

---

# DPZero: Private Fine-Tuning of Language Models without Backpropagation

---

Liang Zhang<sup>1</sup> Bingcong Li<sup>1</sup> Kiran Koshy Thekumparampil<sup>2</sup> Sewoong Oh<sup>3</sup> Niao He<sup>1</sup>

## Abstract

The widespread practice of fine-tuning large language models (LLMs) on domain-specific data faces two major challenges in memory and privacy. First, as the size of LLMs continues to grow, the memory demands of gradient-based training methods via backpropagation become prohibitively high. Second, given the tendency of LLMs to memorize training data, it is important to protect potentially sensitive information in the fine-tuning data from being regurgitated. Zeroth-order methods, which rely solely on forward passes, substantially reduce memory consumption during training. However, directly combining them with standard differentially private gradient descent suffers more as model size grows. To bridge this gap, we introduce DPZERO, a novel private zeroth-order algorithm with nearly dimension-independent rates. The memory efficiency of DPZERO is demonstrated in privately fine-tuning RoBERTa and OPT on several downstream tasks. Our code is available at <https://github.com/Liang137/DPZero>.

## 1. Introduction

Fine-tuning pretrained large language models (LLMs), such as BERT (Devlin et al., 2019; Liu et al., 2019b; Sanh et al., 2019), OPT (Zhang et al., 2022b), LLaMA (Touvron et al., 2023a;b), and GPT (Radford et al., 2018; Brown et al., 2020; Ouyang et al., 2022; OpenAI, 2023), achieves state-of-the-art performance in a wide array of downstream applications. However, two significant challenges persist in practical adoption: memory demands for gradient-based optimizers and the need to safeguard the privacy of domain-specific fine-tuning data.

---

<sup>1</sup> Department of Computer Science, ETH Zurich <sup>2</sup> Amazon Search <sup>3</sup> Paul G. Allen School of Computer Science and Engineering, University of Washington. Correspondence to: Liang Zhang <liang.zhang@inf.ethz.ch>.

As the memory requirement of fine-tuning LLMs is increasingly becoming a bottleneck, various approaches have been proposed, spanning from parameter-efficient fine-tuning (PEFT) (Li & Liang, 2021; Hu et al., 2022) to novel optimization algorithms (Shazeer & Stern, 2018; Anil et al., 2019). Since these methods rely on backpropagation to compute the gradients, which can be memory-intensive, a recent trend has emerged in developing algorithms that do not require backpropagation (Baydin et al., 2022; Silver et al., 2022; Hinton, 2022; Hou et al., 2023; Phang et al., 2023; Chen et al., 2024a). Specifically for LLMs, Malladi et al. (2023) introduced zeroth-order methods for fine-tuning, thereby eliminating the backward pass and freeing up the memory for gradients and activations. Utilizing a single A100 GPU (80 GiB memory), zeroth-order methods are capable of fine-tuning a 30-billion-parameter model, whereas first-order methods, even equipped with PEFT, fail to fit into the memory for a model with more than 6.7 billion parameters. This greatly expands the potential for deploying and fine-tuning LLMs even on personal devices.

On the other hand, empirical studies have highlighted the risk of LLMs inadvertently revealing sensitive information from their fine-tuning datasets (Mireshghallah et al., 2022; Zeng et al., 2023; Mattern et al., 2023; Lukas et al., 2023). Such privacy concerns are pronounced especially when users opt to fine-tune LLMs on datasets of their own. Notably, the expectation that machine learning models should not compromise the confidentiality of their contributing entities is codified into legal frameworks (Voigt & Von dem Bussche, 2017). Differential privacy (DP) (Dwork et al., 2006) is a widely accepted mathematical framework for ensuring privacy by preventing attackers from identifying participating entities (Shokri et al., 2017). Consequently, the development of methods that fine-tune LLMs under differential privacy is of pressing necessity (Li et al., 2022b; Yu et al., 2022; He et al., 2023; Bu et al., 2023b; Du et al., 2023); however, most efforts so far have focused on first-order algorithms.

Motivated by the memory-hungry nature and privacy concerns in fine-tuning LLMs, we investigate zeroth-order methods that guarantee differential privacy for solving the fol-

lowing stochastic optimization problem:

$$\min_{x \in \mathbb{R}^d} F_S(x) := \frac{1}{n} \sum_{i=1}^n f(x; \xi_i), \quad (1)$$

where  $S = \{\xi_i\}_{i=1}^n$  is the training data,  $x \in \mathbb{R}^d$  is the model weight, the loss  $f(x; \xi_i)$  is Lipschitz for each sample  $\xi_i$ , and the averaged loss  $F_S(x)$  is smooth and possibly nonconvex. In theory, previous work on both differentially private optimization (Bassily et al., 2014) and zeroth-order optimization (Duchi et al., 2015) indicated that their convergence guarantees depend explicitly on the dimension  $d$ . Such dimension dependence becomes problematic in the context of LLMs with  $d$  scaling to billions. In practice, and somewhat surprisingly, empirical studies on the fine-tuning of LLMs using zeroth-order methods (Malladi et al., 2023) and DP first-order methods (Yu et al., 2022; Li et al., 2022b;a) have shown that the performance degradation due to the large model size is marginal. For example, Yu et al. (2022) showed that the performance drop due to privacy is smaller for larger architectures. A 345 million-sized GPT-2-Medium, fine-tuned with  $(\epsilon = 6.8, \delta = 10^{-5})$ -DP, showcases a modest drop of 5.1 in BLEU score (Papineni et al., 2002) (compared to a non-private model of the same size and architecture), whereas a larger GPT-2-XL with 1.5 billion parameters exhibits smaller cost in test performance, i.e., 4.3 BLEU score under the same privacy budget.

This gap between theory and practice has been linked to the presence of low-rank structures in the fine-tuning of pretrained LLMs (Malladi et al., 2023; Li et al., 2022a). Empirical evidence suggests that fine-tuning occurs within a low-dimensional subspace (Sagun et al., 2017; Gur-Ari et al., 2018; Ghorbani et al., 2019; Li et al., 2018): 200 dimensions for RoBERTa with 355 million parameters (Aghajanyan et al., 2021) and 100 dimensions for PEFT on DistilRoBERTa with 7 million parameters (Li et al., 2022a). In such cases where the intrinsic dimension is small, zeroth-order methods are known to achieve dimension-independent convergence rate (Malladi et al., 2023) and private first-order methods are also known to achieve dimension-independent guarantees (Ma et al., 2022; Li et al., 2022a).

Given the significance of fine-tuning LLMs on domain-specific datasets, we ask the following fundamental question: *Can we achieve a dimension-independent rate both under differential privacy and with access only to the zeroth-order oracle?* Our contributions are summarized below.

- We first show that the straightforward approach — that combines DP first-order methods with zeroth-order gradient estimators (Algorithm 1) — exhibits an undesirable dimension dependence in the convergence guarantees, even when the effective rank of the problem does not scale with the dimension (Theorems 1 and 2 in Section 3). There are two

root causes. First, the standard practice of choosing the clipping threshold to be the maximum norm of the estimated sample gradient leads to an unnecessarily large threshold. Next, this choice of the clipping threshold forces the addition of a large noise to ensure privacy, and Algorithm 1 adds that noise in all  $d$  directions.

- We present DPZERO (Algorithm 2), the first nearly dimension-independent DP zeroth-order method for stochastic optimization. Its convergence guarantee depends on the effective rank of the problem (specified in Assumption 3.5) and exhibits logarithmic dependence on the dimension  $d$  (Theorem 3 in Section 4). This builds upon two insights. First, the direction of the estimated gradient is a public information and does not need to be private; it is sufficient to make only the magnitude of the estimated gradient private, which is a scalar value. Next, we introduce a tighter analysis that allows us to choose a significantly smaller clipping threshold, leveraging the fact that the typical norm of the estimated gradient is much smaller than its maximum.

- We verify the effectiveness of DPZERO in both synthetic examples and private fine-tuning tasks on RoBERTa (Liu et al., 2019b) and OPT (Zhang et al., 2022b). In contrast to first-order algorithms that demand extensive effort for the efficient implementation of per-sample gradient clipping (Li et al., 2022b; He et al., 2023; Bu et al., 2023b), DPZERO offers the advantage of near-zero additional costs compared to non-private zeroth-order methods (Malladi et al., 2023). Our empirical results validate theoretical findings, revealing only a slight performance decrement for DPZERO even with large model sizes.

## 1.1. Related Works

We build upon exciting advances in zeroth-order optimization and differentially private optimization, which we survey here. Notably, DPZERO is inspired by new empirical and theoretical findings showing that fine-tuning LLMs does not suffer in high-dimensions when using zeroth-order methods in Malladi et al. (2023) or using private first-order optimization in Li et al. (2022a). Due to space limitation, a more comprehensive overview is deferred to Appendix A.

**Zeroth-order optimization.** Nesterov & Spokoiny (2017) pioneered the formal analysis of the convergence rate of zeroth-order methods, i.e., zeroth-order (stochastic) gradient descent (ZO-SGD) that replaces gradients in SGD by their zeroth-order estimators. Their findings are later refined by several works (Ghadimi & Lan, 2013; Shamir, 2017; Lin et al., 2022). These well-established results indicate a runtime complexity  $\mathcal{O}(d)$  worse than first-order methods. Such dimension dependence of zeroth-order methods is proven inevitable without additional structures (Wibisono et al., 2012; Duchi et al., 2015).

There are several recent works that relax the dimension dependence in zeroth-order methods leveraging problem structures. Balasubramanian & Ghadimi (2018) demonstrated that ZO-SGD can directly identify the sparsity of the problem and proved a dimension-independent rate when the support of gradients remains unchanged. Yue et al. (2023) and Malladi et al. (2023) relaxed the dependence on dimension  $d$  to a quantity related to the trace of the loss’s Hessian.

**Differentially private optimization.** Previous works on DP optimization mostly center around first-order methods. When the problem is nonconvex, i.e., the setting of our interest, differentially private (stochastic) gradient descent (DP-GD) achieves a rate of  $\mathcal{O}(\sqrt{d \log(1/\delta)}/(n\varepsilon))$  on the squared norm of the gradient (Wang et al., 2017; Zhou et al., 2020). We show that DPZERO matches this rate with access only to the zeroth-order oracle in Theorem 3. Given access to the first-order oracle, it has been recently shown that such rate can be improved to  $\mathcal{O}((\sqrt{d \log(1/\delta)}/(n\varepsilon))^{4/3})$  leveraging momentum (Tran & Cutkosky, 2022) or variance reduction techniques (Arora et al., 2023).

Early works established dimension-independent rates when the gradients lie in some fixed low-rank subspace (Jain & Thakurta, 2014; Song et al., 2021). Closest to our result is Song et al. (2021), which demonstrated that the rate of DP-GD for smooth nonconvex optimization can be improved to  $\mathcal{O}(\sqrt{r \log(1/\delta)}/(n\varepsilon))$  for generalized linear models (GLMs) with a rank- $r$  feature matrix. DPZERO matches this result with access only to the zeroth-order oracle in Theorem 3 for more general problems beyond low-rank GLMs. Our result is inspired by Li et al. (2022a) that introduced a relaxed Lipschitz condition for the gradients and provided dimension-free bounds when the loss is convex and the relaxed Lipschitz parameters decay rapidly. Similarly, Ma et al. (2022) suggested that the dependence on  $d$  in the utility upper bound for DP stochastic convex optimization can be improved.

Literature on DP optimization beyond first-order methods remains less explored. Recently, Zhang et al. (2024a) studied the problem of private zeroth-order nonsmooth nonconvex optimization and achieved a rate that depends on the dimension  $d$ . As far as we are aware, no prior studies have addressed the challenge of deriving a dimension-independent rate in DP zeroth-order optimization.

After the workshop version of our paper (Zhang et al., 2023) was released, Tang et al. (2024a) concurrently discovered the same algorithm as DPZERO (up to a minor difference in how  $u_t$  is drawn) and showed empirical benefits when applied to fine-tuning OPT models but without theoretical analysis. Also building upon the workshop version of our paper, Liu et al. (2024) introduced DP-ZOSO, a stage-wise zeroth-order method with an additional quadratic regularizer. With

extra hyper-parameters to be tuned, DP-ZOSO demonstrates further empirical gain over DPZERO. However, Liu et al. (2024) only provided *dimension-dependent* guarantees.

## 2. Preliminaries

**Notation.** We use  $\|\cdot\|$  for the Euclidean norm and define  $\|v\|_W^2 = v^\top W v$  for a square matrix  $W$ .  $\mathbb{S}^{d-1} = \{x \in \mathbb{R}^d \mid \|x\| = 1\}$  denotes the unit sphere in  $\mathbb{R}^d$ , and  $\eta \mathbb{S}^{d-1}$  is the sphere of radius  $\eta > 0$ . A function  $p : \mathbb{R}^d \rightarrow \mathbb{R}$  is  $L$ -Lipschitz if  $|p(x_1) - p(x_2)| \leq L\|x_1 - x_2\|, \forall x_1, x_2$ . A function  $q : \mathbb{R}^d \rightarrow \mathbb{R}$  is  $\ell$ -smooth if it is differentiable and  $\|\nabla q(x_1) - \nabla q(x_2)\| \leq \ell\|x_1 - x_2\|$ . The trace of a square matrix  $J$  is denoted by  $\text{Tr}(J)$ . A symmetric real matrix  $M \succeq 0$  if it is positive semi-definite. The clipping operation is defined to be  $\text{clip}_C(x) = x \min\{1, C/\|x\|\}$  given  $C > 0$ . The notation  $\tilde{\mathcal{O}}(\cdot)$  hides additional logarithmic terms.

### 2.1. Differential Privacy

**Definition 2.1** (Differential Privacy (Dwork et al., 2006; 2014)). Two datasets  $S = \{\xi_i\}_{i=1}^n$  and  $S' = \{\xi'_i\}_{i=1}^n$  are *neighboring* if  $\max\{|S \setminus S'|, |S' \setminus S|\} = 1$ , and we denote it by  $S \sim S'$ . For prescribed  $\varepsilon > 0$  and  $\delta \in (0, 1)$ , an algorithm  $\mathcal{A}$  is said to satisfy  $(\varepsilon, \delta)$ -*differential privacy* (DP) if  $\mathbb{P}(\mathcal{A}(S) \in \mathcal{B}) \leq e^\varepsilon \mathbb{P}(\mathcal{A}(S') \in \mathcal{B}) + \delta$  for all  $S \sim S'$  and all measurable set  $\mathcal{B}$  in the range of  $\mathcal{A}$ .

To ensure DP while solving the optimization problem in Eq. (1), first-order approaches, such as DP-GD, update via  $x_{t+1} \leftarrow x_t - \alpha((1/n) \sum_{i=1}^n \text{clip}_C(\nabla f(x_t; \xi_i)) + z_t)$ ; see e.g., (Song et al., 2013; Abadi et al., 2016). Through the following composition lemma (Kairouz et al., 2015, Theorem 4.3), the privacy for entire  $T$  updates is secured by the per-sample clipping operation that ensures finite sensitivity of  $\Delta = 2C/n$  together with the Gaussian noise  $z_t$ .

**Lemma 2.2** (Advanced Composition). *Let  $\mathcal{A}$  be some randomized algorithm operating on a dataset  $S$  and outputting a vector in  $\mathbb{R}^d$ . If  $\mathcal{A}$  has sensitivity  $\Delta := \sup_{S \sim S'} \|\mathcal{A}(S) - \mathcal{A}(S')\|$ , the mechanism that adds Gaussian noise  $\mathcal{N}(0, \sigma^2 \mathbf{I}_d)$  with variance  $\sigma^2 = (2\Delta \sqrt{2T \log(e + (\varepsilon/\delta))})/\varepsilon)^2$  satisfies  $(\varepsilon, \delta)$ -DP under  $T$ -fold adaptive composition for any  $\varepsilon > 0$  and  $\delta \in (0, 1)$ .*

### 2.2. Zeroth-Order Optimization

When the gradient is expensive to compute, zeroth-order methods are useful for optimizing Eq. (1). For example, the two-point gradient estimator below requires only two evaluation of function values (Shamir, 2017)

$$g_\lambda(x; \xi_i) := \frac{f(x + \lambda u; \xi_i) - f(x - \lambda u; \xi_i)}{2\lambda} u, \quad (2)$$

**Algorithm 1** DP-GD with 0th-order gradients (DPGD-0th)

**Input:** Dataset  $S = \{\xi_1, \dots, \xi_n\}$ , initialization  $x_0 \in \mathbb{R}^d$ , number of iterations  $T$ , stepsize  $\alpha > 0$ , smoothing parameter  $\lambda > 0$ , clipping threshold  $C > 0$ , privacy parameters  $\varepsilon > 0, \delta \in (0, 1)$ .

- 1: **for**  $t = 0, 1, \dots, T - 1$  **do**
- 2: Sample  $u_t$  uniformly at random from the Euclidean sphere  $\sqrt{d}\mathbb{S}^{d-1}$  and for all  $i = 1, \dots, n$  compute

$$g_\lambda(x_t; \xi_i) \leftarrow \frac{f(x_t + \lambda u_t; \xi_i) - f(x_t - \lambda u_t; \xi_i)}{2\lambda} u_t.$$

- 3: Sample  $z_t \in \mathbb{R}^d$  randomly from the multivariate Gaussian distribution  $\mathcal{N}(0, \sigma^2 \mathbf{I}_d)$  with variance  $\sigma = 4C\sqrt{2T \log(e + (\varepsilon/\delta))}/(n\varepsilon)$  and update

$$x_{t+1} \leftarrow x_t - \alpha \left( \frac{1}{n} \sum_{i=1}^n \text{clip}_C(g_\lambda(x_t; \xi_i)) + z_t \right).$$

- 4: **end for**

**Output:**  $x_\tau$  for  $\tau$  sampled uniformly at random from  $\{0, 1, \dots, T - 1\}$ .

where  $u$  is sampled uniformly from the Euclidean sphere  $\sqrt{d}\mathbb{S}^{d-1}$  and  $\lambda > 0$  is the smoothing parameter (Yousefian et al., 2012; Duchi et al., 2012). A common approach to generate  $u$  is to set  $u = \sqrt{d}z/\|z\|$ , with  $z$  sampled from the standard multivariate Gaussian  $\mathcal{N}(0, \mathbf{I}_d)$  (Muller, 1959; Marsaglia, 1972). We refer to  $g_\lambda(x; \xi)$  as the *zeroth-order gradient* (estimator) in the sequel. The results in this paper can be directly extended to other zeroth-order gradient estimators, e.g., any  $u$  satisfying  $\mathbb{E}[uu^\top] = \mathbf{I}_d$  (Duchi et al., 2015), the one-point estimator (Flaxman et al., 2005), and the directional derivative (Nesterov & Spokoiny, 2017).

### 3. DP-GD with Zeroth-Order Gradients Suffers in High Dimensions

In this section, we show that the direct integration of zeroth-order gradient estimators in Eq. (2) into DP-GD, which we term DPGD-0th, leads to undesirable dimension dependence in the error rate. Such dependence persists even under a low effective rank assumption.

#### 3.1. Direct Integration Leads to an $\mathcal{O}(d^{3/2})$ Rate

We present in Algorithm 1 the straightforward private zeroth-order approach that substitutes the gradients in DP-GD with zeroth-order estimators  $g_\lambda(x_t; \xi_i)$  in Eq. (2).

The privacy guarantee follows from standard DP-GD analysis, and the utility guarantee on the squared gradient norm is derived from classical techniques for analyzing zeroth-order

methods (Nesterov & Spokoiny, 2017). Before presenting the convergence result, we make the following standard assumption, which is common in nonconvex DP optimization (Wang et al., 2017; 2019; Tran & Cutkosky, 2022).

**Assumption 3.1.** The loss  $f(x; \xi)$  is  $L$ -Lipschitz for every  $\xi$ . The average loss  $F_S(x)$  is  $\ell$ -smooth for every given dataset  $S$ , and its minimum  $F_S^* := \min_{x \in \mathbb{R}^d} F_S(x)$  is finite.

**Theorem 1.** For any  $\varepsilon > 0$  and  $\delta \in (0, 1)$ , Algorithm 1 is  $(\varepsilon, \delta)$ -DP. Under Assumption 3.1, its output  $x_\tau$  satisfies that

$$\mathbb{E}[\|\nabla F_S(x_\tau)\|^2] \leq 16 \left( (F_S(x_0) - F_S^*) \ell + 2L^2 \right) \frac{d\sqrt{d \log(e + (\varepsilon/\delta))}}{n\varepsilon}, \quad (3)$$

with the choice of parameters

$$\alpha = \frac{1}{4\ell d}, \quad T = \frac{n\varepsilon}{\sqrt{d \log(e + (\varepsilon/\delta))}},$$

$$\lambda \leq \frac{4L}{\ell d} \left( \frac{\sqrt{d \log(e + (\varepsilon/\delta))}}{n\varepsilon} \right)^{1/2}, \quad C = Ld.$$

The total number of zeroth-order gradient computations is  $nT = \mathcal{O}(n^2/\sqrt{d})$ .

*Remark 3.2.* Theorem 1 demonstrates that directly combining DP-GD with zeroth-order gradients leads to an  $\mathcal{O}(d^{3/2})$  error complexity, which is  $\mathcal{O}(d)$  worse than first-order DP approaches (Wang et al., 2017).

*Remark 3.3.* Three sources contribute to the dependence in  $d$ : the squared norm of the zeroth-order gradient estimator  $\mathbb{E}[\|(1/n) \sum_{i=1}^n g_\lambda(x, \xi_i)\|^2] = \mathcal{O}(d \|\nabla F_S(x)\|^2)$  when taking  $\lambda \rightarrow 0$  for simplicity, the clipping threshold  $C = \mathcal{O}(d)$ , and the norm of the privacy noise  $\mathbb{E}[\|z_t\|^2] = \mathcal{O}(dC^2) = \mathcal{O}(d^3)$ . The standard analysis of one-step update gives

$$\mathbb{E}[F_S(x_{t+1})] \leq \mathbb{E}[F_S(x_t)] - \frac{\alpha}{2} (1 - 2d\ell\alpha) \mathbb{E}[\|\nabla F_S(x_t)\|^2] + c\alpha^2 d^3, \quad (4)$$

where  $c$  is a constant that depends on problem parameters other than  $\alpha$  and  $d$ ; see Eq. (12) for details. A small enough step size,  $\alpha < 1/(2\ell d)$ , is required to make the second term negative, where the dependence in  $d$  comes from  $\mathbb{E}[\|(1/n) \sum_{i=1}^n g_\lambda(x, \xi_i)\|^2]$ . The dependence on  $d^3$  in the last term arises from  $\mathbb{E}[\|z_t\|^2]$ , which leads to the  $\mathcal{O}(d^{3/2})$  rate in Eq. (3) after balancing error terms. Detailed proofs can be found in Appendix D.

*Remark 3.4.* The choice of the clipping threshold  $C = Ld$  ensures that clipping does not happen with probability one, which is a common choice in the theoretical analysis of private optimization algorithms (Bassily et al., 2014; 2019; Wang et al., 2017). This follows from the fact that,



for  $L$ -Lipschitz  $f(x; \xi)$ , the zeroth-order gradient is upper bounded by  $\|g_\lambda(x; \xi)\| \leq Ld$  almost surely. Selecting the clipping threshold without knowledge of this upper bound remains an active research topic (Chen et al., 2020; Yang et al., 2022; Fang et al., 2023; Koloskova et al., 2023; Zhang et al., 2024b).

### 3.2. Rate Improves to $\mathcal{O}(d)$ under Low Effective Rank

Here, under the low-dimensional structures in fine-tuning LLMs (cf. Section 1), we demonstrate improved performance for Algorithm 1. Unfortunately, a linear dependence in  $d$  still persists even under the low effective rank structure.

