

ENHANCING ZERO-ORDER FINE-TUNING FOR LLMs VIA GRADIENT-GUIDED SUBSPACE SELECTION

Anonymous authors

Paper under double-blind review

ABSTRACT

As a promising memory-efficient technique, zeroth-order (ZO) optimization enables large language models (LLMs) to bypass the costly process of backpropagation during fine-tuning by estimating gradients through function evaluations. However, to minimize approximate variance in high-dimensional parameter spaces, existing ZO methods focus on exploring the estimate of gradients within random subspaces, neglecting the benefits of searching for more accurate subspaces of LLMs on gradient estimates. Due to inaccurate gradient estimates obtained from random spaces, fine-tuning performance is inevitably degraded, thus compromising the performance of downstream tasks. To address the limitation of existing ZO methods, this paper proposes a novel ZO subspace fine-tuning method named *SVD-0*. Based on singular value decomposition (SVD), SVD-0 can effectively obtain more accurate subspace projection matrices, which can be used to improve the accuracy of gradient estimates. Experimental results on various language modeling tasks show that SVD-0 achieves better fine-tuning performance than SOTA ZO methods.

1 INTRODUCTION

Due to the powerful capabilities of language understanding and reasoning, large language models (LLMs) have demonstrated significant performance on a wide range of tasks, such as mathematical reasoning (Guo et al., 2025), creative writing (Shanahan & Clarke, 2023). Currently, fine-tuning (FT) the pre-trained foundation model to adapt to downstream tasks has become the mainstream paradigm for AI application development. However, due to the extremely large number of model parameters, traditional first-order (FO) optimization-based fine-tuning methods face a serious challenge of excessive memory consumption. Typically, since the backpropagation process in FO requires storing activations and optimizer states, the memory requirements of FT are significantly larger than those of reasoning, which severely limits the development of LLM-based applications.

To achieve memory-efficient FT, existing methods can be classified into two categories, i.e., parameter-efficient fine-tuning (PEFT) methods (Liu et al., 2022; Han et al., 2024) and zeroth-order (ZO) optimization methods (Malladi et al., 2023). PEFT methods attempt to reduce the number of trainable parameters to alleviate memory requirements. However, since PEFT methods are still based on FO optimization, they require a significant amount of memory to store intermediate training results, which severely limits the choice of trainable parameters. ZO optimization methods (Malladi et al., 2023) emerge as a promising alternative by estimating gradients through forward-pass perturbations, thereby eliminating the memory overhead associated with backpropagation. However, conventional ZO methods face a critical challenge: the high variance of gradient approximations in billion-parameter spaces severely degrades optimization efficiency and model performance.

Recent advances in ZO optimization for LLMs, such as SubZero (Yu et al., 2024) and LOZO (Chen et al., 2025), attempt to mitigate this issue by constraining perturbations to random low-dimensional subspaces. These methods are based on the finding that gradient matrices become low-rank during LLM training and fine-tuning (Zhao et al., 2024a). While these subspace methods reduce approximation variance, they fundamentally rely on arbitrary projection matrices that fail to match the low-rank structure implied by the gradient. This limitation stems from a fundamental disconnect - the subspace construction process ignores critical gradient information that could guide more effective parameter updates. Therefore, *how to determine the optimal low-dimensional subspaces without relying on first-order optimizers* poses a fundamental challenge.

054 The similarity between the gradient estimated by ZO optimizers and the true gradient has been
055 experimentally demonstrated in (Malladi et al., 2023). We find it feasible to derive the low-rank
056 structure of the true gradient from the estimated gradient. In light of this idea, we conducted a
057 preliminary study (see details in Section 3). The experimental results indicate a significant similarity
058 between the estimated and true gradients, as demonstrated by the resemblance of their singular value
059 vectors. Consequently, we conclude that applying singular value decomposition (SVD) to the gradient
060 estimated by the ZO optimizer allows us to obtain a low-rank structure that closely resembles the
061 low-rank structure of the true gradient.

