

# Zeroth-Order Fine-Tuning of LLMs in Random Subspaces

Ziming Yu  
Beijing Normal University  
zimingyu@mail.bnu.edu.cn

Pan Zhou  
Singapore Management University  
panzhou@smu.edu.sg

Sike Wang  
Beijing Normal University  
sikewang@mail.bnu.edu.cn

Jia Li<sup>†</sup>  
Beijing Normal University  
jiali@bnu.edu.cn

Hua Huang  
Beijing Normal University  
huahuang@bnu.edu.cn

## Abstract

*Fine-tuning Large Language Models (LLMs) has proven effective for a variety of downstream tasks. However, as LLMs grow in size, the memory demands for backpropagation become increasingly prohibitive. Zeroth-order (ZO) optimization methods offer a memory-efficient alternative by using forward passes to estimate gradients, but the variance of gradient estimates typically scales linearly with the model’s parameter dimension—a significant issue for LLMs. In this paper, we propose the random Subspace Zeroth-order (SubZero) optimization to address the challenges posed by LLMs’ high dimensionality. We introduce a low-rank perturbation tailored for LLMs that significantly reduces memory consumption while improving training performance. Additionally, we prove that our gradient estimation closely approximates the backpropagation gradient, exhibits lower variance than traditional ZO methods, and ensures convergence when combined with SGD. Experimental results show that SubZero enhances fine-tuning performance and achieves faster convergence compared to standard ZO approaches like MeZO across various language modeling tasks\*.*

## 1. Introduction

Large Language Models (LLMs), such as the GPT and LLaMA series [54, 60], have recently demonstrated impressive capabilities in natural language processing tasks and beyond [1, 51]. These models utilize deep learning, particularly the transformer architecture [55], to learn complex patterns in language data. However, LLMs can struggle with specialized tasks that require domain-specific knowledge [48]. Fine-tuning presents an effective solution by slightly adjust-

ing pre-trained LLMs with domain data, enabling them to adapt to specific tasks more effectively.

For fine-tuning, first-order (FO) optimizers, such as SGD [3] or Adam [29], are commonly used to achieve promising performance on domain datasets. However, as LLMs grow in size, FO optimizers demand increasingly memory consumption due to the gradient computations required by backpropagation (BP) [63]. To enhance memory efficiency, MeZO [39] first introduces the zeroth-order (ZO) optimizer to LLM fine-tuning without BP. It just needs forward passes and calculates gradient estimates using finite differences of training loss values. Nevertheless, the variance of ZO gradient estimates linearly depends on the perturbation dimension, which corresponds to the number of model parameters. This can become extremely large in LLMs, resulting in significant performance degradation compared to FO optimizers [19, 27, 38].

There are two main attempts to addressing the high variance of ZO gradient estimates. The first approach involves increasing batch size alongside training steps, which reduces gradient noise and variance in ZO gradient estimates [19, 27]. However, this leads to significant runtime and memory costs due to the large batch size in the later training stages. The second approach focuses on perturbing fewer parameters by employing sparse parameter perturbations, such as random and sparse pruning masks [38] and block-coordinate perturbations [61], or by reducing the number of trainable parameters through techniques like parameter-efficient fine-tuning (PEFT) [39, 61] and tensorized adapters [58]. Recent theoretical advancements have proposed using random projections to lessen the dimensionality dependence in ZO optimizers [30, 42, 45] by applying low-dimensional perturbations in random subspaces. Nonetheless, a major drawback of this approach is the need to store a huge projection matrix that scales with model parameter dimensionality, making it impractical for fine-tuning large LLMs.

\*Our code is available at <https://github.com/zimingyy/SubZero>.

<sup>†</sup>Corresponding author.

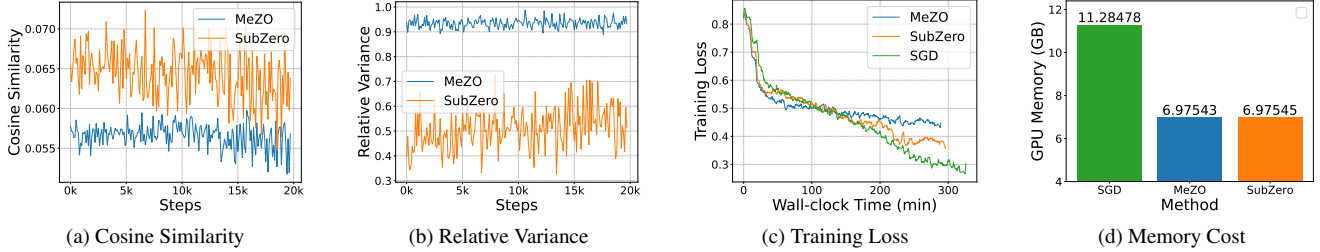


Figure 1. Visualization of cosine similarity  $\mathbb{E}[\text{cosine}(\mathbf{g}, \hat{\mathbf{g}})]$ , relative variance  $\text{Var}[\|\hat{\mathbf{g}}\|] / \|\mathbf{g}\|^2$ , training loss, and peak total GPU memory cost with OPT-1.3B on SST-2 in the prompt tuning scheme. All three methods utilize a batch size of 16 and run for 20K steps. Here,  $\hat{\mathbf{g}}$  represents the gradient estimated by MeZO or our SubZero, and  $\mathbf{g}$  denotes the expected gradient  $\mathbb{E}[\hat{\mathbf{g}}]$ . Theorem 1 (b) ensures that SubZero maintains a small distance between  $\mathbf{g}$  and the BP gradient in a subspace. (a) and (b) demonstrate that SubZero’s estimated gradient  $\hat{\mathbf{g}}$  has lower angle error and variance than MeZO. (c) and (d) indicate that SubZero enhances convergence speed with minimal extra memory usage.

**Contributions.** In this work, we propose the first random Subspace Zeroth-order (SubZero) optimization to tackle the challenges of high-dimensional LLM fine-tuning. We introduce a low-rank perturbation to estimate the gradient, specifically designed for LLM architecture, leading to reduced memory consumption and enhanced training performance. Our main contributions are as follows.

Firstly, we propose a layer-wise low-rank perturbation approach for gradient estimation, specifically designed for fine-tuning LLMs. In each layer, we generate a low-rank perturbation matrix by combining two column-orthogonal matrices with a Gaussian random matrix, which is then used for gradient estimation. Unlike traditional ZO methods like MeZO [39] which apply non-low-rank perturbations to the entire model, our approach significantly reduces the variance of gradient estimates and the angle error between the estimated gradient and its expectation, as respectively shown in Fig. 1 (a) and (b). SubZero also improves upon random subspace ZO methods like S-RGF [42] by using smaller and layer-specific low-rank perturbation matrices instead of a large and model-scale projection matrix, thus cutting memory and computational costs. Additionally, we introduce a lazy update strategy, generating perturbations periodically rather than iteratively, further reducing overhead. Besides, we also successfully apply SubZero to four popular LLM fine-tuning schemes, highlighting the compatibility of SubZero.

Secondly, we provide theoretical guarantees for SubZero. We first convert our gradient estimation into an equivalent formulation, highlighting the key differences between our approach and existing traditional ZO methods [39], as well as random subspace ZO methods [42]. Then, we prove that the gradient estimated by SubZero closely approximates the BP gradient, i.e., the ground-truth gradient, and enjoys significantly lower gradient variance than traditional ZO methods like MeZO. Furthermore, we establish the theoretical convergence of SubZero when combined with the SGD optimizer.

Finally, experimental results demonstrate SubZero’s superior performance and memory efficiency compared to other ZO approaches in both full-parameter tuning and parameter-

efficient fine-tuning (PEFT) schemes, such as LoRA, prefix tuning, and prompt tuning. For instance, SubZero improves upon MeZO by 7.1% on LLaMA-7B and by 3.2% on OPT-1.3B under full-parameter tuning and prompt tuning, while maintaining nearly identical memory costs to MeZO.

## 2. Related Work

**Zeroth-Order Fine-Tuning.** ZO optimizers utilize just two forward passes to estimate gradient without BP. Malladi et al. [39] first used ZO optimization to fine-tune LLMs, significantly lowering the GPU hours and memory usage to levels similar to inference, which offers a considerable advantage over FO optimizers. They demonstrated that LLM fine-tuning benefits from a well-structured loss landscape by introducing suitable task-specific prompt templates. Convergence theories for ZO optimization have been elaborated in both convex [18, 24, 40] and non-convex settings [25, 36]. However, these convergence rates typically increase linearly with the number of trainable parameters [18, 24, 25, 36, 40].

Recently, more work in ZO has focused on improving the convergence rates and reducing gradient estimation variance for LLM fine-tuning. Increasing batch size can diminish noise in ZO gradient estimation [19, 27]. Perturbing a subset of model parameters also lowers gradient variance. This approach induces sparse parameter perturbations through random and sparse pruning masks [38] or block-coordinate perturbations [61]. Additionally, some approaches tried to reduce trainable parameters through PEFT [39, 61] and tensorized adapters [58].

**Random Subspace Optimization.** To lessen dependence on dimensionality, some research utilizes random projections and low-dimensional perturbations in subspaces [30, 42, 45]. However, these methods are hindered by the need to store a large projection matrix that increases with dimensionality, making it impractical for fine-tuning LLMs.

**Memory-Efficient Fine-Tuning.** Fine-tuning generally employs FO optimizers like SGD [3] or Adam [29]. Various approaches have been developed to reduce the memory cost of BP, such as sparsifying gradients [53], projecting gradients into a low-rank subspace [63], and quantizing op-

tokenizer states to lower bits [14, 34]. Additional methods to conserve activation and weight memory during forward and backward passes include gradient checkpointing [8], FlashAttention [11], QLoRA [15], and LLM.int8() [13].

### 3. Preliminaries

Here we introduce the most popular ZO optimization approach and existing random subspace optimization methods. **Notations.** Let non-bold letter like  $a$  and  $A$  denote a scalar, a boldfaced lower-case letter like  $\mathbf{w}$  denote a column vector, and a boldfaced upper-case letter like  $\mathbf{W}$  denote a matrix.  $\mathcal{N}(\mathbf{0}, \mathbf{I})$  is a multivariate normal distribution with a zero mean vector and an identity covariance matrix.  $\text{vec}(\mathbf{W})$  denotes the vectorization of matrix  $\mathbf{W}$  which reshapes  $\mathbf{W}$  into a column vector by stacking the columns of  $\mathbf{W}$  vertically.  $\mathbf{A} \otimes \mathbf{B}$  is the Kronecker product of matrices  $\mathbf{A}$  and  $\mathbf{B}$ .  $\mathbb{E}[x]$  is the expected value of a random variable  $x$ .  $\text{Var}[x]$  is the variance of a random variable  $x$ . The  $\ell_2$ -norm of a vector  $\mathbf{x}$  is  $\|\mathbf{x}\| = \sqrt{\sum_{i=1}^n x_i^2}$ . The spectral norm of a matrix  $\mathbf{A}$  is  $\|\mathbf{A}\|$ . The Frobenius norm of a matrix  $\mathbf{A}$  is  $\|\mathbf{A}\|_F = \sqrt{\langle \mathbf{A}, \mathbf{A} \rangle}$ .  $C_L^{s,p}(\mathcal{S})$  denotes the class of  $s$ -th smooth and  $p$ -th  $L$ -smooth functions over the set  $\mathcal{S}$ .  $\text{bdiag}(\mathbf{A}_1, \mathbf{A}_2, \dots, \mathbf{A}_l)$  is a block diagonal matrix with diagonal blocks  $\mathbf{A}_1, \mathbf{A}_2, \dots, \mathbf{A}_l$ .

We are interested in fine-tuning large LLMs [16]. These models typically comprise multiple layers, with trainable parameter vectors represented as  $\mathbf{w} = [\mathbf{w}_1^\top, \mathbf{w}_2^\top, \dots, \mathbf{w}_l^\top]^\top \in \mathbb{R}^d$ , where  $\mathbf{w}_i$  denotes the flattened parameter vector from the  $i$ -th layer and  $d$  is model parameter dimension. Then training these models involves optimizing the problem:

$$\min_{\mathbf{w}} \mathcal{L}(\mathbf{w}), \quad (1)$$

where  $\mathcal{L}(\cdot)$  denotes the loss function.

**Zeroth-Order Optimization.** ZO optimization is BP-free and estimates gradients via random perturbations. A classical gradient estimator is the simultaneous perturbation stochastic approximation (SPSA) [52], which is defined as

$$\widehat{\nabla} \mathcal{L}(\mathbf{w}; \mathcal{B}) = \frac{\mathcal{L}(\mathbf{w} + \varepsilon \mathbf{z}; \mathcal{B}) - \mathcal{L}(\mathbf{w} - \varepsilon \mathbf{z}; \mathcal{B})}{2\varepsilon} \mathbf{z}, \quad (2)$$

where  $\mathcal{L}(\mathbf{w}; \mathcal{B})$  is the loss on a minibatch  $\mathcal{B}$  of size  $B$  uniformly sampled from the training dataset  $\mathcal{D}$ ,  $\mathbf{z} \in \mathbb{R}^d$  represents a random perturbation sampled from  $\mathcal{N}(\mathbf{0}, \mathbf{I}_d)$ , and  $\varepsilon$  is the perturbation scale.

The SPSA in Eqn. (2) is an unbiased gradient estimator of the desired gradient  $\nabla_{\mathbb{E}_{\mathbf{z}}}[\mathcal{L}(\mathbf{w} + \varepsilon \mathbf{z})]$  [40]. It only requires two forward passes to estimate the gradient and eliminates BP computation, greatly reducing computation cost and GPU memory. With this estimated gradient, one integrate with existing FO optimizers like SGD to develop corresponding ZO optimizers, e.g., ZO-SGD defined as:

$$\mathbf{w}^{t+1} = \mathbf{w}^t - \eta^t \widehat{\nabla} \mathcal{L}(\mathbf{w}^t; \mathcal{B}^t), \quad (3)$$

where  $\eta^t$  is the learning rate at step  $t$ . In practice, MeZO [39] implements ZO-SGD via in-place operations and uses a single random seed to facilitate efficient perturbation regeneration, greatly reducing memory overhead.

**Random Subspace Optimization.** Recent theoretical work [42, 45] has explored using low-dimensional perturbations in random subspaces to reduce gradient variances and hence enhance convergence rates. The key to random subspace methods is the generation of the perturbation vector  $\tilde{\mathbf{z}}$  within a subspace spanned by  $\mathbf{P}$ :

$$\tilde{\mathbf{z}} = \mathbf{P}\mathbf{z}, \quad (4)$$

where  $\mathbf{P} \in \mathbb{R}^{d \times q}$  is a random projection matrix with entries drawn from  $\mathcal{N}(0, 1)$ ,  $\mathbf{z} \in \mathbb{R}^q$  is a low-dimensional random perturbation vector sampled from  $\mathcal{N}(\mathbf{0}, \mathbf{I}_q)$ , and  $q < d$  is the dimension of the subspace. Thus, the gradient estimator in the subspace is given as follows:

$$\widehat{\nabla} \mathcal{L}(\mathbf{w}, \mathbf{P}; \mathcal{B}) = \frac{\mathcal{L}(\mathbf{w} + \varepsilon \mathbf{P}\mathbf{z}; \mathcal{B}) - \mathcal{L}(\mathbf{w} - \varepsilon \mathbf{P}\mathbf{z}; \mathcal{B})}{2\varepsilon} \mathbf{P}\mathbf{z}. \quad (5)$$

LLMs have a large model size, and thus their training and fine-tuning parameters can be very high-dimensional. This results in an excessively large matrix  $\mathbf{P}$  which is  $q$  times larger than the model size  $d$  in full-parameter tuning [2] and is also large in other fine-tuning schemes e.g., LoRA [23]. Consequently, this approach significantly increases memory requirements and computational complexity. Therefore, it is crucial to develop an efficient subspace construction strategy with minimal memory consumption for LLM fine-tuning.

## 4. Methodology

Here we first elaborate on our SubZero, a powerful ZO framework for LLM fine-tuning. Then we present how to integrate SubZero into four representative fine-tuning schemes.

### 4.1. Random Subspace Optimization for LLM Fine-Tuning

Our intuition is that exploring update directions in a low-dimensional subspace may result in a reduced variance of the estimated gradient [42, 45] compared to the estimation in the vanilla space as used in MeZO. Moreover, recent work indicates that BP gradients in LLM fine-tuning rapidly converge to a small subspace [22, 39, 62, 63]. Accordingly, we propose the random Subspace Zeroth-order (SubZero) optimization framework tailored for LLM fine-tuning. This framework reduces gradient estimation variance, and minimizes the memory overhead associated with gradient estimation, such as the memory overhead caused by the projection matrix  $\mathbf{P}$  in Eqn. (5) used in [42, 45].

**Layer-wise Random Subspace Perturbation.** LLMs primarily consist of dense layers that perform matrix multiplication. We denote the trainable parameters of the  $i$ -th layer

in matrix form as  $\mathbf{W}_i \in \mathbb{R}^{m_i \times n_i}$ . Then we will explain how to design its low-rank perturbation  $\tilde{\mathbf{Z}}_i \in \mathbb{R}^{m_i \times n_i}$ .

We propose a low-rank perturbation strategy for model parameter matrix of each layer, contrasting with previous random subspace methods that focus on the entire model’s parameters [42, 45]. At each step, we generate a low-dimensional random matrix  $\mathbf{Z}_i \in \mathbb{R}^{r \times r}$ , where  $r \ll \min\{m_i, n_i\}$ , and perform QR decomposition on two random matrices to create projection matrices  $\mathbf{U}_i \in \mathbb{R}^{m_i \times r}$  and  $\mathbf{V}_i \in \mathbb{R}^{n_i \times r}$  (see Algorithm 1). Both  $\mathbf{U}_i$  and  $\mathbf{V}_i$  are column-orthogonal matrices. Our experiments in Table 6 indicate that using Gaussian random projection matrices yields worse performance than using our designed column-orthogonal matrices. Then we combine these three matrices to yield a low-rank perturbation as follows:

$$\tilde{\mathbf{Z}}_i = \mathbf{U}_i \mathbf{Z}_i \mathbf{V}_i^\top, \quad (6)$$

where  $\tilde{\mathbf{Z}}_i$  is the perturbation matrix in a subspace spanned by  $\mathbf{U}_i$  and  $\mathbf{V}_i$ , and  $\mathbf{Z}_i$  represents the low-dimensional random perturbation matrix with entries sampled from  $\mathcal{N}(0, 1)$ .

