# **Accelerated Parameter-Free Stochastic Optimization**

Itai Kreisler Maor Ivgi Oliver Hinder Yair Carmon KREISLER@MAIL.TAU.AC.IL
MAOR.IVGI@CS.TAU.AC.IL
OHINDER@PITT.EDU
YCARMON@TAUEX.TAU.AC.IL

Editors: Shipra Agrawal and Aaron Roth

#### **Abstract**

We propose a method that achieves near-optimal rates for *smooth* stochastic convex optimization and requires essentially no prior knowledge of problem parameters. This improves on prior work which requires knowing at least the initial distance to optimality  $d_0$ . Our method, U-DoG, combines UniXGrad (Kavis et al. [30]) and DoG (Ivgi et al. [27]) with novel iterate stabilization techniques. It requires only loose bounds on  $d_0$  and the noise magnitude, provides high probability guarantees under sub-Gaussian noise, and is also near-optimal in the non-smooth case. Our experiments show consistent, strong performance on convex problems and mixed results on neural network training.

**Keywords:** Parameter-free, Adaptive, Stochastic convex optimization, Smooth optimization.

# 1. Introduction

We consider the problem of minimizing a smooth convex function using access to an unbiased stochastic gradient oracle. This is a fundamental problem in machine learning, including many important special cases such as logistic and linear regression. Moreover, the smoothness assumption is crucial for developing one of the most widely used improvements for the classical gradient method: Nesterov acceleration [44].

Nesterov acceleration obtains the optimal rate of convergence for this problem but is strongly reliant on knowing the problem parameters. Specifically, Lan [35], who first demonstrated the theoretical value of Nesterov acceleration on smooth *stochastic* convex functions, requires knowledge of the smoothness parameter  $\beta$ , the distance  $d_0$  from the initial point to the optimum, and a value  $\sigma$  for which the noise is  $\sigma$ -sub-Gaussian. Accelerated adaptive methods [14, 30] do not require knowledge of  $\beta$  and  $\sigma$ , but assume knowledge of  $d_0$ . For *non-smooth* stochastic convex optimization, *parameter-free methods* [e.g., 7, 9, 16, 27, 28, 41, 49] require only loose knowledge of problem parameters to obtain near-optimal rates. Finding such parameter-free methods for *smooth* stochastic optimization is a longstanding open problem.

Our contribution. We solve this open problem, designing an accelerated parameter-free method which we call UNIXGRAD-DOG, or U-DOG for short. U-DOG combines the "universal extragradient" (UNIXGRAD) framework [30] with the "distance over gradient" (DOG) technique [27]. More specifically, we replace the domain diameter D in the UNIXGRAD step size numerator with the maximum distance from the initial point, similar to the DoG step size numerator. Furthermore, we use this maximum distance to automatically tune the "momentum" parameter  $\alpha_t$  of UNIXGRAD.

Algorithm name	Unbounded domain?	Insensitive to $d_0/D$ $\beta$ $\sigma$			Rate of convergence	High probability?
U-DoG (this work)	✓	✓	<b>√</b>	<b>✓</b>	$\widetilde{O}\left(\frac{\beta d_0^2}{T^2} + \frac{\sigma d_0}{\sqrt{T}} + \frac{\hat{\sigma} d_0}{T}\right)$	✓
U-DOG (tills work)	×	✓	✓	✓	$\widetilde{O}\left(\frac{\beta D^2}{T^2} + \frac{\sigma D}{\sqrt{T}}\right)$	✓
UNIXGRAD [30]	×	Х	1	✓	$O\left(\frac{\beta D^2}{T^2} + \frac{\sigma D}{\sqrt{T}}\right)$	Х
Cutkosky [14]	✓	X	✓	✓	$\widetilde{O}\left(\frac{\beta d_0^2}{T^2} + \frac{\sigma d_0}{\sqrt{T}}\right)$	X
Lan [35]	✓	X	X	X	$O\left(\frac{\beta d_0^2}{T^2} + \frac{\sigma d_0}{\sqrt{T}}\right)$	✓
DoG [27] / CO [16]	✓	✓	✓	×	$\widetilde{O}\left(\frac{\beta d_0^2}{T} + \frac{\sigma d_0}{\sqrt{T}} + \frac{\hat{L}d_0}{T}\right)$	<b>√</b> / <b>X</b>

Table 1: Comparison of U-DoG and prior work on  $\beta$ -smooth stochastic optimization with  $\sigma$ -sub-Gaussian noise. "Unbounded domain" indicates if the algorithm is defined over the whole Euclidean space or a bounded subspace. In the former case we express rates in terms of the initial distance to optimality  $d_0$  and in the latter case we use the domain diameter D. Under "Insensitive to..." we mark  $\msec I$  if the suboptimality bound grows polynomially with error in the parameter,  $\msec I$  if it only affects logarithmic factors or low order terms, and  $\msec I$  if there is no dependence on the parameter at all. The marker  $\msec I$  indicates algorithms that require an upper bound  $\hat{L}$  on gradient norm, which may be much larger the the upper bound  $\hat{\sigma}$  on the noise. The notation  $\widetilde{O}(\cdot)$  hides polylogarithmic factors.

Finally, we modify the UNIXGRAD step size denominator to ensure the stability of the iterate sequence. U-DoG only requires a loose upper bound  $\hat{\sigma}$  on  $\sigma$  and lower bound  $r_{\epsilon}$  on  $D.^1$  As long as  $\hat{\sigma}$  is loose by at most a  $\sqrt{T}$  factor and  $r_{\epsilon}$  is loose by any poly(T) factor, we obtain a near-optimal, high-probability rate of convergence; Table 1 states U-DoG's guarantees and compares it to prior work. Moreover, U-DoG simultaneously enjoys a near-optimal, parameter-free rate of convergence for *non-smooth* problems.

We conduct preliminary experiments with U-DoG as well as another algorithm, A-DoG, which combines ACCELEGRAD [36] and DoG. On convex optimization problems, both U-DoG and A-DoG often substantially improve over DoG, especially at large batch sizes, with A-DoG outperforming U-DoG, likely due to not requiring an extra-gradient computation at each step. On several problems, A-DoG matches the performance of carefully tuned SGD with Nesterov momentum. On neural network optimization problems, however, we observe that both U-DoG and A-DoG do not consistently improve over DoG.

#### 1.1. Related work

**Non-smooth stochastic optimization.** The majority of tuning-insensitive stochastic optimization methods are developed for online convex optimization. Online regret bounds immediately translate to suboptimality guarantees for non-smooth stochastic optimization using online-to-batch conversion [48, Section 3]. Proposed methods divide roughly into *adaptive* algorithms such as adaptive SGD [22, 38], AdaGrad [21, 40] and variants [e.g., 33, 55, 58], and *parameter-free* methods [7, 15, 16, 28, 39, 41, 47, 49, 59]. Adaptive methods typically require no knowledge of the stochastic gradient bound but need to know the initial distance to optimality (or the domain diameter), while

<sup>1.</sup> In fact, we only require *local* upper bounds of the form  $\hat{\sigma}(x)$  on the noise sub-Gaussianity.

parameter-free methods are robust to uncertainty in the distance but require some (loose) bound on the stochastic gradient norms.

Recent work [9, 27] develops parameter-free methods that hew closer to SGD and eschew online-to-batch conversion for high-probability guarantees in the stochastic setting; U-DoG continues this line. In particular, it extends the core mechanism of DoG [27] wherein iterate movement serves as a proxy for the distance to optimality. D-Adaptation [17], DoWG [32], and Prodigy [42] use a similar mechanism, but only provide guarantees for the non-stochastic setting. Ensuring the validity of the mechanism (i.e., that iterates never move too far away from the optimum) is a key challenge in its analysis. This challenge becomes greater in the smooth setting, where selecting too small of a step size nullifies the benefit of acceleration. Much of our algorithmic and analytical innovation addresses this challenge.

Non-stochastic smooth optimization. Without noise, Nesterov acceleration requires knowledge of the smoothness constant  $\beta$  but not the distance to optimality [44, 45]. The methods [30, 36] reverse this tradeoff, requiring the distance but not  $\beta$ . Line search techniques such as [6, 11] provide much stronger adaptivity, attaining the optimal gradient evaluation complexity up to an additive term that depends logarithmically on the uncertainty in  $\beta$ . However, line search can be challenging to employ efficiently in the stochastic setting as we can no longer accurately evaluate the function. Indeed, there are many works that analyze stochastic line search techniques [e.g., 50, 60] but none have obtained convergence guarantees close to that of Lan [35].

Smooth stochastic optimization. Several adaptive and parameter-free methods [9, 16, 22, 27, 32] converge faster on smooth functions. However, they do not improve all the way to the optimal rate (see Table 1) due to a missing "momentum" component. Cutkosky [14] gives an improved online-to-batch conversion framework that endows adaptive SGD with momentum and accelerated rates in the smooth case, but requires a bound on the distance to optimality. Kavis et al. [30] propose UNIXGRAD, combining ideas from [14] with the mirror-prox/extragradient algorithm [19, 43] and online learning [38, 54] to obtain optimal rates assuming bounded domains of known diameter D and assuming that  $d_0$  is of the order of D. U-DoG modifies UNIXGRAD and removes both assumptions, yielding the first parameter-free accelerated method.

# 2. Preliminaries and algorithmic framework

In this section, we set up our notation and terminology, and use them to present the general U-DoG template (Algorithm 1) defining the algorithm up to the choice of adaptive step sizes, which we gradually develop in the following sections.

**Basic notation and conventions.** Throughout,  $\|\cdot\|$  denotes the Euclidean norm,  $\log$  is base e and  $\log_+(x) \coloneqq 1 + \log(x)$ . The function  $\operatorname{Proj}_{\mathcal{X}}(\cdot)$  denotes Euclidean projection onto set  $\mathcal{X}$ . We say that  $f: \mathcal{K} \to \mathbb{R}$  is  $\beta$ -smooth if  $\nabla f$  is  $\beta$ -Lipschitz, i.e.,  $\|\nabla f(u) - \nabla f(v)\| \le \beta \|u - v\|$  for all  $u, v \in \mathcal{K}$ . We write  $[\cdot]_+ := \max\{\cdot, 0\}$ .

In this work, we minimize an objective function f via queries to a stochastic gradient estimator  $\mathcal{G}$ . We make the following assumption in all of our theoretical analyses.

**Assumption 1 (Made throughout)** The objective function  $f: \mathcal{K} \to \mathbb{R}$  is convex, L-Lipschitz,  $\beta$ smooth, has closed convex domain K, and its minimum is attained at some  $x_* \in \arg\min_{x \in K} f(x)$ . For all  $x \in \mathcal{K}$ , the gradient estimator  $\mathcal{G}$  satisfies  $\mathbb{E}\mathcal{G}(x) = \nabla f(x)$ .

# **Algorithm 1:** U-DoG (UNIXGRAD-DoG) template

**Input:** Initial  $x_0 \in \mathcal{K}$ , iteration budget T, initial movement  $r_{\epsilon}$ , step sizes  $\{\eta_{x,t}, \eta_{u,t}\}$ 

- 1 Set  $y_0 = x_0$
- 2 for  $t = 0, 1, 2, \dots, T 1$  do

3 Set 
$$\alpha_t = \sum_{k=0}^t \bar{r}_k / \bar{r}_t$$
 and  $\omega_t = \alpha_t \bar{r}_t$  for  $\bar{r}_t = \max_{k \le t} \max\{\|y_k - x_0\|, \|x_k - x_0\|, r_\epsilon\}$ 

4 
$$x_{t+1} = \operatorname{Proj}_{\mathcal{K}}(y_t - \alpha_t \eta_{x,t} m_t)$$
 for  $m_t \sim \mathcal{G}(\hat{z}_t)$  and  $\hat{z}_t = \frac{\omega_t y_t + \sum_{k=0}^{t-1} \omega_k x_{k+1}}{\sum_{k=0}^{t} \omega_k}$ 

$$\mathbf{4} \qquad x_{t+1} = \operatorname{Proj}_{\mathcal{K}}(y_t - \alpha_t \eta_{x,t} m_t) \quad \text{for} \quad m_t \sim \mathcal{G}(\hat{z}_t) \quad \text{and} \quad \hat{z}_t = \frac{\omega_t y_t + \sum_{k=0}^{t-1} \omega_k x_{k+1}}{\sum_{k=0}^{t} \omega_k}$$

$$\mathbf{5} \qquad y_{t+1} = \operatorname{Proj}_{\mathcal{K}}(y_t - \alpha_t \eta_{y,t} g_t) \quad \text{for} \quad g_t \sim \mathcal{G}(\hat{x}_t) \quad \text{and} \quad \hat{x}_t = \frac{\omega_t x_{t+1} + \sum_{k=0}^{t-1} \omega_k x_{k+1}}{\sum_{k=0}^{t} \omega_k}$$

6 end

7 return  $\hat{x}_T$ 

**Presenting U-DoG.** Algorithm 1 provides the general template of U-DoG. As in UNIXGRAD [30], each iteration of the algorithm consists of two stochastic gradient steps, with each stochastic gradient queried at a moving average of iterates. Unlike UNIXGRAD, the moving average weights  $\omega_t$  and the step size multipliers  $\alpha_t$  are not fixed in advance, but are instead dynamically set based on the maximum distance moved from the origin, denoted

$$\bar{r}_t \coloneqq \max_{k < t} \max\{\|y_k - x_0\|, \|x_k - x_0\|, r_\epsilon\}.$$

The parameter  $r_{\epsilon}$  serves as a (loose) lower bound on  $||x_0 - x_{\star}||$ ; typically,  $\bar{r}_t$  grows rapidly and then plateaus at a level roughly approximating  $||x_0 - x_*||$ . When that happens, the sequence  $\alpha_t =$  $\sum_{k \le t} \bar{r}_k / \bar{r}_t$  grows linearly in t, similar to  $\alpha_t = t + 1$  in UNIXGRAD.

To complete the specification of U-DoG we must set the step size sequence. UNIXGRAD assumes  $\mathcal K$  the domain has Euclidean diameter D and picks step sizes of the form  $\eta_{x,t}=\eta_{y,t}=0$  $\frac{\sqrt{2}D}{\sqrt{1+Q_{t-1}}}$  where

$$Q_t \coloneqq \sum_{k=0}^t q_k \text{ and } q_t \coloneqq \alpha_t^2 \|g_t - m_t\|^2.$$
 (1)

To handle unknown domain size and unbounded domains, U-DoG follows DoG in using  $\bar{r}_t$  as the step size numerator in lieu of D. Thus, the U-DoG step size admits the general form

$$\eta_{x,t} = \frac{\bar{r}_t}{\sqrt{G_{x,t}}} \text{ and } \eta_{y,t} = \frac{\bar{r}_t}{\sqrt{G_{y,t}}}, \text{ where } G_{x,0} \le G_{y,0} \le G_{x,1} \le \cdots.$$
(2)

<sup>2.</sup> Our results hold in the non-Lipschitz or non-smooth cases by setting  $L=\infty$  or  $\beta=\infty$ , respectively. In the non-smooth case, we define  $\nabla f(x) := \mathbb{E}\mathcal{G}(x)$  and assume it is a subgradient of f.

In the appendix, we also use the notation

$$\tilde{\eta}_{x,t} = \frac{1}{\sqrt{G_{x,t}}} \text{ and } \tilde{\eta}_{y,t} = \frac{1}{\sqrt{G_{y,t}}}.$$
 (3)

For bounded domains, setting  $G_{x,t} = G_{y,t} = 1 + Q_{t-1}$  recovers the UNIXGRAD guarantees up to logarithmic factors. However, for unbounded domains, ensuring the stability of U-DoG (i.e., that  $\bar{r}_t$  never grows much larger than  $\|x_0 - x_\star\|$ ) requires more careful selection of  $G_{x,t}, G_{y,t}$ . Enforcing iterate stability without compromising the rate of convergence is the main challenge we overcome. To that end, we define a few frequently appearing quantities:

$$\begin{split} r_t \coloneqq \max\{\|y_k - x_0\|, \|x_k - x_0\|\} \;,\; d_t \coloneqq \|y_t - x_\star\| \;,\; \bar{d}_t \coloneqq \max_{k \le t} d_k \;, \\ M_t \coloneqq \max_{k \le t} \left\{\alpha_k^2 \|m_k\|^2\right\} \text{ and } \theta_{t,\delta} \coloneqq \log \frac{60 \log(6t)}{\delta}. \end{split}$$

UNIXGRAD as a special case. For a domain with Euclidean diameter D, setting  $r_{\epsilon}=D\sqrt{2}$  and  $G_{x,t}=G_{y,t}=1+Q_{t-1}$  recovers UNIXGRAD (with Euclidean distance generating function) exactly, as it implies  $\bar{r}_t=D\sqrt{2}$  for all t and hence  $\alpha_t=t+1$ .

# 3. Analysis in the noiseless case

We begin our analysis under the simplifying assumption that gradients are computed exactly.

**Assumption 2** In addition to Assumption 1, we assume that  $\mathcal{G}(x) = \nabla f(x)$  with probability 1.

This noiseless setting allows us to isolate and address the key challenges of exploiting smoothness and stabilizing the iterates.

## 3.1. General suboptimally bound

Our first result is a bound on the suboptimality of U-DoG for general step sizes; see Section A.1 for complete proof. To interpret Proposition 3 recall that  $d_0$  is the initial distance to the optimum and the definition of  $Q_t$  given in (1).

**Proposition 3** In the noiseless setting (Assumption 2), suppose the U-DoG step sizes (2) satisfy  $G_{x,t} \ge Q_{t-1}$  for all  $t \ge 0$ . Then for every  $t \ge 0$  and for any number  $s \ge 0$ , we have

$$f(\hat{x}_t) - f(x_\star) \le O\left(\frac{s^{3/2}\beta(\bar{r}_{t+1} + d_0)^2 + (\bar{r}_{t+1} + d_0)\left[\sqrt{\max\{G_{y,t}, Q_t\}} - s\sqrt{Q_t}\right]_+}{\left(\sum_{k=0}^t \bar{r}_k/\bar{r}_{t+1}\right)^2}\right). \tag{4}$$

Before sketching the proof of Proposition 3, let us explain how it yields the desired rates of convergence if we momentarily set aside iterate stability and assume  $\bar{r}_t \leq D$  for all t, e.g., because the domain has diameter D. In this case, we may choose  $G_{x,t} = G_{y,t} = Q_{t-1}$  similarly to UNIXGRAD. Substituting s=1 in eq. (4) guarantees suboptimality  $O\left(\frac{\beta D^2}{\left(\sum_{k=0}^t \bar{r}_k/\bar{r}_{t+1}\right)^2}\right)$ . As shown in [27, Lemma 3], we have  $\max_{t < T} \sum_{k=0}^t \bar{r}_k/\bar{r}_{t+1} = \Omega\left(T\log^{-1}(\bar{r}_T/r_\epsilon)\right)$ , meaning that for some t < T we obtain the near-optimal rate  $O\left(\frac{\beta D^2}{T^2}\log^2\frac{D}{r_\epsilon}\right)$ . Moreover, since  $\alpha_t \leq t+1$  for all t,

when all gradients are bounded by L we have  $Q_t = O(L^2 \sum_{k \le t} \alpha_k^2) = O(L^2 t^3)$ . Substituting s = 0 in eq. (4) and reusing our bound on the denominator gives the near-optimal rate  $O\left(\frac{LD}{\sqrt{T}}\log^2\frac{D}{r_\epsilon}\right)$  in the non-smooth setting. We also see that setting  $r_\epsilon = \Omega(D)$  recovers the UNIXGRAD guarantees in the noiseless setting, which is to be expected since  $r_\epsilon = D\sqrt{2}$  recovers UNIXGRAD itself as explained in the previous section.

Our proof of Proposition 3 combines ideas from the analyses of UNIXGRAD and DOG. It centers on the weighted "regret"  $\mathcal{R}_t := \sum_{k=0}^t \omega_k \langle g_k, x_{k+1} - x_\star \rangle$  where  $\omega_k = \alpha_k \bar{r}_k$ . This is similar to the weighted regret considered for UNIXGRAD with additional weighting by  $\bar{r}_t$  used in the DOG analysis. Algebraic manipulation of  $\mathcal{R}_t$  gives (recall that  $d_t = \|y_t - x_\star\|$ ),

$$\mathcal{R}_t \le O\left(\bar{r}_{t+1}^2 \sqrt{Q_t} + \sum_{k=0}^t (d_k^2 - d_{k+1}^2) \sqrt{G_{y,k}} - \sum_{k=0}^t ||x_{k+1} - y_k||^2 \sqrt{Q_k}\right).$$

We use a telescoping argument from DoG in order to bound  $\sum_{k=0}^{t} (d_k^2 - d_{k+1}^2) \sqrt{G_{y,k}}$  by  $O\left(\bar{r}_{t+1}(\bar{r}_{t+1} + d_0)\sqrt{G_{y,t}}\right)$ . Next, following UNIXGRAD we leverage smoothness to write

$$\|x_{k+1} - y_k\|^2 = \left(\frac{\sum_{i=0}^k \omega_i}{\omega_k}\right)^2 \|\hat{x}_k - \hat{z}_k\|^2 \overset{\text{Lem. 25}}{\geq} \frac{\alpha_k^2}{4} \|\hat{x}_k - \hat{z}_k\|^2 \ge \frac{\alpha_k^2}{4\beta^2} \|\nabla f(\hat{x}_k) - \nabla f(\hat{z}_k)\|^2 = \frac{q_k^2}{4\beta^2},$$

where the last equality is the first time we assumed exact  $\mathcal{G}$ . We then show that, for all  $S \geq 0$ ,

$$\sum_{k=0}^{t} ||x_{k+1} - y_k||^2 \sqrt{Q_k} \ge \sum_{k=0}^{t} \frac{q_k^2}{\beta^2} \sqrt{Q_k} \ge \Omega \left( S \sqrt{Q_t} - S^{3/2} \beta \right); \tag{5}$$

this is a streamlined version of key arguments in [30, 36] where the authors carefully split the sum above based on the value of the adaptive step size. Taking  $S = s \cdot \bar{r}_{t+1}(\bar{r}_{t+1} + d_0)$  and substituting back, we get

$$\mathcal{R}_{t} \leq O\left(s^{3/2}\beta\bar{r}_{t+1}(\bar{r}_{t+1}+d_{0})^{2} + \bar{r}_{t+1}(\bar{r}_{t+1}+d_{0})\left[\sqrt{\max\{G_{y,t},Q_{t}\}} - s\sqrt{Q_{t}}\right]_{+}\right).$$
 (6)

To conclude the proof, we use the following UNIXGRAD "anytime online-to-batch conversion" [14] bound:

$$f(\hat{x}_t) - f(x_\star) \le \sum_{k=0}^t \frac{\omega_k}{\sum_{i=0}^t \omega_i} \langle \nabla f(\hat{x}_k), x_{k+1} - x_\star \rangle = \frac{\mathcal{R}_t}{\sum_{k=0}^t \omega_k},\tag{7}$$

where the last equality is the second and final time the proof uses the noiseless gradient assumption. Dividing eq. (6) by

$$\sum_{k=0}^{t} \omega_k \stackrel{\text{lem. 25}}{\geq} \frac{1}{2} \bar{r}_t \alpha_t^2 = \frac{1}{2} \bar{r}_t \left( \sum_{k=0}^{t} \bar{r}_k / \bar{r}_t \right)^2 \geq \frac{1}{2} \bar{r}_{t+1} \left( \sum_{k=0}^{t} \bar{r}_k / \bar{r}_{t+1} \right)^2, \tag{8}$$

and employing (7) yields the suboptimality bound (4).

## 3.2. Iterate stability

In the discussion following Proposition 3 above, we provisionally imagined that the iterates were bounded ( $\bar{r}_t \leq D$  for all t) and argued that in this case simply setting  $G_{x,t} = Q_{t-1}$  and  $G_{y,t} = Q_t$  suffices for obtaining optimal rates whenever  $D = O(d_0)$ . However, in unconstrained settings this choice of step size is hopeless, as it makes  $\eta_{x,0}$  infinite, implying divergence at the first step!<sup>3</sup>

In the following proposition, we identify two conditions that together guarantee the iterates remain appropriately bounded. The complete proof appears in Section A.2.

**Proposition 4** In the noiseless setting (Assumption 2), let s > 0 and define  $c_t = 12 \log_+^2 \left(\frac{s + Q_t}{s}\right)$ . If  $r_{\epsilon} \leq d_0$  and the U-DoG step sizes (2) satisfy (i)  $G_{y,t} \geq c_t^2 (s + Q_t)$  (with  $G_{x,0} \geq 144s$ ), and (ii)  $\max\{\|x_{t+1} - y_t\|, \|y_{t+1} - x_{t+1}\|\} \leq \frac{2\bar{r}_t}{c_{t-1}}$  for all  $t \geq 0$ , then we have

$$\bar{d}_t \leq 2d_0$$
 and  $\bar{r}_t \leq 4d_0$  for all  $t \geq 0$ .

Let us briefly explain the two requirements in Proposition 4. Requirement (i) folds two conditions into one. The first is that we increase the UNIXGRAD denominator by a logarithmic factor—this is analogous to the step size attenuation necessary to ensure the stability of DoG (i.e., the T-DoG step size [27, Section 3.3]). The second is more subtle, requiring that  $G_{y,t}$  upper bound  $Q_t$  (rather than  $Q_{t-1}$  as in UNIXGRAD and Proposition 3) and hence depend on  $||g_t - m_t||$ . This is essential for guaranteeing stability but is also the cause of considerable technical difficulty in the noisy setting. Requirement (ii) simply asks that U-DoG iterates at time t move by no more than a fraction of the estimated distance to optimality  $\bar{r}_t$ ; a reasonable requirement if the estimate is good.

The proof of Proposition 4 is a careful application of the T-DoG stability proof [27, Proposition 2] to the U-DoG template. The key to the proof is the following modification of the UNIX-GRAD online-to-batch conversion bound (7), which states that for any optimum  $x_*$  we have

$$\mathcal{R}'_{t} := \sum_{k=0}^{t} \eta_{y,k} \alpha_{k} \left\langle g_{k}, x_{k+1} - x_{\star} \right\rangle \stackrel{(\star)}{=} \sum_{k=0}^{t} \eta_{y,k} \alpha_{k} \left\langle \nabla f(\hat{x}_{k}), x_{k+1} - x_{\star} \right\rangle \ge 0, \tag{9}$$

where  $(\star)$  holds only in the noiseless setting. We algebraically manipulate  $\mathcal{R}'_t$  similarly to the weighted regret in the proof of Proposition 3. Writing  $Q'_t = c_{t-1}^2(s+Q_t)$ , we obtain

$$0 \le \mathcal{R}'_t \le \sum_{k=0}^t \left( d_k^2 - d_{k+1}^2 + \frac{q_k \bar{r}_k^2}{\sqrt{G_{y,k} Q'_k}} + \frac{\sqrt{Q'_k} - \sqrt{G_{x,k}}}{\sqrt{G_{y,k}}} (\|x_{k+1} - y_k\|^2 + \|x_{k+1} - y_{k+1}\|^2) \right).$$

Our requirements  $G_{y,k} \geq c_t^2(s+Q_t)$  (which entails  $G_{x,k} \geq G_{y,k-1} \geq c_{t-1}^2(s+Q_{t-1})$ ) and  $\|x_{k+1}-y_k\|^2+\|x_{k+1}-y_{k+1}\|^2 \leq \frac{8\bar{r}_k^2}{c_{k-1}^2}$ , allow us, with some more algebra, to bound the last two summands by  $\frac{9q_k\bar{r}_k^2}{c_k(s+Q_k)}$ . From here, the proof proceeds identically to the T-DoG analysis [27, Section 3.3]: we get that  $\sum_{k=0}^t \frac{9q_k^2\bar{r}_k}{c_k(s+Q_k)} \leq \frac{\bar{r}_t^2}{16}$  by the choice of  $c_t$ , and substituting back obtain that  $d_{t+1}^2 \leq d_0^2 + \frac{\bar{r}_t^2}{16}$ , which by straightforward induction implies the desired bounds on  $\bar{d}_t$  and  $\bar{r}_t$ .

