

# ZEROth-ORDER OPTIMIZATION WITH TRAJECTORY-INFORMED DERIVATIVE ESTIMATION

Yao Shu\*, Zhongxiang Dai\*, Weicong Sng, Arun Verma,

Dept. of Computer Science, National University of Singapore, Republic of Singapore  
 {shuyao, daizhongxiang, sngweicong, arun}@comp.nus.edu.sg

Patrick Jaillet<sup>†</sup> & Bryan Kian Hsiang Low<sup>§</sup>

Dept. of Electrical Engineering and Computer Science, MIT, USA<sup>†</sup>

Dept. of Computer Science, National University of Singapore, Republic of Singapore<sup>§</sup>  
 jaillet@mit.edu, lowkh@comp.nus.edu.sg

## ABSTRACT

Zeroth-order (ZO) optimization, in which the derivative is unavailable, has recently succeeded in many important machine learning applications. Existing algorithms rely on finite difference (FD) methods for derivative estimation and gradient descent (GD)-based approaches for optimization. However, these algorithms suffer from query inefficiency because many additional function queries are required for derivative estimation in their every GD update, which typically hinders their deployment in real-world applications where every function query is expensive. To this end, we propose a *trajectory-informed* derivative estimation method which only employs the optimization trajectory (i.e., the history of function queries during optimization) and hence can eliminate the need for additional function queries to estimate a derivative. Moreover, based on our derivative estimation, we propose the technique of *dynamic virtual updates*, which allows us to reliably perform multiple steps of GD updates without reapplying derivative estimation. Based on these two contributions, we introduce the *zeroth-order* optimization with *trajectory-informed derivative* estimation (ZORD) algorithm for query-efficient ZO optimization. We theoretically demonstrate that our trajectory-informed derivative estimation and our ZORD algorithm improve over existing approaches, which is then supported by our real-world experiments such as black-box adversarial attack, non-differentiable metric optimization, and derivative-free reinforcement learning.

## 1 INTRODUCTION

Zeroth-order (ZO) optimization, in which the objective function to be optimized is only accessible by querying, has received great attention in recent years due to its success in many applications, e.g., black-box adversarial attack (Ru et al., 2020), non-differentiable metric optimization (Hiranandani et al., 2021), and derivative-free reinforcement learning (Salimans et al., 2017). In these problems, the derivative of objective function is either prohibitively costly to obtain or even non-existent, making it infeasible to directly apply standard derivative-based algorithms such as gradient descent (GD). In this regard, existing works have proposed to estimate the derivative using the *finite difference* (FD) methods and then apply GD-based algorithms using the *estimated derivative* for ZO optimization (Nesterov and Spokoiny, 2017; Cheng et al., 2021). These algorithms, which we refer to as *GD with estimated derivatives*, have been the most widely applied approach to ZO optimization especially for problems with high-dimensional input spaces, because of their theoretically guaranteed convergence and competitive practical performance. Unfortunately, these algorithms suffer from query inefficiency, which hinders their real-world deployment especially in applications with expensive-to-query objective functions, e.g., black-box adversarial attack.

---

\* Equal contribution.

Specifically, one of the reasons for the query inefficiency of existing algorithms on GD with estimated derivatives is that in addition to the necessary queries (i.e., the query of every updated input)<sup>1</sup>, the FD methods applied in these algorithms require a large number of *additional queries* to accurately estimate the derivative *at an input* (Berahas et al., 2022). This naturally begs the question: *Can we estimate a derivative without any additional query?* A natural approach to achieve this is to leverage the *optimization trajectory*, which is inherently available as a result of the necessary queries and their observations, to predict the derivatives. However, this requires a non-trivial method to simultaneously (a) predict a derivative using only the optimization trajectory (i.e., the history of updated inputs and their observations), and (b) quantify the uncertainty of this prediction to avoid using inaccurate predicted derivatives. Interestingly, the *Gaussian process* (GP) model satisfies both requirements and is hence a natural choice for such a derivative estimation. Specifically, under the commonly used assumption that the objective function is sampled from a GP (Srinivas et al., 2010), the derivative at *any* input in the domain follows a Gaussian distribution which, surprisingly, can be calculated using only the optimization trajectory. This allows us to (a) employ the mean of this Gaussian distribution as the estimated derivative, and (b) use the covariance matrix of this Gaussian distribution to obtain a principled measure of the predictive uncertainty and the accuracy of this derivative estimation, which together constitute our trajectory-informed derivative estimation (Sec. 3.1).

Another reason for the query inefficiency of the existing algorithms on GD with estimated derivatives is that *every* update in these algorithms requires reapplying derivative estimation and hence necessitates additional queries. This can preclude their adoption of a large number of GD updates since every update requires potentially expensive additional queries. Therefore, another question arises: *Can we perform multiple GD updates without reapplying derivative estimation and hence without any additional query?* To address this question, we propose a technique named *dynamic virtual updates* (Sec. 3.2). Specifically, thanks to the ability of our method to estimate the derivative at *any* input in the domain while only using existing optimization trajectory, we can apply *multi-step* GD updates without the need to reapply derivative estimation and hence without requiring any new query. Moreover, we can dynamically determine the number of steps for these updates by inspecting the aforementioned predictive uncertainty at every step, such that we only perform an update if the uncertainty is small enough (which also indicates that the estimation error is small, see Sec. 4.1).

By incorporating our aforementioned trajectory-informed derivative estimation and dynamic virtual updates into GD-based algorithms, we then introduce the *zeroth-order* optimization with *trajectory-informed derivative* estimation (ZORD) algorithm for query-efficient ZO optimization. We theoretically bound the estimation error of our trajectory-informed derivative estimation and show that this estimation error is non-increasing in the entire domain as the number of queries is increased and can even be exponentially decreasing in some scenarios (Sec. 4.1). Based on this, we prove the convergence of our ZORD algorithm, which improves over the existing ZO optimization algorithms that rely on the FD methods for derivative estimation (Sec. 4.2). Lastly, we use extensive experiments, such as black-box adversarial attack, non-differentiable metric optimization, and derivative-free reinforcement learning, to demonstrate that (a) our trajectory-informed derivative estimation improves over the existing FD methods and that (b) our ZORD algorithm consistently achieves improved query efficiency compared with previous ZO optimization algorithms (Sec. 5).

## 2 PRELIMINARIES

### 2.1 PROBLEM SETUP

Throughout this paper, we use  $\nabla$  and  $\partial_{\mathbf{x}}$  to denote, respectively, the total derivative (i.e., gradient) and partial derivative w.r.t the variable  $\mathbf{x}$ . We consider the minimization of a black-box objective function  $f : \mathcal{X} \rightarrow \mathbb{R}$ , in which  $\mathcal{X} \subset \mathbb{R}^d$  is a convex subset of the  $d$ -dimensional domain:

$$\min_{\mathbf{x} \in \mathcal{X}} f(\mathbf{x}). \quad (1)$$

Since we consider ZO optimization, the derivative information is not accessible and instead, we are only allowed to query the inputs in  $\mathcal{X}$ . For every queried input  $\mathbf{x} \in \mathcal{X}$ , we observe a corresponding noisy output of  $y(\mathbf{x}) = f(\mathbf{x}) + \zeta$ , in which  $\zeta$  is a zero-mean Gaussian noise with a variance of  $\sigma^2$ :

<sup>1</sup>In practice, it is usually necessary to query every updated input to measure the optimization performance and select the best-performing input. We refer to these queries as *necessary queries*.

**Algorithm 1:** Standard (Projected) GD with Estimated Derivatives

---

```

1: Input: Objective function  $f : \mathcal{X} \rightarrow \mathbb{R}$ ,
   initialization  $\mathbf{x}_0$ , iteration number  $T$ ,
   learning rates  $\{\eta_t\}_{t=1}^T$ , projection
   function  $\mathcal{P}_{\mathcal{X}}(\mathbf{x})$ 
2: for iteration  $t = 1, \dots, T$  do
3:    $g(\mathbf{x}_{t-1}) \approx \nabla f(\mathbf{x}_{t-1})$  with (2)
4:    $\mathbf{x}_t \leftarrow \mathcal{P}_{\mathcal{X}}(\mathbf{x}_{t-1} - \eta_{t-1}g(\mathbf{x}_{t-1}))$ 
5:   Query  $\mathbf{x}_t$  to yield  $y(\mathbf{x}_t)$ 
6: end for
7: Return  $\arg \min_{\mathbf{x}_{1:T}} y(\mathbf{x})$ 

```

---

**Algorithm 2:** ZORD (Ours)

---

```

1: Input: In addition to the parameters in Algo. 1, set
   the steps of virtual updates  $\{V_t\}_{t=1}^T$ 
2: for iteration  $t = 1, \dots, T$  do
3:    $\mathbf{x}_{t,0} \leftarrow \mathbf{x}_{t-1}$ 
4:   for iteration  $\tau = 1, \dots, V_t$  do
5:      $\mathbf{x}_{t,\tau} \leftarrow \mathcal{P}_{\mathcal{X}}(\mathbf{x}_{t,\tau-1} - \eta_{t,\tau-1}\nabla\mu_{t-1}(\mathbf{x}_{t,\tau-1}))$ 
6:   end for
7:   Query  $\mathbf{x}_t = \mathbf{x}_{t,\tau}$  to yield  $y(\mathbf{x}_t)$ 
8:   Update (4) using optimization trajectory
9: end for
10: Return  $\arg \min_{\mathbf{x}_{1:T}} y(\mathbf{x})$ 

```

---

$\zeta \sim \mathcal{N}(0, \sigma^2)$ . Besides, we adopt a common assumption on  $f$  which has already been widely used in the literature of *Bayesian optimization* (BO) (Srinivas et al., 2010; Kandasamy et al., 2018): we assume that  $f$  is sampled from a *Gaussian process* (GP). A GP  $\mathcal{GP}(\mu(\cdot), k(\cdot, \cdot))$ , which is characterized by a mean function  $\mu(\cdot)$  and a covariance function  $k(\cdot, \cdot)$ , is a stochastic process in which any finite subset of random variables follows a multi-variate Gaussian distribution (Rasmussen and Williams, 2006). In addition, following the common practice of GP and BO, we assume w.l.o.g. that  $\mu(\mathbf{x}) = 0$  and  $k(\mathbf{x}, \mathbf{x}') \leq 1$  ( $\forall \mathbf{x}, \mathbf{x}' \in \mathcal{X}$ ). We also assume that the kernel function  $k$  is differentiable, and that  $\|\partial_z \partial_{z'} k(\mathbf{z}, \mathbf{z}')|_{\mathbf{z}=\mathbf{z}'=\mathbf{x}}\|_2 \leq \kappa^2$ ,  $\forall \mathbf{x} \in \mathcal{X}$  for some  $\kappa > 0$ . This is satisfied by most commonly used kernels such as the squared exponential (SE) kernel (Rasmussen and Williams, 2006).

## 2.2 ZO OPTIMIZATION WITH ESTIMATED DERIVATIVES

To solve (1), GD with estimated derivatives (e.g., Algo. 1) has been developed (Flaxman et al., 2005; Ghadimi and Lan, 2013; Nesterov and Spokoiny, 2017; Liu et al., 2018a;b). Particularly, these algorithms first estimate the derivative of  $f$  (line 3 of Algo. 1) and then plug the estimated derivative into GD-based methods to obtain the next input for querying (lines 4-5 of Algo. 1). In these algorithms, the derivative is typically estimated by averaging the finite difference approximation of the directional derivatives for  $f$  along certain directions, which we refer to as the *finite difference* (FD) method in this paper. For example, given a parameter  $\lambda$  and directions  $\{\mathbf{u}_i\}_{i=1}^n$ , the derivative  $\nabla f$  at any  $\mathbf{x} \in \mathcal{X}$  can be estimated by the following FD method (Berahas et al., 2022):

$$\nabla f(\mathbf{x}) \approx g(\mathbf{x}) \triangleq \sum_{i=1}^n \frac{y(\mathbf{x} + \lambda \mathbf{u}_i) - y(\mathbf{x})}{\lambda} \mathbf{u}_i. \quad (2)$$

The directions  $\{\mathbf{u}_i\}_{i=1}^n$  are usually sampled from the standard Gaussian distribution (Nesterov and Spokoiny, 2017) or uniformly from the unit sphere (Flaxman et al., 2005), or set as the standard basis vectors with 1 at one of its coordinates and 0 otherwise (Lian et al., 2016). As mentioned before, existing FD methods typically require many additional queries (i.e.,  $\{\mathbf{x} + \lambda \mathbf{u}_i\}_{i=1}^n$ ) to achieve an accurate derivative estimation in every iteration of Algo. 1 (Berahas et al., 2022), making existing ZO optimization algorithms (Flaxman et al., 2005; Nesterov and Spokoiny, 2017) query-inefficient.

## 3 ZO OPTIMIZATION VIA TRAJECTORY-INFORMED DERIVATIVE ESTIMATION

To improve existing GD with estimated derivatives (e.g., Algo. 1), we propose the ZORD algorithm (Algo. 2), which achieves more query-efficient ZO optimization thanks to our two major contributions. Firstly, we propose a derived GP-based derivative estimation method which only uses the optimization trajectory and consequently does not require any additional query for derivative estimation (Sec. 3.1). Secondly, thanks to the ability of our method to estimate the derivative at any input in the domain without any additional query and to measure the estimation error in a principled way, we develop the technique of *dynamic virtual updates* to further improve the query efficiency of our ZORD (Sec. 3.2).

### 3.1 TRAJECTORY-INFORMED DERIVATIVE ESTIMATION

To begin with, if a function  $f$  follows a GP, then its derivative  $\nabla f$  also follows a GP (Rasmussen and Williams, 2006). This is formalized by our Lemma 1 below (proof in Appx. B.1), which then provides us a principled way to estimate the derivative at any input in the domain.

**Lemma 1** (Derived GP for Derivatives). *If a function  $f$  follows a GP:  $f \sim \mathcal{GP}(\mu(\cdot), \sigma^2(\cdot, \cdot))$ , then*

$$\nabla f \sim \mathcal{GP}(\nabla \mu(\cdot), \partial \sigma^2(\cdot, \cdot))$$

where  $\partial \sigma^2(\cdot, \cdot)$  denotes the cross partial derivative w.r.t the first and second arguments of  $\sigma^2(\cdot, \cdot)$ .

**$f$  Follows the Posterior GP.** As discussed in Sec. 2.1, we assume that  $f \sim \mathcal{GP}(\mu(\cdot), k(\cdot, \cdot))$ . So, in every iteration  $t$  of our Algo. 2, conditioned on the current optimization trajectory  $\mathcal{D}_{t-1} \triangleq \{(\mathbf{x}_\tau, y_\tau)\}_{\tau=1}^{t-1}$ ,  $f$  follows the *posterior GP*:  $f \sim \mathcal{GP}(\mu_{t-1}(\cdot), \sigma_{t-1}^2(\cdot, \cdot))$  with the mean function  $\mu_{t-1}(\cdot)$  and the covariance function  $\sigma_{t-1}^2(\cdot, \cdot)$  defined as below (Rasmussen and Williams, 2006):

$$\begin{aligned} \mu_{t-1}(\mathbf{x}) &\triangleq \mathbf{k}_{t-1}(\mathbf{x})^\top (\mathbf{K}_{t-1} + \sigma^2 \mathbf{I})^{-1} \mathbf{y}_{t-1} \\ \sigma_{t-1}^2(\mathbf{x}, \mathbf{x}') &\triangleq k(\mathbf{x}, \mathbf{x}') - \mathbf{k}_{t-1}(\mathbf{x})^\top (\mathbf{K}_{t-1} + \sigma^2 \mathbf{I})^{-1} \mathbf{k}_{t-1}(\mathbf{x}') \end{aligned} \quad (3)$$

where  $\mathbf{y}_{t-1}^\top \triangleq [y_\tau]_{\tau=1}^{t-1}$  and  $\mathbf{k}_{t-1}(\mathbf{x})^\top \triangleq [k(\mathbf{x}, \mathbf{x}_\tau)]_{\tau=1}^{t-1}$  are  $(t-1)$ -dimensional row vectors, and  $\mathbf{K}_{t-1} \triangleq [k(\mathbf{x}_\tau, \mathbf{x}_{\tau'})]_{\tau, \tau'=1}^{t-1}$  is a  $(t-1) \times (t-1)$ -dimensional matrix. Define  $\sigma_{t-1}^2(\mathbf{x}) \triangleq \sigma_{t-1}^2(\mathbf{x}, \mathbf{x})$ , the posterior distribution at  $\mathbf{x}$  is Gaussian with mean  $\mu_{t-1}(\mathbf{x})$  and variance  $\sigma_{t-1}^2(\mathbf{x})$ .