Let the model consist of  $l$  layers, with the parameter matrix set defined as  $\mathcal{W} = \{\mathbf{W}_i\}_{i=1}^l$  and the perturbation matrix set as  $\tilde{\mathcal{Z}} = \{\tilde{\mathbf{Z}}_i\}_{i=1}^l$ . Similar to Eqns. (2) and (5), we compute the loss difference:

$$\rho = \frac{\mathcal{L}(\mathcal{W} + \varepsilon \tilde{\mathcal{Z}}; \mathcal{B}) - \mathcal{L}(\mathcal{W} - \varepsilon \tilde{\mathcal{Z}}; \mathcal{B})}{2\varepsilon}. \quad (7)$$

Note that multiplying a set by a scalar means that the scalar is multiplied by each element in the set. The addition of two sets means that the corresponding elements are added. This is only for mathematical expression, and  $\rho$  in Eqn. (7) can be calculated by two forward passes through all the layers in practice. Then we obtain the gradient estimate for the  $i$ -th layer as

$$\hat{\nabla} \mathcal{L}(\mathbf{W}_i; \mathcal{B}) = \rho \tilde{\mathbf{Z}}_i = \rho \mathbf{U}_i \mathbf{Z}_i \mathbf{V}_i^\top. \quad (8)$$

In Sec. 5, we analyze the effectiveness of this new gradient estimation (8). Specifically, Theorem 1 proves the close distance between our gradient estimate (8) and the vanilla gradient computed by BP in FO methods, while Theorem 2 shows smaller variance and angle error of our gradient estimate in Eqn. (8) compared to the gradient estimate (2) in MeZO [39]. See more theoretical details in Sec. 5.

Then, one can use estimated gradient in (8) to replace the gradient in any FO optimizer such as SGD:

$$\mathbf{W}_i^{t+1} = \mathbf{W}_i^t - \eta^t \hat{\nabla} \mathcal{L}(\mathbf{W}_i^t; \mathcal{B}^t) = \mathbf{W}_i^t - \eta^t \rho^t \mathbf{U}_i^t \mathbf{Z}_i^t \mathbf{V}_i^{t\top}. \quad (9)$$

Here we choose SGD as the default optimizer of SubZero. Theorem 3 in Sec. 5 guarantees the convergence of SubZero with SGD as basic optimizer and gives its convergence rate. The choice of FO optimizers is orthogonal to ZO optimization. Also, some empirical work indicates that adaptive

Table 1. Comparison of memory cost between SubZero and representative optimizers in full-parameter tuning scheme with RoBERTa-large on SST-2.

Method	Total GPU Memory (GB)
SGD	6.063
MeZO [39]	2.683
S-RGF [42]	23.845
SubZero	2.690

optimizers like Adam [29] do not necessarily enhance convergence of ZO approaches during LLM fine-tuning [21, 61]. Consequently, the combination of SubZero and Adam is included in Appendix 8.1 due to the limited space. We apply the primitive ZO approach. There are other ZO optimizers that utilize stochastic momentum [27] and second-order information [64] to facilitate faster convergence. While SubZero can be adapted to these ZO optimizers, we leave a comprehensive evaluation of these approaches for future work.

We compare the memory overhead of SubZero with the existing random subspace method S-RGF [42] using identical experimental settings, including layer-wise perturbation and matching subspace dimension, with all methods utilizing the SGD optimizer. As shown in Table 1, S-RGF’s memory usage is roughly four times greater than SGD and 8.8 times that of MeZO [39], while our SubZero’s memory usage is comparable to MeZO. See more experimental comparisons on OPT-13B in Table 5 of Sec. 6.

**Lazy Low-rank Subspace Update.** According to Eqn. (9), at the  $t$ -th step, the gradient estimate of the parameter matrix in the  $i$ -th layer,  $\hat{\nabla} \mathcal{L}(\mathbf{W}_i^t; \mathcal{B}^t)$ , lies within a subspace defined by the projection matrices  $\mathbf{U}_i^t$  and  $\mathbf{V}_i^t$ . Specifically,  $\mathbf{U}_i^t$  spans the column subspace, while  $\mathbf{V}_i^t$  determines the row subspace, with both matrices generated iteratively, leading to extra computational overhead to LLM fine-tuning.

However, for LLM fine-tuning, enhancing the computational efficiency and the accuracy of gradient subspace approximation is crucial. An excessively short update interval for  $\mathbf{U}_i$  and  $\mathbf{V}_i$ , such as generating them iteratively, can incur high computational costs and limit exploration of the gradient subspace they established. Conversely, a long interval may result in inaccuracies in subspace approximation and fail to capture the evolving nature of the gradient subspace. Accordingly, we propose a lazy subspace update strategy that periodically regenerates the projection matrices  $\mathbf{U}_i$  and  $\mathbf{V}_i$ . Specifically, these matrices are generated at the first step of every  $F > 1$  training steps and remain unchanged for the subsequent  $F - 1$  steps (see lines 4-7 in Algorithm 3). We utilize QR decomposition on two different random matrices for generating the column-orthogonal matrices  $\mathbf{U}_i$  and  $\mathbf{V}_i$ , as summarized in Algorithm 1. This lazy subspace update strategy is both efficient and effective in all our experiments.



---

**Algorithm 1** GenerateProjMatrix( $m, n, r$ )

---

**Input:** size of parameter matrix  $m \times n$ , rank  $r$ .

- 1: Generate random matrices  $\mathbf{R}_1 \in \mathbb{R}^{m \times r}$  and  $\mathbf{R}_2 \in \mathbb{R}^{n \times r}$  whose entries are sampled from  $\mathcal{N}(0, 1)$
- 2:  $\mathbf{U}, \_ \leftarrow \text{QR\_Decomposition}(\mathbf{R}_1)$
- 3:  $\mathbf{V}, \_ \leftarrow \text{QR\_Decomposition}(\mathbf{R}_2)$
- 4: **return**  $\mathbf{U}, \mathbf{V}$

---

---

**Algorithm 2** PerturbParams( $\mathcal{W}, \mathcal{U}, \mathcal{V}, r, \varepsilon, s$ )

---

**Input:** model parameter set  $\mathcal{W}$ , projection matrix sets  $\mathcal{U}$  and  $\mathcal{V}$ , rank  $r$ , perturbation scale  $\varepsilon$ , seed  $s$ .

- 1: Reset random number generator with seed  $s$
- 2: **for**  $i = 1, 2, \dots, l$  **do**
- 3:   Generate the perturbation matrix  $\mathbf{Z}_i \in \mathbb{R}^{r \times r}$  whose entries are sampled from  $\mathcal{N}(0, 1)$
- 4:    $\mathbf{W}_i \leftarrow \mathbf{W}_i + \varepsilon \mathbf{U}_i \mathbf{Z}_i \mathbf{V}_i^\top$
- 5: **return**  $\mathcal{W}$

---

---

**Algorithm 3** SubZero

---

**Input:** parameter matrix in the  $i$ -th layer  $\mathbf{W}_i \in \mathbb{R}^{m_i \times n_i}$ ,  $i = 1, 2, \dots, l$ , loss  $\mathcal{L}$ , step budget  $T$ , perturbation scale  $\varepsilon$ , learning rate schedule  $\{\eta^t\}$ , subspace change frequency  $F$ , rank  $r$ .

- 1: **for**  $t = 0, 1, \dots, T - 1$  **do**
- 2:   Sample a minibatch  $\mathcal{B}^t \subset \mathcal{D}$  and a random seed  $s^t$
- 3:   **for**  $i = 1, 2, \dots, l$  **do**
- 4:     **if**  $t \bmod F \equiv 0$  **then**
- 5:        $\mathbf{U}_i^t, \mathbf{V}_i^t \leftarrow \text{GenerateProjMatrix}(m_i, n_i, r)$
- 6:     **else**
- 7:        $\mathbf{U}_i^t \leftarrow \mathbf{U}_i^{t-1}, \mathbf{V}_i^t \leftarrow \mathbf{V}_i^{t-1}$
- 8:     // Note that  $\mathcal{W}^t = \{\mathbf{W}_i^t\}_{i=1}^l, \mathcal{U}^t = \{\mathbf{U}_i^t\}_{i=1}^l,$   
       $\mathcal{V}^t = \{\mathbf{V}_i^t\}_{i=1}^l$
- 9:      $\mathcal{W}^t \leftarrow \text{PerturbParams}(\mathcal{W}^t, \mathcal{U}^t, \mathcal{V}^t, r, \varepsilon, s^t),$   
       $\ell_+^t \leftarrow \mathcal{L}(\mathcal{W}^t; \mathcal{B}^t)$
- 10:      $\mathcal{W}^t \leftarrow \text{PerturbParams}(\mathcal{W}^t, \mathcal{U}^t, \mathcal{V}^t, r, -2\varepsilon, s^t),$   
       $\ell_-^t \leftarrow \mathcal{L}(\mathcal{W}^t; \mathcal{B}^t)$
- 11:      $\mathcal{W}^t \leftarrow \text{PerturbParams}(\mathcal{W}^t, \mathcal{U}^t, \mathcal{V}^t, r, \varepsilon, s^t)$
- 12:      $\rho^t \leftarrow (\ell_+^t - \ell_-^t) / (2\varepsilon)$
- 13:     Reset random number generator with seed  $s^t$
- 14:     **for**  $i = 1, 2, \dots, l$  **do**
- 15:       Regenerate the perturbation matrix  $\mathbf{Z}_i^t \in \mathbb{R}^{r \times r}$   
      whose entries are sampled from  $\mathcal{N}(0, 1)$
- 16:        $\mathbf{W}_i^{t+1} \leftarrow \mathbf{W}_i^t - \eta^t \rho^t (\mathbf{U}_i^t \mathbf{Z}_i^t \mathbf{V}_i^{t\top})$
- 17: **return**  $\mathcal{W}^{t+1}$

---

SubZero maintains just three small matrices per layer: a perturbation matrix  $\mathbf{Z}_i \in \mathbb{R}^{r \times r}$ , and two column-orthogonal matrices  $\mathbf{U}_i \in \mathbb{R}^{m_i \times r}$  and  $\mathbf{V}_i \in \mathbb{R}^{n_i \times r}$ . This design enhances memory efficiency, as  $r$  is generally much smaller than the size of the corresponding parameter matrix  $\mathbf{W}_i \in \mathbb{R}^{m_i \times n_i}$  (i.e.,  $r \ll \min\{m_i, n_i\}$ ). Moreover, we employ

in-place operations and per-layer parameter updates to estimate gradients and update parameters in parallel (see Appendix 8.2). Consequently, SubZero uses significantly less GPU memory than previous methods while achieving similar or better performance. For example, fine-tuning OPT-1.3B [60] on SST-2 [50] using SGD (without momentum) in full-parameter scheme as shown in Table 3, SubZero requires only 6.8GB GPU memory, compared to 11.5GB for SGD, yielding a  $1.6\times$  improvement in memory efficiency, similar as illustrated in Fig. 1 (d).

Now we are ready to summarize the overall algorithm of SubZero in Algorithm 3. Each training step consists of three sequential phases. First, it obtains the projection matrices  $\mathbf{U}_i^t$  and  $\mathbf{V}_i^t$  using Algorithm 1 or directly adopts previous ones. Next, it computes the loss value difference  $\rho$  with Eqn. (7) by applying Algorithm 2 to perturb all parameter matrices. Finally, SubZero updates all parameter matrices layer by layer, following Eqn. (9).

## 4.2. Integration into Fine-Tuning Schemes

We describe the integration of SubZero into full-parameter tuning [2] and three prominent PEFT schemes: LoRA [23], prefix tuning [35], and prompt tuning [32]. Typically, SubZero can be easily incorporated into these fine-tuning schemes. However, it encounters a challenge with extremely non-square parameter matrices, which have far more rows than columns or vice versa. This issue is particularly prevalent in LoRA, which employs two low-rank matrices  $\mathbf{A}_i \in \mathbb{R}^{m_i \times k}$  and  $\mathbf{B}_i \in \mathbb{R}^{k \times n_i}$  to approximate a full matrix  $\mathbf{W}_i' \in \mathbb{R}^{m_i \times n_i}$ , with  $k \ll \min\{m_i, n_i\}$ , e.g.,  $k = 8$  while  $\min\{m_i, n_i\} = 2048$  used in [61]. Consequently, it is impossible to find a smaller rank  $r \ll k$  to compute the gradient estimates of  $\mathbf{A}_i$  and  $\mathbf{B}_i$  using Eqn. (6), imposing a challenge when applying SubZero to this scenario.

To overcome this limitation, we propose a reshaping strategy that transforms the original non-square matrix into an approximate square matrix. For instance, we reshape  $\mathbf{A}_i \in \mathbb{R}^{m_i \times k}$  into  $\mathbf{A}_i' \in \mathbb{R}^{m_i' \times k'}$  such that  $m_i k = m_i' k'$  and  $m_i'$  is close to  $k'$ . This reshaping allows us to apply Eqn. (6) to find a low-rank perturbation with rank  $r$  significantly smaller than  $\min\{m_i', k'\}$ , demonstrating the applicability of SubZero in the scenario. Table 8 in Sec. 6.4 shows the effectiveness of this reshaping strategy.

## 5. Theoretical Analysis

In this section, we theoretically analyze why SubZero can reduce the variance of gradient estimates and hence accelerate convergence. Before the analysis, we first define some necessary notations:

$$\mathbf{P} = \text{bdiag}(\mathbf{V}_1 \otimes \mathbf{U}_1, \dots, \mathbf{V}_l \otimes \mathbf{U}_l), \quad (10)$$

$$\mathbf{z} = [\text{vec}(\mathbf{Z}_1)^\top, \dots, \text{vec}(\mathbf{Z}_l)^\top]^\top, \quad (11)$$

$$\tilde{\mathbf{z}} = [\text{vec}(\tilde{\mathbf{Z}}_1)^\top, \dots, \text{vec}(\tilde{\mathbf{Z}}_l)^\top]^\top. \quad (12)$$

Then we first state the main theoretical results on our gradient estimation in Eqn. (8).

**Theorem 1.** *For the gradient estimation in Eqn. (8), the following two properties hold.*

*a) By using gradient estimation in (8), our estimated gradient  $\hat{g}_\varepsilon(\mathbf{x}, \mathbf{P}, \mathbf{z})$  is equivalent to*

$$\hat{g}_\varepsilon(\mathbf{x}, \mathbf{P}, \mathbf{z}) = \frac{f(\mathbf{x} + \varepsilon \mathbf{P} \mathbf{z}) - f(\mathbf{x} - \varepsilon \mathbf{P} \mathbf{z})}{2\varepsilon} \mathbf{P} \mathbf{z}, \quad (13)$$

where  $\mathbf{z} \sim \mathcal{N}(\mathbf{0}, \mathbf{I}_q)$ ,  $\varepsilon > 0$ ,  $\mathbf{P} \in \mathbb{R}^{d \times q}$  satisfies  $\mathbf{P}^\top \mathbf{P} = \mathbf{I}_q$  with  $d = \sum_{i=1}^l m_i n_i$  and  $q = lr^2$ .

*b) Let  $\mathbf{z} \sim \mathcal{N}(\mathbf{0}, \mathbf{I}_q)$ , and  $f \in C_{L_2}^{2,2}(\mathbb{R}^d)$ . Then we have*

$$\Phi(\mathbf{x}) = \|\mathbb{E}_{\mathbf{z}}[\hat{g}_\varepsilon(\mathbf{x}, \mathbf{P}, \mathbf{z})] - \mathbf{P} \mathbf{P}^\top \nabla f(\mathbf{x})\|_2 \leq \frac{\varepsilon^2}{6} L_2 (q + 4)^2.$$

See its proof in Appendix 8.6. Theorem 1 (a) provides the equivalent form (13) of our gradient estimation (8). By comparing this with the gradient estimation (5) in random subspace optimization [42, 45], we observe significant differences. First, our gradient estimation (13) accounts for the layer-wise structure of the network, requiring the projection matrix  $\mathbf{P}$  to be block-diagonal, whereas in random subspace optimization,  $\mathbf{P}$  is not. Additionally, our method introduces a layer-wise low-rank perturbation matrix, reflected by the block-diagonal structure of  $\mathbf{P}$ , with lazy updates to the column and row spaces defined by  $\mathbf{U}_i$  and  $\mathbf{V}_i$ . In contrast, random subspace optimization simply requires  $\mathbf{P}$  to be random. These distinctions highlight the key differences between our gradient estimation and existing methods in random subspace optimization.

Theorem 1 (b) guarantees that the distance  $\Phi(\mathbf{x})$  between the expected gradient estimate and the BP gradient in the subspace spanned by  $\mathbf{P}$  is small. Moreover, by setting  $\varepsilon = \frac{1}{q+4}$ , the distance  $\Phi(\mathbf{x})$  is bounded by a constant  $L_2/6$ , independent of the parameter dimension  $d$ . This implies that the error in our gradient estimation does not scale with the extremely high parameter dimensions of LLMs, providing highly accurate gradient estimation—crucial for optimizing LLMs.

Next, we utilize a strictly convex quadratic loss to further analyze our gradient estimation in Eqn. (13). This choice is motivated by the fact that, after pretraining, the LLM parameters tend to converge toward a local minimum within a local basin, which can be well-approximated by a quadratic loss [41].