**Assumption 3.5.** The function  $f(x; \xi)$  is  $L$ -Lipschitz and  $\ell$ -smooth for every  $\xi$ . The average function  $F_S(x)$  is twice differentiable with  $-H \preceq \nabla^2 F_S(x) \preceq H$  for any  $x \in \mathbb{R}^d$ , and its minimum  $F_S^* := \min_{x \in \mathbb{R}^d} F_S(x)$  is finite. Here, the real-valued  $d \times d$  matrix  $H \succeq 0$  satisfies that  $\|H\|_2 \leq \ell$  and  $\text{Tr}(H) \leq r\|H\|_2$ . We refer to  $r$  as the effective rank or the intrinsic dimension of the problem.

Assumption 3.5 boils down to Assumption 3.1 if  $r = d$ . This is because  $-H' \preceq \nabla^2 F_S(x) \preceq H', \forall x \in \mathbb{R}^d$  and  $H' = \ell \mathbf{I}_d$  imply that  $\|H'\|_2 \leq \ell$  and  $\text{Tr}(H') \leq d\|H'\|_2$ . With  $r < d$ , this assumption reflects the additional structures encoded in the Hessian matrix. While Assumption 3.5 naturally holds for low-rank Hessians, it covers more general cases. For example, the assumption is satisfied with  $r = \mathcal{O}(\log d) \ll d$  in the case of a full-rank matrix  $H$ , with its  $i$ -th largest eigenvalue being  $\ell/i$  for  $1 \leq i \leq d$ .

Similar assumptions have been made to relax the dimension dependence in zeroth-order optimization in the limit  $\lambda \rightarrow 0$  (Malladi et al., 2023) and also for DP first-order optimization when the objective is smooth and convex (Ma et al., 2022). However, even under Assumption 3.5, DPGD-0th (Algorithm 1) still suffers from a linear dependence in  $d$  in its error rate, as presented below. A proof is provided in Appendix D.

**Theorem 2.** For any  $\varepsilon > 0$  and  $\delta \in (0, 1)$ , Algorithm 1 is  $(\varepsilon, \delta)$ -DP. Under Assumption 3.5, its output  $x_\tau$  satisfies that

$$\mathbb{E}[\|\nabla F_S(x_\tau)\|^2] \leq 16 \left( (F_S(x_0) - F_S^*) \ell + 2L^2 \right) \frac{d\sqrt{r \log(e + (\varepsilon/\delta))}}{n\varepsilon}, \quad (5)$$

with the choice of parameters

$$\alpha = \frac{1}{4\ell(r+2)}, \quad T = \frac{n(r+2)\varepsilon}{d\sqrt{r \log(e + (\varepsilon/\delta))}},$$

$$\lambda \leq \frac{4L}{\ell d} \left( \frac{\sqrt{r \log(e + (\varepsilon/\delta))}}{n\varepsilon} \right)^{1/2}, \quad C = Ld.$$

The total number of zeroth-order gradient computations is  $nT = \mathcal{O}(n^2\sqrt{r}/d)$ .

*Remark 3.6.* Comparing to Remark 3.3, both the zeroth-order gradient,  $\mathbb{E}[\|(1/n) \sum_{i=1}^n g_\lambda(x_t; \xi_i)\|_H^2]$ , and the DP noise,  $\mathbb{E}[\|z_t\|_H^2]$ , decrease by a factor of  $\mathcal{O}(r/d)$  under low effective rank. This is made precise in Lemma C.1. As a result, the one-step update analysis can be tightened as

$$\mathbb{E}[F_S(x_{t+1})] \leq \mathbb{E}[F_S(x_t)] - \frac{\alpha}{2} (1 - 2(r+2)\ell\alpha) \mathbb{E}[\|\nabla F_S(x_t)\|^2] + c\alpha^2 r d^2. \quad (6)$$

Comparing to the RHS of Eq. (4), it achieves an improved dependence in  $d$ . However, the third term in Eq. (6) is still at  $\mathcal{O}(d^2)$  due to the clipping threshold  $C = \mathcal{O}(d)$ . Consequently, even when the effective rank  $r$  is small, Eq. (5) still grows linearly in  $d$ .

## 4. DPZero: Nearly Dimension-Independent Private Zeroth-Order Optimization

A straightforward combination of DP-GD and zeroth-order methods has a large dimension dependence. Our novel DPZERO overcomes this issue with two key insights elaborated below.

**Scalar privacy noise.** By decoupling zeroth-order gradients in Eq. (2) into direction and magnitude, our key observation is that the direction,  $u_t$ , is public knowledge, and we only need to make the magnitude private. Privacy can be guaranteed by clipping the finite-difference,  $(f(x_t + \lambda u_t; \xi_i) - f(x_t - \lambda u_t; \xi_i))/(2\lambda)$ , and then adding a scalar noise  $z_t$ ; see line 3 of Algorithm 2. This change, when applied to Algorithm 1, can significantly improve the rate in Eq. (5) by a factor of  $d^{1/2}$ .

**Tighter clipping threshold.** Another factor of  $d^{1/2}$  improvement originates from a tighter analysis on the upper bound of the finite-difference term. Although its worst-case upper bound scales with the dimension  $d$ , this only happens with an exponentially small probability over the randomness of  $u_t$ . As proved in Eq. (16) in Appendix E, the size of the finite-difference is

$$\frac{|f(x_t + \lambda u_t; \xi_i) - f(x_t - \lambda u_t; \xi_i)|}{2\lambda} \leq |u_t^\top \nabla f(x_t; \xi_i)| + \frac{\ell}{2} \lambda d,$$

where we use the assumption that each  $f(x; \xi)$  is  $\ell$ -smooth. When  $u_t$  is sampled from the sphere  $\sqrt{d}\mathbb{S}^{d-1}$ , a tail bound (part (ii) of Lemma C.1 in the appendix) implies that

$$\mathbb{P}(|u_t^\top \nabla f(x_t; \xi_i)| \geq C) \leq 2\sqrt{2\pi} \exp\left(-\frac{C^2}{8L^2}\right).$$

By selecting the smoothing parameter  $\lambda$  to be sufficiently small, a careful choice of  $C = \tilde{\mathcal{O}}(L)$ , which is nearly independent of  $d$ , can ensure that clipping does not occur with a

**Algorithm 2** DPZERO

**Input:** Dataset  $S = \{\xi_1, \dots, \xi_n\}$ , initialization  $x_0 \in \mathbb{R}^d$ , number of iterations  $T$ , stepsize  $\alpha > 0$ , smoothing parameter  $\lambda > 0$ , clipping threshold  $C > 0$ , privacy parameters  $\varepsilon > 0, \delta \in (0, 1)$ .

- 1: **for**  $t = 0, 1, \dots, T - 1$  **do**
- 2:   Sample  $u_t$  uniformly at random from the Euclidean sphere  $\sqrt{d}S^{d-1}$ .
- 3:   Sample a scalar  $z_t \in \mathbb{R}$  randomly from the univariate Gaussian distribution  $\mathcal{N}(0, \sigma^2)$  with variance  $\sigma = 4C\sqrt{2T \log(e + (\varepsilon/\delta))}/(n\varepsilon)$  and update the parameter

$$x_{t+1} \leftarrow x_t - \alpha \left( \frac{1}{n} \sum_{i=1}^n \text{clip}_C \left( \frac{f(x_t + \lambda u_t; \xi_i) - f(x_t - \lambda u_t; \xi_i)}{2\lambda} \right) + z_t \right) u_t.$$

4: **end for**

**Output:**  $x_\tau$  for  $\tau$  sampled uniformly at random from  $\{0, 1, \dots, T - 1\}$ .

high probability. This choice is significantly smaller than the worst-case clipping threshold of  $Ld^{1/2}$ . The main technical challenge is that we need to analyze the algorithm given the event that clipping does not happen. The choice of drawing  $u_t$  from the uniform distribution over the sphere, together with corresponding tail bounds in Appendix C, allows us to prove the following nearly dimension-independent bound under the low effective rank structure in Assumption 3.5. A proof is provided in Appendix E.

**Theorem 3.** *For any  $\varepsilon > 0$  and  $\delta \in (0, 1)$ , Algorithm 2 is  $(\varepsilon, \delta)$ -DP. Under Assumption 3.5, suppose  $\max_{0 \leq t \leq T} |F_S(x_t)| \leq B$ , the output  $x_\tau$  satisfies that*

$$\mathbb{E}[\|\nabla F_S(x_\tau)\|^2] \leq \left( 64 \left( (F_S(x_0) - F_S^*) \ell + \tilde{L}^2 \right) + 2L^2 \right) \frac{\sqrt{r \log(e + (\varepsilon/\delta))}}{n\varepsilon}, \quad (7)$$

where we define

$$\tilde{L}^2 = L^2 \log \left( \frac{2\sqrt{2\pi} n^3 \varepsilon^2 (r+2)(d + 8\ell B(r+2)/L^2)}{r \log(e + (\varepsilon/\delta))} \right),$$

and choose the parameters to be

$$\alpha = \frac{1}{4\ell(r+2)}, \quad T = \frac{n(r+2)\varepsilon}{4\sqrt{r \log(e + (\varepsilon/\delta))}}, \quad C = 4\tilde{L},$$

$$\lambda \leq \frac{1}{\ell d} \min \left\{ 4(2 - \sqrt{2})\tilde{L}, \frac{L}{\sqrt{d}} \left( \frac{\sqrt{r \log(e + (\varepsilon/\delta))}}{n\varepsilon} \right)^{\frac{1}{2}} \right\}.$$

The total number of zeroth-order gradient computations is  $nT = \mathcal{O}(n^2\sqrt{r})$ .

**Remark 4.1.** Algorithm 2 is nearly dimension-independent, given its logarithmic dependence on  $d$ . To the best of our knowledge, this is the first zeroth-order DP method that is nearly dimension-independent. This feature is significantly beneficial for fine-tuning pretrained LLMs where the effective rank has been observed to be quite small (Aghajanyan et al., 2021; Li et al., 2022a). When  $r = d$ , our rate in

Table 1. The dependence of the error rate on dimension  $d$  and effective rank  $r$  shows that the proposed DPZERO (Alg. 2) significantly outperforms DPGD-0th (Alg. 1) and achieves performance close to the popular first-order method, DP-GD, on both scenarios with and without a low-effective rank assumption. Note that the error rates of zeroth and first-order DP methods are achieved with different number of iterations.

	w/o Asmp. 3.5	with Asmp. 3.5
DPGD-0th	$\mathcal{O}(d\sqrt{d})$	$\mathcal{O}(d\sqrt{r})$
DPZERO	$\mathcal{O}((\log d)\sqrt{d})$	$\mathcal{O}((\log d)\sqrt{r})$
DP-GD	$\mathcal{O}(\sqrt{d})$	$\mathcal{O}(\sqrt{r})$

Eq. (7) nearly matches that of the best known achievable bound of the first-order method DP-GD for smooth nonconvex losses (Wang et al., 2017). When the effective rank  $r$  is smaller, this algorithm achieves  $\tilde{\mathcal{O}}(\sqrt{r \log(1/\delta)})/(n\varepsilon)$  squared gradient norm. Similar dimension-free error rate is established for DP-GD on unconstrained generalized linear losses (Song et al., 2021), with a dependence on the rank of the feature matrix. Table 1 provides a summary on how DPZERO depends on dimension  $d$  and effective rank  $r$ .

**Remark 4.2.** The RHS of Eq. (7) improves upon Eq. (5) of Algorithm 1 by a factor of  $d$ . Simplifying our analysis in Eq. (22) and conditioned on the event that the clipping does not happen, we get a similar one-step update analysis as Eq. (6) (see Eq. (22) and (23) for a precise inequality). However, since the privacy noise  $z_t$  is a scalar and the clipping threshold has been reduced, we have that  $\mathbb{E}[\|z_t u_t\|_H^2] = \tilde{\mathcal{O}}(r)$  is nearly independent of the dimension  $d$ , and thus the final error scales as  $\tilde{\mathcal{O}}(r^{1/2})$ .

**Remark 4.3.** The strategy of appropriately selecting the clipping threshold to ensure that clipping occurs with low probability is commonly applied in the analysis of private algorithms (Fang et al., 2023; Shen et al., 2023). Adaptive

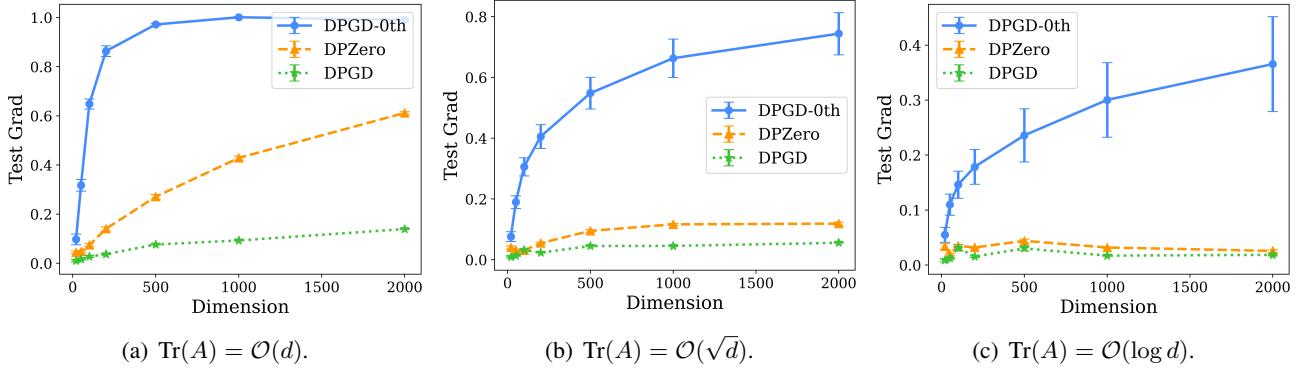


Figure 1. Experiments on the quadratic loss with effective rank  $\text{Tr}(A)$  (Assumption 3.5). For three different modes of the effective rank, we demonstrate how the norm of the test gradient depends on the problem dimension. DPGD-0th (Algorithm 1) has a strong dimension dependence regardless of the effective rank, while DPZERO (Algorithm 2) achieves dimension-independent performance when effective rank is small (right panel), similar to the standard first-order method DP-GD.

choices of clipping thresholds can provably improve error rates for certain problems including PCA (Liu et al., 2022) and linear regression (Liu et al., 2023c). One technical challenge in the choice of the clipping threshold in DPZERO is that we need the *expected* one-step progress to be sufficient in Eq. (22). This requires controlling the progress in the low-probability event that finite difference is clipped. The fact that  $\|u_t\|$  is finite with probability one simplifies the analysis, which is the reason we choose to sample  $u_t$  uniformly at random over the sphere. We believe that the analysis extends to the commonly used spherical Gaussian random vectors, which we leave as a future research direction. Table 7 in the appendix supports our hypothesis that the resulting performances are similar whether Gaussian or spherical random vectors are used. We choose Gaussian vectors for our experiments in Section 5 for simplicity.

**Remark 4.4.** Our theoretical results, including Theorems 1, 2, and 3, can be extended to the setting where the average loss  $F_S(x)$  additionally satisfies the PL inequality (Karimi et al., 2016; Polyak, 1963; Łojasiewicz, 1963). Under Assumption 3.5, DPZERO converges to an optimal solution in a nearly dimension-independent error rate. See more details in Appendix F.

**Remark 4.5.** Per-sample clipping is essential in DP algorithms to ensure bounded sensitivity that determines the magnitude of the DP noise. Besides the dimension-free error rates and memory saving of no backpropagation, another practical merit of DPZERO stems from the significantly simplified clipping compared with DP-GD. In addition to the advantage of clipping a *scalar* function value difference rather than a gradient vector as required by first-order methods, the efficiency of DPZERO is mainly attributed to the low-cost *per-sample* operations. In DP first-order methods, clipping is applied to gradients for every sample in a batch. The straightforward method of performing backward steps

for each sample to compute its gradient loses the benefit of parallelization, leading to significant memory and runtime overhead. Despite extensive effort in improving the efficiency of per-sample gradient clipping (Li et al., 2022b; He et al., 2023; Bu et al., 2023b), these methods still incur extra costs compared to non-DP algorithms. However, the clipping in DPZERO only involves computing the per-sample loss from forward steps and incurs no overhead in memory and runtime. This is straightforward for implementation as it is directly supported by, e.g., PyTorch, and no additional techniques are required. DPZERO is thus the first private method for fine-tuning LLMs that achieves near-zero additional costs compared to non-DP baselines, which is highly preferable especially in resource-constrained scenarios.

## 5. Experiments

We provide empirical results on synthetic problems and private fine-tuning of language models for sentence classification and generation tasks. A thorough description of the experimental settings is available in Appendix B. All experiments are tested on a single NVIDIA GeForce RTX 3090 GPU with 24 GiB memory. Code is available at <https://github.com/Liang137/DPZero>.

**Synthetic example.** Our first evaluation compares the performance of Algorithm 2 (DPZERO) with Algorithm 1 (DPGD-0th) and DP-GD on problems with different effective ranks. In particular, we use a quadratic loss

$$\min_{x \in \mathbb{R}^d} F_S(x) = \frac{1}{2n} \sum_{i=1}^n (x - x_i)^\top A (x - x_i),$$

with three choices of the Hessian matrix,  $A$ , whose effective ranks are designed to be  $\mathcal{O}(d)$ ,  $\mathcal{O}(\sqrt{d})$ , and  $\mathcal{O}(\log d)$ , respectively. All methods are trained with  $(\varepsilon = 2, \delta = 10^{-6})$ -DP on a training set  $\{x_1, \dots, x_n\}$  with  $n = 10,000$  and

Table 2. Experiments on RoBERTa (355M). We report both mean and standard error of the accuracy (%) across three random seeds. Zero-shot results with no fine-tuning provide lower bounds (taken from Malladi et al. (2023)), since they can be achieved with no private data. MeZO is not private and serves as an upper bound of DPZERO. LoRA (Hu et al., 2022) and DP-LoRA adopt AdamW (Loshchilov & Hutter, 2018) as their optimizer. All first-order methods (AdamW, LoRA, and their private versions) utilize the implementation by Li et al. (2022b). Thanks to DPZERO, the performance gaps between zeroth and first-order methods are made smaller in private fine-tuning.

Task	SST-2	SST-5	SNLI	MNLI	RTE	TREC
	— Sentiment —		— Natural Language Inference —			— Topic —
AdamW	93.1 ± 0.3	56.6 ± 0.3	86.4 ± 0.8	81.4 ± 0.9	83.6 ± 1.6	95.9 ± 0.2
DP-AdamW ( $\epsilon = 6$ )	91.6 ± 1.2	49.0 ± 0.3	81.5 ± 1.4	76.3 ± 0.9	77.3 ± 1.1	89.9 ± 0.8
DP-AdamW ( $\epsilon = 2$ )	90.5 ± 1.5	47.5 ± 0.5	74.6 ± 1.0	70.3 ± 0.8	72.8 ± 0.9	85.0 ± 0.5
LoRA	93.3 ± 0.4	55.3 ± 1.0	85.9 ± 0.7	82.2 ± 0.7	84.2 ± 0.4	94.6 ± 0.4
DP-LoRA ( $\epsilon = 6$ )	91.0 ± 1.3	48.8 ± 0.5	81.0 ± 1.5	72.8 ± 1.8	74.7 ± 1.3	89.2 ± 0.8
DP-LoRA ( $\epsilon = 2$ )	90.2 ± 1.2	47.1 ± 0.4	74.7 ± 1.6	65.7 ± 0.9	69.2 ± 1.1	83.2 ± 2.3
MeZO	92.5 ± 0.3	50.8 ± 0.8	80.4 ± 0.6	69.2 ± 0.3	72.8 ± 1.0	88.9 ± 0.1
DPZERO ( $\epsilon = 6$ )	92.2 ± 0.3	49.3 ± 0.6	77.8 ± 1.0	67.4 ± 0.3	71.9 ± 0.9	87.6 ± 0.9
DPZERO ( $\epsilon = 2$ )	91.8 ± 0.1	47.1 ± 0.9	73.6 ± 0.9	62.7 ± 0.9	70.4 ± 0.7	82.0 ± 1.6
Zero-Shot	79.0	35.5	50.2	48.8	51.4	32.0

Table 3. Runtime per iteration (s) and memory consumption (MiB) when fine-tuning RoBERTa (355M) for SST-2. Private methods in the table ensure ( $\epsilon = 2, \delta = 10^{-5}$ )-DP. DPZERO is as memory and runtime efficient as the non-private zeroth-order method MeZO (Malladi et al., 2023). First-order methods DP-AdamW and DP-LoRA (AdamW as the optimizer) both introduce considerable memory and runtime overhead compared to their non-private baselines. All first-order methods use the implementation by Li et al. (2022b). Comparisons with other implementations of DP first-order methods can be found in Table 9 in the appendix.

Method	Time (s/iter)	Memory (MiB)
AdamW	1.25	15820
DP-AdamW	2.12	17126
LoRA	0.821	10366
DP-LoRA	1.05	10496
MeZO	0.345	2668
DPZERO	0.347	2668

evaluated on a test set of the same size. The problem dimension is increased from 20 to 2,000. We perform a parameter search and plot the best gradient norm evaluated on the test set in Figure 1. Every method scales with the dimension  $d$  when the effective rank is  $d$  (as in Figure 1(a)), and DPGD-0th has the worst performance. When the effective rank reduces to  $\log d$  (as in Figure 1(c)), both DP-GD and DPZERO become nearly dimension-independent, which validates the dimension independence of DPZERO. Appendix B.1 includes more results measuring the loss and the gradient norm for both training and test datasets.

**Fine-tuning on RoBERTa.** Next, we follow the experimental setting in Malladi et al. (2023) and evaluate DPZERO on fine-tuning RoBERTa (Liu et al., 2019b) with 355M parameters across six different sentence classification tasks. We consider the few-shot scenario with 512 samples per class. We report the test accuracy for DPZERO trained with ( $\epsilon = \{2, 6\}, \delta = 10^{-5}$ )-DP and non-private zeroth-order baseline MeZO (Malladi et al., 2023) and compare them with first-order methods in Table 2. The memory consumption and per-iteration runtime are shown in Table 3. DP first-order methods introduce additional overhead in both memory and runtime compared to non-DP baselines, with a maximum accuracy drop of 9.5% when  $\epsilon = 6$ . However, DPZERO enjoys the same benefit as MeZO on memory efficiency and achieves near-zero additional costs, with at max only a 2.6% drop in the accuracy. In our experiments, we notice that the clipping threshold of DPZERO is typically larger compared to DP first-order methods; see Figure 4 in the appendix. This is consistent with the results in Theorem 3 regarding the selection of the clipping threshold  $C$ .

Compared with DP first-order methods, the main benefit of DPZERO is memory efficiency. Such memory savings are even greater than those observed in non-DP domains, thanks to DPZERO’s efficient clipping (cf. Remark 4.5). We note that the aim of Table 3 is to explain that DP first-order methods need considerable memory and runtime overhead compared to non-DP methods, while DPZERO does not. Such comparisons happen between DP and non-DP algorithms, respectively. We do not intend to directly compare the runtime of DPZERO to DP first-order methods as it depends on the implementation. In general, zeroth-order methods require more iterations to attain the same level of



Table 4. Experiments on OPT for classification tasks. We report mean and standard error of the accuracy (%) across three random seeds.