**$\nabla f$  Follows the Derived GP for Derivatives.** Substituting (3) into Lemma 1, we have that

$$\nabla f \sim \mathcal{GP}(\nabla \mu_{t-1}(\cdot), \partial \sigma_{t-1}^2(\cdot, \cdot)), \quad (4)$$

in which the mean  $\nabla \mu_{t-1}(\mathbf{x})$  at  $\mathbf{x}$  and the covariance  $\partial \sigma_{t-1}^2(\mathbf{x}, \mathbf{x}')$  at  $\mathbf{x}, \mathbf{x}'$  are

$$\begin{aligned} \nabla \mu_{t-1}(\mathbf{x}) &\triangleq \partial_{\mathbf{z}} \mathbf{k}_{t-1}(\mathbf{z})^\top (\mathbf{K}_{t-1} + \sigma^2 \mathbf{I})^{-1} \mathbf{y}_{t-1} \Big|_{\mathbf{z}=\mathbf{x}}, \\ \partial \sigma_{t-1}^2(\mathbf{x}, \mathbf{x}') &\triangleq \partial_{\mathbf{z}} \partial_{\mathbf{z}'} k(\mathbf{z}, \mathbf{z}') - \partial_{\mathbf{z}} \mathbf{k}_{t-1}(\mathbf{z})^\top (\mathbf{K}_{t-1} + \sigma^2 \mathbf{I})^{-1} \partial_{\mathbf{z}'} \mathbf{k}_{t-1}(\mathbf{z}') \Big|_{\mathbf{z}=\mathbf{x}, \mathbf{z}'=\mathbf{x}'}, \end{aligned} \quad (5)$$

in which  $\partial_{\mathbf{z}} \mathbf{k}_{t-1}(\mathbf{z}) \triangleq [\partial_{\mathbf{z}} k(\mathbf{z}, \mathbf{x}_\tau)]_{\tau=1}^{t-1}$  is a  $(t-1) \times d$ -dimensional matrix and  $\partial_{\mathbf{z}} \partial_{\mathbf{z}'} k(\mathbf{z}, \mathbf{z}')$  is a  $d \times d$ -dimensional matrix. Therefore,  $\nabla \mu_{t-1}(\mathbf{x})$  is a  $d$ -dimensional vector and  $\partial \sigma_{t-1}^2(\mathbf{x}, \mathbf{x}')$  is a  $d \times d$ -dimensional matrix. We refer to this GP (4) followed by  $\nabla f$  as the *derived GP for derivatives*.

So, define  $\partial \sigma_{t-1}^2(\mathbf{x}) \triangleq \partial \sigma_{t-1}^2(\mathbf{x}, \mathbf{x})$ , we have that for any input  $\mathbf{x} \in \mathcal{X}$ , the derivative  $\nabla f(\mathbf{x})$  at  $\mathbf{x}$  follows a  $d$ -dimensional Gaussian distribution:  $\nabla f(\mathbf{x}) \sim \mathcal{N}(\nabla \mu_{t-1}(\mathbf{x}), \partial \sigma_{t-1}^2(\mathbf{x}))$ . This allows us to (a) estimate the derivative  $\nabla f(\mathbf{x})$  at any input  $\mathbf{x} \in \mathcal{X}$  using the posterior mean  $\nabla \mu_{t-1}(\mathbf{x})$  of the derived GP for derivatives (4):

$$\nabla f(\mathbf{x}) \approx \nabla \mu_{t-1}(\mathbf{x}), \quad (6)$$

and (b) employ the posterior covariance matrix  $\partial \sigma_{t-1}^2(\mathbf{x})$  to obtain a principled measure of the uncertainty for this derivative estimation, which together constitute our novel derivative estimation. Remarkably, our derivative estimation only makes use of the naturally available optimization trajectory  $\mathcal{D}_{t-1}$  and *does not need any additional query*, which is in stark contrast to the existing FD methods (e.g., (2)) that require many additional queries for their derivative estimation. Moreover, our principled measure of uncertainty allows us to perform dynamic virtual updates (Sec. 3.2) and theoretically guarantee the quality of our derivative estimation (Sec. 4.1).

### 3.2 DYNAMIC VIRTUAL UPDATES

Note that our derived GP-based derivative estimation (6) can estimate the derivative at *any* input  $\mathbf{x}$  within the domain. As a result, in every iteration  $t$  of our ZORD algorithm, for a step  $\tau \geq 1$ , after performing a GD update using the estimated derivative at  $\mathbf{x}_{t, \tau-1}$  (i.e.,  $\nabla \mu_{t-1}(\mathbf{x}_{t, \tau-1})$ ) to reach the input  $\mathbf{x}_{t, \tau}$  (line 5 of Algo. 2), we can again estimate the derivative at  $\mathbf{x}_{t, \tau}$  (i.e.,  $\nabla \mu_{t-1}(\mathbf{x}_{t, \tau})$ ) and then perform another GD update to reach  $\mathbf{x}_{t, \tau+1}$  without requiring any additional query. This process can be repeated for multiple steps, and can further improve the query efficiency of our ZORD. Formally, given the projection function  $\mathcal{P}_{\mathcal{X}}(\mathbf{x}) \triangleq \arg \min_{\mathbf{z} \in \mathcal{X}} \|\mathbf{x} - \mathbf{z}\|_2^2 / 2$  and learning rates  $\{\eta_{t, \tau}\}_{\tau=0}^{V_t-1}$ , we perform the following virtual updates for  $V_t$  steps (lines 4-6 of Algo. 2):

$$\mathbf{x}_{t, \tau} = \mathcal{P}_{\mathcal{X}}(\mathbf{x}_{t, \tau-1} - \eta_{t, \tau-1} \nabla \mu_{t-1}(\mathbf{x}_{t, \tau-1})) \quad \forall \tau = 1, \dots, V_t \quad (7)$$

and then choose the last  $\mathbf{x}_{t, V_t}$  to query (i.e., line 7 of Algo. 2). Importantly, these multi-step virtual GD updates are only feasible in our ZORD (Algo. 2) because *our derivative estimator* (6) *does not*

require any new query in all these steps, whereas the existing FD methods require additional queries to estimate the derivative in every step.

The number of steps for our virtual updates (i.e.,  $V_t$ ) induces an intriguing trade-off: An overly small  $V_t$  may not be able to fully exploit the benefit of our derivative estimation (6) which is free from the requirement for additional queries, yet an excessively large  $V_t$  may lead to the usage of inaccurate derivative estimations which can hurt the performance (validated in Appx. D.2). Remarkably, (4) allows us to *dynamically* choose  $V_t$  by inspecting our principled measure of the predictive uncertainty (i.e.,  $\partial\sigma_{t-1}^2(\mathbf{x})$ ) for every derivative estimation. Specifically, after reaching the input  $\mathbf{x}_{t,\tau}$ , we continue the virtual updates (to reach  $\mathbf{x}_{t,\tau+1}$ ) if our predictive uncertainty is small, i.e., if  $\|\partial\sigma_{t-1}^2(\mathbf{x}_{t,\tau})\|_2 \leq c$  where  $c$  is a confidence threshold; otherwise, we terminate the virtual updates and let  $V_t = \tau$  since the derivative estimation at  $\mathbf{x}_{t,\tau}$  is likely unreliable.<sup>2</sup>

## 4 THEORETICAL ANALYSIS

### 4.1 DERIVATIVE ESTIMATION ERROR

To begin with, we derive a theoretical guarantee on the error of our derivative estimation at any  $\mathbf{x}$ .

**Theorem 1** (Derivative Estimation Error). *Let  $\delta \in (0, 1)$  and  $\beta \triangleq \sqrt{d + 2(\sqrt{d} + 1) \ln(1/\delta)}$ . For any  $\mathbf{x} \in \mathcal{X}$  and any  $t \geq 1$ , the following holds with probability of at least  $1 - \delta$ ,*

$$\|\nabla f(\mathbf{x}) - \nabla \mu_t(\mathbf{x})\|_2 \leq \beta \sqrt{\|\partial\sigma_t^2(\mathbf{x})\|_2}.$$

Thm. 1 (proof in Appx. B.2) has presented an upper bound on the error of our derivative estimation (6) at any  $\mathbf{x} \in \mathcal{X}$  in terms of  $\sqrt{\|\partial\sigma_t^2(\mathbf{x})\|_2}$ , which is a measure of the uncertainty about our derivative estimation at  $\mathbf{x}$  (Sec. 3.1). This hence implies that the threshold  $c$  applied to our predictive uncertainty  $\|\partial\sigma_t^2(\mathbf{x})\|_2$  (Sec. 3.2) also ensures that the derivative estimation error is small during our dynamic virtual updates. Next, we show in the following theorem (proof in Appx. B.3) that our upper bound on the estimation error from Thm. 1 is non-increasing as the number of function queries is increased.

**Theorem 2** (Non-Increasing Error). *For any  $\mathbf{x} \in \mathcal{X}$  and any  $t \geq 1$ , we have that*

$$\|\partial\sigma_t^2(\mathbf{x})\|_2 \leq \|\partial\sigma_{t-1}^2(\mathbf{x})\|_2.$$

*Let  $\delta \in (0, 1)$ . Define  $r \triangleq \max_{\mathbf{x} \in \mathcal{X}, t \geq 1} \sqrt{\|\partial\sigma_t^2(\mathbf{x})\|_2 / \|\partial\sigma_{t-1}^2(\mathbf{x})\|_2}$ , given the  $\beta$  in Thm. 1, we then have that  $r \in [1/\sqrt{1 + 1/\sigma^2}, 1]$ , and that with a probability of at least  $1 - \delta$ ,*

$$\|\nabla f(\mathbf{x}) - \nabla \mu_t(\mathbf{x})\|_2 \leq \beta \sqrt{\|\partial\sigma_t^2(\mathbf{x})\|_2} \leq \kappa \beta r^t.$$

Thm. 2 shows that our upper bound on the derivative estimation error (i.e.,  $\beta \sqrt{\|\partial\sigma_t^2(\mathbf{x})\|_2}$  from Thm. 1) is guaranteed to be *non-increasing in the entire domain* as the number of function queries is increased. Moreover, in some situations (i.e., when  $r < 1$ ), our upper bound on the estimation error is even exponentially decreasing. Of note,  $r$  characterizes how fast the uncertainty about our derivative estimation (measured by  $\sqrt{\|\partial\sigma_t^2(\mathbf{x})\|_2}$ ) is reduced across the domain. Since GD-based algorithms usually perform a local search in a neighborhood (especially for the problems with high-dimensional input spaces), all the inputs within the local region are expected to be close to each other (measured by the kernel function  $k$ ). Moreover, as the objective function is usually smooth in the local region (i.e., its derivatives are continuous), reducing the uncertainty of the derivative at an input  $\mathbf{x}_t$  (i.e., by querying  $\mathbf{x}_t$ ) is also expected to decrease the uncertainty of the derivatives at the other inputs in the same local region (i.e., decrease  $\sqrt{\|\partial\sigma_t^2(\mathbf{x})\|_2}$ ). So,  $r < 1$  is expected to be a reasonable condition that can be satisfied in practice. This will also be corroborated by our empirical results (e.g., Figs. 1 and 2), which demonstrates that the error of our derivative estimation (6) is indeed reduced very fast.

**Our GP-based Method (6) vs. Existing FD Methods.** Our derivative estimation method based on the derived GP (6) is superior to the traditional FD methods (e.g., (2)) in a number of major aspects. (a) Our derivative estimation error can be exponentially decreasing in some situations (i.e., when  $r < 1$  in Thm. 2), which is unachievable for the existing FD methods since they can only

<sup>2</sup>The first step of GD update to reach  $\mathbf{x}_{t,1}$  is always performed, i.e.,  $V_t \geq 1$ .

attain a polynomial rate of reduction (Berahas et al., 2022). (b) Our method (6) does not need any additional query to estimate the derivative (but only requires the optimization trajectory), whereas the existing FD methods require additional queries for every derivative estimation. (c) Our method (6) is equipped with a principled measure of the predictive uncertainty and hence the estimation error for derivative estimation (i.e., via  $\sqrt{\|\partial\sigma_t^2(\mathbf{x})\|_2}$ , Thm. 1), which is typically unavailable for the existing FD methods. (d) Our method (6), unlike the existing FD methods, makes it possible to apply the technique of dynamic virtual updates (Sec. 3.2) thanks to its capability of estimating the derivative at any input in the domain without requiring any additional query and measuring the estimation error in a principled way (Thm. 1).

## 4.2 CONVERGENCE ANALYSIS

To analyze the convergence of our ZORD, besides our main assumption that  $f$  is sampled from a GP (Sec. 2.1), we assume that  $f$  is  $L_c$ -Lipchitz continuous for  $L_c > 0$ . This is a mild assumption since it has been shown that a function  $f$  sampled from a GP is Lipchitz continuous with high probability for commonly used kernels, e.g., the SE kernel and Matérn kernel with  $\nu > 2$  (Srinivas et al., 2010). We also assume that  $f$  is  $L_s$ -Lipchitz smooth, which is commonly adopted in the analysis GD-based algorithms (J Reddi et al., 2016). We aim to prove the convergence of our ZORD for nonconvex  $f$  by analyzing how fast it converges to a stationary point (Ghadimi and Lan, 2013; Liu et al., 2018a). Specifically, we follow the common practice of previous works (J Reddi et al., 2016; Liu et al., 2018b) to analyze the following derivative mapping:

$$G_{t,\tau} \triangleq (\mathbf{x}_{t,\tau} - \mathcal{P}_{\mathcal{X}}(\mathbf{x}_{t,\tau} - \eta_{t,\tau} \nabla f(\mathbf{x}_{t,\tau}))) / \eta_{t,\tau}. \quad (8)$$

The convergence of our ZORD is formally guaranteed by Thm. 3 below (proof in Appx. B.4).

**Theorem 3 (Convergence of ZORD).** *Let  $\delta \in (0, 1)$ . Suppose our ZORD (Algo. 2) is run with  $V_t = V$  and  $\eta_{t,\tau} = \eta \leq 1/L_s$  for any  $t$  and  $\tau$ . Then with probability of at least  $1 - \delta$ , when  $r < 1$ ,*

$$\min_{t \leq T} \frac{1}{V} \sum_{\tau=0}^{V-1} \|G_{t,\tau}\|_2^2 \leq \underbrace{\frac{2[f(\mathbf{x}_0) - f(\mathbf{x}^*)]/\eta}{TV}}_{\textcircled{1}} + \underbrace{\frac{2\alpha^2 r^2}{T(1-r^2)} + \frac{(2L_c + 1/\eta)\alpha r}{T(1-r)}}_{\textcircled{2}}$$

where  $\alpha \triangleq \kappa \sqrt{d + 2(\sqrt{d} + 1) \ln(VT/\delta)}$ . When  $r = 1$ , we instead have  $\textcircled{2} = 2\alpha^2 + (2L_c + 1/\eta)\alpha$ .

In the upper bound of Thm. 3, the term  $\textcircled{1}$  represents the convergence rate of (projected) GD when the true derivative is used and it asymptotically goes to 0 as  $T$  increases; the term  $\textcircled{2}$  corresponds to the impact of the error of our derivative estimation (6) on the convergence. In situations where  $r < 1$  which is a reasonably achievable condition as we have discussed in Sec. 4.1, the term  $\textcircled{2}$  will also asymptotically approach 0. This, remarkably, suggests that the impact of the derivative estimation error on the convergence vanishes asymptotically and our ZORD algorithm is guaranteed to converge to a stationary point (i.e.,  $\min_{t \leq T} \frac{1}{V} \sum_{\tau=0}^{V-1} \|G_{t,\tau}\|_2^2$  approaches 0) at the rate of  $\mathcal{O}(1/T)$  when  $r < 1$ . This is unattainable by existing ZO optimization algorithms using FD-based derivative estimation (Nesterov and Spokoiny, 2017; Liu et al., 2018b), because these methods typically converge to a stationary point at the rate of  $\mathcal{O}(1/T + \text{const.})$  with a constant learning rate. Even when  $r = 1$  where the term  $\textcircled{2}$  becomes a constant independent of  $T$ , our Thm. 3 is still superior to the convergence of these existing works because our result (Thm. 3) is based on the worst-case analysis whereas these works are typically based on the average-case analysis, i.e., their results only hold in expectation over the randomly sampled directions for derivative estimation. This means that their convergence may become even worse when inappropriate directions are used, e.g., directions that are nearly orthogonal to the true derivative which commonly happens in high-dimensional input spaces. In addition, given a fixed  $T$ , our ZORD enjoys a query complexity (i.e., the number of queries in  $T$  iterations) of  $\mathcal{O}(T)$ , which significantly improves over the  $\mathcal{O}(nT)$  of the existing works based on FD ( $n$  in Sec. 2.2).

The impacts of the number of steps of our virtual updates (i.e.,  $V$ ) are partially reflected in Thm. 3. Specifically, a larger  $V$  improves the reduction rate of the term  $\textcircled{1}$  because a larger number of virtual GD updates (without requiring additional queries) will be applied in our ZORD algorithm. This is also unachievable by existing ZO optimization algorithms using FD-based derivative estimation since they require additional queries for the derivative estimation in their every GD update. Meanwhile, a larger  $V$  may also negatively impact the performance of our ZORD since it may lead to the use of those estimated derivatives with large estimation errors (Sec. 3.2). However, this negative impact has

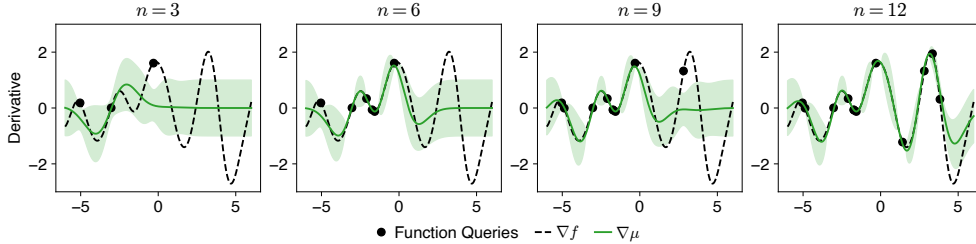


Figure 1: Our derived GP for derivative estimation (4) with different number  $n$  of queries. Green curve and its confidence interval denote the mean  $\nabla\mu(\mathbf{x})$  and standard deviation of the derived GP.