**Theorem 2.** *Let  $f(\mathbf{x}) = \mathbf{x}^\top \mathbf{H} \mathbf{x}$  and  $\mathbf{z} \sim \mathcal{N}(\mathbf{0}, \mathbf{I}_q)$ , where*

*$\mathbf{H} \in \mathbb{R}^{d \times d}$  is positive definite. We have*

$$\mathbb{E}_{\mathbf{z}}[\hat{g}_\varepsilon(\mathbf{x}, \mathbf{P}, \mathbf{z})] = \mathbf{P} \mathbf{P}^\top \nabla f(\mathbf{x}), \quad (14)$$

$$\mathbb{E}_{\mathbf{z}}[\|\hat{g}_\varepsilon(\mathbf{x}, \mathbf{P}, \mathbf{z})\|^2] = (q + 2) \|\mathbf{P}^\top \nabla f(\mathbf{x})\|^2, \quad (15)$$

$$\mathbb{E}_{\mathbf{z}} \left[ \frac{\langle \nabla f(\mathbf{x}), \hat{g}_\varepsilon(\mathbf{x}, \mathbf{P}, \mathbf{z}) \rangle^2}{\|\mathbf{P}^\top \nabla f(\mathbf{x})\|^2 \|\hat{g}_\varepsilon(\mathbf{x}, \mathbf{P}, \mathbf{z})\|^2} \right] = \frac{1}{q}. \quad (16)$$

See its proof in Appendix 8.6. Theorem 2 demonstrates several advantageous properties of our gradient estimation on the quadratic function. First, Eqn. (14) establishes the equivalence between the expected gradient estimation and the BP gradient within the subspace spanned by our projection matrix  $\mathbf{P}$ . Second, Eqn. (15) shows that, in this subspace, the variance of the gradient estimation scales linearly with the subspace dimension  $q$ . In contrast, the variance of gradient estimation (2) in MeZO depends linearly on the model’s parameter dimension  $d$ , which is significantly larger than  $q$ . Finally, Eqn. (16) reveals that the expected cosine similarity between our estimated gradient and the BP gradient within the subspace depends only on the subspace dimension  $q \ll d$ , indicating that our gradient estimation provides a highly accurate parameter update direction.

Building upon the above results, we can prove the convergence of our SubZero.

**Theorem 3.** *Let  $\mathbf{x}^* = \arg \min_{\mathbf{x} \in \mathbb{R}^d} f(\mathbf{x})$ , where  $f \in C_{L_1}^{1,1}(\mathbb{R}^d)$  and  $f$  is non-convex. Suppose  $\mathcal{E}_k = (\mathbf{z}_0, \mathbf{z}_1, \dots, \mathbf{z}_k)$ , where  $\mathbf{z}_k \sim \mathcal{N}(\mathbf{0}, \mathbf{I}_q)$  and  $\eta = \frac{1}{4(q+4)L_1}$ .  $\{\mathbf{x}_k\}_{k>0}$  is the sequence generated by Algorithm 3. For the  $\mathbf{P}$  defined in (10), which is updated lazily at a fixed frequency  $F$ , we have*

$$\frac{1}{T} \sum_{k=0}^{T-1} \mathbb{E}_{\mathcal{E}_k} [\|\nabla f(\mathbf{x}_k)\|^2] \leq \varepsilon$$

for any  $T = \Omega(\frac{d}{\varepsilon})$  if  $\varepsilon \leq \mathcal{O}\left(\frac{\varepsilon^{1/2}}{q^{3/2} d^{1/2} L_1^{3/2}}\right)$ , where  $T = KF$ ,  $K$  represents the total number of subspace updates, and  $\varepsilon$  represents perturbation scale.

See its proof in Appendix 8.6. Theorem 3 guarantees the convergence of our SubZero when the projection matrix  $\mathbf{P}$  is updated at a fixed frequency  $F$ .

## 6. Experiments

In this section, we present comprehensive experiments to evaluate the effectiveness of SubZero. We conduct our experiments using medium-sized masked LLMs (RoBERTa-large [37]) and large-scale autoregressive LLMs (OPT-1.3B and 13B [60], LLaMA2-7B [54], and Mistral-7B [26]). Our exploration covers full-parameter tuning (FT) [2] and three PEFT schemes: LoRA [23], prefix tuning [35], and prompt

Table 2. Performance of fine-tuning OPT-13B on SuperGLUE with various experimental settings (with 1000 examples). AVG: average relative percentage difference with MeZO of all tasks.

Task type Task	classification							- multiple choice -		- generation -		AVG.
	SST-2	RTE	CB	BoolQ	WSC	WIC	MultiRC	COPA	ReCoRD	SQuAD	DROP	
SGD(FT)	94.9	82.3	85.7	78.4	65.3	65.8	74.2	90.0	82.4	88.0	35.5	-
Zero-shot	58.8	59.6	46.4	59.0	38.5	55.0	46.9	80.0	81.2	46.2	14.6	-
ICL	87.0	62.1	57.1	66.9	39.4	50.5	53.1	87.0	82.5	75.9	29.6	-
LP	93.4	68.6	67.9	59.3	63.5	60.2	63.5	55.0	27.1	3.7	11.1	-
MeZO(FT)	92.1	71.5	71.4	74.4	61.5	60.0	60.1	87.0	82.0	84.2	31.2	0%
ZO-AdaMU(FT)	92.1	72.9	67.9	73.0	61.5	60.7	63.0	<b>89.0</b>	<b>83.0</b>	82.4	<b>32.0</b>	0.46%
S-MeZO(FT)	<b>92.3</b>	<b>76.9</b>	<b>75.0</b>	<b>76.5</b>	61.1	58.2	<b>63.3</b>	87.0	71.2	77.9	31.9	-0.10%
HiZOO(FT)	91.3	69.3	69.4	67.3	63.5	59.4	55.5	88.0	81.4	81.9	31.3	-2.15%
SubZero(FT)	92.1	74.0	73.2	75.3	<b>65.4</b>	<b>60.8</b>	61.0	88.0	82.3	<b>84.5</b>	<b>32.0</b>	<b>1.89%</b>
MeZO(LoRA)	92.2	74.4	69.6	75.2	64.4	59.7	58.2	87.0	82.0	82.9	31.0	0%
ZO-AdaMU(LoRA)	88.0	72.0	71.6	72.6	60.1	56.4	58.9	88.0	<b>83.2</b>	76.8	<b>32.4</b>	-1.78%
S-MeZO(LoRA)	90.8	62.2	<b>75.0</b>	72.9	51.9	55.8	56.4	86.0	69.9	76.4	31.7	-5.79%
HiZOO(LoRA)	90.6	67.5	69.6	70.5	63.5	60.2	60.2	87.0	81.9	<b>83.8</b>	31.2	-1.16%
SubZero(LoRA)	<b>93.8</b>	<b>75.5</b>	71.4	<b>76.1</b>	<b>65.4</b>	<b>60.3</b>	<b>60.3</b>	<b>89.0</b>	81.9	83.7	31.3	<b>1.57%</b>

Table 3. Performance of fine-tuning LLaMA2-7B and Mistral-7B on CB, and OPT-1.3B on SST-2.

	LLaMA2-7B				Mistral-7B				OPT-1.3B			
	FT	LoRA	Prefix	Prompt	FT	LoRA	Prefix	Prompt	FT	LoRA	Prefix	Prompt
SGD	69.6	75.0	69.6	69.6	73.2	75.0	69.6	62.5	93.2	93.0	93.1	90.7
MeZO	64.3	73.2	69.6	60.7	62.5	69.6	58.3	57.1	92.3	92.8	91.6	85.9
SubZero	<b>71.4</b>	<b>75.0</b>	<b>76.8</b>	<b>66.1</b>	<b>64.3</b>	<b>73.2</b>	<b>64.3</b>	<b>62.5</b>	<b>93.4</b>	<b>92.9</b>	<b>92.2</b>	<b>89.1</b>

Table 4. Fine-tuning performance comparison between SubZero and MeZO on RoBERTa-large and OPT-13B with non-differentiable objectives.

Model Task	RoBERTa-large				OPT-13B
	SST-2	SST-5	SNLI	MNLI	SQuAD
Zero-shot	79.0	35.5	50.2	48.8	46.2
Cross entropy (Adam)	93.9	55.9	88.7	83.8	84.2
Cross entropy (MeZO)	<b>92.9</b>	53.2	83.0	77.0	84.2
Cross entropy (SubZero)	<b>92.9</b>	<b>54.0</b>	<b>84.7</b>	<b>77.1</b>	<b>84.5</b>
Accuracy/F1 (MeZO)	92.4	46.5	81.9	73.9	80.2
Accuracy/F1 (SubZero)	<b>92.7</b>	<b>47.1</b>	<b>83.0</b>	<b>74.8</b>	<b>81.1</b>

tuning [32]. For comparison, we include leading ZO methods, such as MeZO [39], ZO-AdaMU [27], S-MeZO [38], and HiZOO [64] alongside inference-only memory-efficient baselines like zero-shot, in-context learning (ICL) [7], and linear probing (LP) [31]. As the first and most popular ZO optimizer for LLM fine-tuning, MeZO is considered our primary competitor. We also use the FO optimizer SGD as a benchmark. Since appropriate prompts are critical for ZO optimization [39, 61], all experiments incorporate prompt templates. Since a larger batch size reduces the variance of the estimated gradients, all compared methods use a batch size of 16 unless otherwise specified. All experimental settings are detailed in Appendix 8.2-8.5.

Table 5. Memory usage (GB) and wall-clock time (minutes) of fine-tuning OPT-13B, with SGD’s batch size being 8 for SQuAD and 16 for other tasks.

Task Method	SST-2		WIC		SQuAD	
	Mem.	Time	Mem.	Time	Mem.	Time
Zero-shot/ICL	24.2	0	24.8	0	27.2	0
SGD(FT)	48.9	190.3	48.9	257.3	122.7	623.7
MeZO(FT)	26.1	324.9	26.6	370.5	37.4	670.2
SubZero(FT)	26.5	337.3	27.1	385.3	37.8	690.5
MeZO(LoRA)	26.1	123.9	26.6	171.6	37.4	476.7
SubZero(LoRA)	26.1	130.3	26.6	179.7	37.4	486.5

## 6.1. Results under Different Experimental Settings

Following the settings in MeZO [39], we evaluated SubZero using OPT-13B on the SuperGLUE benchmark [56], which covers a diverse range of tasks, including classification, multiple-choice, and generation, as outlined in Table 2. The ZO methods were applied to both full-parameter tuning (FT) and LoRA fine-tuning schemes. The comparisons with vanilla LoRA and SGD with gradient accumulation are provided in Appendix 8.1.

Table 2 presents the key findings, highlighting the best-performing ZO method in bold. The results show that ZO techniques significantly outperform baseline approaches like zero-shot, in-context learning, and linear probing, underscoring their ability to enhance a pre-trained model’s perfor-

Table 6. Orthogonal or random projection matrix.

Dataset	Ortho.	Accuracy
RTE	✗	67.5
	✓	<b>74.0</b>
WSC	✗	59.6
	✓	<b>65.1</b>

Table 7. Subspace change frequency  $F$  and rank  $r$ .

$F \setminus r$	32	64	128
500	<b>72.6</b>	70.0	72.2
1000	73.6	71.8	<b>74.0</b>
2000	72.2	<b>73.3</b>	72.2
20000	70.4	<b>71.1</b>	68.6

Table 8. Reshaping strategy for non-square matrices on SST-2 with OPT-1.3B in PEFT schemes.

Method	LoRA	Prefix	Prompt
MeZO	92.8	91.6	85.9
SubZero(w/o)	92.1	89.4	74.2
SubZero(w/)	<b>92.9</b>	<b>92.2</b>	<b>89.1</b>

mance on downstream tasks.

From Table 2, one can also observe that MeZO, the first ZO optimizer for LLM fine-tuning, is highly competitive after carefully tuning its hyperparameters. Only ZO-AdaMU, aside from SubZero, outperforms MeZO in FT scheme. SubZero consistently surpasses MeZO across all tasks and fine-tuning schemes. For instance, SubZero boosts MeZO’s accuracy from 61.1% to 65.4% on the WSC task (+4.3%) under FT, and from 58.2% to 60.3% on MultiRC using LoRA (+2.1%). S-MeZO demonstrated competitive performance on several classification tasks in FT scheme. However, SubZero outperformed S-MeZO in 6 out of 11 tasks with FT and 9 out of 11 tasks with LoRA. Additionally, SubZero’s average relative percentage difference with MeZO across all tasks was better than S-MeZO’s, which displayed inconsistent performance due to its selective parameter masking based on pre-determined thresholds—an approach that lacked robustness in practice. Despite tuning S-MeZO’s hyperparameters, its performance on ReCoRD remains unsatisfactory. Excluding ReCoRD, SubZero still outperforms S-MeZO with 2.05% vs. 1.20% in FT scheme and 1.74% vs. -4.90% in LoRA scheme.

We further extended our evaluation of SubZero using OPT-1.3B, LLaMA2-7B, and Mistral-7B in FT and three PEFT schemes: LoRA, prefix tuning, and prompt tuning. As shown in Table 3, SubZero outperformed MeZO across all models and fine-tuning schemes. Notably, while MeZO struggled in the prompt tuning scheme, SubZero excelled, achieving performance levels that closely matched those of the SGD optimizer.

## 6.2. Results on Non-Differentiable Objectives

Following MeZO [39], we respectively apply SubZero to fine-tune RoBERTa-large and OPT-13B using two non-differentiable objectives: accuracy and F1. As a baseline, we also report results using the cross-entropy objective with Adam. As shown in Table 4, SubZero consistently outperforms MeZO across both non-differentiable objectives and the cross-entropy benchmark, demonstrating its effectiveness across varying optimization goals.

## 6.3. Memory Usage and Wall-Clock Time Analysis

Table 5 compares the memory consumption and wall-clock time of ZO methods (MeZO and SubZero), SGD, and inference-only approaches (zero-shot and in-context learn-

ing (ICL)) using OPT-13B. Since inference-only methods do not involve fine-tuning, they have zero wall-clock time and their memory usage reflects only the inference load. For fine-tuning, all methods were run for 20K steps. The ZO methods, including SubZero, achieved over a 1.8× reduction in memory usage compared to SGD. Notably, SubZero’s memory footprint closely aligns with MeZO’s, while offering improved performance. We use per-layer weight updates for MeZO and SubZero (see Appendix 8.2), resulting in nearly identical memory usage for FT and LoRA schemes when one decimal place is reserved.

Although SubZero introduces additional computational overhead for generating projection matrices via QR decomposition, this extra time represents less than 5% of the total wall-clock time. It is important to note that due to differences in how steps are defined between ZO methods and SGD, direct wall-clock time comparisons between the two are not entirely meaningful.

## 6.4. Ablation Study

We conducted a thorough investigation of the effectiveness of our proposed techniques. Table 6 shows that using a column-orthogonal projection matrix significantly outperforms a Gaussian random projection matrix, primarily due to the low-rank structure of the perturbation matrices. This low-rank perturbation is key to improving the quality of gradient estimation.

Next, Table 7 explores the effects of subspace rank  $r$  and update frequency  $F$  in Algorithm 3. The results demonstrate that SubZero is robust to variations in the subspace rank. However, performance drops sharply when the update frequency is too low, as the optimization becomes constrained to a single subspace for too long, limiting its adaptability.

Finally, Table 8 underscores the critical role of the reshaping strategy for handling highly non-square perturbation matrices, essential for ensuring effective perturbations in different layers of the model. Together, these results highlight the improvements brought by our design choices, particularly in terms of projection and reshaping strategies, and their impact on SubZero’s robustness and performance.

Due to limited space, the ablation studies of SubZero on random seed, batch size, and combination with Adam are given in Appendix 8.1.



## 7. Conclusion

We have demonstrated that SubZero effectively fine-tunes large LLMs across various tasks and schemes with a memory cost comparable to that of inference. Extra experiments indicate that SubZero can optimize non-differentiable objectives. Our theory explains how SubZero reduces the variance of gradient estimates and hence accelerates convergence.

**Limitation.** In addition to the representative first-order and primitive zero-order optimizers, we have yet to investigate the combinations of SubZero with other first-order and zero-order optimizers to evaluate the implications on convergence speed. While SubZero is also compatible with memory-efficient techniques like parameter quantization [34], we have not thoroughly explored the practical effects of these combinations. A theoretical analysis of the reshaping strategy is certainly worth exploring. We will leave these explorations for future work.