Model Task	OPT-1.3B		OPT-2.7B		OPT-6.7B	
	SST-2	BoolQ	SST-2	BoolQ	SST-2	BoolQ
MeZO	88.2 ± 0.9	63.2 ± 0.8	91.9 ± 0.5	65.3 ± 1.3	93.0 ± 0.2	67.4 ± 2.3
DPZERO ( $\varepsilon = 6$ )	88.2 ± 1.1	62.4 ± 0.8	91.5 ± 1.7	65.4 ± 1.6	92.6 ± 0.7	66.8 ± 1.6
DPZERO ( $\varepsilon = 2$ )	86.8 ± 1.7	61.6 ± 1.1	90.5 ± 0.9	63.7 ± 0.7	90.6 ± 1.3	63.7 ± 0.7
Zero-Shot	53.6	45.3	56.3	47.7	61.2	59.4

Table 5. Experiments on OPT for generation tasks. We report both mean and standard error of the f1 score (%) across three random seeds.

Model Task	OPT-1.3B		OPT-2.7B		OPT-6.7B	
	SQuAD	DROP	SQuAD	DROP	SQuAD	DROP
MeZO	73.5 ± 1.2	24.4 ± 0.2	76.3 ± 0.8	25.5 ± 1.2	79.7 ± 1.1	28.8 ± 0.7
DPZERO ( $\varepsilon = 6$ )	72.6 ± 0.8	24.7 ± 1.0	75.7 ± 1.5	24.6 ± 0.5	79.5 ± 0.9	28.4 ± 1.3
DPZERO ( $\varepsilon = 2$ )	70.1 ± 1.6	23.9 ± 1.2	71.9 ± 1.2	23.1 ± 0.9	77.1 ± 1.0	27.6 ± 0.7
Zero-Shot	26.8	11.1	29.8	9.7	36.5	17.8

performance as first-order methods (Malladi et al., 2023). In our case, DP first-order methods take 1,000 iterations while DPZERO need 10,000 iterations. This aligns with Theorem 3, which states that DPZERO requires  $\mathcal{O}(r)$  times more iterations than DP-GD to attain the same level of error rate, where  $r$  is the effective rank. However, DPZERO can still be efficient for large models in terms of GPU hours, because first-order methods often require communication-heavy distributed training over more GPUs each with limited memory; see Appendix F.6 of Malladi et al. (2023).

**Fine-tuning on OPT.** We also provide experiments on fine-tuning OPT (Zhang et al., 2022b) in the few-shot setting to illustrate the scalability of DPZERO. On our device (a GPU with 24 GiB memory), the largest model that can fit in for zeroth-order methods is OPT-6.7B, while first-order methods already run out of memory for OPT-1.3B; see Table 11 in the appendix for a detailed comparison of the memory consumption. The results of DPZERO’s test performance on four downstream tasks are reported in Tables 4 and 5. DPZERO demonstrates the same level of scalability as MeZO, with the ability to fine-tune models wherever MeZO is applicable, and experiences only small drops in performance due to privacy (up to 0.9% when  $\varepsilon = 6$ ). Our results indicate the effectiveness of DPZERO for privately fine-tuning pretrained LLMs and confirm that it does not suffer in high dimensions.

## 6. Conclusion

DPZERO is proposed to privately fine-tune language models

in a memory efficient manner by avoiding backpropagation. Theoretically, DPZERO enjoys a provably near dimension-free rate under low-rank structures, clearing the barriers for scaling private fine-tuning of LLMs. When deploying DPZERO, the elimination of gradient computation not only significantly saves memory, but avoids the overhead in gradient clipping as well. Thus the benefit of using zeroth-order method is more significant for private optimization. The theoretical guarantees on scalability and the practical merits of DPZERO are validated on private fine-tuning of RoBERTa and OPT on several downstream tasks.

DPZERO uses the full batch gradient every iteration, and the analysis guarantees an upper bound on the empirical average gradient assuming smooth nonconvex objectives. We defer extensions to the stochastic mini-batch setting, guarantees on the population loss leveraging the stability of zeroth-order methods (Nikolakakis et al., 2022), and considerations of other assumptions on objective functions like convexity or nonsmoothness to future research. We believe this work opens up a plethora of other prospective directions in DP zeroth-order optimization. These include, but are not limited to, understanding advantages of the intrinsic noise in zeroth-order gradient estimators, discovering other structural assumptions like the restricted Lipschitz condition (Li et al., 2022a) for dimension-independent rates, exploring alternative private mechanisms for the privacy guarantees of DPZERO (e.g., the Laplace mechanism for pure DP (Tang et al., 2024a)), and utilizing momentum (Tran & Cutkosky, 2022) or variance reduction (Arora et al., 2023) techniques for an improved rate and computational complexity.

## Acknowledgements

We are grateful to Gavin Brown and Divyansh Pareek for their insightful discussions regarding the proofs. We also thank Fanny Yang for proofreading of the paper. Additionally, we thank all anonymous reviewers for their valuable suggestions. L.Z. gratefully acknowledges funding by the Max Planck ETH Center for Learning Systems (CLS). This work does not relate to the current position of K.T. at Amazon. N.H. is supported by ETH research grant funded through ETH Zurich Foundations and Swiss National Science Foundation Project Funding No. 200021-207343. S.O. is supported in part by the National Science Foundation under grant no. 2019844, 2112471, and 2229876 supported in part by funds provided by the National Science Foundation, by the Department of Homeland Security, and by IBM. Any opinions, findings, and conclusions or recommendations expressed in this material are those of the author(s) and do not necessarily reflect the views of the National Science Foundation or its federal agency and industry partners.

## Impact Statement

A major concern with current use-cases of large language models is privacy of the fine-tuning data. Fine-tuning on in-domain data greatly improves performance and is now a default option. However, in-domain data can contain sensitive information about the participants of the dataset. The proposed solution makes privacy protection easier, consuming less resources, thus democratizing the use of privacy enhancing technology beyond those who have access to large amounts of resources.

## References

- Abadi, M., Chu, A., Goodfellow, I., McMahan, H. B., Mironov, I., Talwar, K., and Zhang, L. Deep learning with differential privacy. In *Proceedings of the ACM SIGSAC Conference on Computer and Communications Security*, pp. 308–318, 2016.
- Aghajanyan, A., Gupta, S., and Zettlemoyer, L. Intrinsic dimensionality explains the effectiveness of language model fine-tuning. In *Proceedings of the Annual Meeting of the Association for Computational Linguistics and the International Joint Conference on Natural Language Processing*, pp. 7319–7328, 2021.
- Anil, R., Gupta, V., Koren, T., and Singer, Y. Memory efficient adaptive optimization. *Advances in Neural Information Processing Systems*, 32, 2019.
- Arora, R., Bassily, R., González, T., Guzmán, C. A., Menart, M., and Ullah, E. Faster rates of convergence to stationary points in differentially private optimization. In *International Conference on Machine Learning*, pp. 1060–1092. PMLR, 2023.
- Asi, H., Feldman, V., Koren, T., and Talwar, K. Private stochastic convex optimization: Optimal rates in  $\ell_1$  geometry. In *International Conference on Machine Learning*, pp. 393–403. PMLR, 2021.
- Auer, P., Cesa-Bianchi, N., and Fischer, P. Finite-time analysis of the multiarmed bandit problem. *Machine learning*, 47:235–256, 2002.
- Balasubramanian, K. and Ghadimi, S. Zeroth-order (non)-convex stochastic optimization via conditional gradient and gradient updates. *Advances in Neural Information Processing Systems*, 31, 2018.
- Bassily, R., Smith, A., and Thakurta, A. Private empirical risk minimization: Efficient algorithms and tight error bounds. In *IEEE Annual Symposium on Foundations of Computer Science*, pp. 464–473. IEEE, 2014.
- Bassily, R., Feldman, V., Talwar, K., and Guha Thakurta, A. Private stochastic convex optimization with optimal rates. *Advances in Neural Information Processing Systems*, 32, 2019.
- Bassily, R., Feldman, V., Guzmán, C., and Talwar, K. Stability of stochastic gradient descent on nonsmooth convex losses. *Advances in Neural Information Processing Systems*, 33, 2020.
- Baydin, A. G., Pearlmutter, B. A., Syme, D., Wood, F., and Torr, P. Gradients without backpropagation. *arXiv preprint arXiv:2202.08587*, 2022.
- Bentivogli, L., Clark, P., Dagan, I., and Giampiccolo, D. The fifth PASCAL recognizing textual entailment challenge, 2009.
- Bowman, S. R., Angeli, G., Potts, C., and Manning, C. D. A large annotated corpus for learning natural language inference. In *Proceedings of the Conference on Empirical Methods in Natural Language Processing*, pp. 632–642, 2015.
- Brown, T., Mann, B., Ryder, N., Subbiah, M., Kaplan, J. D., Dhariwal, P., Neelakantan, A., Shyam, P., Sastry, G., Askell, A., et al. Language models are few-shot learners. *Advances in Neural Information Processing Systems*, 33: 1877–1901, 2020.
- Bu, Z., Chiu, J., Liu, R., Zha, S., and Karypis, G. Zero redundancy distributed learning with differential privacy. *arXiv preprint arXiv:2311.11822*, 2023a.
- Bu, Z., Wang, Y.-X., Zha, S., and Karypis, G. Differentially private optimization on large model at small cost. In

- International Conference on Machine Learning*, pp. 3192–3218. PMLR, 2023b.
- Cai, H., Mckenzie, D., Yin, W., and Zhang, Z. Zeroth-order regularized optimization (ZORO): Approximately sparse gradients and adaptive sampling. *SIAM Journal on Optimization*, 32(2):687–714, 2022.
- Chaudhuri, K., Monteleoni, C., and Sarwate, A. D. Differentially private empirical risk minimization. *Journal of Machine Learning Research*, 12(3), 2011.
- Chen, A., Zhang, Y., Jia, J., Diffenderfer, J., Parasyris, K., Liu, J., Zhang, Y., Zhang, Z., Kailkhura, B., and Liu, S. DeepZero: Scaling up zeroth-order optimization for deep model training. In *International Conference on Learning Representations*, 2024a.
- Chen, P.-Y., Zhang, H., Sharma, Y., Yi, J., and Hsieh, C.-J. ZOO: Zeroth order optimization based black-box attacks to deep neural networks without training substitute models. In *Proceedings of the ACM Workshop on Artificial Intelligence and Security*, pp. 15–26, 2017.
- Chen, T., Da, L., Zhou, H., Li, P., Zhou, K., Chen, T., and Wei, H. Privacy-preserving fine-tuning of large language models through flatness. *arXiv preprint arXiv:2403.04124*, 2024b.
- Chen, X., Liu, S., Xu, K., Li, X., Lin, X., Hong, M., and Cox, D. ZO-AdaMM: Zeroth-order adaptive momentum method for black-box optimization. *Advances in Neural Information Processing Systems*, 32, 2019.
- Chen, X., Wu, S. Z., and Hong, M. Understanding gradient clipping in private SGD: A geometric perspective. *Advances in Neural Information Processing Systems*, 33: 13773–13782, 2020.
- Choromanski, K., Rowland, M., Sindhwani, V., Turner, R., and Weller, A. Structured evolution with compact architectures for scalable policy optimization. In *International Conference on Machine Learning*, pp. 970–978. PMLR, 2018.
- Clark, C., Lee, K., Chang, M.-W., Kwiatkowski, T., Collins, M., and Toutanova, K. BoolQ: Exploring the surprising difficulty of natural yes/no questions. In *Proceedings of the Conference of the North American Chapter of the Association for Computational Linguistics*, pp. 2924–2936, 2019.
- Cramér, H. *Mathematical methods of statistics*, volume 43. Princeton University Press, 1999.
- Dagan, I., Glickman, O., and Magnini, B. The PASCAL recognising textual entailment challenge, 2005.
- Devlin, J., Chang, M.-W., Lee, K., and Toutanova, K. BERT: Pre-training of deep bidirectional Transformers for language understanding. In *Proceedings of the Conference of the North American Chapter of the Association for Computational Linguistics*, pp. 4171–4186, 2019.
- Du, M., Yue, X., Chow, S. S., Wang, T., Huang, C., and Sun, H. DP-Forward: Fine-tuning and inference on language models with differential privacy in forward pass. In *Proceedings of the ACM SIGSAC Conference on Computer and Communications Security*, pp. 2665–2679, 2023.
- Dua, D., Wang, Y., Dasigi, P., Stanovsky, G., Singh, S., and Gardner, M. DROP: A reading comprehension benchmark requiring discrete reasoning over paragraphs. In *Proceedings of the Conference of the North American Chapter of the Association for Computational Linguistics*, pp. 2368–2378, 2019.
- Duchi, J. C., Bartlett, P. L., and Wainwright, M. J. Randomized smoothing for stochastic optimization. *SIAM Journal on Optimization*, 22(2):674–701, 2012.
- Duchi, J. C., Jordan, M. I., Wainwright, M. J., and Wibisono, A. Optimal rates for zero-order convex optimization: The power of two function evaluations. *IEEE Transactions on Information Theory*, 61(5):2788–2806, 2015.
- Dwork, C., McSherry, F., Nissim, K., and Smith, A. Calibrating noise to sensitivity in private data analysis. In *Theory of Cryptography Conference*, pp. 265–284. Springer, 2006.
- Dwork, C., Roth, A., et al. The algorithmic foundations of differential privacy. *Foundations and Trends® in Theoretical Computer Science*, 9(3–4):211–407, 2014.
- Fang, C., Li, C. J., Lin, Z., and Zhang, T. SPIDER: Near-optimal non-convex optimization via stochastic path-integrated differential estimator. *Advances in Neural Information Processing Systems*, 31, 2018.
- Fang, H., Li, X., Fan, C., and Li, P. Improved convergence of differential private SGD with gradient clipping. In *International Conference on Learning Representations*, 2023.
- Fang, W., Yu, Z., Jiang, Y., Shi, Y., Jones, C. N., and Zhou, Y. Communication-efficient stochastic zeroth-order optimization for federated learning. *IEEE Transactions on Signal Processing*, 70:5058–5073, 2022.
- Feldman, V., Koren, T., and Talwar, K. Private stochastic convex optimization: optimal rates in linear time. In *Proceedings of the 52nd Annual ACM SIGACT Symposium on Theory of Computing*, pp. 439–449, 2020.

- Flaxman, A. D., Kalai, A. T., and McMahan, H. B. Online convex optimization in the bandit setting: Gradient descent without a gradient. In *Proceedings of the ACM-SIAM Symposium on Discrete Algorithms*, pp. 385–394, 2005.
- Ganesh, A., Haghifam, M., Nasr, M., Oh, S., Steinke, T., Thakkar, O., Thakurta, A. G., and Wang, L. Why is public pretraining necessary for private model training? In *International Conference on Machine Learning*, pp. 10611–10627. PMLR, 2023a.
- Ganesh, A., Haghifam, M., Steinke, T., and Guha Thakurta, A. Faster differentially private convex optimization via second-order methods. *Advances in Neural Information Processing Systems*, 36, 2023b.
- Gao, T., Fisch, A., and Chen, D. Making pre-trained language models better few-shot learners. In *Proceedings of the Annual Meeting of the Association for Computational Linguistics and the International Joint Conference on Natural Language Processing*, pp. 3816–3830, 2021.
- Ghadimi, S. and Lan, G. Stochastic first-and zeroth-order methods for nonconvex stochastic programming. *SIAM Journal on Optimization*, 23(4):2341–2368, 2013.
- Ghorbani, B., Krishnan, S., and Xiao, Y. An investigation into neural net optimization via Hessian eigenvalue density. In *International Conference on Machine Learning*, pp. 2232–2241. PMLR, 2019.
- Giampiccolo, D., Magnini, B., Dagan, I., and Dolan, W. B. The third PASCAL recognizing textual entailment challenge, 2007.
- Golovin, D., Karro, J., Kochanski, G., Lee, C., Song, X., and Zhang, Q. Gradientless descent: High-dimensional zeroth-order optimization. In *International Conference on Learning Representations*, 2020.
- Gratton, C., Venkatesgowda, N. K., Arablouei, R., and Werner, S. Privacy-preserved distributed learning with zeroth-order optimization. *IEEE Transactions on Information Forensics and Security*, 17:265–279, 2021.
- Grill, J.-B., Valko, M., and Munos, R. Black-box optimization of noisy functions with unknown smoothness. *Advances in Neural Information Processing Systems*, 28, 2015.
- Guha Thakurta, A. and Smith, A. (Nearly) optimal algorithms for private online learning in full-information and bandit settings. *Advances in Neural Information Processing Systems*, 26, 2013.
- Gupta, A. K. and Nadarajah, S. *Handbook of Beta distribution and its applications*. CRC Press, 2004.
- Gur-Ari, G., Roberts, D. A., and Dyer, E. Gradient descent happens in a tiny subspace. *arXiv preprint arXiv:1812.04754*, 2018.
- Haim, R. B., Dagan, I., Dolan, B., Ferro, L., Giampiccolo, D., Magnini, B., and Szpektor, I. The second PASCAL recognising textual entailment challenge, 2006.
- Han, A., Mishra, B., Jawanpuria, P., and Gao, J. Differentially private Riemannian optimization. *Machine Learning*, 113(3):1133–1161, 2024.
- He, J., Li, X., Yu, D., Zhang, H., Kulkarni, J., Lee, Y. T., Backurs, A., Yu, N., and Bian, J. Exploring the limits of differentially private deep learning with group-wise clipping. In *International Conference on Learning Representations*, 2023.
- Hinton, G. The forward-forward algorithm: Some preliminary investigations. *arXiv preprint arXiv:2212.13345*, 2022.
- Hong, J., Wang, J. T., Zhang, C., LI, Z., Li, B., and Wang, Z. DP-OPT: Make large language model your privacy-preserving prompt engineer. In *International Conference on Learning Representations*, 2024.
- Hou, B., O’connor, J., Andreas, J., Chang, S., and Zhang, Y. PromptBoosting: Black-box text classification with ten forward passes. In *International Conference on Machine Learning*, pp. 13309–13324. PMLR, 2023.
- Hu, E. J., Shen, Y., Wallis, P., Allen-Zhu, Z., Li, Y., Wang, S., Wang, L., and Chen, W. LoRA: Low-rank adaptation of large language models. In *International Conference on Learning Representations*, 2022.
- Huang, Z., Hu, R., Guo, Y., Chan-Tin, E., and Gong, Y. DP-ADMM: ADMM-based distributed learning with differential privacy. *IEEE Transactions on Information Forensics and Security*, 15:1002–1012, 2019.
- Jain, P. and Thakurta, A. G. (Near) dimension independent risk bounds for differentially private learning. In *International Conference on Machine Learning*, pp. 476–484. PMLR, 2014.
- Ji, K., Wang, Z., Zhou, Y., and Liang, Y. Improved zeroth-order variance reduced algorithms and analysis for non-convex optimization. In *International Conference on Machine Learning*, pp. 3100–3109. PMLR, 2019.
- Kairouz, P., Oh, S., and Viswanath, P. The composition theorem for differential privacy. In *International Conference on Machine Learning*, pp. 1376–1385. PMLR, 2015.
- Kairouz, P., Diaz, M. R., Rush, K., and Thakurta, A. (Nearly) dimension independent private ERM with adapt rates via publicly estimated subspaces. In *Conference on Learning Theory*, pp. 2717–2746. PMLR, 2021.



- Karimi, H., Nutini, J., and Schmidt, M. Linear convergence of gradient and proximal-gradient methods under the Polyak-Łojasiewicz condition. In *European Conference on Machine Learning and Knowledge Discovery in Databases*, pp. 795–811, 2016.
- Kenthapadi, K., Korolova, A., Mironov, I., and Mishra, N. Privacy via the Johnson-Lindenstrauss transform. *Journal of Privacy and Confidentiality*, 5(1):39–71, 2013.
- Koloskova, A., Hendrikx, H., and Stich, S. U. Revisiting gradient clipping: Stochastic bias and tight convergence guarantees. In *International Conference on Machine Learning*, 2023.
- Kulkarni, J., Lee, Y. T., and Liu, D. Private non-smooth ERM and SCO in subquadratic steps. *Advances in Neural Information Processing Systems*, 34, 2021.
- Li, C., Farkhoor, H., Liu, R., and Yosinski, J. Measuring the intrinsic dimension of objective landscapes. In *International Conference on Learning Representations*, 2018.
- Li, X., Liu, D., Hashimoto, T. B., Inan, H. A., Kulkarni, J., Lee, Y.-T., and Guha Thakurta, A. When does differentially private learning not suffer in high dimensions? *Advances in Neural Information Processing Systems*, 35: 28616–28630, 2022a.
- Li, X., Tramer, F., Liang, P., and Hashimoto, T. Large language models can be strong differentially private learners. In *International Conference on Learning Representations*, 2022b.
- Li, X. L. and Liang, P. Prefix-tuning: Optimizing continuous prompts for generation. In *Proceedings of the Annual Meeting of the Association for Computational Linguistics and the International Joint Conference on Natural Language Processing*, pp. 4582–4597, 2021.
- Lian, X., Zhang, H., Hsieh, C.-J., Huang, Y., and Liu, J. A comprehensive linear speedup analysis for asynchronous stochastic parallel optimization from zeroth-order to first-order. *Advances in Neural Information Processing Systems*, 29, 2016.
- Lin, T., Zheng, Z., and Jordan, M. Gradient-free methods for deterministic and stochastic nonsmooth nonconvex optimization. *Advances in Neural Information Processing Systems*, 35:26160–26175, 2022.
- Liu, D., Ganesh, A., Oh, S., and Guha Thakurta, A. Private (stochastic) non-convex optimization revisited: Second-order stationary points and excess risks. *Advances in Neural Information Processing Systems*, 36, 2023a.
- Liu, S., Kailkhura, B., Chen, P.-Y., Ting, P., Chang, S., and Amini, L. Zeroth-order stochastic variance reduction for nonconvex optimization. *Advances in Neural Information Processing Systems*, 31, 2018.
- Liu, S., Chen, P.-Y., Chen, X., and Hong, M. SignSGD via zeroth-order oracle. In *International Conference on Learning Representations*, 2019a.
- Liu, T., Tang, J., Vietri, G., and Wu, S. Generating private synthetic data with genetic algorithms. In *International Conference on Machine Learning*, pp. 22009–22027. PMLR, 2023b.
- Liu, X., Kong, W., Jain, P., and Oh, S. DP-PCA: Statistically optimal and differentially private PCA. *Advances in Neural Information Processing Systems*, 35:29929–29943, 2022.
- Liu, X., Jain, P., Kong, W., Oh, S., and Suggala, A. Label robust and differentially private linear regression: Computational and statistical efficiency. *Advances in Neural Information Processing Systems*, 36, 2023c.
- Liu, Y., Ott, M., Goyal, N., Du, J., Joshi, M., Chen, D., Levy, O., Lewis, M., Zettlemoyer, L., and Stoyanov, V. RoBERTa: A robustly optimized BERT pretraining approach. *arXiv preprint arXiv:1907.11692*, 2019b.
- Liu, Z., Lou, J., Bao, W., Hu, Y., Li, B., Qin, Z., and Ren, K. Differentially private zeroth-order methods for scalable large language model finetuning. *arXiv preprint arXiv:2402.07818*, 2024.
- Łojasiewicz, S. A topological property of real analytic subsets. *Coll. du CNRS, Les équations aux dérivées partielles*, 117(87-89):2, 1963.
- Loshchilov, I. and Hutter, F. Decoupled weight decay regularization. In *International Conference on Learning Representations*, 2018.
- Lowy, A., Li, Z., Huang, T., and Razaviyayn, M. Optimal differentially private learning with public data. *arXiv preprint arXiv:2306.15056*, 2023.
- Lukas, N., Salem, A., Sim, R., Tople, S., Wutschitz, L., and Zanella-Béguelin, S. Analyzing leakage of personally identifiable information in language models. In *IEEE Symposium on Security and Privacy*, pp. 346–363. IEEE, 2023.
- Ma, Y.-A., Marinov, T. V., and Zhang, T. Dimension independent generalization of DP-SGD for overparameterized smooth convex optimization. *arXiv preprint arXiv:2206.01836*, 2022.