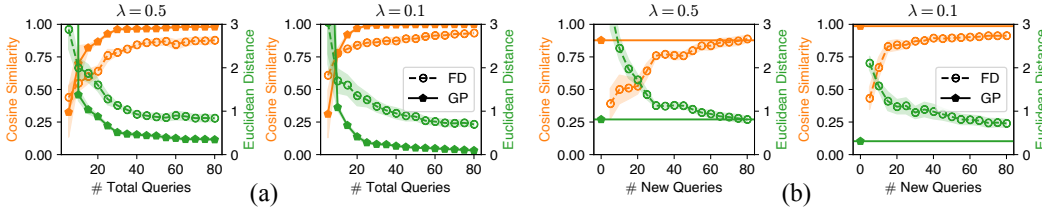


Figure 2: Comparison of the derivative estimation errors of our derived GP-based estimator (6) (GP) and the FD estimator, measured by cosine similarity (larger is better) and Euclidean distance (smaller is better). Each curve is the mean  $\pm$  standard error from five independent runs.

only been implicitly accounted for by the term ② because this term comes from our Thm. 2, which is based on a worst-case analysis and gives a uniform upper bound on the derivative estimation error *for all inputs in the domain*  $\mathcal{X}$ .

## 5 EXPERIMENTS

In this section, we firstly empirically verify the efficacy of our derived GP-based derivative estimator (6) in Sec. 5.1, and then demonstrate that our ZORD outperforms existing baseline methods for ZO optimization using synthetic experiments (Sec. 5.2) and real-world experiments (Secs. 5.3, 5.4).

### 5.1 DERIVATIVE ESTIMATION

Here we investigate the efficacy of our derivative estimator (6) based on the derived GP for derivatives (4). Specifically, we sample a function  $f$  (defined on a one-dimensional domain) from a GP using the SE kernel, and then use a set of randomly selected inputs as well as their noisy observations (as optimization trajectory) to calculate our derived GP for derivatives. The results (Fig. 1) illustrate a number of interesting insights. Firstly, in regions where (even only a few) function queries are performed (e.g., in the region of  $[-3, 0]$ ), our estimated derivative (i.e.,  $\nabla\mu_{t-1}(\mathbf{x})$ ) (6) generally aligns with the groundtruth derivative (i.e.,  $\nabla f(\mathbf{x})$ ) and our estimation uncertainty (i.e., characterized by  $\sqrt{\|\partial\sigma_{t-1}^2(\mathbf{x})\|_2}$ ) shrinks compared with other un-queried regions. These results hence demonstrate that our (4) is able to accurately estimate derivatives and reliably quantify the uncertainty of these estimations within the regions where function queries are performed. Secondly, as more input queries are collected (i.e., from left to right in Fig. 1), the uncertainty  $\sqrt{\|\partial\sigma_{t-1}^2(\mathbf{x})\|_2}$  in the entire domain is decreased in general. This provides an empirical justification for our Thm. 2 which guarantees non-increasing uncertainty and hence non-increasing estimation error. Lastly, note that with only 12 queries (rightmost figure), our derivative estimator is already able to accurately estimate the derivative in the entire domain, which represents a remarkable reduction rate of our derivative estimation error.

Next, we compare our derivative estimator (6) with the FD estimator (Sec. 2.2). Specifically, using the Ackley function with  $d = 10$  (see Appx. C.2), we firstly select an input  $\mathbf{x}_0$  and then follow the FD method (2) to randomly sample  $n$  directions  $\{\mathbf{u}_i\}_{i=1}^n$  from the standard Gaussian distribution, to construct input queries  $\{\mathbf{x}_0 + \lambda\mathbf{u}_i\}_{i=1}^n$  (see Sec. 2.2). Next, these queries and their observations are (a) used as the optimization trajectory to apply our derivative estimator (6), and (b) used by the FD method to estimate the derivative following (2). The results are shown in Fig. 2a (for two different values of  $\lambda$ ), in which for both our derived GP-based estimator (6) and the FD estimator, we measure the cosine similarity (larger is better) and Euclidean distance (smaller is better) between the estimated

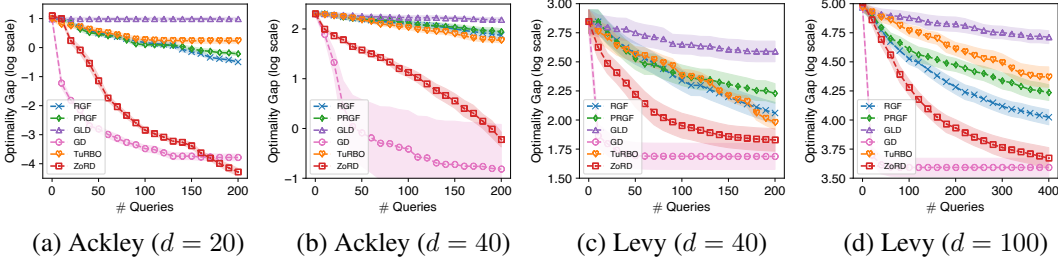


Figure 3: Optimization of Ackley and Levy functions with different dimensions. The  $x$ -axis and  $y$ -axis denote the number of queries and log-scaled optimality gap (i.e.,  $\log(f(\mathbf{x}_T) - f(\mathbf{x}^*))$ ) achieved after this number of queries. Each curve is the mean  $\pm$  standard error from ten independent runs.

Table 1: Comparison of the number of required queries to achieve a successful black-box adversarial attack. Every entry represents mean  $\pm$  standard deviation from five independent runs.

| Dataset  | Metric    | GLD            | RGF             | PRGF            | TuRBO-1       | TuRBO-10      | ZoRD                          |
|----------|-----------|----------------|-----------------|-----------------|---------------|---------------|-------------------------------|
| MNIST    | # Queries | 1780 $\pm$ 222 | 1192 $\pm$ 260  | 1236 $\pm$ 145  | 654 $\pm$ 70  | 747 $\pm$ 60  | <b>248<math>\pm</math>50</b>  |
|          | Speedup   | 7.2 $\times$   | 4.8 $\times$    | 5.0 $\times$    | 2.6 $\times$  | 3.0 $\times$  | <b>1.0<math>\times</math></b> |
| CIFAR-10 | # Queries | 964 $\pm$ 175  | 3622 $\pm$ 1155 | 4133 $\pm$ 1525 | 638 $\pm$ 108 | 708 $\pm$ 105 | <b>384<math>\pm</math>59</b>  |
|          | Speedup   | 2.5 $\times$   | 9.4 $\times$    | 10.8 $\times$   | 1.7 $\times$  | 1.8 $\times$  | <b>1.0<math>\times</math></b> |

derivative and the true derivative at  $\mathbf{x}_0$ . The figures show that our derivative estimation error enjoys a faster rate of reduction compared with the FD method, which corroborates our theoretical insights from Thm. 2 (Sec. 4.1) positing that our estimation error can be rapidly decreasing. Subsequently, to further highlight our advantage of being able to exploit the optimization trajectory and hence to eliminate the need for additional function queries (Sec. 4.1), we perform another comparison where our derived GP-based estimator (6) only utilizes 20 queries from the optimization trajectory (sampled using the same method above) for derivative estimation. The results (Fig. 2b) show that even with only these 20 queries (without any additional function query), our derivative estimator (6) achieves comparable or better estimation errors than FD using as many as 80 additional queries. Overall, the results in Fig. 2 have provided empirical supports for the superiority of our derived GP-based derivative estimation (6), which substantiates our theoretical justifications in Sec. 4.1.

## 5.2 SYNTHETIC EXPERIMENTS

Here we adopt the widely used Ackley and Levy functions with various dimensions (Eriksson et al., 2019) to show the superiority of our ZoRD. We compare ZoRD with a number of representative baselines for ZO optimization, e.g., RGF (Nesterov and Spokoiny, 2017) which uses FD for derivative estimation, PRGF (Cheng et al., 2021) which is a recent extension of RGF, GLD (Golovin et al., 2020) which is a recent ZO optimization algorithm based on direct search, and TuRBO (Eriksson et al., 2019) which is a highly performant Bayesian optimization (BO) algorithm. We also evaluate the performance of a first-order optimization algorithm, i.e., GD with true derivatives. More details are in Appx. C.2. The results are shown in Fig. 3, where ZoRD outperforms all other ZO optimization algorithms. Particularly, ZoRD considerably outperforms both RGF and PRGF, which can be attributed to our two major contributions. Firstly, our derivative estimator (6) used by ZoRD is more accurate and more query-efficient than the FD method adopted by RGF and PRGF, as theoretically justified in Sec. 4.1 and empirically demonstrated in Sec. 5.1. Secondly, our dynamic virtual updates (Sec. 3.2) can perform multi-step GD updates without requiring any additional query, which further improves the performance of ZoRD (validated in Appx. D.2). Moreover, ZoRD is the only ZO optimization algorithm that is able to converge to a comparable final performance to that of the GD with true derivatives in every figure of Fig. 3.

## 5.3 BLACK-BOX ADVERSARIAL ATTACK

We further compare our ZoRD with other ZO optimization algorithms in the problem of black-box adversarial attack on images, which is one of the most important applications of ZO optimization in recent years. In black-box adversarial attack (Ru et al., 2020), given a fully trained ML model and an image  $\mathbf{z}$ , we intend to find (through only function queries) a small perturbation  $\mathbf{x}$  to be added to  $\mathbf{z}$



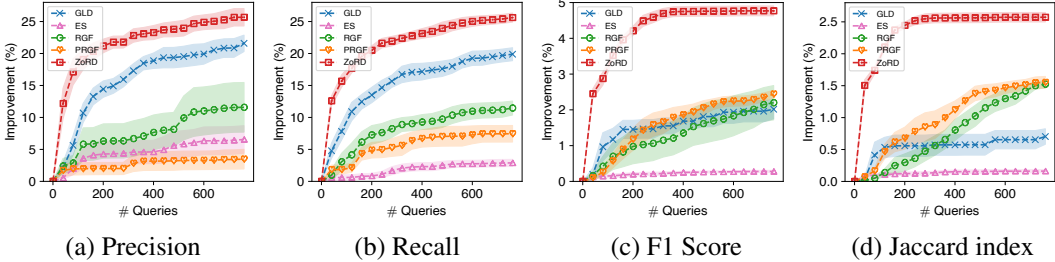


Figure 4: Optimization of different non-differentiable metrics on the Covertypes dataset. The  $x$ -axis and  $y$ -axis denote, respectively, the number of queries and the improvement on the non-differentiable metric. Each curve is the mean  $\pm$  standard error from five independent experiments.

such that the perturbed image  $z + x$  will be incorrectly classified by the ML model. Following the practice from (Cheng et al., 2021), we randomly select an image from MNIST (Lecun et al., 1998) ( $d = 28 \times 28$ ) or CIFAR-10 (Krizhevsky et al., 2009) ( $d = 32 \times 32$ ), and aim to add a perturbation with an  $L_\infty$  constraint to make a trained deep neural network misclassify the image (more details in Appx. C.3). Tab. 1 summarizes the number of required queries to achieve a successful attack by different algorithms (see results on multiple images in Appx. D.3). The results show that in such high-dimensional ZO optimization problems, our ZORD again significantly outperforms the other algorithms since it requires a considerably smaller number of queries to achieve a successful attack. Particularly, our ZORD is substantially more query-efficient than RGF and PRGF which rely on the FD methods for derivative estimation, e.g., for CIFAR-10, the number of queries required by RGF and PRGF are  $9.4\times$  and  $10.8\times$  of that required by ZORD. This further verifies the advantages of our trajectory-informed derivative estimation (as justified theoretically in Sec. 4.1 and empirically in Sec. 5.1) and dynamic virtual updates (as demonstrated in Appx. D.2). Remarkably, our ZORD also outperforms BO (i.e., TuRBO-1/10 which correspond to two versions of the TuRBO algorithm (Eriksson et al., 2019)) which has been widely shown to be query-efficient in black-box adversarial attack (Ru et al., 2020). Overall, these results showcase the ability of our ZORD to advance the other ZO optimization algorithms in challenging real-world ZO optimization problems.

#### 5.4 NON-DIFFERENTIABLE METRIC OPTIMIZATION

Non-differentiable metric optimization (Hiranandani et al., 2021; Huang et al., 2021), which has received a surging interest recently, can also be cast as a ZO optimization problem. We therefore use it to further demonstrate the superiority of our ZORD to other ZO optimization algorithms. Specifically, we firstly train a multilayer perceptron (MLP) ( $d = 2189$ ) on the Covertypes (Dua and Graff, 2017) dataset with the cross-entropy loss function. Then, we use the same dataset to fine-tune this MLP model by exploiting ZO optimization algorithms to optimize a non-differentiable metric, such as precision, recall, F1 score and Jaccard index (see more details in Appx. C.4). Here we additionally compare with the evolutionary strategy (ES) which has been previously applied for non-differentiable metric optimization (Huang et al., 2021). Fig. 4 illustrates the percentage improvements achieved by different algorithms during the fine-tuning process (i.e.,  $(f(x_0) - f(x_T)) \times 100\% / f(x_0)$ ). The results show that our ZORD again consistently outperforms the other ZO optimization algorithms in terms of both the query efficiency and the final converged performance. These results therefore further substantiate the superiority of ZORD in optimizing high-dimensional non-differentiable functions.

## 6 CONCLUSION

We have introduced the ZORD algorithm, which achieves query-efficient ZO optimization through two major contributions. Firstly, we have proposed a novel derived GP-based method (6) which only uses the optimization trajectory and hence eliminates the requirement for additional queries (Sec. 3.1) to estimate derivatives. Secondly, we have introduced a novel technique, i.e., dynamic virtual updates, which is made possible by our GP-based derivative estimation, to further improve the performance of our ZORD (Sec. 3.2). Through theoretical justifications (Sec. 4) and empirical demonstrations (Sec. 5), we show that our derived GP-based derivative estimation improve over existing FD methods and that our ZORD outperforms various ZO optimization baselines.

## 7 REPRODUCIBILITY STATEMENT

For our theoretical results, we have discussed all our assumptions in Sec. 2.1 & Sec. 4.2, and provided our complete proofs in Appx. B. For our empirical results, we have provided our detailed experimental settings in Appx. C and included our codes in the supplementary materials (i.e., the zip file).

### ACKNOWLEDGMENTS

This research is part of the programme DesCartes and is supported by the National Research Foundation, Prime Minister’s Office, Singapore under its Campus for Research Excellence and Technological Enterprise (CREATE) programme.

### REFERENCES

- Binxin Ru, Adam D. Cobb, Arno Blaas, and Yarin Gal. Bayesopt adversarial attack. In *Proc. ICLR*, 2020.
- Gaurush Hiranandani, Jatin Mathur, Harikrishna Narasimhan, Mahdi Milani Fard, and Sanmi Koyejo. Optimizing black-box metrics with iterative example weighting. In *Proc. ICML*, 2021.
- Tim Salimans, Jonathan Ho, Xi Chen, and Ilya Sutskever. Evolution strategies as a scalable alternative to reinforcement learning. arXiv:1703.03864, 2017.
- Yurii E. Nesterov and Vladimir G. Spokoiny. Random gradient-free minimization of convex functions. *Found. Comput. Math.*, 17(2):527–566, 2017.
- Shuyu Cheng, Guoqiang Wu, and Jun Zhu. On the convergence of prior-guided zeroth-order optimization algorithms. In *Proc. NeurIPS*, 2021.
- Albert S. Berahas, Liyuan Cao, Krzysztof Choromanski, and Katya Scheinberg. A theoretical and empirical comparison of gradient approximations in derivative-free optimization. *Found. Comput. Math.*, 22(2):507–560, 2022.
- Niranjan Srinivas, Andreas Krause, Sham M. Kakade, and Matthias W. Seeger. Gaussian process optimization in the bandit setting: No regret and experimental design. In *Proc. ICML*, 2010.
- Kirthevasan Kandasamy, Akshay Krishnamurthy, Jeff Schneider, and Barnabás Póczos. Parallelised Bayesian optimisation via Thompson sampling. In *Proc. AISTATS*, 2018.
- Carl Edward Rasmussen and Christopher K. I. Williams. *Gaussian processes for machine learning*. Adaptive computation and machine learning. MIT Press, 2006.
- Abraham Flaxman, Adam Tauman Kalai, and H. Brendan McMahan. Online convex optimization in the bandit setting: Gradient descent without a gradient. In *Proc. SODA*, 2005.
- Saeed Ghadimi and Guanghui Lan. Stochastic first- and zeroth-order methods for nonconvex stochastic programming. *SIAM Journal on Optimization*, 23(4):2341–2368, 2013.
- Sijia Liu, Bhavya Kailkhura, Pin-Yu Chen, Pai-Shun Ting, Shiyu Chang, and Lisa Amini. Zeroth-order stochastic variance reduction for nonconvex optimization. In *Proc. NeurIPS*, 2018a.
- Sijia Liu, Xingguo Li, Pin-Yu Chen, Jarvis D. Haupt, and Lisa Amini. Zeroth-order stochastic projected gradient descent for nonconvex optimization. In *Proc. GlobalSIP*, 2018b.
- Xiangru Lian, Huan Zhang, Cho-Jui Hsieh, Yijun Huang, and Ji Liu. A comprehensive linear speedup analysis for asynchronous stochastic parallel optimization from zeroth-order to first-order. In *Proc. NIPS*, 2016.
- Sashank J Reddi, Suvrit Sra, Barnabas Póczos, and Alexander J Smola. Proximal stochastic methods for nonsmooth nonconvex finite-sum optimization. In *Proc. NIPS*, 2016.
- David Eriksson, Michael Pearce, Jacob R. Gardner, Ryan Turner, and Matthias Poloczek. Scalable global optimization via local Bayesian optimization. In *Proc. NeurIPS*, 2019.