## References

- [1] Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, Florencia Leoni Aleman, Diogo Almeida, Janko Altschmidt, Sam Altman, Shyamal Anadkat, et al. GPT-4 technical report. *arXiv:2303.08774*, 2023. 1
- [2] Armen Aghajanyan, Sonal Gupta, and Luke Zettlemoyer. 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*, pages 7319–7328, 2021. 3, 5, 6
- [3] Shun-ichi Amari. Backpropagation and stochastic gradient descent method. *Neurocomputing*, 5(4-5):185–196, 1993. 1, 2
- [4] Roy Bar Haim, Ido Dagan, Bill Dolan, Lisa Ferro, Danilo Giampiccolo, Bernardo Magnini, and Idan Szpektor. The second PASCAL recognising textual entailment challenge. In *Proceedings of the Second PASCAL Challenges Workshop on Recognising Textual Entailment*, 2006. 14
- [5] Luisa Bentivogli, Peter Clark, Ido Dagan, and Danilo Giampiccolo. The fifth PASCAL recognizing textual entailment challenge. In *Proceedings of the Second Text Analysis Conference*, 2009. 14
- [6] Samuel R. Bowman, Gabor Angeli, Christopher Potts, and Christopher D. Manning. A large annotated corpus for learning natural language inference. In *Proceedings of the Conference on Empirical Methods in Natural Language Processing*, pages 632–642, 2015. 14
- [7] Tom B Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, et al. Language models are few-shot learners. *Advances in Neural Information Processing Systems*, pages 1877–1901, 2020. 7
- [8] Tianqi Chen, Bing Xu, Chiyuan Zhang, and Carlos Guestrin. Training deep nets with sublinear memory cost. *arXiv:1604.06174*, 2016. 3
- [9] Christopher Clark, Kenton Lee, Ming-Wei Chang, Tom Kwiatkowski, Michael Collins, and Kristina Toutanova. 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: Human Language Technologies*, 2019. 14
- [10] Ido Dagan, Oren Glickman, and Bernardo Magnini. The PASCAL recognising textual entailment challenge. In *Proceedings of the International Conference on Machine Learning Challenges: Evaluating Predictive Uncertainty Visual Object Classification, and Recognizing Textual Entailment*, 2005. 14
- [11] Tri Dao, Dan Fu, Stefano Ermon, Atri Rudra, and Christopher Ré. FlashAttention: Fast and memory-efficient exact attention with IO-awareness. *Advances in Neural Information Processing Systems*, 35:16344–16359, 2022. 3
- [12] Marie-Catherine de Marneffe, Mandy Simons, and Judith Tonhauser. The commitmentbank: Investigating projection in naturally occurring discourse. In *Proceedings of Sinn und Bedeutung 23*, 2019. 14
- [13] Tim Dettmers, Mike Lewis, Younes Belkada, and Luke Zettlemoyer. GPT3.int8 (): 8-bit matrix multiplication for transformers at scale. *Advances in Neural Information Processing Systems*, 35:30318–30332, 2022. 3
- [14] Tim Dettmers, Mike Lewis, Sam Shleifer, and Luke Zettlemoyer. 8-bit optimizers via block-wise quantization. In *Proceedings of the International Conference on Learning Representations*, 2022. 3
- [15] Tim Dettmers, Artidoro Pagnoni, Ari Holtzman, and Luke Zettlemoyer. QLoRA: Efficient finetuning of quantized LLMs. *Advances in Neural Information Processing Systems*, 36, 2024. 3
- [16] Ning Ding, Yujia Qin, Guang Yang, Fuchao Wei, Zonghan Yang, Yusheng Su, Shengding Hu, Yulin Chen, Chi-Min Chan, Weize Chen, et al. Parameter-efficient fine-tuning of large-scale pre-trained language models. *Nature Machine Intelligence*, 5(3):220–235, 2023. 3
- [17] Dheeru Dua, Yizhong Wang, Pradeep Dasigi, Gabriel Stanovsky, Sameer Singh, and Matt Gardner. 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: Human Language Technologies*, pages 2368–2378, 2019. 14
- [18] John C Duchi, Michael I Jordan, Martin J Wainwright, and Andre Wibisono. Optimal rates for zero-order convex optimization: The power of two function evaluations. *IEEE Transactions on Information Theory*, 61(5):2788–2806, 2015. 2
- [19] Tanmay Gautam, Youngsuk Park, Hao Zhou, Parameswaran Raman, and Wooseok Ha. Variance-reduced zeroth-order methods for fine-tuning language models. In *Proceedings of the International Conference on Machine Learning*, 2024. 1, 2, 14
- [20] Danilo Giampiccolo, Bernardo Magnini, Ido Dagan, and Bill Dolan. The third PASCAL recognizing textual entailment

- challenge. In *Proceedings of the ACL-PASCAL Workshop on Textual Entailment and Paraphrasing*, 2007. 14
- [21] Wentao Guo, Jikai Long, Yimeng Zeng, Zirui Liu, Xinyu Yang, Yide Ran, Jacob R Gardner, Osbert Bastani, Christopher De Sa, Xiaodong Yu, Beidi Chen, and Zhaozhuo Xu. Zeroth-order fine-tuning of LLMs with extreme sparsity. *arXiv:2406.02913*, 2024. 4, 13
- [22] Yongchang Hao, Yanshuai Cao, and Lili Mou. Flora: Low-rank adapters are secretly gradient compressors. In *Proceedings of the International Conference on Machine Learning*, 2024. 3
- [23] Edward J Hu, Yelong Shen, Phillip Wallis, Zeyuan Allen-Zhu, Yanzhi Li, Shean Wang, Lu Wang, and Weizhu Chen. LoRA: Low-rank adaptation of large language models. In *Proceedings of the International Conference on Learning Representations*, 2022. 3, 5, 6
- [24] Kevin G Jamieson, Robert Nowak, and Ben Recht. Query complexity of derivative-free optimization. *Advances in Neural Information Processing Systems*, 25, 2012. 2
- [25] Kaiyi Ji, Zhe Wang, Yi Zhou, and Yingbin Liang. Improved zeroth-order variance reduced algorithms and analysis for nonconvex optimization. In *Proceedings of the International Conference on Machine Learning*, pages 3100–3109. PMLR, 2019. 2
- [26] Albert Q Jiang, Alexandre Sablayrolles, Arthur Mensch, Chris Bamford, Devendra Singh Chaplot, Diego de las Casas, Florian Bressand, Gianna Lengyel, Guillaume Lample, Lucile Saulnier, et al. Mistral 7B. *arXiv:2310.06825*, 2023. 6
- [27] Shuoran Jiang, Qingcai Chen, Youcheng Pan, Yang Xiang, Yukang Lin, Xiangping Wu, Chuanyi Liu, and Xiaobao Song. ZO-AdaMU optimizer: Adapting perturbation by the momentum and uncertainty in zeroth-order optimization. In *Proceedings of the AAAI Conference on Artificial Intelligence*, pages 18363–18371, 2024. 1, 2, 4, 7
- [28] Daniel Khashabi, Snigdha Chaturvedi, Michael Roth, Shyam Upadhyay, and Dan Roth. Looking beyond the surface: A challenge set for reading comprehension over multiple sentences. In *Proceedings of the Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 252–262, 2018. 14
- [29] Diederik P Kingma and Jimmy Ba. Adam: A method for stochastic optimization. In *Proceedings of the International Conference on Learning Representations*, 2015. 1, 2, 4
- [30] David Kozak, Stephen Becker, Alireza Doostan, and Luis Tenorio. A stochastic subspace approach to gradient-free optimization in high dimensions. *Computational Optimization and Applications*, 79(2):339–368, 2021. 1, 2
- [31] Ananya Kumar, Aditi Raghunathan, Robbie Jones, Tengyu Ma, and Percy Liang. Fine-tuning can distort pretrained features and underperform out-of-distribution. In *Proceedings of the International Conference on Learning Representations*, 2022. 7
- [32] Brian Lester, Rami Al-Rfou, and Noah Constant. The power of scale for parameter-efficient prompt tuning. In *Proceedings of the Conference on Empirical Methods in Natural Language Processing*, pages 3045–3059, 2021. 5, 7
- [33] Hector Levesque, Ernest Davis, and Leora Morgenstern. The winograd schema challenge. In *Proceedings of the International Conference on the Principles of Knowledge Representation and Reasoning*, 2012. 14
- [34] Bingrui Li, Jianfei Chen, and Jun Zhu. Memory efficient optimizers with 4-bit states. *Advances in Neural Information Processing Systems*, 36, 2024. 3, 9
- [35] Xiang Lisa Li and Percy Liang. 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*, pages 4582–4597, 2021. 5, 6
- [36] Sijia Liu, Bhavya Kailkhura, Pin-Yu Chen, Paishun Ting, Shiyu Chang, and Lisa Amini. Zeroth-order stochastic variance reduction for nonconvex optimization. *Advances in Neural Information Processing Systems*, 31, 2018. 2
- [37] Yinhan Liu, Myle Ott, Naman Goyal, Jingfei Du, Mandar Joshi, Danqi Chen, Omer Levy, Mike Lewis, Luke Zettlemoyer, and Veselin Stoyanov. RoBERTa: A robustly optimized BERT pretraining approach. *arXiv:1907.11692*, 2019. 6
- [38] Yong Liu, Zirui Zhu, Chaoyu Gong, Minhao Cheng, Chou-Jui Hsieh, and Yang You. Sparse MeZO: Less parameters for better performance in zeroth-order LLM fine-tuning. *arXiv:2402.15751*, 2024. 1, 2, 7, 14, 15
- [39] Sadhika Malladi, Tianyu Gao, Eshaan Nichani, Alex Damian, Jason D Lee, Danqi Chen, and Sanjeev Arora. Fine-tuning language models with just forward passes. *Advances in Neural Information Processing Systems*, 36:53038–53075, 2023. 1, 2, 3, 4, 7, 8, 14, 15, 16
- [40] Yuri Nesterov and Vladimir Spokoiny. Random gradient-free minimization of convex functions. *Foundations of Computational Mathematics*, 17:527–566, 2017. 2, 3, 18, 20, 21
- [41] Behnam Neyshabur, Hanie Sedghi, and Chiyuan Zhang. What is being transferred in transfer learning? *Advances in Neural Information Processing Systems*, 33:512–523, 2020. 6
- [42] Ryota Nozawa, Pierre-Louis Poirion, and Akiko Takeda. Zeroth-order random subspace algorithm for non-smooth convex optimization. *arXiv:2401.13944*, 2024. 1, 2, 3, 4, 6
- [43] Mohammad Taher Pilehvar and Jose Camacho-Collados. WiC: the word-in-context dataset for evaluating context-sensitive meaning representations. In *Proceedings of the Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 1267–1273, 2019. 14
- [44] Pranav Rajpurkar, Jian Zhang, Konstantin Lopyrev, and Percy Liang. SQuAD: 100,000+ questions for machine comprehension of text. In *Proceedings of the Conference on Empirical Methods in Natural Language Processing*, pages 2383–2392, 2016. 14
- [45] Lindon Roberts and Clément W Royer. Direct search based on probabilistic descent in reduced spaces. *SIAM Journal on Optimization*, 33(4):3057–3082, 2023. 1, 2, 3, 4, 6
- [46] Melissa Roemmele, Cosmin Adrian Bejan, and Andrew S. Gordon. Choice of plausible alternatives: An evaluation of commonsense causal reasoning. In *Proceedings of the AAAI Spring Symposium Series*, 2011. 14

- [47] Keisuke Sakaguchi, Ronan Le Bras, Chandra Bhagavatula, and Yejin Choi. Winogrande: An adversarial winograd schema challenge at scale. *Communications of the ACM*, 64(9):99–106, 2021. [13](#)
- [48] Junhong Shen, Neil Tenenholtz, James Brian Hall, David Alvarez-Melis, and Nicolo Fusi. Tag-LLM: Repurposing general-purpose LLMs for specialized domains. In *Proceedings of the International Conference on Machine Learning*, pages 44759–44773, 2024. [1](#)
- [49] Richard Socher, Alex Perelygin, Jean Wu, Jason Chuang, Christopher D. Manning, Andrew Ng, and Christopher Potts. Recursive deep models for semantic compositionality over a sentiment treebank. In *Proceedings of the Conference on Empirical Methods in Natural Language Processing*, 2013. [14](#)
- [50] Richard Socher, Alex Perelygin, Jean Wu, Jason Chuang, Christopher D Manning, Andrew Y Ng, and Christopher Potts. Recursive deep models for semantic compositionality over a sentiment treebank. In *Proceedings of the Conference on Empirical Methods in Natural Language Processing*, pages 1631–1642, 2013. [5](#)
- [51] Irene Solaiman, Miles Brundage, Jack Clark, Amanda Askell, Ariel Herbert-Voss, Jeff Wu, Alec Radford, Gretchen Krueger, Jong Wook Kim, Sarah Kreps, et al. Release strategies and the social impacts of language models. *arXiv:1908.09203*, 2019. [1](#)
- [52] James C Spall. Multivariate stochastic approximation using a simultaneous perturbation gradient approximation. *IEEE Transactions on Automatic Control*, 37(3):332–341, 1992. [3](#)
- [53] Xu Sun, Xuancheng Ren, Shuming Ma, and Houfeng Wang. meProp: Sparsified back propagation for accelerated deep learning with reduced overfitting. In *Proceedings of the International Conference on Machine Learning*, pages 3299–3308, 2017. [2](#)
- [54] Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier Martinet, Marie-Anne Lachaux, Timothée Lacroix, Baptiste Rozière, Naman Goyal, Eric Hambro, Faisal Azhar, et al. LLaMA: Open and efficient foundation language models. *arXiv preprint arXiv:2302.13971*, 2023. [1](#), [6](#)
- [55] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N. Gomez, Lukasz Kaiser, and Illia Polosukhin. Attention is all you need. *Advances in Neural Information Processing Systems*, 2017. [1](#)
- [56] Alex Wang, Yada Pruksachatkun, Nikita Nangia, Amanpreet Singh, Julian Michael, Felix Hill, Omer Levy, and Samuel R. Bowman. SuperGLUE: A stickier benchmark for general-purpose language understanding systems. *arXiv: 1905.00537*, 2019. [7](#), [14](#)
- [57] Adina Williams, Nikita Nangia, and Samuel Bowman. 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: Human Language Technologies*, 2018. [14](#)
- [58] Yifan Yang, Kai Zhen, Ershad Banijamal, Athanasios Mouchtaris, and Zheng Zhang. AdaZeta: Adaptive zeroth-order tensor-train adaption for memory-efficient large language models fine-tuning. *arXiv:2406.18060*, 2024. [1](#), [2](#), [14](#), [15](#)
- [59] Sheng Zhang, Xiaodong Liu, Jingjing Liu, Jianfeng Gao, Kevin Duh, and Benjamin Van Durme. ReCoRD: Bridging the gap between human and machine commonsense reading comprehension. *arXiv preprint 1810.12885*, 2018. [14](#)
- [60] Susan Zhang, Stephen Roller, Naman Goyal, Mikel Artetxe, Moya Chen, Shuohui Chen, Christopher Dewan, Mona Diab, Xian Li, Xi Victoria Lin, et al. OPT: Open pre-trained transformer language models. *arXiv:2205.01068*, 2022. [1](#), [5](#), [6](#)
- [61] Yihua Zhang, Pingzhi Li, Junyuan Hong, Jiayang Li, Yimeng Zhang, Wenqing Zheng, Pin-Yu Chen, Jason D. Lee, Wotao Yin, Mingyi Hong, Zhangyang Wang, Sijia Liu, and Tianlong Chen. Revisiting zeroth-order optimization for memory-efficient LLM fine-tuning: A benchmark. In *Proceedings of the International Conference on Machine Learning*, pages 59173–59190, 2024. [1](#), [2](#), [4](#), [5](#), [7](#), [13](#), [15](#)
- [62] Zhong Zhang, Bang Liu, and Junming Shao. Fine-tuning happens in tiny subspaces: Exploring intrinsic task-specific subspaces of pre-trained language models. In *Proceedings of the Annual Meeting of the Association for Computational Linguistics*, pages 1701–1713, 2023. [3](#)
- [63] Jiawei Zhao, Zhenyu Zhang, Beidi Chen, Zhangyang Wang, Anima Anandkumar, and Yuandong Tian. Galore: Memory-efficient LLM training by gradient low-rank projection. In *Proceedings of the International Conference on Machine Learning*, 2024. [1](#), [2](#), [3](#)
- [64] Yanjun Zhao, Sizhe Dang, Haishan Ye, Guang Dai, Yi Qian, and Ivor W Tsang. Second-order fine-tuning without pain for LLMs: A Hessian informed zeroth-order optimizer. *arXiv:2402.15173*, 2024. [4](#), [7](#)

## 8. Appendix

### 8.1. Additional Results

#### More Comparisons

In the main manuscript, we use the identical batch size for FO and ZO optimizers. Here, we adjust SGD with gradient accumulation to match the memory usage of ZO optimizers, and then compare their convergence speed and performance. The experimental settings are the same as those in Figure 1, and the experimental results are shown in Figure 2. With similar memory usage, SubZero attains a convergence rate nearly on par with SGD, surpasses MeZO, and achieves test accuracy comparable to that of SGD.

The vanilla LoRA is fine-tuned by Adam. We compare SubZero with SGD in the FT and LoRA schemes with vanilla LoRA using the pretrained OPT-1.3B model on SST-2. For Adam and SubZero with SGD, we apply the constant learning rate schedule. The results are given in Table 9. We can see that SubZero with SGD in the FT scheme outperforms vanilla LoRA in both test accuracy and memory usage. SubZero with SGD in the LoRA scheme also achieves comparable test accuracy while maintaining minimal memory usage.

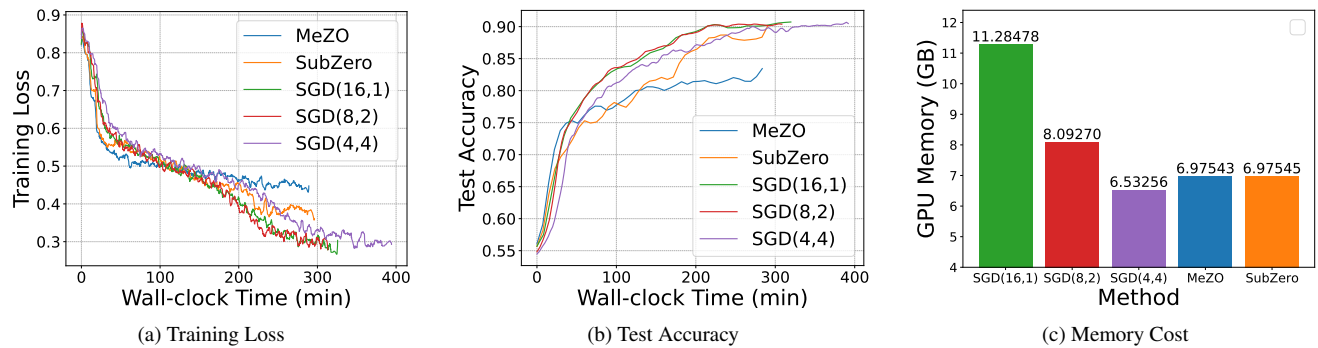


Figure 2. Visualization of training loss, test accuracy, and peak total GPU memory usage with OPT-1.3B on SST-2 in prompt tuning scheme. SGD(BS, GA) refers to SGD with a batch size of BS and GA times of gradient accumulation. All ZO methods utilize a batch size of 16, while SGD(BS, GA) applies gradient accumulation to ensure its memory usage aligns with that of the ZO optimizers. All methods are executed for 20K steps.

Table 9. Comparison with vanilla LoRA using the pretrained OPT-1.3B model on SST-2.