<sup>3.</sup> For constrained domains, however, this choice results in a valid scheme where the first step jumps to the domain boundary. Indeed, UNIXGRAD also behaves this way for sufficiently scaled-up instances since it uses a fixed, arbitrary value for  $\eta_{x,0}$ . This underscores UNIXGRAD's strong reliance on the bounded domain assumption.

## 3.3. Rate of convergence in the noiseless case

With the conditional stability guarantee of Proposition 8 in place, we are ready to face a central challenge: finding step sizes  $\eta_{x,t}, \eta_{y,t}$  that satisfy the proposition's conditions but still lead to good rates of convergence in the smooth case. Our solution is (recalling the notation  $M_t$  $\max_{k < t} \{ \alpha_k^2 || m_k ||^2 \} )$ :

$$\eta_{x,t} = \frac{\bar{r}_t}{12 \log_+^2 \left(\frac{\|m_0\|^2 + Q_{t-1}}{\|m_0\|^2}\right) \sqrt{\max\{\|m_0\|^2 + Q_{t-1}, M_t\}}}$$

$$\eta_{y,t} = \frac{\bar{r}_t}{12 \log_+^2 \left(\frac{\|m_0\|^2 + Q_t}{\|m_0\|^2}\right) \sqrt{\max\{\|m_0\|^2 + Q_t, M_t\}}}.$$
(10)

Clearly, the step sizes (10) satisfy the first condition in Proposition 4 with  $s=\|m_0\|^2$ . To see why the second condition holds, note that, since  $\sqrt{M_t} \ge \alpha_t ||m_t||$ , we have  $\eta_{x,t} \le \frac{\bar{r}_t}{c_t \alpha_t ||m_t||}$ . By the contractive property of projections, we therefore have

$$||x_{t+1} - y_t|| \le \eta_{x,t} \alpha_t ||m_t|| \le \frac{\bar{r}_t}{c_t} \le \frac{2\bar{r}_t}{c_t}.$$

A similar argument also shows that  $||x_{t+1} - y_{t+1}|| \leq \frac{2\bar{r}_t}{c_t}$ , fulfilling the conditions of Proposition 4 (see Lemma 19).

Now the question becomes: how does the introduction of  $M_t$  into the step size affect suboptimality? In the non-smooth case the effect is minimal, as we anyway bound  $Q_t$  with  $O(L^2t^3)$ , and  $M_t = O(L^2 t^2)$  is of a lower order. In the smooth case, however,  $M_t$  is potentially more harmful, since while Proposition 3 allows us to cancel the dependence on  $Q_t$  by setting  $s = c_t$ , it leaves  $M_t$ hanging in the numerator, yielding  $f(\hat{x}_t) - f(x_\star) \leq O\left(\frac{1}{\alpha_t^2}\left(c_t^{3/2}\beta d_0^2 + c_t d_0\sqrt{M_t}\right)\right)$ . Fortunately, smoothness allows us to relate  $M_t$  back to the optimality gap  $f(\hat{x}_t) - f(x_\star)$ . In

particular, in the unconstrained setting  $\mathcal{K} = \mathbb{R}^n$  we have

$$||m_t||^2 \le 2||g_t - m_t||^2 + 2||g_t||^2 \le 2Q_t/\alpha_t^2 + 4\beta[f(\hat{x}_t) - f(x_\star)],$$

where the last transition used that  $g_t = \nabla f(\hat{x}_t)$  in the noiseless setting. Combining this bound with Proposition 3, we obtain

$$f(\hat{x}_t) - f(x_\star) \le O\left(\frac{c_t^{3/2}\beta d_0^2 + \sqrt{c_t^2\beta d_0^2 \max_{k \le t} \alpha_k^2 [f(\hat{x}_k) - f(x_\star)]}}{\alpha_t^2}\right),$$

from which  $f(\hat{x}_t) - f(x_\star) \le O\left(\frac{c_t^2 \beta d_0^2}{\alpha_t^2}\right)$  follows by induction. Thus we arrive at our final guarantee in the noiseless case: Theorem 5 (see full proof in Section A.3).

**Theorem 5** In the noiseless setting (Assumption 2) with  $K = \mathbb{R}^n$  and  $r_{\epsilon} \leq d_0$ , using the step sizes eq. (10), we get that  $\bar{d}_T \leq 2d_0$ ,  $\bar{r}_T \leq 4d_0$  and, for  $\tau = \arg\max_{t < T} \sum_{i \leq t} \frac{\bar{r}_i}{\bar{r}_{t+1}}$ , the suboptimality is

$$f(\hat{x}_{\tau}) - f(x_{\star}) \le O\left(c_{r_{\epsilon},T} \min\left\{\frac{\beta d_0^2}{T^2}, \frac{Ld_0}{\sqrt{T}}\right\}\right),$$

where 
$$c_{r_{\epsilon},T} = \log_{+}^{4} \left( 1 + \frac{T \min\{\beta d_{0}^{2}, Ld_{0}\}}{f(x_{0}) - f(x_{\star})} \right) \log_{+}^{2} \left( \frac{d_{0}}{r_{\epsilon}} \right)$$
.

# 4. Analysis in the stochastic case

In this section, we extend the U-DoG guarantees to the noisy case. We start by assuming that the gradient noise is bounded, a setting that captures most of the remaining technical challenges. We then generalize our results to sub-Gaussian noise by means of a black-box reduction [3]. Finally, we specialize the U-DoG guarantee for mini-batches of bounded gradient estimates and conclude with a discussion of the (weak) dependence of our result on problem parameter bounds. Throughout this section, we denote the empirical variance at time t by

$$V_t := \frac{1}{t+1} \sum_{k=0}^{t} (\|g_t - \nabla f(\hat{x}_t)\|^2 + \|m_t - \nabla f(\hat{z}_t)\|^2).$$
 (11)

We also recall the notation

$$\theta_{t,\delta} \coloneqq \log \frac{60 \log(6t)}{\delta}.$$

# 4.1. Analysis with bounded noise

We formalize the bounded noise assumption as follows.

**Assumption 6** In addition to Assumption 1, we assume that  $\|\mathcal{G}(x) - \nabla f(x)\| \leq \mathfrak{b}(x)$  with probability 1 for all  $x \in \mathcal{K}$ , for some (known<sup>4</sup>) function  $\mathfrak{b} : \mathcal{K} \to \mathbb{R}_+$ .

For the iterates of U-DoG we define

$$\mathfrak{b}_t \coloneqq \mathfrak{b}(\hat{x}_t) \text{ and } \bar{\mathfrak{b}}_t \coloneqq \max \Big\{ \max_{i \le t} \mathfrak{b}_i, \mathfrak{b}(\hat{z}_0) \Big\}.$$
 (12)

With the assumption and notation in place, we state the stochastic equivalent of Proposition 3 in the following (see proof in Section B.1).

**Proposition 7** In the bounded noise setting (Assumption 6), suppose the U-DoG step sizes (2) satisfy  $G_{x,t} \ge Q_{t-1}$  for every  $t \ge 0$ . Then for any  $\mathfrak{B} > 0$ ,  $T \in \mathbb{N}$ , and  $\delta \in (0,1)$ , with probability at least  $1 - \delta - \mathbb{P}[\bar{\mathfrak{b}}_{T-1} > \mathfrak{B}]$  we have, for all t < T and  $s \ge 0$ ,

$$f(\hat{x}_t) - f(x_\star) \le O\left(RHS_{eq. (4)} + \frac{(1+s)(\bar{r}_{t+1} + d_0)\sqrt{t^3\theta_{t+1,\delta}V_t + (t\theta_{t+1,\delta}\mathfrak{B})^2}}{\left(\sum_{k=0}^t \bar{r}_k/\bar{r}_{t+1}\right)^2}\right)$$

where RHS<sub>eq. (4)</sub> = 
$$\frac{s^{3/2}\beta(\bar{r}_{t+1}+d_0)^2 + (\bar{r}_{t+1}+d_0)\left[\sqrt{\max\{G_{y,t},Q_t\}} - s\sqrt{Q_t}\right]_+}{\left(\sum_{k=0}^t \bar{r}_k/\bar{r}_{t+1}\right)^2} \text{ as in Proposition 3.}$$

Proposition 7 is a fairly straightforward extension of its noiseless counterpart. The bound (5) continues to hold if we replace  $Q_t$  with  $\hat{Q}_t = \sum_{k=0}^t \alpha_k^2 \min\{\|g_k - m_k\|^2, \|\nabla f(\hat{x}_k) - \nabla f(\hat{z}_k)\|^2\}$ . Proceeding as in the proof of Proposition 3, we conclude that

$$f(\hat{x}_t) - f(x_\star) \le O\left(\text{RHS}_{eq. (4)} + \frac{s(\bar{r}_{t+1} + d_0)(Q_t^{1/2} - 2\hat{Q}_t^{1/2}) + \sum_{k=0}^t \omega_k \langle \nabla f(\hat{x}_k) - g_k, x_{k+1} - x_\star \rangle}{\left(\sum_{k=0}^t \bar{r}_k / \bar{r}_{t+1}\right)^2}\right).$$

<sup>4.</sup> We may view b as a coarse upper bound on the true noise magnitude, as it only affects low order terms in our bounds.

We show that  $Q_t^{1/2} - 2\hat{Q}_t^{1/2} \leq O\left(\sqrt{t^3V_t}\right)$  by straightforward manipulation. Furthermore, using a time-uniform empirical-Bernstein-type concentration bound [26, 27] (Lemma 21) to show that (with the appropriate high probability) the martingale difference sum  $\sum_{k=0}^t \omega_k \left\langle \nabla f(\hat{x}_k) - g_k, x_{k+1} - x_\star \right\rangle$ is bounded by  $O\left(\bar{r}_t \bar{d}_{t+1} \sqrt{t^3 \theta_{t+1,\delta} V_t + (t \theta_{t+1,\delta} \mathfrak{B})^2}\right)$ . Next, we extend our iterate stability guarantee to the stochastic setting (see proof in section B.2).

**Proposition 8** In the bounded noise setting Assumption 6, let s > 0,  $T \in \mathbb{N}$  and  $\delta \in (0,1)$ , and define  $c_t = 400\theta_{T,\delta} \log_+^2 \left(\frac{s+Q_t}{s}\right)$ . Suppose that  $r_{\epsilon} \leq d_0$  and the U-DoG step sizes (2) satisfy, with probability 1, for all  $t \geq 0$ : (i)  $G_{y,t} \geq c_t^2(s+Q_t)$  (with  $G_{x,0} \geq 400^2 \theta_{T,\delta}^2 s$ ), (ii)  $\max\{\|x_{t+1} - x_t\|_{t=0}^2 s \leq t \leq t \}$  $|y_t||, ||y_{t+1} - x_{t+1}|| \le \frac{2\bar{r}_t}{c_{t-1}}, (iii) \sqrt{G_{y,t}} \ge c_t \alpha_t \max\{\|\nabla f(\hat{x}_t) - g_t\|, \|\nabla f(\hat{x}_t) - m_t\|\}, and (iv)$  $\eta_{y,t}$  is independent of  $g_t$  given  $x_0, \ldots, x_t$ . Then, we have with probability of at least  $1 - \delta$ ,

$$\bar{d}_t \leq 2d_0$$
 and  $\bar{r}_t \leq 4d_0$  for all  $t < T$ .

Conditions (i) and (ii) of Proposition 8 are identical to their noiseless counterparts in Proposition 4, while conditions (iii) and (iv) are new and facilitate the application of a concentration bound to the weighed regret  $\mathcal{R}'_t$  defined in eq. (9). In particular, the condition (iv) ensures that  $\sum_{k=0}^{t} \eta_{y,k} \alpha_k \langle g_k - \nabla f(\hat{x}_k), x_{k+1} - x_{\star} \rangle$  is a martingale difference sequence, and condition (iii)guarantees boundedness required by our concentration bound (Lemma 22). With this high probability bound in place, the proof continues in the same vein as the noiseless case.

When searching for step sizes meeting the conditions of Proposition 8 we encounter two challenges. First, condition (iii) asks  $G_{y,t}$  to be large compared to a quantity depending on the exact gradient  $\nabla f(\hat{x}_t)$ , which we cannot access directly. We solve it using the bounds given in (12). Simply adding  $c_t^2(t+1)^2\bar{\mathfrak{b}}_t^2 \geq c_t^2\alpha_t^2\mathfrak{b}_t^2$  to  $G_{y,t}$  guarantees that  $\sqrt{G_{y,t}} \geq c_t\alpha_t\|\nabla f(\hat{x}_t) - g_t\|$ . Moreover, using  $\|u\|^2 + \|v\|^2 \geq \frac{1}{2}\|v + u\|^2$ , we have

$$||g_t - m_t||^2 + \bar{\mathfrak{b}}_t^2 \ge ||g_t - m_t||^2 + ||\nabla f(\hat{x}_t) - g_t||^2 \ge \frac{1}{2} ||\nabla f(\hat{x}_t) - m_t||^2.$$

Therefore, taking  $G_{y,t} = c_t^2(s + 2Q_t + 2(t+1)^2\bar{\mathfrak{b}}_t^2)$  fulfills condition (iii). However, it violates condition (iv) which leads us to the second challenge: how to avoid dependence on  $q_t$ ? To address this challenge, we employ the somewhat unusual trick of drawing a fresh stochastic gradient  $\tilde{q}_t \sim$  $\mathcal{G}(\hat{x}_t)$  which is, by construction, independent of  $g_t$  given  $\hat{x}_t$ . We can now replace the forbidden  $||g_t - m_t||$  with the valid upper bound  $2||\tilde{g}_t - m_t|| + 8\bar{\mathfrak{b}}_t$  and thus satisfy conditions (i) and (iii) without violating condition (iv).

To satisfy condition (ii) we introduce  $M_t$  to  $G_{y,t}$  as done in the noiseless setting and make another modification to ensure the monotonicity required in (2). Writing,

$$\tilde{q}_t \coloneqq 2\alpha_t^2 \|\tilde{g}_t - m_t\|^2 \ , \ \bar{Q}_t \coloneqq \sum_{k=0}^t \max\{q_k, \tilde{q}_k\} \ \text{and} \ p_t \coloneqq 8(t+1)^2 \bar{\mathfrak{b}}_t^2,$$
 (13)

our final step sizes are:

$$\eta_{x,t} = \frac{\bar{r}_t}{400\theta_{T,\delta} \log_+^2 \left(1 + \frac{p_{t-1} + \bar{Q}_{t-1}}{\|m_0\|^2 + p_0}\right) \sqrt{\max\{\|m_0\|^2 + p_0 + p_{t-1} + \bar{Q}_{t-1}, M_t\}}}{\bar{r}_t}$$

$$\eta_{y,t} = \frac{\bar{r}_t}{400\theta_{T,\delta} \log_+^2 \left(1 + \frac{p_t + \tilde{q}_t + \bar{Q}_{t-1}}{\|m_0\|^2 + p_0}\right) \sqrt{\max\{\|m_0\|^2 + p_0 + p_t + \tilde{q}_t + \bar{Q}_{t-1}, M_t\}}}.$$
(14)

Similar to the T-DoG step sizes [27, Section 3.3], our step sizes depend logarithmically on the desired confidence level  $\delta$  and double-logarithmically on the maximum iteration budget T.

With all the pieces in place, we now state our main result (see proof in Section B.3).

**Theorem 9** In the bounded noise setting (Assumption 6) with  $K = \mathbb{R}^n$ , for any  $T \in \mathbb{N}$  and  $\delta \in (0, \frac{1}{5})$ , consider U-DoG with step sizes (14). With probability at least  $1-5\delta$ , we have  $\bar{d}_T \leq 2d_0$ ,  $\bar{r}_T \leq 4d_0$  and for  $\tau = \arg\max_{t < T} \sum_{i \leq t} \frac{\bar{r}_i}{\bar{r}_{t+1}}$  and  $\mathfrak{b}_{\star} := \max_{x: ||x-x_{\star}|| \leq 2d_0} \{\mathfrak{b}(x)\}$  we have

$$f(\hat{x}_{\tau}) - f(x_{\star}) \le O\left(c_{\delta, r_{\epsilon}, T}\left(\min\left\{\frac{\beta d_0^2}{T^2}, \frac{Ld_0}{\sqrt{T}}\right\} + \frac{d_0\sqrt{V_T}}{\sqrt{T}} + \frac{d_0\mathfrak{b}_{\star}}{T}\right)\right),\tag{15}$$

where 
$$c_{\delta,r_{\epsilon},T} = \log^2\left(\frac{\log_+(T)}{\delta}\right)\log_+^4\left(1 + T\frac{\mathfrak{b}_{\star} + \min\left\{\beta d_0^2, L d_0\right\}}{f(x_0) - f(x_{\star})}\right)\log_+^2\left(\frac{d_0}{r_{\epsilon}}\right)$$
 and  $V_t$ , defined in eq. (11), is the empirical noise variance.

We remark that under our assumptions it is straightforward to replace the empirical variance  $V_t$  in eq. (15) with its expectation without altering other non-logarithmic terms in the bound, e.g., via Hoeffding's inequality.

#### 4.2. From bounded to sub-Gaussian noise

The bounded noise assumption makes analysis convenient but is not entirely satisfactory since averaging independent bounded-noise estimators does not reduce the probability 1 noise bound, preventing us from making statements about mini-batch scaling. To address this issue, we consider the following standard assumption.

**Assumption 10** In addition to Assumption 1, we assume that  $\|\mathcal{G}(x) - \nabla f(x)\|$  is  $\sigma^2(x)$ -sub-Gaussian for all  $x \in \mathcal{K}$ , for some (known)  $\sigma : \mathcal{K} \to \mathbb{R}_+$ . That is,

$$\mathbb{P}(\|\mathcal{G}(x) - \nabla f(x)\| \ge z) \le 2\exp(-z^2/\sigma^2(x))$$

for all  $z \geq 0$  and  $x \in \mathcal{K}$ .

To move from bounded to sub-Gaussian we utilize a reduction due to Attia and Koren [3] that allows us to essentially replace  $\mathfrak{b}(\cdot)$  with  $\sigma(\cdot)$  in Theorem 9 at the cost of additional logarithmic factors. To that end, we define  $\bar{\sigma}_t \coloneqq \max\{\max_{i \le t} \sigma(\hat{x}_k), \sigma(\hat{z}_0)\}$ , as well as  $\sigma_\star \coloneqq \max_{x:\|x-x_\star\| \le 2d_0} \sigma(x)$  and  $\varsigma_{t,\delta} \coloneqq 3\log^{1/2}(\frac{15(t+1)^2}{\delta})$ . With this notation in hand, we state our guarantee for the sub-Gaussian setting (see proof in Section B.4).

**Corollary 11** Consider the sub-Gaussian noise setting (Assumption 10) with  $K = \mathbb{R}^n$  and  $\delta \in (0, \frac{1}{6})$ , using the step sizes (14) with  $\bar{\mathfrak{b}}_t = \bar{\sigma}_t \varsigma_{t,\delta}$ , then with probability at least  $1 - 6\delta$  we get that  $\bar{d}_T \leq 2d_0$ ,  $\bar{r}_T \leq 4d_0$ , and the suboptimality bound (15) holds for  $\mathfrak{b}_\star = \sigma_\star \varsigma_{T-1,\delta}$ .

# 4.3. Corollary: mini-batch of bounded noise

Finally, we leverage our result for sub-Gaussian noise to demonstrate that U-DoG automatically benefits from increasing mini-batch size (see proof in Section B.5).

**Assumption 12** In addition to Assumption 1, we assume that G(x) is the average of B unbiased estimates of  $\nabla f(x)$ , each bounded by L with a known upper bound  $\hat{L} \geq L$ .

**Corollary 13** In the mini-batch setting (Assumption 12) with  $\mathcal{K} = \mathbb{R}^n$ , for any  $T \in \mathbb{N}$  and  $\delta \in (0, \frac{1}{6})$ , consider U-DoG with step sizes (14) where  $\bar{\mathfrak{b}}_t = \sqrt{2} \frac{\hat{L}}{\sqrt{B}} \varsigma_{t,\delta}$ . With probability at least  $1 - 6\delta$  we have  $\bar{d}_T \leq 2d_0$ ,  $\bar{r}_T \leq 4d_0$  and, for  $\tau = \arg\max_{t < T} \sum_{i < t} \frac{\bar{r}_i}{\bar{r}_{t+1}}$ ,

$$f(\hat{x}_{\tau}) - f(x_{\star}) \le O\left(c_{\delta, r_{\epsilon}, T} \left(\frac{\beta d_0^2}{T^2} + \frac{\left(L + \hat{L}/\sqrt{T}\right) d_0}{\sqrt{TB}}\right)\right),$$

where 
$$c_{\delta,r_{\epsilon},T} = \sqrt{\log_{+}(\frac{T}{\delta})} \log^{2}(\frac{\log_{+}(T)}{\delta}) \log_{+}^{4}(1 + T\frac{\hat{L} + \min\{\beta d_{0}^{2}, Ld_{0}\}}{f(x_{0}) - f(x_{\star})}) \log_{+}^{2}(\frac{d_{0}}{r_{\epsilon}}).$$

## 4.4. Discussion: how parameter-free is our algorithm?

With our results established, we now discuss in more detail the extent to which our algorithms and complexity bounds are free of a-priori knowledge of problem parameters. U-DoG requires a lower bound  $r_{\epsilon}$  on the initial distance to the optimum  $d_0$ , and pointwise upper bounds  $\mathfrak{b}_t$  on the noise magnitude at each iteration. Theorem 9 provides suboptimality bounds that depend poly-logarithmically on  $\frac{d_0}{r_{\epsilon}}$  which quantifies how  $r_{\epsilon}$  underestimates  $d_0$ . Many works [e.g., 7, 9, 15, 16, 27, 28, 39, 41, 49] treat such logarithmic dependence as the *definition* of a parameter-free algorithm, and in that strict sense our method is certainly parameter-free. The noise bounds impact our suboptimality guarantees polynomially via the additive term  $\mathfrak{b}_{\star}/T$  where  $\mathfrak{b}_{\star} := \max_{x:\|x-x_{\star}\| \leq 2d_0} \{\mathfrak{b}(x)\}$ , potentially implying greater sensitivity to problem parameters. Neverthless, we argue that our method fully deserve the title "parameter-free" for the following reasons.

- 1. The noise bound only contributes a low-order error term. To see why  $\mathfrak{b}_{\star}/T$  is low-order, let  $\hat{\mathfrak{b}}_T$  be the largest stochastic gradient error in the first T iterations of U-DoG. Then the empirical variance satisfies  $V_T = O\left(\hat{\mathfrak{b}}_T^2\right)$  and the noise-dependent part of Theorem 9 is  $\widetilde{O}\left(\frac{d_0\hat{\mathfrak{b}}_T}{\sqrt{T}} + \frac{d_0\mathfrak{b}_{\star}}{T}\right)$ . Therefore, as long as  $\mathfrak{b}_{\star}/\hat{\mathfrak{b}}_T = O(\sqrt{T})$  (i.e. the noise bound is loose by  $\sqrt{T}$  or less) we get the near-optimal dependence on the unknown, true noise magnitude  $\hat{\mathfrak{b}}_T$ .
- 2. The low-order term and noise bound assumption are unavoidable. Recent work [4, Theorem 6] proves that any algorithm with logarithmic dependence on uncertainty in distance to optimality must suffer the low-order error term b<sub>\*</sub>/T, and hence also require an a-priori noise bound (concurrent work [10, 31] also shows similar results). Moreover, prior parameter-free algorithms assume known bounds on stochastic gradient magnitude, which is stronger than assuming noise bounds. In this sense, our method is as parameter-free as it gets.
- 3. The noise bound is often easy to obtain and vanishes as batch size grows. Corollary 13 and Assumption 12 give a general setting where a noise bound is readily available. For a concrete instantiation, consider logistic regression with normalized covariates. In this case  $\hat{L} = 1$  and the noise bound at batch size B is O(1/B), which decreases as the batch size grows.

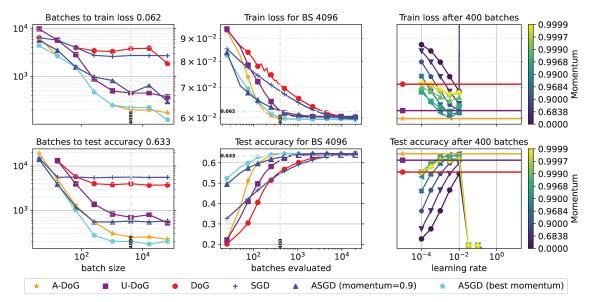


Figure 1: Training a linear model with ViT-32 features and least-squares loss on SVHN. Top: Train loss. Bottom: Test accuracy after iterate averaging. First column: Batch size scaling of complexity to reach target performance. Second column: Learning curves. Third column: ASGD performance at all learning rates and momenta, contrasted with DoG variants.

# 5. Experiments

We test U-DoG on a suite of experiments on convex and non-convex learning problems. We also heuristically derive and experiment with an algorithm we call A-DoG, which integrates ideas from ACCELEGRAD [36] and DoG. Namely, it uses the ACCELEGRAD step with DoG numerator and  $\alpha_t$  as in U-DoG. The pseudocode for A-DoG is given in Algorithm 2 in Section G.2.

We compare our algorithms to DoG as well as carefully tuned SGD with constant Nesterov momentum (ASGD for short) across a wide range of batch sizes. Detailed experimental results and analyses, as well as implementation details, are presented in Appendix G.

Our testbed consists of multiple classification problems based on the VTAB benchmark [67] and libsvm datasets [12], which we solve with both multiclass log loss and least squares loss, as well as a synthetic noiseless linear regression problem (see Section G.3). In addition, we perform preliminary experiments in the non-convex setting, including training neural networks from scratch on CIFAR-10 and VTAB datasets, and fine-tuning a CLIP model on ImageNet (see Section G.4).