- Daniel Golovin, John Karro, Greg Kochanski, Chansoo Lee, Xingyou Song, and Qiuyi (Richard) Zhang. Gradientless descent: High-dimensional zeroth-order optimization. In *Proc. ICLR*, 2020.
- Y. Lecun, L. Bottou, Y. Bengio, and P. Haffner. Gradient-based learning applied to document recognition. *Proceedings of the IEEE*, pages 2278–2324, 1998.
- Alex Krizhevsky, Geoffrey Hinton, et al. Learning multiple layers of features from tiny images. Technical report, Citeseer, 2009.
- Chen Huang, Shuangfei Zhai, Pengsheng Guo, and Josh M. Susskind. Metricopt: Learning to optimize black-box evaluation metrics. In *Proc. CVPR*, 2021.
- Dheeru Dua and Casey Graff. UCI machine learning repository, 2017. URL <http://archive.ics.uci.edu/ml>.
- Sebastian U Stich, Christian L Muller, and Bernd Gartner. Optimization of convex functions with random pursuit. *SIAM Journal on Optimization*, 23(2):1284–1309, 2013.
- Sayak Ray Chowdhury and Aditya Gopalan. On kernelized multi-armed bandits. In *Proc. ICML*, 2017.
- Zhongxiang Dai, Haibin Yu, Bryan Kian Hsiang Low, and Patrick Jaillet. Bayesian optimization meets Bayesian optimal stopping. In *Proc. ICML*, 2019.
- Zhongxiang Dai, Bryan Kian Hsiang Low, and Patrick Jaillet. Federated bayesian optimization via thompson sampling. In *Proc. NeurIPS*, 2020.
- Benjamin Letham, Roberto Calandra, Akshara Rai, and Eytan Bakshy. Re-examining linear embeddings for high-dimensional Bayesian optimization. In *Proc. NeurIPS*, 2020.
- Andrew Ilyas, Logan Engstrom, and Aleksander Madry. Prior convictions: Black-box adversarial attacks with bandits and priors. In *Proc. ICLR*, 2019.
- Florian Meier, Asier Mujika, Marcelo Matheus Gauy, and Angelika Steger. Improving gradient estimation in evolutionary strategies with past descent directions. arXiv:1910.05268, 2019.
- Niru Maheswaranathan, Luke Metz, George Tucker, Dami Choi, and Jascha Sohl-Dickstein. Guided evolutionary strategies: Augmenting random search with surrogate gradients. In *Proc. ICML*, 2019.
- Shuyu Cheng, Yinpeng Dong, Tianyu Pang, Hang Su, and Jun Zhu. Improving black-box adversarial attacks with a transfer-based prior. In *NeurIPS*, 2019.
- Beatrice Laurent and Pascal Massart. Adaptive estimation of a quadratic functional by model selection. *Annals of Statistics*, pages 1302–1338, 2000.
- Sayak Ray Chowdhury and Aditya Gopalan. No-regret algorithms for multi-task Bayesian optimization. In *Proc. AISTATS*, 2021.
- Stephen P. Boyd and Lieven Vandenbergh. *Convex Optimization*. Cambridge University Press, 2014.
- Diederik P. Kingma and Jimmy Ba. Adam: A method for stochastic optimization. In *Proc. ICLR*, 2015.
- Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *Proc. CVPR*, 2016.
- Greg Brockman, Vicki Cheung, Ludwig Pettersson, Jonas Schneider, John Schulman, Jie Tang, and Wojciech Zaremba. OpenAI Gym. arXiv:1606.01540, 2016.
- M. D. McKay, R. J. Beckman, and W. J. Conover. A comparison of three methods for selecting values of input variables in the analysis of output from a computer code. *Technometrics*, 21(2): 239–245, 1979.
- Jian Tan, Niv Nayman, and Mengchang Wang. CobBO: Coordinate backoff Bayesian optimization with two-stage kernels. arXiv:2101.05147, 2021.
- Hong Qian and Yang Yu. Derivative-free reinforcement learning: A review. *Frontiers Comput. Sci.*, 15(6):156336, 2021.

## APPENDIX A RELATED WORK

Various types of algorithms have been proposed in the literature to solve ZO optimization problems, e.g., direct search, Bayesian optimization (BO) and GD-based algorithms with estimated derivatives. Particularly, direct search, e.g., (Stich et al., 2013; Golovin et al., 2020), relies on the comparison of function values at different inputs for the updates, which can be query-inefficient in practice owing to its indirect utilization of function values. In contrast, Bayesian optimization (BO) directly utilizes the function values to model the objective function using a Gaussian process (GP) and iteratively selects the inputs to query by trading off sampling potentially optimal inputs (i.e., exploitation) and inputs that can improve the GP belief of the objective function over the entire input domain (i.e., exploration) (Chowdhury and Gopalan, 2017; Srinivas et al., 2010; Dai et al., 2019; 2020). However, in ZO optimization problems with high-dimensional input spaces, BO algorithms typically suffer from query inefficiency and large computational complexity (Rasmussen and Williams, 2006; Letham et al., 2020; Eriksson et al., 2019), which significantly hinders their real-world applications. Therefore, GD-based algorithms with estimated derivatives, which inherit the advantage of GD-based algorithms in optimizing functions with high-dimensional input spaces, have been more widely applied in practice. For these algorithms, the derivatives are commonly estimated using the finite difference (FD) approximation (which requires additional function queries) of the directional derivatives along selected directions, in which the directions can be randomly sampled unit vectors Flaxman et al. (2005), Gaussian vectors (Nesterov and Spokoiny, 2017), or standard bases (Lian et al., 2016) (Sec. 2.2). More recently, some works have incorporated a time-dependent prior (i.e., the estimated derivative in the previous iteration) into existing FD methods to improve the quality of its derivative estimation (Ilyas et al., 2019; Meier et al., 2019; Cheng et al., 2021). Nevertheless, such a prior is also estimated by the FD method (i.e., in the previous iteration) and can hence be biased owing to its estimation error, which may even lead to larger derivative estimation errors in practice due to compounding errors. Another line of work has taken the surrogate derivatives from other sources to help reduce the derivative estimation error of existing FD methods (Maheswaranathan et al., 2019; Cheng et al., 2019). However, these surrogate derivatives may generally be unavailable in practice. Importantly, these existing FD methods require additional function queries for every derivation estimation during optimization, which will significantly increase the query complexity of ZO optimization algorithms which employ these FD methods for derivative estimation.

## APPENDIX B PROOFS

### B.1 PROOF OF LEMMA 1

According to Rasmussen and Williams (2006), if a function  $f$  follows from a Gaussian process, its derivative also follows a Gaussian process determined by its mean  $\mathbb{E}[\cdot]$  and covariance  $\text{Cov}(\cdot, \cdot)$ , i.e.,

$$\nabla f \sim \mathcal{GP}(\mathbb{E}[\nabla f], \text{Cov}(\nabla f, \nabla f)) . \quad (9)$$

So, to prove Lemma 1, we only need to derive the mean and the covariance of the Gaussian process above for a function  $f$  that is sampled from another Gaussian process, i.e.,  $f \sim \mathcal{GP}(\mu(\cdot), \sigma^2(\cdot, \cdot))$ . Specifically, for the mean  $\mathbb{E}[\nabla f]$ , we have

$$\mathbb{E}[\nabla f] = \nabla \mathbb{E}[f] = \nabla \mu . \quad (10)$$

where the first equality derives from the interchangeability of the expectation and derivative operation based on the Leibniz integral rule. The second equality comes from the fact that  $\mathbb{E}[f] = \mu$ .

For the covariance  $\text{Cov}(\nabla f, \nabla f)$ , we have

$$\begin{aligned} \text{Cov}(\nabla f(\mathbf{z}), \nabla f(\mathbf{z}')) &\stackrel{(a)}{=} \mathbb{E} \left[ (\nabla f(\mathbf{z}) - \mathbb{E}[\nabla f(\mathbf{z})])^\top (\nabla f(\mathbf{z}') - \mathbb{E}[\nabla f(\mathbf{z}')]) \right] \\ &\stackrel{(b)}{=} \mathbb{E} \left[ \nabla (f(\mathbf{z}) - \mathbb{E}[f(\mathbf{z})])^\top \nabla (f(\mathbf{z}') - \mathbb{E}[f(\mathbf{z}')]) \right] \\ &\stackrel{(c)}{=} \mathbb{E} \left[ \partial_{\mathbf{z}} \partial_{\mathbf{z}'} (f(\mathbf{z}) - \mathbb{E}[f(\mathbf{z})])^\top (f(\mathbf{z}') - \mathbb{E}[f(\mathbf{z}')]) \right] \\ &\stackrel{(d)}{=} \partial_{\mathbf{z}} \partial_{\mathbf{z}'} \mathbb{E} \left[ (f(\mathbf{z}) - \mathbb{E}[f(\mathbf{z})])^\top (f(\mathbf{z}') - \mathbb{E}[f(\mathbf{z}')]) \right] \\ &\stackrel{(e)}{=} \partial_{\mathbf{z}} \partial_{\mathbf{z}'} \sigma_f^2(\mathbf{z}, \mathbf{z}') . \end{aligned} \quad (11)$$

Notably, (b) and (d) also derive from the interchangeability of the expectation and derivative operation based on the Leibniz integral rule. Besides, (e) is obtained based on  $\text{Cov}(f, f) = \sigma^2(\cdot, \cdot)$ . This finally completes our proof.

## B.2 PROOF OF THEOREM 1

To begin with, we introduce the following concentration inequality for standard multi-variate Gaussian distribution:

**Lemma B.1** (Laurent and Massart (2000)). *Let  $\zeta \sim \mathcal{N}(\mathbf{0}, \mathbf{I}_m)$  and  $\delta \in (0, 1)$  then*

$$\mathbb{P}\left(\|\zeta\|_2 \leq \sqrt{m + 2(\sqrt{m} + 1)\ln(1/\delta)}\right) \geq 1 - \delta. \quad (12)$$

Define  $\zeta \triangleq (\partial\sigma_t^2(\mathbf{x}))^{-1/2} (\nabla f(\mathbf{x}) - \nabla\mu_t(\mathbf{x}))$ , according to Lemma 1, we then have that  $\zeta$  follows a standard multi-variate Gaussian distribution, i.e.,

$$\zeta \sim \mathcal{N}(\mathbf{0}, \mathbf{I}_d). \quad (13)$$

Let  $\delta \in (0, 1)$ . By substituting the result above into Lemma B.1, the following holds with probability of at least  $1 - \delta$ :

$$\begin{aligned} \|\nabla f(\mathbf{x}) - \nabla\mu_t(\mathbf{x})\|_2 &= \left\| (\partial\sigma_t^2(\mathbf{x}))^{-1/2} \zeta \right\|_2 \\ &\leq \sqrt{\|\partial\sigma_t^2(\mathbf{x})\|_2} \|\zeta\|_2 \\ &\leq \sqrt{d + 2(\sqrt{d} + 1)\ln(1/\delta)} \sqrt{\|\partial\sigma_t^2(\mathbf{x})\|_2} \\ &= \beta \sqrt{\|\partial\sigma_t^2(\mathbf{x})\|_2} \end{aligned} \quad (14)$$

with  $\beta \triangleq \sqrt{d + 2(\sqrt{d} + 1)\ln(1/\delta)}$  and the first inequality is from the Cauchy-Schwarz inequality, which completes our proof.

## B.3 PROOF OF THEOREM 2

We first introduce the following lemmas.

**Lemma B.2** (Chowdhury and Gopalan (2021)). *For any  $\sigma \in \mathbb{R}$  and any matrix  $\mathbf{A}$ , the following hold*

$$\mathbf{I} - \mathbf{A}^\top (\mathbf{A}\mathbf{A}^\top + \sigma^2\mathbf{I})^{-1} \mathbf{A} = \sigma^2 (\mathbf{A}^\top \mathbf{A} + \sigma^2\mathbf{I})^{-1}. \quad (15)$$

**Lemma B.3** (Sherman-Morrison formula). *For any invertible square matrix  $\mathbf{A}$  and column vectors  $\mathbf{u}, \mathbf{v}$ , suppose  $\mathbf{A} + \mathbf{u}\mathbf{v}^\top$  is invertible, then the following holds*

$$(\mathbf{A} + \mathbf{u}\mathbf{v}^\top)^{-1} = \mathbf{A}^{-1} - \frac{\mathbf{A}^{-1}\mathbf{u}\mathbf{v}^\top\mathbf{A}^{-1}}{1 + \mathbf{v}^\top\mathbf{A}^{-1}\mathbf{u}}. \quad (16)$$

**Preparation.** We then introduce some additional notations and representations for our proof of Theorem 2. Following the common practice in (Chowdhury and Gopalan, 2021), we let the kernel  $k$  be defined by  $\psi(\mathbf{x})$ , i.e.,  $k(\mathbf{x}, \mathbf{x}') = \psi(\mathbf{x})^\top \psi(\mathbf{x}')$ , and  $\phi(\mathbf{x}) \triangleq \nabla\psi(\mathbf{x})$ . We then further define the  $(t \times d)$ -dimensional Jacobian matrix  $\phi_t(\mathbf{x}) \triangleq [\phi(\mathbf{x})^\top \psi(\mathbf{x}_\tau)]_{\tau=1}^t$  and  $\Psi_t \triangleq [\psi(\mathbf{x}_\tau)]_{\tau=1}^t$ . The matrix  $\mathbf{K}_t$  and the covariance matrix  $\partial\sigma_t^2(\mathbf{x})$  defined on the optimization trajectory  $\mathcal{D}_t$  in our Sec. 3.1 can be reformulated as

$$\begin{aligned} \mathbf{K}_t &= \Psi_t^\top \Psi_t, \\ \partial\sigma_t^2(\mathbf{x}) &= \phi(\mathbf{x})^\top \phi(\mathbf{x}) - \phi_t(\mathbf{x})^\top (\mathbf{K}_t + \sigma^2\mathbf{I})^{-1} \phi_t(\mathbf{x}). \end{aligned} \quad (17)$$

Based on the reformulation above, define  $\mathbf{V}_t \triangleq \Psi_t \Psi_t^\top + \sigma^2 \mathbf{I}$ , we can further reformulate  $\partial \sigma_t^2(\mathbf{x})$  as below

$$\begin{aligned}
\partial \sigma_t^2(\mathbf{x}) &\stackrel{(a)}{=} \phi(\mathbf{x})^\top \phi(\mathbf{x}) - \phi_t(\mathbf{x})^\top (\mathbf{K}_t + \sigma^2 \mathbf{I})^{-1} \phi_t(\mathbf{x}) \\
&\stackrel{(b)}{=} \phi(\mathbf{x})^\top \phi(\mathbf{x}) - \phi(\mathbf{x})^\top \Psi_t (\Psi_t^\top \Psi_t + \sigma^2 \mathbf{I})^{-1} \Psi_t^\top \phi(\mathbf{x}) \\
&\stackrel{(c)}{=} \phi(\mathbf{x})^\top \left( \mathbf{I} - \Psi_t (\Psi_t^\top \Psi_t + \sigma^2 \mathbf{I})^{-1} \Psi_t^\top \right) \phi(\mathbf{x}) \\
&\stackrel{(d)}{=} \sigma^2 \phi(\mathbf{x})^\top (\Psi_t \Psi_t^\top + \sigma^2 \mathbf{I})^{-1} \phi(\mathbf{x}) \\
&\stackrel{(e)}{=} \sigma^2 \phi(\mathbf{x})^\top \mathbf{V}_t^{-1} \phi(\mathbf{x}).
\end{aligned} \tag{18}$$

Note that (b) is obtained by exploiting the fact that  $\mathbf{K}_t = \Psi_t^\top \Psi_t$  and  $\phi_t(\mathbf{x}) = \phi(\mathbf{x})^\top \Psi_t$ . In addition, (d) comes from Lemma B.2 by replacing the matrix  $\mathbf{A}$  in Lemma B.2 with the matrix  $\Psi_t^\top$ .