Method	Test Accuracy(%)	Total Memory (GB)
LoRA (Adam)	93.2	10.75
SubZero (FT)	<b>93.4</b>	6.88
SubZero (LoRA)	92.9	<b>6.80</b>

#### More Ablation Studies

We investigate the effects of random seed, batch size, and combination with Adam for SubZero. We also provide more results on the reshaping strategy.

We first fine-tune the OPT-1.3B model on the SST-2 dataset in prompt tuning scheme with three random seeds. We present the results in Table 10, and hyperparameters are presented in Table 15. For various random seeds, the variance of MeZO is quite large, whereas the variance of SubZero is small, and its average performance is superior.

Then we examine the impact of batch size for ZO optimizers using the RoBERTa-large model on SST-2 in full-parameter tuning scheme. The results are shown in Table 11. The training epochs are 100K in Table 4, while they are 20K in Table 11. The remaining hyperparameters are consistent with Table 4, as detailed in Appendix 8.4. For ZO optimizers, a large batch size always gets better performance. Across various batch sizes, SubZero demonstrates better fine-tuning performance compared to MeZO.

Next, we assess the impact of the Adam optimizer. We fine-tune the OPT-1.3B model on the SST-2 dataset, and the experimental results are displayed in Table 12. For ZO optimizers with Adam, we perform a grid search on the hyperparameters and find that keeping the learning rate and perturbation scale consistent with those with SGD resulted in good convergence, as



Table 10. The impact of random seed with the pretrained OPT-1.3B model on SST-2 in prompt tuning scheme.

Seed	42	0	1234	AVG.
MeZO	85.9	83.3	80.7	83.3
SubZero	<b>89.1</b>	<b>89.4</b>	<b>89.2</b>	<b>89.2</b>

Table 11. The impact of batch size with the pretrained RoBERTa-large model in full-parameter tuning scheme.

Batch Size	Method	SST-2	SST-5	SNLI	MNLI	AVG.
16	MeZO	91.7	44.7	77.3	<b>53.0</b>	66.7
	SubZero	<b>91.9</b>	<b>45.9</b>	<b>77.5</b>	52.8	<b>67.0</b>
32	MeZO	92.9	45.4	78.3	53.2	67.5
	SubZero	<b>93.0</b>	<b>45.5</b>	<b>79.6</b>	<b>54.0</b>	<b>68.0</b>

Table 12. Comparison of test accuracy (%) for the pretrained OPT-1.3B model fine-tuned on SST-2 with SGD and Adam.

Method	FT	LoRA	Prefix	Prompt	AVG.
SGD	<b>93.2</b>	93.0	<b>93.1</b>	90.7	92.5
Adam	92.6	<b>93.2</b>	92.9	<b>93.3</b>	<b>93.0</b>
MeZO_SGD	92.3	92.8	91.6	85.9	90.7
SubZero_SGD	<b>93.4</b>	<b>92.9</b>	<b>92.2</b>	<b>89.1</b>	<b>91.9</b>
MeZO_Adam(constant)	92.3	<b>93.3</b>	90.7	84.6	90.2
SubZero_Adam(constant)	<b>93.2</b>	92.4	<b>90.9</b>	<b>89.3</b>	<b>91.5</b>
MeZO_Adam(cosine)	<b>91.9</b>	<b>93.1</b>	86.1	78.7	87.5
SubZero_Adam(cosine)	91.7	92.0	<b>86.6</b>	<b>83.4</b>	<b>88.4</b>

Table 13. Reshaping strategy for non-square matrices with the pretrained OPT-1.3B model fine-tuned on Winogrande in the PEFT schemes.

Method	LoRA	Prefix	Prompt	AVG.
SGD	58.3	56.9	58.4	57.9
SubZero(w/o)	56.6	56.6	56.5	56.6
SubZero(w/)	<b>57.8</b>	<b>57.3</b>	<b>57.6</b>	<b>57.6</b>

detailed in Table 15. We utilize the linear and the constant learning rate schedules for SGD and Adam, respectively. For all ZO optimizers with SGD, we apply the constant learning rate schedule. For all ZO optimizers with Adam, we test the constant and the cosine annealing schedules. We note that SubZero surpasses MeZO when employing the Adam optimizer with both constant and cosine annealing schedules. Also, Adam does not provide an advantage over SGD for ZO optimization, which aligns with the conclusions of previous studies [21, 61].

Finally, we provide more ablations on the reshaping strategy with OPT-1.3B on Winogrande in the PEFT schemes. The Winogrande dataset [47] is a benchmark for commonsense reasoning and available at <https://winogrande.allenai.org/>. The results are shown in Table 13. We can see that the reshaping strategy clearly enhances performance, aligning with the conclusion presented in Table 8.

## 8.2. Implementation Details

We use one A800 GPU with the PyTorch 2.1.0+CUDA 11.8 framework for ZO methods and, if needed, two A800 GPUs for SGD.

The gradient estimation in SubZero is applicable to parameter matrices, while LLMs mainly consist of dense layers. For other trainable parameters, such as biases and layer normalization parameters, we recommend using the gradient estimation in MeZO [39], as these layers contain fewer parameters.

We introduce two useful strategies to implement our SubZero efficiently in memory.

**In-place Operation.** As indicated in Eqn. (7), directly computing the loss difference  $\rho$  requires twice the memory of inference, as it must store both the parameter matrix set  $\mathcal{W}$  and the perturbation matrix set  $\hat{\mathcal{Z}}$ . To mitigate this, we draw inspiration from MeZO and utilize in-place operations. By employing the random seed trick, we store a random seed to compute  $\rho$  (see lines 9-12 in Algorithm 3 and Algorithm 2) and regenerate the low-dimensional perturbation matrices  $\mathbf{Z}_1, \mathbf{Z}_2, \dots, \mathbf{Z}_l$  (see line 15 in Algorithm 3). Consequently, the memory cost for fine-tuning with SubZero is nearly equivalent to that of inference (see Table 1 and Table 5).

**Per-layer Weight Update.** FO optimizers update all model parameters after BP by storing the entire gradients in memory. In contrast, ZO optimizers like SubZero calculate gradient estimates by first determining the loss value difference from two forward passes, then calculating the gradient estimate for each layer using this difference along with the layer’s perturbation. To reduce memory usage during training, we can implement the parameter update with `optimizer.step()` after calculating the gradient estimate for each layer.

SubZero significantly reduces GPU memory consumption with the two implementation strategies. It should note that we use the per-layer weight update strategy for MeZO in all experiments.

To simplify hyperparameter tuning, we employ a norm alignment trick, allowing SubZero to directly utilize hyperparameter settings, such as the learning rate, from MeZO [39]. For a random perturbation matrix  $\mathbf{Z} \in \mathbb{R}^{m \times n}$ , and its low-rank approximation is  $\hat{\mathbf{Z}} = \mathbf{U}\mathbf{Z}'\mathbf{V}^\top$ , where  $\mathbf{U} \in \mathbb{R}^{m \times r}$ ,  $\mathbf{V} \in \mathbb{R}^{n \times r}$ , and  $\mathbf{Z}' \in \mathbb{R}^{r \times r}$ . If  $\mathbf{Z}$  and  $\mathbf{Z}'$  are Gaussian random matrices, and  $\mathbf{U}$  and  $\mathbf{V}$  are column-orthogonal matrices, then we have:

$$\mathbb{E}[\|\mathbf{Z}\|_F] = \sqrt{\frac{m \times n}{r^2}} \mathbb{E}[\|\hat{\mathbf{Z}}\|_F]. \quad (17)$$

Define  $\mu = \sqrt{\frac{m \times n}{r^2}}$ . Let MeZO’s learning rate be  $\eta$  and perturbation scale be  $\varepsilon$ . There are two equivalent approaches to obtain the perturbation for SubZero. The first approach involves multiplying the random low-dimensional perturbation matrix by  $\mu$ , with SubZero adopting MeZO’s hyperparameters directly:  $\eta' = \eta$  and  $\varepsilon' = \varepsilon$ . The second approach keeps the random low-dimensional perturbation matrix fixed and sets SubZero’s learning rate and perturbation scale as follows:

$$\eta' = \eta\mu^2, \varepsilon' = \varepsilon\mu.$$

We argue that norm alignment is crucial for SubZero, as changing the rank  $r$  affects the norm of the gradient estimate, complicating the fine-tuning of the associated learning rate.

S-MeZO [38], a new ZO method, aims to improve MeZO’s performance and convergence speed. However, its source code and detailed layer-wise hyperparameter configurations have not been released. Yang et al. [58] reproduce S-MeZO using a fixed sparsity ratio for each layer, selected based on the best overall result shown in Fig. 6 of their paper. So we perform S-MeZO with this non-official implementation code available at <https://github.com/yifanycc/AdaZeta>.

### 8.3. Datasets

Following [39], we use SuperGLUE [56] for OPT experiments, including BoolQ [9], CB [12], COPA [46], MultiRC [28], ReCoRD [59], RTE [4, 5, 10, 20], WiC [43], and WSC [33]. We also utilize SST-2 [49] and two question answering (QA) datasets, SQuAD [44] and DROP [17]. For each task, we randomly sampled 1000 examples for training, 500 for validation, and 1000 for testing.

For LLama2-7B and Mistral-7B, we use CB [12] in the full-parameter tuning and three PEFT schemes. For OPT-1.3B, we utilize SST-2 [49] in the full-parameter tuning and three PEFT schemes.

For RoBERTa-large, we consider classification datasets: SST-2 [49], SST-5 [49], MNLI [57], and SNLI [6]. Following [39], the test set has 1000 examples for fast iteration, while we have 512 examples per class for both training and validation.

### 8.4. Hyperparameters

Using a larger batch size can consistently reduce the variance in ZO optimization, thus enhancing fine-tuning performance [19, 39, 58]. However, this increase in batch size also raises the time for forward passes and significantly elevates memory usage. We focus on developing ZO methods that minimize variance and improve performance with small batch sizes, with a default

Table 14. The hyperparameter search grids for OPT-13B. For each task, we run 20K steps for ZO methods (MeZO, S-MeZO, and SubZero) and SGD. We record the best model checkpoint based on the validation loss every 500 training steps.

Experiment	Hyperparameter	Value
MeZO(FT)	batch size	16
	learning rate	{1e-7, 2e-7, 5e-7, 1e-6}
	$\epsilon$	1e-3
MeZO(LoRA)	batch size	16
	learning rate	{1.5e-5, 3e-5, 5e-5}
	$\epsilon$	1e-3
S-MeZO(FT)	batch size	16
	learning rate	{1e-6, 5e-6}
	$\epsilon$	1e-3
S-MeZO(LoRA)	batch size	16
	learning rate	{5e-5, 1e-4, 1e-3}
	$\epsilon$	1e-3
SubZero(FT)	batch size	16
	learning rate	{1e-7, 2e-7, 5e-7, 1e-6}
	$\epsilon$	1e-3
SubZero(LoRA)	batch size	16
	learning rate	{1.5e-5, 3e-5, 5e-5}
	$\epsilon$	1e-3
SGD(FT)	batch size	16
	learning rate	{1e-4, 1e-3, 5e-3}
	rank	{32, 64, 128, 256}
SubZero(FT)	rank	{32, 64, 128, 256}
	subspace change frequency	{500, 1000, 2000}
	rank	{4, 8, 16}
SubZero(LoRA)	rank	{4, 8, 16}
	subspace change frequency	{500, 1000, 2000}
	rank	{4, 8, 16}
SGD(FT)	batch size	16
	learning rate	{1e-4, 1e-3, 5e-3}
	rank	{4, 8, 16}

setting of 16. In some SGD experiments, like on MultiRC and SQuAD, the batch size is reduced to 8 due to limited GPU resources.

Consistent with previous studies [38, 39, 58, 61], we employ SGD without momentum by default to maintain memory efficiency. SGD utilizes the linear learning rate schedule, while all ZO methods with SGD apply a constant learning rate schedule, with weight decay set to 0.

For RoBERTa, we run Adam for 1K steps and ZO methods for 100K steps. In the rest experiments, we run Adam for 5 epochs and SGD and ZO methods for 20K steps.

We follow previous work to set the hyperparameters in the PEFT schemes [39, 61]. For LoRA, the rank is set to 8 and  $\alpha$  is set to 16. For prefix tuning, the length of prefix tokens is set to 5, and we initialize these tunable representations by randomly sampling tokens from the vocabulary and then passing them through the LLM to get their keys and values at different attention layers. For prompt tuning, the length of prompt virtual tokens is set to 10, and the prompt tokens are initialized with actual token values from the model’s embedding.

We present the hyperparameter search grids in Tables 14 and 15 to assist with result reproduction. For OPT-1.3B, we utilize the same hyperparameter settings as in Table 15. For RoBERTa-large, we use a learning rate of {1e-6, 5e-6} and  $\epsilon=1e-3$  for MeZO and SubZero, with a batch size of 64. The rank for SubZero is set to {8, 16, 24}, and subspace change frequency is adjusted to {1000, 2000}.

For the ablation study, we evaluate the effectiveness of the orthogonal projection matrix using the OPT-13B model in full-parameter tuning scheme on the RTE and WSC datasets, and the results are presented in Table 6. The hyperparameter settings are consistent with those in Table 2, and further details are available in Table 14. The subspace dimensionality remains fixed across all experiments. It is noteworthy that both orthogonal and non-orthogonal projection matrices can utilize the same

Table 15. The hyperparameter search grids for LLama2-7B and Mistral-7B. For each task, we run 20K steps for ZO methods (MeZO and SubZero) and SGD. We record the best model checkpoint based on the validation loss every 500 training steps.

Experiment	Hyperparameter	Value
MeZO(FT)	batch size	16
	learning rate	{ 1e-7, 5e-7, 1e-6}
	$\epsilon$	1e-3
MeZO(LoRA)	batch size	16
	learning rate	{ 1e-6, 5e-6, 1e-5, 3e-5}
	$\epsilon$	1e-3
MeZO(Prefix)	batch size	16
	learning rate	{ 1e-3, 5e-3, 1e-2}
	$\epsilon$	1e-1
MeZO(Prompt)	batch size	16
	learning rate	{ 1e-3, 5e-3, 1e-2}
	$\epsilon$	1e-2
SubZero(FT)	batch size	16
	learning rate	{ 1e-7, 5e-7, 1e-6}
	$\epsilon$	1e-3
	rank	{24, 48}
SubZero(LoRA)	subspace change frequency	1000
	batch size	16
	learning rate	{ 1e-6, 5e-6, 1e-5, 3e-5}
	$\epsilon$	1e-3
SubZero(Prefix)	rank	{4, 8}
	subspace change frequency	1000
	batch size	16
	learning rate	{ 1e-3, 5e-3, 1e-2}
SubZero(Prompt)	$\epsilon$	1e-1
	rank	{4, 8}
	subspace change frequency	1000
	learning rate	{ 1e-3, 5e-3, 1e-2}
SubZero(LoRA)	batch size	16
	learning rate	{ 1e-3, 5e-3, 1e-2}
	$\epsilon$	1e-2
	rank	{16, 24}
SubZero(Prompt)	subspace change frequency	1000
	batch size	16
	learning rate	{ 1e-5, 1e-4, 1e-3, 5e-3}
SGD(FT)	learning rate	{ 1e-5, 1e-4, 1e-3, 5e-3}

learning rate and perturbation scale. This is because the overall perturbation matrix is scaled by a factor of  $\frac{1}{r}$ , following a similar norm alignment strategy as detailed in Eqn. (17). We also perform ablation studies on the rank and subspace update frequency for SubZero, with results shown in Table 7. Full-parameter tuning scheme is conducted on the RTE dataset using the OPT-13B model, with specific experimental settings outlined in Table 14. All experiments employ the same learning rate and perturbation scale, enabled by the norm alignment technique described in Eqn. (17).

## 8.5. Prompt Templates

For autoregressive LLMs, we have three task types: classification, multiple-choice, and question answering. We adopt the prompt templates for various tasks in [39], which are summarized in Table 16. For masked LLMs, we also adopt the prompt templates in [39] and present them in Table 17.



Table 16. The prompt templates used in the OPT-1.3B, OPT-13B, LLama2-7B, and Mistral-7B experiments.

Task	Type	Prompt
SST-2	cls.	<text> It was <b>terrible/great</b>
RTE	cls.	<premise> Does this mean that "<hypothesis>" is true? Yes or No? <b>Yes or No</b>
CB	cls.	Does this mean that "<hypothesis>" is true? Yes or No? <b>Yes/No/Maybe</b>
BoolQ	cls.	<passage> <question>? <b>Yes/No</b>
WSC	cls.	<text> In the previous sentence, does the pronoun "<span2>" refer to <span1>? Yes or No? <b>Yes/No</b>
WIC	cls.	Does the word "<word>" have the same meaning in these two sentences? Yes, No? <sentence1> <sentence2> <b>Yes/No</b>
MultiRC	cls.	<paragraph> Question: <question> I found this answer "<answer>". Is that correct? Yes or No? <b>Yes/No</b>
COPA	mch.	<premise> so/because <candidate>
ReCoRD	mch.	<passage> <query>.replace("@placeholder", <candidate>)
SQuAD	QA	Title: <title> Context: <context> Question: <question> Answer:
DROP	QA	Passage: <context> Question: <question> Answer:

Table 17. The prompt templates used in RoBERTa-large experiments.  $C$  is the number of classification categories.

Task	$C$	Type	Prompt
SST-2	2	sentiment cls.	<sentence1> It was <b>great/terrible</b>
SST-5	5	sentiment cls.	<sentence1> It was <b>great/good/okay/bad/terrible</b>
MNLI	3	NLI	<sentence1> ? <b>Yes/Maybe/No</b> , <sentence2>
SNLI	3	NLI	<sentence1> ? <b>Yes/Maybe/No</b> , <sentence2>

## 8.6. Proofs

In practice, SubZero employs smaller and layer-specific low-rank perturbation matrices instead of a large model-scale projection matrix. However, it is more convenient to prove SubZero’s properties using a model-scale projection. Fortunately, the following lemma shows that the low-rank perturbation matrix for each layer can be represented as a layer-scale projection matrix, which is column orthogonal.