- Malladi, S., Gao, T., Nichani, E., Damian, A., Lee, J. D., Chen, D., and Arora, S. Fine-tuning language models with just forward passes. *Advances in Neural Information Processing Systems*, 36:53038–53075, 2023.
- Mania, H., Guy, A., and Recht, B. Simple random search of static linear policies is competitive for reinforcement learning. *Advances in Neural Information Processing Systems*, 31, 2018.
- Marsaglia, G. Choosing a point from the surface of a sphere. *The Annals of Mathematical Statistics*, 43(2):645–646, 1972.
- Mattern, J., Mireshghallah, F., Jin, Z., Schölkopf, B., Sachan, M., and Berg-Kirkpatrick, T. Membership inference attacks against language models via neighbourhood comparison. *arXiv preprint arXiv:2305.18462*, 2023.
- Mireshghallah, F., Uniyal, A., Wang, T., Evans, D., and Berg-Kirkpatrick, T. Memorization in NLP fine-tuning methods. *arXiv preprint arXiv:2205.12506*, 2022.
- Muller, M. E. A note on a method for generating points uniformly on  $n$ -dimensional spheres. *Communications of the ACM*, 2(4):19–20, 1959.
- Nesterov, Y. *Introductory lectures on convex optimization: A basic course*, volume 87. Springer Science & Business Media, 2003.
- Nesterov, Y. and Spokoiny, V. Random gradient-free minimization of convex functions. *Foundations of Computational Mathematics*, 17:527–566, 2017.
- Nikolakakis, K., Haddadpour, F., Kalogerias, D., and Karbasi, A. Black-box generalization: Stability of zeroth-order learning. *Advances in Neural Information Processing Systems*, 35:31525–31541, 2022.
- OpenAI. GPT-4 Technical Report, 2023.
- Ouyang, L., Wu, J., Jiang, X., Almeida, D., Wainwright, C., Mishkin, P., Zhang, C., Agarwal, S., Slama, K., Ray, A., et al. Training language models to follow instructions with human feedback. *Advances in Neural Information Processing Systems*, 35:27730–27744, 2022.
- Papineni, K., Roukos, S., Ward, T., and Zhu, W.-J. BLEU: A method for automatic evaluation of machine translation. In *Proceedings of the Annual Meeting of the Association for Computational Linguistics*, pp. 311–318, 2002.
- Phang, J., Mao, Y., He, P., and Chen, W. HyperTuning: Toward adapting large language models without backpropagation. In *International Conference on Machine Learning*, pp. 27854–27875. PMLR, 2023.
- Polyak, B. T. Gradient methods for the minimisation of functionals. *USSR Computational Mathematics and Mathematical Physics*, 3(4):864–878, 1963.
- Radford, A., Narasimhan, K., Salimans, T., Sutskever, I., et al. Improving language understanding by generative pre-training. *OpenAI*, 2018.
- Rajbhandari, S., Rasley, J., Ruwase, O., and He, Y. ZeRO: Memory optimizations toward training trillion parameter models. In *International Conference for High Performance Computing, Networking, Storage and Analysis*, pp. 1–16. IEEE, 2020.
- Rajpurkar, P., Zhang, J., Lopyrev, K., and Liang, P. SQuAD: 100,000+ questions for machine comprehension of text. In *Proceedings of the Conference on Empirical Methods in Natural Language Processing*, pp. 2383–2392, 2016.
- Reimherr, M., Bharath, K., and Soto, C. Differential privacy over Riemannian manifolds. *Advances in Neural Information Processing Systems*, 34:12292–12303, 2021.
- Sagun, L., Evci, U., Guney, V. U., Dauphin, Y., and Bottou, L. Empirical analysis of the Hessian of over-parametrized neural networks. *arXiv preprint arXiv:1706.04454*, 2017.
- Salimans, T., Ho, J., Chen, X., Sidor, S., and Sutskever, I. Evolution strategies as a scalable alternative to reinforcement learning. *arXiv preprint arXiv:1703.03864*, 2017.
- Sanh, V., Debut, L., Chaumond, J., and Wolf, T. DistilBERT, a distilled version of BERT: smaller, faster, cheaper and lighter. *arXiv preprint arXiv:1910.01108*, 2019.
- Shamir, O. An optimal algorithm for bandit and zero-order convex optimization with two-point feedback. *The Journal of Machine Learning Research*, 18(1):1703–1713, 2017.
- Shariff, R. and Sheffet, O. Differentially private contextual linear bandits. *Advances in Neural Information Processing Systems*, 31, 2018.
- Shazeer, N. and Stern, M. Adafactor: Adaptive learning rates with sublinear memory cost. In *International Conference on Machine Learning*, pp. 4596–4604. PMLR, 2018.
- Shen, Z., Ye, J., Kang, A., Hassani, H., and Shokri, R. Share your representation only: Guaranteed improvement of the privacy-utility tradeoff in federated learning. In *International Conference on Learning Representations*, 2023.
- Shokri, R., Stronati, M., Song, C., and Shmatikov, V. Membership inference attacks against machine learning models. In *IEEE Symposium on Security and Privacy*, pp. 3–18. IEEE, 2017.

- Silver, D., Goyal, A., Danihelka, I., Hessel, M., and van Hasselt, H. Learning by directional gradient descent. In *International Conference on Learning Representations*, 2022.
- Socher, R., Perelygin, A., Wu, J., Chuang, J., Manning, C. D., Ng, A. Y., and Potts, C. Recursive deep models for semantic compositionality over a sentiment treebank. In *Proceedings of the Conference on Empirical Methods in Natural Language Processing*, pp. 1631–1642, 2013.
- Song, S., Chaudhuri, K., and Sarwate, A. D. Stochastic gradient descent with differentially private updates. In *IEEE Global Conference on Signal and Information Processing*, pp. 245–248. IEEE, 2013.
- Song, S., Steinke, T., Thakkar, O., and Thakurta, A. Evading the curse of dimensionality in unconstrained private GLMs. In *International Conference on Artificial Intelligence and Statistics*, pp. 2638–2646. PMLR, 2021.
- Tang, X., Panda, A., Nasr, M., Mahloujifar, S., and Mittal, P. Private fine-tuning of large language models with zeroth-order optimization. *arXiv preprint arXiv:2401.04343*, 2024a.
- Tang, X., Shin, R., Inan, H. A., Manoel, A., Mireshghallah, F., Lin, Z., Gopi, S., Kulkarni, J., and Sim, R. Privacy-preserving in-context learning with differentially private few-shot generation. In *International Conference on Learning Representations*, 2024b.
- Tossou, A. and Dimitrakakis, C. Algorithms for differentially private multi-armed bandits. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 30, 2016.
- Touvron, H., Lavril, T., Izacard, G., Martinet, X., Lachaux, M.-A., Lacroix, T., Rozière, B., Goyal, N., Hambro, E., Azhar, F., et al. LLaMA: Open and efficient foundation language models. *arXiv preprint arXiv:2302.13971*, 2023a.
- Touvron, H., Martin, L., Stone, K., Albert, P., Almahairi, A., Babaei, Y., Bashlykov, N., Batra, S., Bhargava, P., Bhosale, S., et al. LLaMA 2: Open foundation and fine-tuned chat models. *arXiv preprint arXiv:2307.09288*, 2023b.
- Tran, H. and Cutkosky, A. Momentum aggregation for private non-convex ERM. *Advances in Neural Information Processing Systems*, 35:10996–11008, 2022.
- Utpala, S., Han, A., Jawanpuria, P., and Mishra, B. Improved differentially private Riemannian optimization: Fast sampling and variance reduction. *Transactions on Machine Learning Research*, 2023a. ISSN 2835-8856.
- Utpala, S., Vepakomma, P., and Miolane, N. Differentially private Fréchet mean on the manifold of symmetric positive definite (SPD) matrices with log-Euclidean metric. *Transactions on Machine Learning Research*, 2023b. ISSN 2835-8856.
- Vershynin, R. *High-dimensional probability: An introduction with applications in data science*, volume 47. Cambridge University Press, 2018.
- Voigt, P. and Von dem Bussche, A. The EU general data protection regulation (GDPR). *A Practical Guide, 1st Ed., Cham: Springer International Publishing*, 10(3152676): 10–5555, 2017.
- Voorhees, E. M. and Tice, D. M. Building a question answering test collection. In *Proceedings of the Annual International ACM SIGIR Conference on Research and Development in Information Retrieval*, pp. 200–207, 2000.
- Wainwright, M. J. *High-dimensional statistics: A non-asymptotic viewpoint*, volume 48. Cambridge University Press, 2019.
- Wang, A., Singh, A., Michael, J., Hill, F., Levy, O., and Bowman, S. R. GLUE: A multi-task benchmark and analysis platform for natural language understanding. In *International Conference on Learning Representations*, 2018a.
- Wang, D., Ye, M., and Xu, J. Differentially private empirical risk minimization revisited: Faster and more general. *Advances in Neural Information Processing Systems*, 30, 2017.
- Wang, D., Chen, C., and Xu, J. Differentially private empirical risk minimization with non-convex loss functions. In *International Conference on Machine Learning*, pp. 6526–6535. PMLR, 2019.
- Wang, Y., Du, S., Balakrishnan, S., and Singh, A. Stochastic zeroth-order optimization in high dimensions. In *International Conference on Artificial Intelligence and Statistics*, pp. 1356–1365. PMLR, 2018b.
- Wang, Z., Balasubramanian, K., Ma, S., and Razaviyayn, M. Zeroth-order algorithms for nonconvex–strongly-concave minimax problems with improved complexities. *Journal of Global Optimization*, pp. 1–32, 2022.
- Wibisono, A., Wainwright, M. J., Jordan, M., and Duchi, J. C. Finite sample convergence rates of zero-order stochastic optimization methods. *Advances in Neural Information Processing Systems*, 25, 2012.
- Williams, A., Nangia, N., and Bowman, S. A broad-coverage challenge corpus for sentence understanding through inference. In *Proceedings of the Conference*

- of the North American Chapter of the Association for Computational Linguistics, pp. 1112–1122, 2018.
- Wu, X., Li, F., Kumar, A., Chaudhuri, K., Jha, S., and Naughton, J. Bolt-on differential privacy for scalable stochastic gradient descent-based analytics. In *Proceedings of ACM International Conference on Management of Data*, pp. 1307–1322, 2017.
- Xu, M., Wu, Y., Cai, D., Li, X., and Wang, S. Federated fine-tuning of billion-sized language models across mobile devices. *arXiv preprint arXiv:2308.13894*, 2023.
- Yang, X., Zhang, H., Chen, W., and Liu, T.-Y. Normalized/Clipped SGD with perturbation for differentially private non-convex optimization. *arXiv preprint arXiv:2206.13033*, 2022.
- Yousefian, F., Nedić, A., and Shanbhag, U. V. On stochastic gradient and subgradient methods with adaptive steplength sequences. *Automatica*, 48(1):56–67, 2012.
- Yu, D., Naik, S., Backurs, A., Gopi, S., Inan, H. A., Kamath, G., Kulkarni, J., Lee, Y. T., Manoel, A., Wutschitz, L., Yekhanin, S., and Zhang, H. Differentially private fine-tuning of language models. In *International Conference on Learning Representations*, 2022.
- Yue, P., Yang, L., Fang, C., and Lin, Z. Zeroth-order optimization with weak dimension dependency. In *Annual Conference on Learning Theory*, pp. 4429–4472. PMLR, 2023.
- Zelikman, E., Huang, Q., Liang, P., Haber, N., and Goodman, N. D. Just one byte (per gradient): A note on low-bandwidth decentralized language model finetuning using shared randomness. *arXiv preprint arXiv:2306.10015*, 2023.
- Zeng, S., Li, Y., Ren, J., Liu, Y., Xu, H., He, P., Xing, Y., Wang, S., Tang, J., and Yin, D. Exploring memorization in fine-tuned language models. *arXiv preprint arXiv:2310.06714*, 2023.
- Zhang, J., Zheng, K., Mou, W., and Wang, L. Efficient private ERM for smooth objectives. In *Proceedings of the 26th International Joint Conference on Artificial Intelligence*, pp. 3922–3928, 2017.
- Zhang, L., Thekumparampil, K. K., Oh, S., and He, N. Bring your own algorithm for optimal differentially private stochastic minimax optimization. *Advances in Neural Information Processing Systems*, 35:35174–35187, 2022a.
- Zhang, L., Thekumparampil, K. K., Oh, S., and He, N. DPZero: Dimension-independent and differentially private zeroth-order optimization. *International Workshop on Federated Learning in the Age of Foundation Models in Conjunction with NeurIPS*, 2023.
- Zhang, Q., Tran, H., and Cutkosky, A. Private zeroth-order nonsmooth nonconvex optimization. In *International Conference on Learning Representations*, 2024a.
- Zhang, S., Roller, S., Goyal, N., Artetxe, M., Chen, M., Chen, S., Dewan, C., Diab, M., Li, X., Lin, X. V., et al. OPT: Open pre-trained Transformer language models. *arXiv preprint arXiv:2205.01068*, 2022b.
- Zhang, X., Bu, Z., Wu, S., and Hong, M. Differentially private SGD without clipping bias: An error-feedback approach. In *International Conference on Learning Representations*, 2024b.
- Zhou, Y., Chen, X., Hong, M., Wu, Z. S., and Banerjee, A. Private stochastic non-convex optimization: Adaptive algorithms and tighter generalization bounds. *arXiv preprint arXiv:2006.13501*, 2020.
- Zhou, Y., Wu, S., and Banerjee, A. Bypassing the ambient dimension: Private SGD with gradient subspace identification. In *International Conference on Learning Representations*, 2021.



## A. Additional Related Works

**Zeroth-order optimization.** Nesterov & Spokoiny (2017) pioneered the formal analysis of the convergence rate of zeroth-order methods, i.e., zeroth-order (stochastic) gradient descent (ZO-SGD) that replaces gradients in SGD by their zeroth-order estimators. This is motivated by renewed interest in adopting zeroth-order methods in industry due to, for example, fast differentiation techniques that require storing all intermediate computations reaching the memory limitations. Their findings on nonsmooth convex functions are later refined by Shamir (2017). Lin et al. (2022) contributed to further advancements on nonsmooth nonconvex functions recently. Additionally, Ghadimi & Lan (2013) extended the results for smooth functions into the stochastic setting. Zeroth-order methods have also been expanded to incorporate approaches such as coordinate descent (Lian et al., 2016), conditional gradient descent (Balasubramanian & Ghadimi, 2018), variance reduction techniques (Liu et al., 2018; Fang et al., 2018; Ji et al., 2019), SignSGD (Liu et al., 2019a), and minimax optimization (Wang et al., 2022). Additionally, zeroth-order methods find applications in fields such as black-box machine learning (Grill et al., 2015; Chen et al., 2017; 2019), bandit optimization (Flaxman et al., 2005; Shamir, 2017), reinforcement learning (Salimans et al., 2017; Choromanski et al., 2018; Mania et al., 2018), and distributed learning (Fang et al., 2022; Zelikman et al., 2023; Xu et al., 2023) to reduce communication overhead.

These well-established results indicate that the norm of the zeroth-order gradient scales with the dimension  $d$  and the required stepsize is  $d$ -times smaller than that in first-order gradient-based methods, leading to a  $d$ -times increase in the final time complexity. For example, the convergence rate of gradient descent for minimizing a smooth convex function  $f(x)$  is  $f(\bar{x}_T) - \min_{x \in \mathbb{R}^d} f(x) \leq \mathcal{O}(1/T)$  where  $\bar{x}_T$  is the average of  $T$  iterates (Nesterov, 2003), while the zeroth-order method only achieves a rate  $\mathcal{O}(d/T)$ . It has been shown that such dimension dependence of zeroth-order methods is inevitable without additional structures (Wibisono et al., 2012; Duchi et al., 2015).

There are several recent works that relax the dimension dependence in zeroth-order methods leveraging problem structures. Wang et al. (2018b) and Cai et al. (2022) assumed certain sparsity structure in the problem and applied sparse recovering algorithms, e.g. LASSO, to obtain sparse gradients from zeroth-order observations. Golovin et al. (2020) analyzed the case when the objective function is  $f(Px)$  for some low-rank projection matrix  $P$ . These works either require the objective or the algorithm to be modified to have a dimension-independent guarantee. Balasubramanian & Ghadimi (2018) demonstrated that ZO-SGD can directly identify the sparsity of the problem and proved a dimension-independent rate when the support of gradients remains unchanged (Cai et al., 2022). Recently, Yue et al. (2023) and Malladi et al. (2023) relaxed the dependence on dimension  $d$  to a quantity related to the trace of the loss’s Hessian.

**Differentially private optimization.** Previous works on DP optimization mostly center around first-order methods. For constrained convex problems, tight utility guarantees on both excess empirical (Chaudhuri et al., 2011; Bassily et al., 2014; Wu et al., 2017; Zhang et al., 2017; Wang et al., 2017) and population (Bassily et al., 2019; 2020; Feldman et al., 2020; Asi et al., 2021; Kulkarni et al., 2021; Zhang et al., 2022a) losses are well-understood. As an example, a typical result states that the optimal rate on the excess empirical loss for convex objectives is  $\Theta(\sqrt{d \log(1/\delta)}/(n\varepsilon))$ , where  $(\varepsilon, \delta)$  are privacy parameters,  $n$  is the number of samples, and  $d$  is the dimension. The dimension dependence is fundamental as both the upper bound (Bassily et al., 2014), using differentially private (stochastic) gradient descent (DP-GD) introduced in (Song et al., 2013), and the lower bound (Bassily et al., 2014), using a reduction to finger printing codes, have the same dependence.

When the problem is nonconvex, i.e., the setting of our interest, DP-GD achieves a rate of  $\mathcal{O}(\sqrt{d \log(1/\delta)}/(n\varepsilon))$  on the squared norm of the gradient (Wang et al., 2017; Zhou et al., 2020). We show that DPZERO matches this rate with access only to the zeroth-order oracle in Theorem 3. Given access to the first-order oracle, it has been recently shown that such rate can be improved to  $\mathcal{O}((\sqrt{d \log(1/\delta)}/(n\varepsilon))^{4/3})$  leveraging momentum (Tran & Cutkosky, 2022) or variance reduction techniques (Arora et al., 2023). Further, the convergence to second-order stationary points in nonconvex DP optimization is studied in (Liu et al., 2023a). Recent advancements in DP optimization have also delved into the understanding of the potential of public data (Ganesh et al., 2023a; Lowy et al., 2023), the convergence properties of per-sample gradient clipping (Yang et al., 2022; Fang et al., 2023; Koloskova et al., 2023; Zhang et al., 2024b), and the relaxation of the dimension dependence in the utility upper bound (Ma et al., 2022; Li et al., 2022a).

Early works established that dimension-independent rates can be attained when the gradients lie in some fixed low-rank subspace (Jain & Thakurta, 2014; Song et al., 2021). By first identifying this gradient subspace, dimension-independent algorithms can be designed (Zhou et al., 2021; Kairouz et al., 2021). Closest to our result is Song et al. (2021), which demonstrated that the rate of DP-GD for smooth nonconvex optimization can be improved to  $\mathcal{O}(\sqrt{r \log(1/\delta)}/(n\varepsilon))$  under certain structural assumptions, i.e., for generalized linear models (GLMs) with a rank- $r$  feature matrix. DPZERO matches

this result with access only to the zeroth-order oracle in Theorem 3 for more general problems beyond low-rank GLMs. Our result is inspired by Li et al. (2022a) that introduced a relaxed Lipschitz condition for the gradients and provided dimension-free bounds when the loss is convex and the relaxed Lipschitz parameters decay rapidly. Similarly, Ma et al. (2022) suggested that the dependence on  $d$  in the utility upper bound for DP stochastic convex optimization can be improved to a dependence on the trace of the Hessian. There is also a line of work on DP Riemannian optimization that achieves utility bounds dependent on the intrinsic dimension of the manifold (Reimherr et al., 2021; Utpala et al., 2023b;a; Han et al., 2024). Further exploration of its connection to the low-rank structure in this work is reserved for future.

Literature on DP optimization beyond first-order methods remains less explored. Ganesh et al. (2023b) investigated the potential of second-order methods for DP convex optimization. Gratton et al. (2021) proposed to use zeroth-order methods for DP-ADMM (Huang et al., 2019) in distributed learning. They state that the noise intrinsic in zeroth-order methods is enough to provide privacy guarantee and rely on the output of zeroth-order methods being Gaussian, which is unverified to the best of our knowledge. Liu et al. (2023b) proposed a private genetic algorithm based on zeroth-order optimization heuristics for private synthetic data generation. Recently, Zhang et al. (2024a) studied the problem of private zeroth-order nonsmooth nonconvex optimization and achieved a rate that depends on the dimension  $d$ . After the workshop version of our paper (Zhang et al., 2023) was released, Tang et al. (2024a) concurrently discovered the same algorithm as DPZERO (up to a minor difference in how  $u_t$  is drawn) and showed empirical benefits when applied to fine-tuning OPT models but without theoretical analysis. Also building upon the workshop version of our paper, Liu et al. (2024) introduced DP-ZOSO, a stage-wise zeroth-order method with an additional quadratic regularizer. With extra hyper-parameters to be tuned, DP-ZOSO demonstrates further empirical gain over DPZERO. However, Liu et al. (2024) only provided dimension-dependent guarantees. As far as we are aware, no prior studies have addressed the challenge of deriving a dimension-independent rate in DP zeroth-order optimization.

**Other relevant works.** Du et al. (2023) introduced a novel noise adding mechanism that happens in the forward pass of training. Although the algorithm is termed “DP-Forward”, it still requires backpropagation for training. In a separate context, Bu et al. (2023a) coincidentally proposed DP-ZERO, a term identical to ours, denoting a private version of the zero redundancy optimizer (ZeRO) by Rajbhandari et al. (2020) that aims at enhancing memory efficiency in data and model parallelisms. DP prompt tuning (Hong et al., 2024) and DP in-context learning (Tang et al., 2024b) provide resource-efficient alternatives compared to private fine-tuning, enabling the private adaptation of pretrained LLMs to specific tasks without extensive computational demands. Investigating how DPZERO performs relative to these methods and whether different techniques can be integrated is an interesting research problem. More recently, Chen et al. (2024b) proposed differentially private algorithms that enforce weight flatness to improve generalization, which can also handle zeroth-order oracles. There is also another line of research (Guha Thakurta & Smith, 2013; Tossou & Dimitrakakis, 2016; Shariff & Sheffet, 2018) on the design of differentially private algorithms for the stochastic bandit problem based on upper confidence bound (Auer et al., 2002). Their algorithms are not directly applicable to our setting.

## B. Additional Experiment Details

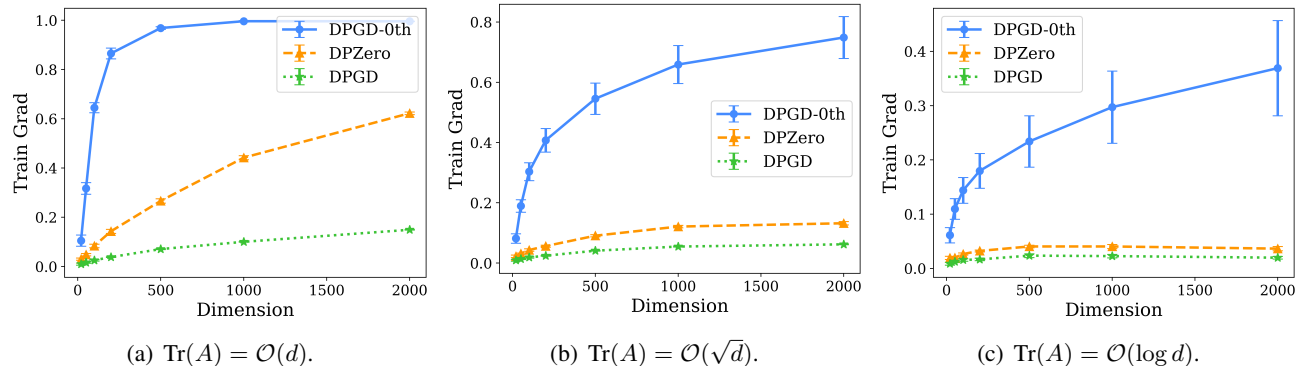


Figure 2. Experiments on the quadratic loss with effective rank  $\text{Tr}(A)$ . For three different modes of the effective rank, we increase the problem dimension and report the best gradient norm evaluated on the training set. Insights for the saturation of DPGD-0th when the dimension increases can be found in Remark F.5.