062 Motivated by the above findings, we propose SVD-0, a novel gradient-guided subspace optimization
063 framework that combines zeroth-order efficiency with principled subspace discovery. Our key
064 insight is that, while exact first-order gradients remain inaccessible due to memory constraints, ZO
065 gradient estimates contain sufficient directional information to reconstruct high-fidelity subspaces.
066 Specifically, SVD-0 periodically performs singular value decomposition (SVD) on ZO gradient
067 estimates to derive layer-wise projection matrices that capture dominant optimization directions. By
068 preserving the intrinsic structure of the subspace, our method effectively enhances the performance
069 of subspace-based ZO methods. The contributions of this work are summarized as follows:

- 070 • We propose a novel method for exploring more accurate subspace projection matrices and
071 conducting layer-wise perturbations on low-rank matrices. With periodic updates of the
072 projection matrices, our method continuously captures the subspaces of the parameters.
- 073 • We develop a novel gradient-guided ZO method to approximate these two projection matri-
074 ces, ensuring low memory usage throughout the entire fine-tuning process, to overcome the
075 paradox that obtaining subspace projection matrices requires FO gradients.
- 076 • We conduct comprehensive experiments on various model scales and language modeling
077 tasks. The corresponding results demonstrate the superiority of our method over various ZO
078 optimization methods specifically designed for LLM fine-tuning.

080 2 RELATED WORK

082 **Memory-efficient fine-tuning for LLMs.** Recent work has concentrated on exploring memory-
083 efficient fine-tuning methods to enable LLM fine-tuning on memory-intensive hardware. A critical
084 line of research centers on Parameter-Efficient Fine-Tuning (PEFT) methods (Liu et al., 2022; Han
085 et al., 2024) by freezing the backbone of LLMs while only tuning a small group of parameters.
086 For instance, LoRA (Hu et al., 2022) only updates parameters based on low-rank structures while
087 being competitive with full-parameter fine-tuning. LISA (Pan et al., 2024) distinguishes trainable
088 layers based on their contribution to task-specific performance and freezes other layers to reduce
089 the memory footprint. Further, parameter quantization (Lin et al., 2024; Frantar et al., 2022) has
090 played a pivotal role in enhancing memory efficiency. By discretizing model parameters (e.g., from
091 32-bit to 8-bit or lower precision), quantization methods such as QLoRA (Dettmers et al., 2023)
092 and LLM.int8() (Dettmers et al., 2022) reduce storage requirements without significant degradation
093 in task performance. Complementary to PEFT and the quantization method, subspace projection
094 techniques have emerged as a powerful strategy to reduce the dimensionality of the optimization
095 space. Galora (Zhao et al., 2024a) and FLORA (Hao et al., 2024) both leverage the low-rank property
096 of gradients to constrain updates on a compact subspace of the full parameter space (Huang et al.,
097 2025). By discovering the projection matrices of low-rank subspaces, the memory costs for storing
098 gradients and optimizer states (e.g., the first and second order states in Adam optimizer (Kingma &
Ba, 2014)) are greatly reduced.

099 **Zero-order optimization.** ZO approaches enable backpropagation-free optimization by approxim-
100 ating exact gradients through finite differences. This flexibility has driven interest in ZO for solving a
101 range of machine learning problems, including on-chip learning, black-box adversarial strategies,
102 and memory-efficient LLMs fine-tuning (Malladi et al., 2023; Zhang et al., 2024). Despite these
103 strengths, the practical application of ZO is primarily limited to smaller-scale tasks and models. A
104 critical limitation stems from the high error in its gradient approximations (Park et al., 2025), which
105 becomes more pronounced as problems grow larger and more complex, making scaling particularly
106 challenging. To address this issue, approaches such as MeZO-SVRG (Gautam et al., 2024) and
107 DiZO (Tan et al., 2025) utilize variance-reduction methodologies (Ma & Huang, 2025) to mitigate
gradient divergence. Furthermore, methods including SparseMezo (Liu et al., 2024), TeZO (Sun et al.,