**Lemma 1.** Let  $\tilde{\mathbf{Z}} = \mathbf{UZV}^T$ , where  $\mathbf{U} \in \mathbb{R}^{m \times r}$ ,  $\mathbf{Z} \in \mathbb{R}^{r \times r}$ ,  $\mathbf{V} \in \mathbb{R}^{n \times r}$ , and  $\mathbf{U}^T \mathbf{U} = \mathbf{V}^T \mathbf{V} = \mathbf{I}_r$ . Then we have  $\text{vec}(\tilde{\mathbf{Z}}) = \mathbf{P} \text{vec}(\mathbf{Z})$  and  $\mathbf{P}^T \mathbf{P} = \mathbf{I}_{r^2}$ , where  $\mathbf{P} = \mathbf{V} \otimes \mathbf{U}$ .

*Proof.* Since  $\text{vec}(U\mathbf{Z}\mathbf{V}^\top) = (\mathbf{V} \otimes \mathbf{U})\text{vec}(\mathbf{Z})$ , we only need to show  $(\mathbf{V} \otimes \mathbf{U})^\top(\mathbf{V} \otimes \mathbf{U}) = \mathbf{I}_{r^2}$ . In fact

$$(\mathbf{V} \otimes \mathbf{U})^\top(\mathbf{V} \otimes \mathbf{U}) = (\mathbf{V}^\top \otimes \mathbf{U}^\top)(\mathbf{V} \otimes \mathbf{U}) = (\mathbf{V}^\top \mathbf{V}) \otimes (\mathbf{U}^\top \mathbf{U}) = \mathbf{I}_r \otimes \mathbf{I}_r = \mathbf{I}_{r^2}.$$

The proof is completed.  $\square$

We can also demonstrate that the low-rank perturbation matrices across all layers can be represented as a model-scale projection matrix. We first give the following lemma.

**Lemma 2.** *Let a block diagonal matrix  $\mathbf{P} = \text{bdiag}(\mathbf{P}_1, \mathbf{P}_2, \dots, \mathbf{P}_l)$  and  $\tilde{\mathbf{z}}_i = \mathbf{P}_i \mathbf{z}_i$ , where  $\mathbf{P}_i^\top \mathbf{P}_i = \mathbf{I}_{r^2}$  and  $i = 1, 2, \dots, l$ . Then we have  $\tilde{\mathbf{z}} = \mathbf{P}\mathbf{z}$ , where  $\tilde{\mathbf{z}} = [\tilde{\mathbf{z}}_1^\top, \dots, \tilde{\mathbf{z}}_l^\top]^\top$ ,  $\mathbf{z} = [\mathbf{z}_1^\top, \dots, \mathbf{z}_l^\top]^\top$  and  $\mathbf{P}^\top \mathbf{P} = \mathbf{I}_{lr^2}$ .*

*Proof.* It is easy to check that  $\tilde{\mathbf{z}} = \mathbf{P}\mathbf{z}$ . Besides, we have

$$\mathbf{P}^\top \mathbf{P} = \text{bdiag}(\mathbf{P}_1^\top, \dots, \mathbf{P}_l^\top) \text{bdiag}(\mathbf{P}_1, \dots, \mathbf{P}_l) = \text{bdiag}(\mathbf{P}_1^\top \mathbf{P}_1, \dots, \mathbf{P}_l^\top \mathbf{P}_l) = \mathbf{I}_{lr^2}.$$

The proof is completed.  $\square$

We may define  $\mathbf{P} = \text{bdiag}(\mathbf{V}_1 \otimes \mathbf{U}_1, \mathbf{V}_2 \otimes \mathbf{U}_2, \dots, \mathbf{V}_l \otimes \mathbf{U}_l)$  that satisfies  $\mathbf{P}^\top \mathbf{P} = \mathbf{I}$ ,  $\mathbf{z} = [\text{vec}(\mathbf{Z}_1)^\top, \text{vec}(\mathbf{Z}_2)^\top, \dots, \text{vec}(\mathbf{Z}_l)^\top]^\top$ , and  $\tilde{\mathbf{z}} = [\text{vec}(\tilde{\mathbf{Z}}_1)^\top, \text{vec}(\tilde{\mathbf{Z}}_2)^\top, \dots, \text{vec}(\tilde{\mathbf{Z}}_l)^\top]^\top$ . Then according to Lemma 2, the perturbation vector of SubZero is  $\tilde{\mathbf{z}} = \mathbf{P}\mathbf{z}$ , which is similar as existing random subspace methods in Eqn. (4), but with SubZero's projection matrix being block diagonal and column orthogonal.

To prove Theorem 1 and Theorem 2, we first introduce some definitions and lemmas about Gaussian distribution.

**Definition 1.** *We say  $\mathbf{z}$  is a standard  $n$ -dimensional Gaussian vector (denote by  $\mathbf{z} \sim \mathcal{N}(\mathbf{0}, \mathbf{I}_n)$ ), if its probability density function  $p(\mathbf{z}) = \frac{1}{\kappa} e^{-\frac{1}{2}\|\mathbf{z}\|^2}$ , where  $\kappa > 0$  satisfies  $\int_{\mathbb{R}^n} \frac{1}{\kappa} e^{-\frac{1}{2}\|\mathbf{z}\|^2} d\mathbf{z} = 1$ .*

**Definition 2.** *Let  $\mathbf{z} \sim \mathcal{N}(\mathbf{0}, \mathbf{I}_n)$ . We say  $x$  is a chi-square random variable with degrees of freedom  $n$  (denote by  $x \sim \chi^2(n)$ ), if  $x = \|\mathbf{z}\|^2$ .*

**Lemma 3.** *Let  $\mathbf{z} \sim \mathcal{N}(\mathbf{0}, \mathbf{I}_n)$ . For any orthogonal  $(n \times n)$ -matrix  $\mathbf{Q}$  and continuous function  $f$ , we have  $\mathbb{E}_{\mathbf{z}}[f(\mathbf{z})] = \mathbb{E}_{\mathbf{z}}[f(\mathbf{Q}\mathbf{z})]$ .*

**Lemma 4.** *If  $x \sim \chi^2(n)$ , then we have*

$$\mathbb{E}_x[x] = n, \quad \text{Var}_x[x] = 2n.$$

**Lemma 5.** [40] *Let  $f \in C_{L_2}^{2,2}(\mathbb{R}^n)$ . Then for all  $\mathbf{x}, \mathbf{y} \in \mathbb{R}^n$ , we have*

$$|f(\mathbf{y}) - f(\mathbf{x}) - \langle \nabla f(\mathbf{x}), \mathbf{y} - \mathbf{x} \rangle - \frac{1}{2} \langle \nabla^2 f(\mathbf{x})(\mathbf{y} - \mathbf{x}), \mathbf{y} - \mathbf{x} \rangle| \leq \frac{L_2}{6} \|\mathbf{y} - \mathbf{x}\|^3.$$

**Lemma 6.** [40] *Let  $\mathbf{z} \sim \mathcal{N}(\mathbf{0}, \mathbf{I}_n)$ . For  $0 \leq t \leq 2$ , we have*

$$\mathbb{E}_{\mathbf{z}}[\|\mathbf{z}\|^t] \leq n^{t/2}.$$

For  $t \geq 2$ , we have

$$n^{t/2} \leq \mathbb{E}_{\mathbf{z}}[\|\mathbf{z}\|^t] \leq (n+t)^{t/2}.$$

**Lemma 7.** *Let  $\mathbf{z} \sim \mathcal{N}(\mathbf{0}, \mathbf{I}_n)$ . For all  $\mathbf{y} \in \mathbb{R}^n$ , we have*

$$\mathbb{E}_{\mathbf{z}}[\|\langle \mathbf{y}, \mathbf{z} \rangle \mathbf{z}\|^2] = (n+2)\|\mathbf{y}\|^2.$$

*Proof.* Note that for any orthogonal  $(n \times n)$ -matrix  $\mathbf{Q}$ , we have

$$\|\langle \mathbf{y}, \mathbf{Q}\mathbf{z} \rangle \mathbf{Q}\mathbf{z}\|^2 = \|\langle \mathbf{Q}^\top \mathbf{y}, \mathbf{z} \rangle \mathbf{z}\|^2, \quad \|\mathbf{Q}^\top \mathbf{y}\| = \|\mathbf{y}\|.$$

In accordance with Lemma 3, we can set  $\mathbf{y} = [1, 0, \dots, 0]^\top$ , and only need to prove  $\mathbb{E}_{\mathbf{z}}[\|\langle \mathbf{y}, \mathbf{z} \rangle \mathbf{z}\|^2] = n + 2$ . Equipped with Lemma 4, we get

$$\mathbb{E}_{\mathbf{z}}[\|\langle \mathbf{y}, \mathbf{z} \rangle \mathbf{z}\|^2] = \mathbb{E}_{\mathbf{z}} \left[ \sum_{i=1}^n z_1^2 z_i^2 \right] = \sum_{i=1}^n \mathbb{E}_{\mathbf{z}}[z_1^2 z_i^2] = \mathbb{E}_{z_1}[z_1^4] + \mathbb{E}_{z_1}[z_1^2] \sum_{i=2}^n \mathbb{E}_{z_i}[z_i^2] = n + 2.$$

The proof is completed. □

**Theorem 1.** For the gradient estimation in Eqn. (8), the following two properties hold.

a) By using gradient estimation in (8), our estimated gradient  $\hat{g}_\varepsilon(\mathbf{x}, \mathbf{P}, \mathbf{z})$  is equivalent to

$$\hat{g}_\varepsilon(\mathbf{x}, \mathbf{P}, \mathbf{z}) = \frac{f(\mathbf{x} + \varepsilon \mathbf{P}\mathbf{z}) - f(\mathbf{x} - \varepsilon \mathbf{P}\mathbf{z})}{2\varepsilon} \mathbf{P}\mathbf{z}, \quad (13)$$

where  $\mathbf{z} \sim \mathcal{N}(\mathbf{0}, \mathbf{I}_q)$ ,  $\varepsilon > 0$ ,  $\mathbf{P} \in \mathbb{R}^{d \times q}$  satisfies  $\mathbf{P}^\top \mathbf{P} = \mathbf{I}_q$  with  $d = \sum_{i=1}^l m_i n_i$  and  $q = lr^2$ .

b) Let  $\mathbf{z} \sim \mathcal{N}(\mathbf{0}, \mathbf{I}_q)$ , and  $f \in C_{L_2}^{2,2}(\mathbb{R}^d)$ . Then we have

$$\Phi(\mathbf{x}) = \|\mathbb{E}_{\mathbf{z}}[\hat{g}_\varepsilon(\mathbf{x}, \mathbf{P}, \mathbf{z})] - \mathbf{P}\mathbf{P}^\top \nabla f(\mathbf{x})\|_2 \leq \frac{\varepsilon^2}{6} L_2 (q + 4)^2.$$

*Proof.* a) Evidently, the conclusion is established based on Lemma 1 and Lemma 2.

b) Let  $a_{\mathbf{z}}(\tau) = f(\mathbf{x} + \tau \mathbf{z}) - f(\mathbf{x}) - \tau \langle \nabla f(\mathbf{x}), \mathbf{z} \rangle - \frac{\tau^2}{2} \langle \nabla^2 f(\mathbf{x}) \mathbf{z}, \mathbf{z} \rangle$ . Lemma 5 implies that

$$|a_{\mathbf{z}}(\pm\varepsilon)| \leq \frac{\varepsilon^3}{6} L_2 \|\mathbf{z}\|^3.$$

Note that

$$\begin{aligned} & \mathbb{E}_{\mathbf{z}}[\hat{g}_\varepsilon(\mathbf{x}, \mathbf{P}, \mathbf{z})] - \mathbf{P}\mathbf{P}^\top \nabla f(\mathbf{x}) \\ &= \frac{\mathbf{P}}{2\kappa\varepsilon} \int_{\mathbb{R}^q} [f(\mathbf{x} + \varepsilon \mathbf{P}\mathbf{z}) - f(\mathbf{x} - \varepsilon \mathbf{P}\mathbf{z}) - 2\varepsilon \langle \nabla f(\mathbf{x}), \mathbf{P}\mathbf{z} \rangle] \mathbf{z} e^{-\frac{1}{2}\|\mathbf{z}\|^2} d\mathbf{z}. \end{aligned}$$

Therefore, in accordance with Lemma 6, we have

$$\begin{aligned} & \|\mathbb{E}_{\mathbf{z}}[\hat{g}_\varepsilon(\mathbf{x}, \mathbf{P}, \mathbf{z})] - \mathbf{P}\mathbf{P}^\top \nabla f(\mathbf{x})\| \\ & \leq \frac{1}{2\kappa\varepsilon} \int_{\mathbb{R}^q} |f(\mathbf{x} + \varepsilon \mathbf{P}\mathbf{z}) - f(\mathbf{x} - \varepsilon \mathbf{P}\mathbf{z}) - 2\varepsilon \langle \nabla f(\mathbf{x}), \mathbf{P}\mathbf{z} \rangle| \|\mathbf{z}\| e^{-\frac{1}{2}\|\mathbf{z}\|^2} d\mathbf{z} \\ & = \frac{1}{2\kappa\varepsilon} \int_{\mathbb{R}^q} |a_{\mathbf{P}\mathbf{z}}(\varepsilon) - a_{\mathbf{P}\mathbf{z}}(-\varepsilon)| \|\mathbf{z}\| e^{-\frac{1}{2}\|\mathbf{z}\|^2} d\mathbf{z} \\ & \leq \frac{\varepsilon^2 L_2}{6\kappa} \int_{\mathbb{R}^q} \|\mathbf{z}\|^4 e^{-\frac{1}{2}\|\mathbf{z}\|^2} d\mathbf{z} \leq \frac{\varepsilon^2}{6} L_2 (q + 4)^2. \end{aligned}$$

The proof is completed. □

**Theorem 2.** Let  $f(\mathbf{x}) = \mathbf{x}^\top \mathbf{H}\mathbf{x}$  and  $\mathbf{z} \sim \mathcal{N}(\mathbf{0}, \mathbf{I}_q)$ , where  $\mathbf{H} \in \mathbb{R}^{d \times d}$  is positive definite. We have

$$\mathbb{E}_{\mathbf{z}}[\hat{g}_\varepsilon(\mathbf{x}, \mathbf{P}, \mathbf{z})] = \mathbf{P}\mathbf{P}^\top \nabla f(\mathbf{x}), \quad (14)$$

$$\mathbb{E}_{\mathbf{z}}[\|\hat{g}_\varepsilon(\mathbf{x}, \mathbf{P}, \mathbf{z})\|^2] = (q + 2) \|\mathbf{P}^\top \nabla f(\mathbf{x})\|^2, \quad (15)$$

$$\mathbb{E}_{\mathbf{z}} \left[ \frac{\langle \nabla f(\mathbf{x}), \hat{g}_\varepsilon(\mathbf{x}, \mathbf{P}, \mathbf{z}) \rangle^2}{\|\mathbf{P}^\top \nabla f(\mathbf{x})\|^2 \|\hat{g}_\varepsilon(\mathbf{x}, \mathbf{P}, \mathbf{z})\|^2} \right] = \frac{1}{q}. \quad (16)$$

*Proof.* It is easy to check that  $\hat{g}_\varepsilon(\mathbf{x}, \mathbf{P}, \mathbf{z}) = \mathbf{P}\langle \mathbf{P}^\top \nabla f(\mathbf{x}), \mathbf{z} \rangle \mathbf{z}$ . Thus we have  $\mathbb{E}_z[\hat{g}_\varepsilon(\mathbf{x}, \mathbf{P}, \mathbf{z})] = \mathbf{P}\mathbf{P}^\top \nabla f(\mathbf{x})$ . Combined with Lemma 7, we get  $\mathbb{E}_z[\|\hat{g}_\varepsilon(\mathbf{x}, \mathbf{P}, \mathbf{z})\|^2] = (q+2)\|\mathbf{P}^\top \nabla f(\mathbf{x})\|^2$ . Note that for any orthogonal  $(q \times q)$ -matrix  $\mathbf{Q}$ , we have

$$\begin{aligned} \mathbb{E}_z \left[ \frac{\langle \nabla f(\mathbf{x}), \hat{g}_\varepsilon(\mathbf{x}, \mathbf{P}, \mathbf{z}) \rangle^2}{\|\mathbf{P}^\top \nabla f(\mathbf{x})\|^2 \|\hat{g}_\varepsilon(\mathbf{x}, \mathbf{P}, \mathbf{z})\|^2} \right] &= \mathbb{E}_z \left[ \frac{\langle \mathbf{P}^\top \nabla f(\mathbf{x}), \mathbf{z} \rangle^2}{\|\mathbf{P}^\top \nabla f(\mathbf{x})\|^2 \|\mathbf{z}\|^2} \right] \\ &= \mathbb{E}_z \left[ \frac{\langle \mathbf{P}^\top \nabla f(\mathbf{x}), \mathbf{Q}\mathbf{z} \rangle^2}{\|\mathbf{P}^\top \nabla f(\mathbf{x})\|^2 \|\mathbf{Q}\mathbf{z}\|^2} \right] \\ &= \mathbb{E}_z \left[ \frac{\langle \mathbf{Q}^\top \mathbf{P}^\top \nabla f(\mathbf{x}), \mathbf{z} \rangle^2}{\|\mathbf{Q}^\top \mathbf{P}^\top \nabla f(\mathbf{x})\|^2 \|\mathbf{z}\|^2} \right]. \end{aligned}$$

In accordance with Lemma 3, we can set  $\mathbf{P}^\top \nabla f(\mathbf{x}) = [1, 0, \dots, 0]^\top$ . Thus we have

$$\mathbb{E}_z \left[ \frac{\langle \nabla f(\mathbf{x}), \hat{g}_\varepsilon(\mathbf{x}, \mathbf{P}, \mathbf{z}) \rangle^2}{\|\mathbf{P}^\top \nabla f(\mathbf{x})\|^2 \|\hat{g}_\varepsilon(\mathbf{x}, \mathbf{P}, \mathbf{z})\|^2} \right] = \mathbb{E}_z \left[ \frac{z_1^2}{\|\mathbf{z}\|^2} \right] = \frac{1}{q}.$$

The proof is completed.  $\square$

To illustrate the convergence of Subzero with SGD, our analysis is divided into two main segments. We first investigate the convergence behavior of SubZero solution process while keeping the projection matrix  $\mathbf{P}$  constant. Next, we evaluate the effects of the lazy updates to  $\mathbf{P}$ . Based on these evaluations, we establish the global convergence of Subzero. Without loss of generality, we concentrate on the scenario where the number of layers is 1.