On convex optimization problems, both U-DoG and A-DoG often substantially improve over DoG, with A-DoG achieving results comparable to well-tuned ASGD and outperforming U-DoG, likely by avoiding extra-gradient computations. Figure 1 illustrates these results on a particular dataset and least-squares loss function configuration and Section G.3 repeats this figure for additional configurations. The left panels in the figure show the rate of convergence of A-DoG, U-DoG and ASGD plateaus at a larger batch size compared to DoG and SGD without momentum. This is the typical effect of acceleration in stochastic optimization [56], and is also supported by Corollary 13 which shows that, for sufficiently large batch size, U-DoG converges at rate scaling as  $1/T^2$ . In contrast, non-accelerated methods like DoG and SGD converge with rate scaling as 1/T. The right panels of the figure show that, at a tight computational budget, the performance

of ASGD is very sensitive to the tuning of both step size and momentum, with only the very best values matching the performance of A-DoG. When using logarithmic instead of least-squares loss, the test accuracy becomes more robust to large step size choices (see Figure 2 in the appendix). This is partly because the log loss is Lispchitz which prevents complete divergence at any fixed step size.

In our preliminary non-convex experiments on neural network models (reported in detail in Sections G.3 and G.4), we find that U-DoG often fails to converge to competitive results, while A-DoG is competitive with DoG on most VTAB tasks, but under-performs it for CIFAR-10 and ImageNet fine-tuning, indicating that it is not a yet a viable general-purpose neural network optimizer.

# Acknowledgments

We thank Konstantin Mishchenko for helpful discussion. This work was supported by the NSF-BSF program, under NSF grant #2239527 and BSF grant #2022663. MI acknowledges support from the Israeli Council of Higher Education. OH acknowledges support from Pitt Momentum Funds, and AFOSR grant #FA955023-1-0242. YC acknowledges support from the Israeli Science Foundation (ISF) grant no. 2486/21 and the Alon Fellowship.

#### References

- [1] Martín Abadi, Ashish Agarwal, Paul Barham, Eugene Brevdo, Zhifeng Chen, Craig Citro, Greg S. Corrado, Andy Davis, Jeffrey Dean, Matthieu Devin, Sanjay Ghemawat, Ian Goodfellow, Andrew Harp, Geoffrey Irving, Michael Isard, Yangqing Jia, Rafal Jozefowicz, Lukasz Kaiser, Manjunath Kudlur, Josh Levenberg, Dandelion Mané, Rajat Monga, Sherry Moore, Derek Murray, Chris Olah, Mike Schuster, Jonathon Shlens, Benoit Steiner, Ilya Sutskever, Kunal Talwar, Paul Tucker, Vincent Vanhoucke, Vijay Vasudevan, Fernanda Viégas, Oriol Vinyals, Pete Warden, Martin Wattenberg, Martin Wicke, Yuan Yu, and Xiaoqiang Zheng. TensorFlow: Large-scale machine learning on heterogeneous systems, 2015.
- [2] E. Alpaydin and Fevzi. Alimoglu. Pen-Based Recognition of Handwritten Digits. UCI Machine Learning Repository, 1998. DOI: https://doi.org/10.24432/C5MG6K.
- [3] Amit Attia and Tomer Koren. SGD with AdaGrad stepsizes: Full adaptivity with high probability to unknown parameters, unbounded gradients and affine variance. In *International Conference on Machine Learning (ICML)*, 2023.
- [4] Amit Attia and Tomer Koren. How free is parameter-free stochastic optimization? In *International Conference on Machine Learning (ICML)*, 2024.
- [5] Charles Beattie, Joel Z Leibo, Denis Teplyashin, Tom Ward, Marcus Wainwright, Heinrich Küttler, Andrew Lefrancq, Simon Green, Víctor Valdés, Amir Sadik, et al. Deepmind lab. *arXiv:1612.03801*, 2016.
- [6] Amir Beck and Marc Teboulle. A fast iterative shrinkage-thresholding algorithm for linear inverse problems. *SIAM journal on imaging sciences*, 2(1):183–202, 2009.
- [7] Aditya Bhaskara, Ashok Cutkosky, Ravi Kumar, and Manish Purohit. Online learning with imperfect hints. In *International Conference on Machine Learning (ICML)*, 2020.
- [8] Jock Blackard. Covertype. UCI Machine Learning Repository, 1998. DOI: https://doi.org/10.24432/C50K5N.
- [9] Yair Carmon and Oliver Hinder. Making SGD parameter-free. In *Conference on Learning Theory (COLT)*, 2022.
- [10] Yair Carmon and Oliver Hinder. The price of adaptivity in stochastic convex optimization. In *Conference on Learning Theory (COLT)*, 2024.
- [11] Yair Carmon, Danielle Hausler, Arun Jambulapati, Yujia Jin, and Aaron Sidford. Optimal and adaptive monteiro-svaiter acceleration. In *Advances in Neural Information Processing Systems* (NeurIPS), 2022.
- [12] Chih-Chung Chang and Chih-Jen Lin. LIBSVM: a library for support vector machines. *ACM Transactions on Intelligent Systems and Technology*, 2011.
- [13] Gong Cheng, Junwei Han, and Xiaoqiang Lu. Remote sensing image scene classification: Benchmark and state of the art. *Proceedings of the IEEE*, 105(10):1865–1883, 2017.

- [14] Ashok Cutkosky. Anytime online-to-batch, optimism and acceleration. In *International Conference on Machine Learning (ICML)*, pages 1446–1454, 2019.
- [15] Ashok Cutkosky. Artificial constraints and hints for unbounded online learning. In *Conference on Learning Theory (COLT)*, 2019.
- [16] Ashok Cutkosky and Francesco Orabona. Black-box reductions for parameter-free online learning in Banach spaces. In *Conference on Learning Theory (COLT)*, 2018.
- [17] Aaron Defazio and Konstantin Mishchenko. Learning-rate-free learning by D-adaptation. In *International Conference on Machine Learning (ICML)*, 2023.
- [18] Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. ImageNet: A large-scale hierarchical image database. In Conference on Computer Vision and Pattern Recognition (CVPR), 2009.
- [19] Jelena Diakonikolas and Lorenzo Orecchia. Accelerated extra-gradient descent: A novel accelerated first-order method. In *Innovations in Theoretical Computer Science (ITCS)*, 2018.
- [20] Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, Jakob Uszkoreit, and Neil Houlsby. An image is worth 16x16 words: Transformers for image recognition at scale. In *International Conference on Learning Representations (ICLR)*, 2021.
- [21] John Duchi, Elad Hazan, and Yoram Singer. Adaptive subgradient methods for online learning and stochastic optimization. *Journal of Machine Learning Research*, 12(7), 2011.
- [22] Vineet Gupta, Tomer Koren, and Yoram Singer. A unified approach to adaptive regularization in online and stochastic optimization. *arXiv:1706.06569*, 2017.
- [23] Charles R. Harris, K. Jarrod Millman, Stéfan J. van der Walt, Ralf Gommers, Pauli Virtanen, David Cournapeau, Eric Wieser, Julian Taylor, Sebastian Berg, Nathaniel J. Smith, Robert Kern, Matti Picus, Stephan Hoyer, Marten H. van Kerkwijk, Matthew Brett, Allan Haldane, Jaime Fernández del Río, Mark Wiebe, Pearu Peterson, Pierre Gérard-Marchant, Kevin Sheppard, Tyler Reddy, Warren Weckesser, Hameer Abbasi, Christoph Gohlke, and Travis E. Oliphant. Array programming with NumPy. Nature, 585(7825):357–362, 2020.
- [24] Kaiming He, X. Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. *Conference on Computer Vision and Pattern Recognition (CVPR)*, 2016.
- [25] Steven R Howard, Aaditya Ramdas, Jon McAuliffe, and Jasjeet Sekhon. Time-uniform chernoff bounds via nonnegative supermartingales. *Probability Surveys*, 17:257–317, 2020.
- [26] Steven R Howard, Aaditya Ramdas, Jon McAuliffe, and Jasjeet Sekhon. Time-uniform, non-parametric, nonasymptotic confidence sequences. *The Annals of Statistics*, 49(2):1055–1080, 2021.
- [27] Maor Ivgi, Oliver Hinder, and Yair Carmon. DoG is SGD's best friend: A parameter-free dynamic step size schedule. In *International Conference on Machine Learning (ICML)*, 2023. We refer to the latest arXiv version: https://arxiv.org/abs/2302.12022.

- [28] Andrew Jacobsen and Ashok Cutkosky. Parameter-free mirror descent. In *Conference on Learning Theory (COLT)*, 2022.
- [29] Justin Johnson, Bharath Hariharan, Laurens Van Der Maaten, Li Fei-Fei, C Lawrence Zitnick, and Ross Girshick. CLEVR: A diagnostic dataset for compositional language and elementary visual reasoning. In Conference on Computer Vision and Pattern Recognition (CVPR), 2017.
- [30] Ali Kavis, Kfir Y Levy, Francis Bach, and Volkan Cevher. UniXGrad: A universal, adaptive algorithm with optimal guarantees for constrained optimization. *Advances in Neural Information Processing Systems (NeurIPS)*, 2019.
- [31] Ahmed Khaled and Chi Jin. Tuning-free stochastic optimization. In *International Conference on Machine Learning (ICML)*, 2024.
- [32] Ahmed Khaled, Konstantin Mishchenko, and Chi Jin. DoWG unleashed: An efficient universal parameter-free gradient descent method. In *Advances in Neural Information Processing Systems (NeurIPS)*, 2023.
- [33] Diederik P Kingma and Jimmy Ba. ADAM: A method for stochastic optimization. In *International Conference on Learning Representations (ICLR)*, 2015.
- [34] Alex Krizhevsky. Learning multiple layers of features from tiny images. Technical report, University of Toronto, 2009.
- [35] Guanghui Lan. An optimal method for stochastic composite optimization. *Mathematical Programming*, 133(1):365–397, 2012.
- [36] Kfir Y Levy, Alp Yurtsever, and Volkan Cevher. Online adaptive methods, universality and acceleration. *Advances in Neural Information Processing Systems (NeurIPS)*, 2018.
- [37] Ilya Loshchilov and Frank Hutter. SGDR: Stochastic gradient descent with warm restarts. *International Conference on Learning Representations*, 2017.
- [38] H Brendan McMahan. A survey of algorithms and analysis for adaptive online learning. *The Journal of Machine Learning Research*, 18(1):3117–3166, 2017.
- [39] H Brendan McMahan and Francesco Orabona. Unconstrained online linear learning in Hilbert spaces: Minimax algorithms and normal approximations. In *Conference on Learning Theory* (*COLT*), 2014.
- [40] H Brendan McMahan and Matthew Streeter. Adaptive bound optimization for online convex optimization. *arXiv:1002.4908*, 2010.
- [41] Zakaria Mhammedi and Wouter M Koolen. Lipschitz and comparator-norm adaptivity in online learning. In *Conference on Learning Theory (COLT)*, 2020.
- [42] Konstantin Mishchenko and Aaron Defazio. Prodigy: An expeditiously adaptive parameter-free learner. *arXiv*:2306.06101, 2023.

- [43] Arkadi Nemirovski. Prox-method with rate of convergence o(1/t) for variational inequalities with Lipschitz continuous monotone operators and smooth convex-concave saddle point problems. SIAM Journal on Optimization, 15(1):229–251, 2004.
- [44] Yurii Nesterov. A method of solving a convex programming problem with convergence rate  $O(1/k^2)$ . Soviet Mathematics Doklady, 27(2):372–376, 1983.
- [45] Yurii Nesterov. *Introductory Lectures on Convex Optimization: A Basic Course*, volume 87. Springer Science & Business Media, 2013.
- [46] Yuval Netzer, Tao Wang, Adam Coates, A. Bissacco, Bo Wu, and A. Ng. Reading digits in natural images with unsupervised feature learning. In *NIPS Workshop on Deep Learning and Unsupervised Feature Learning 2011*, 2011.
- [47] Francesco Orabona. Dimension-free exponentiated gradient. *Advances in Neural Information Processing Systems (NeurIPS)*, 2013.
- [48] Francesco Orabona. A modern introduction to online learning. arXiv:1912.13213, 2021.
- [49] Francesco Orabona and Dávid Pál. Coin betting and parameter-free online learning. In *Advances in Neural Information Processing Systems (NeurIPS)*, 2016.
- [50] Courtney Paquette and Katya Scheinberg. A stochastic line search method with expected complexity analysis. *SIAM Journal on Optimization*, 30(1):349–376, 2020.
- [51] Adam Paszke, Sam Gross, Francisco Massa, Adam Lerer, James Bradbury, Gregory Chanan, Trevor Killeen, Zeming Lin, Natalia Gimelshein, Luca Antiga, Alban Desmaison, Andreas Kopf, Edward Yang, Zachary DeVito, Martin Raison, Alykhan Tejani, Sasank Chilamkurthy, Benoit Steiner, Lu Fang, Junjie Bai, and Soumith Chintala. PyTorch: An imperative style, high-performance deep learning library. In Advances in Neural Information Processing Systems (NeurIPS), 2019.
- [52] Fabian Pedregosa, Gaël Varoquaux, Alexandre Gramfort, Vincent Michel, Bertrand Thirion, Olivier Grisel, Mathieu Blondel, Peter Prettenhofer, Ron Weiss, Vincent Dubourg, et al. Scikitlearn: Machine learning in Python. *Journal of Machine Learning Research*, 12:2825–2830, 2011.
- [53] Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, et al. Learning transferable visual models from natural language supervision. In *International Conference on Machine Learning (ICML)*, 2021.
- [54] Sasha Rakhlin and Karthik Sridharan. Optimization, learning, and games with predictable sequences. In *Advances in Neural Information Processing Systems (NeurIPS)*, 2013.
- [55] Sashank J Reddi, Satyen Kale, and Sanjiv Kumar. On the convergence of Adam and beyond. In *International Conference on Learning Representations (ICLR)*, 2018.
- [56] Christopher J Shallue, Jaehoon Lee, Joseph Antognini, Jascha Sohl-Dickstein, Roy Frostig, and George E Dahl. Measuring the effects of data parallelism neural network training. *Journal of Machine Learning Research*, 20:1–49, 2019.

- [57] Ohad Shamir and Tong Zhang. Stochastic gradient descent for non-smooth optimization: Convergence results and optimal averaging schemes. In *International Conference on Machine Learning (ICML)*, 2013.
- [58] Noam Shazeer and Mitchell Stern. Adafactor: Adaptive learning rates with sublinear memory cost. In *International Conference on Machine Learning (ICML)*, 2018.
- [59] Matthew Streeter and H Brendan McMahan. No-regret algorithms for unconstrained online convex optimization. In *Advances in Neural Information Processing Systems (NeurIPS)*, 2012.
- [60] Sharan Vaswani, Aaron Mishkin, Issam Laradji, Mark Schmidt, Gauthier Gidel, and Simon Lacoste-Julien. Painless stochastic gradient: Interpolation, line-search, and convergence rates. In *Advances in Neural Information Processing Systems (NeurIPS)*, 2019.
- [61] Pauli Virtanen, Ralf Gommers, Travis E. Oliphant, Matt Haberland, Tyler Reddy, David Cournapeau, Evgeni Burovski, Pearu Peterson, Warren Weckesser, Jonathan Bright, Stéfan J. van der Walt, Matthew Brett, Joshua Wilson, K. Jarrod Millman, Nikolay Mayorov, Andrew R. J. Nelson, Eric Jones, Robert Kern, Eric Larson, C J Carey, İlhan Polat, Yu Feng, Eric W. Moore, Jake VanderPlas, Denis Laxalde, Josef Perktold, Robert Cimrman, Ian Henriksen, E. A. Quintero, Charles R. Harris, Anne M. Archibald, Antônio H. Ribeiro, Fabian Pedregosa, Paul van Mulbregt, and SciPy 1.0 Contributors. SciPy 1.0: Fundamental Algorithms for Scientific Computing in Python. *Nature Methods*, 17:261–272, 2020.
- [62] Wes McKinney. Data Structures for Statistical Computing in Python. In *Proceedings of the 9th Python in Science Conference*, 2010.
- [63] Ross Wightman. PyTorch image models. https://github.com/rwightman/pytorch-image-models, 2019.
- [64] Jianxiong Xiao, James Hays, Krista A Ehinger, Aude Oliva, and Antonio Torralba. Sun database: Large-scale scene recognition from abbey to zoo. In *Conference on Computer Vision and Pattern Recognition (CVPR)*, 2010.
- [65] Jianxiong Xiao, Krista A Ehinger, James Hays, Antonio Torralba, and Aude Oliva. Sun database: Exploring a large collection of scene categories. *International Journal of Computer Vision*, 119(1):3–22, 2016.
- [66] Sergey Zagoruyko and Nikos Komodakis. Wide residual networks. In *British Machine Vision Conference (BMVC)*, 2016.
- [67] Xiaohua Zhai, Joan Puigcerver, Alexander Kolesnikov, Pierre Ruyssen, Carlos Riquelme, Mario Lucic, Josip Djolonga, André Susano Pinto, Maxim Neumann, Alexey Dosovitskiy, Lucas Beyer, Olivier Bachem, Michael Tschannen, Marcin Michalski, Olivier Bousquet, Sylvain Gelly, and Neil Houlsby. A large-scale study of representation learning with the visual task adaptation benchmark. arXiv:1910.04867, 2019.

# **Contents**

1	<b>Intr</b> 1.1		<b>1</b> 2
2	Prel	iminaries and algorithmic framework	3
3	Ana	v	5
	3.1	General suboptimally bound	5
	3.2	Iterate stability	7
	3.3	Rate of convergence in the noiseless case	8
4	Ana	v	9
	4.1	,	9
	4.2	From bounded to sub-Gaussian noise	.1
	4.3	,	.1
	4.4	Discussion: how parameter-free is our algorithm?	2
5	Exp	eriments 1	3
A	Proc	of for Section 3 (the noiseless setting)	22
		· · · · · · · · · · · · · · · · · · ·	22
	A.2	Proof of Proposition 4	25
	A.3		27
В	Proc	ofs for Section 4 (the stochastic setting)	80
	B.1	Proof of Proposition 7	80
	B.2	Proof of Proposition 8	32
	B.3	Proof of Theorem 9	32
	B.4	Proof of Corollary 11	35
	B.5	Proof of Corollary 13	86
C	Sub	optimality lemmas 3	36
	<b>C</b> .1	Weighted regret to suboptimality conversion (Lemma 14)	6
	C.2	Inductive suboptimality bound (Lemma 15)	37
	C.3	General regret bound (Lemma 16)	9
D	Itera	ate stability lemmas 4	1
	D.1	A weighted regret bound (Lemma 17)	1
	D.2		12
	D.3	Single-step iterate stability (Lemma 19)	13
E	Con	centration bounds 4	4
	E.1	An empirical-Bernstein-type time uniform concentration bound (Lemma 20) 4	4
	E.2		4
	E.3		15

# ACCELERATED PARAMETER-FREE STOCHASTIC OPTIMIZATION

	E.4	Relating $Q_t$ to $Q_t$ (Lemma 23)	46
	E.5	Concentration inequality for bounded random vectors (Lemma 24)	47
F Aux		iliary lemmas	48
	F.1	The growth rate of $\sum_{k} \bar{r}_{k} \alpha_{k}$ (Lemma 25)	48
	F.2	Discrete derivative lemma (Lemma 26)	48
	F.3	Discrete integral lemma (Lemma 27)	49
	F.4	Additional lemmas from prior work	50
G Exp		erimental details	50
	G.1	U-DoG step sizes	50
	G.2	ACCELEGRAD-DOG (A-DOG)	51
	G.3	Convex experiments	51
	G.4	Non-convex experiments	52
	G.5	Implementation details	52

# Appendix A. Proof for Section 3 (the noiseless setting)

## A.1. Proof of Proposition 3

**Proof** Define

$$\begin{split} \rho_t &\coloneqq \frac{1}{\sqrt{Q_t}} \text{ and } \\ \hat{Q}_t &\coloneqq \sum_{k=0}^t \alpha_k^2 \min \big\{ \|\nabla f(\hat{x}_k) - \nabla f(\hat{z}_k)\|^2, \|g_k - m_k\|^2 \big\}. \end{split}$$

Note that in the noiseless setting  $\hat{Q}_T = Q_T$ . However, most of the proof carries over to the noisy setting as well. Therefore, until a later stage of the proof, we do not use that  $m_t = \nabla f(\hat{z}_t)$ ,  $g_t = \nabla f(\hat{x}_t)$  and  $\hat{Q}_t = Q_t$  in the noiseless setting.

Recall the notation  $\tilde{\eta}_{x,t}=\frac{1}{\sqrt{G_{x,t}}}$  and  $\tilde{\eta}_{y,t}=\frac{1}{\sqrt{G_{y,t}}}$ . Algebraic manipulation gives us that for all  $k\geq 0$ 

$$\begin{aligned} \bar{r}_{k}\alpha_{k} \langle g_{k}, x_{k+1} - x_{\star} \rangle &\leq \frac{\bar{r}_{k}^{2}\alpha_{k}^{2}\rho_{k}}{2} \|g_{k} - m_{k}\|^{2} - \sum_{k=0}^{t} \frac{1}{2\rho_{k}} \|x_{k+1} - y_{k}\|^{2} \\ &+ \left(\frac{1}{2\rho_{k}} - \frac{1}{2\tilde{\eta}_{x,k}}\right) \left(\|x_{k+1} - y_{k}\|^{2} + \|x_{k+1} - y_{k+1}\|^{2}\right) \\ &+ \frac{1}{2\tilde{\eta}_{y,k}} \left(\|x_{\star} - y_{k}\|^{2} - \|x_{\star} - y_{k+1}\|^{2}\right); \end{aligned}$$

see Lemma 16 for a proof. Therefore, by summing over both sides of the inequality we get that for all  $t \ge 0$ 

$$\sum_{k=0}^{t} \bar{r}_{k} \alpha_{k} \langle g_{k}, x_{k+1} - x_{\star} \rangle \leq \underbrace{\frac{\bar{r}_{t}^{2}}{2} \sum_{k=0}^{t} \frac{\alpha_{k}^{2} \|g_{k} - m_{k}\|^{2}}{\sqrt{\sum_{j=0}^{t} \alpha_{j}^{2} \|g_{j} - m_{j}\|^{2}}}_{(A)} \underbrace{-\sum_{k=0}^{t} \frac{1}{2\rho_{k}} \|x_{k+1} - y_{k}\|^{2}}_{(B)} + \underbrace{4\bar{r}_{t+1}^{2} \sum_{k=0}^{t} \left[\frac{1}{\rho_{k}} - \frac{1}{\tilde{\eta}_{x,k}}\right]_{+}}_{(C)} + \underbrace{\frac{1}{2} \sum_{k=0}^{t} \frac{1}{\tilde{\eta}_{y,k}} (d_{k}^{2} - d_{k+1}^{2})}_{(D)}.$$

**Bounding** (A): We have  $\sum_{k=0}^{t} \frac{\alpha_k^2 \|g_k - m_k\|^2}{\sqrt{\sum_{j=0}^{k} \alpha_j^2 \|g_j - m_j\|^2}} \le 2\sqrt{\sum_{k=0}^{t} \alpha_k^2 \|g_k - m_k\|^2}$ ; see Lemma 28 with  $s_k = \alpha_k^2 \|g_k - m_k\|^2$ , and therefore

$$\frac{\bar{r}_t^2}{2} \sum_{k=0}^t \frac{\alpha_k^2 \|g_k - m_k\|^2}{\sqrt{\sum_{j=0}^k \alpha_j^2 \|g_j - m_j\|^2}} = \frac{\bar{r}_t^2}{\rho_t}.$$

**Bounding** (B): We have that for all  $k \ge 0$ 

$$\|\nabla f(\hat{x}_{k}) - \nabla f(\hat{z}_{k})\|^{2} \stackrel{(1)}{\leq} \beta^{2} \|\hat{x}_{k} - \hat{z}_{k}\|^{2}$$

$$= \frac{\beta^{2} \bar{r}_{k}^{2} \alpha_{k}^{2}}{\left(\sum_{0=1}^{k} \bar{r}_{i} \alpha_{i}\right)^{2}} \|x_{k+1} - y_{k}\|^{2}$$

$$\stackrel{(2)}{\leq} \frac{4\beta^{2}}{\alpha_{k}^{2}} \|x_{k+1} - y_{k}\|^{2},$$

where (1) is from the  $\beta$ -smoothness of f, and (2) is because  $\bar{r}_k \alpha_k^2 \leq 2 \sum_{0=1}^k \bar{r}_i \alpha_i$  by Lemma 25 . Therefore,

$$-\|x_{k+1} - y_k\|^2 \le -\frac{\alpha_k^2 \|\nabla f(\hat{x}_k) - \nabla f(\hat{z}_k)\|^2}{4\beta^2}.$$

Thus,

$$-\sum_{k=0}^{t} \frac{1}{2\rho_k} \|x_{k+1} - y_k\|^2 \le -\sum_{k=0}^{t} \frac{\alpha_k^2 \|\nabla f(\hat{x}_k) - \nabla f(\hat{z}_k)\|^2}{8\beta^2 \rho_k}$$

**Bounding** (C): As  $\frac{1}{\rho_k} - \frac{1}{\tilde{\eta}_{x,k}}$  is not necessarily non-negative for all  $k \in \{0, \dots, t\}$  we define the set of indices for which it is non-negative as

$$I \triangleq \left\{ k \in \{0, 1, \dots, t\} : \frac{1}{\rho_t} - \frac{1}{\tilde{\eta}_{x,t}} \ge 0 \right\}.$$

Define  $i_k$  as the k-th smallest index in I, and define  $i_{|I|+1} := t+1$ . We note that for all  $k \in I$  then  $i_k \le i_{k+1} - 1 \le t$ . Therefore,

$$\begin{split} (C) &= 4\bar{r}_{t+1}^2 \sum_{k=1}^{|I|} \left(\frac{1}{\rho_{i_k}} - \frac{1}{\tilde{\eta}_{x,i_k}}\right) \leq 4\bar{r}_{t+1}^2 \sum_{k=1}^{|I|} \left(\frac{1}{\rho_{[i_{k+1}-1]}} - \frac{1}{\tilde{\eta}_{x,i_k}}\right) \\ &\leq \frac{4\bar{r}_{t+1}^2}{\rho_t} + 4\bar{r}_{t+1}^2 \sum_{k=2}^{|I|} \left(\frac{1}{\rho_{[i_k-1]}} - \frac{1}{\tilde{\eta}_{x,i_k}}\right) \leq \frac{4\bar{r}_{t+1}^2}{\rho_t} + 4\bar{r}_{t+1}^2 \sum_{k=0}^{t-1} \left[\frac{1}{\rho_k} - \frac{1}{\tilde{\eta}_{x,k+1}}\right]_+. \end{split}$$

