} \|\tilde{g}_i\|^2 \lesssim dG^2, \tag{60}$$

$$\mathbb{E} \|\tilde{g}^j - \mathbb{E}_e \tilde{g}^j\|^2 \leq \frac{dG^2}{2^j l}, \tag{61}$$

$$\mathbb{E} \|\mathbb{E}_e \tilde{g}^j - \nabla f_t\|^2 \leq \frac{4C_1 \tau G^2}{2^j l}, \tag{62}$$

$$\|\mathbb{E} \hat{g}^j - \nabla f_t\|^2 \lesssim \frac{d^2 \Delta^2}{t^2} + \frac{\tau G^2}{2^j l}, \tag{63}$$

$$\mathbb{E} \|\hat{g}^j - \nabla f_t\|^2 \lesssim \frac{d^2 \Delta^2}{t^2} + \frac{(d+\tau)G^2}{2^j l}. \tag{64}$$

738 *Proof.*

739 Proof of (60):

$$\mathbb{E} \|\tilde{g}_i\|^2 \stackrel{(11)}{=} \frac{d^2}{4t^2} \mathbb{E} |F(x + te_i, Z_i) - F(x - te_i, Z_i)|^2 \stackrel{\text{like (40)}}{\lesssim} dG^2.$$

740 Proof of (61):

$$\mathbb{E} \|\tilde{g}^j - \mathbb{E}_e \tilde{g}^j\|^2 \stackrel{\text{like (35)}}{\leq} \frac{1}{2^{2j} l^2} \sum_{i=1}^{2^j l} \mathbb{E}_Z \mathbb{E}_e \|\tilde{g}_i\|^2 \stackrel{(60)}{\leq} \frac{dG^2}{2^j l}.$$

741 Proof of (62):

$$\begin{aligned} \mathbb{E} \|\mathbb{E}_e \tilde{g}^j - \nabla f_t\|^2 &\stackrel{(15)}{=} \mathbb{E} \left\| \mathbb{E}_e \left[\frac{1}{2^j l} \sum_{i=1}^{2^j l} \tilde{g}_i \right] - \nabla f_t \right\|^2 \\ &\stackrel{(24)}{=} \mathbb{E} \left\| \mathbb{E}_e \left[\frac{1}{2^j l} \sum_{i=1}^{2^j l} \nabla F_t(x, Z_i) \right] - \nabla f_t \right\|^2 \\ &\stackrel{\textcircled{1}}{\leq} \frac{4C_1 \tau G^2}{2^j l}, \end{aligned}$$

742 where $\textcircled{1}$ uses (22) with $\sigma_2^2 = 4G^2$. Let us verify that Assumption 4' holds:

743 Unbiasedness:

$$\mathbb{E}_Z [\nabla F_t(x, Z)] \stackrel{(24)}{=} \mathbb{E}_Z \mathbb{E}_e \left[d \frac{F(x + te, Z) - F(x - te, Z)}{2t} e \right]$$

$$\begin{aligned}
&= \mathbb{E}_e \mathbb{E}_Z \left[d \frac{F(x+te, Z) - F(x-te, Z)}{2t} e \right] \\
&\stackrel{(7)}{=} \mathbb{E}_e \left[d \frac{f(x+te) - f(x-te)}{2t} e \right] \stackrel{(24)}{=} \nabla f_t(x).
\end{aligned}$$

744 Variance: $\|\nabla F_t(x, Z) - \nabla f_t(x)\| \stackrel{(77)}{\leq} 2\|\nabla F_t(x, Z)\|^2 + 2\|\nabla f_t(x)\|^2 \stackrel{\textcircled{2}}{\leq} 4G^2$,

745 where $\textcircled{2}$ uses that from Lemma 4 the smoothed f_t and F_t are differentiable, G -Lipshitz and
746 thus have norm of their gradients bounded by G .

747 Proof of (63):

$$\begin{aligned}
\|\mathbb{E}\hat{g}^j - \nabla f_t\|^2 &\stackrel{(77), (80)}{\leq} 2 \left[\mathbb{E}\|\hat{g}^j - \tilde{g}^j\|^2 + \mathbb{E}_Z \|\mathbb{E}_e \tilde{g}^j - \nabla f_t\|^2 \right] \\
&\stackrel{(31), (61)}{\lesssim} \frac{d^2 \Delta^2}{t^2} + \frac{\tau G^2}{2^j l}.
\end{aligned}$$

748 Proof of (64):

$$\begin{aligned}
\mathbb{E}\|\hat{g}^j - \nabla f_t\|^2 &\stackrel{(77)}{\leq} 3\mathbb{E} \left[\|\hat{g}^j - \tilde{g}^j\|^2 + \|\tilde{g}^j - \mathbb{E}_e \tilde{g}^j\|^2 + \|\mathbb{E}_e \tilde{g}^j - \nabla f_t\|^2 \right] \\
&\stackrel{(31), (61), (62)}{\lesssim} \frac{d^2 \Delta^2}{t^2} + \frac{(d + \tau)G^2}{2^j l}.
\end{aligned}$$

749 □

750 **Lemma 7'.** *Let Assumptions 3 and 4' hold. For any initial distribution ξ on (Z, \mathcal{Z}) the*
751 *following inequalities hold:*

752 Under Assumption 1:

$$\mathbb{E}[\|\nabla f_t(x) - \hat{g}_{ml}\|^2] \lesssim \frac{d^2 \Delta^2}{t^2} + \frac{d^2 t^2 L^2}{B} + \frac{d + \tau}{B} \sigma_2^2 + \frac{d}{B} \|\nabla f\|^2. \quad (65)$$

$$\|\nabla f_t(x) - \mathbb{E}[\hat{g}_{ml}]\|^2 \lesssim \frac{d^2 \Delta^2}{t^2} + \frac{\tau}{MB} \sigma_2^2. \quad (66)$$

753 Under Assumption 6:

$$\mathbb{E}[\|\nabla f_t(x) - \hat{g}_{ml}\|^2] \lesssim \frac{d^2 \Delta^2}{t^2} + \frac{(d + \tau)G^2}{2^j l}. \quad (67)$$

$$\|\nabla f_t(x) - \mathbb{E}[\hat{g}_{ml}]\|^2 \lesssim \frac{d^2 \Delta^2}{t^2} + \frac{\tau}{MB} G^2. \quad (68)$$

754 *Proof.* The proof is almost identical to Lemma 7, so we will leave the calculations only.