First, when the subspace  $\mathbf{P}$  is fixed, the original problem of SubZero can be reformulated as an optimization problem within the subspace. Define  $h(\mathbf{y}) = f(\mathbf{x} + \mathbf{P}\mathbf{y})$ ,  $h_\varepsilon(\mathbf{y}) = \mathbb{E}_z[h(\mathbf{y} + \varepsilon\mathbf{z})]$ , and  $g_\varepsilon(\mathbf{y}) = \frac{h(\mathbf{y} + \varepsilon\mathbf{z}) - f(\mathbf{y})}{\varepsilon} \mathbf{z}$ . According to Lemma 8, if  $f$  is first  $L_1$ -smooth, then  $h$  is also first  $L_1$ -smooth.

**Lemma 8.** *Let  $h(\mathbf{y}) = f(\mathbf{x} + \mathbf{P}\mathbf{y})$ , where  $f \in C_{L_1}^{1,1}(\mathbb{R}^d)$ , and  $\mathbf{P}^\top \mathbf{P} = \mathbf{I}$ , then we have  $h \in C_{L_1}^{1,1}(\mathbb{R}^q)$ .*

*Proof.* The following proves that if  $f$  is first  $L_1$ -smooth, then  $h$  is also first  $L_1$ -smooth. For any  $\mathbf{y}_1 \in \mathbb{R}^q$  and  $\mathbf{y}_2 \in \mathbb{R}^q$ , we have

$$\begin{aligned} \|\nabla h(\mathbf{y}_1) - \nabla h(\mathbf{y}_2)\| &= \|\mathbf{P}^\top \nabla(f(\mathbf{x} + \mathbf{P}\mathbf{y}_1) - \mathbf{P}^\top \nabla(f(\mathbf{x} + \mathbf{P}\mathbf{y}_2)))\| \\ &\leq \|\mathbf{P}^\top\| \|\nabla(f(\mathbf{x} + \mathbf{P}\mathbf{y}_1) - \nabla(f(\mathbf{x} + \mathbf{P}\mathbf{y}_2)))\| \\ &\leq L_1 \|\mathbf{P}(\mathbf{y}_1 - \mathbf{y}_2)\| \\ &= L_1 \|\mathbf{y}_1 - \mathbf{y}_2\|. \end{aligned}$$

The proof is completed.  $\square$

Now, we can analyze the convergence of SubZero when fixing the subspace.

**Lemma 9.** [40] *Let  $f \in C_{L_1}^{1,1}(\mathbb{R})$ . Then, for any  $\mathbf{x} \in \mathbb{R}$ , we have*

$$E_z[\|g_\varepsilon(\mathbf{x})\|^2] = E_z \left[ \left\| \frac{f(\mathbf{x} + \varepsilon\mathbf{z}) - f(\mathbf{x})}{\varepsilon} \right\|^2 \right] \leq 4(n+4)\|\nabla f_\varepsilon(\mathbf{x})\|^2 + 3\varepsilon^2 L_1^2(f)(n+4)^3, \quad (18)$$

and

$$\|\nabla f(\mathbf{x})\|^2 \leq 2\|\nabla f_\varepsilon(\mathbf{x})\|^2 + \frac{\varepsilon^2}{2} L_1^2(f)(n+6)^3, \quad (19)$$

where  $f_\varepsilon(\mathbf{x}) = \mathbb{E}_z[f(\mathbf{x} + \varepsilon\mathbf{z})]$ .

**Lemma 10.** *Let  $\mathbf{y}^* = \arg \min_{\mathbf{x} \in \mathbb{R}^q} h(\mathbf{y})$ , where  $h \in C_{L_1}^{1,1}(\mathbb{R}^q)$  and  $h$  is non-convex. Suppose  $\mathcal{E}_k = (\mathbf{z}_0, \mathbf{z}_1, \dots, \mathbf{z}_{k-1}, \mathbf{z}_k)$ , where  $\mathbf{z}_k \sim \mathcal{N}(0, \mathbf{I}_q)$  and  $\eta = \frac{1}{4(q+4)L_1}$ .  $\{\mathbf{y}_k\}_{k>0}$  is the sequence generated by Algorithm 3. Let  $\phi_0 = h(\mathbf{y}_0)$ , and for  $k \geq 1$ ,  $\phi_k = \mathbb{E}_{\mathcal{E}_{k-1}}[h(\mathbf{y}_k)]$ . For the  $\mathbf{P}$  defined in (10), which is fixed, we have*

$$\phi_{k+1} - \phi_k \leq -\frac{1}{4}\eta \mathbb{E}_{\mathcal{E}_k}[\|\nabla h(\mathbf{y}_k)\|^2] + \frac{\varepsilon^2(q+6)^3}{8} L_1^2 + \frac{3\varepsilon^2(q+4)}{32} L_1 \quad (20)$$



*Proof.* If a subspace  $\mathbf{P} \in \mathbb{R}^{d \times q}$  is fixed, the optimization objective can be reformulated as

$$\min_{\mathbf{y} \in \mathbb{R}^q} h(\mathbf{y}) := f(\mathbf{x} + \mathbf{P}\mathbf{y}),$$

Let  $\mathbf{y}_0$  be an initial point and  $\{\eta_k\}_{k \geq 0}$  a sequence of positive real numbers. Consider the randomized gradient search algorithm  $\mathcal{R}\mathcal{G}_\varepsilon(\varepsilon > 0)$ :

- 1) Generate  $\mathbf{z}_k$  and the corresponding  $g_\varepsilon(\mathbf{y}_k)$ , where  $\mathbf{z}_k \sim \mathcal{N}(\mathbf{0}, \mathbf{I}_q)$ .
- 2) Update  $\mathbf{y}_{k+1} = \mathbf{y}_k - \eta_k g_\varepsilon(\mathbf{y}_k)$ .

We aim to estimate the evolution of the function  $h_\varepsilon$  after one iteration of this algorithm.

Given that  $h$  is  $L_1$ -Lipschitz continuous for the first derivative, and  $h_\varepsilon$  is  $L_\varepsilon$ -Lipschitz continuous for the first derivative (where  $L_\varepsilon \leq L_1$ ) [40]. Thus, we have

$$h_\varepsilon(\mathbf{y}_{k+1}) \leq h_\varepsilon(\mathbf{y}_k) - \eta_k \langle \nabla h_\varepsilon(\mathbf{y}_k), g_\varepsilon(\mathbf{y}_k) \rangle + \frac{1}{2} \eta_k^2 L_\varepsilon \|g_\varepsilon(\mathbf{y}_k)\|^2.$$

Taking expectation with respect to  $\mathbf{z}_k$ , we obtain

$$\mathbb{E}_{\mathbf{z}_k}[h_\varepsilon(\mathbf{y}_{k+1})] \leq h_\varepsilon(\mathbf{y}_k) - \eta_k \|\nabla h_\varepsilon(\mathbf{y}_k)\|^2 + \frac{1}{2} \eta_k^2 L_\varepsilon \mathbb{E}_{\mathbf{z}_k}[\|g_\varepsilon(\mathbf{y}_k)\|^2].$$

Since  $h \in C^{1,1}(\mathbb{R}^q)$ , from Lemma 9, we have

$$\begin{aligned} \mathbb{E}_{\mathbf{z}_k}[h_\varepsilon(\mathbf{y}_{k+1})] &\leq h_\varepsilon(\mathbf{y}_k) - \eta_k \|\nabla h_\varepsilon(\mathbf{y}_k)\|^2 \\ &\quad + \frac{1}{2} \eta_k^2 L_1 (4(q+4) \|\nabla h_\varepsilon(\mathbf{y}_k)\|^2 + 3\varepsilon^2 L_1^2 (q+4)^3). \end{aligned}$$

Setting  $\eta_k = \hat{\eta} = \frac{1}{4(q+4)L_1}$ , we get

$$\mathbb{E}_{\mathbf{z}_k}[h_\varepsilon(\mathbf{y}_{k+1})] \leq h_\varepsilon(\mathbf{y}_k) - \frac{1}{2} \hat{\eta} \|\nabla h_\varepsilon(\mathbf{y}_k)\|^2 + \frac{3\varepsilon^2}{32} L_1 (q+4).$$

Taking the expectation with respect to  $\mathcal{E}_k$ , we get

$$\phi_{k+1} \leq \phi_k - \frac{1}{2} \hat{\eta} \mathbb{E}_{\mathcal{E}_k}[\|\nabla h_\varepsilon(\mathbf{y}_k)\|^2] + \frac{3\varepsilon^2(q+4)}{32} L_1,$$

From Lemma 9, we have  $\mathbb{E}_{\mathcal{E}_k}[\|\nabla h(\mathbf{y}_k)\|^2] \leq 2\mathbb{E}_{\mathcal{E}_k}[\|\nabla h_\varepsilon(\mathbf{y}_k)\|^2] + \frac{\varepsilon^2(q+6)^3}{2} L_1^2$ . Therefore,

$$\phi_{k+1} - \phi_k \leq -\frac{1}{4} \hat{\eta} \mathbb{E}_{\mathcal{E}_k}[\|\nabla h(\mathbf{y}_k)\|^2] + \frac{\varepsilon^2(q+6)^3}{8} L_1^2 + \frac{3\varepsilon^2(q+4)}{32} L_1. \quad (21)$$

The proof is completed.  $\square$

Next, we need to measure the randomness of our random subspace. From Lemma 14, if the projection matrix is obtained by Algorithm 1, we have  $\mathbb{E}[\mathbf{P}\mathbf{P}^T] = \frac{q}{d}\mathbf{I}$ , where  $q$  represents the dimension of the subspace,  $d$  represents the dimension of the origin space, and  $\mathbf{P} = \mathbf{V} \otimes \mathbf{U}$  (see Lemma 1).

**Lemma 11.** *Let matrix  $\mathbf{A} = (\mathbf{a}_1, \mathbf{a}_2, \dots, \mathbf{a}_r) \in \mathbb{R}^{n \times r}$  be composed of column vectors  $\mathbf{a}_k$  which are mutually independent and  $\mathbf{a}_k \in \mathcal{N}(0, \mathbf{I}_n)$ . Suppose Gram-Schmidt process  $\mathbf{u}_k = \mathbf{a}_k - \sum_{s=1}^{k-1} \langle \mathbf{a}_k, \mathbf{e}_s \rangle \mathbf{e}_s$  and  $\mathbf{e}_k = \frac{\mathbf{u}_k}{\|\mathbf{u}_k\|}$ .  $[\mathbf{a}_k]_i \leftrightarrow [\mathbf{a}_k]_j$  represents the exchange of the  $i$ -th element and the  $j$ -th element of  $\mathbf{a}_k$ , while all other elements remain unchanged.  $[\mathbf{a}_k]_i = -1 \times [\mathbf{a}_k]_i$  signifies that only the  $i$ -th element of  $\mathbf{a}_k$  is multiplied by  $-1$ , while all other elements remain unchanged. Suppose  $f(\mathbf{A}, \mathbf{U}, \mathbf{E})$  be a function of the matrix  $\mathbf{A}$ ,  $\mathbf{U} = (\mathbf{u}_1, \mathbf{u}_2, \dots, \mathbf{u}_r)$  and  $\mathbf{E} = (\mathbf{e}_1, \mathbf{e}_2, \dots, \mathbf{e}_r)$ , then*

- (1) if  $[\mathbf{a}_k]_i \leftrightarrow [\mathbf{a}_k]_j$  or  $[\mathbf{a}_k]_i = -1 \times [\mathbf{a}_k]_i$ ,  $\mathbb{E}[f]$  remain unchanged.
- (2) if  $[\mathbf{a}_k]_i \leftrightarrow [\mathbf{a}_k]_j \Rightarrow [\mathbf{u}_k]_i \leftrightarrow [\mathbf{u}_k]_j$  and  $[\mathbf{e}_k]_i \leftrightarrow [\mathbf{e}_k]_j$ .
- (3) if  $[\mathbf{a}_k]_i = -1 \times [\mathbf{a}_k]_i \Rightarrow [\mathbf{u}_k]_i = -1 \times [\mathbf{u}_k]_i$ ,  $[\mathbf{e}_k]_i = -1 \times [\mathbf{e}_k]_i$ ,  $[\mathbf{u}_k]_j = 1 \times [\mathbf{u}_k]_j$ , and  $[\mathbf{e}_k]_j = 1 \times [\mathbf{e}_k]_j$ , where  $i \neq j$ .

$$(4) \mathbb{E} \left[ \frac{[\mathbf{u}_k]_i^2}{\langle \mathbf{u}_k, \mathbf{u}_k \rangle} \right] = \frac{1}{n}.$$

$$(5) \mathbb{E} \left[ \frac{[\mathbf{u}_k]_i [\mathbf{u}_k]_j}{\langle \mathbf{u}_k, \mathbf{u}_k \rangle} \right] = 0, \text{ where } i \neq j.$$

*Proof.* According to real analysis, the matrix  $\mathbf{A}$  is full rank almost everywhere under a Gaussian distribution, and both  $\mathbf{u}_k$  and  $\mathbf{e}_k$  are non-zero almost everywhere.

(1) Since  $\mathbf{a}_k$  is independently and identically distributed, it obviously holds.

(2) For base case  $k = 1$ , it obviously holds. Assume the result holds for all  $k = 1, 2, \dots, k-1$ , where  $k \geq 2$ , then  $[\mathbf{a}_k]_i \leftrightarrow [\mathbf{a}_k]_j \Rightarrow [\mathbf{u}_k]_i = [\mathbf{a}_k]_i - \sum_{s=1}^{k-1} \langle \mathbf{a}_k, \mathbf{e}_s \rangle [\mathbf{e}_s]_i$ ,  $[\mathbf{u}_k]_j = [\mathbf{a}_k]_j - \sum_{s=1}^{k-1} \langle \mathbf{a}_k, \mathbf{e}_s \rangle [\mathbf{e}_s]_j$ ,  $[\mathbf{e}_k]_i = \frac{[\mathbf{u}_k]_i}{\|\mathbf{u}_k\|}$ , and  $[\mathbf{e}_k]_j = \frac{[\mathbf{u}_k]_j}{\|\mathbf{u}_k\|}$ .

Thus, by strong induction, we have  $[\mathbf{u}_k]_i \leftrightarrow [\mathbf{u}_k]_j$  and  $[\mathbf{e}_k]_i \leftrightarrow [\mathbf{e}_k]_j$ .

(3) For base case  $k = 1$ , it obviously holds. Assume the result holds for all  $k = 1, 2, \dots, k-1$ , where  $k \geq 2$ , then

$$\begin{aligned} [\mathbf{a}_k]_i = -1 \times [\mathbf{a}_k]_i &\Rightarrow \begin{cases} [\mathbf{u}_k]_i = [\mathbf{a}_k]_i \times (-1) - \sum_{s=1}^{k-1} \langle \mathbf{a}_k, \mathbf{e}_s \rangle [\mathbf{e}_s]_i \times (-1) = [\mathbf{u}_k]_i \times (-1) \\ [\mathbf{u}_k]_j = [\mathbf{u}_k]_j \times 1, i \neq j \end{cases} \\ &\Rightarrow \begin{cases} [\mathbf{e}_k]_i \times (-1) = \frac{[\mathbf{u}_k]_i}{\|\mathbf{u}_k\|} \times (-1) \\ [\mathbf{e}_k]_j = [\mathbf{e}_k]_j \times 1, j \neq i \end{cases} \end{aligned}$$

By strong induction, we have  $[\mathbf{u}_k]_i = -1 \times [\mathbf{u}_k]_i$ ,  $[\mathbf{e}_k]_i = -1 \times [\mathbf{e}_k]_i$ ,  $[\mathbf{u}_k]_j = 1 \times [\mathbf{u}_k]_j$ , and  $[\mathbf{e}_k]_j = 1 \times [\mathbf{e}_k]_j$ , where  $i \neq j$ .  $\square$

(4) Since  $\left| \frac{[\mathbf{u}_k]_i^2}{\langle \mathbf{u}_k, \mathbf{u}_k \rangle} \right| \leq 1$ ,  $\mathbb{E} \left[ \frac{[\mathbf{u}_k]_i^2}{\langle \mathbf{u}_k, \mathbf{u}_k \rangle} \right]$  exists.  $[\mathbf{a}_k]_i \leftrightarrow [\mathbf{a}_k]_j \Rightarrow \frac{[\mathbf{u}_k]_i^2}{\langle \mathbf{u}_k, \mathbf{u}_k \rangle} \leftrightarrow \frac{[\mathbf{u}_k]_j^2}{\langle \mathbf{u}_k, \mathbf{u}_k \rangle}$ .

Thus,  $\mathbb{E} \left[ \frac{[\mathbf{u}_k]_i^2}{\langle \mathbf{u}_k, \mathbf{u}_k \rangle} \right] \times n = \sum_{s=1}^n \mathbb{E} \left[ \frac{[\mathbf{u}_k]_s^2}{\langle \mathbf{u}_k, \mathbf{u}_k \rangle} \right] = \mathbb{E} \left[ \frac{\langle \mathbf{u}_k, \mathbf{u}_k \rangle}{\langle \mathbf{u}_k, \mathbf{u}_k \rangle} \right] = 1 \Rightarrow \mathbb{E} \left[ \frac{[\mathbf{u}_k]_i^2}{\langle \mathbf{u}_k, \mathbf{u}_k \rangle} \right] = \frac{1}{n}$ .