2025), and AdaZeta (Yang et al., 2024) have been proposed to diminish approximation errors by reducing dependence on the parameter dimension through parameter sparsification and tensorization. Subspace methods (Nozawa et al., 2025), including SubZero (Yu et al., 2024) and LOZO (Chen et al., 2025), are explored to leverage low-rank structures for decreasing the error. Although they effectively alleviate the variance of gradient approximation, the randomly generated projection matrices cannot precisely reflect the transformation between the subspace and the full space, leading to degradation in model performance.

3 PRESTUDY

In exploring the alignment between estimated ZO and true FO gradients in the parameter spaces of large language models, we perform a targeted analysis using the OPT-1.3B model (Zhang et al., 2022) on the RTE task (Dagan et al., 2005; Bar Haim et al., 2006; Giampiccolo et al., 2007; Bentivogli et al., 2009). For every 50 training steps, we determine the exact FO gradients through backpropagation with a batch size of 16 and ZO gradient estimates via MeZO’s simultaneous perturbation method (Malladi et al., 2023). Subsequently, we apply singular value decomposition (SVD) to both gradient matrices. We then assess the cosine similarity between the singular value vectors.

Figure 1 illustrates that the singular vectors demonstrate high cosine similarity. This finding indicates that the ZO gradients maintain critical optimization directions and exhibit a similar low-rank structure. This supports our main hypothesis that ZO gradient estimates contain sufficient spectral information to reconstruct low-rank subspaces guided by FO methods. The preserved accuracy in directional estimates suggests that by limiting ZO perturbations to the primary gradient subspaces, we can reduce approximation variance while still achieving effective updates. These concepts form the foundation of our SVD-0 optimization framework, which systematically leverages the inherent structure in ZO gradient estimates to achieve FO-guided efficiency without the computational overhead associated with backpropagation. Additional experiments can be found in Appendix B.

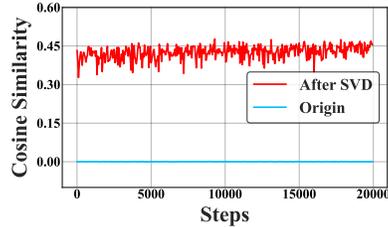


Figure 1: Cosine similarities between estimated ZO gradients and true gradients for “Origin” and “After SVD”.

4 METHODOLOGY

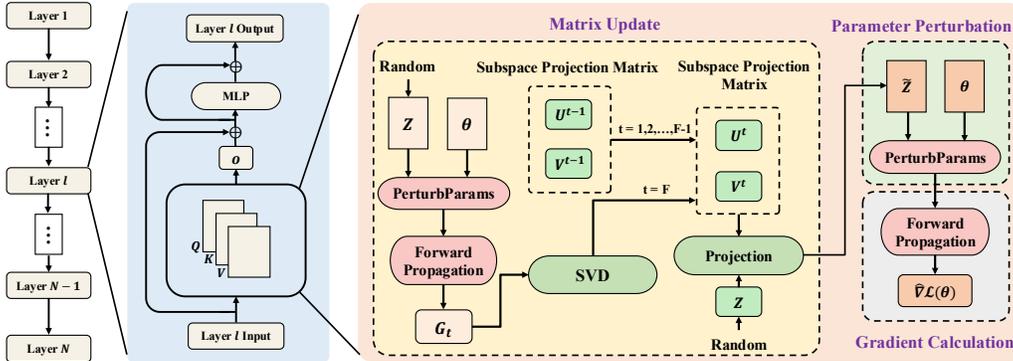


Figure 2: Framework and workflow of our SVD-0 method.