755 Proof of (66):

$$\begin{aligned}
\|\nabla f_t(x) - \mathbb{E}[\hat{g}_{ml}]\|^2 &\stackrel{(47)}{=} \left\| \nabla f_t(x) - \mathbb{E}[\hat{g}^{\lfloor \log_2 M \rfloor}] \right\|^2 \\
&\stackrel{(57)}{\lesssim} \frac{d^2 \Delta^2}{t^2} + \frac{\tau}{MB} \sigma_2^2.
\end{aligned}$$

756 Proof of (68):

$$\begin{aligned}
\|\nabla f_t(x) - \mathbb{E}[\hat{g}_{ml}]\|^2 &\stackrel{(47)}{=} \left\| \nabla f_t(x) - \mathbb{E}[\hat{g}^{\lfloor \log_2 M \rfloor}] \right\|^2 \\
&\stackrel{(63)}{\lesssim} \frac{d^2 \Delta^2}{t^2} + \frac{\tau}{MB} G^2.
\end{aligned}$$

757 Proof of (65):

$$\mathbb{E}[\|\nabla f_t(x) - \hat{g}_{ml}\|^2]$$

$$\begin{aligned}
&\leq 2\mathbb{E}[\|\nabla f_t(x) - \hat{g}^0\|^2] + 2 \sum_{j=1}^{\lfloor \log_2 M \rfloor} 2^j \mathbb{E}[\|\tilde{g}^j - \tilde{g}^{j-1}\|^2] \\
&\leq 2\mathbb{E}[\|\nabla f_t(x) - \hat{g}^0\|^2] + 4 \sum_{j=1}^{\lfloor \log_2 M \rfloor} 2^j (\mathbb{E}\|\tilde{g}^j - \nabla f_t(x)\|^2 + \mathbb{E}\|\nabla f_t(x) - \tilde{g}^{j-1}\|^2) \\
&\leq 2\mathbb{E}[\|\nabla f_t(x) - \hat{g}^0\|^2] + 16 \sum_{j=0}^{\lfloor \log_2 M \rfloor} 2^j \mathbb{E}[\|\nabla f_t(x) - \hat{g}^j\|^2] \\
&\stackrel{(59),(58)}{\lesssim} \frac{d^2 \Delta^2}{t^2} + \frac{d^2 t^2 L^2}{B} + \frac{d + \tau}{B} \sigma_2^2 + \frac{d}{B} \|\nabla f\|^2.
\end{aligned}$$

758 Proof of (67):

$$\mathbb{E}[\|\nabla f_t(x) - \hat{g}_{ml}\|^2] \stackrel{(64)}{\lesssim} \frac{d^2 \Delta^2}{t^2} + \frac{(d + \tau)G^2}{2^j l}.$$

759

□

760 E.2 Proof of Theorem 1'

761 **Theorem 1'.** *Let Assumptions 1' to 4' hold, and consider problem (4) solved by Algorithm 1.*
762 *Then, for a suitable choice of hidden parameters (with $p \simeq \frac{B}{B+d}$) and arbitrary choice of free*
763 *parameters (see Table 2), it holds that:*

$$\mathbb{E}r^N \lesssim \exp\left(-\sqrt{\frac{p^2 \mu \gamma N^2}{3}}\right) r^0 + \frac{p\sqrt{\gamma}}{\mu^{3/2}} \cdot \left[\sigma_2^2 \frac{d + \tau}{B} + t^2 \frac{L^2 d^2}{B}\right] + \frac{\Delta^2 d^2}{\mu^2 t^2} + \frac{Lt^2}{\mu}$$

764 Moreover, for arbitrary $\varepsilon \gtrsim \frac{d\Delta\sqrt{L}}{\mu^{3/2}}$ and an appropriate choice of t and γ , the number of oracle
765 calls required to ensure $r^N \lesssim \varepsilon$ is bounded by

$$B \cdot \tilde{\mathcal{O}}\left(\max\left[1, \frac{d}{B}\right] \sqrt{\frac{L}{\mu}} \log \frac{1}{\varepsilon} + \frac{(d + \tau)\sigma_2^2}{B\mu^2 \varepsilon}\right) \quad \text{two-point oracle calls.}$$

766 *Proof.* Replacing Lemma 7 with Lemma 7' in the proof of Theorem 1, we get:

$$\begin{aligned}
\mathbb{E}r^N &\lesssim \exp\left(-\sqrt{\frac{p^2 \mu \gamma N^2}{3}}\right) r^0 \\
&\quad + \frac{p\sqrt{\gamma}}{\mu^{3/2}} \cdot \left[\sigma_2^2 \frac{d + \tau}{B} + t^2 \frac{L^2 d^2}{B}\right] \\
&\quad + \frac{\Delta^2 d^2}{\mu^2 t^2} + \frac{Lt^2}{\mu}
\end{aligned}$$

767 Applying Lemma 10 with:

$$\begin{aligned}
\Gamma &= \sqrt{\gamma} \\
u &\simeq \sqrt{L} \\
a &\simeq p\sqrt{\mu} \\
b &\simeq \frac{p}{\mu^{3/2}} \cdot \left[\sigma_2^2 \frac{d + \tau}{B} + t^2 \frac{L^2 d^2}{B}\right]
\end{aligned}$$

768 We get that $r^N \lesssim \varepsilon$ takes N iterations:

$$N = \tilde{\mathcal{O}}\left(\frac{1}{p} \sqrt{\frac{L}{\mu}} \log \frac{1}{\varepsilon} + \frac{1}{B\mu^2 \varepsilon} [(d + \tau)\sigma_2^2 + L^2 t^2 d^2]\right)$$

769 Bounds on Δ , t and p are inherited: $\varepsilon \gtrsim \frac{d\Delta\sqrt{L}}{\mu^{3/2}} \Leftrightarrow \Delta \lesssim \frac{\varepsilon \mu^{3/2}}{d\sqrt{L}}$, $t \simeq \frac{\sqrt{\mu\varepsilon}}{\sqrt{L}}$, $p \simeq \frac{B}{B+d}$. Thus the
770 total number of iterations is:

$$\tilde{\mathcal{O}}\left(\left[1 + \frac{d}{B}\right] \sqrt{\frac{L}{\mu}} \log \frac{1}{\varepsilon} + \frac{(d + \tau)\sigma_2^2}{B\mu^2 \varepsilon}\right)$$

771 Finally, the oracle complexity is:

$$B \cdot \tilde{\mathcal{O}} \left(\max \left[1, \frac{d}{B} \right] \sqrt{\frac{L}{\mu}} \log \frac{1}{\varepsilon} + \frac{(d + \tau)\sigma_2^2}{B\mu^2\varepsilon} \right) \quad \text{two-point oracle calls.}$$

772

□

773 E.3 Proof of Theorem 3'

774 **Theorem 3'.** *Let Assumptions 1' to 4' hold, and consider problem (4) solved by Algorithm 1.*
 775 *Then, for a suitable choice of hidden parameters (with $p \simeq 1$) and arbitrary choice of free*
 776 *parameters (see Table 2), it holds that:*

$$\mathbb{E}r^N \lesssim \exp \left(-\sqrt{\frac{p^2\mu\gamma N^2}{3}} \right) r^0 + \frac{p\sqrt{\gamma}}{\mu^{3/2}} \cdot G^2 \frac{d + \tau}{B} + \frac{\Delta^2 d^2}{\mu^2 t^2} + \frac{Gt}{\mu}$$

777 Moreover, for arbitrary $\varepsilon \gtrsim \frac{d\Delta\sqrt{L}}{\mu^{3/2}}$ and an appropriate choice of t and γ , the number of oracle
 778 calls required to ensure $r^N \lesssim \varepsilon$ is bounded by

$$B \cdot \tilde{\mathcal{O}} \left(\sqrt{\frac{\sqrt{d}G^2}{\mu^2\varepsilon}} \log \frac{1}{\varepsilon} + \frac{(d + \tau)G^2}{B\mu^2\varepsilon} \right) \quad \text{two-point oracle calls.}$$