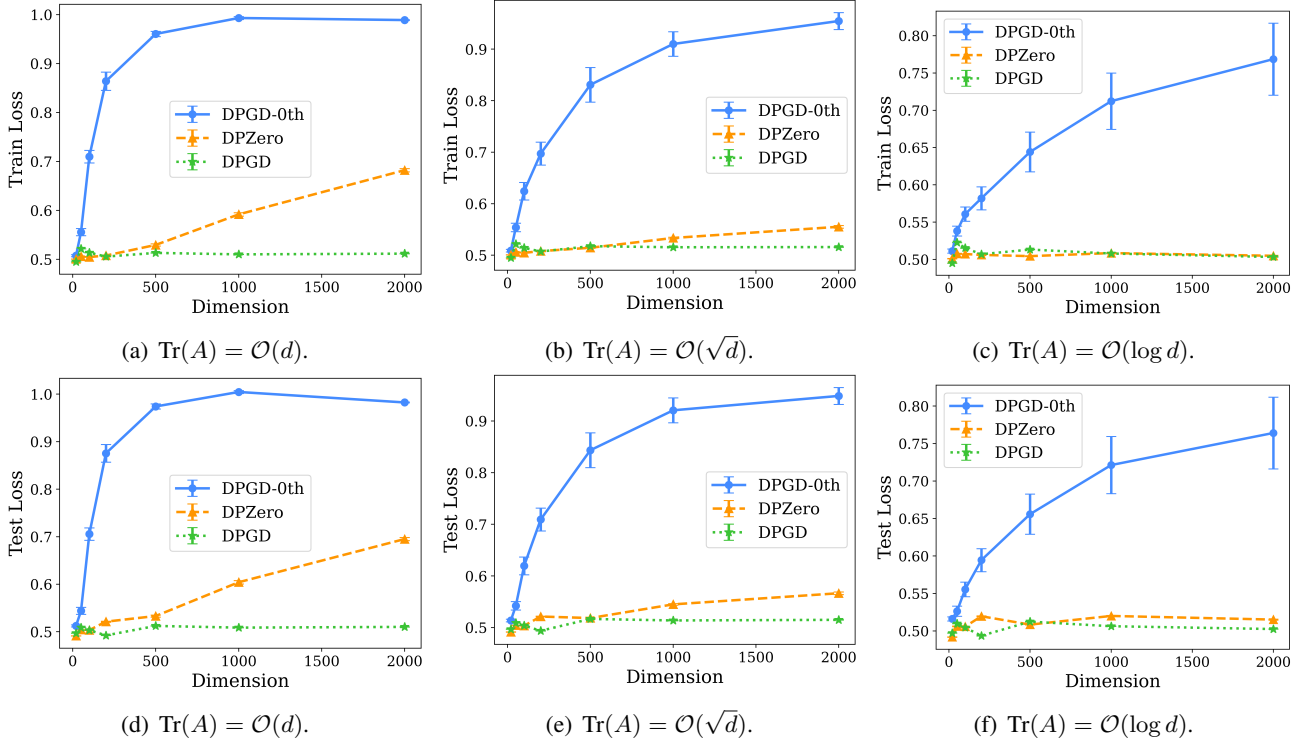


Figure 3. Experiments on the quadratic loss with effective rank  $\text{Tr}(A)$ . For three different modes of the effective rank, we increase the problem dimension and report the best loss evaluated on both training set ((a), (b), and (c)) and test set ((d), (e), and (f)).

In this section, we discuss our experimental setups in detail.

### B.1. Synthetic Example on a Quadratic Loss

Given a training dataset  $S = \{x_1, \dots, x_n\}$  with each coordinate of  $x_i \in \mathbb{R}^d$  sampled independently from the Gaussian  $\mathcal{N}(1, 1)$ , we implement DPZERO on the quadratic loss

$$\min_{x \in \mathbb{R}^d} F_S(x) = \frac{1}{2n} \sum_{i=1}^n (x - x_i)^\top A (x - x_i),$$

with a fixed Hessian  $A \in \mathbb{R}^{d \times d}$  that can be designed to implement different effective ranks  $r = \text{Tr}(A)/\|A\|_2$  according to Assumption 3.5. We compare DPZERO (Algorithm 2) with DPGD-0th (Algorithm 1) and first-order algorithm DP-GD on three patterns of the effective rank

- (a)  $\text{Tr}(A) = \mathcal{O}(d) : A = \text{diag}\{1, 1, \dots, 1\}$ ;
- (b)  $\text{Tr}(A) = \mathcal{O}(\sqrt{d}) : A = \text{diag}\{1, 1/\sqrt{2}, \dots, 1/\sqrt{d}\}$ ;
- (c)  $\text{Tr}(A) = \mathcal{O}(\log d) : A = \text{diag}\{1, 1/2, \dots, 1/d\}$ .

Since  $\|A\|_2 = 1$  in all cases, the effective rank  $r = \text{Tr}(A)$ . For each mode of the effective rank, we increase the problem dimension  $d$  from 20 to 2000. We perform a parameter search and plot the best gradient norm evaluated on the training set (see Figure 2) and a test set that follows the same distribution of the training set (see Figure 1). For completeness, we also plot both training and test loss in Figure 3. The key hyper-parameters used for the experiments are summarized in Table 6.

In all figures, we observe that the performance of each method is improved with smaller effective rank. For each pattern of the effective rank, DPGD-0th (Algorithm 1) has the worst performance, while DP-GD consistently achieves the best results. When the effective rank is  $d$ , every method scales with the dimension. When the effective rank improves to  $\log d$ ,

Table 6. Hyper-parameters used for the synthetic example on the quadratic loss. The number of iterations, stepsize, and clipping threshold are optimized through a grid search using given values. Other parameters are fixed to the listed values.

Hyper-parameters	Values
Number of training samples	10000
Number of test samples	10000
Dimension $d$	{20, 50, 100, 200, 500, 1000, 2000}
Privacy	( $\epsilon = 2, \delta = 10^{-6}$ )
Smoothing $\lambda$ (DPZERO and DPGD-0th)	$10^{-4}$
Number of iterations	{10, 20, 40, 80, 160, 320, 640, 1280, 2560, 5120}
Stepsize	{ $10^{-5}, 3 \times 10^{-5}, 10^{-4}, 3 \times 10^{-4}, 0.001, 0.003, 0.01, 0.03, 0.1, 0.3, 1$ }
Clipping	{0.1, 0.3, 1, 3, 10, 30, 100, 300}

DPZERO and DP-GD become nearly dimension-independent, and DPZERO matches the performance of the first-order method DP-GD. This validates our theoretical findings, as summarized in Table 1, and demonstrates the effectiveness of DPZERO. We want to mention that a similar set of experiments to verify the performance of DP-GD when dimension increases was also provided by Li et al. (2022a). Our implementation of this synthetic example is based on their code.

## B.2. Private Fine-Tuning of the Language Model RoBERTa

We follow experiment settings in Malladi et al. (2023) to evaluate the performance of DPZERO in the private fine-tuning of RoBERTa (Liu et al., 2019b) across six sentence classification datasets: SST-2 and SST-5 (Socher et al., 2013) for sentiment classification, SNLI (Bowman et al., 2015), MNLI (Williams et al., 2018), and RTE (Dagan et al., 2005; Haim et al., 2006; Giampiccolo et al., 2007; Bentivogli et al., 2009; Wang et al., 2018a) for natural language inference tasks, and TREC (Voorhees & Tice, 2000) for topic classification. In our experiments, we employ the same prompts as used in Malladi et al. (2023), which are adapted from Gao et al. (2021).

**Implementation details.** Our implementation of DPZERO utilizes the codebase provided by Malladi et al. (2023). For easier implementation and better memory efficiency, we follow Malladi et al. (2023) to sample the zeroth-order direction  $u_t$  from the Gaussian distribution  $\mathcal{N}(0, I_d)$  instead of the sphere as stated in Algorithm 2. Table 7 compares the performance of DPZERO on SST-2 and SST-5 when  $u_t$  is sampled from Gaussian and sphere. Given the negligible differences between the two sampling strategies, we continue with the Gaussian sampling for its simplicity. Another strategy in the implementation to further save memory involves storing only the random seed for the generation of the zeroth-order direction  $u_t$ , rather than the complete vector, and regenerating this direction whenever it’s used. Although DPZERO is stated for the full-batch case in Algorithm 2, we adopt a mini-batch setting in the experiments.

Table 7. Test accuracy (mean %  $\pm$  standard error %) of DPZERO when fine-tuning RoBERTa (355M) for SST-2 and SST-5 with ( $\epsilon = \{2, 6\}, \delta = 10^{-5}$ )-DP and using different sampling strategies of the zeroth-order update direction  $u_t$ . No notable difference is observed when  $u_t$  is sampled from either the Gaussian distribution or the Euclidean sphere.

Randomness	Gaussian		Sphere	
	$\epsilon = 6$	$\epsilon = 2$	$\epsilon = 6$	$\epsilon = 2$
SST-2	92.2 $\pm$ 0.3	91.8 $\pm$ 0.1	91.8 $\pm$ 0.1	91.5 $\pm$ 0.5
SST-5	49.3 $\pm$ 0.6	47.1 $\pm$ 0.9	49.9 $\pm$ 1.3	47.4 $\pm$ 1.3

**Hyper-parameter selection.** For all experiments, we employ a few-shot setting, utilizing 512 samples per class in the training set, randomly selected from the original dataset. The test set is also composed of 1000 randomly selected samples from the original test dataset. We fix the total number of iterations to be 10000, the batch size to be 64, and the smoothing parameter  $\lambda = 10^{-3}$  for both DPZERO and the non-private zeroth-order baseline MeZO (Malladi et al., 2023). Note that the original results of MeZO reported in Malladi et al. (2023) run for 100000 iterations. A parameter search of the learning



Table 8. Hyper-parameters used in DPZERO for fine-tuning RoBERTa (355M). We only optimize the clipping threshold through a grid search from 50 to 400. Other parameters are fixed to the listed values.

Hyper-parameters	Values
Number of training samples	512 per class
Number of test samples	1000
Number of iterations	10000
Batch size	64
Privacy	$(\epsilon = \{2, 6\}, \delta = 10^{-5})$
Smoothing $\lambda$	$10^{-3}$
Stepsize	$10^{-6}$
Clipping	{50, 100, 150, 200, 250, 300, 400}

rate for MeZO is performed, and it turns out  $10^{-6}$  consistently yields the best performance. We then fix the learning rate to be  $10^{-6}$  for DPZERO and only search for the clipping threshold for different tasks. There is potential for improved performance by well-optimizing other hyper-parameters, such as the learning rate and the number of iterations. All results are averaged through three different random seeds {42, 13, 21} for selecting the few-shot datasets. The hyper-parameters used for our language model fine-tuning experiments are summarized in Table 8.

**Comparison with first-order methods.** Regarding the first-order methods, we use the same few-shot setting as before, and the results are averaged over three different random seeds {42, 13, 21}. The number of iterations is set to be 1000, and the batch size is fixed to be 64. The learning rate is optimized by a grid search over  $\{5 \times 10^{-5}, 10^{-4}, 5 \times 10^{-4}, 10^{-3}\}$ , and the clipping threshold is optimized by a grid search over {0.1, 0.5, 1, 10}. In the experiments for LoRA, we set the rank to be 8 and the LoRA  $\alpha = 16$ , which remain the same as in the original paper (Hu et al., 2022). All other parameters are fixed to their default values. In addition to Li et al. (2022b) in Tables 2 and 3, we also compare the performance of DPZERO to two other implementations of DP first-order methods, Yu et al. (2022) and Bu et al. (2023b), in Table 9. DPZERO achieves similar performance on SST-2 as DP first-order methods, while saving a significant amount of memory. Such memory savings are greater than the savings of MeZO (Malladi et al., 2023) over AdamW (Loshchilov & Hutter, 2018) and LoRA (Hu et al., 2022) (AdamW as the optimizer), due to DPZero’s simpler clipping (cf. Remark 4.5).

Table 9. Test accuracy (%), runtime per iteration (s), and memory consumption (MiB) when fine-tuning RoBERTa (355M) for SST-2. Private methods in the table guarantee  $(\epsilon = 2, \delta = 10^{-5})$ -DP. A fair comparison is ensured among Li et al. (2022b) and Bu et al. (2023b), as they are implemented using the same codebase. It is important to note, however, that they cannot be directly compared with those of Yu et al. (2022), due to differences in implementations. LoRA (Hu et al., 2022) and DP-LoRA use the first-order method AdamW (Loshchilov & Hutter, 2018) as the optimizer. DP first-order methods introduce considerable overheads in both memory and runtime compared to their non-DP baselines, while DPZERO does not, thanks to its novel design of the efficient clipping. Also note that such comparisons between DP and non-DP algorithms are fair since they use the same codebase.

Method	Acc.	Time (s/iter)	Memory (MiB)
AdamW (Li et al., 2022b)	93.1	1.25	15820
DP-AdamW (Li et al., 2022b)	90.5	2.12	17126
DP-AdamW (Bu et al., 2023b)	91.1	1.55	18372
AdamW (Yu et al., 2022)	94.4	0.425	16960
DP-AdamW (Yu et al., 2022)	92.3	2.33	21494
LoRA (Li et al., 2022b)	93.3	0.821	10366
DP-LoRA (Li et al., 2022b)	90.2	1.05	10496
LoRA (Yu et al., 2022)	94.3	0.301	11512
DP-LoRA (Yu et al., 2022)	91.3	0.332	11522
MeZO	92.5	0.345	2668
DPZERO	91.8	0.347	2668

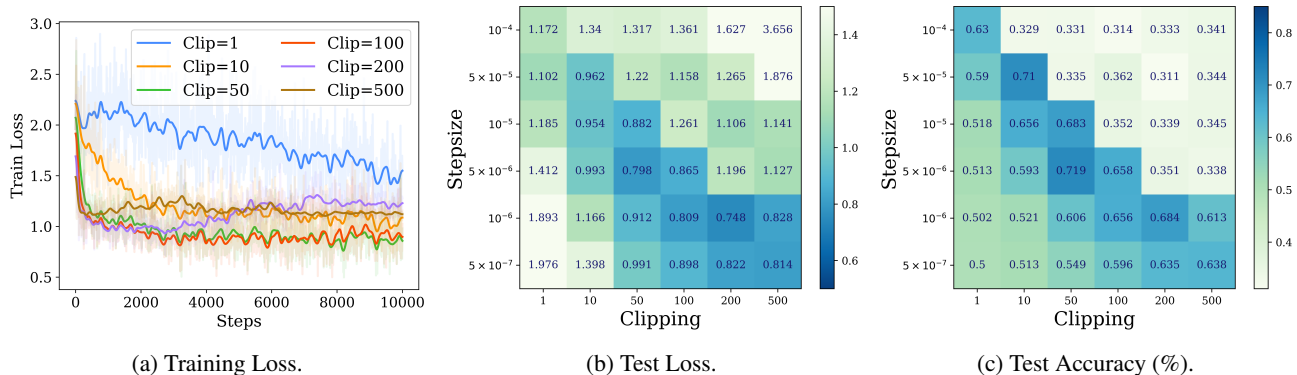


Figure 4. Experiments on private fine-tuning RoBERTa (125M) for SNLI with DPZERO. (a) (Smoothed) training curves when fixing the stepsize to be  $5 \times 10^{-6}$  and varying the clipping threshold from 1 to 500. In the choice of clipping, a tradeoff emerges; larger clipping values result in unnecessarily high privacy noise, while smaller values can induce increased bias in the optimization process. (b) and (c) Test loss and accuracy (%) when varying the stepsize and clipping threshold together. Consistent with first-order methods (Li et al., 2022b), we observe that larger clipping necessitates smaller stepsizes, whereas smaller clipping favors larger stepsizes.

**Comparison with DPGD-0th.** In the previous synthetic example, DPGD-0th suffers from worse performance in larger dimensions. To provide a more complete comparison, we also evaluate the performance of DPGD-0th (Algorithm 1) for fine-tuning RoBERTa-large on the dataset TREC with a privacy budget of  $\epsilon = 2$  (the same setting as Table 2). DPGD-0th only achieves a test accuracy of 67.0, while DPZERO attains 82.0. Moreover, DPGD-0th still requires per-sample clipping of the gradient estimator, which is costly in both memory and runtime compared to DPZERO.

**Clipping threshold.** Our findings indicate that the optimal clipping threshold for DPZERO tends to be higher than that for first-order methods. This observation aligns with the theoretical outcomes presented in Theorem 3, where the clipping threshold for DPZERO is  $C = \mathcal{O}(L\sqrt{\log(nd)})$ , in contrast to the  $\mathcal{O}(L)$  threshold adequate for first-order methods. In the concurrent study by (Tang et al., 2024a), the chosen clipping threshold is 0.05. However, their implementation applies the clipping to the term  $f(x + \lambda u; \xi) - f(x - \lambda u; \xi)$ . After normalization by  $\lambda = 10^{-3}$ , it aligns with the order of magnitude used in our method. The validity of opting for a larger clipping threshold in DPZERO is further confirmed through the private fine-tuning of RoBERTa (125M) on the SNLI dataset in Figure 4. An additional observation from our experiments is that the non-private baseline MeZO also appears to benefit from clipping. For instance, without clipping, the original MeZO encounters non-convergence issues at a stepsize of  $5 \times 10^{-6}$ . Conversely, incorporating clipping permits the use of larger stepsizes and yields better results. A thorough investigation of this phenomenon is reserved for future research.

### B.3. Private Fine-Tuning of the Language Model OPT

Table 10. Hyper-parameters used for fine-tuning OPT. We randomly sample 1000 samples for training and 1000 samples for testing. Stepsize and clipping are optimized through a grid search over the listed values. Other parameters are fixed to the values provided.

Hyper-parameters	Values
Number of training samples	1000
Number of test samples	1000
Number of iterations	20000
Batch size	8
Privacy	$(\epsilon = \{2, 6\}, \delta = 10^{-5})$
Smoothing $\lambda$	$10^{-3}$
Stepsize	$\{10^{-6}, 10^{-7}\}$
Clipping	$\{10, 50, 100, 200\}$

We follow experiment settings in Malladi et al. (2023) to evaluate the performance of DPZERO in the private fine-tuning of OPT (Zhang et al., 2022b) across four different datasets: SST-2 (Socher et al., 2013) for sentiment classification and

Table 11. Memory consumption (MiB) when fine-tuning OPT for BoolQ with batch size 8. All experiments are tested on a single GPU with 24 GiB memory. ‘-’ in the table denotes out of memory. MeZO and DPZERO can fit models up to OPT-6.7B, while the first-order method AdamW already runs out of memory on OPT-1.3B.

Method	OPT-1.3B	OPT-2.7B	OPT-6.7B	OPT-13B
AdamW	-	-	-	-
MeZO	7866	11602	20548	-
DPZERO	7866	11602	20548	-

BoolQ (Clark et al., 2019), SQuAD (Rajpurkar et al., 2016), and DROP (Dua et al., 2019) for question answering. In our experiments, we employ the same prompts as used in Malladi et al. (2023) and use the same implementation as explained before. All results are averaged over three random seeds  $\{0, 29, 83\}$ . The hyper-parameters used for our experiments are summarized in Table 10, and the memory usages on the dataset BoolQ are reported in Table 11.

### C. Technical Lemmas

**Lemma C.1.** *Let  $u$  be uniformly sampled from the Euclidean sphere  $\sqrt{d}\mathbb{S}^{d-1}$ ,  $a \in \mathbb{R}^d$  be some fixed vector independent of  $u$ , and  $H \in \mathbb{R}^{d \times d}$  be some fixed matrix independent of  $u$ . We have that*

(i)  $\mathbb{E}[u] = 0$  and  $\mathbb{E}[uu^\top] = \mathbf{I}_d$ .

(ii)  $\mathbb{E}_u[u^\top a] = 0$ ,  $\mathbb{E}_u[(u^\top a)^2] = \|a\|^2$  and  $\forall C \geq 0$ ,

$$\mathbb{P}(|u^\top a| \geq C) \leq 2\sqrt{2\pi} \exp\left(-\frac{C^2}{8\|a\|^2}\right).$$

(iii)  $\mathbb{E}_u[(u^\top a)u] = a$  and

$$\begin{aligned} \mathbb{E}_u[(u^\top a)^2 \|u\|^2] &= d\|a\|^2, \\ \mathbb{E}_u[(u^\top a)^2 uu^\top] &= \frac{d}{d+2} (2aa^\top + \|a\|^2 \mathbf{I}_d). \end{aligned}$$

(iv)  $\mathbb{E}_u[u^\top H u] = \text{Tr}(H)$  and

$$\mathbb{E}_u[(u^\top a)^2 u^\top H u] = \frac{d}{d+2} (2a^\top H a + \|a\|^2 \text{Tr}(H)).$$

*Proof.* (i) is a standard result, e.g., in Duchi et al. (2015), and follows by the symmetry of the sphere. For any  $u \in \sqrt{d} \cdot \mathbb{S}^{d-1}$ , it must be the case that  $-u \in \sqrt{d} \cdot \mathbb{S}^{d-1}$  as well, which suggests that  $\mathbb{E}[u] = 0$ . Since  $\mathbb{E}[\sum_{i=1}^d u_i^2] = \mathbb{E}\|u\|^2 = d$ , we immediately have that  $\mathbb{E}[u_i^2] = 1$  for every  $i$  by symmetry. Then for the off-diagonal terms, since for any  $u = (u_1, \dots, u_i, \dots, u_j, \dots, u_d) \in \sqrt{d} \cdot \mathbb{S}^{d-1}$ , it must be the case that  $(u_1, \dots, u_i, \dots, -u_j, \dots, u_d) \in \sqrt{d} \cdot \mathbb{S}^{d-1}$  as well, which suggests that  $\mathbb{E}[u_i u_j] = 0$  when  $i \neq j$ . As a result, we can conclude that the matrix  $\mathbb{E}[uu^\top] = \mathbf{I}_d$ .