**First Part.** We then prove the first half part of our Theorem 2, i.e., the following Lemma B.4.

**Lemma B.4** (Non-Increasing Variance Norm). *For any  $\mathbf{x} \in \mathcal{X}$  and any  $t \geq 1$ , we have that*

$$\|\partial \sigma_t^2(\mathbf{x})\|_2 \leq \|\partial \sigma_{t-1}^2(\mathbf{x})\|_2. \tag{19}$$

*Proof.* Based on our additional notations and representations, we have

$$\begin{aligned}
\partial \sigma_t^2(\mathbf{x}) &\stackrel{(a)}{=} \sigma^2 \phi(\mathbf{x})^\top \mathbf{V}_t^{-1} \phi(\mathbf{x}) \\
&\stackrel{(b)}{=} \sigma^2 \phi(\mathbf{x})^\top (\Psi_{t-1} \Psi_{t-1}^\top + \sigma^2 \mathbf{I} + \psi(\mathbf{x}_t) \psi(\mathbf{x}_t)^\top)^{-1} \phi(\mathbf{x}) \\
&\stackrel{(c)}{=} \sigma^2 \phi(\mathbf{x})^\top (\mathbf{V}_{t-1} + \psi(\mathbf{x}_t) \psi(\mathbf{x}_t)^\top)^{-1} \phi(\mathbf{x}) \\
&\stackrel{(d)}{=} \sigma^2 \phi(\mathbf{x})^\top \mathbf{V}_{t-1}^{-1} \phi(\mathbf{x}) - \sigma^2 (1 + \psi(\mathbf{x}_t)^\top \mathbf{V}_{t-1}^{-1} \psi(\mathbf{x}_t))^{-1} \phi(\mathbf{x})^\top \mathbf{V}_{t-1}^{-1} \psi(\mathbf{x}_t) \psi(\mathbf{x}_t)^\top \mathbf{V}_{t-1}^{-1} \phi(\mathbf{x}) \\
&\stackrel{(e)}{=} \partial \sigma_{t-1}^2(\mathbf{x}) - \sigma^2 (1 + \psi(\mathbf{x}_t)^\top \mathbf{V}_{t-1}^{-1} \psi(\mathbf{x}_t))^{-1} \phi(\mathbf{x})^\top \mathbf{V}_{t-1}^{-1} \psi(\mathbf{x}_t) \psi(\mathbf{x}_t)^\top \mathbf{V}_{t-1}^{-1} \phi(\mathbf{x}) \\
&\stackrel{(f)}{\preceq} \partial \sigma_{t-1}^2(\mathbf{x}).
\end{aligned} \tag{20}$$

Note that (a) follows from the aforementioned definition of  $\mathbf{V}_t$  and (b) comes from the fact that  $\Psi_t \Psi_t^\top = \Psi_{t-1} \Psi_{t-1}^\top + \psi(\mathbf{x}_t) \psi(\mathbf{x}_t)^\top$ . Similarly, (c) uses the definition of  $\mathbf{V}_{t-1}$ . In addition, equality (d) derives from Lemma B.3 by letting  $\mathbf{A} = \mathbf{V}_{t-1}$  and  $\mathbf{u} = \mathbf{v} = \psi(\mathbf{x}_t)$  and (e) follows from the reformulation of  $\partial \sigma_{t-1}^2(\mathbf{x})$  in (18). Finally, (f) derives from the positive semi-definite property of  $\phi(\mathbf{x})^\top \mathbf{V}_{t-1}^{-1} \psi(\mathbf{x}_t) \psi(\mathbf{x}_t)^\top \mathbf{V}_{t-1}^{-1} \phi(\mathbf{x})$  as well as the fact that  $1 + \psi(\mathbf{x}_t)^\top \mathbf{V}_{t-1}^{-1} \psi(\mathbf{x}_t) > 0$ . That is, for any column vector  $\mathbf{z}$  we have that

$$\begin{aligned}
\mathbf{z}^\top \phi(\mathbf{x})^\top \mathbf{V}_{t-1}^{-1} \psi(\mathbf{x}_t) \psi(\mathbf{x}_t)^\top \mathbf{V}_{t-1}^{-1} \phi(\mathbf{x}) \mathbf{z} &= (\phi(\mathbf{x}_t)^\top \mathbf{V}_{t-1}^{-1} \phi(\mathbf{x}) \mathbf{z})^\top (\phi(\mathbf{x}_t)^\top \mathbf{V}_{t-1}^{-1} \phi(\mathbf{x}) \mathbf{z}) \\
&= \|\phi(\mathbf{x}_t)^\top \mathbf{V}_{t-1}^{-1} \phi(\mathbf{x}) \mathbf{z}\|_2^2 \\
&\geq 0.
\end{aligned} \tag{21}$$

So,  $\phi(\mathbf{x})^\top \mathbf{V}_{t-1}^{-1} \psi(\mathbf{x}_t) \psi(\mathbf{x}_t)^\top \mathbf{V}_{t-1}^{-1} \phi(\mathbf{x})$  is positive semi-definite. Following a similar way, we are also able to verify that  $1 + \psi(\mathbf{x}_t)^\top \mathbf{V}_{t-1}^{-1} \psi(\mathbf{x}_t) > 0$  by showing that  $\psi(\mathbf{x}_t)^\top \mathbf{V}_{t-1}^{-1} \psi(\mathbf{x}_t) \geq 0$  using the decomposition of  $\mathbf{V}_{t-1}^{-1}$  from the Principle Component Analysis (PCA). Since  $\partial \sigma_t^2(\mathbf{x}) \preceq \partial \sigma_{t-1}^2(\mathbf{x})$  is equivalent to  $\|\partial \sigma_t^2(\mathbf{x})\|_2 \leq \|\partial \sigma_{t-1}^2(\mathbf{x})\|_2$ , we then complete the proof of first half part of our Theorem 2.  $\square$

**Second Part.** To prove the rest of our Theorem 2, we firstly introduce the following lemmas.

**Lemma B.5.** *For any  $\mathbf{x} \in \mathcal{X}$  and any  $t \geq 1$ , the following holds*

$$\mathbf{V}_t^{-1} \preceq \mathbf{V}_{t-1}^{-1}. \tag{22}$$

*Proof.* For any column vector  $\mathbf{z}$ , we have

$$\begin{aligned}
\mathbf{z}^\top (\mathbf{V}_t - \mathbf{V}_{t-1}) \mathbf{z} &= \mathbf{z}^\top \psi(\mathbf{x}_t) \psi(\mathbf{x}_t)^\top \mathbf{z} \\
&= (\psi(\mathbf{x}_t)^\top \mathbf{z})^\top (\psi(\mathbf{x}_t)^\top \mathbf{z}) \\
&= \|\psi(\mathbf{x}_t)^\top \mathbf{z}\|_2^2 \\
&\geq 0.
\end{aligned} \tag{23}$$

The first equality comes from the intermediate result in (20). So,  $\mathbf{V}_t - \mathbf{V}_{t-1}$  is positive semi-definite, i.e.,  $\mathbf{V}_{t-1} \preceq \mathbf{V}_t$ . This can also indicate that  $\mathbf{V}_t^{-1} \preceq \mathbf{V}_{t-1}^{-1}$ , which thus completes our proof.  $\square$

**Lemma B.6** (Lower Bound of Variance Norm). *For any  $\mathbf{x} \in \mathcal{X}$  and any  $t \geq 1$ , the following holds*

$$1/(1 + 1/\sigma^2) \|\partial\sigma_{t-1}^2(\mathbf{x})\|_2 \leq \|\partial\sigma_t^2(\mathbf{x})\|_2. \tag{24}$$

*Proof.* We firstly show that

$$\begin{aligned}
\left\| \mathbf{V}_t^{-1/2} \psi(\mathbf{x}) \psi(\mathbf{x})^\top \mathbf{V}_t^{-1/2} \right\|_2 &\stackrel{(a)}{\leq} \left\| \mathbf{V}_t^{-1/2} \psi(\mathbf{x}) \right\|_2 \left\| \psi(\mathbf{x})^\top \mathbf{V}_t^{-1/2} \right\|_2 \\
&\stackrel{(b)}{=} \left\| \psi(\mathbf{x})^\top \mathbf{V}_t^{-1/2} \right\|_2^2 \\
&\stackrel{(c)}{=} \psi(\mathbf{x})^\top \mathbf{V}_t^{-1/2} \mathbf{V}_t^{-1/2} \psi(\mathbf{x}) \\
&\stackrel{(d)}{=} \psi(\mathbf{x})^\top \mathbf{V}_t^{-1} \psi(\mathbf{x}) \\
&\stackrel{(e)}{\leq} \psi(\mathbf{x})^\top \mathbf{V}_{t-1}^{-1} \psi(\mathbf{x}) \\
&\stackrel{(f)}{\leq} \psi(\mathbf{x})^\top \mathbf{V}_0^{-1} \psi(\mathbf{x}) \\
&\stackrel{(g)}{=} \psi(\mathbf{x})^\top \psi(\mathbf{x}) / \sigma^2 \\
&\stackrel{(h)}{=} 1/\sigma^2.
\end{aligned} \tag{25}$$

Note that (a) derives from the Cauchy-Schwarz inequality. As for (b) and (c), they have exploited the fact that  $(\mathbf{V}_t^{-1/2} \psi(\mathbf{x}))^\top = \psi(\mathbf{x})^\top \mathbf{V}_t^{-1/2}$  and  $\psi(\mathbf{x})^\top \mathbf{V}_t^{-1/2}$  is a row vector. In addition, (e) follows from Lemma B.5. Finally, (g) results from  $\mathbf{V}_0^{-1} = \mathbf{I}/\sigma^2$  and (h) derives from the assumption that  $k(\mathbf{x}, \mathbf{x}) \leq 1$  ( $\forall \mathbf{x} \in \mathcal{X}$ ) in Sec. 2.1. Alternatively, we can restate the result above as

$$\mathbf{V}_t^{-1/2} \psi(\mathbf{x}) \psi(\mathbf{x})^\top \mathbf{V}_t^{-1/2} \preceq \sigma^{-2} \mathbf{I}. \tag{26}$$

We then complete our proof on the first inequality in Lemma B.6 using the following inequality:

$$\begin{aligned}
\partial\sigma_t^2(\mathbf{x}) &\stackrel{(a)}{=} \sigma^2 \phi(\mathbf{x})^\top (\mathbf{V}_{t-1} + \psi(\mathbf{x}_t) \psi(\mathbf{x}_t)^\top)^{-1} \phi(\mathbf{x}) \\
&\stackrel{(b)}{=} \sigma^2 \phi(\mathbf{x})^\top \left[ \mathbf{V}_{t-1}^{1/2} \left( \mathbf{I} + \mathbf{V}_{t-1}^{-1/2} \psi(\mathbf{x}_t) \psi(\mathbf{x}_t)^\top \mathbf{V}_{t-1}^{-1/2} \right) \mathbf{V}_{t-1}^{1/2} \right]^{-1} \phi(\mathbf{x}) \\
&\stackrel{(c)}{=} \sigma^2 \phi(\mathbf{x})^\top \mathbf{V}_{t-1}^{-1/2} \left( \mathbf{I} + \mathbf{V}_{t-1}^{-1/2} \psi(\mathbf{x}_t) \psi(\mathbf{x}_t)^\top \mathbf{V}_{t-1}^{-1/2} \right)^{-1} \mathbf{V}_{t-1}^{-1/2} \phi(\mathbf{x}) \\
&\stackrel{(d)}{\succcurlyeq} \sigma^2 \phi(\mathbf{x})^\top \mathbf{V}_{t-1}^{-1} \phi(\mathbf{x}) / (1 + 1/\sigma^2) \\
&\stackrel{(e)}{=} \partial\sigma_{t-1}^2(\mathbf{x}) / (1 + 1/\sigma^2)
\end{aligned} \tag{27}$$

where (a) derives from (20) and (c) comes from the inversion of matrix product. Finally (d) follows from the result in (26) and (e) exploits the reformulation of  $\partial\sigma_{t-1}^2(\mathbf{x})$ .  $\square$

According to Lemma B.4 and Lemma B.6, the following holds for any  $\mathbf{x} \in \mathcal{X}$  and any  $t \geq 1$ ,

$$\frac{1}{1 + 1/\sigma^2} \leq \frac{\|\partial\sigma_t^2(\mathbf{x})\|_2}{\|\partial\sigma_{t-1}^2(\mathbf{x})\|_2} \leq 1. \tag{28}$$

Based on the definition of  $r$  in our Theorem 2, we therefore also have

$$r \triangleq \max_{\mathbf{x} \in \mathcal{X}, t \geq 1} \sqrt{\|\partial\sigma_t^2(\mathbf{x})\|_2 / \|\partial\sigma_{t-1}^2(\mathbf{x})\|_2} \in \left[1/\sqrt{1+1/\sigma^2}, 1\right]. \quad (29)$$

As a result, for every iteration  $t$  of our Algo. 2, we have

$$\begin{aligned} \sqrt{\|\partial\sigma_t^2(\mathbf{x})\|_2} &\leq r \sqrt{\|\partial\sigma_{t-1}^2(\mathbf{x})\|_2} \\ &\leq r^t \sqrt{\|\partial\sigma_0^2(\mathbf{x})\|_2} \\ &= r^t \sqrt{\|\partial_z \partial_{z'} k(\mathbf{z}, \mathbf{z}')|_{\mathbf{z}=\mathbf{z}'=\mathbf{x}}\|_2} \\ &\leq r^t \kappa \end{aligned} \quad (30)$$

where the last inequality derives from our assumption of  $\|\partial_z \partial_{z'} k(\mathbf{z}, \mathbf{z}')|_{\mathbf{z}=\mathbf{z}'=\mathbf{x}}\|_2 \leq \kappa^2$  ( $\forall \mathbf{x} \in \mathcal{X}$ ) in our Sec. 2.1. By substituting the result above into our Theorem 1, we complete our proof of Theorem 2.

#### B.4 PROOF OF THEOREM 3

**Preparation.** Following the definition of the derivative mapping on the true derivative  $\nabla f(\mathbf{x}_{t,\tau})$  in (8), we defined the following derivative mapping on our estimated derivative  $\nabla \mu_{t-1}(\mathbf{x}_{t,\tau})$ :

$$\widehat{G}_{t,\tau} \triangleq \frac{\mathbf{x}_{t,\tau} - \mathbf{x}_{t,\tau+1}}{\eta_{t,\tau}} = \frac{\mathbf{x}_{t,\tau} - \mathcal{P}_{\mathcal{X}}(\mathbf{x}_{t,\tau} - \eta_{t,\tau} \nabla \mu_t(\mathbf{x}_{t,\tau}))}{\eta_{t,\tau}}. \quad (31)$$

By re-arranging it, we have the following update rule that has reformulated (7):

$$\mathbf{x}_{t,\tau+1} = \mathbf{x}_{t,\tau} - \eta_{t,\tau} \widehat{G}_{t,\tau}. \quad (32)$$

Based on our definition of the derivative mappings in (31) and (8), we introduce the following lemmas:

**Lemma B.7** (General Projection Inequalities). *Given  $\mathcal{P}_{\mathcal{X}}(\mathbf{x}) = \arg \min_{\mathbf{z} \in \mathcal{X}} \|\mathbf{x} - \mathbf{z}\|_2^2/2$  and domain  $\mathcal{X}$ , for any  $\mathbf{x}, \mathbf{x}'$ , we have*

$$\|\mathbf{x} - \mathcal{P}_{\mathcal{X}}(\mathbf{x})\|_2 \leq \|\mathbf{x} - \mathcal{P}_{\mathcal{X}}(\mathbf{x}')\|_2, \quad (33)$$

$$\|\mathcal{P}_{\mathcal{X}}(\mathbf{x}) - \mathcal{P}_{\mathcal{X}}(\mathbf{x}')\|_2 \leq \|\mathbf{x} - \mathbf{x}'\|_2. \quad (34)$$

*Proof.* For (33), as  $\mathcal{P}_{\mathcal{X}}(\mathbf{x}') \in \mathcal{X}$  ( $\forall \mathbf{x}'$ ) and  $\mathcal{P}_{\mathcal{X}}(\mathbf{x}) = \arg \min_{\mathbf{z} \in \mathcal{X}} \|\mathbf{x} - \mathbf{z}\|_2^2/2$ , we then naturally have (33).