779 *Proof.* Replacing (65) and (66) with (67) and (68) in the proof of the smooth case we get:

$$\begin{aligned} \mathbb{E}r^N &\lesssim \exp \left(-\sqrt{\frac{p^2\mu\gamma N^2}{3}} \right) r^0 \\ &\quad + \frac{p\sqrt{\gamma}}{\mu^{3/2}} \cdot G^2 \frac{d + \tau}{B} \\ &\quad + \frac{\Delta^2 d^2}{\mu^2 t^2} + \frac{Gt}{\mu} \end{aligned}$$

780 Applying Lemma 10 with:

$$\begin{aligned} \Gamma &= \sqrt{\gamma} \\ u &\simeq \sqrt{L} \stackrel{\text{Lemma 4}}{\simeq} \sqrt{\frac{\sqrt{d}G}{t}} \\ a &\simeq p\sqrt{\mu} \\ b &\simeq \frac{p(d + \tau)G^2}{\mu^{3/2}B} \end{aligned}$$

781 We get that $r^N \lesssim \varepsilon$ takes N iterations:

$$N = \tilde{\mathcal{O}} \left(\sqrt{\frac{\sqrt{d}G}{t\mu}} \log \frac{1}{\varepsilon} + \frac{(d + \tau)G^2}{B\mu^2\varepsilon} \right)$$

782 Bounds on Δ , t and p are inherited: $\varepsilon \gtrsim \left[\frac{d\Delta G}{\mu^2} \right]^{2/3} \Leftrightarrow \Delta \lesssim \frac{\varepsilon^{3/2}\mu^2}{dG}$, $t \simeq \frac{\mu\varepsilon}{G}$, $p \simeq 1$. Thus the
 783 total number of iterations is:

$$\tilde{\mathcal{O}} \left(\sqrt{\frac{\sqrt{d}G^2}{\mu^2\varepsilon}} \log \frac{1}{\varepsilon} + \frac{(d + \tau)G^2}{B\mu^2\varepsilon} \right)$$

784 Finally, the oracle complexity is:

$$B \cdot \tilde{\mathcal{O}} \left(\sqrt{\frac{\sqrt{d}G^2}{\mu^2\varepsilon}} \log \frac{1}{\varepsilon} + \frac{(d+\tau)G^2}{B\mu^2\varepsilon} \right) \quad \text{two-point oracle calls.}$$

785

□

786 F Lower Bounds

787 F.1 Main theorems

788 First, we introduce the results that confirm the optimality of our analysis with a second
789 moment bounds. By this we mean that we check

$$\mathbb{E}_\pi |F(x, Z) - f(x)|^2 < \sigma_1^2$$

790 instead of Assumption 4 and

$$\mathbb{E}_\pi \|\nabla F(x, Z) - \nabla f(x)\|^2 < \sigma_2^2$$

791 instead of Assumption 4'.

792 Then, we show how to use clipping technique in the construction of the hard instance
793 problems to preserve the lower bounds up to logarithmic factors.

794 Our main results here are the following two theorems. They show theoretical optimality of
795 our method and analysis in both one-point and two-point regimes.

796 **Theorem 6** (one-point feedback). *For any (possibly randomized) algorithm that solves the*
797 *problem (4), there exists a function f that satisfies Assumptions 1 to 4, s.t.*

$$\mathbb{E} \|\hat{x}_N - x^*\|^2 \gtrsim \frac{\sqrt{d(d+\tau)\sigma_1^2}}{\mu\sqrt{N}} \quad \text{as } N \rightarrow \infty.$$

798 *Consequently, to get to the ε -neighborhood of the solution with one-point feedback the algorithm*
799 *needs at least*

$$N = \Omega \left(\frac{d(d+\tau)\sigma_1^2}{\mu^2\varepsilon^2} \right) \quad \text{one-point oracle calls.}$$

800 **Theorem 7** (two-point feedback). *For any (possibly randomized) algorithm that solves the*
801 *problem (4), there exists a function f that satisfies Assumptions 1' to 4', s.t.*

$$\mathbb{E} \|\hat{x}_N - x^*\|^2 \gtrsim \frac{(d+\tau)\sigma_2^2}{\mu^2 N} \quad \text{as } N \rightarrow \infty.$$

802 *Consequently, to get to the ε -neighborhood of the solution with two-point feedback one needs*
803 *at least*

$$N = \Omega \left(\frac{(d+\tau)\sigma_2^2}{\mu^2\varepsilon} \right) \quad \text{two-point oracle calls.}$$

804 We note that due to the two-part structure of the optimal rates, it is natural to prove both
805 parts separately in a regime where the part becomes dominant. We introduce those regimes:

- 806 • $\tau \geq d$ — high-correlation regime
- 807 • $\tau \leq d$ — high-dimensional regime

808 Next, we summarize the lower bounds that we claim to hold in each regime:

Table 3: Strongly convex case, lower bounds

	high-correlation		high-dimensional	
ZO 1-point	$\frac{d\tau\sigma_1^2}{\mu^2\varepsilon^2}$	(New , Theorem 8)	$\frac{d^2\sigma_1^2}{\mu^2\varepsilon^2}$	Akhavan et al. [2] (our Theorem 9)
ZO 2-point	$\frac{\tau\sigma_2^2}{\mu^2\varepsilon}$	Beznosikov et al. [5] (even for FO)	$\frac{d\sigma_2^2}{\mu^2\varepsilon}$	Duchi et al. [15] (our Theorem 10)

It becomes obvious that only 1 out of 4 bounds depend on dimension and mixing time simultaneously. For other cases, we can use existing constructions which deal with mixing and zero-order information separately and adapt them to our assumptions. Combining all four bounds, we come up with tight lower bounds in both one-point and two-point settings. Let us discuss the important related results.

Akhavan et al. [2] work with a special case of one-point feedback when noise variables do not depend on query points — this makes their lower bound applicable to our case. The only factor they do not consider is σ_1^2 , which, however, appears from their proof if used with scaled Gaussian noise, as well as additional μ^2 factor; see our Theorem 9 for the result. In the work of Beznosikov et al. [5], a first-order oracle is considered, but the hard instance problem is a 1-dimensional quadratic problem, which makes first-order and zero-order information equivalent. Duchi et al. [15] consider a general convex case of a two-point setting and provide a tight lower bound. However, their proof can be translated for strongly convex problems using the trick of adding a common quadratic part to each of the linear functions from the hard-to-distinguish family. For a more formal reduction, see Theorem 10.

Finally, we provide a, to the best of our knowledge, novel lower bound in one-point feedback and high-correlation regime.