**Bounding** (D):

$$\frac{1}{2} \sum_{k=0}^{t} \frac{1}{\tilde{\eta}_{y,t}} (d_k^2 - d_{k+1}^2) = \frac{d_0^2}{2\tilde{\eta}_{y,0}} - \frac{d_{t+1}^2}{2\tilde{\eta}_{y,t}} + \frac{1}{2} \sum_{k=0}^{t} \left( \frac{1}{\tilde{\eta}_{y,k}} - \frac{1}{\tilde{\eta}_{y,k-1}} \right) d_k^2 \\
\leq \frac{d_0^2}{2\tilde{\eta}_{y,0}} - \frac{d_{t+1}^2}{2\tilde{\eta}_{y,t}} + \frac{d_t^2}{2} \sum_{k=0}^{t} \left( \frac{1}{\tilde{\eta}_{y,k}} - \frac{1}{\tilde{\eta}_{y,k-1}} \right).$$

By performing telescopic summation we obtain

$$(D) \le \frac{\bar{d}_{t+1}^2 - d_{t+1}^2}{2\tilde{\eta}_{u,t}}.$$

Let  $s \in \arg\max_{k \le t+1} d_k$ , we have that  $\bar{d}_{t+1}^2 - d_{t+1}^2 = \bar{d}_s^2 - d_{t+1}^2 = (\bar{d}_s - d_{t+1})(\bar{d}_s + d_{t+1}) \le \|y_s - y_{t+1}\|(\bar{d}_s + \bar{d}_{t+1}) \le (\bar{r}_s + r_{t+1})(\bar{d}_s + d_{t+1}) \le 4\bar{r}_{t+1}\bar{d}_{t+1}$ . Thus,

$$(D) \le \frac{2\bar{r}_{t+1}\bar{d}_{t+1}}{\tilde{\eta}_{y,t}}.$$

**Bounding** (A) + (B) + (C) + (D): Combining all of the above, we obtain that

$$\sum_{k=0}^{t} \bar{r}_{k} \alpha_{k} \langle g_{k}, x_{k+1} - x_{\star} \rangle \leq 5 \bar{r}_{t+1} (\bar{r}_{t+1} + \bar{d}_{t+1}) \max \left\{ \frac{1}{\rho_{t}}, \frac{1}{\tilde{\eta}_{y,t}} \right\} + 4 \bar{r}_{t+1}^{2} \sum_{k=0}^{t-1} \left[ \frac{1}{\rho_{k}} - \frac{1}{\tilde{\eta}_{x,k+1}} \right]_{+} - \sum_{k=0}^{t} \frac{\alpha_{k}^{2} ||\nabla f(\hat{x}_{k}) - \nabla f(\hat{z}_{k})||^{2}}{8\beta^{2} \rho_{k}}$$

Therefore, as for any we have that  $G_{x,k} \ge Q_{k-1}$ ,

$$\sum_{k=0}^{t} \bar{r}_{k} \alpha_{k} \langle g_{k}, x_{k+1} - x_{\star} \rangle \leq 5 \bar{r}_{t+1} (\bar{r}_{t+1} + \bar{d}_{t+1}) \sqrt{\max\{G_{y,t}, Q_{t}\}} - \sum_{k=0}^{t} \frac{\alpha_{k}^{2} \|\nabla f(\hat{x}_{k}) - \nabla f(\hat{z}_{k})\|^{2}}{8\beta^{2} \rho_{k}}.$$

Let  $s \geq 0$  and recall that  $\frac{1}{\rho_k} = \sqrt{Q_k}$ . We get that

$$\sum_{k=0}^{t} \bar{r}_{k} \alpha_{k} \langle g_{k}, x_{k+1} - x_{\star} \rangle \leq 10 s \bar{r}_{t+1} (\bar{r}_{t+1} + \bar{d}_{t+1}) \sqrt{\hat{Q}_{t}} - \sum_{k=0}^{t} \frac{\alpha_{k}^{2} \|\nabla f(\hat{x}_{k}) - \nabla f(\hat{z}_{k})\|^{2}}{8\beta^{2}} \sqrt{Q_{k}} + 5 \bar{r}_{t+1} (\bar{r}_{t+1} + \bar{d}_{t+1}) \left( s \sqrt{Q_{t}} - 2s \sqrt{\hat{Q}_{t}} \right) + 5 \bar{r}_{t+1} (\bar{r}_{t+1} + \bar{d}_{t+1}) \left( \sqrt{\max\{G_{y,t}, Q_{t}\}} - s \sqrt{Q_{t}} \right). \tag{16}$$

We have that

$$10s\bar{r}_{t+1}(\bar{r}_{t+1} + \bar{d}_{t+1})\sqrt{\hat{Q}_t} - \sum_{k=0}^t \frac{\alpha_t^2 \|\nabla f(\hat{x}_t) - \nabla f(\hat{z}_t)\|^2}{8\beta^2} \sqrt{Q_k}$$

$$\leq 10s\bar{r}_{t+1}(\bar{r}_{t+1} + \bar{d}_{t+1})\sqrt{\hat{Q}_t} - \sum_{k=0}^t \frac{\alpha_k^2 \min\{\|\nabla f(\hat{x}_k) - \nabla f(\hat{z}_k)\|^2, \|g_k - m_k\|^2\}}{8\beta^2} \sqrt{\hat{Q}_k}.$$

Define  $B_k^2 = \alpha_k^2 \min \left\{ \|\nabla f(\hat{x}_k) - \nabla f(\hat{z}_k)\|^2, \|g_k - m_k\|^2 \right\}, c_1 = 10 s \bar{r}_{t+1} (\bar{r}_{t+1} + \bar{d}_{t+1}), \text{ and } c_2 = 8\beta^2.$  Lemma 27 gives us that for all  $t \geq 0$ 

$$c_1 \sqrt{\sum_{k=0}^t B_k^2 - \sum_{k=0}^t \frac{B_k^2}{c_2}} \sqrt{\sum_{j=0}^k B_j^2} \le 2c_1^{3/2} c_2^{1/2}.$$

Therefore,

$$10s\bar{r}_{t+1}(\bar{r}_{t+1} + \bar{d}_{t+1})\sqrt{\hat{Q}_t} - \sum_{k=0}^t \frac{\alpha_t^2 \|\nabla f(\hat{x}_t) - \nabla f(\hat{z}_t)\|^2}{8\beta^2} \sqrt{Q_k}$$

$$\leq 2(10s\bar{r}_{t+1}(\bar{r}_{t+1} + \bar{d}_{t+1}))^{3/2}(8\beta)^{1/2} \leq 180s^{3/2}\bar{r}_{t+1}(\bar{r}_{t+1} + \bar{d}_{t+1})^2\beta.$$

Combining this result with eq. (16) yields that for all  $t \ge 0$  and  $s \ge 0$ 

$$\sum_{k=0}^{t} \bar{r}_{k} \alpha_{k} \langle g_{k}, x_{k+1} - x_{\star} \rangle \leq 180 s^{3/2} \bar{r}_{t+1} (\bar{r}_{t+1} + \bar{d}_{t+1})^{2} \beta 
+ 5 \bar{r}_{t+1} (\bar{r}_{t+1} + \bar{d}_{t+1}) \left( s \sqrt{Q_{t}} - 2s \sqrt{\hat{Q}_{t}} \right) 
+ 5 \bar{r}_{t+1} (\bar{r}_{t+1} + \bar{d}_{t+1}) \left( \sqrt{\max\{G_{y,t}, Q_{t}\}} - s \sqrt{Q_{t}} \right).$$
(17)

Lemma 14 gives us that

$$f(\hat{x}_t) - f(x_\star) \le \frac{1}{\sum_{k=0}^t \bar{r}_k \alpha_k} \sum_{k=0}^t \bar{r}_k \alpha_k \left\langle \nabla f(\hat{x}_k), x_{k+1} - x_\star \right\rangle.$$

Now, by additionally using the fact that in the noiseless setting

$$\hat{Q}_t = Q_t \ \text{ and }$$
 
$$\sum_{k=0}^t \bar{r}_k \alpha_k \left\langle \nabla f(\hat{x}_k), x_{k+1} - x_\star \right\rangle = \sum_{k=0}^t \bar{r}_k \alpha_k \left\langle g_k, x_{k+1} - x_\star \right\rangle$$

we get that

$$f(\hat{x}_t) - f(x_\star) \le 180s^{3/2} \frac{\bar{r}_{t+1}}{\sum_{k=0}^t \bar{r}_k \alpha_k} \beta \left(\bar{r}_{t+1} + \bar{d}_{t+1}\right)^2 + 5 \frac{\bar{r}_{t+1}}{\sum_{k=0}^t \bar{r}_k \alpha_k} \left(\bar{r}_{t+1} + \bar{d}_{t+1}\right) \left(\sqrt{\max\{G_{y,t}, Q_t\}} - s\sqrt{Q_t}\right).$$

Finally, by using the fact that  $\bar{d}_{t+1} \leq d_0 + \bar{r}_{t+1}$  and because  $\bar{r}_k \alpha_k^2 \leq 2 \sum_{0=1}^k \bar{r}_i \alpha_i$  for all  $k \geq 0$  (Lemma 25), we obtain that

$$f(\hat{x}_t) - f(x_\star) \le O\left(\frac{s^{3/2}\beta(\bar{r}_{t+1} + d_0)^2 + (\bar{r}_{t+1} + d_0)\left[\sqrt{\max\{G_{y,t}, Q_t\}} - s\sqrt{Q_t}\right]_+}{\left(\sum_{k=0}^t \bar{r}_k/\bar{r}_{t+1}\right)^2}\right).$$

# A.2. Proof of Proposition 4

**Proof** For any h > 0 (in this case h = 12), define

$$c_t \coloneqq h \log_+^2 \left( \frac{s + Q_t}{s} \right)$$
 and 
$$\rho_t \coloneqq \frac{1}{c_{t-1} \sqrt{s + Q_t}}.$$

Lemma 16 gives us that, for all  $t \ge 0$ ,

$$\bar{r}_{t}\alpha_{t} \langle g_{t}, x_{t+1} - x_{\star} \rangle \leq \frac{\bar{r}_{t}^{2}\alpha_{t}^{2}\rho_{t}}{2} \|g_{t} - m_{t}\|^{2} + \left(\frac{1}{2\rho_{t}} - \frac{1}{2\tilde{\eta}_{x,t}}\right) (\|x_{t+1} - y_{t}\|^{2} + \|x_{t+1} - y_{t+1}\|^{2}) + \frac{1}{2\tilde{\eta}_{y,t}} (\|x_{\star} - y_{t}\|^{2} - \|x_{\star} - y_{t+1}\|^{2}).$$

From the definitions of  $\rho_t$  and  $\tilde{\eta}_{x,t}=1/\sqrt{G_{x,t}}\leq 1/\sqrt{G_{y,t-1}}\leq 1/\rho_{t-1}$  we obtain that

$$\frac{1}{2\rho_t} - \frac{1}{2\tilde{\eta}_{x,t}} \le \frac{c_{t-1}^2}{2} \rho_t(Q_t - Q_{t-1});$$

See proof in Lemma 26. Now, because we also have that  $\max\{\|x_{t+1}-y_t\|, \|y_{t+1}-x_{t+1}\|\} \le \frac{2\overline{r}_t}{c_{t-1}}$ , we get

$$\bar{r}_t \alpha_t \langle g_t, x_{t+1} - x_{\star} \rangle \leq \frac{9}{2} \bar{r}_t^2 \rho_t \alpha_t^2 ||g_t - m_t||^2 + \frac{1}{2\tilde{\eta}_{u,t}} (d_t^2 - d_{t+1}^2).$$

Thus,

$$2\tilde{\eta}_{y,t}\bar{r}_t\alpha_t \langle g_t, x_{t+1} - x_{\star} \rangle \le 9\bar{r}_t^2\tilde{\eta}_{y,t}\rho_t\alpha_t^2 \|g_t - m_t\|^2 + (d_t^2 - d_{t+1}^2).$$

Consequentially, by summing the two sides of the inequality, we get that for all  $t \ge 0$ 

$$2\sum_{k=0}^{t} \tilde{\eta}_{y,k} \bar{r}_{k} \alpha_{k} \langle g_{k}, x_{k+1} - x_{\star} \rangle \leq 9\sum_{k=0}^{t} \bar{r}_{k}^{2} \tilde{\eta}_{y,k} \rho_{k} \alpha_{k}^{2} \|g_{k} - m_{k}\|^{2} + \sum_{k=0}^{t} (d_{k}^{2} - d_{k+1}^{2})$$

$$\leq \frac{9\bar{r}_{t}^{2}}{h^{2}} \sum_{k=0}^{t} \frac{Q_{k} - Q_{k-1}}{(s + Q_{k}) \log_{+}^{2} \left(\frac{s + Q_{k}}{s}\right)} + \sum_{k=0}^{t} (d_{k}^{2} - d_{k+1}^{2}).$$

Lemma 30 gives us that

$$\sum_{k=0}^{t} \frac{Q_k - Q_{k-1}}{(s + Q_k) \log_+^2 \left(\frac{s + Q_k}{s}\right)} \le 1.$$

Therefore, we obtain that

$$2\sum_{k=0}^{t} \tilde{\eta}_{y,k} \bar{r}_{k} \alpha_{k} \langle g_{k}, x_{k+1} - x_{\star} \rangle \leq \frac{9\bar{r}_{t}^{2}}{h^{2}} + \sum_{k=0}^{t} (d_{k}^{2} - d_{k+1}^{2}).$$

Thus,

$$2\sum_{k=0}^{t} \tilde{\eta}_{y,k} \bar{r}_{k} \alpha_{k} \left\langle \nabla f(\hat{x}_{k}), x_{k+1} - x_{\star} \right\rangle \leq \frac{9\bar{r}_{t}^{2}}{h^{2}} + 2\sum_{k=0}^{t} \tilde{\eta}_{y,k} \bar{r}_{k} \alpha_{k} \left\langle \nabla f(\hat{x}_{k}) - g_{k}, x_{k+1} - x_{\star} \right\rangle + \sum_{k=0}^{t} \left( d_{k}^{2} - d_{k+1}^{2} \right).$$

Consequentially, as Lemma 17 gives us that

$$\sum_{k=0}^{t} \tilde{\eta}_{y,k} \bar{r}_k \alpha_k \left\langle \nabla f(\hat{x}_k), x_{k+1} - x_{\star} \right\rangle \ge 0,$$

we get that

$$0 \le \frac{9\bar{r}_t^2}{h^2} + 2\sum_{k=0}^t \eta_{y,k} \alpha_k \left\langle \nabla f(\hat{x}_k) - g_k, x_{k+1} - x_\star \right\rangle + \sum_{k=0}^t \left( d_k^2 - d_{k+1}^2 \right).$$

Therefore, we get that for all  $t \ge 0$ 

$$d_{t+1}^2 \le \frac{9\bar{r}_t^2}{h^2} + 2\sum_{k=0}^t \eta_{y,k} \alpha_k \left\langle \nabla f(\hat{x}_k) - g_k, x_{k+1} - x_\star \right\rangle + d_0^2. \tag{18}$$

As we are in the noiseless case, and h = 12, we get that for all  $t \ge 0$ 

$$d_{t+1}^{2} \leq \frac{\bar{r}_{t}^{2}}{16} + d_{0}^{2}$$

$$\leq \left(d_{0} + \frac{1}{4}\bar{r}_{t}\right)^{2}.$$

Finally, Lemma 18 now gives us that for all  $t \ge 0$ 

$$d_t < 2d_0$$
 and  $r_t < 4d_0$ .

## A.3. Proof of Theorem 5

**Proof** Define

$$c_t = 12 \log_+^2 \left( \frac{\|m_0\|^2 + Q_t}{\|m_0\|^2} \right).$$

From Lemma 19, we get that for all  $t \ge 0$  the distance between iterates is not large:

$$\max\{\|x_{t+1} - y_t\|, \|x_{t+1} - y_{t+1}\|\} \le \frac{2\bar{r}_t}{c_{t-1}}.$$

Now, we fulfill all the conditions for Proposition 4 and therefore, for all  $t \ge 0$ 

$$\bar{d}_t \leq 2d_0$$
 and  $\bar{r}_t \leq 4d_0$ .

Proposition 3 gives that for all  $t \ge 0$  and for all  $s \ge 0$ 

$$f(\hat{x}_t) - f(x_\star) \le O\left(\frac{s^{3/2}\beta(\bar{r}_{t+1} + d_0)^2 + (\bar{r}_{t+1} + d_0)\left[\sqrt{G_{y,t}} - s\sqrt{Q_t}\right]_+}{\left(\sum_{k=0}^t \bar{r}_k/\bar{r}_{t+1}\right)^2}\right).$$

By using the fact that  $\bar{r}_t \leq 4d_0$ , we get that for all  $t \geq 0$ 

$$f(\hat{x}_t) - f(x_\star) \le O\left(\frac{s^{3/2}\beta d_0^2 + d_0\left[\sqrt{G_{y,t}} - s\sqrt{Q_t}\right]_+}{\left(\sum_{k=0}^t \bar{r}_k/\bar{r}_{t+1}\right)^2}\right). \tag{19}$$

Recall that

$$\tau = \operatorname*{arg\,max}_{t < T} \sum_{i < t} \frac{\bar{r}_i}{\bar{r}_{t+1}}.$$

To show the non-smooth rate, we set s=0 and obtain

$$\sqrt{G_{y,\tau}} \le c_T \sqrt{\max_{k \le T-1} \left\{ \alpha_k^2 \|m_k\|^2 \right\}} + \sum_{k=0}^{T-1} \alpha_k^2 \|g_k - m_k\|^2 \le c_T \sqrt{T^2 L^2 + T^3 L^2} \le 2L T^{3/2} c_T.$$

This result, with eq. (19), gives us that

$$f(\hat{x}_{\tau}) - f(x_{\star}) \le O\left(\frac{Ld_0 T^{3/2}}{\left(\sum_{k=0}^{\tau} \bar{r}_k / \bar{r}_{t+1}\right)^2} c_T\right).$$
 (20)

To show the smooth rate, setting  $s = 2c_{t+1}$  yields

$$\sqrt{G_{y,t}} - s\sqrt{Q_t} \le c_{t+1} \left(\sqrt{Q_t + M_t} - 2\sqrt{Q_t}\right) \le c_{t+1} \left(\sqrt{M_t} - \sqrt{Q_t}\right).$$

For some  $\kappa_t \leq t$  we have that  $\sqrt{M_t} = \alpha_{\kappa_t} \|m_{\kappa_t}\|$ . In addition, the smoothness of f implies that  $\|\nabla f(z)\|^2 \leq 2\beta [f(z) - f(x_\star)]$  for all  $z \in \mathcal{X}$ . Combining this fact with the triangle inequality gives us that, in the noiseless setting,

$$\alpha_{\kappa_t} \| m_{\kappa_t} \| = \alpha_{\kappa_t} \| \nabla f(\hat{z}_{\kappa_t}) \| \le \alpha_{\kappa_t} \| \nabla f(\hat{z}_{\kappa_t}) - \nabla f(\hat{z}_{\kappa_t}) \| + \alpha_{\kappa_t} \sqrt{2\beta} \sqrt{f(\hat{x}_{\kappa_t}) - f(x_\star)}.$$

Thus,

$$\sqrt{M_t} \le \sqrt{Q_t} + \alpha_{\kappa_t} \sqrt{2\beta} \sqrt{f(\hat{x}_{\kappa_t}) - f(x_\star)}.$$

Therefore,

$$\sqrt{G_{y,t}} - s\sqrt{Q_t} \le \alpha_{\kappa_t} \sqrt{2c_{t+1}^2 \beta} \sqrt{f(\hat{x}_{\kappa_t}) - f(x_\star)}.$$

This result, together with eq. (19), give us that for all  $t \ge 0$ , there exist  $\kappa_t \le t$  such as

$$f(\hat{x}_t) - f(x_\star) \le O\left(\frac{c_{t+1}^{3/2}\beta d_0^2 + \alpha_{\kappa_t}\sqrt{c_{t+1}^2\beta d_0^2}\sqrt{f(\hat{x}_{\kappa_t}) - f(x_\star)}}{\left(\sum_{k=0}^t \bar{r}_k/\bar{r}_{t+1}\right)^2}\right).$$

Using the previous inequality and Lemma 15 we obtain that for all  $t \ge 0$  that

$$f(\hat{x}_t) - f(x_\star) \le O\left(\frac{\beta d_0^2}{\left(\sum_{k=0}^t \bar{r}_k / \bar{r}_{t+1}\right)^2} c_{t+1}^2\right). \tag{21}$$

Combining the result from eq. (20) and eq. (21) gives

$$f(\hat{x}_{\tau}) - f(x_{\star}) \le O\left(\frac{\min\{\beta d_0^2, L d_0 T^{3/2}\}}{\left(\sum_{k=0}^{\tau} \bar{r}_k / \bar{r}_{t+1}\right)^2} c_T^2\right). \tag{22}$$

Lemma 29 gives us that

$$\sum_{k=0}^{\tau} \bar{r}_k / \bar{r}_{t+1} \ge \frac{1}{e} \left( \frac{T}{\log_+(\bar{r}_T / r_{\epsilon})} - 1 \right).$$

Thus, if  $T \geq 2 \log_+(\bar{r}_T/r_\epsilon)$  then

$$\frac{1}{\sum_{k=0}^{\tau} \bar{r}_k / \bar{r}_{t+1}} \le O\left(\frac{1}{T} \log_+ \left(\frac{\bar{r}_T}{r_{\epsilon}}\right)\right).$$

Therefore, from eq. (22), we obtain

$$f(\hat{x}_{\tau}) - f(x_{\star}) \le O\left(\frac{\min\{\beta d_0^2, L d_0 T^{3/2}\}}{T^2} c_T^2 \log_+^2 \left(\frac{\bar{r}_T}{r_{\epsilon}}\right)\right).$$
 (23)

We have that

$$c_{T} \leq O\left(\log_{+}^{2}\left(\frac{\|m_{0}\|^{2} + Q_{T-1}}{\|m_{0}\|^{2}}\right)\right)$$

$$\stackrel{(i)}{\leq} O\left(\log_{+}^{2}\left(1 + \frac{T^{3}\min\{\beta d_{0}, L\}}{\|\nabla f(\hat{z}_{0})\|^{2}}\right)\right) \leq O\left(\log_{+}^{2}\left(1 + T\frac{\min\{\beta d_{0}^{2}, Ld_{0}\}}{f(x_{0}) - f(x_{\star})}\right)\right),$$

due to (i) the noiseless setting and f being  $\beta$ -smooth and L-Lipschitz, and (ii) convexity, which implies  $f(x_0) - f(x_\star) \le d_0 \|\nabla f(\hat{z}_0)\|$  Finally, from eq. (23), we obtain

$$f(\hat{x}_{\tau}) - f(x_{\star}) \le O\left(\frac{\min\{\beta d_0^2, Ld_0 T^{3/2}\}}{T^2} \log_+^4 \left(1 + T \frac{\min\{\beta d_0^2, Ld_0\}}{f(x_0) - f(x_{\star})}\right) \log_+^2 \left(\frac{d_0}{r_{\epsilon}}\right)\right). \tag{24}$$

Finally, for  $T \leq 2\log_+(\bar{r}_T/r_\epsilon)$  the theorem holds trivially since  $f(\hat{x}_\tau) - f(x_\star) \leq \min\{\beta \bar{d}_\tau^2, L\bar{d}_\tau\}$  and  $\bar{d}_\tau \leq 2d_0$  by Proposition 8. Therefore,

$$f(\hat{x}_{\tau}) - f(x_{\star}) \le O\left(\min\{\beta d_0^2, Ld_0\}\right) \le O\left(\frac{\min\{\beta \bar{d}_0^2, L\bar{d}_0\}}{T^2}\log_+^2(\bar{r}_T/r_{\epsilon})\right),$$

and so the bound Equation (24) holds in all cases, concluding the proof.