For (34), since  $\mathcal{P}_{\mathcal{X}}(\mathbf{x})$  is the optimum of  $h(\mathbf{z}) = \|\mathbf{x} - \mathbf{z}\|_2^2/2$ , according to the optimality condition of the convex projection function  $h(\mathbf{z})$  within the domain  $\mathbf{z} \in \mathcal{X}$  (Boyd and Vandenberghe, 2014), we then have the following inequality for any  $\mathcal{P}_{\mathcal{X}}(\mathbf{x}') \in \mathcal{X}$ :

$$\nabla h(\mathbf{z})^\top (\mathcal{P}_{\mathcal{X}}(\mathbf{x}') - \mathbf{z}) \geq 0. \quad (35)$$

By taking  $\nabla h(\mathbf{z}) = \mathbf{z} - \mathbf{x}$  with  $\mathbf{z} = \mathcal{P}_{\mathcal{X}}(\mathbf{x})$  into the inequality above, we have

$$(\mathcal{P}_{\mathcal{X}}(\mathbf{x}) - \mathbf{x})^\top (\mathcal{P}_{\mathcal{X}}(\mathbf{x}') - \mathcal{P}_{\mathcal{X}}(\mathbf{x})) \geq 0. \quad (36)$$

By exchanging  $\mathbf{x}$  and  $\mathbf{x}'$  in the result above, we achieve the following similar result:

$$(\mathcal{P}_{\mathcal{X}}(\mathbf{x}') - \mathbf{x}')^\top (\mathcal{P}_{\mathcal{X}}(\mathbf{x}) - \mathcal{P}_{\mathcal{X}}(\mathbf{x}')) \geq 0. \quad (37)$$

By summing (36) and (37),

$$(\mathbf{x} - \mathbf{x}')^\top (\mathcal{P}_{\mathcal{X}}(\mathbf{x}) - \mathcal{P}_{\mathcal{X}}(\mathbf{x}')) \geq \|\mathcal{P}_{\mathcal{X}}(\mathbf{x}) - \mathcal{P}_{\mathcal{X}}(\mathbf{x}')\|_2^2. \quad (38)$$

Based on the Cauchy-Schwarz inequality, we finally achieve (34) using

$$\begin{aligned} \|\mathcal{P}_{\mathcal{X}}(\mathbf{x}) - \mathcal{P}_{\mathcal{X}}(\mathbf{x}')\|_2^2 &\leq (\mathbf{x} - \mathbf{x}')^\top (\mathcal{P}_{\mathcal{X}}(\mathbf{x}) - \mathcal{P}_{\mathcal{X}}(\mathbf{x}')) \\ &\leq \|\mathbf{x} - \mathbf{x}'\|_2 \|\mathcal{P}_{\mathcal{X}}(\mathbf{x}) - \mathcal{P}_{\mathcal{X}}(\mathbf{x}')\|_2 \end{aligned} \quad (39)$$

where both sides need to be divided by  $\|\mathcal{P}_{\mathcal{X}}(\mathbf{x}) - \mathcal{P}_{\mathcal{X}}(\mathbf{x}')\|_2$  to complete our proof.  $\square$



**Lemma B.8** (Inequalities for Derivative Mappings). *Given (31) and (8), for every  $t$  and  $\tau$ , we have*

$$\left\| \widehat{G}_{t,\tau} \right\|_2^2 \leq \nabla \mu_{t-1}(\mathbf{x}_{t,\tau})^\top \widehat{G}_{t,\tau}, \quad (40)$$

$$\|G_{t,\tau}\|_2 \leq \|\nabla f(\mathbf{x}_{t,\tau})\|_2, \quad (41)$$

$$\left\| \widehat{G}_{t,\tau} - G_{t,\tau} \right\|_2 \leq \|\nabla \mu_{t-1}(\mathbf{x}_{t,\tau}) - \nabla f(\mathbf{x}_{t,\tau})\|_2. \quad (42)$$

*Proof.* For (40), let  $\widehat{\mathbf{x}}_{t,\tau} = \mathbf{x}_{t,\tau} - \eta_{t,\tau} \nabla \mu_{t-1}(\mathbf{x}_{t,\tau})$ , we then have

$$\begin{aligned} & \|\mathcal{P}_{\mathcal{X}}(\mathbf{x}_{t,\tau}) - \mathcal{P}_{\mathcal{X}}(\widehat{\mathbf{x}}_{t,\tau})\|_2^2 - (\mathbf{x}_{t,\tau} - \widehat{\mathbf{x}}_{t,\tau})^\top (\mathcal{P}_{\mathcal{X}}(\mathbf{x}_{t,\tau}) - \mathcal{P}_{\mathcal{X}}(\widehat{\mathbf{x}}_{t,\tau})) \\ \stackrel{(a)}{=} & \|\mathbf{x}_{t,\tau} - \mathbf{x}_{t,\tau+1}\|_2^2 - \eta_{t,\tau} \nabla \mu_{t-1}(\mathbf{x}_{t,\tau})^\top (\mathbf{x}_{t,\tau} - \mathbf{x}_{t,\tau+1}) \\ \stackrel{(b)}{=} & \eta_{t,\tau}^2 \left\| \widehat{G}_{t,\tau} \right\|_2^2 - \eta_{t,\tau}^2 \nabla \mu_{t-1}(\mathbf{x}_{t,\tau})^\top \widehat{G}_{t,\tau} \\ \stackrel{(c)}{\leq} & 0 \end{aligned} \quad (43)$$

where (a) results from the fact that  $\mathbf{x}_{t,\tau+1} = \mathcal{P}_{\mathcal{X}}(\mathbf{x}_{t,\tau} - \eta_{t,\tau} \nabla \mu_{t-1}(\mathbf{x}_{t,\tau}))$  based on our (7) and (b) derives from the definition of  $\widehat{G}_{t,\tau}$  in (31). In addition, (c) is based on the following result by substituting  $\mathbf{x} = \mathbf{x}_{t,\tau}$  and  $\mathbf{x}' = \widehat{\mathbf{x}}_{t,\tau}$  into (38):

$$\|\mathcal{P}_{\mathcal{X}}(\mathbf{x}_{t,\tau}) - \mathcal{P}_{\mathcal{X}}(\widehat{\mathbf{x}}_{t,\tau})\|_2^2 - (\mathbf{x}_{t,\tau} - \widehat{\mathbf{x}}_{t,\tau})^\top (\mathcal{P}_{\mathcal{X}}(\mathbf{x}_{t,\tau}) - \mathcal{P}_{\mathcal{X}}(\widehat{\mathbf{x}}_{t,\tau})) \leq 0. \quad (44)$$

Finally, by dividing  $\eta_{t,\tau}^2$  on the both sides of the last inequality in (43), we finish the proof for (40).

For (41), following the same proof above, we can also obtain achieve the following inequality for the projected derivative  $G_{t,\tau}$ :

$$\|G_{t,\tau}\|_2^2 \leq \nabla f(\mathbf{x}_{t,\tau})^\top G_{t,\tau} \leq \|\nabla f(\mathbf{x}_{t,\tau})\|_2 \|G_{t,\tau}\|_2. \quad (45)$$

We complete the proof for (41) by dividing  $\|G_{t,\tau}\|_2$  on the both sides of the inequality above.

For (42), define  $\mathbf{x}'_{t,\tau+1} \triangleq \mathbf{x}_{t,\tau} - \eta_{t,\tau} G_{t,\tau}$ , we have

$$\begin{aligned} \left\| \widehat{G}_{t,\tau} - G_{t,\tau} \right\|_2 & \stackrel{(a)}{=} \frac{1}{\eta_{t,\tau}} \left\| \mathbf{x}_{t,\tau} - \mathbf{x}_{t,\tau+1} - (\mathbf{x}_{t,\tau} - \mathbf{x}'_{t,\tau+1}) \right\|_2 \\ & \stackrel{(b)}{=} \frac{1}{\eta_{t,\tau}} \left\| \mathbf{x}_{t,\tau+1} - \mathbf{x}'_{t,\tau+1} \right\|_2 \\ & \stackrel{(c)}{=} \frac{1}{\eta_{t,\tau}} \left\| \mathcal{P}_{\mathcal{X}}(\mathbf{x}_{t,\tau} - \eta_{t,\tau} \nabla \mu_{t-1}(\mathbf{x}_{t,\tau})) - \mathcal{P}_{\mathcal{X}}(\mathbf{x}_{t,\tau} - \eta_{t,\tau} \nabla f(\mathbf{x}_{t,\tau})) \right\|_2 \\ & \stackrel{(d)}{\leq} \frac{1}{\eta_{t,\tau}} \left\| \mathbf{x}_{t,\tau} - \eta_{t,\tau} \nabla \mu_{t-1}(\mathbf{x}_{t,\tau}) - (\mathbf{x}_{t,\tau} - \eta_{t,\tau} \nabla f(\mathbf{x}_{t,\tau})) \right\|_2 \\ & \stackrel{(e)}{=} \|\nabla \mu_{t-1}(\mathbf{x}_{t,\tau}) - \nabla f(\mathbf{x}_{t,\tau})\|_2 \end{aligned} \quad (46)$$

where (a) comes from the definition of  $\widehat{G}_{t,\tau}$  and  $G_{t,\tau}$  in (31) and (8), respectively. In addition, (c) derives from (7) and (8). Finally, (d) results from (34).  $\square$

**Proof.** Since the objective function  $f$  is assumed to be  $L_s$ -Lipschitz smooth (Sec. 4.2), we have the following inequality for any  $\mathbf{x}_{t,\tau} \in \mathcal{X}$  in our ZORD algorithm:

$$f(\mathbf{x}_{t,\tau+1}) - f(\mathbf{x}_{t,\tau}) \leq \nabla f(\mathbf{x}_{t,\tau})^\top (\mathbf{x}_{t,\tau+1} - \mathbf{x}_{t,\tau}) + \frac{L_s}{2} \|\mathbf{x}_{t,\tau+1} - \mathbf{x}_{t,\tau}\|_2^2. \quad (47)$$

Let  $\delta' \in (0, 1)$ . Define  $\beta \triangleq \sqrt{d + 2(\sqrt{d} + 1) \ln(1/\delta')}$ , by substituting (32) into the inequality above, the following inequality holds with probability of at least  $1 - \delta'$ :

$$\begin{aligned}
& f(\mathbf{x}_{t,\tau+1}) - f(\mathbf{x}_{t,\tau}) \\
\stackrel{(a)}{\leq} & -\eta_{t,\tau} \nabla f(\mathbf{x}_{t,\tau})^\top \widehat{G}_{t,\tau} + \frac{L_s \eta_{t,\tau}^2}{2} \|\widehat{G}_{t,\tau}\|_2^2 \\
\stackrel{(b)}{=} & \eta_{t,\tau} (\nabla \mu_{t-1}(\mathbf{x}_{t,\tau}) - \nabla f(\mathbf{x}_{t,\tau}))^\top \widehat{G}_{t,\tau} - \eta_{t,\tau} \nabla \mu_{t-1}(\mathbf{x}_{t,\tau})^\top \widehat{G}_{t,\tau} + \frac{L_s \eta_{t,\tau}^2}{2} \|\widehat{G}_{t,\tau}\|_2^2 \\
\stackrel{(c)}{=} & \eta_{t,\tau} \left[ (\nabla \mu_{t-1}(\mathbf{x}_{t,\tau}) - \nabla f(\mathbf{x}_{t,\tau}))^\top (\widehat{G}_{t,\tau} - G_{t,\tau}) + (\nabla \mu_{t-1}(\mathbf{x}_{t,\tau}) - \nabla f(\mathbf{x}_{t,\tau}))^\top G_{t,\tau} \right] \\
& - \eta_{t,\tau} \nabla \mu_{t-1}(\mathbf{x}_{t,\tau})^\top \widehat{G}_{t,\tau} + \frac{L_s \eta_{t,\tau}^2}{2} \|\widehat{G}_{t,\tau}\|_2^2 \\
\stackrel{(d)}{\leq} & \eta_{t,\tau} \left[ \|\nabla \mu_{t-1}(\mathbf{x}_{t,\tau}) - \nabla f(\mathbf{x}_{t,\tau})\|_2 \|\widehat{G}_{t,\tau} - G_{t,\tau}\|_2 + \|\nabla \mu_{t-1}(\mathbf{x}_{t,\tau}) - \nabla f(\mathbf{x}_{t,\tau})\|_2 \|G_{t,\tau}\|_2 \right] \\
& - \eta_{t,\tau} \nabla \mu_{t-1}(\mathbf{x}_{t,\tau})^\top \widehat{G}_{t,\tau} + \frac{L_s \eta_{t,\tau}^2}{2} \|\widehat{G}_{t,\tau}\|_2^2 \\
\stackrel{(e)}{\leq} & \eta_{t,\tau} \left[ \|\nabla \mu_{t-1}(\mathbf{x}_{t,\tau}) - \nabla f(\mathbf{x}_{t,\tau})\|_2^2 + \|\nabla \mu_{t-1}(\mathbf{x}_{t,\tau}) - \nabla f(\mathbf{x}_{t,\tau})\|_2 \|\nabla f(\mathbf{x}_{t,\tau})\|_2 \right] \\
& - \frac{2\eta_{t,\tau} - L_s \eta_{t,\tau}^2}{2} \|\widehat{G}_{t,\tau}\|_2^2 \\
\stackrel{(f)}{\leq} & \eta_{t,\tau} \kappa^2 \beta^2 r^{2t} + \eta_{t,\tau} L_c \kappa \beta r^t - \frac{\eta_{t,\tau}}{2} \|\widehat{G}_{t,\tau}\|_2^2
\end{aligned} \tag{48}$$

where (d) derives from the Cauchy-Schwarz inequality and (e) follows from the Lemma B.7. Finally, (f) result from the bounded derivative estimation error in Theorem 2 and the fact that  $f$  is  $L_c$ -Lipschitz continuous (i.e.,  $\|\nabla f(\mathbf{x})\|_2 \leq L_c$  for any  $\mathbf{x} \in \mathcal{X}$ ) and  $\eta_{t,\tau} \leq 1/L_s$  ( $\forall \tau$ ).

For every iteration  $t$  our ZORD algorithm, we in fact will apply the virtual updates (7) for  $V_t$  times (see Algo. 2). Therefore, for probability  $\geq 1 - V_t \delta'$ , we have

$$\begin{aligned}
\frac{1}{V_t} \sum_{\tau=0}^{V_t-1} \eta_{t,\tau} \|\widehat{G}_{t,\tau}\|_2^2 & \leq \frac{2}{V_t} \sum_{\tau=0}^{V_t-1} [f(\mathbf{x}_{t,\tau}) - f(\mathbf{x}_{t,\tau+1}) + \eta_{t,\tau} (\kappa^2 \beta^2 r^{2t} + L_c \kappa \beta r^t)] \\
& = \frac{2}{V_t} [f(\mathbf{x}_{t-1}) - f(\mathbf{x}_t)] + \left( \frac{2}{V_t} \sum_{\tau=0}^{V_t-1} \eta_{t,\tau} \right) (\kappa^2 \beta^2 r^{2t} + L_c \kappa \beta r^t)
\end{aligned} \tag{49}$$

where the first inequality results from (48) by re-arranging it and then sum it up over  $\tau$ .

However, in order to prove the convergence of our ZORD algorithm to a stationary point, we need to consider the derivative mapping of  $G_{t,\tau}$  instead (refer to our Sec. 4.2). So, for any  $\tau$ , we propose the following inequality:

$$\begin{aligned}
\|G_{t,\tau}\|_2 & = \|G_{t,\tau} - \widehat{G}_{t,\tau} + \widehat{G}_{t,\tau}\|_2 \\
& \leq \|G_{t,\tau} - \widehat{G}_{t,\tau}\|_2 + \|\widehat{G}_{t,\tau}\|_2 \\
& \leq \|\nabla \mu_{t-1}(\mathbf{x}_{t,\tau}) - \nabla f(\mathbf{x}_{t,\tau})\|_2 + \|\widehat{G}_{t,\tau}\|_2 \\
& \leq \kappa \beta r^t + \|\widehat{G}_{t,\tau}\|_2
\end{aligned} \tag{50}$$

where the first inequality is from the Cauchy-Schwarz inequality and the second inequality comes from (42). Finally, by taking the result above into (49), we have

$$\frac{1}{V_t} \sum_{\tau=0}^{V_t-1} \eta_{t,\tau} \|G_{t,\tau}\|_2^2 \leq \frac{2}{V_t} [f(\mathbf{x}_{t-1}) - f(\mathbf{x}_t)] + \left( \frac{2}{V_t} \sum_{\tau=0}^{V_t-1} \eta_{t,\tau} \right) (\kappa^2 \beta^2 r^{2t} + L_c \kappa \beta r^t) + \kappa \beta r^t. \tag{51}$$

Then, substituting  $V_t = V$  and  $\eta_{t,\tau} = \eta$  for any  $t, \tau$  into the result above, the following inequality holds with probability of at least  $1 - VT\delta'$  when  $r < 1$ :

$$\begin{aligned}
\frac{1}{T} \sum_{t=1}^T \frac{1}{V} \sum_{\tau=0}^{V-1} \eta \|G_{t,\tau}\|_2^2 &\stackrel{(a)}{\leq} \frac{1}{T} \sum_{t=1}^T \left( \frac{2(f(\mathbf{x}_{t-1}) - f(\mathbf{x}_t))}{V} + 2\eta\kappa^2\beta^2r^{2t} + (2\eta L_c + 1)\kappa\beta r^t \right) \\
&\stackrel{(b)}{\leq} \frac{2}{TV} [f(\mathbf{x}_0) - f(\mathbf{x}_T)] + \frac{2\eta(1 - r^{2T})}{T(1 - r^2)} \kappa^2\beta^2r^2 \\
&\quad + \frac{(2\eta L_c + 1)(1 - r^T)}{T(1 - r)} \kappa\beta r \\
&\stackrel{(c)}{\leq} \frac{2}{TV} [f(\mathbf{x}_0) - f(\mathbf{x}^*)] + \frac{2\eta\kappa^2\beta^2r^2}{T(1 - r^2)} + \frac{(2\eta L_c + 1)\kappa\beta r}{T(1 - r)}.
\end{aligned} \tag{52}$$