(5) Since  $\left| \frac{[\mathbf{u}_k]_i [\mathbf{u}_k]_j}{\langle \mathbf{u}_k, \mathbf{u}_k \rangle} \right| \leq \left| \frac{[\mathbf{u}_k]_i^2 + [\mathbf{u}_k]_j^2}{2\langle \mathbf{u}_k, \mathbf{u}_k \rangle} \right| \leq 1$ ,  $\mathbb{E} \left[ \frac{[\mathbf{u}_k]_i [\mathbf{u}_k]_j}{\langle \mathbf{u}_k, \mathbf{u}_k \rangle} \right]$  exists.

$[\mathbf{a}_k]_i = [\mathbf{a}_k]_i \times -1 \Rightarrow \mathbb{E} \left[ \frac{[\mathbf{u}_k]_i [\mathbf{u}_k]_j}{\langle \mathbf{u}_k, \mathbf{u}_k \rangle} \right] = \mathbb{E} \left[ \frac{-[\mathbf{u}_k]_i [\mathbf{u}_k]_j}{\langle \mathbf{u}_k, \mathbf{u}_k \rangle} \right] = 0$ , where  $i \neq j$ .

**Lemma 12.** Let  $\mathbf{A} \in \mathbb{R}^{n \times r}$  be a matrix with independent standard normal entries, i.e., each element of  $\mathbf{A}$  is an i.i.d.  $\mathcal{N}(0, 1)$  random variable. Suppose  $\mathbf{A}$  undergoes QR decomposition via the Gram-Schmidt process to yield a column-orthogonal matrix  $\mathbf{Q} \in \mathbb{R}^{n \times r}$  with orthonormal columns  $\mathbf{e}_1, \mathbf{e}_2, \dots, \mathbf{e}_r$  and an upper triangular matrix  $\mathbf{R} \in \mathbb{R}^{r \times r}$ . Then, for each  $k = 1, 2, \dots, r$ , the expected value of the outer product of the  $k$ -th orthonormal column vector  $\mathbf{e}_k$  of  $\mathbf{Q}$  is given by:

$$\mathbb{E}[\mathbf{e}_k \mathbf{e}_k^T] = \frac{1}{n} \mathbf{I},$$

where  $\mathbf{I}$  is the  $n \times n$  identity matrix.

*Proof.* By the Gram-Schmidt process, we have  $\mathbf{e}_k = \frac{\mathbf{u}_k}{\|\mathbf{u}_k\|}$ , where  $\mathbf{u}_k = \mathbf{a}_k - \sum_{s=1}^{k-1} \langle \mathbf{a}_k, \mathbf{e}_s \rangle \mathbf{e}_s$ . Thus,  $\mathbf{e}_k \mathbf{e}_k^T = \frac{\mathbf{u}_k \mathbf{u}_k^T}{\langle \mathbf{u}_k, \mathbf{u}_k \rangle}$ .

The  $(i, j)$ -th entry of  $\mathbb{E}[\mathbf{e}_k \mathbf{e}_k^T]$  can be written as:

$$\mathbb{E}[(\mathbf{e}_k \mathbf{e}_k^T)_{ij}] = \mathbb{E} \left[ \frac{[\mathbf{u}_k]_i [\mathbf{u}_k]_j}{\langle \mathbf{u}_k, \mathbf{u}_k \rangle} \right].$$

For diagonal entries ( $i = j$ ): When  $i = j$ , from Lemma 11(4), we have:

$$\mathbb{E}[(\mathbf{e}_k \mathbf{e}_k^T)_{ii}] = \mathbb{E} \left[ \frac{[\mathbf{u}_k]_i^2}{\langle \mathbf{u}_k, \mathbf{u}_k \rangle} \right] = \frac{1}{n}.$$

For off-diagonal entries ( $i \neq j$ ): When  $i \neq j$ , from Lemma 11(5), we have:

$$\mathbb{E}[(\mathbf{e}_k \mathbf{e}_k^T)_{ij}] = \mathbb{E} \left[ \frac{[\mathbf{u}_k]_i [\mathbf{u}_k]_j}{\langle \mathbf{u}_k, \mathbf{u}_k \rangle} \right] = 0.$$

Combining these two cases, we conclude that  $\mathbb{E}[\mathbf{e}_k \mathbf{e}_k^T]$  is a diagonal matrix with all diagonal entries equal to  $\frac{1}{n}$ . Thus,

$$\mathbb{E}[\mathbf{e}_k \mathbf{e}_k^T] = \frac{1}{n} \mathbf{I},$$

where  $\mathbf{I}$  is the  $n \times n$  identity matrix. The proof is completed.  $\square$

**Lemma 13.** Let  $\mathbf{A} \in \mathbb{R}^{n \times r}$  be a matrix with independent standard normal entries, i.e., each element of  $\mathbf{A}$  is an i.i.d.  $\mathcal{N}(0, 1)$  random variable. Suppose  $\mathbf{A}$  undergoes QR decomposition to yield an orthogonal matrix  $\mathbf{Q} \in \mathbb{R}^{n \times r}$  with orthonormal columns and an upper triangular matrix  $\mathbf{R} \in \mathbb{R}^{r \times r}$ . Then, the expected value of the outer product of the matrix  $\mathbf{Q}$  with itself is given by:

$$\mathbb{E}[\mathbf{Q}\mathbf{Q}^T] = \frac{r}{n}\mathbf{I}$$

where  $\mathbf{I}$  is the  $n \times n$  identity matrix.

*Proof.* The QR decomposition of  $\mathbf{A}$  is given by  $\mathbf{A} = \mathbf{Q}\mathbf{R}$ , where  $\mathbf{Q}$  is an orthogonal matrix with columns  $\mathbf{e}_1, \mathbf{e}_2, \dots, \mathbf{e}_r$  and  $\mathbf{R}$  is an upper triangular matrix. Since  $\mathbf{Q}$  is orthogonal,  $\mathbf{Q}\mathbf{Q}^T = \mathbf{I}_r$ , where  $\mathbf{I}_r$  is the  $r \times r$  identity matrix. We aim to compute  $\mathbb{E}[\mathbf{Q}\mathbf{Q}^T]$ . By linearity of expectation and the fact that the columns of  $\mathbf{Q}$  are orthonormal, we have:

$$\mathbb{E}[\mathbf{Q}\mathbf{Q}^T] = \mathbb{E}\left[\sum_{k=1}^r \mathbf{e}_k \mathbf{e}_k^T\right] = \sum_{k=1}^r \mathbb{E}[\mathbf{e}_k \mathbf{e}_k^T].$$

From Lemma 12, we know that  $\mathbb{E}[\mathbf{e}_k \mathbf{e}_k^T] = \frac{1}{n}\mathbf{I}$  for each  $k$ . Therefore:

$$\mathbb{E}[\mathbf{Q}\mathbf{Q}^T] = \sum_{k=1}^r \frac{1}{n}\mathbf{I} = \frac{r}{n}\mathbf{I}.$$

The proof is completed. □

**Lemma 14.** Let  $\mathbf{A}_1 \in \mathbb{R}^{m \times r}$  and  $\mathbf{A}_2 \in \mathbb{R}^{n \times r}$  be matrices with independent standard normal entries, i.e., each element of  $\mathbf{A}_1$  and  $\mathbf{A}_2$  is an i.i.d.  $\mathcal{N}(0, 1)$  random variable. Suppose  $\mathbf{A}_1$  and  $\mathbf{A}_2$  undergo QR decomposition to yield orthogonal matrices  $\mathbf{Q}_1 \in \mathbb{R}^{m \times r}$  and  $\mathbf{Q}_2 \in \mathbb{R}^{n \times r}$  with orthonormal columns, respectively. Define  $\mathbf{P} = \mathbf{Q}_2 \otimes \mathbf{Q}_1$ , where  $\otimes$  denotes the Kronecker product. Then, the expected value of the outer product of the matrix  $\mathbf{P}$  with itself is given by:

$$\mathbb{E}[\mathbf{P}\mathbf{P}^T] = \frac{r^2}{mn}\mathbf{I},$$

where  $\mathbf{I}$  is the  $mn \times mn$  identity matrix.

*Proof.* The Kronecker product  $\mathbf{P} = \mathbf{Q}_2 \otimes \mathbf{Q}_1$  results in a matrix  $\mathbf{P} \in \mathbb{R}^{mn \times r^2}$ . From Lemma 13, we have  $\mathbf{Q}_1\mathbf{Q}_1^T = \frac{r}{m}\mathbf{I}$  and  $\mathbf{Q}_2\mathbf{Q}_2^T = \frac{r}{n}\mathbf{I}$ . We aim to compute  $\mathbb{E}[\mathbf{P}\mathbf{P}^T]$ . Using the properties of the Kronecker product, we have:

$$\mathbb{E}[\mathbf{P}\mathbf{P}^T] = \mathbb{E}[(\mathbf{Q}_2 \otimes \mathbf{Q}_1)(\mathbf{Q}_2^T \otimes \mathbf{Q}_1^T)] = \mathbb{E}[(\mathbf{Q}_2\mathbf{Q}_2^T)] \otimes \mathbb{E}[(\mathbf{Q}_1\mathbf{Q}_1^T)] = \frac{r^2}{mn}\mathbf{I} \otimes \mathbf{I} = \frac{r^2}{mn}\mathbf{I}$$

The proof is completed. □

Now we can assess the impact of the lazy updates to  $\mathbf{P}$ .

**Theorem 3.** Let  $\mathbf{x}^* = \arg \min_{\mathbf{x} \in \mathbb{R}^d} f(\mathbf{x})$ , where  $f \in C_{L_1}^{1,1}(\mathbb{R}^d)$  and  $f$  is non-convex. Suppose  $\mathcal{E}_k = (\mathbf{z}_0, \mathbf{z}_1, \dots, \mathbf{z}_k)$ , where  $\mathbf{z}_k \sim \mathcal{N}(\mathbf{0}, \mathbf{I}_q)$  and  $\eta = \frac{1}{4(q+4)L_1}$ .  $\{\mathbf{x}_k\}_{k>0}$  is the sequence generated by Algorithm 3. For the  $\mathbf{P}$  defined in (10), which is updated lazily at a fixed frequency  $F$ , we have

$$\frac{1}{T} \sum_{k=0}^{T-1} \mathbb{E}_{\mathcal{E}_k} [\|\nabla f(\mathbf{x}_k)\|^2] \leq \epsilon$$

for any  $T = \Omega\left(\frac{d}{\epsilon}\right)$  if  $\epsilon \leq \mathcal{O}\left(\frac{\epsilon^{1/2}}{q^{3/2}d^{1/2}L_1^{3/2}}\right)$ , where  $T = KF$ ,  $K$  represents the total number of subspace updates, and  $\epsilon$  represents perturbation scale.

*Proof.* Suppose  $\mathcal{P}_j = (\mathbf{P}_0, \mathbf{P}_1, \dots, \mathbf{P}_j)$ , where  $\mathbf{P}_j$  is the sequence generated by Eqn. (10) and  $j \leq K$ . In accordance with Lemma 8 and Lemma 10, if the subspace is fixed, we can transform the original problem  $f \in C_{L_1}^{1,1}(\mathbb{R}^d)$  into  $h \in C_{L_1}^{1,1}(\mathbb{R}^q)$  through transformation  $h(\mathbf{y}) = f(\mathbf{x} + \mathbf{P}_j \mathbf{y})$ . Consider the update rule:

$$\mathbf{y}_{j,0} = 0, h_j(\mathbf{y}) = f(\mathbf{x}_{jF} + \mathbf{P}_j \mathbf{y}), \forall j \in 0, 1, \dots, K-1 \quad (22)$$

$$\mathbf{y}_{j,k} = \mathbf{y}_{j,k-1} - \eta \widehat{\nabla} h_j(\mathbf{y}_{j,k-1}), \forall k \in 0, 1, \dots, F \quad (23)$$

$$\mathbf{x}_{jF+k} = \mathbf{x}_{jF} + \mathbf{P}_j \mathbf{y}_k, \quad (24)$$

In the  $j$ -th subspace, the projection matrix  $\mathbf{P}_j$  remains constant, hence we can accumulate the changes of  $\phi$  within the current subspace. Using Lemma 10, we have

$$\phi_{(j+1)F} - \phi_{jF} \leq -\frac{1}{4} \hat{\eta} \sum_{i=0}^{K-1} \mathbb{E}_{\mathcal{E}_{jF+i}} [\|\nabla h_j(\mathbf{y}_{j,i})\|^2] + \frac{\varepsilon^2(q+6)^3}{8} KL_1^2 + \frac{3\varepsilon^2(q+4)}{32} KL_1 \quad (25)$$

$$\leq -\frac{1}{4} \hat{\eta} \mathbb{E}_{\mathcal{E}_{jF}} [\|\nabla h_j(\mathbf{y}_{j,0})\|^2] + \frac{\varepsilon^2(q+6)^3}{8} KL_1^2 + \frac{3\varepsilon^2(q+4)}{32} KL_1. \quad (26)$$

Additionally, we note that  $\nabla h_j(\mathbf{y}_{j,0}) = (\mathbf{P}_j)^\top \nabla f(\mathbf{x}_{jF})$ . Taking expectations over the overall historical projection matrix  $\mathcal{P}_j$ , and noting Lemma 14,  $\mathbb{E}[\mathbf{P}_j (\mathbf{P}_j)^\top] = \frac{q}{d} \mathbf{I}$ , with  $\mathbf{P}_j$  independent of  $\mathbf{x}_{jF}$ , we get

$$\mathbb{E}_{\mathcal{P}_{j+1}} [\phi_{(j+1)F}] - \mathbb{E}_{\mathcal{P}_j} [\phi_{jF}] \leq -\frac{1}{4} \hat{\eta} \mathbb{E}_{\mathcal{E}_{jF}, \mathcal{P}_j} [\|(\mathbf{P}_j)^\top \nabla f(\mathbf{x}_{jF})\|^2] + \frac{\varepsilon^2(q+6)^3}{8} KL_1^2 + \frac{3\varepsilon^2(q+4)}{32} KL_1 \quad (27)$$

$$= -\frac{q}{4d} \hat{\eta} \mathbb{E}_{\mathcal{E}_{jF}, \mathcal{P}_j} [\|\nabla f(\mathbf{x}_{jF})\|^2] + \frac{\varepsilon^2(q+6)^3}{8} KL_1^2 + \frac{3\varepsilon^2(q+4)}{32} KL_1. \quad (28)$$

Assuming  $f(\mathbf{x}) \geq f^*$  holds for all  $\mathbf{x} \in \mathbb{R}^d$ , and letting  $T = KF$ , summing the inequality yields

$$\mathbb{E}_{\mathcal{P}_{K-1}} [\phi_T] \leq \mathbb{E}_{\mathcal{P}_0} [\phi_0] - \frac{q}{4d} \hat{\eta} \sum_{j=0}^{K-1} \mathbb{E}_{\mathcal{E}_{jF}, \mathcal{P}_j} [\|\nabla f(\mathbf{x}_{jF})\|^2] + T \frac{\varepsilon^2(q+6)^3}{8} L_1^2 + T \frac{3\varepsilon^2(q+4)}{32} L_1. \quad (29)$$

Since  $\mathbb{E}_{\mathcal{P}_{K-1}} [\phi_T] \geq f^*$ , we have:

$$f^* \leq \mathbb{E}_{\mathcal{P}_0} [\phi_0] - \frac{q}{4d} \hat{\eta} \sum_{j=0}^{K-1} \mathbb{E}_{\mathcal{E}_{jF}, \mathcal{P}_j} [\|\nabla f(\mathbf{x}_{jF})\|^2] + T \frac{\varepsilon^2(q+6)^3}{8} L_1^2 + T \frac{3\varepsilon^2(q+4)}{32} L_1. \quad (30)$$

Rearranging the inequality, we get

$$\frac{q}{4d} \hat{\eta} \sum_{j=0}^{K-1} \mathbb{E}_{\mathcal{E}_{jF}, \mathcal{P}_j} [\|\nabla f(\mathbf{x}_{jF})\|^2] \leq \mathbb{E}_{\mathcal{P}_0} [\phi_0] - f^* + T \frac{\varepsilon^2(q+6)^3}{8} L_1^2 + T \frac{3\varepsilon^2(q+4)}{32} L_1. \quad (31)$$

Substituting  $\hat{\eta} = \frac{1}{4(q+4)L_1}$ , we obtain:

$$\frac{q}{16d(q+4)L_1} \sum_{j=0}^{K-1} \mathbb{E}_{\mathcal{E}_{jF}, \mathcal{P}_j} [\|\nabla f(\mathbf{x}_{jF})\|^2] \leq \mathbb{E}_{\mathcal{P}_0} [\phi_0] - f^* + T \frac{\varepsilon^2(q+6)^3}{8} L_1^2 + T \frac{3\varepsilon^2(q+4)}{32} L_1. \quad (32)$$

Thus, we have

$$\frac{1}{T} \sum_{k=0}^T \mathbb{E}_{\mathcal{E}_k, \mathcal{P}_{\lfloor k/F \rfloor}} [\|\nabla f(\mathbf{x}_k)\|^2] \leq \frac{16(q+4)dL_1(\mathbb{E}_{\mathcal{P}_0} [\phi_0] - f^*)}{qT} + \frac{2\varepsilon^2(q+6)^3(q+4)d}{q} L_1^3 + \frac{3\varepsilon^2(q+4)^2d}{2q} L_1^2. \quad (33)$$

To ensure  $\sum_{k=0}^T \mathbb{E}_{\mathcal{E}_k, \mathcal{P}_{\lfloor k/F \rfloor}} [\|\nabla f(\mathbf{x}_k)\|^2] \leq \epsilon$ , we can choose

$$\varepsilon \leq O\left(\frac{\epsilon^{1/2}}{q^{3/2}d^{1/2}L_1^{3/2}}\right).$$

So that after  $\Omega\left(\frac{d}{\epsilon}\right)$  iterations,  $\sum_{k=0}^T \mathbb{E}_{\mathcal{E}_k, \mathcal{P}_{\lfloor k/F \rfloor}} [\|\nabla f(\mathbf{x}_k)\|^2] \leq \epsilon$ . The proof is completed.  $\square$