Theorem 8 (one-point, high-correlation). *Under the conditions of Theorem 6 the following bound holds:*

$$\mathbb{E}\|\hat{x}_N - x^*\|^2 \gtrsim \frac{\sqrt{d\tau\sigma_1^2}}{\mu\sqrt{N}}.$$

Proof. Let's consider family of functions

$$f_\omega(x) = \frac{\mu}{2}\|x\|^2 + \langle S(x), \omega \rangle$$

with $\omega \in \{\pm 1\}^d$ and $S : \mathbb{R}^d \rightarrow \mathbb{R}^d$ to be chosen later. For the same values of ω , consider zeroth order oracles

$$F_\omega(x, Z) = \frac{\mu}{2}\|x\|^2 + \langle S(x), Z + \omega \rangle = f_\omega(x) + \langle S(x), Z \rangle$$

and discrete-time Markov process with transition probabilities determined by the formula

$$Z_{t+1} = \begin{cases} \xi_{t+1}, & \text{w.p. } 1/\tau, \\ Z_t, & \text{w.p. } 1 - 1/\tau, \end{cases}$$

where $\{\xi_t\}_{t=1}^\infty$ are independent and

$$\xi_t \sim \pi = \mathcal{N}(0, s^2 I_d).$$

With such pick of Z_t it is clear that Assumption 3 is satisfied and

$$\mathbb{E}_\pi F_\omega(x, Z) = f_\omega(x).$$

Now, we will prove that all algorithms fail at distinguishing between f_ω in a short amount of time. First, note that

$$\|\hat{x} - x_\omega^*\|^2 \geq \frac{1}{4}\|x_{\omega'}^* - x_\omega^*\|^2 \quad (69)$$

836 where $\omega' = \arg \min_{\tilde{\omega}} \|\hat{x} - x_{\tilde{\omega}}^*\|^2$. We will later bound $\|x_{\omega'}^* - x_{\omega}^*\|$ using Hamming distance
837 $\rho(\omega', \omega)$. But first, we bound the distance itself.
838 Applying Assouad's Lemma [55] we get

$$\max_{\omega} \mathbb{E}_{\omega} \rho(\omega', \omega) \geq \frac{d}{2} \left(1 - \min_{\rho(\omega_1, \omega_2)=1} \|P_{\omega_1} - P_{\omega_2}\|_{TV} \right) \quad (70)$$

839 where P_{ω} denotes joint distribution of outputs of F_{ω} on sequential queries produced by the
840 algorithm. And \hat{x} is the output of the algorithm after N steps. Now we bound the total
841 variation between neighbouring distributions. First, we use Pinsker's inequality:

$$2\|P_{\omega_1} - P_{\omega_2}\|_{TV}^2 \leq D_{KL}(\text{Law}(\{\omega_1 + Z_i\}_{i=1}^N), \text{Law}(\{\omega_2 + Z_i\}_{i=1}^N)) =$$

842 Then, using law of total probability, we consider a conditional KL -divergence for a fixed set
843 of indices that introduce new samples. The one step KL equals 0 if it is known that the
844 chain's state did not change. On other steps it equals to the KL between Gaussians with
845 mean ω_1 and ω_2 . We group the terms by the number of state switches k .

$$= \sum_{k=0}^N D_{KL}(\text{Law}(\{\mathcal{N}(\omega_1, s^2 I)\}_{i=1}^k), \text{Law}(\{\mathcal{N}(\omega_2, s^2 I)\}_{i=1}^k)) \mathbb{P}(|\{1 \leq t \leq N : Z_t = \xi_t\}| = k)$$

846 Using $\rho(\omega_1, \omega_2) = 1$, we simplify

$$\begin{aligned} &= \sum_{k=0}^N D_{KL}(\mathcal{N}(1, s^2 I_k), \mathcal{N}(-1, s^2 I_k)) \mathbb{P}(|\{1 \leq t \leq N : Z_t = \xi_t\}| = k) = \\ &\sum_{k=0}^N \frac{4k}{s^2} \mathbb{P}(|\{1 \leq t \leq N : Z_t = \xi_t\}| = k) = \frac{4}{s^2} \sum_{k=0}^N k \mathbb{P}(|\{1 \leq t \leq N : Z_t = \xi_t\}| = k) = \\ &\frac{4}{s^2} \mathbb{E}(|\{1 \leq t \leq N : Z_t = \xi_t\}|) = \frac{4}{s^2} \sum_{t=1}^N \mathbb{E}(I_{Z_t=\xi_t}) = \frac{4N}{s^2 \tau}. \end{aligned}$$

847 Choosing $s^2 = \frac{8N}{\tau}$ we get

$$\|P_{\omega_1} - P_{\omega_2}\|_{TV} \leq \sqrt{\frac{4N\tau}{16N\tau}} = \frac{1}{2}. \quad (71)$$

848 Now we claim that there exists such pick of $S(x)$, that satisfies Assumptions 1 to 4 and

$$\|x_{\omega'}^* - x_{\omega}^*\|^2 \geq \frac{1}{2} \frac{\sqrt{\frac{\sigma_1^2}{9}}}{\sqrt{4\mu^2 dN}} \rho(\omega', \omega) = \frac{1}{12} \frac{\sqrt{\sigma_1^2 \tau}}{\sqrt{\mu^2 dN}} \rho(\omega', \omega). \quad (72)$$

849 Combining (69), (72), (70), (71), we conclude

$$\max_{\omega} \mathbb{E}_{\omega} \|\hat{x} - x_{\omega}^*\|^2 \geq \frac{1}{48} \frac{d\sqrt{\sigma_1^2 \tau}}{\sqrt{\mu^2 dN}} = \frac{1}{48} \frac{\sqrt{d\tau\sigma_1^2}}{\mu\sqrt{N}}.$$

850 Now we should introduce $S(x)$ and check (72) and Assumptions 1 to 4.

851 Denote $\delta = \sqrt[4]{\frac{\sigma_1^2 \tau}{\mu^2 dN}}$. Let $S(x)$ be separable and

$$S(x)_i = \frac{\mu}{4} s(x_i) = \frac{\mu}{4} \cdot \begin{cases} 2\delta x_i, & 0 \leq x_i \leq \delta, \\ 3\delta^2 - (x_i - 2\delta)^2, & \delta \leq x_i \leq 2\delta, \\ 3\delta^2, & 2\delta \leq x_i. \end{cases}$$

852 And $s(x_i)$ is symmetric around zero. It is straightforward to verify that $s(x_i)$ is 2-smooth.
853 To check strong convexity and smoothness of f_{ω} we note that

$$\nabla f_{\omega}(x) = \mu x + \nabla \langle S(x), \omega \rangle = \mu x + \nabla S(x) \odot \omega,$$

where \odot is a coordinate-wise product. The Lipschitz constant of the second term is bounded

$$\|\nabla S(x) \odot \omega - \nabla S(y) \odot \omega\| = \|\nabla S(x) - \nabla S(y)\| \leq \frac{\mu}{2} \|x - y\|.$$

It means that the strong convexity constant μ and gradient Lipschitz constant L of the function f_ω are in range $[\frac{\mu}{2}; \frac{3\mu}{2}]$. Therefore, for a completely rigorous bound, we use 2μ in (72) instead of μ .

It is also straightforward to verify by stationarity condition that $x_\omega^* = -\frac{1}{2}\omega\delta$ and (72) follows. Here we also note that $\|x_\omega^*\|^2 = \frac{1}{2}d\delta^2 < 1$ for big enough N , therefore the minimizer of the function lies in the standard unit ball when the desired accuracy is small enough.