We then show (ii). Applying (i), we have that  $\mathbb{E}_u[u^\top a] = 0$ , and that

$$\begin{aligned} \mathbb{E}_u[(u^\top a)^2] &= \sum_{i=1}^d a_i^2 \mathbb{E}[u_i^2] + \sum_{i \neq j} a_i a_j \mathbb{E}[u_i u_j] \\ &= \|a\|^2. \end{aligned}$$

The tail bound follows from Example 3.12 in Wainwright (2019), where they showed that for any function  $h : \mathbb{S}^{d-1} \rightarrow \mathbb{R}$  such that  $\forall x, y \in \mathbb{S}^{d-1}$ ,

$$|h(x) - h(y)| \leq \arccos(x^\top y),$$

when  $x$  is uniformly sampled from  $\mathbb{S}^{d-1}$ , it holds that  $\forall \gamma \geq 0$ ,

$$\mathbb{P}(|h(x) - \mathbb{E}[h(x)]| \geq \gamma) \leq 2\sqrt{2\pi} \exp\left(-\frac{d\gamma^2}{8}\right). \quad (8)$$

Let  $h(x) = x^\top a / \|a\|$  for  $x \in \mathbb{S}^{d-1}$ . First, we have that  $\forall x, y \in \mathbb{S}^{d-1}$ ,

$$\begin{aligned} |h(x) - h(y)|^2 &= \frac{|(x - y)^\top a|^2}{\|a\|^2} \\ &\leq \|x - y\|^2 \\ &= 2(1 - x^\top y) \\ &\leq (\arccos(x^\top y))^2, \end{aligned}$$

where we use the inequality that  $\theta^2/2 + \cos(\theta) - 1 \geq 0$  for  $\theta \in [0, \pi]$  and let  $x^\top y = \cos(\theta)$  such that  $\arccos(x^\top y) = \theta$  for some  $\theta \in [0, \pi]$ . When  $u$  is uniformly sampled from  $\sqrt{d} \cdot \mathbb{S}^{d-1}$ , we know  $u/\sqrt{d}$  is uniformly from  $\mathbb{S}^{d-1}$ . Applying (8) for  $h(x) = x^\top a / \|a\|$  where  $x \in \mathbb{S}^{d-1}$ , we obtain that

$$\mathbb{P}\left(\left|\frac{u^\top a}{\sqrt{d}\|a\|} - \frac{\mathbb{E}[u^\top a]}{\sqrt{d}\|a\|}\right| \geq \gamma\right) \leq 2\sqrt{2\pi} \exp\left(-\frac{d\gamma^2}{8}\right).$$

Setting  $C = \gamma\sqrt{d}\|a\|$ , the proof is complete since  $\mathbb{E}[u^\top a] = 0$ . Similar results also exist in Theorem 5.1.4 of [Vershynin \(2018\)](#), with all constants hidden behind some absolute  $c$ .

Next, we prove (iii). Applying (i), we have that

$$\begin{aligned} \mathbb{E}_u[(u^\top a)u_i] &= a_i \mathbb{E}[u_i^2] + \sum_{j \neq i} a_j \mathbb{E}[u_i u_j] \\ &= a_i. \end{aligned}$$

This implies that  $\mathbb{E}_u[(u^\top a)u] = a$ . Applying (ii), we obtain that

$$\begin{aligned} \mathbb{E}_u[(u^\top a)^2 \|u\|^2] &= d \cdot \mathbb{E}_u[(u^\top a)^2] \\ &= d\|a\|^2. \end{aligned}$$

For the expectation of the matrix, we start from the diagonal terms.

$$\begin{aligned} \mathbb{E}_u[(u^\top a)^2 u_i^2] &= \sum_{j=1}^d a_j^2 \mathbb{E}[u_j^2 u_i^2] + \sum_{j \neq k} a_j a_k \mathbb{E}[u_j u_k u_i^2] \\ &= a_i^2 \mathbb{E}[u_i^4] + \sum_{j \neq i} a_j^2 \mathbb{E}[u_j^2 u_i^2]. \end{aligned} \quad (9)$$

Here, we use the property that  $\mathbb{E}[u_j u_k u_i^2] = 0$  for every  $i$  when  $j \neq k$ . This follows from symmetry of the sphere such that for any  $u = (u_1, \dots, u_j, \dots, u_k, \dots, u_d) \in \sqrt{d} \cdot \mathbb{S}^{d-1}$ , it must be the case that  $(u_1, \dots, u_j, \dots, -u_k, \dots, u_d) \in \sqrt{d} \cdot \mathbb{S}^{d-1}$  as well. Again by symmetry, we have that  $\mathbb{E}[u_i^4]$  remains the same for every  $i$ , and  $\mathbb{E}[u_i^2 u_j^2]$  remains the same for every  $i \neq j$ . Denote  $w_1 = \mathbb{E}[u_i^4]$  and  $w_2 = \mathbb{E}[u_i^2 u_j^2]$ . Since it holds that

$$\begin{aligned} \sum_{i=1}^d \mathbb{E}_u[(u^\top a)^2 u_i^2] &= \mathbb{E}_u[(u^\top a)^2 \|u\|^2] \\ &= d\|a\|^2, \end{aligned}$$



taking summation over (9), we can have that

$$\begin{aligned}
 d\|a\|^2 &= \sum_{i=1}^d a_i^2 \mathbb{E}[u_i^4] + \sum_{i=1}^d \sum_{j=1, j \neq i}^d a_j^2 \mathbb{E}[u_j^2 u_i^2] \\
 &= w_1 \|a\|^2 + w_2 \sum_{i=1}^d (\|a\|^2 - a_i^2) \\
 &= w_1 \|a\|^2 + (d-1)w_2 \|a\|^2.
 \end{aligned}$$

This holds for arbitrary  $a \in \mathbb{R}^d$ , and thus we obtain that

$$w_1 + (d-1)w_2 = d. \quad (10)$$

We only compute  $w_1 = \mathbb{E}[u_i^4]$  by showing that  $u_i^2/d$  actually follows the Beta distribution, and the value of  $w_2$  can be derived from (10). First,  $z/\|z\|$  is uniformly distributed on the unit sphere  $\mathbb{S}^{d-1}$  for  $z \in \mathbb{R}^d$  sampled from the standard multivariate Gaussian  $\mathcal{N}(0, I_d)$  (Muller, 1959; Marsaglia, 1972). This means that  $z_i^2$  is distributed according to the  $\chi^2$ -distribution with 1 degree of freedom, and  $\bar{z}_i^2 := \sum_{j \neq i} z_j^2$  is distributed according to the  $\chi^2$ -distribution with degree  $(d-1)$ . Since  $\chi^2$ -distribution is a special case of the Gamma distribution and  $z_i^2, \bar{z}_i^2$  are independent, we conclude that  $z_i^2/(z_i^2 + \bar{z}_i^2)$  has the Beta distribution with parameters  $1/2$  and  $(d-1)/2$  (Cramér, 1999; Gupta & Nadarajah, 2004). Finally, since  $u/\sqrt{d}$  is uniformly distributed on  $\mathbb{S}^{d-1}$ , by symmetry of the sphere, we know that  $u_i^2/d$  has the same Beta distribution as  $z_i^2/(z_i^2 + \bar{z}_i^2)$ . The mean and variance of Beta( $1/2, (d-1)/2$ ) is  $1/d$  and  $2(d-1)/(d^2(d+2))$ . This suggests that  $\mathbb{E}[u_i^2] = 1$ , as already proved in (i), and that

$$\begin{aligned}
 w_1 &= \mathbb{E}[(u_i^2 - \mathbb{E}[u_i^2])^2] + (\mathbb{E}[u_i^2])^2 \\
 &= d^2 \left( \frac{2(d-1)}{d^2(d+2)} + \frac{1}{d^2} \right) \\
 &= \frac{3d}{d+2}.
 \end{aligned}$$

By (10), we know  $w_2 = d/(d+2)$ . According to (9), we have that the diagonal terms

$$\begin{aligned}
 \mathbb{E}_u[(u^\top a)^2 u_i^2] &= w_1 a_i^2 + w_2 (\|a\|^2 - a_i^2) \\
 &= \frac{2d}{d+2} a_i^2 + \frac{d}{d+2} \|a\|^2.
 \end{aligned}$$

Then we compute the off-diagonal entries for  $i \neq j$ . By the same reasoning as (9), we have that

$$\begin{aligned}
 \mathbb{E}_u[(u^\top a)^2 u_i u_j] &= \sum_{i \neq j} a_i a_j \mathbb{E}[u_i^2 u_j^2] \\
 &= \frac{2d}{d+2} a_i a_j.
 \end{aligned}$$

All other terms equal to 0 by symmetry of the sphere. Combining both diagonal and off-diagonal elements, we have that  $\mathbb{E}_u[(u^\top a)^2 u u^\top] = (d/(d+2))(2a a^\top + \|a\|^2 I_d)$ . Similar results are also shown in Appendix F of Malladi et al. (2023).

Finally, we give the proof of (iv). For the first statement, applying (i) in this lemma, we have that

$$\begin{aligned}
 \mathbb{E}_u[u^\top H u] &= \mathbb{E}[\text{Tr}(u u^\top H)] \\
 &= \text{Tr}(\mathbb{E}[u u^\top] \cdot H) \\
 &= \text{Tr}(H).
 \end{aligned}$$

Similarly for the second statement, we apply (iii) in this lemma and obtain that

$$\begin{aligned}
 \mathbb{E}_u [(u^\top a)^2 u^\top H u] &= \mathbb{E} \left[ (u^\top a)^2 \cdot \text{Tr}(u u^\top H) \right] \\
 &= \mathbb{E} \left[ \text{Tr} \left( (u^\top a)^2 u u^\top \cdot H \right) \right] \\
 &= \text{Tr} \left( \mathbb{E} \left[ (u^\top a)^2 u u^\top \right] \cdot H \right) \\
 &= \frac{2d}{d+2} \text{Tr}(a a^\top H) + \frac{d}{d+2} \|a\|^2 \text{Tr}(H) \\
 &= \frac{2d}{d+2} a^\top H a + \frac{d}{d+2} \|a\|^2 \text{Tr}(H).
 \end{aligned}$$

This concludes the proof.  $\square$

**Lemma C.2.** *Let  $u$  be uniformly sampled from the Euclidean sphere  $\sqrt{d}\mathbb{S}^{d-1}$  and  $v$  be uniformly sampled from the Euclidean ball  $\sqrt{d}\mathbb{B}^d = \{x \in \mathbb{R}^d \mid \|x\| \leq \sqrt{d}\}$ . For any function  $f(x) : \mathbb{R}^d \rightarrow \mathbb{R}$  and  $\lambda > 0$ , we define its zeroth-order gradient estimator as  $g_\lambda(x) = ((f(x + \lambda u) - f(x - \lambda u)) / (2\lambda))u$  and the smoothed function as  $f_\lambda(x) = \mathbb{E}_v[f(x + \lambda v)]$ . The following properties hold:*

(i)  $f_\lambda(x)$  is differentiable and  $\mathbb{E}_u[g_\lambda(x)] = \nabla f_\lambda(x)$ .

(ii) If  $f(x)$  is  $\ell$ -smooth, then we have that

$$\begin{aligned}
 \|\nabla f(x) - \nabla f_\lambda(x)\| &\leq \frac{\ell}{2} \lambda d^{3/2}, \\
 \mathbb{E}_u[\|g_\lambda(x)\|^2] &\leq 2d \cdot \|\nabla f(x)\|^2 + \frac{\ell^2}{2} \lambda^2 d^3.
 \end{aligned}$$

The above results are consistent with (iii) in Lemma C.1 when  $\lambda \rightarrow 0$  and  $f(x)$  is differentiable such that the two-point estimator reduces to the directional derivative  $g_0(x) = u^\top \nabla f(x)u$ .

*Proof.* We first show (i). Similarly to Lemma 10 in Shamir (2017), we have that

$$\mathbb{E}_{u \in \sqrt{d}\mathbb{S}^{d-1}}[g_\lambda(x)] = \mathbb{E}_{u \in \sqrt{d}\mathbb{S}^{d-1}} \left[ \frac{f(x + \lambda u)u}{\lambda} \right].$$

Applying Lemma 2.1 in Flaxman et al. (2005), we know

$$\mathbb{E}_{u' \in \mathbb{S}^{d-1}}[f(x + \lambda' u')u'] = \frac{\lambda'}{d} \nabla \mathbb{E}_{v' \in \mathbb{B}^d}[f(x + \lambda' v')].$$

Introducing  $u = \sqrt{d}u'$ ,  $v = \sqrt{d}v'$  and  $\lambda = \lambda' / \sqrt{d}$ , we thus obtain

$$\begin{aligned}
 \mathbb{E}_{u \in \sqrt{d}\mathbb{S}^{d-1}} \left[ \frac{f(x + \lambda u)u}{\lambda} \right] &= \mathbb{E}_{u' \in \mathbb{S}^{d-1}} \left[ \frac{f(x + \lambda' u')u'd}{\lambda'} \right] \\
 &= \nabla \mathbb{E}_{v' \in \mathbb{B}^d}[f(x + \lambda' v')] \\
 &= \nabla \mathbb{E}_{v \in \sqrt{d}\mathbb{B}^d}[f(x + \lambda v)].
 \end{aligned}$$

The proof of (ii) mostly follows from Nesterov & Spokoiny (2017), where the results are originally obtained for the case that  $u$  is sampled from the standard multivariate Gaussian distribution. By (iii) in Lemma C.1 and (i) here, we have that for

$u$  uniformly sampled from  $\sqrt{d} \cdot \mathbb{S}^{d-1}$ ,

$$\begin{aligned}
 \|\nabla f(x) - \nabla f_\lambda(x)\| &= \left\| \mathbb{E}_u[(u^\top \nabla f(x))u] - \mathbb{E}_u \left[ \frac{f(x + \lambda u) - f(x - \lambda u)}{2\lambda} u \right] \right\| \\
 &\leq \mathbb{E}_u \left\| \left( \frac{f(x + \lambda u) - f(x - \lambda u)}{2\lambda} - u^\top \nabla f(x) \right) u \right\| \\
 &\leq \frac{\sqrt{d}}{2\lambda} \mathbb{E}_u |f(x + \lambda u) - f(x) - \lambda u^\top \nabla f(x)| \\
 &\quad + \frac{\sqrt{d}}{2\lambda} \mathbb{E}_u |f(x) - f(x - \lambda u) - \lambda u^\top \nabla f(x)| \\
 &\leq \frac{\ell}{2} \lambda d^{3/2},
 \end{aligned}$$

where in the last step we use smoothness of  $f(x)$  such that  $|f(x + \lambda u) - f(x) - \lambda u^\top \nabla f(x)| \leq \ell \lambda^2 d/2$  and the same holds for  $|f(x) - f(x - \lambda u) - \lambda u^\top \nabla f(x)| = |f(x - \lambda u) - f(x) + \lambda u^\top \nabla f(x)|$ . The last statement holds similarly:

$$\begin{aligned}
 \mathbb{E}_u[\|g_\lambda(x)\|^2] &= \frac{d}{4\lambda^2} \mathbb{E}_u[(f(x + \lambda u) - f(x - \lambda u))^2] \\
 &\leq 2d \cdot \mathbb{E}_u[(u^\top \nabla f(x))^2] + \frac{d}{2\lambda^2} \mathbb{E}_u[(f(x + \lambda u) - f(x - \lambda u) - 2\lambda u^\top \nabla f(x))^2] \\
 &\leq 2d \cdot \mathbb{E}_u[(u^\top \nabla f(x))^2] + \frac{d}{\lambda^2} \mathbb{E}_u[(f(x + \lambda u) - f(x) - \lambda u^\top \nabla f(x))^2] \\
 &\quad + \frac{d}{\lambda^2} \mathbb{E}_u[(f(x) - f(x - \lambda u) - \lambda u^\top \nabla f(x))^2] \\
 &\leq 2d \cdot \|\nabla f(x)\|^2 + \frac{\ell^2}{2} \lambda^2 d^3,
 \end{aligned} \tag{11}$$

where in the last step we use Lemma C.1 and smoothness of  $f(x)$ .  $\square$

## D. Detailed Proof and Analysis of DPGD-0th (Algorithm 1)

*Proof of Theorem 1.* The privacy guarantees directly follow from Lemma 2.2 noticing that the sensitivity is  $2C/n$ . Note that the original advanced composition theorem in Kairouz et al. (2015) is stated for the case where the output of  $\mathcal{A}$  is a scalar. Given the spherical symmetry properties of Gaussian noise, the results can be readily extended to multiple dimensions, as outlined in Lemma 1 of Kenthapadi et al. (2013) where the basis can be selected in a way such that  $\mathcal{A}(S)$  and  $\mathcal{A}(S')$  differ in exactly one dimension.

We then focus on the utility guarantee on  $\mathbb{E}[\|\nabla F_S(x_\tau)\|^2]$ . Since  $f(x; \xi)$  is  $L$ -Lipschitz for every  $\xi$  by Assumption 3.1 and  $\|u_t\| = \sqrt{d}$  by its construction, we have that

$$\begin{aligned}
 \|g_\lambda(x_t; \xi_i)\| &= \frac{|f(x_t + \lambda u_t; \xi_i) - f(x_t - \lambda u_t; \xi_i)|}{2\lambda} \|u_t\| \\
 &\leq L \|u_t\|^2 \\
 &= Ld.
 \end{aligned}$$

This means  $\text{clip}_C(g_\lambda(x_t; \xi_i)) = g_\lambda(x_t; \xi_i)$  when setting  $C = Ld$ . For notation simplicity, we let

$$\begin{aligned}
 G_\lambda(x_t) &:= \frac{1}{n} \sum_{i=1}^n g_\lambda(x_t; \xi_i) \\
 &= \frac{1}{n} \sum_{i=1}^n \frac{f(x_t + \lambda u_t; \xi_i) - f(x_t - \lambda u_t; \xi_i)}{2\lambda} u_t \\
 &= \frac{F_S(x_t + \lambda u_t) - F_S(x_t - \lambda u_t)}{2\lambda} u_t.
 \end{aligned}$$

Algorithm 1 reduces to  $x_{t+1} = x_t - \alpha(G_\lambda(x_t) + z_t)$ . By smoothness of  $F_S(x)$ , we have that

$$\begin{aligned} F_S(x_{t+1}) &\leq F_S(x_t) + \nabla F_S(x_t)^\top (x_{t+1} - x_t) + \frac{\ell}{2} \|x_{t+1} - x_t\|^2 \\ &= F_S(x_t) - \alpha \nabla F_S(x_t)^\top (G_\lambda(x_t) + z_t) + \frac{\ell}{2} \alpha^2 \|G_\lambda(x_t)\|^2 + \frac{\ell}{2} \alpha^2 \|z_t\|^2 + \ell \alpha^2 z_t^\top G_\lambda(x_t). \end{aligned}$$

Since  $z_t$  is sampled from  $\mathcal{N}(0, \sigma^2 \mathbf{I}_d)$  and is independent of  $x_t$ ,  $u_t$  and  $S$ , we have that

$$\mathbb{E}_{z_t}[F_S(x_{t+1})] \leq F_S(x_t) - \alpha \nabla F_S(x_t)^\top G_\lambda(x_t) + \frac{\ell}{2} \alpha^2 \|G_\lambda(x_t)\|^2 + \frac{\ell}{2} \alpha^2 d \sigma^2.$$

Define  $F_\lambda(x) := \mathbb{E}_v[F_S(x + \lambda v)]$  for  $v$  sampled uniformly from the Euclidean ball  $\sqrt{d} \cdot \mathbb{B}^d$ . By Lemma C.2, we know  $\mathbb{E}_{u_t}[G_\lambda(x_t)] = \nabla F_\lambda(x_t)$ . Since  $u_t$  is independent of  $x_t$  and  $S$ , taking expectation with respect to  $u_t$  and applying (ii) in Lemma C.2, we obtain that

$$\begin{aligned} \mathbb{E}_{z_t, u_t}[F_S(x_{t+1})] &\leq F_S(x_t) - \alpha \nabla F_S(x_t)^\top \nabla F_\lambda(x_t) + \frac{\ell}{2} \alpha^2 \mathbb{E}_{u_t}[\|G_\lambda(x_t)\|^2] + \frac{\ell}{2} \alpha^2 d \sigma^2 \\ &= F_S(x_t) - \frac{\alpha}{2} \|\nabla F_S(x_t)\|^2 - \frac{\alpha}{2} \|\nabla F_\lambda(x_t)\|^2 + \frac{\alpha}{2} \|\nabla F_\lambda(x_t) - \nabla F_S(x_t)\|^2 \\ &\quad + \frac{\ell}{2} \alpha^2 \mathbb{E}_{u_t}[\|G_\lambda(x_t)\|^2] + \frac{\ell}{2} \alpha^2 d \sigma^2 \\ &\leq F_S(x_t) - \frac{\alpha}{2} (1 - 2d\ell\alpha) \|\nabla F_S(x_t)\|^2 + \frac{\ell^2}{8} \alpha (1 + 2\ell\alpha) \lambda^2 d^3 + \frac{\ell}{2} \alpha^2 d \sigma^2. \end{aligned} \quad (12)$$

Choosing  $\alpha = 1/(4\ell d)$  such that  $1 - 2d\ell\alpha = 1/2$  and  $2\ell\alpha < 1$ , we obtain that

$$\begin{aligned} \mathbb{E}[\|\nabla F_S(x_t)\|^2] &< \frac{4 \mathbb{E}[F_S(x_t) - F_S(x_{t+1})]}{\alpha} + \ell^2 \lambda^2 d^3 + 2\ell\alpha d \sigma^2 \\ &= \frac{4 \mathbb{E}[F_S(x_t) - F_S(x_{t+1})]}{\alpha} + \ell^2 \lambda^2 d^3 + \frac{64\ell C^2 \alpha T d \log(e + (\varepsilon/\delta))}{n^2 \varepsilon^2} \\ &= \frac{4 \mathbb{E}[F_S(x_t) - F_S(x_{t+1})]}{\alpha} + \ell^2 \lambda^2 d^3 + \frac{64\ell L^2 \alpha T d^3 \log(e + (\varepsilon/\delta))}{n^2 \varepsilon^2}. \end{aligned}$$

As a result, taking summation from  $t = 0$  to  $T - 1$  and dividing both sides by  $T$ , we have that

$$\begin{aligned} \mathbb{E}[\|\nabla F_S(x_\tau)\|^2] &= \frac{1}{T} \sum_{t=0}^{T-1} \mathbb{E}[\|\nabla F_S(x_t)\|^2] \\ &\leq \frac{4(F_S(x_0) - F_S^*)}{\alpha T} + \ell^2 \lambda^2 d^3 + \frac{64\ell L^2 \alpha T d^3 \log(e + (\varepsilon/\delta))}{n^2 \varepsilon^2} \\ &\leq \frac{16(\ell(F_S(x_0) - F_S^*) + 2L^2)d\sqrt{d\log(e + (\varepsilon/\delta))}}{n\varepsilon}, \end{aligned}$$

with the choice of parameters

$$\alpha T = \frac{n\varepsilon}{4\ell d \sqrt{d\log(e + (\varepsilon/\delta))}}, \quad \lambda \leq \frac{4L}{\ell d} \left( \frac{\sqrt{d\log(e + (\varepsilon/\delta))}}{n\varepsilon} \right)^{1/2}.$$

This suggests that the total number of iteration is  $T = n\varepsilon/\sqrt{d\log(e + (\varepsilon/\delta))}$  and the total number of zeroth-order gradient computations is  $nT = n^2\varepsilon/\sqrt{d\log(e + (\varepsilon/\delta))}$ . Note that the above selection of parameters ensures scale invariance.  $\square$

*Proof of Theorem 2.* The privacy analysis remains the same as before, and we focus on the utility analysis on  $\mathbb{E}\|\nabla F_S(x_\tau)\|^2$ . By the same reasoning, when setting  $C = Ld$ , Algorithm 1 reduces to  $x_{t+1} = x_t - \alpha(G_\lambda(x_t) + z_t)$  where  $G_\lambda(x_t) =$

$(F_S(x_t + \lambda u_t) - F_S(x_t - \lambda u_t))u_t/(2\lambda)$ . By Taylor's theorem with remainder, for some  $\theta \in (0, 1)$ , we have that

$$\begin{aligned} F_S(x_{t+1}) &= F_S(x_t) + \nabla F_S(x_t)^\top (x_{t+1} - x_t) + \frac{1}{2}(x_{t+1} - x_t)^\top \nabla^2 F_S(x_t + \theta(x_{t+1} - x_t))(x_{t+1} - x_t) \\ &\leq F_S(x_t) - \alpha \nabla F_S(x_t)^\top (G_\lambda(x_t) + z_t) + \frac{\alpha^2}{2} G_\lambda(x_t)^\top H G_\lambda(x_t) + \frac{\alpha^2}{2} z_t^\top H z_t \\ &\quad + \frac{\alpha^2}{2} (G_\lambda(x_t)^\top H z_t + z_t^\top H G_\lambda(x_t)). \end{aligned}$$