# Appendix B. Proofs for Section 4 (the stochastic setting)

# **B.1. Proof of Proposition 7**

**Proof** Define

$$\hat{Q}_t := \sum_{k=0}^t \alpha_k^2 \min \{ \|\nabla f(\hat{x}_k) - \nabla f(\hat{z}_k)\|^2, \|g_k - m_k\|^2 \}.$$

Our proof continues from eq. (17) in the proof Proposition 3, which also holds for stochastic gradients.

$$\sum_{k=0}^{t} \bar{r}_{k} \alpha_{k} \langle g_{k}, x_{k+1} - x_{\star} \rangle \leq 180 s^{3/2} \bar{r}_{t+1} (\bar{r}_{t+1} + \bar{d}_{t+1})^{2} \beta 
+ 5 \bar{r}_{t+1} (\bar{r}_{t+1} + \bar{d}_{t+1}) \left( s \sqrt{Q_{t}} - 2s \sqrt{\hat{Q}_{t}} \right) 
+ 5 \bar{r}_{t+1} (\bar{r}_{t+1} + \bar{d}_{t+1}) \left( \sqrt{\max\{G_{y,t}, Q_{t}\}} - s \sqrt{Q_{t}} \right).$$

For all k > 0

$$||g_k - m_k||^2 \le 2||\nabla f(\hat{x}_k) - \nabla f(\hat{z}_k)||^2 + 2||(g_k - \nabla f(\hat{x}_k)) - (m_k - \nabla f(\hat{z}_k))||^2$$

$$\le 2\min\{||\nabla f(\hat{x}_k) - \nabla f(\hat{z}_k)||^2, ||g_k - m_k||^2\} + 4||m_t - \nabla f(\hat{z}_t)||^2 + 4||g_t - \nabla f(\hat{x}_t)||^2.$$

Thus, for all k > 0

$$||g_k - m_k||^2 \le 2\min\{||\nabla f(\hat{x}_k) - \nabla f(\hat{z}_k)||^2, ||g_k - m_k||^2\} + 4||m_t - \nabla f(\hat{z}_t)||^2 + 4||g_t - \nabla f(\hat{x}_t)||^2.$$

Multiplying by  $\alpha_k^2$ , summing and recalling that  $\alpha_k \leq k+1$  implies  $Q_t \leq 2\hat{Q}_t + 4(t+1)^3V_t$ , where  $V_t = \frac{1}{t+1}\sum_{k=0}^t \left(\|g_t - \nabla f(\hat{x}_t)\|^2 + \|m_t - \nabla f(\hat{z}_t)\|^2\right)$  is the empirical variance. Substituting into eq. (17), we get that

$$\sum_{k=0}^{t} \bar{r}_{k} \alpha_{k} \langle g_{k}, x_{k+1} - x_{\star} \rangle \leq 180 s^{3/2} \bar{r}_{t+1} (\bar{r}_{t+1} + \bar{d}_{t+1})^{2} \beta 
+ 5 \bar{r}_{t+1} (\bar{r}_{t+1} + \bar{d}_{t+1}) \left( \sqrt{\max\{G_{y,t}, Q_{t}\}} - s \sqrt{Q_{t}} \right) 
+ 10 s \bar{r}_{t+1} (\bar{r}_{t+1} + \bar{d}_{t+1}) \sqrt{(t+1)^{3} V_{t}}.$$
(25)

 $\text{Lemma 21 gives us that with probability of at least } 1-\delta-\mathbb{P}\big[\bar{\mathfrak{b}}_{T-1}>\mathfrak{B}\big], \text{ for all } t\in\{0,1,\ldots,T-1\},$ 

$$\left| \sum_{k=0}^{t} \bar{r}_{t} \alpha_{k} \left\langle \nabla f(\hat{x}_{k}) - g_{k}, x_{k+1} - x_{\star} \right\rangle \right| \leq 8\alpha_{t} \bar{r}_{t} (\bar{r}_{t+1} + d_{0}) \sqrt{\theta_{t+1,\delta} \sum_{k=0}^{t} \|\nabla f(\hat{x}_{k}) - g_{k}\|^{2} + (\theta_{t+1,\delta} \mathfrak{B})^{2}}.$$

Using the previous equality and the definition of  $V_t$  we obtain that

$$\sum_{k=0}^{t} \bar{r}_{k} \alpha_{k} \left\langle \nabla f(\hat{x}_{k}), x_{k+1} - x_{\star} \right\rangle$$

$$= \sum_{k=0}^{t} \bar{r}_{k} \alpha_{k} \left\langle g_{k}, x_{k+1} - x_{\star} \right\rangle + \sum_{k=0}^{t} \bar{r}_{k} \alpha_{k} \left\langle \nabla f(\hat{x}_{k}) - g_{k}, x_{k+1} - x_{\star} \right\rangle$$

$$\leq \sum_{k=0}^{t} \bar{r}_{k} \alpha_{k} \left\langle g_{k}, x_{k+1} - x_{\star} \right\rangle + 8\alpha_{t} \bar{r}_{t} (\bar{r}_{t+1} + d_{0}) \sqrt{(t+1)\theta_{t+1,\delta} V_{t} + (\theta_{t+1,\delta} \mathfrak{B})^{2}}. \tag{26}$$

Lemma 14 gives us that

$$f(\hat{x}_t) - f(x_\star) \le \frac{1}{\sum_{k=0}^t \bar{r}_k \alpha_k} \sum_{k=0}^t \bar{r}_k \alpha_k \left\langle \nabla f(\hat{x}_k), x_{k+1} - x_\star \right\rangle.$$

By combining the above inequality with eq. (25) and eq. (26), we obtain

$$f(\hat{x}_{t}) - f(x_{\star}) \leq 180s^{3/2} \frac{\bar{r}_{t+1}}{\sum_{k=0}^{t} \bar{r}_{k} \alpha_{k}} \beta \left(\bar{r}_{t+1} + \bar{d}_{t+1}\right)^{2}$$

$$+ 5 \frac{\bar{r}_{t+1}}{\sum_{k=0}^{t} \bar{r}_{k} \alpha_{k}} \left(\bar{r}_{t+1} + \bar{d}_{t+1}\right) \left(\sqrt{\max\{G_{y,t}, Q_{t}\}} - s\sqrt{Q_{t}}\right)$$

$$+ 10(1+s) \frac{\bar{r}_{t+1}}{\sum_{k=0}^{t} \bar{r}_{k} \alpha_{k}} \left(\bar{r}_{t+1} + \bar{d}_{t+1}\right) \sqrt{(t+1)^{3} V_{t} + (\theta_{t+1}, \delta \mathfrak{B})^{2}}.$$

Now, as Lemma 25 gives us that  $\bar{r}_t \alpha_t^2 \leq 2 \sum_{k=0}^t \bar{r}_k \alpha_k$ , we obtain that

$$f(\hat{x}_{t}) - f(x_{\star}) \leq 360s^{3/2} \frac{1}{\left(\sum_{k=0}^{t} \bar{r}_{k}/\bar{r}_{t+1}\right)^{2}} \beta(\bar{r}_{t+1} + \bar{d}_{t+1})^{2} + 10 \frac{1}{\left(\sum_{k=0}^{t} \bar{r}_{k}/\bar{r}_{t+1}\right)^{2}} (\bar{r}_{t+1} + \bar{d}_{t+1}) \left(\sqrt{\max\{G_{y,t}, Q_{t}\}} - s\sqrt{Q_{t}}\right) + 20 \frac{1}{\left(\sum_{k=0}^{t} \bar{r}_{k}/\bar{r}_{t+1}\right)^{2}} (\bar{r}_{t+1} + \bar{d}_{t+1}) \sqrt{(t+1)^{3}V_{t} + (\theta_{t+1,\delta}\mathfrak{B})^{2}}.$$

Finally, because that  $\bar{d}_{t+1} \leq d_0 + \bar{r}_{t+1}$ , we get that for any  $\mathfrak{B} > 0$  with probability of at least  $1 - \delta - \mathbb{P}[\bar{\mathfrak{b}}_{T-1} > \mathfrak{B}]$  we have that for all t < T and for any number  $s \geq 0$ 

$$f(\hat{x}_t) - f(x_\star) \le O\left(\text{RHS}_{eq. (4)} + \frac{(1+s)(\bar{r}_{t+1} + d_0)\sqrt{t^3\theta_{t,\delta}V_t + (t\theta_{t,\delta}\mathfrak{B})^2}}{\left(\sum_{k=0}^t \bar{r}_k/\bar{r}_{t+1}\right)^2}\right)$$

where RHS<sub>eq. (4)</sub> =  $\frac{s^{3/2}\beta(\bar{r}_{t+1}+d_0)^2 + (\bar{r}_{t+1}+d_0)\left[\sqrt{\max\{G_{y,t},Q_t\}} - s\sqrt{Q_t}\right]_+}{\left(\sum_{k=0}^t \bar{r}_k/\bar{r}_{t+1}\right)^2} \text{ is the error term appearing in Proposition 3.}$ 

## **B.2. Proof of Proposition 8**

**Proof** The proof continues from eq. (18) in the proof of Proposition 4, which also holds for stochastic gradients. Substituting h = 400 in eq. (18) gives, for all  $t \ge 0$ ,

$$d_{t+1}^2 \le \frac{9\bar{r}_t^2}{400^2} + 2\sum_{k=0}^t \eta_{y,k} \alpha_k \left\langle \nabla f(\hat{x}_k) - g_k, x_{k+1} - x_\star \right\rangle + d_0^2.$$

Now, Lemma 22 gives us that with probability at least  $1 - \delta$ , for all t < T

$$\left| \sum_{k=0}^{t} \alpha_k \eta_{y,k} \left\langle g_k - \nabla f(\hat{x}_k), x_{k+1} - x_{\star} \right\rangle \right| \leq \frac{12\theta_{t+1,\delta}}{400\theta_{T,\delta}} \bar{r}_t (\bar{r}_{t+1} + d_0)$$

$$\leq \frac{12\theta_{t+1,\delta}}{400\theta_{T,\delta}} (\bar{r}_t \bar{r}_{t+1} + \bar{r}_t d_0) \leq \frac{12}{400} \left( 1 + \frac{3}{400} \right) \bar{r}_t^2 + \frac{12}{400} \bar{r}_t d_0.$$

Therefore,

$$d_{t+1}^2 \le \frac{81\bar{r}_t^2}{400^2} + \frac{24}{400}\bar{r}_t^2 + \frac{24}{400}\bar{r}_t d_0 + d_0^2 \le \frac{\bar{r}_t^2}{16} + \frac{\bar{r}_t d_0}{2} + d_0.$$

Thus, with probability of at least  $1 - \delta$ , for all t < T

$$d_{t+1}^2 \le \left(d_0 + \frac{1}{4}\bar{r}_t\right)^2.$$

Finally, Lemma 18 gives us that with probability of at least  $1 - \delta$  for all t < T

$$d_t \leq 2d_0$$
 and  $r_t \leq 4d_0$ .

# **B.3. Proof of Theorem 9**

**Proof** Recall the notation

$$\tilde{q}_t \coloneqq 2\alpha_t^2 \|\tilde{g}_t - m_t\|^2$$
,  $\bar{Q}_t \coloneqq \sum_{k=0}^t \max\{q_k, \tilde{q}_k\}$  and  $p_t \coloneqq 8(t+1)^2 \bar{\mathfrak{b}}_t^2$ ,

and that our step sizes are of the form (2) with

$$G_{y,t} = \hat{c}_t^2 \max\{\|m_0\|^2 + p_0 + p_t + \tilde{q}_t + \bar{Q}_{t-1}, M_t\},\$$

where

$$\hat{c}_t = 400\theta_{T,\delta} \log_+^2 \left( 1 + \frac{p_t + \tilde{q}_t + \bar{Q}_{t-1}}{\|m_0\|^2 + p_0} \right).$$

We begin by verifying the conditions of Proposition 8 with  $s = ||m_0||^2 + p_0$ , where condition (iv) holds by construction. By Assumption 6 we have

$$||q_t - \tilde{q}_t||^2 < 2||q_t - \nabla f(\hat{x}_t)||^2 + 2||\tilde{q}_t - \nabla f(\hat{x}_t)||^2 < 4\bar{\mathfrak{b}}_t^2$$

Therefore, since  $t + 1 \ge \alpha_t$ , we have

$$\tilde{q}_t + p_t \ge \alpha_t^2 (2\|\tilde{g}_t - m_t\|^2 + 2\|g_t - \tilde{g}_t\|^2) \ge \alpha_t^2 \|g_t - m_t\|^2 = q_t,$$

and consequently

$$\tilde{q}_t + p_t + \bar{Q}_{t-1} \ge Q_t.$$

Defining

$$c_t = 400\theta_{T,\delta} \log_+^2 \left(1 + \frac{Q_t}{\|m_0\|^2 + p_0}\right),$$

we conclude that

$$G_{y,t} \ge c_t^2 \max\{\|m_0\|^2 + p_0 + Q_t, M_t\} \ge c_t^2(\|m_0\|^2 + p_0 + Q_t)$$

so that condition (i) of Proposition 8 holds. Next, since

$$G_{y,t} \ge c_t^2 \max\{Q_t, M_t\} \ge c_t^2 \alpha_t^2 \max\{\|g_t - m_t\|^2, \|m_t\|^2\},$$

Lemma 19 guarantees condition (ii) of Proposition 8. Finally, we note that

$$p_t \ge 8\alpha_t^2 \max\{\|g_t - \nabla f(\hat{x}_t)\|^2, \|\tilde{g}_t - \nabla f(\hat{x}_t)\|^2\}$$

and

$$p_t + \tilde{q}_t \ge \alpha_t^2 (2||m_t - \tilde{g}_t||^2 + 2||\tilde{g}_t - \nabla f(\hat{x}_t)||^2) \ge \alpha_t^2 ||m_t - \nabla f(\hat{x}_t)||^2.$$

Therefore, as  $\sqrt{G_{y,t}} \ge c_t \sqrt{p_t + \tilde{q}_t}$ , condition (iii) of Proposition 8 holds.

As all the conditions for Proposition 8 hold, with probability of at least  $1 - \delta$ , for all  $t \ge 0$ 

$$\bar{d}_t \leq 2d_0$$
 and  $\bar{r}_t \leq 4d_0$ .

Recalling that  $\mathfrak{b}_{\star} \coloneqq \max_{x:\|x-x_{\star}\| \leq 2d_0} \{\mathfrak{b}(x)\}$ , this also implies that  $\mathbb{P}[\bar{\mathfrak{b}}_{T-1} > \mathfrak{b}_{\star}] \leq \delta$ .

We now combine the conclusions of Proposition 8 with Proposition 3 to obtain a suboptimality bound for U-DoG. Substituting  $\mathbb{P}(\bar{r}_T \leq 4d_0) \leq \delta$  and  $\mathbb{P}[\bar{\mathfrak{b}}_{T-1} > \mathfrak{b}_{\star}] \leq \delta$  into Proposition 7 we get that, with probability at least  $1 - 3\delta$ , for all t < T and  $s \geq 0$ ,

$$f(\hat{x}_{t}) - f(x_{\star}) \leq O\left(\frac{s^{3/2}\beta d_{0}^{2} + d_{0}\left[\sqrt{G_{y,t}} - s\sqrt{Q_{t}}\right]_{+} + (1+s)d_{0}\sqrt{t^{3}\theta_{t+1,\delta}V_{t} + (t\theta_{t+1,\delta}\mathfrak{b}_{\star})^{2}}}{\left(\sum_{k=0}^{t} \bar{r}_{k}/\bar{r}_{t+1}\right)^{2}}\right).$$
(27)

To simplify  $G_{y,t}$  in the bound above, we invoke Lemma 23 which gives that, with probability at least  $1 - \delta - \mathbb{P}[\bar{\mathfrak{b}}_{T-1} > \mathfrak{b}_{\star}] \geq 1 - 2\delta$ , for all t < T,

$$\bar{Q}_t \leq 5Q_t + 80(t+1)^3 \sqrt{\theta_{t+1,\delta}} V_t + 2(t+1)^2 \theta_{t+1,\delta} \mathfrak{b}_{+}^2$$

and hence

$$\sqrt{G_{y,t}} \le \hat{c}_t \sqrt{\bar{Q}_t + 2M_t + 2p_t} = O\left(\hat{c}_t \sqrt{Q_t} + \hat{c}_t \sqrt{M_t} + \hat{c}_t \theta_{T,\delta} \sqrt{t^3 V_t + t^2 \mathfrak{b}_{\star}^2}\right).$$

Combining this with the bound (27) and replacing s with  $s\hat{c}_t\sqrt{3}$ , we get that with probability at least  $1-5\delta$ , for all t< T and  $s\geq 0$ ,

$$f(\hat{x}_t) - f(x_\star) \le O\left(\frac{s^{3/2}\hat{c}_t^{3/2}\beta d_0^2 + \hat{c}_t d_0 \left(\left[(1-s)\sqrt{Q_t} + \sqrt{M_t}\right]_+ + (1+s)\theta_{T,\delta}\sqrt{t^3V_t + t^2\mathfrak{b}_\star^2}\right)}{\left(\sum_{k=0}^t \bar{r}_k/\bar{r}_{t+1}\right)^2}\right)$$
(28)

The remainder of the proof parallels the proof of Theorem 5, where we specialize our bound to the Lipschitz and smooth cases by choosing different values of s. For the Lipschitz case, we use the facts that

$$Q_t \le 4 \sum_{k \le T} \alpha_t^2 (\|\nabla f(\hat{x}_k)\|^2 + \|\nabla f(\hat{z}_k)\|^2 + \|g_k - \nabla f(\hat{x}_k)\|^2 + \|m_k - \nabla f(\hat{z}_k)\|^2) = O(L^2 T^3 + V_T T^3)$$

$$Q_t = O(L^2T^3)$$
 and  $M_t \leq O(L^2T^2)$  and (under the event  $\bar{d}_T \leq 2d_0$ )

$$M_t \le \max_{k \le T} \left\{ 2\alpha_t^2 (\|\nabla f(\hat{z}_k)\|^2 + \|m_k - \nabla f(\hat{z}_k)\|^2) \right\} = O(L^2 T^2 + \mathfrak{b}_{\star}^2 T^2),$$

giving the suboptimality bound. Substituting these expression and s=0 into (28) we get, for all t < T,

$$f(\hat{x}_t) - f(x_\star) \le O\left(\hat{c}_t \frac{Ld_0 T^{3/2} + d_0 \theta_{T,\delta} \sqrt{T^3 V_T + T^2 \mathfrak{b}_\star^2}}{\left(\sum_{k=0}^t \bar{r}_k / \bar{r}_{t+1}\right)^2}\right). \tag{29}$$

For the smooth case and any t < T, let  $\kappa_t \le t$  be such that For some  $\kappa_t \le t$  we have that

$$\sqrt{M_t} = \alpha_{\kappa_t} || m_{\kappa_t} ||.$$

The smoothness of f implies that  $\|\nabla f(z)\|^2 \leq 2\beta [f(z) - f(x_\star)]$  for all  $z \in \mathcal{X}$ . Combining this fact with the triangle inequality gives us that

$$\alpha_{\kappa_t} \| m_{\kappa_t} \| \le \alpha_{\kappa_t} \| \nabla f(\hat{x}_{\kappa_t}) - \nabla f(\hat{z}_{\kappa_t}) \| + \alpha_{\kappa_t} \| m_{\kappa_t} - \nabla f(\hat{z}_{\kappa_t}) \| + \alpha_{\kappa_t} \sqrt{2\beta} \sqrt{f(\hat{x}_{\kappa_t}) - f(x_\star)}$$

and therefore,

$$\sqrt{M_t} < \sqrt{Q_t} + \sqrt{(t+1)^3 V_t} + \alpha_{\kappa_t} \sqrt{2\beta} \sqrt{f(\hat{x}_{\kappa_t}) - f(x_\star)}$$

Substituting into eq. (28) and taking s = 2, we get, for all t < T,

$$f(\hat{x}_t) - f(x_\star) \le O\left(\frac{\hat{c}_t^{3/2}\beta d_0^2 + \hat{c}_t d_0 \theta_{T,\delta} \sqrt{T^3 V_T + T^2 \mathfrak{b}_\star^2} + \alpha_{\kappa_t} \sqrt{\hat{c}_{t+1}^2 \beta d_0^2} \sqrt{f(\hat{x}_{\kappa_t}) - f(x_\star)}}{\left(\sum_{k=0}^t \bar{r}_k / \bar{r}_{t+1}\right)^2}\right).$$

Applying Lemma 15 and noting that  $\theta_{T,\delta} \leq \hat{c}_t$  simplifies the bound to

$$f(\hat{x}_t) - f(x_\star) \le O\left(\frac{\hat{c}_T^2 \beta d_0^2 + \hat{c}_T \theta_{T,\delta} d_0 \sqrt{T^3 V_T + T^2 \mathfrak{b}_\star^2}}{\left(\sum_{k=0}^t \bar{r}_k / \bar{r}_{t+1}\right)^2}\right). \tag{30}$$

Combining the bounds eq. (29) and eq. (30) and noting that  $\theta_{T,\delta} \leq \hat{c}_T$ , we conclude that, with probability at least  $1 - 5\delta$ , for all t < T,

$$f(\hat{x}_t) - f(x_\star) \le O\left(\hat{c}_T^2 \cdot \frac{\min\left\{\beta d_0^2, L d_0 T^{3/2}\right\} + d_0 \sqrt{T^3 V_{T-1} + T^2 \mathfrak{b}_\star^2}}{\left(\sum_{k=0}^{\tau} \bar{r}_k / \bar{r}_{t+1}\right)^2}\right).$$

For  $\tau = \arg\max_{t < T} \sum_{i \le t} \frac{\bar{r}_i}{\bar{r}_{t+1}}$ , Lemma 29 gives us that

$$\sum_{k=0}^{\tau} \bar{r}_k / \bar{r}_{t+1} \ge \frac{1}{e} \left( \frac{T}{\log_+(\bar{r}_T / r_\epsilon)} - 1 \right).$$

Thus, for  $T \ge 2\log_+(\bar{r}_T/r_\epsilon)$  we get (under the event  $\bar{r}_T \le 4d_0$ )

$$f(\hat{x}_{\tau}) - f(x_{\star}) \le O\left(\hat{c}_{T}^{2} \log_{+}^{2} \left(\frac{d_{0}}{r_{\epsilon}}\right) \cdot \frac{\min\left\{\beta d_{0}^{2}, L d_{0} T^{3/2}\right\} + d_{0} \sqrt{T^{3} V_{T-1} + T^{2} \mathfrak{b}_{\star}^{2}}}{T^{2}}\right),$$

which establishes the theorem, since

$$\hat{c}_{T} \stackrel{(i)}{\leq} O\left(\log_{+}^{2} \left(1 + \frac{T^{2}\mathfrak{b}_{\star}^{2} + \bar{Q}_{T-1}}{\|\nabla f(\hat{z}_{0})\|^{2}}\right)\right) \leq O\left(\log_{+}^{2} \left(1 + \frac{T^{3}\mathfrak{b}_{\star}^{2} + T^{3}\sum_{k=0}^{T-1}\|\nabla f(\hat{x}_{k}) - \nabla f(\hat{z}_{k})\|^{2}}{\|\nabla f(\hat{z}_{0})\|^{2}}\right)\right) \\
\leq O\left(\log_{+}^{2} \left(1 + \frac{T^{3}\mathfrak{b}_{\star}^{2} + T^{3}\min\{\beta d_{0}, L\}}{\|\nabla f(\hat{z}_{0})\|^{2}}\right)\right) \stackrel{(ii)}{\leq} O\left(\log_{+}^{2} \left(1 + \frac{T^{3}\mathfrak{b}_{\star}^{2} d_{0}^{2} + T^{3}\min\{\beta d_{0}^{3}, L d_{0}^{2}\}}{f(x_{0}) - f(x_{\star})}\right)\right) \\
= O\left(\log_{+}^{2} \left(1 + T\frac{\mathfrak{b}_{\star} d_{0} + \min\{\beta d_{0}^{2}, L d_{0}\}}{f(x_{0}) - f(x_{\star})}\right)\right),$$

where (i) is because  $\|\nabla f(\hat{z}_0)\|^2 \le \|\nabla f(\hat{z}_0) - m_0 + m_0\|^2 \le 2\|m_0\|^2 + p_0$ , and (ii) is from convexity:  $f(x_0) - f(x_\star) \le d_0 \|\nabla f(\hat{z}_0)\|$ .

Finally, when  $T \leq 2 \log_+(\bar{r}_T/r_\epsilon)$  the required bound is immediate from problem geometry, as explained at the end of the proof of Theorem 5.

#### **B.4. Proof of Corollary 11**

**Proof** Define

$$\delta_t' = \frac{\delta}{5(t+1)^2}.$$

A black-box reduction from sub-Gaussian to bounded stochastic gradient (Lemma 31) shows that at each iteration t, with probability at least  $1-\delta'_t$ , a call to a  $\sigma^2$ -sub-Gaussian subgradient oracle produces an identical result to a call to an alternative stochastic gradient that is bounded by  $3\sigma\sqrt{\log(3/\delta'_t)}$ .

We apply Theorem 9 to U-DoG with the alternative, bounded stochastic gradient oracle. Thus, for this setting, with probability at least  $1-5\delta$ , we have  $\bar{d}_T \leq 2d_0$ ,  $\bar{r}_T \leq 4d_0$ , and the suboptimality bound (15) holds for  $\mathfrak{b}_\star = \sigma_\star\varsigma_{T-1,\delta}$ . To conclude the proof we use Lemma 31 to show that the

algorithm described above produces output different than U-DoG with the original sub-Gaussian oracle as at most

$$3\sum_{t=0}^{\infty} \delta_t' \le \frac{3\delta}{5}\sum_{t=1}^{\infty} \frac{1}{t^2} \le \frac{3 \cdot \pi^2}{5 \cdot 6} \delta \le \delta,$$

where the factor of 3 comes from the fact that every U-DoG iteration involves 3 stochastic gradient queries.

## **B.5. Proof of Corollary 13**

**Proof** A mini-batch of B gradient oracle results, each with noise bounded by L, is a  $\frac{2L^2}{B}$ -sub-Gaussian (see Lemma 24), and we can therefore apply Corollary 11 with  $\sigma_t^2 = \frac{2L^2}{B}$ . Moreover, reusing the sub-Gaussian-to-bounded reduction in the proof of Corollary 11 (Section B.4) we get that, with probability at least  $1-6\delta$ ,

$$\sqrt{V_T} \le \frac{\sqrt{2}L}{\sqrt{B}} \varsigma_{T,\delta}$$

holds in addition to the suboptimality bound given by Corollary 11. Substituting the above bound on  $\sqrt{V_T}$  along with  $\mathfrak{b}_{\star} \leq \sqrt{2} \frac{L}{\sqrt{B}} \zeta_{T,\delta}$  concludes the proof.

# Appendix C. Suboptimality lemmas

# C.1. Weighted regret to suboptimality conversion (Lemma 14)

The following lemma is a straightforward reproduction of Lemma 1 from Kavis et al. [30] with minor changes. In addition, we use the proof of the following lemma as a starting point for the proof of Lemma 17.

**Lemma 14 (Kavis et al. [30])** For any sequence of positive numbers  $\omega_0, \omega_1, \omega_2, \ldots$ , define

$$\hat{x}_t := \frac{\sum_{k=0}^t \omega_k x_{k+1}}{\sum_{k=0}^t \omega_k}.$$

We have that for any T > 0

$$f(\hat{x}_{T-1}) - f(x_{\star}) \le \frac{1}{\sum_{t=0}^{T-1} \omega_t} \sum_{t=0}^{T-1} \omega_t \langle \nabla f(\hat{x}_t), x_{t+1} - x_{\star} \rangle.$$

**Proof** For any  $t \ge 0$  we have that

$$\omega_{t} \langle \nabla f(\hat{x}_{t}), x_{t+1} - x_{\star} \rangle = \omega_{t} \left\langle \nabla f(\hat{x}_{t}), \frac{\sum_{k=0}^{t} \omega_{k}}{\omega_{t}} \hat{x}_{t} - \frac{\sum_{k=0}^{t-1} \omega_{k}}{\omega_{t}} \hat{x}_{t-1} - x_{\star} \right\rangle$$

$$= \omega_{t} \left\langle \nabla f(\hat{x}_{t}), \frac{\sum_{k=0}^{t} \omega_{k}}{\omega_{t}} (\hat{x}_{t} - x_{\star}) - \frac{\sum_{k=0}^{t-1} \omega_{k}}{\omega_{t}} (\hat{x}_{t-1} - x_{\star}) \right\rangle$$

$$= \sum_{k=0}^{t} \omega_{k} \langle \nabla f(\hat{x}_{t}), \hat{x}_{t} - x_{\star} \rangle - \sum_{k=0}^{t-1} \omega_{k} \langle \nabla f(\hat{x}_{t}), \hat{x}_{t-1} - x_{\star} \rangle$$

$$= \omega_{t} \langle \nabla f(\hat{x}_{t}), \hat{x}_{t} - x_{\star} \rangle + \sum_{k=0}^{t-1} \omega_{k} \langle \nabla f(\hat{x}_{t}), \hat{x}_{t} - \hat{x}_{t-1} \rangle.$$

By using the convexity of f, we get

$$\omega_t \langle \nabla f(\hat{x}_t), x_{t+1} - x_{\star} \rangle \ge \omega_t (f(\hat{x}_t) - f(x_{\star})) + \sum_{k=0}^{t-1} \omega_k (f(\hat{x}_t) - f(\hat{x}_{t-1})). \tag{31}$$

Therefore, for any T > 0

$$\sum_{t=0}^{T-1} \omega_t \left\langle \nabla f(\hat{x}_t), x_{t+1} - x_\star \right\rangle \ge \sum_{t=0}^{T-1} \omega_t (f(\hat{x}_t) - f(x_\star)) + \sum_{t=0}^{T-1} \sum_{k=0}^{t-1} \omega_k (f(\hat{x}_t) - f(\hat{x}_{t-1}))$$

$$= \sum_{t=0}^{T-1} \omega_t (f(\hat{x}_t) - f(x_\star)) + \sum_{k=0}^{T-2} \sum_{t=k+1}^{T-1} \omega_k (f(\hat{x}_t) - f(\hat{x}_{t-1})).$$

By performing a telescopic summation, we obtain

$$\sum_{t=0}^{T-1} \omega_t \langle \nabla f(\hat{x}_t), x_{t+1} - x_{\star} \rangle \ge \sum_{t=0}^{T-1} \omega_t (f(\hat{x}_t) - f(x_{\star})) + \sum_{t=0}^{T-2} \omega_t (f(\hat{x}_{T-1}) - f(\hat{x}_t))$$

$$= \omega_{T-1} (f(\hat{x}_{T-1}) - f(x_{\star})) + \sum_{t=0}^{T-2} \omega_t (f(\hat{x}_t) - f(x_{\star}) + f(\hat{x}_{T-1}) - f(\hat{x}_t))$$

$$= \sum_{t=0}^{T-1} \omega_t (f(\hat{x}_{T-1}) - f(x_{\star})).$$

Dividing both sides by  $\sum_{t=0}^{T-1} \omega_t$  concludes the proof.