Lastly, we need to check the bounded noise assumption (4). With our current setup we can guarantee bounded variance with respect to stationary distribution

$$\mathbb{E}_\pi h^2(x, Z) = \mathbb{E}_\pi \langle S(x), Z \rangle^2 = s^2 \|S(x)\|^2 \leq \frac{9s^2\mu^2 d\delta^4}{16} = \frac{9N}{2\tau} \frac{\sigma_1^2 \tau}{N} \leq 9\sigma_1^2. \quad (73)$$

Therefore, for a completely rigorous bound, we use $\sigma_1^2/9$ in (72) instead of σ_1^2 . And a proper uniform bound is achieved via clipping, see Appendix F.2. \square

F.2 Remarks on clipping

There is, however, another problem we have to deal with — for now there is only a second-moment bound on the noise, just as in other lower bounds used that work with i.i.d. noise instead of Markovian. Tackling uniform boundness of an i.i.d. noise is straightforward — since the noise distribution is Gaussian, we can use tail bounds to clip the noise within $[-\sigma \log N; \sigma \log N]$ for all querying points with probability $1 - o(1/N)$. It gives the desired bounds up to logarithmic factors for Theorems 9 and 10.

However, in the settings of Theorem 8, this trick will not work as the algorithm can deliberately call the oracle at a point that would produce high noise on the next step. To deal with this, we clip the oracle rather than noise.

For some $t > 1$ (t is going to be logarithmic in N) we introduce

$$\hat{F}(x, Z) = \max \left(\min_\omega f_\omega(x) - t\sigma, \min(F(x, Z), \max_\omega f_\omega(x) + t\sigma) \right).$$

By construction

$$\begin{aligned} |\hat{F}(x, Z) - \mathbb{E}_\pi \hat{F}(x, Z)|^2 &\leq 2t\sigma_1^2 + 2|\max_\omega f_\omega(x) - \min_\omega f_\omega(x)|^2 = \\ &2t\sigma_1^2 + 2\|S(x)\|_1^2 \leq 2t\sigma_1^2 + 8d^2\mu^2\delta^4 = 2t\sigma_1^2 + \frac{8d^2\sigma_1^2\tau}{N}. \end{aligned}$$

Note that for big enough N , the second term becomes negligible. Now, the clipping introduces bias of the form

$$\begin{aligned} |\mathbb{E}_\pi F(x, Z) - \mathbb{E}_\pi \hat{F}(x, Z)| &\leq |\mathbb{E}_\pi h(x, Z) I_{h(x, Z) > t\sigma}| \leq \\ &\leq \int_{t\sigma}^{\infty} x e^{-\frac{x^2}{2\sigma_1^2}} dx = \sigma_1^2 \int_t^{\infty} x e^{-\frac{x^2}{2}} dx = \sigma_1^2 e^{-\frac{t^2}{2}}. \end{aligned}$$

Choosing $t \sim \log N$ makes this bias superpolynomially small in N i.e. $\lesssim \text{poly}(\frac{1}{N})$, making it within an admissible level of adversarial bias $\Delta \lesssim \frac{\varepsilon\mu^{3/2}}{dL}$. This last step, which introduces a bias, can be avoided through a careful adjustments of the Gaussian distributions used in the proof so that the mutual truncation would not result in a change of expected value. This is possible since the total probability mass that is affected by the truncation is exponentially small, therefore the total variation distance remains large after any transformations with this mass.

886 F.3 One-point high dimensional regime

887 An i.i.d. one-point setup is covered by Akhavan et al. [1], where authors considered a more
 888 general case of high-order smoothness of the objective and provided a lower bound for *any*
 889 distribution of the additive noise. Our point of view is different – we work with usual
 890 smooth functions, consider a limiting behavior when $N \rightarrow \infty$ and are free to choose the
 891 noise structure. However, we also claim stronger result - our bound shows additional $\mu^2\sigma_1^2$
 892 scaling and is asymptotically tight, according to the Theorem 1.

893 **Theorem 9** (one-point, high-dimensional). *Under the conditions of Theorem 6 the following*
 894 *bound holds:*

$$\mathbb{E}\|\hat{x}_N - x^*\|^2 \gtrsim \frac{\sqrt{d^2\sigma_1^2}}{\mu\sqrt{N}}.$$

895 *Proof.* Under closer consideration, the proof repeats, simplifies and extends the construction
 896 of Akhavan et al. [1], using our assumptions. But it will be easier for presentation to build
 897 on our own notation from Theorem 8.

898 We consider the same family of functions f_ω , but the noise is i.i.d. and point-independent
 899 Gaussian with variance σ_1^2 . This requires redefining δ and revising (71) and (72). With
 900 this noise, we use bound on the KL divergence between neighboring distributions similar to
 901 Akhavan et al. [1, Theorem 6.1]. We also use that $I_0 = \frac{1}{2\sigma_1^2}$ for Gaussian distributions. We
 902 get

$$D_{KL}(P_{\omega_1}, P_{\omega_2}) \leq \frac{N}{2\sigma_1^2} \|f_{\omega_1} - f_{\omega_2}\|_\infty^2 < \frac{N\mu^2\delta^4}{2\sigma_1^2}.$$

903 Redefining $\delta = \sqrt[4]{\frac{\sigma_1^2}{\mu^2 N}}$ we check that (71) holds. The (72) then transforms into

$$\|x_{\omega'}^* - x_\omega^*\|^2 \geq \frac{1}{2} \sqrt{\frac{\sigma_1^2}{4\mu^2 N}} \rho(\omega', \omega).$$

904 Combining (69), (72), (70), (71), we conclude

$$\max_{\omega} \mathbb{E}_{\omega} \|\hat{x} - x_{\omega}^*\|^2 \geq \frac{1}{16} \frac{d\sqrt{\sigma_1^2}}{\sqrt{\mu^2 N}}.$$

905

□

906 F.4 Two-point high dimensional regime

907 Theorem 10 below shows a reduction from the lower bound by Duchi et al. [15] to a strongly
 908 convex objectives. Coupled with the clipping technique discussed above, it concludes all the
 909 proofs of the section.

910 **Theorem 10** (two-point, high-dimensional). *Under the conditions of Theorem 7 the following*
 911 *bound holds:*

$$\mathbb{E}\|\hat{x}_N - x^*\|^2 \gtrsim \frac{d\sigma_2^2}{\mu^2 N}.$$

912 *Proof.* Let's consider family of functions for $v \in \{\pm 1\}^d$

$$f_v(x) = \frac{\mu}{2} \|x\|^2 + \delta \langle x, v \rangle$$

913 and corresponding oracles

$$F_v(x, Z) = \frac{\mu}{2} \|x\|^2 + \langle x, \delta v + Z \rangle.$$

914 The noise sequence Z_i is not given any Markovianity, instead we choose it to be i.i.d.
 915 $\sim \mathcal{N}(0, s^2 I_d)$. This family readily satisfies Assumptions 1' to 4' with the parameter $\sigma_2^2 \geq$
 916 $\mathbb{E}\|Z\|^2 = ds^2$. Again, here we consider only a second moment bound, as discussed above.