Here in the inequality, we use Assumption 3.5 such that  $\nabla^2 F_S(x) \preceq H$  for any  $x \in \mathbb{R}^d$ . Similarly to (iv) in Lemma C.1, we have that  $\mathbb{E}[z_t^\top H z_t] = \text{Tr}(\mathbb{E}[z_t z_t^\top] H) = \sigma^2 \text{Tr}(H)$ . Since  $z_t$  is sampled from  $\mathcal{N}(0, \sigma^2 I_d)$  and is independent of  $u_t, x_t$  and the dataset  $S$ , taking expectation with respect to  $z_t$ , we can then obtain that

$$\begin{aligned} \mathbb{E}_{z_t}[F_S(x_{t+1})] &\leq F_S(x_t) - \alpha \nabla F_S(x_t)^\top G_\lambda(x_t) + \frac{\alpha^2}{2} G_\lambda(x_t)^\top H G_\lambda(x_t) + \frac{\alpha^2}{2} \mathbb{E}_{z_t}[z_t^\top H z_t] \\ &= F_S(x_t) - \alpha \nabla F_S(x_t)^\top G_\lambda(x_t) + \frac{\alpha^2}{2} G_\lambda(x_t)^\top H G_\lambda(x_t) + \frac{\alpha^2 \sigma^2}{2} \text{Tr}(H). \end{aligned} \quad (13)$$

Assumption 3.5 implies  $F_S(x)$  is also  $\ell$ -smooth. By a similar argument as (11) in the proof of (ii) in Lemma C.2, we have

$$\left( \frac{F_S(x_t + \lambda u_t) - F_S(x_t - \lambda u_t)}{2\lambda} \right)^2 \leq 2 (u_t^\top \nabla F_S(x_t))^2 + \frac{\ell^2}{2} \lambda^2 d^2. \quad (14)$$

As  $u_t^\top H u_t \geq 0$ , by (iv) in Lemma C.1 and Assumption 3.5, we have that

$$\begin{aligned} \mathbb{E}[G_\lambda(x_t)^\top H G_\lambda(x_t)] &= \mathbb{E} \left[ \left( \frac{F_S(x_t + \lambda u_t) - F_S(x_t - \lambda u_t)}{2\lambda} \right)^2 u_t^\top H u_t \right] \\ &\leq 2 \mathbb{E} \left[ (u_t^\top \nabla F_S(x_t))^2 u_t^\top H u_t \right] + \frac{\ell^2}{2} \lambda^2 d^2 \mathbb{E}[u_t^\top H u_t] \\ &= \frac{2d}{d+2} (2 \nabla F_S(x_t)^\top H \nabla F_S(x_t) + \|\nabla F_S(x_t)\|^2 \text{Tr}(H)) + \frac{\ell^2}{2} \lambda^2 d^2 \text{Tr}(H) \\ &\leq 2\ell(r+2) \|\nabla F_S(x_t)\|^2 + \frac{\ell^3}{2} \lambda^2 d^2 r. \end{aligned}$$

Taking expectation of (13) with respect to  $u_t$ , by Lemma C.2 for  $F_\lambda(x) = \mathbb{E}_v[F_S(x + \lambda v)]$  with  $v$  uniformly sampled from  $\sqrt{d} \cdot \mathbb{B}^d$ , we have that

$$\begin{aligned} \mathbb{E}[F_S(x_{t+1})] &\leq F_S(x_t) - \alpha \nabla F_S(x_t)^\top \nabla F_\lambda(x_t) + \ell \alpha^2 (r+2) \|\nabla F_S(x_t)\|^2 + \frac{\ell^3 \alpha^2 \lambda^2 d^2 r}{4} + \frac{\ell \alpha^2 r \sigma^2}{2} \\ &\leq F_S(x_t) - \frac{\alpha}{2} (1 - 2(r+2)\ell \alpha) \|\nabla F_S(x_t)\|^2 + \frac{\alpha}{2} \|\nabla F_S(x_t) - \nabla F_\lambda(x_t)\|^2 + \frac{\ell^3 \alpha^2 \lambda^2 d^2 r}{4} + \frac{\ell \alpha^2 r \sigma^2}{2} \\ &\leq F_S(x_t) - \frac{\alpha}{2} (1 - 2(r+2)\ell \alpha) \|\nabla F_S(x_t)\|^2 + \frac{\ell^2 \alpha \lambda^2 d^2 (d + 2r\ell \alpha)}{8} + \frac{\ell \alpha^2 r \sigma^2}{2}. \end{aligned} \quad (15)$$

Choosing  $\alpha = 1/(4\ell(r+2))$  such that  $1 - 2(r+2)\ell \alpha = 1/2$  and  $2\ell \alpha r < 1 \leq d$ , we have that

$$\begin{aligned} \mathbb{E}[\|\nabla F_S(x_t)\|^2] &< \frac{4 \mathbb{E}[F_S(x_t) - F_S(x_{t+1})]}{\alpha} + \ell^2 \lambda^2 d^3 + 2\ell \alpha r \sigma^2 \\ &= \frac{4 \mathbb{E}[F_S(x_t) - F_S(x_{t+1})]}{\alpha} + \ell^2 \lambda^2 d^3 + \frac{64\ell C^2 \alpha T r \log(e + (\varepsilon/\delta))}{n^2 \varepsilon^2} \\ &= \frac{4 \mathbb{E}[F_S(x_t) - F_S(x_{t+1})]}{\alpha} + \ell^2 \lambda^2 d^3 + \frac{64\ell L^2 \alpha T d^2 r \log(e + (\varepsilon/\delta))}{n^2 \varepsilon^2}. \end{aligned}$$



As a result, taking summation from  $t = 0$  to  $T - 1$  and dividing both sides by  $T$ , we have that

$$\begin{aligned}\mathbb{E}[\|\nabla F_S(x_\tau)\|^2] &= \frac{1}{T} \sum_{t=0}^{T-1} \mathbb{E}[\|\nabla F_S(x_t)\|^2] \\ &\leq \frac{4(F_S(x_0) - F_S^*)}{\alpha T} + \ell^2 \lambda^2 d^3 + \frac{64\ell L^2 \alpha T d^2 r \log(e + (\varepsilon/\delta))}{n^2 \varepsilon^2} \\ &\leq \frac{16(\ell(F_S(x_0) - F_S^*) + 2L^2)d\sqrt{r \log(e + (\varepsilon/\delta))}}{n\varepsilon},\end{aligned}$$

with the choice of parameters

$$\alpha T = \frac{n\varepsilon}{4\ell d\sqrt{r \log(e + (\varepsilon/\delta))}}, \quad \lambda \leq \frac{4L}{\ell d} \left( \frac{\sqrt{r \log(e + (\varepsilon/\delta))}}{n\varepsilon} \right)^{1/2}.$$

This suggests that the total number of iteration is  $T = n(r+2)\varepsilon/(d\sqrt{r \log(e + (\varepsilon/\delta))})$  and the total number of zeroth-order gradient computations is  $nT = n^2(r+2)\varepsilon/(d\sqrt{r \log(e + (\varepsilon/\delta))})$ . The above selection ensures scale invariance.  $\square$

## E. Detailed Proof and Analysis of DPZero (Algorithm 2)

**Privacy guarantee.** Since  $u_t$  is independent of the dataset  $S$ , the privacy guarantees directly follow from Lemma 2.2 and post-processing (Dwork et al., 2014) noticing that the sensitivity is  $2C/n$ . We want to emphasize that the randomness of  $u_t$  is never used for the privacy guarantee, and the analysis holds for any  $u_t$  as long as it is independent of the dataset.

**Utility guarantee.** We then focus on the utility guarantee on  $\mathbb{E}\|\nabla F_S(x_\tau)\|^2$ . Since  $f(x; \xi)$  is  $\ell$ -smooth for every  $\xi$  by Assumption 3.5, we have that

$$\begin{aligned}\frac{|f(x_t + \lambda u_t; \xi_i) - f(x_t - \lambda u_t; \xi_i)|}{2\lambda} &\leq |u_t^\top \nabla f(x_t; \xi_i)| + \frac{|f(x_t + \lambda u_t; \xi_i) - f(x_t; \xi_i) - \lambda u_t^\top \nabla f(x_t; \xi_i)|}{2\lambda} \\ &\quad + \frac{|f(x_t - \lambda u_t; \xi_i) - f(x_t; \xi_i) + \lambda u_t^\top \nabla f(x_t; \xi_i)|}{2\lambda} \\ &\leq |u_t^\top \nabla f(x_t; \xi_i)| + \frac{\ell}{2} \lambda d.\end{aligned}\tag{16}$$

Therefore, by (ii) in Lemma C.1 and Lipschitzness of  $f(x; \xi)$ , we have that

$$\begin{aligned}\mathbb{P}\left(\frac{|f(x_t + \lambda u_t; \xi_i) - f(x_t - \lambda u_t; \xi_i)|}{2\lambda} \geq C_0 + \frac{\ell}{2} \lambda d\right) &\leq \mathbb{P}(|u_t^\top \nabla f(x_t; \xi_i)| \geq C_0) \\ &\leq 2\sqrt{2\pi} \exp\left(-\frac{C_0^2}{8\|\nabla f(x_t; \xi_i)\|^2}\right) \\ &\leq 2\sqrt{2\pi} \exp\left(-\frac{C_0^2}{8L^2}\right).\end{aligned}$$

We define  $Q_{t,i}$  to be the event that the clipping does not happen at iteration  $t$  for sample  $\xi_i$  and  $\bar{Q}_{t,i}$  to be the event that the clipping does happen. The above equation implies that if the clipping threshold  $C \geq C_0 + \ell\lambda d/2$ , then we have that  $\mathbb{P}(Q_{t,i}) \leq 2\sqrt{2\pi} \exp(-C_0^2/(8L^2))$ . Let  $Q_t$  denote the event that the clipping does not happen at iteration  $t$  for every sample  $1 \leq i \leq n$ , and let  $\bar{Q}_t$  be the event that there exist some  $i$  such that the clipping does happen at iteration  $t$ . We also denote  $Q$  as the event that the clipping does not happen for every iteration  $t = 0, 1, \dots, T - 1$  and every sample  $1 \leq i \leq n$  and  $\bar{Q}$  as the event that there exist some  $t$  and  $i$  such that the clipping does happen. By the union bound, we have that

$$\begin{aligned}\mathbb{P}(\bar{Q}) &= \mathbb{P}\left(\bigcup_{t=0}^{T-1} \bigcup_{i=1}^n \bar{Q}_{t,i}\right) \\ &\leq 2\sqrt{2\pi} nT \exp\left(-\frac{C_0^2}{8L^2}\right).\end{aligned}$$

To simplify the notation, we let

$$\begin{aligned} G_\lambda(x_t) &= \frac{1}{n} \sum_{i=1}^n \frac{f(x_t + \lambda u_t; \xi_i) - f(x_t - \lambda u_t; \xi_i)}{2\lambda} u_t \\ &= \frac{F_S(x_t + \lambda u_t) - F_S(x_t - \lambda u_t)}{2\lambda} u_t, \end{aligned}$$

and its per-sample clipped version as

$$\hat{G}_\lambda(x_t) = \frac{1}{n} \sum_{i=1}^n \text{clip}_C \left( \frac{f(x_t + \lambda u_t; \xi_i) - f(x_t - \lambda u_t; \xi_i)}{2\lambda} \right) u_t.$$

Algorithm 2 becomes  $x_{t+1} = x_t - \alpha(\hat{G}_\lambda(x_t) + z_t u_t)$  under the above notation. By Taylor's theorem with remainder, for some  $\theta \in (0, 1)$ , we have that

$$\begin{aligned} F_S(x_{t+1}) &= F_S(x_t) + \nabla F_S(x_t)^\top (x_{t+1} - x_t) + \frac{1}{2} (x_{t+1} - x_t)^\top \nabla^2 F_S(x_t + \theta(x_{t+1} - x_t)) (x_{t+1} - x_t) \\ &\leq F_S(x_t) - \alpha \nabla F_S(x_t)^\top (\hat{G}_\lambda(x_t) + z_t u_t) + \frac{\alpha^2}{2} \hat{G}_\lambda(x_t)^\top H \hat{G}_\lambda(x_t) + \frac{\alpha^2}{2} z_t^2 u_t^\top H u_t \\ &\quad + \frac{\alpha^2}{2} z_t (\hat{G}_\lambda(x_t)^\top H u_t + u_t^\top H \hat{G}_\lambda(x_t)). \end{aligned}$$

Here in the inequality, we use Assumption 3.5 such that  $\nabla^2 F_S(x) \preceq H$  for any  $x \in \mathbb{R}^d$ . The event  $Q_t$  depends on the randomness in  $u_{<(t+1)} := \{u_0, u_1, \dots, u_t\}$  and  $z_{<t} := \{z_0, z_1, \dots, z_{t-1}\}$ . Note that the scalar noise  $z_t$  sampled from  $\mathcal{N}(0, \sigma^2)$  is independent of  $u_{<(t+1)}$ ,  $z_{<t}$ ,  $x_t$ , and the dataset  $S$ . Conditioned on the event  $Q_t$  and taking expectation with respect to  $z_{<(t+1)}$  and  $u_{<(t+1)}$ , we have that

$$\begin{aligned} \mathbb{E}_{z_{<(t+1)}, u_{<(t+1)}} [F_S(x_{t+1}) | Q_t] &\leq \mathbb{E}_{z_{<t}, u_{<t}} [F_S(x_t) | Q_t] - \alpha \mathbb{E}_{z_{<t}, u_{<(t+1)}} [\nabla F_S(x_t)^\top \hat{G}_\lambda(x_t) | Q_t] \\ &\quad + \frac{\alpha^2}{2} \mathbb{E}_{z_{<t}, u_{<(t+1)}} [\hat{G}_\lambda(x_t)^\top H \hat{G}_\lambda(x_t) | Q_t] + \frac{\alpha^2 \sigma^2}{2} \mathbb{E}_{z_{<t}, u_{<(t+1)}} [u_t^\top H u_t | Q_t]. \end{aligned} \quad (17)$$

Let  $\mathbb{E}_t := \mathbb{E}_{z_{<t}, u_{<(t+1)}}$  for simplicity. Given the condition that  $Q_t$  happens, we know that  $\hat{G}_\lambda(x_t) = G_\lambda(x_t)$  and thus

$$\mathbb{E}_t \left[ \hat{G}_\lambda(x_t)^\top H \hat{G}_\lambda(x_t) \mid Q_t \right] = \mathbb{E}_t \left[ \left( \frac{F_S(x_t + \lambda u_t) - F_S(x_t - \lambda u_t)}{2\lambda} \right)^2 u_t^\top H u_t \mid Q_t \right].$$

Since  $H \succeq 0$ , we have that  $u_t^\top H u_t \geq 0$ . By the law of total probability, we obtain

$$\begin{aligned} &\mathbb{E}_t \left[ \left( \frac{F_S(x_t + \lambda u_t) - F_S(x_t - \lambda u_t)}{2\lambda} \right)^2 u_t^\top H u_t \right] \\ &= \mathbb{E}_t \left[ \left( \frac{F_S(x_t + \lambda u_t) - F_S(x_t - \lambda u_t)}{2\lambda} \right)^2 u_t^\top H u_t \mid Q_t \right] \mathbb{P}(Q_t) \\ &\quad + \mathbb{E}_t \left[ \left( \frac{F_S(x_t + \lambda u_t) - F_S(x_t - \lambda u_t)}{2\lambda} \right)^2 u_t^\top H u_t \mid \bar{Q}_t \right] \mathbb{P}(\bar{Q}_t) \\ &\geq \mathbb{E}_t \left[ \left( \frac{F_S(x_t + \lambda u_t) - F_S(x_t - \lambda u_t)}{2\lambda} \right)^2 u_t^\top H u_t \mid Q_t \right] \mathbb{P}(Q_t). \end{aligned} \quad (18)$$

Assumption 3.5 implies  $F_S(x)$  is also  $\ell$ -smooth. Similarly to the proof of Theorem 2, by (14) and the fact that  $u_t^\top H u_t \geq 0$ , applying (iv) in Lemma C.1 and Assumption 3.5, we can then obtain that

$$\begin{aligned}
 \mathbb{E}_t \left[ \hat{G}_\lambda(x_t)^\top H \hat{G}_\lambda(x_t) \mid Q_t \right] &\leq \frac{\mathbb{E}_t \left[ (F_S(x_t + \lambda u_t) - F_S(x_t - \lambda u_t))^2 u_t^\top H u_t \right]}{4\lambda^2 \cdot \mathbb{P}(Q_t)} \\
 &\leq \frac{\mathbb{E}_t \left[ 2(u_t^\top \nabla F_S(x_t))^2 u_t^\top H u_t \right]}{\mathbb{P}(Q_t)} + \frac{\ell^2 \lambda^2 d^2}{2\mathbb{P}(Q_t)} \mathbb{E}_t \left[ u_t^\top H u_t \right] \\
 &= \frac{2d \mathbb{E}_{z < t, u < t} \left[ 2\nabla F_S(x_t)^\top H \nabla F_S(x_t) + \|\nabla F_S(x_t)\|^2 \text{Tr}(H) \right]}{(d+2)\mathbb{P}(Q_t)} + \frac{\ell^2 \lambda^2 d^2 \text{Tr}(H)}{2\mathbb{P}(Q_t)} \\
 &\leq \frac{2\ell(r+2)}{\mathbb{P}(Q_t)} \mathbb{E}_{z < t, u < t} \|\nabla F_S(x_t)\|^2 + \frac{\ell^3 \lambda^2 d^2 r}{2\mathbb{P}(Q_t)}. \tag{19}
 \end{aligned}$$

The same as (18), we can also get that

$$\begin{aligned}
 \mathbb{E}_t \left[ u_t^\top H u_t \mid Q_t \right] &\leq \frac{\mathbb{E}_t \left[ u_t^\top H u_t \right]}{\mathbb{P}(Q_t)} \\
 &\leq \frac{r\ell}{\mathbb{P}(Q_t)}. \tag{20}
 \end{aligned}$$

For the inner-product term, we have that

$$\mathbb{E}_t \left[ \nabla F_S(x_t)^\top \hat{G}_\lambda(x_t) \mid Q_t \right] = \mathbb{E}_t \left[ \nabla F_S(x_t)^\top G_\lambda(x_t) \mid Q_t \right].$$

By the law of total probability, since  $u_t$  is independent of  $x_t$ , we know that

$$\begin{aligned}
 \mathbb{E}_t \left[ \nabla F_S(x_t)^\top G_\lambda(x_t) \mid Q_t \right] \mathbb{P}(Q_t) + \mathbb{E}_t \left[ \nabla F_S(x_t)^\top G_\lambda(x_t) \mid \bar{Q}_t \right] \mathbb{P}(\bar{Q}_t) &= \mathbb{E}_t \left[ \nabla F_S(x_t)^\top G_\lambda(x_t) \right] \\
 &= \mathbb{E}_{z < t, u < t} \left[ \nabla F_S(x_t)^\top \nabla F_\lambda(x_t) \right],
 \end{aligned}$$

where we use Lemma C.2 for  $F_\lambda(x) = \mathbb{E}_v[F_S(x + \lambda v)]$  with  $v$  uniformly sampled from  $\sqrt{d}\mathbb{B}^d$ . Rearranging terms, we thus obtain that

$$\begin{aligned}
 \mathbb{E}_t \left[ \nabla F_S(x_t)^\top G_\lambda(x_t) \mid Q_t \right] &= \frac{\mathbb{E}_{z < t, u < t} \left[ \nabla F_S(x_t)^\top \nabla F_\lambda(x_t) \right]}{\mathbb{P}(Q_t)} - \frac{\mathbb{E}_t \left[ \nabla F_S(x_t)^\top G_\lambda(x_t) \mid \bar{Q}_t \right] \mathbb{P}(\bar{Q}_t)}{\mathbb{P}(Q_t)} \\
 &= \frac{\mathbb{E}_{z < t, u < t} \|\nabla F_S(x_t)\|^2}{2\mathbb{P}(Q_t)} + \frac{\mathbb{E}_{z < t, u < t} \|\nabla F_\lambda(x_t)\|^2}{2\mathbb{P}(Q_t)} - \frac{\mathbb{E}_{z < t, u < t} \|\nabla F_S(x_t) - \nabla F_\lambda(x_t)\|^2}{2\mathbb{P}(Q_t)} \\
 &\quad - \frac{\mathbb{E}_t \left[ \nabla F_S(x_t)^\top G_\lambda(x_t) \mid \bar{Q}_t \right] \mathbb{P}(\bar{Q}_t)}{\mathbb{P}(Q_t)} \\
 &\geq \frac{\mathbb{E}_{z < t, u < t} \|\nabla F_S(x_t)\|^2}{2\mathbb{P}(Q_t)} - \frac{\ell^2 \lambda^2 d^3}{8\mathbb{P}(Q_t)} - \frac{\mathbb{E}_t \left[ \nabla F_S(x_t)^\top G_\lambda(x_t) \mid \bar{Q}_t \right] \mathbb{P}(\bar{Q}_t)}{\mathbb{P}(Q_t)},
 \end{aligned}$$

where we apply (ii) in Lemma C.2. Assumption 3.5 implies that  $F_S(x)$  is also Lipschitz, and thus

$$\begin{aligned}
 \nabla F_S(x_t)^\top G_\lambda(x_t) &\leq \|\nabla F_S(x_t)\| \|G_\lambda(x_t)\| \\
 &\leq L^2 \|u_t\|^2 \\
 &= L^2 d.
 \end{aligned}$$

As a result, we obtain that

$$\mathbb{E}_t \left[ \nabla F_S(x_t)^\top \hat{G}_\lambda(x_t) \mid Q_t \right] \geq \frac{\mathbb{E}_{z < t, u < t} \|\nabla F_S(x_t)\|^2}{2\mathbb{P}(Q_t)} - \frac{\ell^2 \lambda^2 d^3}{8\mathbb{P}(Q_t)} - \frac{L^2 d \mathbb{P}(\bar{Q}_t)}{\mathbb{P}(Q_t)}. \tag{21}$$

Plugging (21), (19) and (20) back into (17), we obtain that

$$\begin{aligned} \mathbb{E}_{z_{<(t+1)}, u_{<(t+1)}} [F_S(x_{t+1}) | Q_t] &\leq \mathbb{E}_{z_{<t}, u_{<t}} [F_S(x_t) | Q_t] - \frac{\alpha}{2} (1 - 2(r+2)\ell\alpha) \frac{\mathbb{E}_{z_{<t}, u_{<t}} \|\nabla F_S(x_t)\|^2}{\mathbb{P}(Q_t)} + \frac{\ell\alpha^2 r\sigma^2}{2\mathbb{P}(Q_t)} \\ &\quad + \frac{\ell^2\alpha(d+2\ell\alpha r)\lambda^2 d^2}{8\mathbb{P}(Q_t)} + \frac{\alpha L^2 d \mathbb{P}(\bar{Q}_t)}{\mathbb{P}(Q_t)}. \end{aligned} \quad (22)$$