Figure 2 illustrates our approach, which focuses on two main components: the matrix update module and the parameter perturbation module. The matrix update module is for computing and adjusting the projection matrices, represented as $U \in \mathbb{R}^{m \times r}$ and $V \in \mathbb{R}^{n \times r}$. Together with a low-dimensional random matrix $Z \in \mathbb{R}^{r \times r}$, these matrices are used to generate a low-rank perturbation \tilde{Z} .

Within the first module (i.e., the matrix update module), we introduce an innovative and precise approach to acquire the matrices U and V , as detailed in Algorithm 1. Traditional approaches often utilize random low-rank perturbation matrices (Chen et al., 2025; Yu et al., 2024). This randomness contributed to uncertainty in the gradient update process during training. In contrast,

our approach computes the U and V matrices based on the gradient information derived using the MeZO method (Malladi et al., 2023) before each update.

The second module serves to perturb the parameters, as described in Algorithm 2. Common enhancements, such as SubZero (Yu et al., 2024) and the SVD-0 approach proposed here, reformulate the update mechanism by adopting a low-rank perturbation method. As illustrated in Figure 2, the low-rank perturbation $\tilde{Z} \in \mathbb{R}^{m \times n}$ is determined in the following manner:

$$\tilde{Z} = UZV^T, \quad (1)$$

where $Z \in \mathbb{R}^{r \times r}$ is a random perturbation matrix sampled from $\mathcal{N}(0, 1)$. Consequently, the parameter $\theta_t \in \mathbb{R}^{m \times n}$ during the t^{th} iteration is determined by $\theta_t^\pm = \theta \pm \tilde{Z} = \theta \pm UZV^T$. Thus, the gradient is approximated using two forward evaluations as expressed below:

$$\hat{\nabla} \mathcal{L}(\theta_t^\pm) = \frac{\mathcal{L}(\theta_t^+; \mathcal{B}) - \mathcal{L}(\theta_t^-; \mathcal{B})}{2\epsilon} UZV^T. \quad (2)$$

4.1 GRADIENT-GUIDED SUBSPACE PROJECTION MATRIX ACQUISITION

Existing approaches to projection matrix construction consist of a spectrum of techniques, ranging from randomized sampling methods (Chen et al., 2025; Yu et al., 2024) to computationally intensive deterministic algorithms (Zhao et al., 2024b). Although the former is computationally efficient, it has the drawback of insufficient approximation accuracy due to its reliance on randomness. The latter introduces significant computational overhead while not significantly improving the approximation accuracy. To address this limitation, we propose a balancing strategy based on adaptive subspace decomposition, as shown in lines 4-7 of Algorithm 3.

To retain the advantage of memory efficiency of zero-order optimizations, we calculate the gradient using the MeZO (Malladi et al., 2023) method, as shown in lines 5-6 of Algorithm 3. Before calculating the projection matrix each time, the gradient calculation is required. Then, as shown in the algorithm 1, the U and V matrices are updated according to the gradient obtained this time. We use the SVD method to calculate the projection matrix. Through this method, the original gradient is projected onto a compact space $\mathbf{R} \in \mathbb{R}^{r \times r}$: $\mathbf{R} = U^T G V$. After that, we can generate a low-rank perturbation Z in this space, as shown in lines 3-4 of the Algorithm 2, and then use the previously calculated U and V matrices to restore this low-rank perturbation to the original high-rank space. In this way, we can successfully apply gradient-based low-rank perturbations to the parameters, and this process introduces no additional overhead compared to the traditional ZO method (i.e., MeZO).

4.2 PERIODICAL SUBSPACE UPDATE

As mentioned above, we obtain the gradient using the MeZO (Malladi et al., 2023) method and then calculate the projection matrices U and V via SVD. These two projection matrices jointly determine the gradient approximation and the parameter update of the t^{th} step. However, this iterative update method presents a critical trade-off between computational efficiency and subspace adaptability. High-frequency updates restrict the complete evolution of the gradient subspace while incurring substantial computational costs, particularly due to the need for gradient recomputation before each projection matrix update. In contrast, low-frequency updates may fail to capture the dynamic variations in the gradient subspace throughout the training process.