Note that (b) derives from the summation of the geometric sequence about  $r$  and (c) comes from  $\mathbf{x}^* \triangleq \arg \min_{\mathbf{x} \in \mathcal{X}} f(\mathbf{x})$ . When  $r = 1$ , the following holds with probability of at least  $\geq 1 - VT\delta'$  accordingly:

$$\begin{aligned}
\frac{1}{T} \sum_{t=1}^T \frac{1}{V} \sum_{\tau=0}^{V-1} \eta \|G_{t,\tau}\|_2^2 &\leq \frac{1}{T} \sum_{t=1}^T \left( \frac{2(f(\mathbf{x}_{t-1}) - f(\mathbf{x}_t))}{V} + 2\eta\kappa^2\beta^2r^{2t} + (2\eta L_c + 1)\kappa\beta r^t \right) \\
&= \frac{2}{TV} [f(\mathbf{x}_0) - f(\mathbf{x}_T)] + 2\eta\kappa^2\beta^2 + (2\eta L_c + 1)\kappa\beta.
\end{aligned} \tag{53}$$

Finally, let  $\delta = VT\delta' \in (0, 1)$ , the following holds with probability of at least  $1 - \delta$ ,

$$\begin{aligned}
\min_{t \leq T} \frac{1}{V} \sum_{\tau=0}^{V-1} \|G_{t,\tau}\|_2^2 &\leq \frac{1}{T} \sum_{t=1}^T \frac{1}{V} \sum_{\tau=0}^{V-1} \|G_{t,\tau}\|_2^2 \\
&\leq \textcircled{1} + \textcircled{2}
\end{aligned} \tag{54}$$

where  $\textcircled{1}$  and  $\textcircled{2}$  can be defined as below with  $\alpha \triangleq \kappa\sqrt{d + 2(\sqrt{d} + 1) \ln(VT/\delta)}$ :

$$\begin{aligned}
\textcircled{1} &= \frac{2/\eta}{TV} [f(\mathbf{x}_0) - f(\mathbf{x}_T)] \\
\textcircled{2} &= \begin{cases} 2\alpha^2r^2 / [T(1 - r^2)] + (2L_c + 1/\eta)\alpha r / [T(1 - r)] & (r < 1), \\ 2\alpha^2 + (2L_c + 1/\eta)\alpha & (r = 1). \end{cases}
\end{aligned} \tag{55}$$

## APPENDIX C EXPERIMENTAL SETTINGS

### C.1 GENERAL SETTINGS

**Derived GP.** Among all our experiments in Sec. 5, to apply the derivative estimation in Sec. 3.1 for every iteration  $t$  and every step  $\tau$  of our ZORD algorithm, we use the derived GP (4) based on the Matérn kernel with  $\nu = 2.5$  and fit this derived GP using 150 queries that achieves the smallest Euclidean distance with input  $\mathbf{x}_{t,\tau}$  from the optimization trajectory. This is because we only need to model the objective function  $f$  in the vicinity of input  $\mathbf{x}_{t,\tau}$  precisely rather than the entire domain, so as to achieve an accurate derivative estimation at input  $\mathbf{x}_{t,\tau}$ .

**Confidence Threshold.** Among all our experiments in Sec. 5, the confidence threshold  $c$  of our dynamic virtual updates (Sec. 3.2) is set to be 0.35 in order to realize a good trade-off between query efficiency and accurate derivative estimation in practice, which can already allow our ZORD to achieve compelling empirical results consistently (see our Sec. 5). In light of this,  $c = 0.35$  would be a reasonably good choice in practice, especially when there is no prior knowledge about the objective functions. When we have prior knowledge about the smoothness of the objective functions, we can likely make a better choice for  $c$ : Intuitively, smooth objective functions usually can be modeled by the Gaussian process effectively (Rasmussen and Williams, 2006), so an accurate derivative estimation from our derived GP is also likely to be achieved. In this scenario, a large confidence threshold can be applied to fully exploit the benefit of our derivative estimation that is free from the requirement for additional queries and consequently results in an improved query efficiency in practice.

**Baselines.** In addition, among all our experiments in Sec. 5, we consistently use  $n = 10$ ,  $\lambda = 0.01$  and directions  $\{\mathbf{u}_i\}_{i=1}^n$  that are randomly sampled from a unit sphere for the derivative estimation of the FD method (2) applied in the RGF and PRGF algorithm. Moreover, following the common practice of (Berahas et al., 2022; Cheng et al., 2021), we conduct orthogonalization on these randomly selected directions via the Gram-Schmidt procedure. As for the ES algorithm (e.g., the one applied in (Salimans et al., 2017)), we apply the same  $n$ ,  $\lambda$  and  $\{\mathbf{u}_i\}_{i=1}^n$  in RGF and PRGF for their update in every iteration.

**Domain Transformation.** Following the practice that has been used in (Eriksson et al., 2019), for all our experiments, we firstly re-scale the input domains into  $[0, 10]^d$  to ease the optimization and then re-scale the updated inputs back to the original domains for querying.

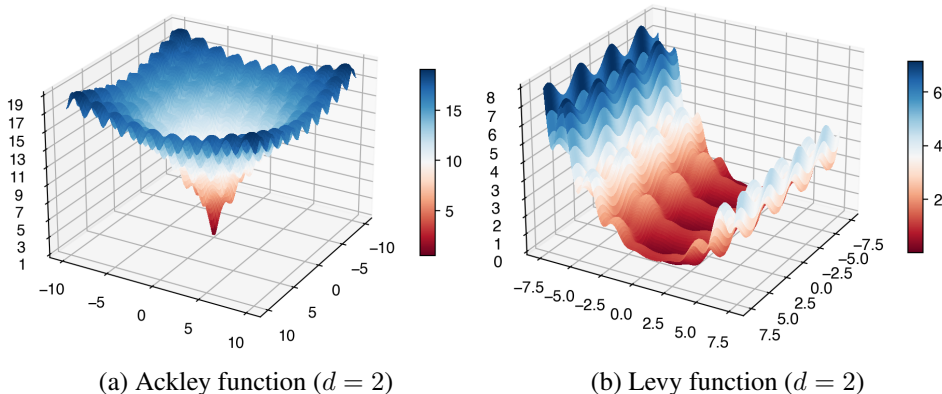
### C.2 SYNTHETIC EXPERIMENTS

Let input  $\mathbf{x} = [x_i]_{i=1}^d$ , the Ackley and Levy function applied in our synthetic experiments are given below,

$$\begin{aligned}
 f(\mathbf{x}) &= -20 \exp \left( -0.2 \sqrt{\frac{1}{d} \sum_{i=1}^d x_i^2} \right) - \exp \left( \frac{1}{d} \sum_{i=1}^d \cos(2\pi x_i) \right) + 20 + \exp(1), \text{ (Ackley)} \\
 f(\mathbf{x}) &= \sin^2(\pi w_1) + \sum_{i=1}^{d-1} (w_i - 1)^2 [1 + 10 \sin^2(\pi w_i + 1)] + (w_d - 1)^2 [1 + \sin^2(2\pi w_d)] \text{ (Levy)}
 \end{aligned}
 \tag{56}$$

where  $w_i = 1 + (x_i - 1)/4$  for any  $i = 1, \dots, d$ , Ackley function achieves its minimum (i.e.,  $\min f(\mathbf{x}) = 0$ ) at  $\mathbf{x}^* = \mathbf{0}$ , and Levy function achieves its minimum (i.e.,  $\min f(\mathbf{x}) = 0$ ) at  $\mathbf{x}^* = \mathbf{1}$ . Note that the Ackley and Levy function for the synthetic experiments in our Sec. 5.2 are defined within the domain  $[-20, 20]^d$  and  $[-7.5, 7.5]^d$ , respectively. To give a better understanding of these two synthetic functions, we provide a 3D illustration of these two synthetic functions with  $d = 2$  in our Fig. 5. As shown in Fig. 5, these two synthetic functions are highly nonconvex and therefore have local minimums within their domains.

To compare our ZORD algorithm with other ZO/FO optimization baselines in Sec. 5.2, we firstly employ TURBO of 300 queries to find a good initialization for all other ZO/FO optimization algorithms in Fig. 3 because of the nonconvexity of these two synthetic functions as shown in Fig. 5. We then

Figure 5: The 3D illustration of Ackley and Levy synthetic function with  $d = 2$ .

apply these ZO/FO optimization algorithms with a query budget of 200 for  $d = 20, 40$ , and a query budget of 400 for  $d = 100$  to compare their query efficiency. We use the same Adam optimizer (Kingma and Ba, 2015) with a learning rate of 0.1 and exponential decay rates of 0.9, 0.999 for RGF, PRGF, GD, and our ZORD algorithm, for faster convergence compared with standard GD.

### C.3 BLACK-BOX ADVERSARIAL ATTACK

For the black-box adversarial attack experiment on the MNIST dataset, we use the same fully trained deep neural networks from (Cheng et al., 2021) and adopt a  $L_\infty$  constraint of  $\|\mathbf{x}\|_\infty \leq 0.3$  on the input perturbation  $\mathbf{x}$ . For the black-box adversarial attack experiment on the CIFAR-10 dataset, we fully train a ResNet-18 (He et al., 2016) on CIFAR-10 using stochastic gradient descent (SGD) with a cosine annealed learning rate from 0.1 to 0, a momentum of 0.9 and a weight decay of  $5 \times 10^{-4}$  for 200 epochs, and adopt a  $L_\infty$  constraint of  $\|\mathbf{x}\|_\infty \leq 0.2$  on the input perturbation  $\mathbf{x}$ . Note that we use the same loss function as (Cheng et al., 2021) for these two experiments. Meanwhile, to apply RGF, PRGF and our ZORD, we adopt Adam optimizer with the same learning rate of 0.5 and the same exponential decay rates of 0.9, 0.999.

### C.4 NON-DIFFERENTIABLE METRIC OPTIMIZATION

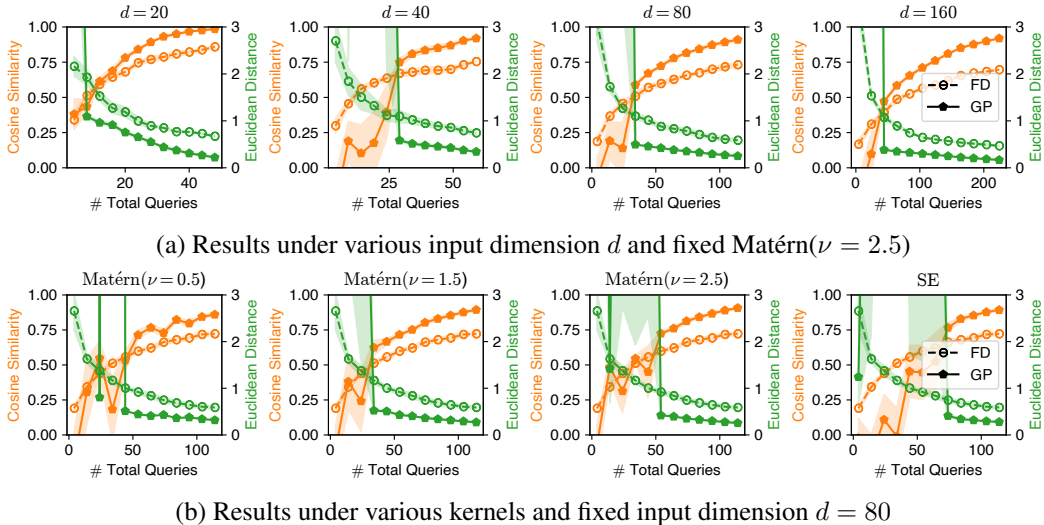
The Covertypes dataset used in Sec. 5.4 is a classification dataset consisting of 581,012 samples from 7 different categories. Each sample from this dataset is a 54-dimensional vector of integers. In this experiment, we randomly split the dataset into training and test sets with each containing 290,506 samples. The MLP classifier applied in Sec. 5.4 consists of 2 layers with 30 and 14 hidden neurons respectively, leading to 2189 parameters in total (i.e.,  $d = 2189$ ). We first train this MLP classifier on the training dataset of Covertypes using the L-BFGS algorithm with the cross-entropy loss function for 300 epochs, and then apply ZO optimization algorithms to fine-tune our trained MLP directly on the non-differentiable metrics (i.e., using these metrics as the new loss functions), including precision, recall, F1 score and Jaccard index. To obtain the results of ES, RGF, PRGF and our ZORD algorithm in Sec. 5.4, we apply the same Adam optimizer with a learning rate of 0.2 (for precision and recall) or 0.01 (for F1 score and Jaccard index) and exponential decay rates of 0.9, 0.999. Note that standard BO algorithms (including TuRBO) fail to achieve any percentage improvements (i.e., achieving 0% in the  $y$ -axis of Fig. 4) in this experiment according to our five independent runs, which is likely due to their aggressive exploration in the input domain of such a high dimension. In light of this, we do not include them in our comparison since all other methods are able to achieve certain improvements.

### C.5 DERIVATIVE-FREE REINFORCEMENT LEARNING

Our derivative-free RL experiments aim to learn controllers (which outputs policies) that maximize the rewards/return for several environments in the OpenAI Gym (Brockman et al., 2016) without using true derivatives. Specifically, we need to optimize the parameters (i.e.,  $\mathbf{x}$ ) of our neural network

Table 2: OpenAI Gym environment properties and their respective network dimensions.

|       | Acrobot | Swimmer | Lunar | BipedalWalker | Walker2D | HalfCheetah |
|-------|---------|---------|-------|---------------|----------|-------------|
| $ S $ | 6       | 8       | 8     | 24            | 17       | 17          |
| $ A $ | 3       | 2       | 4     | 4             | 6        | 6           |
| $d$   | 213     | 222     | 244   | 404           | 356      | 356         |

Figure 6: Comparison of the derivative estimation errors of our derived GP-based estimator (GP) and the FD estimator under various input dimensions and kernels. Similarly, each result is reported with the mean  $\pm$  standard error from five independent runs.

(MLP) controller with 2 hidden layers, where each hidden layer has 10 hidden neurons and one bias term. We adopt a  $L_\infty$  constraint of  $\|x\|_\infty \leq 1$  on the parameters  $x$ . We use a softmax output layer for the policies that deal with discrete action spaces, and a tanh output layer for the policies that deal with continuous action spaces. The dimension of neural network parameters (represented as a column vector)  $d$  is determined by the dimensions of both the observation  $|S|$  and the action space  $|A|$  of an environment, as detailed in Tab. 2.

In order to search for policies that are robust to different random state initializations, we use the vectorized API of OpenAI Gym, and our observed function value  $y(x)$  given the network parameters  $x$  is an averaged return of 32 parallel environments. We also fix the seed of OpenAI Gym for all queries, which ensures that we are evaluating on a fixed set of 32 state initializations and that our results can be reproduced. We first initialize a sample of 500 points from a Latin Hypercube (McKay et al., 1979) to find a good initial input, and then proceed to apply ZO optimization algorithms (i.e., ES, RGF, PRGF, and our ZORD) with the same query budget of 1000 on this initial input. For all these ZO optimization algorithms, we employ the same Adam optimizer with a learning rate of 1.0 and exponential decay rates of 0.9, 0.999. Considering the prohibitive noise in RL experiments, we use 300 queries from the optimization trajectory that has the smallest Euclidean distance with an input needing to be updated. Of note, we conduct 10 trials in total where each trial differs from each other by both the OpenAI Gym seed and the Latin Hypercube initializations.