Choosing  $\alpha = 1/(4\ell(r+2))$  such that  $1 - 2(r+2)\ell\alpha = 1/2$  and  $2\ell\alpha r < 1 \leq d$ , we have that

$$\begin{aligned} \mathbb{E}_{z_{<t}, u_{<t}} \|\nabla F_S(x_t)\|^2 &\leq \frac{4\mathbb{E}_{z_{<(t+1)}, u_{<(t+1)}} [F_S(x_t) - F_S(x_{t+1}) | Q_t] \mathbb{P}(Q_t)}{\alpha} + 2\ell\alpha r\sigma^2 + \ell^2 d^3 \lambda^2 + 4L^2 d \mathbb{P}(\bar{Q}_t) \\ &\leq \frac{4\mathbb{E}_{z_{<(t+1)}, u_{<(t+1)}} [F_S(x_t) - F_S(x_{t+1}) | Q_t] \mathbb{P}(Q_t)}{\alpha} + 2\ell\alpha r\sigma^2 + \ell^2 d^3 \lambda^2 + 4L^2 d \mathbb{P}(\bar{Q}). \end{aligned}$$

Recall  $Q_t$  is the event that clipping does not happen at iteration  $t$  and  $Q$  is the event that clipping does not happen for every iteration. By the law of total probability and the assumption that  $|F_S(x_t)| \leq B$  for every  $t$ , we have that

$$\begin{aligned} \mathbb{E}_{z_{<(t+1)}, u_{<(t+1)}} [F_S(x_t) - F_S(x_{t+1}) | Q_t] \mathbb{P}(Q_t) &= \mathbb{E}_{z_{<T}, u_{<T}} [F_S(x_t) - F_S(x_{t+1}) | Q_t] \mathbb{P}(Q_t) \\ &= \mathbb{E}_{z_{<T}, u_{<T}} [F_S(x_t) - F_S(x_{t+1}) | Q_t \cap Q] \mathbb{P}(Q_t \cap Q) \\ &\quad + \mathbb{E}_{z_{<T}, u_{<T}} [F_S(x_t) - F_S(x_{t+1}) | Q_t \cap \bar{Q}] \mathbb{P}(Q_t \cap \bar{Q}) \\ &\leq \mathbb{E}_{z_{<T}, u_{<T}} [F_S(x_t) - F_S(x_{t+1}) | Q] \mathbb{P}(Q) + 2B \mathbb{P}(\bar{Q}). \end{aligned}$$

As a result, we have that

$$\mathbb{E}_{z_{<t}, u_{<t}} \|\nabla F_S(x_t)\|^2 \leq \frac{4\mathbb{E}_{z_{<T}, u_{<T}} [F_S(x_t) - F_S(x_{t+1}) | Q] \mathbb{P}(Q)}{\alpha} + 2\ell\alpha r\sigma^2 + \ell^2 d^3 \lambda^2 + \left(4L^2 d + \frac{8B}{\alpha}\right) \mathbb{P}(\bar{Q}). \quad (23)$$

Taking expectation with respect to all randomness, i.e.,  $\mathbb{E} = \mathbb{E}_{z_{<T}, u_{<T}}$ , summing up from  $t = 0$  to  $T - 1$ , and dividing both sides by  $T$ , we have that

$$\begin{aligned} \mathbb{E} \|\nabla F_S(x_T)\|^2 &= \frac{1}{T} \sum_{t=0}^{T-1} \mathbb{E}_{z_{<t}, u_{<t}} \|\nabla F_S(x_t)\|^2 \\ &\leq \frac{4\mathbb{E}[F_S(x_0) - F_S(x_T) | Q] \mathbb{P}(Q)}{\alpha T} + \frac{64\ell C^2 \alpha T r \log(e + (\varepsilon/\delta))}{n^2 \varepsilon^2} + \ell^2 d^3 \lambda^2 \\ &\quad + 8\sqrt{2\pi} n T (L^2 d + 8\ell B(r+2)) \exp\left(-\frac{C_0^2}{8L^2}\right) \\ &\leq \left(64\ell[F_S(x_0) - F_S^*] + 4C^2\right) \frac{\sqrt{r \log(e + (\varepsilon/\delta))}}{n\varepsilon} + \ell^2 d^3 \lambda^2 \\ &\quad + \frac{2\sqrt{2\pi} n^2 \varepsilon (r+2) (L^2 d + 8\ell B(r+2))}{\sqrt{r \log(e + (\varepsilon/\delta))}} \exp\left(-\frac{C_0^2}{8L^2}\right), \end{aligned}$$

with the choice of parameters to be

$$\alpha T = \frac{n\varepsilon}{16\ell\sqrt{r \log(e + (\varepsilon/\delta))}}, \quad \alpha = \frac{1}{4\ell(r+2)}, \quad T = \frac{n(r+2)\varepsilon}{4\sqrt{r \log(e + (\varepsilon/\delta))}}.$$

When selecting  $\lambda \leq 2(\sqrt{2} - 1)C_0/(\ell d)$ , we can set  $C = \sqrt{2}C_0$  such that  $C \geq C_0 + \ell\lambda d/2$  is satisfied. If  $C_0$  and  $\lambda$  further satisfy the conditions that

$$C_0^2 = 8L^2 \log\left(\frac{2\sqrt{2\pi} n^3 \varepsilon^2 (r+2) (d + 8\ell B(r+2)/L^2)}{r \log(e + (\varepsilon/\delta))}\right), \quad \lambda \leq \frac{L}{\ell d^{3/2}} \left(\frac{\sqrt{r \log(e + (\varepsilon/\delta))}}{n\varepsilon}\right)^{1/2},$$

we can then obtain that

$$\begin{aligned} \mathbb{E}\|\nabla F_S(x_\tau)\|^2 &\leq \left(64\ell[F_S(x_0) - F_S^*] + 4C^2 + 2L^2\right) \frac{\sqrt{r \log(e + (\varepsilon/\delta))}}{n\varepsilon} \\ &= \left(64\ell[F_S(x_0) - F_S^*] + 64L^2 \log\left(\frac{2\sqrt{2\pi} n^3 \varepsilon^2 (r+2)(d+8\ell B(r+2)/L^2)}{r \log(e + (\varepsilon/\delta))}\right) + 2L^2\right) \frac{\sqrt{r \log(e + (\varepsilon/\delta))}}{n\varepsilon}. \end{aligned}$$

We conclude that the clipping threshold  $C$  and smoothing parameter  $\lambda$  should satisfy that

$$\begin{aligned} C &= 4L \sqrt{\log\left(\frac{2\sqrt{2\pi} n^3 \varepsilon^2 (r+2)(d+8\ell B(r+2)/L^2)}{r \log(e + (\varepsilon/\delta))}\right)}, \\ \lambda &\leq \frac{L}{\ell d} \min \left\{ 4(2 - \sqrt{2}) \sqrt{\log\left(\frac{2\sqrt{2\pi} n^3 \varepsilon^2 (r+2)(d+8\ell B(r+2)/L^2)}{r \log(e + (\varepsilon/\delta))}\right)}, \frac{1}{\sqrt{d}} \left(\frac{\sqrt{r \log(e + (\varepsilon/\delta))}}{n\varepsilon}\right)^{1/2} \right\}. \end{aligned}$$

The total number of zeroth-order gradient computations is  $nT = n^2(r+2)\varepsilon/(4\sqrt{r \log(e + (\varepsilon/\delta))})$ .

## F. Extension to the PL Setting

**Assumption F.1.** The average loss  $F_S(x)$  satisfies the PL inequality with parameter  $\mu > 0$ . That is, it holds that  $\forall x \in \mathbb{R}^d$ ,

$$\|\nabla F_S(x)\|^2 \geq 2\mu(F_S(x) - F_S^*).$$

**Corollary F.2.** Under the same setting of Theorem 1, when Assumption F.1 is also met, let  $\kappa = \ell/\mu$  be the condition number, the last iterate of Algorithm 1 satisfies that

$$\mathbb{E}[F_S(x_T) - F_S^*] \leq \left(\ell(F_S(x_0) - F_S^*) + 64L^2\kappa \log\left(\frac{n^2\varepsilon^2}{\kappa d^3 \log(e + (\varepsilon/\delta))}\right) + 2L^2\right) \frac{d^3 \log(e + (\varepsilon/\delta))}{\mu n^2 \varepsilon^2},$$

with the choice of parameters

$$\alpha = \frac{1}{4\ell d}, \quad T = 8\kappa d \log\left(\frac{n^2\varepsilon^2}{\kappa d^3 \log(e + (\varepsilon/\delta))}\right), \quad \lambda \leq \frac{2L}{\ell} \frac{\sqrt{\log(e + (\varepsilon/\delta))}}{n\varepsilon}, \quad C = Ld.$$

The total number of zeroth-order gradient computations is  $nT = \tilde{O}(nd\kappa)$ .

*Proof.* Starting from (12) in the proof of Theorem 1, with the choice that  $\alpha = 1/(4\ell d)$ , we have that

$$\begin{aligned} \mathbb{E}_{z_t, u_t}[F_S(x_{t+1})] &\leq F_S(x_t) - \frac{\alpha}{4} \|F_S(x_t)\|^2 + \frac{\alpha}{4} \ell^2 \lambda^2 d^3 + \frac{\ell}{2} \alpha^2 d \sigma^2 \\ &\leq F_S(x_t) - \frac{\mu\alpha}{2} (F_S(x_t) - F_S^*) + \frac{\alpha}{4} \ell^2 \lambda^2 d^3 + \frac{\ell}{2} \alpha^2 d \sigma^2. \end{aligned}$$

This gives the recursion that

$$\mathbb{E}[F_S(x_{t+1}) - F_S^*] \leq \left(1 - \frac{\mu\alpha}{2}\right) \mathbb{E}[F_S(x_t) - F_S^*] + \frac{\alpha}{4} \ell^2 \lambda^2 d^3 + \frac{\ell}{2} \alpha^2 d \sigma^2.$$

Resolving the recursion, we obtain that

$$\begin{aligned}
 \mathbb{E}[F_S(x_T) - F_S^*] &\leq \left(1 - \frac{\mu\alpha}{2}\right)^T (F_S(x_0) - F_S^*) + \left(\frac{\alpha}{4} \ell^2 \lambda^2 d^3 + \frac{\ell}{2} \alpha^2 d \sigma^2\right) \left(\left(1 - \frac{\mu\alpha}{2}\right)^{T-1} + \dots + \left(1 - \frac{\mu\alpha}{2}\right) + 1\right) \\
 &\leq \exp\left(-\frac{\mu\alpha T}{2}\right) (F_S(x_0) - F_S^*) + \frac{\ell^2 \lambda^2 d^3}{2\mu} + \frac{\ell \alpha d \sigma^2}{\mu} \\
 &= \exp\left(-\frac{\mu\alpha T}{2}\right) (F_S(x_0) - F_S^*) + \frac{32\ell L^2 \alpha T d^3 \log(e + (\varepsilon/\delta))}{\mu n^2 \varepsilon^2} + \frac{\ell^2 \lambda^2 d^3}{2\mu} \\
 &= \left(\ell(F_S(x_0) - F_S^*) + 64L^2 \kappa \log\left(\frac{n^2 \varepsilon^2}{\kappa d^3 \log(e + (\varepsilon/\delta))}\right) + 2L^2\right) \frac{d^3 \log(e + (\varepsilon/\delta))}{\mu n^2 \varepsilon^2},
 \end{aligned}$$

with the choice of parameters

$$\alpha T = \frac{2}{\mu} \log\left(\frac{n^2 \varepsilon^2}{\kappa d^3 \log(e + (\varepsilon/\delta))}\right), \quad \lambda \leq \frac{2L}{\ell} \frac{\sqrt{\log(e + (\varepsilon/\delta))}}{n\varepsilon}.$$

The total number of iteration is  $T = \tilde{O}(\kappa d)$ .  $\square$

**Corollary F.3.** *Under the same setting of Theorem 2, when Assumption F.1 is also met, let  $\kappa = \ell/\mu$  be the condition number, the last iterate of Algorithm 1 satisfies that*

$$\mathbb{E}[F_S(x_T) - F_S^*] \leq \left(\ell(F_S(x_0) - F_S^*) + 64L^2 \kappa \log\left(\frac{n^2 \varepsilon^2}{\kappa r d^2 \log(e + (\varepsilon/\delta))}\right) + 2L^2\right) \frac{r d^2 \log(e + (\varepsilon/\delta))}{\mu n^2 \varepsilon^2},$$

with the choice of parameters

$$\alpha = \frac{1}{4\ell(r+2)}, \quad T = 8\kappa(r+2) \log\left(\frac{n^2 \varepsilon^2}{\kappa r d^2 \log(e + (\varepsilon/\delta))}\right), \quad \lambda \leq \frac{2L}{\ell\sqrt{d}} \frac{\sqrt{r \log(e + (\varepsilon/\delta))}}{n\varepsilon}, \quad C = Ld.$$

The total number of zeroth-order gradient computations is  $nT = \tilde{O}(nr\kappa)$ .

*Proof.* Starting from (15) in the proof of Theorem 2, with the choice that  $\alpha = 1/(4\ell(r+2))$ , we have that

$$\begin{aligned}
 \mathbb{E}_{z_t, u_t}[F_S(x_{t+1})] &\leq F_S(x_t) - \frac{\alpha}{4} \|F_S(x_t)\|^2 + \frac{\alpha}{4} \ell^2 \lambda^2 d^3 + \frac{\ell}{2} \alpha^2 r \sigma^2 \\
 &\leq F_S(x_t) - \frac{\mu\alpha}{2} (F_S(x_t) - F_S^*) + \frac{\alpha}{4} \ell^2 \lambda^2 d^3 + \frac{\ell}{2} \alpha^2 r \sigma^2.
 \end{aligned}$$

This gives the recursion that

$$\mathbb{E}[F_S(x_{t+1}) - F_S^*] \leq \left(1 - \frac{\mu\alpha}{2}\right) \mathbb{E}[F_S(x_t) - F_S^*] + \frac{\alpha}{4} \ell^2 \lambda^2 d^3 + \frac{\ell}{2} \alpha^2 r \sigma^2.$$

Resolving the recursion, we obtain that

$$\begin{aligned}
 \mathbb{E}[F_S(x_T) - F_S^*] &\leq \left(1 - \frac{\mu\alpha}{2}\right)^T (F_S(x_0) - F_S^*) + \left(\frac{\alpha}{4} \ell^2 \lambda^2 d^3 + \frac{\ell}{2} \alpha^2 r \sigma^2\right) \left(\left(1 - \frac{\mu\alpha}{2}\right)^{T-1} + \dots + \left(1 - \frac{\mu\alpha}{2}\right) + 1\right) \\
 &\leq \exp\left(-\frac{\mu\alpha T}{2}\right) (F_S(x_0) - F_S^*) + \frac{\ell^2 \lambda^2 d^3}{2\mu} + \frac{\ell \alpha r \sigma^2}{\mu} \\
 &= \exp\left(-\frac{\mu\alpha T}{2}\right) (F_S(x_0) - F_S^*) + \frac{32\ell L^2 \alpha T r d^2 \log(e + (\varepsilon/\delta))}{\mu n^2 \varepsilon^2} + \frac{\ell^2 \lambda^2 d^3}{2\mu} \\
 &= \left(\ell(F_S(x_0) - F_S^*) + 64L^2 \kappa \log\left(\frac{n^2 \varepsilon^2}{\kappa r d^2 \log(e + (\varepsilon/\delta))}\right) + 2L^2\right) \frac{r d^2 \log(e + (\varepsilon/\delta))}{\mu n^2 \varepsilon^2},
 \end{aligned}$$

with the choice of parameters

$$\alpha T = \frac{2}{\mu} \log\left(\frac{n^2 \varepsilon^2}{\kappa r d^2 \log(e + (\varepsilon/\delta))}\right), \quad \lambda \leq \frac{2L}{\ell\sqrt{d}} \frac{\sqrt{r \log(e + (\varepsilon/\delta))}}{n\varepsilon}.$$

The total number of iteration is  $T = \tilde{O}(\kappa r)$ .  $\square$



**Corollary F.4.** Under the same setting of Theorem 3, when Assumption F.1 is also met, let  $\kappa = \ell/\mu$  be the condition number, suppose  $\max_{0 \leq t \leq T} |F_S(x_t)| \leq B$  and  $|F_S^*| \leq B$ , the last iterate of Algorithm 2 satisfies that

$$\mathbb{E}[F_S(x_T) - F_S^*] \leq \left( \ell(F_S(x_0) - F_S^*) + \log \left( \frac{n^2 \varepsilon^2}{\kappa r \log(e + (\varepsilon/\delta))} \right) \right) \left( L^2 + 16\tilde{L}^2 \kappa \right) + 2L^2 \frac{r \log(e + (\varepsilon/\delta))}{\mu n^2 \varepsilon^2},$$

where we define

$$\tilde{L}^2 = 64L^2 \log \left( \frac{32\sqrt{2\pi} \kappa n^3 \varepsilon^2 (r+2)(d + (8\ell(r+2) + \mu)B/L^2)}{r \log(e + (\varepsilon/\delta))} \right),$$

and choose the parameters to be

$$\alpha = \frac{1}{4\ell(r+2)}, \quad T = 8\kappa(r+2) \log \left( \frac{n^2 \varepsilon^2}{\kappa r \log(e + (\varepsilon/\delta))} \right), \quad C = \frac{\tilde{L}}{2},$$

$$\lambda \leq \frac{1}{2\ell d} \min \left\{ (2 - \sqrt{2})\tilde{L}, \frac{4L}{\sqrt{d}} \frac{\sqrt{r \log(e + (\varepsilon/\delta))}}{n\varepsilon} \right\}.$$

The total number of zeroth-order gradient computations is  $nT = \tilde{O}(nr\kappa)$ .

*Remark F.5.* A more precise expression of our theoretical results, including Theorems 1, 2, and 3 and their corresponding Corollaries F.2, F.3, and F.4, is to cover cases where  $T$  may be less than 1. Considering Theorem 3 as an example, a more accurate statement is

$$T = \max \left\{ \frac{n(r+2)\varepsilon}{4\sqrt{r \log(e + (\varepsilon/\delta))}}, 1 \right\}, \quad \mathbb{E}[\|\nabla F_S(x_\tau)\|^2] \leq \min \left\{ \tilde{O} \left( \frac{\sqrt{r \log(e + (\varepsilon/\delta))}}{n\varepsilon} \right), L^2 \right\}.$$

For the sake of clarity and simplicity in presentation, this detail is omitted in the main results.

*Proof.* Starting from (23) in the proof of Theorem 3 with the choice  $\alpha = 1/(4\ell(r+2))$  and using Assumption F.1 such that

$$\begin{aligned} \mathbb{E}\|\nabla F_S(x_t)\|^2 &\geq 2\mu \mathbb{E}[F_S(x_t) - F_S^*] \\ &= 2\mu \mathbb{E}[F_S(x_t) - F_S^*|Q] \mathbb{P}(Q) + 2\mu \mathbb{E}[F_S(x_t) - F_S^*|\bar{Q}] \mathbb{P}(\bar{Q}) \\ &\geq 2\mu \mathbb{E}[F_S(x_t) - F_S^*|Q] \mathbb{P}(Q), \end{aligned}$$

we have the recursion that

$$\mathbb{E}[F_S(x_{t+1}) - F_S^*|Q] \mathbb{P}(Q) \leq \left(1 - \frac{\mu\alpha}{2}\right) \mathbb{E}[F_S(x_t) - F_S^*|Q] \mathbb{P}(Q) + \frac{\alpha}{4} \ell^2 \lambda^2 d^3 + \frac{\ell}{2} \alpha^2 r \sigma^2 + (L^2 d \alpha + 2B) \mathbb{P}(\bar{Q}).$$

Resolving the recursion, we obtain that

$$\begin{aligned} \mathbb{E}[F_S(x_T) - F_S^*|Q] \mathbb{P}(Q) &\leq \left(1 - \frac{\mu\alpha}{2}\right)^T (F_S(x_0) - F_S^*) + \frac{2}{\mu\alpha} \left( \frac{\alpha}{4} \ell^2 \lambda^2 d^3 + \frac{\ell}{2} \alpha^2 r \sigma^2 + (L^2 d \alpha + 2B) \mathbb{P}(\bar{Q}) \right) \\ &\leq \exp \left( -\frac{\mu\alpha T}{2} \right) (F_S(x_0) - F_S^*) + \frac{\ell^2 \lambda^2 d^3}{2\mu} + \frac{\ell \alpha r \sigma^2}{\mu} + \frac{(2L^2 d + 4B/\alpha) \mathbb{P}(\bar{Q})}{\mu}. \end{aligned}$$

Since the event  $Q$  happens with high probability, the above results can be refined to

$$\begin{aligned} \mathbb{E}[F_S(x_T) - F_S^*] &= \mathbb{E}[F_S(x_T) - F_S^*|Q] \mathbb{P}(Q) + \mathbb{E}[F_S(x_T) - F_S^*|\bar{Q}] \mathbb{P}(\bar{Q}) \\ &\leq \mathbb{E}[F_S(x_T) - F_S^*|Q] \mathbb{P}(Q) + 2B \mathbb{P}(\bar{Q}). \end{aligned}$$

Therefore, we can obtain that

$$\begin{aligned}
 \mathbb{E}[F_S(x_T) - F_S^*] &\leq \exp\left(-\frac{\mu\alpha T}{2}\right) (F_S(x_0) - F_S^*) + \frac{32\ell C^2 \alpha T r \log(e + (\varepsilon/\delta))}{\mu n^2 \varepsilon^2} \\
 &\quad + \frac{4\sqrt{2\pi} n T (L^2 d + 2B/\alpha + B\mu)}{\mu} \exp\left(-\frac{C_0^2}{8L^2}\right) + \frac{\ell^2 \lambda^2 d^3}{2\mu} \\
 &= \left( \ell(F_S(x_0) - F_S^*) + L^2 \log\left(\frac{n^2 \varepsilon^2}{\kappa r \log(e + (\varepsilon/\delta))}\right) \right) \frac{r \log(e + (\varepsilon/\delta))}{\mu n^2 \varepsilon^2} \\
 &\quad + \frac{32\ell C^2 \alpha T r \log(e + (\varepsilon/\delta))}{\mu n^2 \varepsilon^2} + \frac{\ell^2 \lambda^2 d^3}{2\mu},
 \end{aligned}$$

with the choice of parameters

$$\alpha T = \frac{2}{\mu} \log\left(\frac{n^2 \varepsilon^2}{\kappa r \log(e + (\varepsilon/\delta))}\right), \quad C_0^2 = 8L^2 \log\left(\frac{32\sqrt{2\pi} \kappa n^3 \varepsilon^2 (r+2)(d + (8\ell(r+2) + \mu)B/L^2)}{r \log(e + (\varepsilon/\delta))}\right).$$

When selecting  $\lambda$  to be

$$\lambda \leq \min\left\{ \frac{2(\sqrt{2}-1)C_0}{\ell d}, \frac{2L}{\ell d^{3/2}} \frac{\sqrt{r \log(e + (\varepsilon/\delta))}}{n\varepsilon} \right\},$$

we can set  $C = \sqrt{2}C_0$  such that  $C \geq C_0 + \ell\lambda d/2$  is satisfied, and thus

$$\mathbb{E}[F_S(x_T) - F_S^*] \leq \left( \ell(F_S(x_0) - F_S^*) + \log\left(\frac{n^2 \varepsilon^2}{\kappa r \log(e + (\varepsilon/\delta))}\right) \right) (L^2 + 16\tilde{L}^2 \kappa) + 2L^2 \frac{r \log(e + (\varepsilon/\delta))}{\mu n^2 \varepsilon^2},$$

where we define

$$\tilde{L}^2 = 64L^2 \log\left(\frac{32\sqrt{2\pi} \kappa n^3 \varepsilon^2 (r+2)(d + (8\ell(r+2) + \mu)B/L^2)}{r \log(e + (\varepsilon/\delta))}\right).$$

The total number of iteration is  $T = \tilde{O}(\kappa r)$ . □