Algorithm 1 GenerateProjMatrix(G, r)

Input: i) G , estimated gradient of parameter matrix; ii) r , rank.

Output: U, V , projection matrices.

- 1: $(P, S, Q) \leftarrow \text{SVD}(G)$
 - 2: $U \leftarrow P[:, :r]$
 - 3: $V \leftarrow Q[:, :r]$
 - 4: **return** U, V
-

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

Input: i) W , model parameter set; ii) \mathcal{U} and \mathcal{V} , projection matrix sets; iii) r , rank; iv) ϵ , perturbation scale; v) s , seed.

Output: Model parameter set after perturbation.

- 1: ResetGenerator(s)
 - 2: **for** $i = 1, 2, \dots, l$ **do**
 - 3: $Z_i \leftarrow \text{GeneratePerturbMatrix}(r)$
 - 4: $W_i \leftarrow W_i + \epsilon U_i Z_i V_i^T$
 - 5: **end for**
 - 6: **return** W
-

Therefore, we propose a periodic subspace update strategy. As presented in lines 4-10 in Algorithm 3, we use the MeZO method to calculate the gradient once at the start step and every F steps thereafter. Then the obtained gradient is used to update the projection matrices \mathbf{U} and \mathbf{V} , and keep them unchanged in the subsequent steps. We have experimentally proved the effectiveness and necessity of this strategy. As shown in Table 4, the appropriate update frequency can not only ensure efficiency but also bring significant improvements to model performance.

Algorithm 3 SVD-0

Input: i) $\mathbf{W}_i \in \mathbb{R}^{m_i \times n_i}$, $i = 1, \dots, l$, parameter matrix in the i -th layer; ii) \mathcal{L} , loss; iii) T , step budget; iv) ϵ , perturbation scale; v) $\{\eta^t\}$, learning rate schedule; vi) F , subspace update frequency; vii) r , rank.

```

1: for  $t = 1, \dots, T$  in parallel do
2:    $\mathcal{B}^t \leftarrow \text{SampleMinbatch}(s^t)$  {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:        $G_i \leftarrow \text{EstimateGradient}(\mathbf{W}_i^t, \epsilon)$  {Estimate the gradient of  $\mathbf{W}_i^t$  using MeZO}
6:        $\mathbf{U}_i^t, \mathbf{V}_i^t \leftarrow \text{GenerateProjMatrix}(G_i, r)$ 
7:     else
8:        $\mathbf{U}_i^t \leftarrow \mathbf{U}_i^{t-1}, \mathbf{V}_i^t \leftarrow \mathbf{V}_i^{t-1}$ 
9:     end if
10:  end for
11:   $\{\mathbf{W}^t = \{\mathbf{W}_i^t\}_{i=1}^l, \mathcal{U}^t = \{\mathbf{U}_i^t\}_{i=1}^l, \mathcal{V}^t = \{\mathbf{V}_i^t\}_{i=1}^l\}$ 
12:   $\ell_+^t \leftarrow \mathcal{L}(\mathbf{W}^t; \mathcal{B}^t)$ 
13:   $\mathbf{W}^t \leftarrow \text{PerturbParams}(\mathbf{W}^t, \mathcal{U}^t, \mathcal{V}^t, r, \epsilon, s^t)$ 
14:   $\ell_-^t \leftarrow \mathcal{L}(\mathbf{W}^t; \mathcal{B}^t)$ 
15:   $\mathbf{W}^t \leftarrow \text{PerturbParams}(\mathbf{W}^t, \mathcal{U}^t, \mathcal{V}^t, r, \epsilon, s^t)$ 
16:   $\rho^t \leftarrow (\ell_+^t - \ell_-^t) / (2\epsilon)$ 
17:   $\text{ResetGenerator}(s^t)$  {Reset random number generator with seed  $s^t$ }
18:  for  $i = 1, 2, \dots, l$  do
19:     $\mathbf{Z}_i^t \leftarrow \text{GeneratePerturbMatrix}(r)$  {Regenerate the perturbation matrix  $\mathbf{Z}_i^t \in \mathbb{R}^{r \times r}$  whose entries are sampled from  $\mathcal{N}(0, 1)$ }
20:     $\mathbf{W}_i^{t+1} \leftarrow \mathbf{W}_i^t - \eta^t \rho^t (\mathbf{U}_i^t \mathbf{Z}_i^t \mathbf{V}_i^{tT})$ 
21:  end for
22: end for