## APPENDIX D MORE RESULTS

### D.1 MORE RESULTS ON DERIVATIVE ESTIMATION

Besides the comparison in Fig. 2, we provide additional comparison between our derived GP-based estimator (6) and the FD estimator (2) under various input dimensions in Fig. 6(a) and various kernels

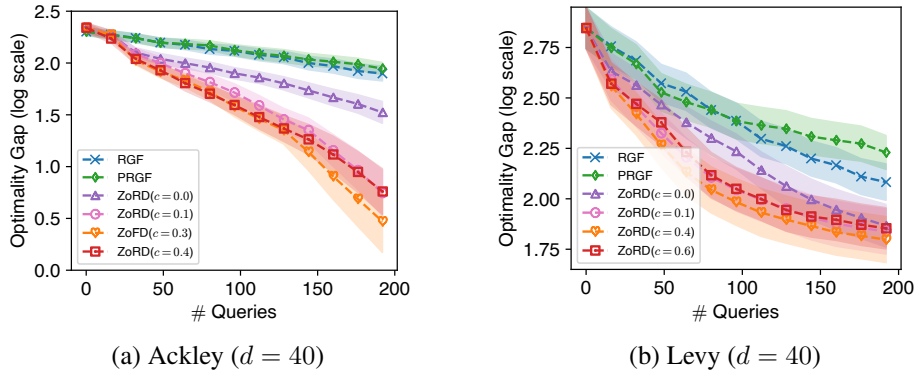


Figure 7: Comparison of our ZORD algorithm using different confidence thresholds  $c$  for its dynamic virtual updates, where the  $x$ -axis and the  $y$ -axis denote the number of function queries and the log-scaled optimality gap (i.e.,  $\log(f(\mathbf{x}_T) - f(\mathbf{x}^*))$ ) achieved with this number of queries, respectively.

in Fig. 6(b) using the Ackley function. We adopt the same setting in Sec. 5.2. Interestingly, Fig. 6(a)(b) show that under various input dimensions and GP kernels, our derived GP-based estimator (6) is still able to achieve faster reduction rates compared with the FD estimator. Of note, all the function queries applied in our derived GP-based estimator is from the optimization trajectory whereas the FD estimator requires additional function queries for its derivative estimation. So, Fig. 6(a)(b) also show that our derived GP method is still able to achieve improved query efficiency for accurate derivative estimation than FD method under various input dimensions and GP kernels because our method avoids the requirement of additional queries for derivative estimation. Interestingly, the objective function (i.e., the Ackley function) is not truly sampled from the GPs based on these kernels. This therefore means that though we have assumed that we need the prior knowledge about the GP in which the objective function is sampled from (Sec. 2.1), such an assumption does not really need to be satisfied for our derived GP-based method to achieve accurate derivative estimation in practice. More interestingly, we notice that Matérn( $\nu = 0.5$ ) and SE kernel will achieve slightly worse derivative estimation, indicating that the choice of GP kernels may impact the quality of our derived GP-based derivative estimation. However, in practice, our derived GP method based on Matérn( $\nu = 2.5$ ) kernel, which has been widely adopted in our experiments, is already able to provide us with good derivative estimation for ZO optimization as confirmed by the results in our other experiments.

## D.2 MORE RESULTS ON SYNTHETIC EXPERIMENTS

In this section, we compare ZORD with more baselines in Fig. 8. Notably, we mainly compare our ZORD with CobBO (based on the code implementation provided by (Tan et al., 2021)) since CobBO generally performs better than other baselines, e.g., TPE, ATPE, and BADS according to (Tan et al., 2021). As shown in the results in Fig. 8, our ZoRD algorithm is still able to outperform the other benchmark BO algorithm (i.e., CobBO).

We then investigate the impacts of the dynamic virtual updates (Sec. 3.2) on our ZORD algorithm. In particular, we apply the same setting in Appx. C.2 to optimize the Ackley and Levy function with  $d = 40$  under various confidence thresholds  $c$  for our dynamic virtual updates. Fig. 7 illustrates the results. As shown in both Fig. 7(a) and (b), our ZORD algorithm using the technique of dynamic virtual updates (i.e.,  $c > 0$ ) can consistently achieve improved query efficiency compared with the one not using the technique of dynamic virtual updates (i.e.,  $c = 0$ ). This indicates the essence of dynamic virtual updates in helping improve the query efficiency of our ZORD algorithm. Such a result actually corroborates our theoretical insights about virtual updates (Sec. 4.2). Remarkably, our ZORD algorithm without the technique of dynamic virtual updates (i.e.,  $c = 0$ ) is still able to achieve both improved query efficiency and better converged performance compared with RGF and PRGF, which further verifies the superiority of our derived GP-based derivative estimation. More interestingly, both Fig. 7(a) and Fig. 7(b) have verified that there indeed exists a trade-off for the confidence threshold  $c$  as we have discussed in Sec. 3.2: The confidence threshold  $c$  can not be overly

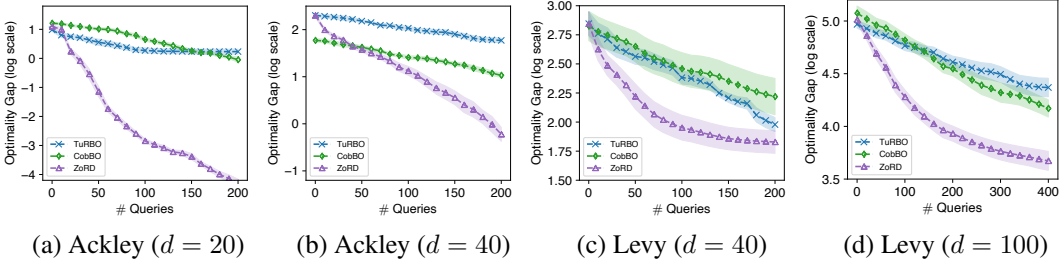


Figure 8: Additional comparison between our ZORD and other baselines. The  $x$ -axis and  $y$ -axis denote the number of queries and log-scaled optimality gap (i.e.,  $\log(f(x_T) - f(x^*))$ ) achieved after this number of queries. Each curve is the mean  $\pm$  standard error from ten independent runs.

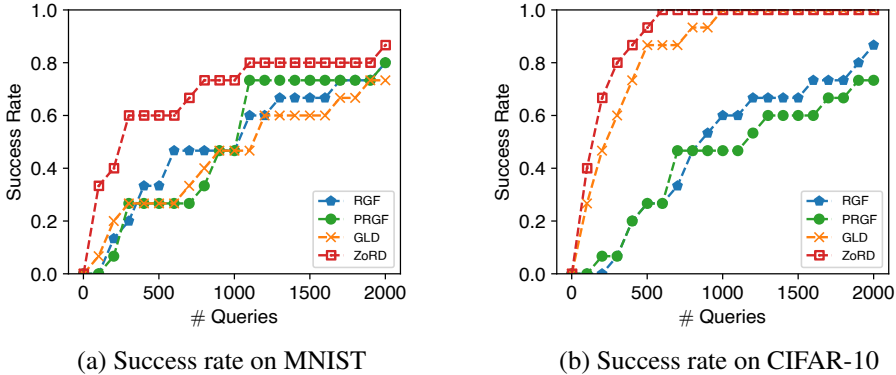


Figure 9: Comparison of the success rate achieved by various ZO optimization algorithms on the 15 images selected from MNIST and CIFAR-10 dataset. Note that the  $x$ -axis and the  $y$ -axis denote the number of queries and the success rate (within the range of  $[0, 1]$ ) achieved after this number of queries, respectively.

small or excessively large in order to achieve the best query efficiency of our ZORD algorithm, e.g.,  $c = 0.3$  for Ackley ( $d = 40$ ) and  $c = 0.4$  for Levy ( $d = 40$ ).

### D.3 MORE RESULTS ON BLACK-BOX ADVERSARIAL ATTACK

Besides the comparison in our Sec. 5.3, we also compare the success rate achieved by different ZO optimization algorithms on the 15 images selected from MNIST or CIFAR-10 in Fig. 9. Note that we adopt the same settings in Appx. C.3 for this comparison. Considering the large computational complexity of Turbo-1/10 algorithm for hard-to-attack images<sup>3</sup> which is usually undesirable in practice, we drop the comparison with them in this experiment. Fig. 9 shows that under the same query budget, our ZORD algorithm is able to achieve considerably improved success rate over other ZO optimization algorithms. These results therefore further support the superior query efficiency of our ZORD algorithm in real-world challenging problems.

### D.4 MORE RESULTS FOR DERIVATIVE-FREE REINFORCEMENT LEARNING

Recent years have also witnessed a surging interest in derivative-free reinforcement learning (Salimans et al., 2017; Qian and Yu, 2021), where ZO optimization algorithms are widely applied. In light of this, we also demonstrate the superiority of our ZORD algorithm in the problem of derivative-free reinforcement learning. Specifically, we adopt the setting in Sec. C.5 to experiment in different RL environments. Tab. 3 summarizes the comparison among different ZO optimization algorithms under

<sup>3</sup>Bayesian optimization algorithms, including Turbo-1/10, are widely known to suffer from the prohibitive computational complexity when they need a large number of function queries for optimization, e.g.,  $T > 1000$  (Rasmussen and Williams, 2006).



Table 3: Comparison of the rewards (larger is better) achieved by various ZO optimization algorithms in different RL environments. Each result is reported with the mean  $\pm$  standard deviation from ten independent runs.

| Algorithm | Acrobot                         | Swimmer                          | Lunar                            | BipedalWalker                   | Walker2D                          | HalfCheetah                        |
|-----------|---------------------------------|----------------------------------|----------------------------------|---------------------------------|-----------------------------------|------------------------------------|
| ES        | -86.2 $\pm$ 11.0                | 176.0 $\pm$ 56.8                 | -94.7 $\pm$ 24.4                 | -34.7 $\pm$ 27.3                | 340.4 $\pm$ 143.0                 | 1042.4 $\pm$ 753.9                 |
| RGF       | -83.0 $\pm$ 5.6                 | 213.2 $\pm$ 65.1                 | -93.8 $\pm$ 19.1                 | -30.3 $\pm$ 40.3                | 368.4 $\pm$ 223.1                 | 1129.3 $\pm$ 748.5                 |
| PRGF      | -86.3 $\pm$ 9.9                 | 218.6 $\pm$ 66.2                 | -100.1 $\pm$ 16.0                | -29.9 $\pm$ 35.2                | 344.6 $\pm$ 152.3                 | 1083.3 $\pm$ 722.2                 |
| ZoRD      | <b>-73.3<math>\pm</math>2.4</b> | <b>280.5<math>\pm</math>77.6</b> | <b>-45.1<math>\pm</math>38.3</b> | <b>12.9<math>\pm</math>37.8</b> | <b>729.1<math>\pm</math>304.2</b> | <b>1950.5<math>\pm</math>576.1</b> |

the same query budget of 1000. As BO algorithms usually suffer from the prohibitive computational complexity for a large  $T$  (Rasmussen and Williams, 2006) and GLD has never been applied in RL, we mainly compare our ZoRD algorithm with ES, RGF and PRGF, which also belongs to the same type of ZO optimization algorithm: GD with estimated derivative. Remarkably, Tab. 3 shows that under the same query budget, our ZoRD algorithm can consistently enjoy improved performance (i.e., highest rewards) than the other ZO optimization algorithms in different RL environments. This further supports the superiority of our ZoRD algorithm to other FD-based ZO optimization algorithms.

## APPENDIX E DISCUSSIONS

### E.1 ZoRD vs. FD-BASED ZO OPTIMIZATION

Of note, the novelty of our work in fact lies in its way of exploiting the GP assumption to help design an improved derivative estimation and hence an improved ZO optimization algorithm, which to the best of our knowledge has not been explored theoretically yet in the field of ZO optimization via GD with estimated derivative. That is, at this moment, it is still not known in the literature how existing FD methods can utilize such an assumption to achieve better derivative estimation (i.e., their derivative estimation quality will remain the same), even when they make the same assumption as us. In light of this, the comparison between our derived GP method and the FD method in Sec. 4 is not only necessary but also meaningful to show the advantage of exploiting such an assumption in ZO derivative estimation. Importantly, our empirical results further show that such an assumption is in fact not restrictive for our ZoRD to achieve compelling performance in practice. For example, our Fig. 2 and Fig. 6 have shown that our derived GP-based method is able to achieve smaller derivative estimation error than the FD method when the objective functions are not designed to be sampled from a GP with the kernel that we had applied for our derivative estimation. Moreover, the results in our Sec. 5.2, 5.3, 5.4 have shown that our ZoRD is capable of achieving competitive optimization performance for real-world optimization problems where the objective functions are also not designed to be sampled from a GP with the kernel that we had used for our ZoRD.

Meanwhile, the theoretical challenges of our work lie in the theoretical guarantee on the derivative estimation error of our unique derived GP-based method for any input in the domain as well as the convergence analysis based on such a unique derivative estimation, which to the best of our knowledge have not been studied in the literature. This means that our Thm. 1 and Thm. 2 have provided new developments in the analysis of gradient estimation error and our Thm. 3 will be the first convergence result for GD using our unique derivative estimation method. Interestingly, the bound in our Thm. 3 also improves over the standard ones from (Nesterov and Spokoiny, 2017; Liu et al., 2018b) in several aspects, as discussed in our Sec. 4.2.

### E.2 ZoRD vs. BO

Our ZoRD algorithm and standard BO algorithms (e.g., GP-UCB) have in fact applied the same GP assumption for their algorithm design. That is, however, where the similarity ends. Of note, our ZoRD exploits such an assumption to derive a specific GP (i.e., (4)) for derivative estimation, which is then employed for local exploitation via (projected) GD update. In contrast, BO algorithms utilize such an assumption to construct their acquisition functions for a global optimization that can trade off between exploitation and exploration. In practice, the exploration of BO algorithms is usually query-inefficient, especially for problems with high-dimensional input spaces, and therefore GD with

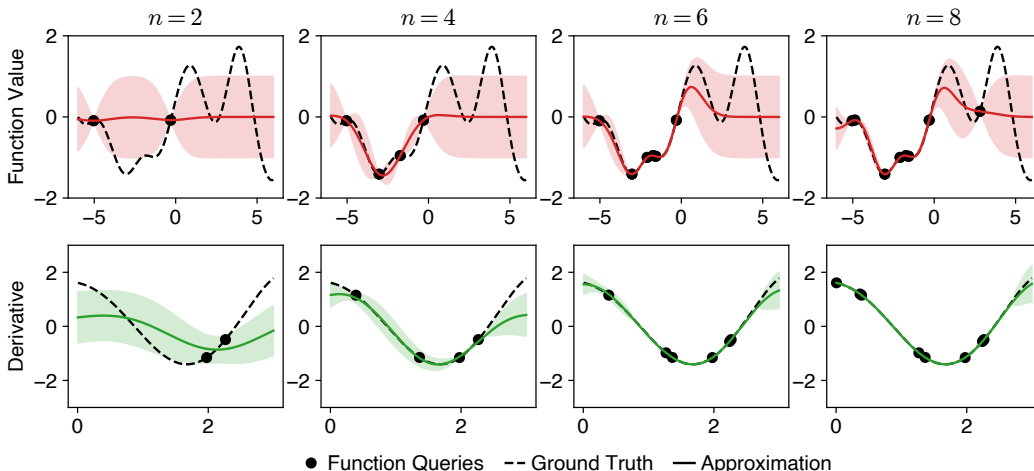


Figure 10: Comparison of local derivative estimation (in the input domain of  $[0, 3]$ ) in our ZORD and global function approximation (in the input domain of  $[-6, 6]$ ) in BO under various number of random function queries.

estimated derivatives (especially our ZORD) is preferred to realize better optimization performances in these problems (see our Sec. 5.2). So, our ZORD and BO algorithms belong to two different types of ZO optimization algorithms (i.e., GD-type vs. BO-type), where their theoretical analyses are in fact not comparable. In particular, GD-type and BO-type ZO optimization algorithms apply different metrics for their theoretical analyses, e.g., the derivative estimation error as well as the convergence to a stationary point (in the nonconvex case) for GD-type ZO optimization algorithms vs. the global asymptotic convergence in terms of the regret for BO-type ZO optimization algorithms. So, it is more reasonable to compare the theory (including the theoretical challenge, the new developments, and the novelty of the convergence result) of our ZORD with other GD-type ZO optimization algorithms, e.g., the ones using FD methods for their derivative estimation (Nesterov and Spokoiny, 2017; Liu et al., 2018b), as what we have discussed in Sec. E.1.

In addition, in contrast to using the GP to model the objective function within the *entire* domain for global *exploration* in BO, our derived GP in ZORD will be applied to estimate the derivative of the objective function for local *exploitation* by GD as shown in Sec. 3.1. As GD typically optimizes in a local region, our derived GP only needs to estimate the derivative *locally*, which is known to be much simpler than modeling the objective function within the *entire* domain in BO especially for objective functions in high-dimensional input spaces. In light of this, the derived GP for derivative estimation (4) in our ZORD algorithm advances the standard GP in BO in the following aspects:

1. **Improved Query Efficiency for Estimation.** The derived GP in our ZORD algorithm requires fewer function queries to provide accurate derivative estimation. We provide a visual example in Fig. 10, in which we sample a one-dimensional function  $f$  from a GP prior  $\mathcal{GP}(0, k(x, x))$  using the standard SE kernel and then randomly select the same number of queries from the input domain of  $[-6, 6]$  and  $[0, 3]$  for standard GP and our derived GP, respectively. As illustrated in Fig. 10, function in a local region (i.e.,  $x \in [0, 3]$ ) is usually smoother than its counterpart in the entire domain (i.e.,  $x \in [-6, 6]$ ). As a result, with only 4 function queries, our derived GP can already provide accurate estimation to the derivative of this objective function whereas standard GP requires more than 8 function queries to model this objective function accurately in the entire domain.
2. **Reduced Computational Complexity.** Comparing (3) and (5), both the derived GP for derivative estimation in our ZORD algorithm and the standard GP in BO enjoy a computational complexity of  $\mathcal{O}(n^3)$  with  $n$  function queries. However, as a consequence of the improved query efficiency of our derived GP, it is able to require fewer function queries (i.e.,

smaller  $n$ ) for accurate derivative estimation<sup>4</sup> and hence can enjoy a reduced computational complexity in practice especially when a large number of queries (e.g.,  $n > 1000$ ) are applied to the standard GP in BO.

---

<sup>4</sup>As introduced in our Appx. C, 150 function queries for our derived GP can already help our ZORD algorithm to achieve remarkable results in practice (refer to the experiments in our Sec. 5).