#### C.2. Inductive suboptimality bound (Lemma 15)

**Lemma 15** Let  $s_0, s_1, \ldots, s_{T-1}$  and  $h_0, h_1, \ldots, h_{T-1}$  be non-negative non-decreasing sequences. Let b > 1 such that  $\bar{r}_{t+1}/\bar{r}_t \leq b$  for any  $t \in \{0, 1, 2 \ldots, T-1\}$ . If for all  $t \in \{0, 1, 2 \ldots, T-1\}$  there exist  $\kappa_t \in \{0, 1, 2 \ldots, t\}$  such that

$$f(\hat{x}_t) - f(x_\star) \le \frac{\alpha_{\kappa_t} \sqrt{s_t} \sqrt{f(\hat{x}_{\kappa_t}) - f(x_\star)} + h_t}{\left(\sum_{k=0}^t \bar{r}_k / \bar{r}_{t+1}\right)^2},$$

then for all  $t \in \{0, 1, 2, \dots, T-1\}$  we have that

$$f(\hat{x}_t) - f(x_\star) \le \frac{4b^2(s_t + h_t)}{\left(\sum_{k=0}^t \bar{r}_k/\bar{r}_{t+1}\right)^2}.$$

**Proof** We prove by induction that

$$f(\hat{x}_t) - f(x_\star) \le \frac{4b^2(s_t + h_t)}{\left(\sum_{k=0}^t \bar{r}_k / \bar{r}_{t+1}\right)^2}.$$

We will only use the induction assumption for the case where  $\kappa_t < t$ .

If  $\kappa_t = t$ : We have that

$$f(\hat{x}_{t}) - f(x_{\star}) \leq \frac{\alpha_{\kappa_{t}} \sqrt{s_{t}} \sqrt{f(\hat{x}_{\kappa_{t}}) - f(x_{\star})} + h_{t}}{\left(\sum_{k=0}^{t} \bar{r}_{k} / \bar{r}_{t+1}\right)^{2}}$$

$$\leq \frac{\frac{\bar{r}_{t+1}}{\bar{r}_{t}} \sqrt{s_{t}} \sqrt{f(\hat{x}_{\kappa_{t}}) - f(x_{\star})}}{\sum_{k=0}^{t} \bar{r}_{k} / \bar{r}_{t+1}} + \frac{h_{t}}{\left(\sum_{k=0}^{t} \bar{r}_{k} / \bar{r}_{t+1}\right)^{2}}$$

$$\leq \frac{b\sqrt{s_{t}} \sqrt{f(\hat{x}_{\kappa_{t}}) - f(x_{\star})}}{\sum_{k=0}^{t} \bar{r}_{k} / \bar{r}_{t+1}} + \frac{h_{t}}{\left(\sum_{k=0}^{t} \bar{r}_{k} / \bar{r}_{t+1}\right)^{2}}.$$

Thus,

$$f(\hat{x}_t) - f(x_\star) - \frac{b\sqrt{s_t}\sqrt{f(\hat{x}_{\kappa_t}) - f(x_\star)}}{\sum_{k=0}^t \bar{r}_k/\bar{r}_{t+1}} \le \frac{h_t}{\left(\sum_{k=0}^t \bar{r}_k/\bar{r}_{t+1}\right)^2}.$$

If

$$\frac{f(\hat{x}_t) - f(x_{\star})}{2} \le f(\hat{x}_t) - f(x_{\star}) - \frac{b\sqrt{s_t}\sqrt{f(\hat{x}_{\kappa_t}) - f(x_{\star})}}{\sum_{k=0}^{t} \bar{r}_k/\bar{r}_{t+1}},$$

then

$$f(\hat{x}_t) - f(x_\star) \le \frac{2h_t}{\left(\sum_{k=0}^t \bar{r}_k / \bar{r}_{t+1}\right)^2}.$$

Otherwise,

$$\frac{f(\hat{x}_t) - f(x_{\star})}{2} \le \frac{b\sqrt{s_t}\sqrt{f(\hat{x}_{\kappa_t}) - f(x_{\star})}}{\sum_{k=0}^{t} \bar{r}_k/\bar{r}_{t+1}}.$$

Therefore,

$$\sqrt{f(\hat{x}_t) - f(x_\star)} \le \frac{2b\sqrt{s_t}}{\sum_{k=0}^t \bar{r}_k/\bar{r}_{t+1}}.$$

Consequentially,

$$f(\hat{x}_t) - f(x_\star) \le \frac{4b^2 s_t}{\left(\sum_{k=0}^t \bar{r}_k / \bar{r}_{t+1}\right)^2}.$$

In either case, we obtain that

$$f(\hat{x}_t) - f(x_\star) \le \frac{4b^2(s_t + h_t)}{\left(\sum_{k=0}^t \bar{r}_k / \bar{r}_{t+1}\right)^2}.$$

If  $\kappa_t < t$ : We assume by induction that

$$f(\hat{x}_{\kappa_t}) - f(x_\star) \le \frac{4b^2(s_{\kappa_t} + h_{\kappa_t})}{\left(\sum_{k=0}^{\kappa_t} \bar{r}_k / \bar{r}_{\kappa_t + 1}\right)^2}.$$

Therefore,

$$\alpha_{\kappa_t} \sqrt{s_t} \sqrt{f(\hat{x}_{\kappa_t}) - f(x_\star)} \le 2b \sqrt{s_t} \sqrt{s_{\kappa_t} + h_{\kappa_t}} \frac{\alpha_{\kappa_t}}{\sum_{k=0}^{\kappa_t} \bar{r}_k / \bar{r}_{\kappa_t + 1}}$$

$$\le 2b \frac{\bar{r}_{t+1}}{\bar{r}_t} \sqrt{s_t} \sqrt{s_t + h_t}$$

$$\le 2b^2 (s_t + h_t).$$

Thus,

$$f(\hat{x}_t) - f(x_\star) \le \frac{2b^2(s_t + h_t) + h_t}{\left(\sum_{k=0}^t \bar{r}_k / \bar{r}_{t+1}\right)^2}$$
$$\le \frac{4b^2(s_t + h_t)}{\left(\sum_{k=0}^t \bar{r}_k / \bar{r}_{t+1}\right)^2}.$$

**Finalizing the induction:** For t = 0 we have  $\kappa_t = 0 = t$ . For the case  $\kappa_t = t$  we did not use the induction assumption, and therefore we have the base of the induction:

$$f(\hat{x}_0) - f(x_\star) \le \frac{4b^2(s_0 + h_0)}{(\bar{r}_0/\bar{r}_1)^2}.$$

Thus, by induction we get that for all  $t \in \{0, 1, 2, \dots, T-1\}$ ,

$$f(\hat{x}_t) - f(x_\star) \le \frac{4b^2(s_t + h_t)}{\left(\sum_{k=0}^t \bar{r}_k / \bar{r}_{t+1}\right)^2}.$$

### C.3. General regret bound (Lemma 16)

The following lemma is inspired by the regret analysis of UNIXGRAD [30].

**Lemma 16** Using Algorithm 1, eq. (2) and eq. (3), for any  $t \ge 0$ ,  $\rho_t > 0$ , we have that

$$\bar{r}_{t}\alpha_{t} \langle g_{t}, x_{t+1} - x_{\star} \rangle \leq \frac{\bar{r}_{t}^{2}\alpha_{t}^{2}\rho_{t}}{2} \|g_{t} - m_{t}\|^{2} - \frac{1}{2\rho_{t}} \|x_{t+1} - y_{t}\|^{2} + \left(\frac{1}{2\rho_{t}} - \frac{1}{2\tilde{\eta}_{x,t}}\right) (\|x_{t+1} - y_{t}\|^{2} + \|x_{t+1} - y_{t+1}\|^{2}) + \frac{1}{2\tilde{\eta}_{y,t}} (\|y_{t} - x_{\star}\|^{2} - \|y_{t+1} - x_{\star}\|^{2}).$$

#### **Proof** We have

$$\bar{r}_t \alpha_t \langle g_t, x_{t+1} - x_{\star} \rangle 
= \bar{r}_t \alpha_t \langle g_t - m_t, x_{t+1} - y_{t+1} \rangle + \bar{r}_t \alpha_t \langle m_t, x_{t+1} - y_{t+1} \rangle + \bar{r}_t \alpha_t \langle g_t, y_{t+1} - x_{\star} \rangle.$$
(32)

In addition

$$\bar{r}_{t}\alpha_{t} \langle g_{t} - m_{t}, x_{t+1} - y_{t+1} \rangle \stackrel{(i)}{\leq} \bar{r}_{t}\alpha_{t} \|g_{t} - m_{t}\| \|x_{t+1} - y_{t+1}\|$$

$$\stackrel{(ii)}{\leq} \frac{\rho_{t}\bar{r}_{t}\alpha_{t}^{2}}{2} \|g_{t} - m_{t}\|^{2} + \frac{1}{2\rho_{t}} \|x_{t+1} - y_{t+1}\|,$$

$$(33)$$

where (i) is from Holder's Inequality and (ii) is due to Young's Inequality.

For the Euclidean Bregman divergence  $\mathcal{D}_{\mathcal{R}}(x,y) = \frac{1}{2}\|x-y\|^{\overline{2}}$  we have that the update rule  $x_{t+1} = \operatorname{Proj}_{\mathcal{K}}(y_t - \alpha_t \eta_{x,t} m_t) = \operatorname{Proj}_{\mathcal{K}}(y_t - \bar{r}_t \alpha_t \tilde{\eta}_{x,t} m_t)$  is equivalent to the update rule  $x_{t+1} = \arg\min_{x \in \mathcal{K}} \left\{ \bar{r}_t \alpha_t \left\langle x, m_t \right\rangle + \frac{1}{\tilde{\eta}_{x,t}} \mathcal{D}_{\mathcal{R}}(x,y_t) \right\}$ . Therefore, from the optimality of  $x_{t+1}$  we get

$$\bar{r}_{t}\alpha_{t} \langle m_{t}, x_{t+1} - y_{t+1} \rangle \leq \frac{1}{\tilde{\eta}_{x,t}} \langle \nabla_{x} \mathcal{D}_{\mathcal{R}}(x_{t+1}, y_{t}), x_{t+1} - y_{t+1} \rangle 
= \frac{1}{\tilde{\eta}_{x,t}} (\mathcal{D}_{\mathcal{R}}(y_{t+1}, y_{t}) - \mathcal{D}_{\mathcal{R}}(x_{t+1}, y_{t}) - \mathcal{D}_{\mathcal{R}}(y_{t+1}, x_{t+1})).$$
(34)

Similarly,  $y_{t+1} = \arg\min_{y \in \mathcal{K}} \left\{ \bar{r}_t \alpha_t \langle y, g_t \rangle + \frac{1}{\tilde{\eta}_{y,t}} \mathcal{D}_{\mathcal{R}}(y, y_t) \right\}$ . Therefore, from the optimality of  $y_{t+1}$  we get

$$\bar{r}_{t}\alpha_{t} \langle g_{t}, y_{t+1} - x_{\star} \rangle \leq \frac{1}{\tilde{\eta}_{y,t}} \langle \nabla_{x} \mathcal{D}_{\mathcal{R}}(y_{t+1}, y_{t}), x_{\star} - y_{t+1} \rangle$$

$$= \frac{1}{\tilde{\eta}_{y,t}} (\mathcal{D}_{\mathcal{R}}(x_{\star}, y_{t}) - \mathcal{D}_{\mathcal{R}}(y_{t+1}, y_{t}) - \mathcal{D}_{\mathcal{R}}(x_{\star}, y_{t+1})). \tag{35}$$

By combining eqs. (34), (35) and (33) into eq. (32) we obtain that

$$\begin{split} \bar{r}_{t}\alpha_{t}\left\langle g_{t},x_{t+1}-x_{\star}\right\rangle &\leq \frac{\bar{r}_{t}^{2}\alpha_{t}^{2}\rho_{t}}{2}\|g_{t}-m_{t}\|^{2}+\frac{1}{2\rho_{t}}\|x_{t+1}-y_{t+1}\|^{2}\\ &+\frac{1}{2\tilde{\eta}_{x,t}}\left(\|y_{t+1}-y_{t}\|^{2}-\|x_{t+1}-y_{t}\|^{2}-\|y_{t+1}-x_{t+1}\|^{2}\right)\\ &+\frac{1}{2\tilde{\eta}_{y,t}}\left(\|x_{\star}-y_{t}\|^{2}-\|y_{t+1}-y_{t}\|^{2}-\|x_{\star}-y_{t+1}\|^{2}\right)\\ &=\frac{\bar{r}_{t}^{2}\alpha_{t}^{2}\rho_{t}}{2}\|g_{t}-m_{t}\|^{2}-\frac{1}{2\rho_{t}}\|x_{t+1}-y_{t}\|^{2}\\ &+\left(\frac{1}{2\rho_{t}}-\frac{1}{2\tilde{\eta}_{x,t}}\right)\left(\|x_{t+1}-y_{t}\|^{2}+\|x_{t+1}-y_{t+1}\|^{2}\right)\\ &+\frac{1}{2\tilde{\eta}_{y,t}}\left(\|x_{\star}-y_{t}\|^{2}-\|x_{\star}-y_{t+1}\|^{2}\right)+\left(\frac{1}{2\tilde{\eta}_{x,t}}-\frac{1}{2\tilde{\eta}_{y,t}}\right)\|y_{t+1}-y_{t}\|^{2}. \end{split}$$

Since  $\tilde{\eta}_{y,t} \leq \tilde{\eta}_{x,t}$ , we may drop the final term in the above display, completing the proof.

## Appendix D. Iterate stability lemmas

## D.1. A weighted regret bound (Lemma 17)

**Lemma 17** For any sequence of positive numbers  $\omega_0, \omega_1, \omega_2, \ldots$ , define

$$\hat{x}_t := \frac{\sum_{k=0}^t \omega_k x_{k+1}}{\sum_{k=0}^t \omega_k}.$$

Let  $\tilde{\eta}_0, \tilde{\eta}_1, \tilde{\eta}_2, \dots$  be a non-increasing sequence of positive numbers. We have that for any T > 0,

$$\sum_{t=0}^{T-1} \omega_t \tilde{\eta}_t \left\langle \nabla f(\hat{x}_t), x_{t+1} - x_{\star} \right\rangle \ge 0.$$

**Proof** Define

$$\tilde{f}(x) = f(x) - f(x_{\star}).$$

We start from eq. (31) inside the proof of Lemma 14, which says that for all  $t \ge 0$ 

$$\omega_t \langle \nabla f(\hat{x}_t), x_{t+1} - x_{\star} \rangle \ge \omega_t (f(\hat{x}_t) - f(x_{\star})) + \sum_{k=0}^{t-1} \omega_k (f(\hat{x}_t) - f(\hat{x}_{t-1})).$$

Multiplying each side by  $\tilde{\eta}_t$  and summing, we obtain

$$\sum_{t=0}^{T-1} \omega_{t} \tilde{\eta}_{t} \left\langle \nabla f(\hat{x}_{t}), x_{t+1} - x_{\star} \right\rangle \geq \sum_{t=0}^{T-1} \omega_{t} \tilde{\eta}_{t} (f(\hat{x}_{t}) - f(x_{\star})) + \sum_{t=0}^{T-1} \sum_{k=0}^{t-1} \omega_{k} \tilde{\eta}_{t} (f(\hat{x}_{t}) - f(\hat{x}_{t-1}))$$

$$= \sum_{t=0}^{T-1} \omega_{t} \tilde{\eta}_{t} \tilde{f}(\hat{x}_{t}) + \sum_{t=0}^{T-1} \sum_{k=0}^{t-1} \omega_{k} \tilde{\eta}_{t} \left( \tilde{f}(\hat{x}_{t}) - \tilde{f}(\hat{x}_{t-1}) \right)$$

$$\stackrel{(\star)}{\geq} \sum_{t=0}^{T-1} \omega_{t} \tilde{\eta}_{t} \tilde{f}(\hat{x}_{t}) + \sum_{t=0}^{T-1} \sum_{k=0}^{t-1} \omega_{k} \left( \tilde{\eta}_{t} \tilde{f}(\hat{x}_{t}) - \tilde{\eta}_{t-1} \tilde{f}(\hat{x}_{t-1}) \right)$$

$$= \sum_{t=0}^{T-1} \omega_{t} \tilde{\eta}_{t} \tilde{f}(\hat{x}_{t}) + \sum_{k=0}^{T-2} \sum_{t=k+1}^{T-1} \omega_{k} \left( \tilde{\eta}_{t} \tilde{f}(\hat{x}_{t}) - \tilde{\eta}_{t-1} \tilde{f}(\hat{x}_{t-1}) \right),$$

where  $(\star)$  is because that  $\tilde{f}(\hat{x}_{t-1}) \geq 0$  and  $\tilde{\eta}_{t-1} \geq \tilde{\eta}_t > 0$ .

We can now perform a telescopic summation and obtain

$$\sum_{t=0}^{T-1} \omega_t \tilde{\eta}_t \left\langle \nabla f(\hat{x}_t), x_{t+1} - x_{\star} \right\rangle \ge \sum_{t=0}^{T-1} \omega_t \tilde{\eta}_t \tilde{f}(\hat{x}_t) + \sum_{t=0}^{T-2} \omega_t \left( \tilde{\eta}_{T-1} \tilde{f}(\hat{x}_{T-1}) - \tilde{\eta}_t \tilde{f}(\hat{x}_t) \right) \\
= \omega_{T-1} \tilde{\eta}_{T-1} \tilde{f}(\hat{x}_{T-1}) + \sum_{t=0}^{T-2} \omega_t \left( \tilde{\eta}_t \tilde{f}(\hat{x}_t) + \tilde{\eta}_{T-1} \tilde{f}(\hat{x}_{T-1}) - \tilde{\eta}_t \tilde{f}(\hat{x}_t) \right) \\
= \omega_{T-1} \tilde{\eta}_{T-1} \tilde{f}(\hat{x}_{T-1}) + \sum_{t=1}^{T-1} \omega_t \tilde{\eta}_{T-1} \tilde{f}(\hat{x}_{T-1}).$$

Thus, because  $\tilde{f}(\hat{x}_{T-1}) \geq 0$ , we obtain that

$$\sum_{t=0}^{T-1} \omega_t \tilde{\eta}_t \langle \nabla f(\hat{x}_t), x_{t+1} - x_{\star} \rangle \ge 0.$$

## D.2. Inductive stability bound (Lemma 18)

**Lemma 18** If  $r_{\epsilon} = r_0 \le d_0$ , and for all  $t \ge 1$  we have that

$$||y_t - x_t|| \le \frac{\bar{r}_{t-1}}{4} \text{ and}$$

$$d_t^2 \le \left(d_0 + \frac{1}{4}\bar{r}_{t-1}\right)^2,$$

then for all  $t \ge 0$  we get that

$$d_t \leq 2d_0$$
 and  $r_t \leq 4d_0$ .

**Proof** We prove this lemma by induction. The basis of the induction is that for t=0 we get that  $d_0 \le 2d_0$  and  $r_0 \le d_0 \le 4d_0$ .

For any  $t \ge 1$ , we assume that  $\bar{d}_{t-1} \le 2d_0$  and  $\bar{r}_{t-1} \le 4d_0$ . Thus,

$$d_t \le d_0 + \frac{1}{4}\bar{r}_{t-1} \le 2d_0.$$

Also,

$$||y_t - x_0|| < ||y_t - x_\star|| + ||x_0 - x_\star|| = d_t + d_0 < 3d_0.$$

In addition,

$$||x_{t} - x_{0}|| \leq ||y_{t} - x_{0}|| + ||x_{t} - y_{t}||$$

$$\stackrel{(\star)}{\leq} 3d_{0} + \frac{\bar{r}_{t-1}}{4}$$

$$\leq 4d_{0}.$$

where  $(\star)$  is because  $||x_t - y_t|| \le \frac{\bar{r}_{t-1}}{4}$ . As a result,

$$d_t \leq 2d_0$$
 and  $r_t \leq 4d_0$ .

Finally, by induction, we get that for all  $t \ge 0$ 

$$d_t \leq 2d_0$$
 and  $r_t \leq 4d_0$ .

## D.3. Single-step iterate stability (Lemma 19)

**Lemma 19** Let c be a positive number. Using Algorithm l, for any  $t \geq 0$ , if  $\eta_{x,t} \leq \frac{\bar{r}_t}{c\alpha_t \|m_t\|}$ ,  $\eta_{y,t} \leq \frac{\bar{r}_t}{c\alpha_t \|g_t - m_t\|}$  and  $\eta_{y,t} \leq \eta_{x,t}$  then

$$||x_{t+1} - y_t|| \le \frac{r_t}{c}$$

$$||y_{t+1} - y_t|| \le \frac{2\bar{r}_t}{c}$$

$$||x_{t+1} - y_{t+1}|| \le \frac{2\bar{r}_t}{c}$$

$$\bar{r}_{t+1} \le \bar{r}_t \left(1 + \frac{2}{c}\right).$$

**Proof** First, by definition of the iterates and the fact that K is convex (and projection onto a closed convex set is nonexpansive) we have

$$||x_{t+1} - y_t|| = ||\operatorname{Proj}_{\mathcal{K}}(y_t - \alpha_t \eta_{x,t} m_t) - y_t|| \le \alpha_t \eta_{x,t} ||m_t|| \le \frac{\bar{r}_t}{c}.$$
 (36)

Second, by definition of the iterates and the fact that K is convex, we also have

$$||y_{t+1} - y_t|| = ||\operatorname{Proj}_{\mathcal{K}}(y_t - \alpha_t \eta_{y,t} g_t) - y_t|| \le \alpha_t \eta_{y,t} ||g_t||$$

$$\le \alpha_t \eta_{y,t} ||g_t - m_t|| + \alpha_t \eta_{y,t} ||m_t|| \le \frac{2\bar{r}_t}{\epsilon}.$$
(37)

Third, by definition of the iterates, the fact that K is convex, the fact  $\eta_{y,t} \leq \eta_{x,t}$ , and the assumed upper bounds on  $\eta_{y,t}$  and  $\eta_{x,t}$  in the premise of this lemma we have

$$||x_{t+1} - y_{t+1}|| = ||\operatorname{Proj}_{\mathcal{K}}(y_t - \alpha_t \eta_{x,t} m_t) - \operatorname{Proj}_{\mathcal{K}}(y_t - \alpha_t \eta_{y,t} g_t)||$$

$$\leq \alpha_t ||\eta_{x,t} m_t - \eta_{y,t} g_t|| \leq \alpha_t \eta_{y,t} ||g_t - m_t|| + \alpha_t (\eta_{x,t} - \eta_{y,t}) ||m_t||$$

$$\leq \alpha_t \eta_{y,t} ||g_t - m_t|| + \alpha_t \eta_{x,t} ||m_t|| \leq \frac{2\bar{r}_t}{c}.$$

Finally,

$$r_{t+1} \le r_t + \max(\|x_{t+1} - y_t\|, \|y_{t+1} - y_t\|).$$

Therefore, using eq. (36) and eq. (37) we obtain

$$\bar{r}_{t+1} = \max(\bar{r}_t, r_{t+1}) \le \bar{r}_t + \max(\|x_{t+1} - y_t\|, \|y_{t+1} - y_t\|) \le \bar{r}_t \left(1 + \frac{2}{c}\right).$$

# Appendix E. Concentration bounds

## E.1. An empirical-Bernstein-type time uniform concentration bound (Lemma 20)

**Lemma 20 (From Ivgi et al. [27])** Let S be the set of nonnegative and nondecreasing sequences. Let  $C_t \in \mathcal{F}_{t-1}$  and let  $X_t$  be a martingale difference sequence adapted to  $\mathcal{F}_t$  such that  $|X_t| \leq C_t$  with probability I for all t. Then, for all  $\delta \in (0,1)$ , c > 0, and  $\hat{X}_t \in \mathcal{F}_{t-1}$  such that  $|\hat{X}_t| \leq C_t$  with probability I,

$$\mathbb{P}\left(\exists t \leq T, \exists \{y_i\}_{i=1}^{\infty} \in S : \left| \sum_{i=1}^{t} y_i X_i \right| \geq 8y_t \sqrt{\theta_{t,\delta} \sum_{i=1}^{t} \left( X_i - \hat{X}_i \right)^2 + c^2 \theta_{t,\delta}^2} \right) \leq \delta + \mathbb{P}(\exists t \leq T : C_t > c).$$

### E.2. Concentration bound for suboptimally proof (Lemma 21)

**Lemma 21** Let  $\mathfrak{B} > 0$  and  $\delta \in (0,1)$ . In the bounded noise setting (Assumption 6), using Algorithm 1 and eq. (12), with probability of at least  $1 - \delta - \mathbb{P}[\bar{\mathfrak{b}}_{T-1} > \mathfrak{B}]$  we get that for all  $t \in \{0,1,\ldots,T-1\}$  then

$$\left| \sum_{k=0}^{t} \bar{r}_{t} \alpha_{k} \left\langle \nabla f(\hat{x}_{k}) - g_{k}, x_{k+1} - x_{\star} \right\rangle \right| \leq 8\alpha_{t} \bar{r}_{t} (\bar{r}_{t+1} + d_{0}) \sqrt{\theta_{t+1,\delta} \sum_{k=0}^{t} \|\nabla f(\hat{x}_{k}) - g_{k}\|^{2} + (\theta_{t+1,\delta} \mathfrak{B})^{2}}.$$

**Proof** For  $k \in \{0, 1, \dots, T\}$  define

$$\tilde{d}_k = \max_{i \le k} ||x_k - x_\star||.$$

For  $k \in \{0, 1, \dots, T-1\}$  define the random variables:

$$Y_k = \alpha_k \bar{r}_k \tilde{d}_{k+1}$$
, and  $X_k = \left\langle \nabla f(\hat{x}_k) - g_k, \frac{x_{k+1} - x_{\star}}{\tilde{d}_{k+1}} \right\rangle$ .