917 This construction is similar to the one used in a proof by Duchi et al. [15, Proposition 1],
 918 but here we add a deterministic quadratic part, as we work with a strongly convex problems.
 919 Therefore, there is always a global minimizer of the function

$$x_v^* = \arg \min f_v(x) = -\frac{\delta}{\mu}v.$$

920 As usual, we can bound distance to the optima with the Hamming distance between the
 921 signs of the estimate and the optima

$$\max_v \mathbb{E}\|\hat{x}_N - x_v^*\|^2 \geq \frac{\delta^2}{\mu^2} \sum_{i=1}^d \mathbb{P}(\text{sign}(\hat{x}_N^i) \neq -\text{sign}(v^i)).$$

922 Duchi et al. [15] prove a lower bound on the sum of such probabilities

$$\sum_{i=1}^d \mathbb{P}(\text{sign}(\hat{x}_N^i) \neq -\text{sign}(v^i)) \geq d \left(1 - \sqrt{\frac{2N\delta^2}{ds^2}}\right).$$

923 This inequality also applies to our set of functions as they differ only by a common deter-
 924 ministic function. Therefore, we get

$$\max_v \mathbb{E}\|\hat{x}_N - x_v^*\|^2 \geq \frac{d\delta^2}{\mu^2} \left(1 - \sqrt{\frac{2N\delta^2}{ds^2}}\right).$$

925 Choosing $s^2 = \frac{\sigma_2^2}{d}$ and $\delta^2 = \frac{\sigma_2^2}{4N}$ gives the desired result

$$\mathbb{E}\|\hat{x}_N - x^*\|^2 \gtrsim \frac{d\sigma_2^2}{\mu^2 N}.$$

926

□

927 G Basic Facts

928 **Lemma 12.** *If f is L -smooth in \mathbb{R}^d , then for any $x, y \in \mathbb{R}^d$*

$$f(x) - f(y) - \langle \nabla f(y), x - y \rangle \leq \frac{L}{2} \|x - y\|^2. \quad (74)$$

929 **Lemma 13** (Cauchy Schwartz inequality). *For any $a, b, x_1, \dots, x_n \in \mathbb{R}^d$ and $c > 0$ the*
 930 *following inequalities hold:*

$$2\langle a, b \rangle \leq \frac{\|a\|^2}{c} + c\|b\|^2, \quad (75)$$

931

$$\|a + b\|^2 \leq \left(1 + \frac{1}{c}\right) \|a\|^2 + (1 + c)\|b\|^2, \quad (76)$$

932

$$\left\| \sum_{i=1}^n x_i \right\|^2 \leq n \cdot \sum_{i=1}^n \|x_i\|^2. \quad (77)$$

933 **Lemma 14.** *For a random variable ξ with a finite second moment:*

$$\mathbb{E}\|\xi - \mathbb{E}\xi\|^2 \leq \mathbb{E}\|\xi\|^2. \quad (78)$$

934 **Lemma 15** (Jensen's inequality). *If f is a convex function, then for any $n \in \mathbb{N}^*$ and*
 935 *$x_1, \dots, x_n \in \mathbb{R}^d$ the following inequality holds:*

$$f\left(\frac{1}{n} \sum_{i=1}^n x_i\right) \leq \frac{1}{n} \sum_{i=1}^n f(x_i). \quad (79)$$

936 *Probabilistic form:*

$$f(\mathbb{E}[X]) \leq \mathbb{E}[f(X)].$$

937 *Applied to $f(X) = \|X\|^2$:*

$$\|\mathbb{E}[X]\|^2 \leq \mathbb{E}[\|X\|^2]. \quad (80)$$

938 **Lemma 16** (Norm of random projection). *For $e \sim RS_2^d(1)$ the following equality holds:*

$$\mathbb{E}_e \langle v, e \rangle^2 = \|v\|^2 \cdot 1/d. \quad (81)$$

Proof.

$$\mathbb{E} \langle v, e \rangle^2 = \|v\|^2 \mathbb{E} \langle v/\|v\|, e \rangle^2 = \|v\|^2 \mathbb{E} \langle (1, 0, \dots, 0), \tilde{e} \rangle^2 = \|v\|^2 \mathbb{E} [\tilde{e}_1]^2 \stackrel{\textcircled{1}}{=} \|v\|^2 \cdot 1/d,$$

939 where $\textcircled{1}$ uses $\sum_i \mathbb{E} [\tilde{e}_i]^2 = 1$ and $E [\tilde{e}_1]^2 = \mathbb{E} [\tilde{e}_2]^2 = \dots$ \square

NeurIPS Paper Checklist

1. Claims

Question: Do the main claims made in the abstract and introduction accurately reflect the paper’s contributions and scope?

Answer: [\[Yes\]](#)

Justification: main contributions of this paper are described accurately in a dedicated subsection (Section 1.2) of the introduction.

Guidelines:

- The answer NA means that the abstract and introduction do not include the claims made in the paper.
- The abstract and/or introduction should clearly state the claims made, including the contributions made in the paper and important assumptions and limitations. A No or NA answer to this question will not be perceived well by the reviewers.
- The claims made should match theoretical and experimental results, and reflect how much the results can be expected to generalize to other settings.
- It is fine to include aspirational goals as motivation as long as it is clear that these goals are not attained by the paper.

2. Limitations

Question: Does the paper discuss the limitations of the work performed by the authors?

Answer: [\[Yes\]](#)

Justification: assumptions we use to prove the main results are presented in Section 2. The motivation for these assumptions as well their limitations are also described there.

Guidelines:

- The answer NA means that the paper has no limitation while the answer No means that the paper has limitations, but those are not discussed in the paper.
- The authors are encouraged to create a separate "Limitations" section in their paper.
- The paper should point out any strong assumptions and how robust the results are to violations of these assumptions (e.g., independence assumptions, noiseless settings, model well-specification, asymptotic approximations only holding locally). The authors should reflect on how these assumptions might be violated in practice and what the implications would be.
- The authors should reflect on the scope of the claims made, e.g., if the approach was only tested on a few datasets or with a few runs. In general, empirical results often depend on implicit assumptions, which should be articulated.
- The authors should reflect on the factors that influence the performance of the approach. For example, a facial recognition algorithm may perform poorly when image resolution is low or images are taken in low lighting. Or a speech-to-text system might not be used reliably to provide closed captions for online lectures because it fails to handle technical jargon.
- The authors should discuss the computational efficiency of the proposed algorithms and how they scale with dataset size.
- If applicable, the authors should discuss possible limitations of their approach to address problems of privacy and fairness.
- While the authors might fear that complete honesty about limitations might be used by reviewers as grounds for rejection, a worse outcome might be that reviewers discover limitations that aren’t acknowledged in the paper. The authors should use their best judgment and recognize that individual actions in favor of transparency play an important role in developing norms that preserve the integrity of the community. Reviewers will be specifically instructed to not penalize honesty concerning limitations.

3. Theory assumptions and proofs

Question: For each theoretical result, does the paper provide the full set of assumptions and a complete (and correct) proof?

Answer: [\[Yes\]](#)

Justification: all assumptions and definitions are carefully stated. The complete proofs appear in the supplemental material and are properly referenced in the main part.