```

Table 1: Computational cost (in minutes) comparison.

Method	WIC	ReCoRD	FiQA-SA	TFNS	MultiRC
MeZO (Malladi et al., 2023)	114.7	211.4	71.5	113.3	1125.3
SubZero (Yu et al., 2024)	114.5	207.5	71.3	124.9	1149.5
SVD-0	116.1	220.6	76.3	116.2	1154.8

Table 2: Memory cost comparison.

Method	RoBERTa-large	OPT-1.3B	OPT-13B
MeZO (Malladi et al., 2023)	2.042GB	4.732GB	27.693GB
LOZO (Chen et al., 2025)	2.042GB	4.732GB	27.789GB
SVD-0	2.562GB	4.891GB	28.767GB

As shown in Table 1, we compared two representative ZO variants (i.e., MeZO (Malladi et al., 2023) and SubZero (Yu et al., 2024)) and our SVD-0 method. The WIC (Pilehvar & Camacho-Collados, 2019) and ReCoRD (Zhang et al., 2018) datasets were paired with the OPT-1.3B model, the FiQA-SA (Maia et al., 2018) and TFNS (Magic, 2022) datasets were paired with the Qwen-1.8B model, whereas the MultiRC (Khashabi et al., 2018) dataset was paired with the OPT-13B model. Each model was fine-tuned on its respective task for 20,000 steps. The findings show that SVD-0 requires a marginally longer training period, approximately 7% longer than the two ZO variants. It’s worth mentioning that the time complexity for SVD processes remains at $O(n^3)$, where n is the matrix’s dimension, ensuring that the upper bound of the training time complexity remains unchanged. In practice, the extra time required by SVD operations is negligible compared to the benefits gained in classification, multiple-choice, and generation tasks.

Despite reducing computational complexity, this strategy will result in minimal additional memory usage, as shown in Table 2. We adopt a layer-wise parameter update strategy, where we update only the parameters of a specific layer of the model simultaneously. This means that during the entire training process, we only need to store two additional small matrices at the same time, including the projection matrices $\mathbf{U} \in \mathbb{R}^{m \times r}$ and $\mathbf{V} \in \mathbb{R}^{n \times r}$, where r is much smaller than the dimension

of the parameter matrix $\theta \in \mathbb{R}^{m \times n}$. Therefore, the memory usage introduced by the two matrices remains at the same low level as that introduced in (Yu et al., 2024). This strategy makes our method almost consistent with the memory required by the MeZO (Malladi et al., 2023) method without any performance loss, and maintains the memory-saving advantage of the ZO method.

5 CONVERGENCE ANALYSIS

In this section, we analyze the convergence of our proposed SVD-0. Following the derivations in (Yu et al., 2024; Nozawa et al., 2025) and (Zhao et al., 2024a), we first present our proposition along with the corresponding lemma.