From these definitions we get

$$\sum_{k=0}^{t} Y_k X_k = \sum_{k=0}^{t} \bar{r}_t \alpha_k \left\langle \nabla f(\hat{x}_k) - g_k, x_{k+1} - x_{\star} \right\rangle,$$

and that  $\{Y_k\}_{k=0}^{T-1}$  is a non-decreasing sequence of non-negative numbers. In addition, as  $x_{k+1}$  and  $\tilde{d}_{k+1}$  are independent of the noise of  $g_k$  then  $X_k$  is a martingale difference sequence. Therefore, as  $|X_k| \leq \bar{\mathfrak{b}}_k$  with probability of 1, Lemma 20 gives us that

$$\mathbb{P}\left(\exists t < T : \left| \sum_{k=0}^{t} Y_k X_k \right| \ge 8Y_t \sqrt{\theta_{t+1,\delta} \sum_{k=0}^{t} (X_k - 0)^2 + (\theta_{t+1,\delta} \mathfrak{B})^2} \right) \le \delta + \mathbb{P}\left[\bar{\mathfrak{b}}_{T-1} > \mathfrak{B}\right].$$

Therefore, by using the Cauchy–Schwarz inequality, we obtain that, with a probability of at least  $1 - \delta - \mathbb{P}[\bar{\mathfrak{b}}_{T-1} > \mathfrak{B}]$ , for all  $t \in \{0, 1, \dots, T-1\}$ 

$$\left| \sum_{k=0}^{t} \bar{r}_{t} \alpha_{k} \left\langle \nabla f(\hat{x}_{k}) - g_{k}, x_{k+1} - x_{\star} \right\rangle \right| \leq 8\alpha_{t} \bar{r}_{t} \tilde{d}_{t+1} \sqrt{\theta_{t+1,\delta} \sum_{k=0}^{t} \|\nabla f(\hat{x}_{k}) - g_{k}\|^{2} + (\theta_{t+1,\delta} \mathfrak{B})^{2}}$$

$$\leq 8\alpha_{t} \bar{r}_{t} \left( \bar{r}_{t+1} + \tilde{d}_{0} \right) \sqrt{\theta_{t+1,\delta} \sum_{k=0}^{t} \|\nabla f(\hat{x}_{k}) - g_{k}\|^{2} + (\theta_{t+1,\delta} \mathfrak{B})^{2}}.$$

Thus,

$$\left| \sum_{k=0}^{t} \bar{r}_{t} \alpha_{k} \left\langle \nabla f(\hat{x}_{k}) - g_{k}, x_{k+1} - x_{\star} \right\rangle \right| \leq 8\alpha_{t} \bar{r}_{t} (\bar{r}_{t+1} + d_{0}) \sqrt{\theta_{t+1,\delta} \sum_{k=0}^{t} \|\nabla f(\hat{x}_{k}) - g_{k}\|^{2} + (\theta_{t+1,\delta} \mathfrak{B})^{2}}.$$

### E.3. Concentration bound for iterate stability proof (Lemma 22)

**Lemma 22** Let  $\tilde{\eta}_{y,t}$  be such that, for some c, s > 0 we have

$$\frac{1}{\tilde{\eta}_{y,t}} \ge c \max \left\{ \sqrt{s + Q_t} \log_+ \left( \frac{s + Q_t}{s} \right), \alpha_t \|\nabla f(\hat{x}_t) - m_t\|, \alpha_t \|\nabla f(\hat{x}_t) - g_t\| \right\}.$$

If for all  $t \ge 0$  we have that  $\eta_{y,t} = \bar{r}_t \tilde{\eta}_{y,t}$  is independent of  $g_t$  given  $x_0, \dots, x_t$ , then, with probability of at least  $1 - \delta$ , for all  $t \ge 0$ ,

$$\left| \sum_{k=0}^{t} \alpha_k \eta_{y,k} \left\langle g_k - \nabla f(\hat{x}_k), x_{k+1} - x_{\star} \right\rangle \right| \leq \frac{12\theta_{t+1,\delta}}{c} \bar{r}_t (\bar{r}_{t+1} + d_0).$$

**Proof** For  $t \in \{0, 1, \dots, T\}$  define

$$\tilde{d}_t = \max_{k < t} ||x_t - x_\star||.$$

For  $t \in \{0, 1, ..., T - 1\}$  define

$$X_t = lpha_t ilde{\eta}_{y,t} \left\langle g_t - 
abla f(\hat{x}_t), rac{x_{t+1} - x_\star}{ ilde{d}_{t+1}} 
ight
angle \; ,$$
 $\hat{X}_t = lpha_t ilde{\eta}_{y,t} \left\langle 
abla f(\hat{x}_t) - m_t, rac{x_{t+1} - x_\star}{ ilde{d}_{t+1}} 
ight
angle \; ext{ and }$ 
 $Y_t = ar{r}_t ilde{d}_{t+1}.$ 

The assumption  $\frac{1}{\tilde{\eta}_{y,t}} \ge c\alpha_t \max\{\|\nabla f(\hat{x}_t) - m_t\|, \|g_t - \nabla f(\hat{x}_t)\|\}$  implies that  $\max\{|X_t|, |\hat{X}_t|\} \le \frac{1}{c}$ . In addition, as  $m_t$ ,  $x_{t+1}$  and  $\tilde{d}_{t+1}$  are independent of the noise of  $g_t$  then  $\hat{X}_t$  is independent of the noise of  $g_t$  and  $X_t$  is a martingale difference sequence. Thus, Lemma 20 gives us that

$$\mathbb{P}\left(\forall t \in \{0, 1, \ldots\} : \left| \sum_{k=0}^{t} Y_k X_k \right| < 8\bar{r}_t \tilde{d}_{t+1} \sqrt{\theta_{t+1, \delta} \sum_{k=0}^{t} \left( X_k - \hat{X}_k \right)^2 + \frac{1}{c^2} \theta_{t+1, \delta}^2} \right) \ge 1 - \delta.$$

Furthermore, we have

$$\sum_{k=0}^{t} \left( X_{t} - \hat{X}_{t} \right)^{2} = \sum_{k=0}^{t} \left( \alpha_{k} \tilde{\eta}_{y,k} \left\langle g_{k} - m_{k}, \frac{x_{k+1} - x_{\star}}{\tilde{d}_{t+1}} \right\rangle \right)^{2} \leq \sum_{k=0}^{t} \alpha_{t}^{2} \tilde{\eta}_{y,k}^{2} \|g_{k} - m_{k}\|^{2}$$

$$\stackrel{(i)}{\leq} \frac{1}{c^{2}} \sum_{k=0}^{t} \frac{\alpha_{t}^{2} \|g_{k} - m_{k}\|^{2}}{\left( s + \sum_{k=0}^{t} \alpha_{k}^{2} \|g_{k} - m_{k}\|^{2} \right) \log_{+}^{2} \left( \frac{s + \sum_{k=0}^{t} \alpha_{k}^{2} \|g_{k} - m_{k}\|^{2}}{s} \right)} \stackrel{(ii)}{\leq} \frac{1}{c^{2}},$$

where (i) follows from the assumption that  $\frac{1}{\tilde{\eta}_{y,t}} \geq c\sqrt{s+Q_t}\log_+\left(\frac{s+Q_t}{s}\right)$  and the definition of  $Q_t$ , and (ii) is a direct result of Lemma 30 with  $a_k = s + \sum_{k=0}^t \alpha_k^2 \|g_k - m_k\|^2$ . In addition, we have that

$$Y_t X_t = \alpha_t \eta_{y,t} \langle g_t - \nabla f(\hat{x}_t), x_{t+1} - x_{\star} \rangle.$$

Therefore, with probability of at least  $1 - \delta$ , for all  $t \ge 0$  we have that

$$\left| \sum_{k=0}^{t} \alpha_k \eta_{y,k} \left\langle g_k - \nabla f(\hat{x}_k), x_{k+1} - x_{\star} \right\rangle \right| \leq 8\bar{r}_t \tilde{d}_{t+1} \sqrt{\frac{\theta_{t+1,\delta}}{c^2} + \frac{\theta_{t+1,\delta}^2}{c^2}}$$

$$\leq \frac{12\theta_{t+1,\delta}}{c} \bar{r}_t \Big( \bar{r}_{t+1} + \tilde{d}_0 \Big)$$

$$\leq \frac{12\theta_{t+1,\delta}}{c} \bar{r}_t (\bar{r}_{t+1} + d_0).$$

# **E.4.** Relating $\bar{Q}_t$ to $Q_t$ (Lemma 23)

**Lemma 23** Let  $\mathfrak{B} > 0$  and  $\delta \in (0,1)$ . In the bounded noise setting (Assumption 6), using Algorithm 1 and the step sizes (14), with probability of at least  $1 - \delta - \mathbb{P}[\bar{\mathfrak{b}}_{T-1} > \mathfrak{B}]$  we get that, for all  $t \in \{0,1,\ldots,T-1\}$ ,

$$\bar{Q}_t \leq 5Q_t + 80(t+1)^3 \sqrt{\theta_{t+1,\delta}} V_t + 2(t+1)^2 \theta_{t+1,\delta} \mathfrak{B}^2.$$

**Proof** For all  $k \ge 0$  we have

$$\|\tilde{g}_k - m_k\|^2 \le 2\|g_k - m_k\|^2 + 2\|g_k - \tilde{g}_t\|^2$$

$$\le 2\|g_k - m_k\|^2 + 4\|g_k - \nabla f(\hat{x}_k)\|^2 + 4\|\tilde{g}_k - \nabla f(\hat{x}_k)\|^2.$$

Therefore, since  $\alpha_k \leq k+1$ ,

$$\sum_{k=0}^{t} \alpha_k^2 \|\tilde{g}_k - m_k\|^2 \le 2 \sum_{k=0}^{t} \alpha_k^2 \|g_k - m_k\|^2 + 8 \sum_{k=0}^{t} (k+1)^2 \|g_k - \nabla f(\hat{x}_k)\|^2 
+ 4 \sum_{k=0}^{t} (k+1)^2 (\|\tilde{g}_k - \nabla f(\hat{x}_k)\|^2 - \|g_k - \nabla f(\hat{x}_k)\|^2).$$
(38)

We now bound 
$$\sum_{k=0}^t (k+1)^2 (\|\tilde{g}_k - \nabla f(\hat{x}_k)\|^2 - \|g_k - \nabla f(\hat{x}_k)\|^2)$$
. Define  $X_t = (\|\tilde{g}_t - \nabla f(\hat{x}_t)\|^2 - \|g_t - \nabla f(\hat{x}_t)\|^2)$ ,  $\hat{X}_t = \|\tilde{g}_t - \nabla f(\hat{x}_t)\|^2$  and  $Y_t = (t+1)^2$ .

We have that for all  $t \geq 0$  then  $|X_t| \leq \bar{\mathfrak{b}}_t^2$  and  $|\hat{X}_t| \leq \bar{\mathfrak{b}}_t^2$  with probability 1. Therefore, Lemma 20 gives us that

$$\mathbb{P}\left(\forall t \in \{0, 1, \dots, T - 1\} : \left| \sum_{k=0}^{t} Y_k X_k \right| < 8Y_t \sqrt{\theta_{t+1,\delta} \sum_{k=0}^{t} \left(X_k - \hat{X}_k\right)^2 + \theta_{t+1,\delta}^2 \mathfrak{B}^4} \right) \\
\geq 1 - \delta - \mathbb{P}(\bar{\mathfrak{b}}_{T-1} > \mathfrak{B}).$$

Consequentially, by combining this result with eq. (38), we get that with probability at least  $1 - \delta - \mathbb{P}(\bar{\mathfrak{b}}_{T-1} > \mathfrak{B})$  that for all  $t \in \{0, 1, \dots, T-1\}$  we have that

$$\sum_{k=0}^{t} \alpha_k^2 \|\tilde{g}_k - m_k\|^2 \le 2 \sum_{k=0}^{t} \alpha_k^2 \|g_k - m_k\|^2 + 40(t+1)^2 \sqrt{\theta_{t+1,\delta}} \sum_{k=0}^{t} \|g_k - \nabla f(\hat{x}_k)\|^2 + (t+1)^2 \theta_{t+1,\delta} \mathfrak{B}^2.$$

Substituting into the above equation the definition of  $Q_t$  and  $V_t$  given in eq. (1) and eq. (11), respectively, and recalling the definition of  $\bar{Q}_t$  given in eq. (13)

$$\bar{Q}_t = \sum_{k=0}^t \alpha_k^2 \max\{\|g_k - m_k\|^2, 2\|\tilde{g}_k - m_k\|^2\} \le Q_t + 2\sum_{k=0}^t \alpha_k^2 \|\tilde{g}_k - m_k\|^2$$

completes the proof.

## E.5. Concentration inequality for bounded random vectors (Lemma 24)

**Lemma 24 (Howard et al. [25])** For  $T \in \mathbb{N}$ , let  $\{U_t\}_{t \in [T]}$  be a sequence of mean zero random vectors in  $\mathbb{R}^d$  with  $||U_t|| \le c$  almost surely. Then

$$\mathbb{P}\left(\left\|\sum_{t=1}^{T} U_{t}\right\| \ge x\right) \le 2\exp\left(-\frac{x^{2}}{2c^{2}T}\right).$$

**Proof** This result follows from Howard et al. [25, Corollary 10.a] with  $Y_t = \sum_{k=1}^t U_k$ ,  $\Psi(\cdot) = \|\cdot\|$ ,  $c_t = c$  and  $m = c^2 T$ . The selection of  $\Psi(\cdot) = \|\cdot\|$  yields  $D_\star = 1$  (see discussion preceding [25, Corollary 10.a]). Setting  $c_t = c$  yields  $V_t = c^2 t$ . Hence  $\frac{D_\star^2}{2m}(V_T - m) \leq 0$  and Howard et al. [25, eq. (4.28)] gives the desired result.

# Appendix F. Auxiliary lemmas

# **F.1.** The growth rate of $\sum_k \bar{r}_k \alpha_k$ (Lemma 25)

We note that in accelerated optimization algorithms we normally have that  $\alpha_t = \Theta(t)$ . Even though this is not the case for U-DoG,  $\alpha_t$  is roughly similar to t. First, it is easy to see that  $1 \le \alpha_t \le t$ . Secondly, the running sum of  $\bar{r}_t \alpha_t$  grows roughly quadratically. This is shown in the following lemma, in which we replace  $\alpha_t$  and  $\bar{r}_t$  with  $a_t$  and  $s_t$ , respectively

**Lemma 25** Let  $s_0, s_1, \ldots, s_t$  be a non-decreasing sequence of positive numbers. Define  $a_k := \sum_{i=0}^k \frac{s_i}{s_k}$ , then

$$s_t a_t^2 \le 2 \sum_{k=0}^t s_k a_k.$$

Proof We have

$$\sum_{k=0}^{t} s_k a_k = \sum_{k=0}^{t} \sum_{i=0}^{k} s_i = \sum_{k=0}^{t} (t - k + 1) s_k.$$

And,

$$s_t a_t^2 = \frac{1}{s_t} \sum_{k=0}^t \sum_{i=0}^t s_k s_i = \frac{2}{s_t} \sum_{k=0}^t s_k \sum_{i=k}^t s_i - \frac{1}{s_t} \sum_{k=0}^t s_k^2 \le 2 \sum_{k=0}^t s_k \sum_{i=k}^t \frac{s_i}{s_t} \le 2 \sum_{k=0}^t (t-k+1) s_k.$$

Thus,

$$s_t a_t^2 \le 2 \sum_{k=0}^t s_k a_k.$$

## F.2. Discrete derivative lemma (Lemma 26)

**Lemma 26** Let c be a positive number, and let  $s_0, s_1, s_2, \ldots$  be a sequence of positive numbers. For every  $t \ge 0$  define

$$\rho_t = \frac{1}{c\sqrt{\sum_{k=0}^t s_k}}.$$

We have that for every  $t \geq 0$ 

$$\frac{1}{\rho_{t+1}} - \frac{1}{\rho_t} \le c^2 \rho_{t+1} s_{t+1}.$$

**Proof** For every  $t \geq 0$  we have that

$$s_{t+1} = \sum_{k=0}^{t+1} s_k - \sum_{k=0}^{t} s_k \ge \sqrt{\sum_{k=0}^{t+1} s_k} \left( \sqrt{\sum_{k=0}^{t+1} s_k} - \sqrt{\sum_{k=0}^{t} s_k} \right) = \frac{1}{c^2 \rho_{t+1}} \left( \frac{1}{\rho_{t+1}} - \frac{1}{\rho_t} \right).$$

Thus,

$$\frac{1}{\rho_{t+1}} - \frac{1}{\rho_t} \le c^2 \rho_{t+1} s_{t+1}.$$

### F.3. Discrete integral lemma (Lemma 27)

**Lemma 27** For any positive numbers  $c_1, c_2$ , for any  $t \ge 0$ , and for any sequence of non-negative numbers  $B_0, B_1, B_2, \ldots, B_t$  we have that

$$c_1 \sqrt{\sum_{k=0}^{t} B_k^2} - \sum_{k=0}^{t} \frac{B_k^2}{c_2} \sqrt{\sum_{j=0}^{k} B_j^2} \le 2c_1^{3/2} c_2^{1/2}.$$

**Proof** Define

$$\eta_{B,k} = \frac{1}{\sqrt{\sum_{j=1}^{k} B_j^2}}.$$

Lemma 28 gives us that

$$\sqrt{\sum_{k=0}^{t} B_k^2} \le \sum_{k=0}^{t} \frac{B_k^2}{\sqrt{\sum_{j=0}^{k} B_j^2}}.$$

Therefore, we obtain

$$c_{1}\sqrt{\sum_{k=0}^{t}B_{k}^{2}} - \sum_{k=0}^{t}\frac{B_{k}^{2}}{c_{2}}\sqrt{\sum_{j=0}^{k}B_{j}^{2}} \le c_{1}\sum_{k=0}^{t}\frac{B_{k}^{2}}{\sqrt{\sum_{j=0}^{k}B_{j}^{2}}} - \sum_{k=0}^{t}\frac{B_{k}^{2}}{c_{2}}\sqrt{\sum_{j=0}^{k}B_{j}^{2}}$$
$$= \sum_{k=0}^{t}\left(c_{1}\eta_{B,k} - \frac{1}{c_{2}\eta_{B,k}}\right)B_{k}^{2}.$$

Define

$$\kappa = \max \left[ \left\{ t \in \{0, 1, \dots, t\} : 2c_1 \eta_{B, t} - \frac{1}{c_2 \eta_{B, t}} > 0 \right\} \cup \{-1\} \right].$$

We have,

$$c_1\sqrt{\sum_{k=0}^{t}B_k^2} - \sum_{k=0}^{t}\frac{B_k^2}{c_2}\sqrt{\sum_{j=0}^{k}B_j^2} \leq \sum_{k=0}^{\kappa}c_1\eta_{B,k}B_k^2 = c_1\sum_{k=0}^{\kappa}\frac{B_k^2}{\sqrt{\sum_{j=0}^{k}B_j^2}} \overset{(\star)}{\leq} 2c_1\sqrt{\sum_{k=0}^{\kappa}B_k^2} = \frac{2c_1}{\eta_{B,\kappa}}\mathbb{1}_{\{\kappa\geq 0\}},$$

where  $(\star)$  is because of Lemma 28. From the definition of  $\kappa$ , we obtain that

$$c_1 \eta_{B,\kappa} > \frac{1}{c_2 \eta_{B,\kappa}}.$$

Thus,

$$c_1 \sqrt{\sum_{k=0}^t B_k^2} - \sum_{k=0}^t \frac{B_k^2}{c_2} \sqrt{\sum_{j=0}^k B_j^2} \le \frac{2c_1}{\eta_{B,\kappa}} \mathbb{1}_{\{\kappa \ge 0\}} \le 2c_1^{3/2} c_2^{1/2}.$$

## F.4. Additional lemmas from prior work

**Lemma 28 (e.g., Levy et al. [36])** For any  $k \ge 0$  and for any sequence on non-negative numbers  $s_0, s_1, s_2, \ldots, s_k$  the following holds:

$$\sqrt{\sum_{i=0}^{k} s_i} \le \sum_{i=0}^{k} \frac{s_i}{\sqrt{\sum_{j=0}^{i} s_j}} \le 2\sqrt{\sum_{i=0}^{k} s_i}.$$

**Lemma 29 (Ivgi et al. [27, Lemma 3])** Let  $s_0, s_1, \ldots, s_T$  be a positive nondecreasing sequence. Then

$$\max_{t \le T} \sum_{i < t} \frac{s_i}{s_t} \ge \frac{1}{e} \left( \frac{T}{\log_+(s_T/s_0)} - 1 \right).$$

**Lemma 30 (Ivgi et al. [27, Lemma 6])** Let  $a_{-1}, a_0, a_1, \ldots, a_t$  be a non-decreasing sequence of non-negative numbers, then

$$\sum_{k=0}^{t} \frac{a_k - a_{k-1}}{a_k \log_+^2(a_k/a_{-1})} \le 1.$$

**Lemma 31 (Attia and Koren [3, Lemma 15])** Let X be a  $\sigma^2(x)$ -sub-Gaussian. For and  $\delta \in (0,1)$  here exist a random variable  $\bar{X}$  such that:

- 1.  $\bar{X}$  is zero-mean:  $\mathbb{E}\bar{X}=0$ .
- 2.  $\bar{X}$  is equal to X w.h.p:  $\mathbb{P}(\bar{X} = X) \geq 1 \delta$ .
- 3.  $\bar{X}$  is bounded with probability 1:  $\mathbb{P}\Big(\|\bar{X}\| = 3\sigma\sqrt{\log(4/\delta)}\Big) = 1$ .

### Appendix G. Experimental details

#### G.1. U-DoG step sizes

In the experiments, we use the following step sizes for U-DoG

$$\eta_{x,t} = \frac{\bar{r}_t}{\sqrt{\max\{Q_{t-1}, M_t\}}} \text{ and } \eta_{y,t} = \frac{\bar{r}_t}{\sqrt{\max\{Q_t, M_t\}}},$$

with  $\bar{r}_t$ ,  $Q_t$ , and  $M_t$  as defined in Section 2. This step size is similar to the choice in eq. (10), which enjoys proven stability in the noiseless case, except we replace the logarithmic factor in the denominator with 1; preliminary experiments indicated 1 was the smallest value for which the algorithm was stable in practice. This difference between practical and theoretical algorithms is analogous to the difference between DoG and its theoretically stable variant T-DoG [27]. However, we maintain the maximization with  $M_t$  in the denominator, mainly in order to ensure that  $\eta_{x,t}$  and  $\eta_{y,t}$  are not too large early in the training. As with DoG, the additional step size adjustments necessary for the stochastic setting (given in eq. (14)) do not appear to be useful in practical settings.

### G.2. ACCELEGRAD-DOG (A-DOG)

While U-DoG enjoys strong theoretical guarantees, it requires an extra-gradient computation at each step, which can be expensive in practice. To address this, we propose an alternative algorithm, A-DoG, which combines ACCELEGRAD [36] and DoG. To complete the combination we set  $\alpha_t$  in the same way as it is calculated in U-DoG (algorithm 1). A-DoG is a simple algorithm that does not require an extra-gradient computation at each step and is presented in Algorithm 2. While we do not provide theoretical guarantees for A-DoG, our experiments demonstrate its efficacy in practice. The main challenge in proving guarantees for A-DoG appears to lie in deriving a suboptimality bound akin to Proposition 3, whose proof strongly leverages U-DoG's extra-gradient structure.

# Algorithm 2: ACCELEGRAD-DOG (A-DOG)

```
Input: Initialization z^{(0)} \in \mathcal{K}, positive constant r_{\epsilon} and number of iterations T.
```

```
1 Set y_0 = x_0 = z_0 and \bar{r}_0 = r_\epsilon

2 for t = 0, 1, \dots, T - 1 do

3 x_t = \sum_{k=0}^t \bar{r}_k / \bar{r}_t

4 y_t \sim \mathcal{G}(x_{t+1})

5 y_t = \frac{\bar{r}_t}{\sqrt{\sum_{k=0}^t \alpha_k^2 ||g_k||^2}}

6 x_{t+1} = \frac{\alpha_t}{\sum_{k=0}^t \alpha_k} z_t + \left(1 - \frac{\alpha_t}{\sum_{k=0}^t \alpha_k}\right) y_t

7 y_{t+1} = x_{t+1} - \eta_t g_t

8 z_{t+1} = \Pi_{\mathcal{K}}(z_t - \alpha_t \eta_t g_t)

9 \bar{r}_{t+1} = \max\{\bar{r}_t, ||z_{t+1} - z_0||\}
```

10 end

11 **return**  $x_T$   $\triangleright$  returning  $y_T$  gives similar results in practice

## **G.3.** Convex experiments

The bulk of our experiments focus on smooth stochastic convex optimization problems, matching our theoretical assumptions.

**Multiclass logistic regression.** We experiment with multi-class logistic regression on multiple tasks from the VTAB benchmark and the LIBSVM [12] suite (a full list is given in Section G.5). For VTAB tasks we use features obtained from a pretrained ViT-B/32 [20] model (i.e., perform linear probes), and for LIBSVM tasks we use apply logistic regression directly on the features provided.

Figures 2, 4, 6, 8, 10, 12, 14 and 16 show a view of the results for different datasets analogous to Figure 1. Figures 3, 5, 7, 9, 11, 13, 15 and 17 give a complementary view by providing training curves at different batch sizes. As discussed in Section 5, we find that both U-DoG and A-DoG are competitive with well-tuned accelerated SGD (ASGD) and often significantly outperform DoG and tuned SGD. This is especially true for the training loss (for which our theory directly holds) and at large batch sizes, with A-DoG outperforming U-DoG in most cases, as both algorithms take advantage of the reduced variance in the gradient estimates to scale effectively with the batch size, as the theory suggests. In most experiments A-DoG attain and tuned ASGD attain superior convergence rate in terms of test accuracy as well as train loss; the only exception is CIFAR-100 (Figures 4 and 5, bottom rows) where the test accuracy does not closely track the train loss.

**Least-squares.** We modify the loss on a subset of the previous experiments to least squares, learned over a one-hot encoding of the features. We use features obtained from a pretrained ViT-B/32, similar to what we used for the multiclass logistic regression. We find that our algorithms perform well in this setting as well. In comparison, while SGD and ASGD can perform well when tuned correctly, they become more sensitive to the choice of step size and momentum, performing poorly when not properly tuned and sometimes diverging completely. Similar to the other experiments, the results are given in Figures 18 to 21.

**Noiseless quadratic experiments.** As a final experiment, we compare the performance of the different algorithms on the quadratic function  $f(x) = \sum_{i=1}^n \left(\frac{i}{2n}x_i^2 + x_i\right)$  with  $n=10^4$ . The results agree with the theoretical analysis, with all algorithms reaching the optimal solution or very close to it, barring GD and AGD with excessively high momentum and learning rate. Results are depicted in Figure 22.