Guidelines:

- The answer NA means that the paper does not include theoretical results.
- All the theorems, formulas, and proofs in the paper should be numbered and cross-referenced.
- All assumptions should be clearly stated or referenced in the statement of any theorems.
- The proofs can either appear in the main paper or the supplemental material, but if they appear in the supplemental material, the authors are encouraged to provide a short proof sketch to provide intuition.
- Inversely, any informal proof provided in the core of the paper should be complemented by formal proofs provided in appendix or supplemental material.
- Theorems and Lemmas that the proof relies upon should be properly referenced.

4. Experimental result reproducibility

Question: Does the paper fully disclose all the information needed to reproduce the main experimental results of the paper to the extent that it affects the main claims and/or conclusions of the paper (regardless of whether the code and data are provided or not)?

Answer: [\[Yes\]](#)

Justification: see Section 3. The setup is fully disclosed.

Guidelines:

- The answer NA means that the paper does not include experiments.
- If the paper includes experiments, a No answer to this question will not be perceived well by the reviewers: Making the paper reproducible is important, regardless of whether the code and data are provided or not.
- If the contribution is a dataset and/or model, the authors should describe the steps taken to make their results reproducible or verifiable.
- Depending on the contribution, reproducibility can be accomplished in various ways. For example, if the contribution is a novel architecture, describing the architecture fully might suffice, or if the contribution is a specific model and empirical evaluation, it may be necessary to either make it possible for others to replicate the model with the same dataset, or provide access to the model. In general, releasing code and data is often one good way to accomplish this, but reproducibility can also be provided via detailed instructions for how to replicate the results, access to a hosted model (e.g., in the case of a large language model), releasing of a model checkpoint, or other means that are appropriate to the research performed.
- While NeurIPS does not require releasing code, the conference does require all submissions to provide some reasonable avenue for reproducibility, which may depend on the nature of the contribution. For example
 - (a) If the contribution is primarily a new algorithm, the paper should make it clear how to reproduce that algorithm.
 - (b) If the contribution is primarily a new model architecture, the paper should describe the architecture clearly and fully.
 - (c) If the contribution is a new model (e.g., a large language model), then there should either be a way to access this model for reproducing the results or a way to reproduce the model (e.g., with an open-source dataset or instructions for how to construct the dataset).

1047 (d) We recognize that reproducibility may be tricky in some cases, in which
1048 case authors are welcome to describe the particular way they provide for
1049 reproducibility. In the case of closed-source models, it may be that access to
1050 the model is limited in some way (e.g., to registered users), but it should be
1051 possible for other researchers to have some path to reproducing or verifying
1052 the results.

1053 5. Open access to data and code

1054 Question: Does the paper provide open access to the data and code, with sufficient
1055 instructions to faithfully reproduce the main experimental results, as described in
1056 supplemental material?

1057 Answer: [No]

1058 Justification: our experiments are rather a practical confirmation of theoretical
1059 results, and these experiments can be easily reproduced.

1060 Guidelines:

- 1061 • The answer NA means that paper does not include experiments requiring code.
- 1062 • Please see the NeurIPS code and data submission guidelines ([https://nips.](https://nips.cc/public/guides/CodeSubmissionPolicy)
1063 [cc/public/guides/CodeSubmissionPolicy](https://nips.cc/public/guides/CodeSubmissionPolicy)) for more details.
- 1064 • While we encourage the release of code and data, we understand that this might
1065 not be possible, so “No” is an acceptable answer. Papers cannot be rejected
1066 simply for not including code, unless this is central to the contribution (e.g., for
1067 a new open-source benchmark).
- 1068 • The instructions should contain the exact command and environment needed
1069 to run to reproduce the results. See the NeurIPS code and data submis-
1070 sion guidelines (<https://nips.cc/public/guides/CodeSubmissionPolicy>)
1071 for more details.
- 1072 • The authors should provide instructions on data access and preparation, in-
1073 cluding how to access the raw data, preprocessed data, intermediate data, and
1074 generated data, etc.
- 1075 • The authors should provide scripts to reproduce all experimental results for
1076 the new proposed method and baselines. If only a subset of experiments are
1077 reproducible, they should state which ones are omitted from the script and why.
- 1078 • At submission time, to preserve anonymity, the authors should release
1079 anonymized versions (if applicable).
- 1080 • Providing as much information as possible in supplemental material (appended
1081 to the paper) is recommended, but including URLs to data and code is permitted.

1082 6. Experimental setting/details

1083 Question: Does the paper specify all the training and test details (e.g., data splits,
1084 hyperparameters, how they were chosen, type of optimizer, etc.) necessary to
1085 understand the results?

1086 Answer: [Yes]

1087 Justification: see Section 3, all parameters are described there.

1088 Guidelines:

- 1089 • The answer NA means that the paper does not include experiments.
- 1090 • The experimental setting should be presented in the core of the paper to a level
1091 of detail that is necessary to appreciate the results and make sense of them.
- 1092 • The full details can be provided either with the code, in appendix, or as
1093 supplemental material.

1094 7. Experiment statistical significance

1095 Question: Does the paper report error bars suitably and correctly defined or other
1096 appropriate information about the statistical significance of the experiments?

1097 Answer: [No]

Justification: we use experiments to verify the theoretical rates and have no statistical effects associated with running the experiments.

Guidelines:

- The answer NA means that the paper does not include experiments.
- The authors should answer "Yes" if the results are accompanied by error bars, confidence intervals, or statistical significance tests, at least for the experiments that support the main claims of the paper.
- The factors of variability that the error bars are capturing should be clearly stated (for example, train/test split, initialization, random drawing of some parameter, or overall run with given experimental conditions).
- The method for calculating the error bars should be explained (closed form formula, call to a library function, bootstrap, etc.)
- The assumptions made should be given (e.g., Normally distributed errors).
- It should be clear whether the error bar is the standard deviation or the standard error of the mean.
- It is OK to report 1-sigma error bars, but one should state it. The authors should preferably report a 2-sigma error bar than state that they have a 96% CI, if the hypothesis of Normality of errors is not verified.
- For asymmetric distributions, the authors should be careful not to show in tables or figures symmetric error bars that would yield results that are out of range (e.g. negative error rates).
- If error bars are reported in tables or plots, The authors should explain in the text how they were calculated and reference the corresponding figures or tables in the text.

8. Experiments compute resources

Question: For each experiment, does the paper provide sufficient information on the computer resources (type of compute workers, memory, time of execution) needed to reproduce the experiments?

Answer: [No]

Justification: the experiments performed are not computationally heavy and can be reproduced on an average machine in a fairly reasonable amount of time.

Guidelines:

- The answer NA means that the paper does not include experiments.
- The paper should indicate the type of compute workers CPU or GPU, internal cluster, or cloud provider, including relevant memory and storage.
- The paper should provide the amount of compute required for each of the individual experimental runs as well as estimate the total compute.
- The paper should disclose whether the full research project required more compute than the experiments reported in the paper (e.g., preliminary or failed experiments that didn't make it into the paper).

9. Code of ethics

Question: Does the research conducted in the paper conform, in every respect, with the NeurIPS Code of Ethics <https://neurips.cc/public/EthicsGuidelines>?

Answer: [Yes]

Justification: the research follows the NeurIPS Code of Ethics.

Guidelines:

- The answer NA means that the authors have not reviewed the NeurIPS Code of Ethics.
- If the authors answer No, they should explain the special circumstances that require a deviation from the Code of Ethics.
- The authors should make sure to preserve anonymity (e.g., if there is a special consideration due to laws or regulations in their jurisdiction).

10. Broader impacts

Question: Does the paper discuss both potential positive societal impacts and negative societal impacts of the work performed?

Answer: [NA]

Justification: there is no societal impact of the work performed – we only develop the theoretical understanding of Optimization.

Guidelines:

- The answer NA means that there is no societal impact of the work performed.
- If the authors answer NA or No, they should explain why their work has no societal impact or why the paper does not address societal impact.
- Examples of negative societal impacts include potential malicious or unintended uses (e.g., disinformation, generating fake profiles, surveillance), fairness considerations (e.g., deployment of technologies that could make decisions that unfairly impact specific groups), privacy considerations, and security considerations.
- The conference expects that many papers will be foundational research and not tied to particular applications, let alone deployments. However, if there is a direct path to any negative applications, the authors should point it out. For example, it is legitimate to point out that an improvement in the quality of generative models could be used to generate deepfakes for disinformation. On the other hand, it is not needed to point out that a generic algorithm for optimizing neural networks could enable people to train models that generate Deepfakes faster.
- The authors should consider possible harms that could arise when the technology is being used as intended and functioning correctly, harms that could arise when the technology is being used as intended but gives incorrect results, and harms following from (intentional or unintentional) misuse of the technology.
- If there are negative societal impacts, the authors could also discuss possible mitigation strategies (e.g., gated release of models, providing defenses in addition to attacks, mechanisms for monitoring misuse, mechanisms to monitor how a system learns from feedback over time, improving the efficiency and accessibility of ML).

11. Safeguards

Question: Does the paper describe safeguards that have been put in place for responsible release of data or models that have a high risk for misuse (e.g., pretrained language models, image generators, or scraped datasets)?

Answer: [NA]

Justification: the paper poses no such risks.

Guidelines:

- The answer NA means that the paper poses no such risks.
- Released models that have a high risk for misuse or dual-use should be released with necessary safeguards to allow for controlled use of the model, for example by requiring that users adhere to usage guidelines or restrictions to access the model or implementing safety filters.
- Datasets that have been scraped from the Internet could pose safety risks. The authors should describe how they avoided releasing unsafe images.
- We recognize that providing effective safeguards is challenging, and many papers do not require this, but we encourage authors to take this into account and make a best faith effort.

12. Licenses for existing assets

Question: Are the creators or original owners of assets (e.g., code, data, models), used in the paper, properly credited and are the license and terms of use explicitly mentioned and properly respected?

Answer: [NA]

1203 Justification: the paper does not use existing assets.

1204 Guidelines:

- 1205 • The answer NA means that the paper does not use existing assets.
- 1206 • The authors should cite the original paper that produced the code package or
- 1207 dataset.
- 1208 • The authors should state which version of the asset is used and, if possible,
- 1209 include a URL.
- 1210 • The name of the license (e.g., CC-BY 4.0) should be included for each asset.
- 1211 • For scraped data from a particular source (e.g., website), the copyright and
- 1212 terms of service of that source should be provided.
- 1213 • If assets are released, the license, copyright information, and terms of use in
- 1214 the package should be provided. For popular datasets, [paperswithcode.com/](https://paperswithcode.com/datasets)
- 1215 [datasets](https://paperswithcode.com/datasets) has curated licenses for some datasets. Their licensing guide can help
- 1216 determine the license of a dataset.
- 1217 • For existing datasets that are re-packaged, both the original license and the
- 1218 license of the derived asset (if it has changed) should be provided.
- 1219 • If this information is not available online, the authors are encouraged to reach
- 1220 out to the asset’s creators.

1221 **13. New assets**

1222 Question: Are new assets introduced in the paper well documented and is the

1223 documentation provided alongside the assets?

1224 Answer: [NA]

1225 Justification: the paper does not propose new assets.

1226 Guidelines:

- 1227 • The answer NA means that the paper does not release new assets.
- 1228 • Researchers should communicate the details of the dataset/code/model as part
- 1229 of their submissions via structured templates. This includes details about
- 1230 training, license, limitations, etc.
- 1231 • The paper should discuss whether and how consent was obtained from people
- 1232 whose asset is used.
- 1233 • At submission time, remember to anonymize your assets (if applicable). You
- 1234 can either create an anonymized URL or include an anonymized zip file.

1235 **14. Crowdsourcing and research with human subjects**

1236 Question: For crowdsourcing experiments and research with human subjects, does

1237 the paper include the full text of instructions given to participants and screenshots,

1238 if applicable, as well as details about compensation (if any)?

1239 Answer: [NA]

1240 Justification: the paper does not involve crowdsourcing nor research with human

1241 subjects.

1242 Guidelines:

- 1243 • The answer NA means that the paper does not involve crowdsourcing nor
- 1244 research with human subjects.
- 1245 • Including this information in the supplemental material is fine, but if the main
- 1246 contribution of the paper involves human subjects, then as much detail as
- 1247 possible should be included in the main paper.
- 1248 • According to the NeurIPS Code of Ethics, workers involved in data collection,
- 1249 curation, or other labor should be paid at least the minimum wage in the
- 1250 country of the data collector.

1251 **15. Institutional review board (IRB) approvals or equivalent for research**

1252 **with human subjects**

1253 Question: Does the paper describe potential risks incurred by study participants,
 1254 whether such risks were disclosed to the subjects, and whether Institutional Review
 1255 Board (IRB) approvals (or an equivalent approval/review based on the requirements
 1256 of your country or institution) were obtained?

1257 Answer: [NA]

1258 Justification: the paper does not involve crowdsourcing nor research with human
 1259 subjects.

1260 Guidelines:

- 1261 • The answer NA means that the paper does not involve crowdsourcing nor
 1262 research with human subjects.
- 1263 • Depending on the country in which research is conducted, IRB approval (or
 1264 equivalent) may be required for any human subjects research. If you obtained
 1265 IRB approval, you should clearly state this in the paper.
- 1266 • We recognize that the procedures for this may vary significantly between insti-
 1267 tutions and locations, and we expect authors to adhere to the NeurIPS Code of
 1268 Ethics and the guidelines for their institution.
- 1269 • For initial submissions, do not include any information that would break
 1270 anonymity (if applicable), such as the institution conducting the review.

1271 **16. Declaration of LLM usage**

1272 Question: Does the paper describe the usage of LLMs if it is an important, original,
 1273 or non-standard component of the core methods in this research? Note that if
 1274 the LLM is used only for writing, editing, or formatting purposes and does not
 1275 impact the core methodology, scientific rigorousness, or originality of the research,
 1276 declaration is not required.

1277 Answer: [NA]

1278 Justification: LLMs were used only for editing.

1279 Guidelines:

- 1280 • The answer NA means that the core method development in this research does
 1281 not involve LLMs as any important, original, or non-standard components.
- 1282 • Please refer to our LLM policy (<https://neurips.cc/Conferences/2025/LLM>)
 1283 for what should or should not be described.