#### **G.4.** Non-convex experiments

While we mainly focus on demonstrating the effectiveness of U-DoG and A-DoG in settings that match our theoretical analysis, we also perform preliminary experimentation in practical scenarios, namely training neural networks on datasets of moderate scales. In particular, we train a ResNet-50 [24] from scratch on a subset of the VTAB benchmark (Figures 23 to 27). Additionally, we repeat two experiments from [27]: fine-tuning a CLIP model [53] on ImageNet (Figure 28), and training a WideResnet-28-10 [66] model from scratch on CIFAR-10 (Figure 29). We observe that U-DoG often fails to converge to competitive results, while A-DoG is quite competitive with DoG on the VTAB tasks, but under-performs it for CIFAR-10 and ImageNet fine-tuning, indicating that it is not a yet a viable general-purpose neural network optimizer.

# **G.5.** Implementation details

**Environment settings.** All of our experiments were based on PyTorch [51] (version 1.12.0). For DoG and the implementation of polynomial-decay model averaging [57], we used the dog-optimizer package (version 1.0.3) [27]. For ASGD, we used the native PyTorch SGD<sup>5</sup> with the Nesterov option enabled.

VTAB experiments were based on the PyTorch Image Models (timm, version0.7.0dev0) repository [63], with TensorFlow datasets (version 4.6.0) as a dataset backend [1]. LIBSVM [12] experiments were based on the libsymdata (version 0.4.1) package.

<sup>5.</sup> https://pytorch.org/docs/stable/generated/torch.optim.SGD.html

To support the training and analysis of the results, we used numpy [23], scipy [61], pandas [62] and scikit-learn [52].

As much as possible, we leveraged existing recipes as provided by timm to train the models.

**Datasets.** The subset of datasets used in our VTAB experiments are: **CIFAR-100** [34], **CLEVR-Dist** [29], **DMLab** [5], **Resisc45** [13], **Sun397** [64, 65], and **SVHN** [46]. From LIBSVM, we used the **Pendigits** [2] and **Covertype** [8] datasets, where cover covertype we used the scaled features version (i.e., covtype.scale). We also experiment with **CIFAR-10** [34] and **ImageNet** [18].

**Models.** The computer vision pre-trained models were accessed via timm. The strings used to load the models were: 'resnet50', 'vit\_base\_patch32\_224\_in21k'.

Complexity measure. To fairly compare all algorithms, we measure complexity by the number of batches evaluated, i.e., the number of stochastic gradient queries performed by the algorithm. U-DoG requires two batches per iteration while the rest of the algorithms we consider require only one. We note that the algorithms we compare also have different memory footprints and runtimes per iteration (by constant factors). We focus on the number of batches as our complexity metric since it is most relevant to our theory. Memory and per-iteration runtime optimizations are potentially possible for U-DoG and A-DoG; we leave investigating those to future work.

**ASGD model selection.** In the convex optimization experiments, we run (A)SGD over a wide range of momentum and learning rate parameters. For the batch size scaling figures (e.g., the left panels in Figure 1), we pick the parameters that reach the target metric in the smallest number of batches, providing a conservative upper bound on the performance obtainable with a very carefully tuned algorithm. The learning curve figures adjacent to the batch size scaling figures (e.g., the middle panels in Figure 1) show the learning curve for the (A)SGD run attaining the best target performance at the batch size indicated. For plots of learning curves at different batch sizes (e.g., Figure 19), we select the (A)SGD parameters that are the first to reach 95% of the best metric attained by A-DoG. If no such parameters exist, we take the parameters that reach the best performance within the iteration budget.

**Iterate averaging.** When evaluating test accuracy, we follow Ivgi et al. [27] and apply polynomial-decay weight averaging [57] with parameter 8. We did not tune this parameter or comprehensively check how beneficial the averaging is. Nevertheless, a cursory examination of our data suggests that averaging is mostly helpful across the board, but much more so for DoG and SGD than their accelerated counterparts. This is in line with the theory, which provides guarantees on (essentially) the last iterate of U-DoG, but only the averaged iterate of DoG.

**Learning rate schedule.** We use a constant learning rate schedule for (A)SGD. We do not use a decaying schedule such as cosine decay [37] as it would complicate comparing the smallest number of steps required to reach a target metric, since a decaying schedule requires knowing the number of steps in advance. Preliminary experiments indicate that, in the settings we study, cosine decay is not significantly better than a constant schedule combined with iterate averaging.

**Setting**  $r_{\epsilon}$ . Similarly to Ivgi et al. [27] we set  $r_{\epsilon} = \gamma(1 + ||x_0||)$  with  $\gamma = 10^{-6}$ . Our theoretical analysis suggests that the particular choice of  $r_{\epsilon}$  does not matter as long as it is sufficiently small relative to the distance between the weight initialization x0 and the optimum.

#### KREISLER IVGI HINDER CARMON

**Weight decay.** We do not use weight decay in most experiments, except for training from scratch on CIFAR-10 (Figure 29), where we use a weight decay of  $5 \cdot 10^{-4}$ . For DoG we decay the parameters toward zero, while for U-DoG and A-DoG we decay the parameters toward the initial point  $x_0$ . That is, for DoG we add  $5 \cdot 10^{-4}x$  to the stochastic gradient evaluated at x, while for U-DoG and A-DoG we add  $5 \cdot 10^{-4}(x - x_0)$ .

**Gradient accumulation.** Due to GPU memory limitations, in the non-convex experiments, for large batch sizes we divide each batch into smaller sub-batches of size of either 128 or 256 samples. We calculate the gradient for each sub-batch and average those into a single gradient which we then use to perform a single step. When batch normalization is used (that is, for ResNet50), this is not mathematically identical to computing the gradient in one large batch.

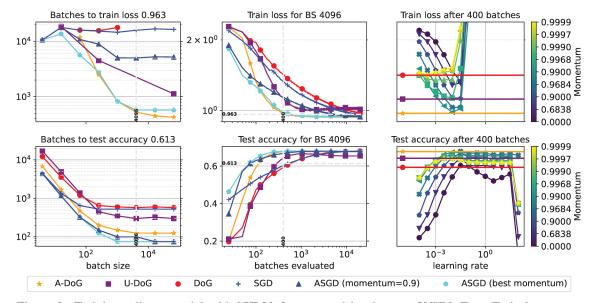


Figure 2: Training a linear model with ViT-32 features and log loss on SVHN. Top: Train loss. Bottom: Test accuracy after iterate averaging. First column: Batch size scaling of complexity to reach target performance. Second column: Learning curves. Third column: ASGD performance at all learning rates and momenta, contrasted with DoG variants.

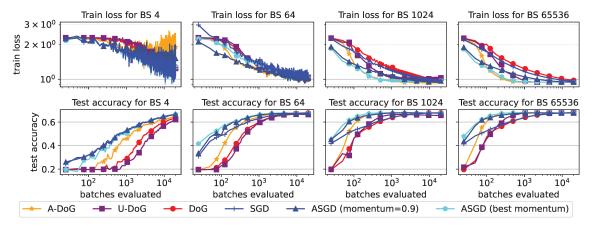


Figure 3: Training a linear model with ViT-32 features and log loss on SVHN. Top: Loss vs. batches processed training for different batch sizes. Bottom: Test accuracy of averaged model vs. batches processed for different batch sizes.

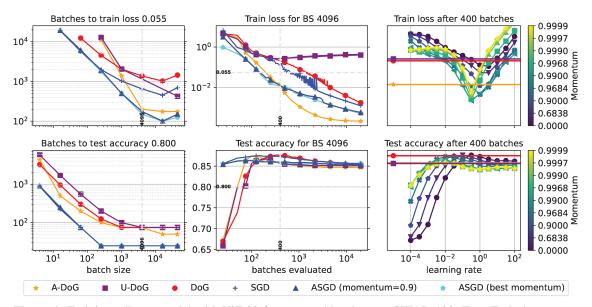


Figure 4: Training a linear model with ViT-32 features and log loss on CIFAR-100. Top: Train loss. Bottom: Test accuracy after iterate averaging. First column: Batch size scaling of complexity to reach target performance. Second column: Learning curves. Third column: ASGD performance at all learning rates and momenta, contrasted with DoG variants.

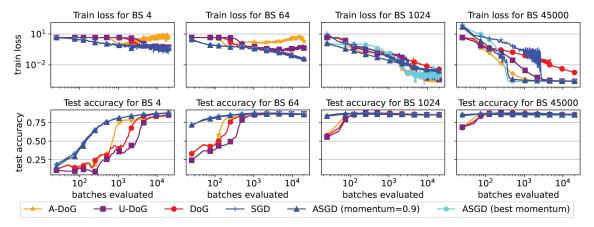


Figure 5: Training a linear model with ViT-32 features and log loss on CIFAR-100. Top: Loss vs. batches processed training for different batch sizes. Bottom: Test accuracy of averaged model vs. batches processed for different batch sizes.

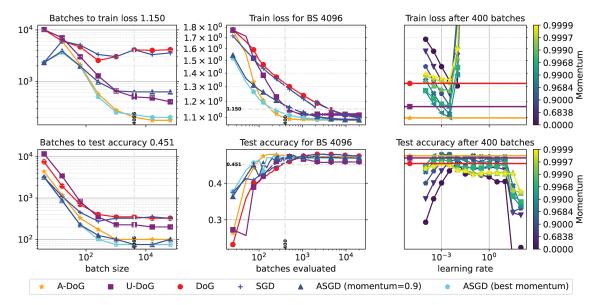


Figure 6: Training a linear model with ViT-32 features and log loss on DMLab. Top: Train loss. Bottom: Test accuracy after iterate averaging. First column: Batch size scaling of complexity to reach target performance. Second column: Learning curves. Third column: ASGD performance at all learning rates and momenta, contrasted with DoG variants.

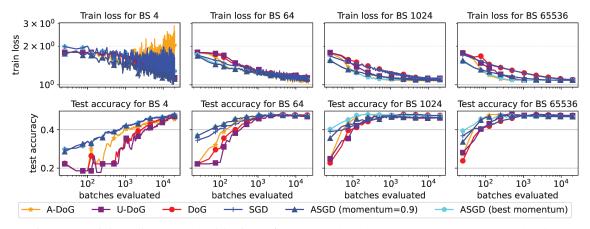


Figure 7: Training a linear model with ViT-32 features and log loss on DMLab. Top: Loss vs. batches processed training for different batch sizes. Bottom: Test accuracy of averaged model vs. batches processed for different batch sizes.

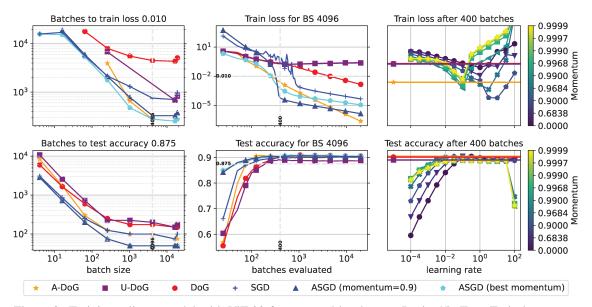


Figure 8: Training a linear model with ViT-32 features and log loss on Resisc45. Top: Train loss. Bottom: Test accuracy after iterate averaging. First column: Batch size scaling of complexity to reach target performance. Second column: Learning curves. Third column: ASGD performance at all learning rates and momenta, contrasted with DoG variants.

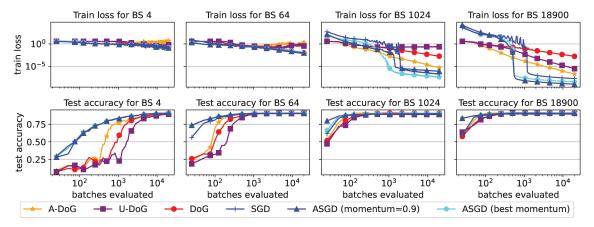


Figure 9: Training a linear model with ViT-32 features and log loss on Resisc45. Top: Loss vs. batches processed training for different batch sizes. Bottom: Test accuracy of averaged model vs. batches processed for different batch sizes.

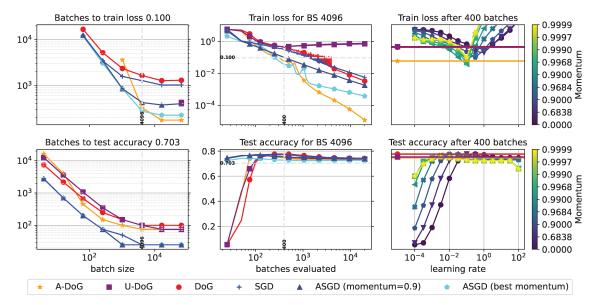


Figure 10: Training a linear model with ViT-32 features and log loss on Sun397. Top: Train loss. Bottom: Test accuracy after iterate averaging. First column: Batch size scaling of complexity to reach target performance. Second column: Learning curves. Third column: ASGD performance at all learning rates and momenta, contrasted with DoG variants.

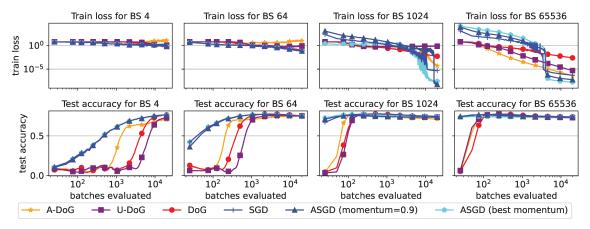


Figure 11: Training a linear model with ViT-32 features and log loss on Sun397. Top: Loss vs. batches processed training for different batch sizes. Bottom: Test accuracy of averaged model vs. batches processed for different batch sizes.

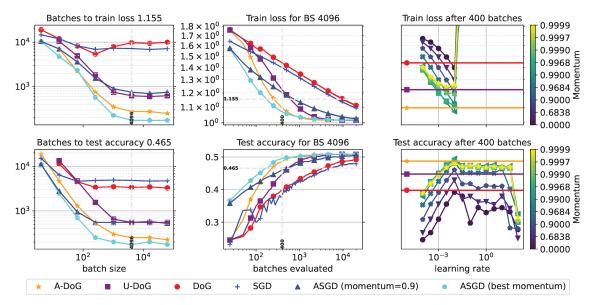


Figure 12: Training a linear model with ViT-32 features and log loss on CLEVR-Dist. Top: Train loss. Bottom: Test accuracy after iterate averaging. First column: Batch size scaling of complexity to reach target performance. Second column: Learning curves. Third column: ASGD performance at all learning rates and momenta, contrasted with DoG variants.

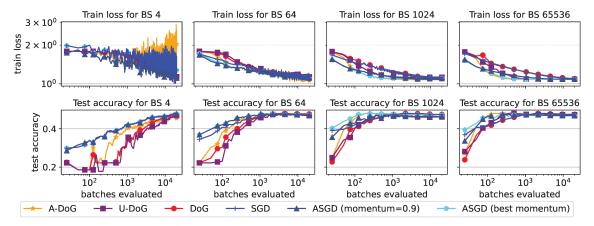


Figure 13: Training a linear model with ViT-32 features and log loss on CLEVR-Dist. Top: Loss vs. batches processed training for different batch sizes. Bottom: Test accuracy of averaged model vs. batches processed for different batch sizes.

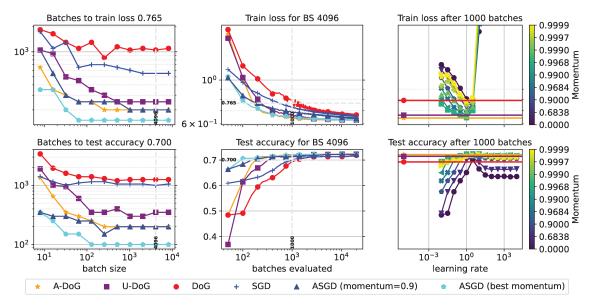


Figure 14: Training a linear model with log loss on LIBSVM/CovertypeScale. Top: Train loss. Bottom: Test accuracy after iterate averaging. First column: Batch size scaling of complexity to reach target performance. Second column: Learning curves. Third column: ASGD performance at all learning rates and momenta, contrasted with DoG variants.

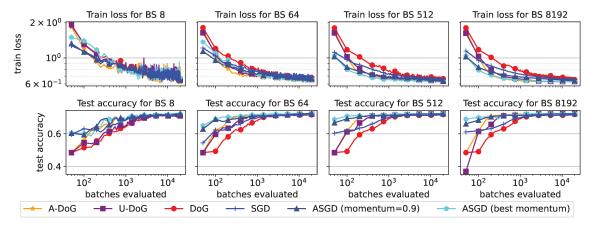


Figure 15: Training a linear model with log loss on LIBSVM/CovertypeScale. Top: Loss vs. batches processed training for different batch sizes. Bottom: Test accuracy of averaged model vs. batches processed for different batch sizes.

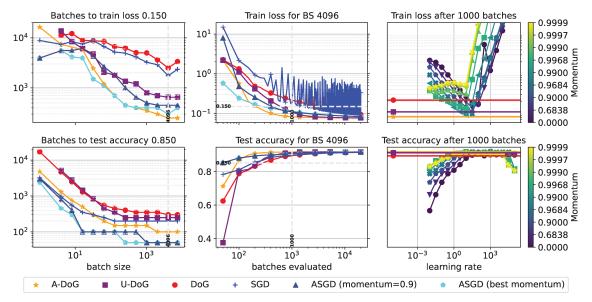


Figure 16: Training a linear model with log loss on LIBSVM/Pendigits. Top: Train loss. Bottom: Test accuracy after iterate averaging. First column: Batch size scaling of complexity to reach target performance. Second column: Learning curves. Third column: ASGD performance at all learning rates and momenta, contrasted with DoG variants.

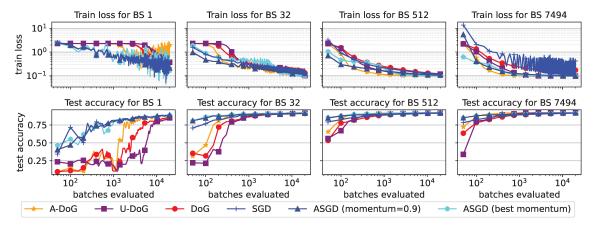


Figure 17: Training a linear model with log loss on LIBSVM/Pendigits. Top: Loss vs. batches processed training for different batch sizes. Bottom: Test accuracy of averaged model vs. batches processed for different batch sizes. As most algorithms here fail to converge at reasonable rate, we use significantly lower targets to choose hyper-parameters.

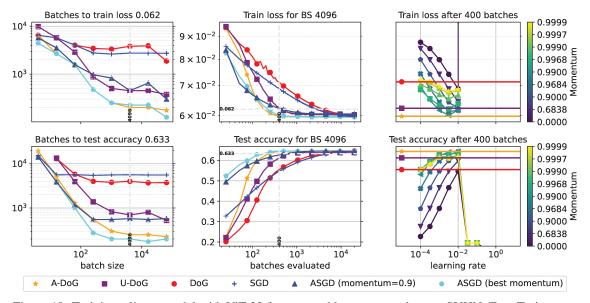


Figure 18: Training a linear model with ViT-32 features and least-squares loss on SVHN. Top: Train loss. Bottom: Test accuracy after iterate averaging. First column: Batch size scaling of complexity to reach target performance. Second column: Learning curves. Third column: ASGD performance at all learning rates and momenta, contrasted with DoG variants. This is the same as Figure 1.

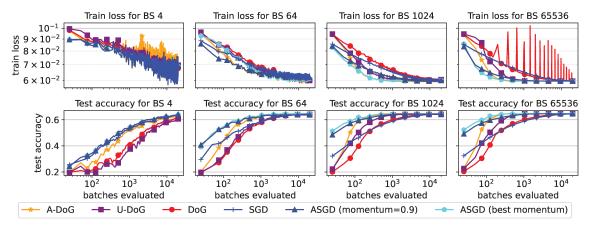


Figure 19: Training a linear model with ResNet50 features and least-squares loss on SVHN. Top: Loss vs. batches processed training for different batch sizes. Bottom: Test accuracy of averaged model vs. batches processed for different batch sizes.

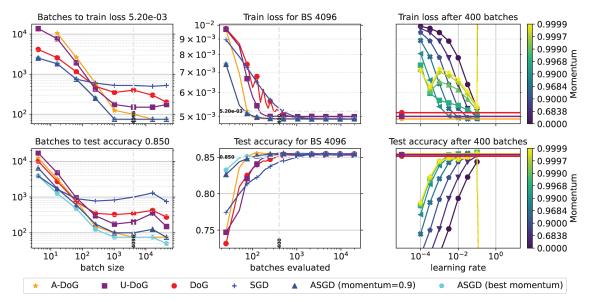


Figure 20: Training a linear model with ViT-32 features and least-squares loss on CIFAR-100. Top: Train loss. Bottom: Test accuracy after iterate averaging. First column: Batch size scaling of complexity to reach target performance. Second column: Learning curves. Third column: ASGD performance at all learning rates and momenta, contrasted with DoG variants.

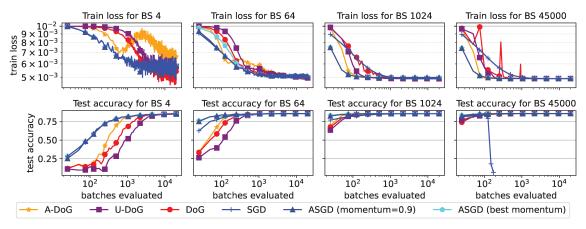


Figure 21: Training a linear model with ResNet50 features and least-squares loss on CIFAR-100. Top: Loss vs. batches processed training for different batch sizes. Bottom: Test accuracy of averaged model vs. batches processed for different batch sizes.

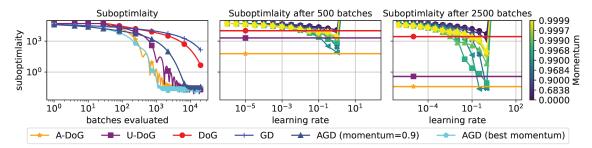


Figure 22: Training a model on a noiseless quadratic problem. At larger base learning rates, all AGD variants diverge while DoG variants remain stable, and U-DoG and A-DoG perform especially well.

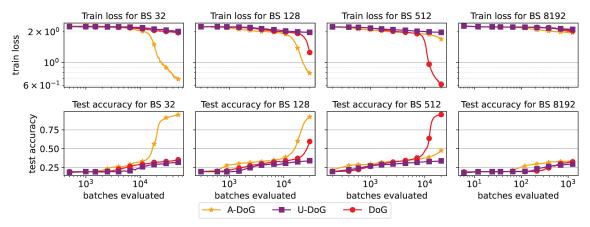


Figure 23: Training a ResNet50 model from scratch on SVHN. Top: Loss vs. batches processed training for different batch sizes. Bottom: Test accuracy vs. batches processed for averaged iterates at varied batch sizes.

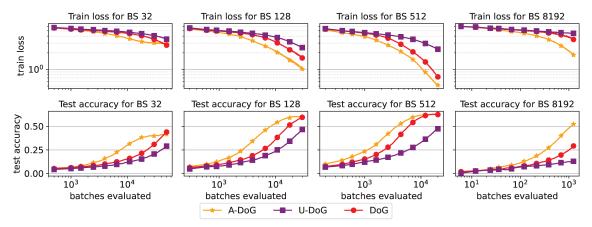


Figure 24: Training a ResNet50 model from scratch on Sun397. Top: Loss vs. batches processed training for different batch sizes. Bottom: Test accuracy vs. batches processed for averaged iterates at varied batch sizes.

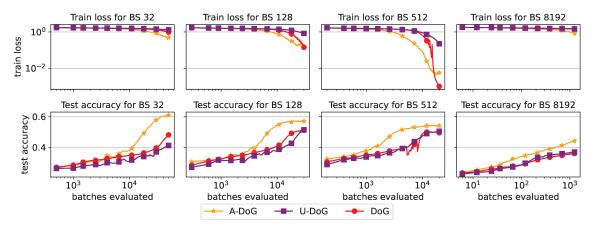


Figure 25: Training a ResNet50 model from scratch on DMLab. Top: Loss vs. batches processed training for different batch sizes. Bottom: Test accuracy vs. batches processed for averaged iterates at varied batch sizes.

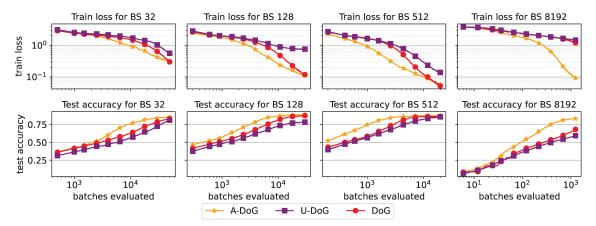


Figure 26: Training a ResNet50 model from scratch on Resisc45. Top: Loss vs. batches processed training for different batch sizes. Bottom: Test accuracy vs. batches processed for averaged iterates at varied batch sizes.

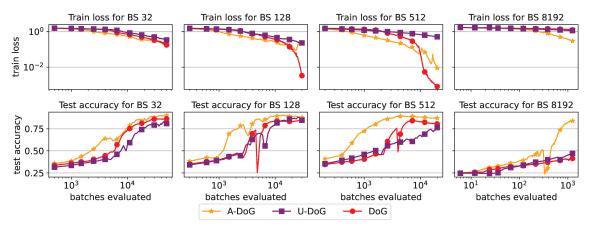


Figure 27: Training a ResNet50 model from scratch on CLEVR-Dist. Top: Loss vs. batches processed training for different batch sizes. Bottom: Test accuracy vs. batches processed for averaged iterates at varied batch sizes.

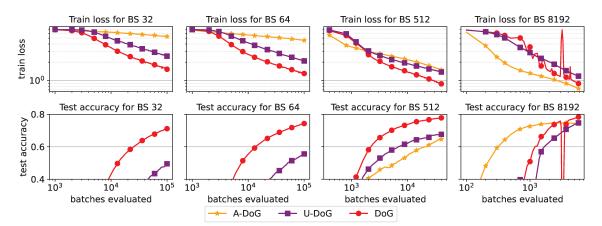


Figure 28: Fine-tuning a Clip-ViT-B/32 model on ImageNet, at different batch sizes. Top: Loss vs. step training curve for different batch sizes. Bottom: Test accuracy vs. step curve for averaged iterates at varied batch sizes.

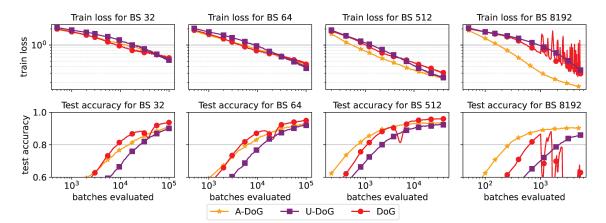


Figure 29: Training a Wide-ResNet-28-10 model on CIFAR-10 from scratch, at different batch sizes. Top: Loss vs. step training curve for different batch sizes. Bottom: Test accuracy vs. step curve for averaged iterates at varied batch sizes.

## KREISLER IVGI HINDER CARMON