Lemma 1. (Low-rank subspace of weight matrices (Zhao et al., 2024a)). *Gradient matrices become low-rank during fine-tuning. The weight matrix update can be formed as:*

$$\theta_T = \theta_0 + \eta \sum_{t=0}^{T-1} \tilde{\nabla} f(\theta_t), \quad \tilde{\nabla} f(\theta_t) = \mathbf{U}_t (\mathbf{U}_t^\top \nabla f(\theta_t) \mathbf{V}_t) \mathbf{V}_t^\top, \quad (3)$$

where η is the learning rate, $\mathbf{U}_t \in \mathbb{R}^{m \times r}$ and $\mathbf{V}_t \in \mathbb{R}^{n \times r}$ are projection matrices and can be approximated by the spectrum of $\nabla f(\theta_t)$ through $(\mathbf{U}, \mathbf{V}) = \text{SVD}(\nabla f(\theta_t))$.

Lemma 1 shows that subspace projection matrices can be approximated by adopting SVD on gradients. Given that the SPSA is an unbiased approximation of the exact gradient $\nabla f(\theta)$, we can use the SPSA gradient to compute the two projection matrices.

Proposition 1. (Block-diagonal matrix based on SVD). *The singular matrices \mathbf{U} and \mathbf{V} are column-orthogonal. Therefore, we can similarly define the following notations based on Equation 1:*

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

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

Proposition 1 indicates that the projection matrices in our method exhibit the same properties as the column-orthogonal matrices discussed in (Yu et al., 2024). Consequently, the subsequent theoretical analysis can follow the same approach as that demonstrated in (Yu et al., 2024).

Lemma 2. (Bounded gradient estimation error (Yu et al., 2024)). *For the gradient estimation in Equation 2, the following two properties hold.*

i) *By using gradient estimation in Equation 2, the estimated gradient $\hat{\nabla} f(\theta)$ is equivalent to:*

$$\hat{\nabla} f(\theta) = \frac{f(\theta + \varepsilon \mathbf{P} \mathbf{z}) - f(\theta - \varepsilon \mathbf{P} \mathbf{z})}{2\varepsilon} \mathbf{P} \mathbf{z}, \quad (4)$$

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 $q = lr^2$.

ii) *Let $\mathbf{z} \sim \mathcal{N}(\mathbf{0}, \mathbf{I}_q)$, and $f \in C_{L_2}^{2,2}(\mathbb{R}^d)$. Based on Equation 4 whose properties have been analyzed in (Nozawa et al., 2025), our method has the same bounded gradient estimation error as that in (Yu et al., 2024):*

$$\left\| \mathbb{E}_{\mathbf{z}} \left[\hat{\nabla} f(\theta) \right] - \mathbf{P} \mathbf{P}^\top \nabla f(\theta) \right\|_2 \leq \frac{\varepsilon^2}{6} L_2 (q + 4)^2. \quad (5)$$

Note that $f \in C_L^{s,p}(S)$ denotes the class of s -th smooth and p -th L -smooth functions over the set S .

Theorem 1. (Convergence of SVD-0). *Consider the optimization problem $\mathbf{x}^* = \arg \min_{\mathbf{x} \in \mathbb{R}^d} f(\mathbf{x})$,*

in which $f \in C_{L_1}^{1,1}(\mathbb{R}^d)$ and f exhibits non-convex behavior. Define the stochastic sequence $\mathcal{E}_k = (\mathbf{z}_0, \mathbf{z}_1, \dots, \mathbf{z}_k)$, where each \mathbf{z}_k follows the normal distribution $\mathcal{N}(\mathbf{0}, \mathbf{I}_q)$. Set the step-size parameter as $\eta = \frac{1}{4(q+4)L_1}$. Let $\{\mathbf{x}_k\}_{k>0}$ denote the iterates produced via Algorithm 3. For SVD-0, we establish its convergence rate as:

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

324 Under the scaling $T = \Omega\left(\frac{d}{\epsilon^2}\right)$ for $\epsilon \leq \mathcal{O}\left(\frac{1}{q^{3/2}d^{1/2}L_1^{3/2}}\right)$, this aligns with prior theoretical
 325 derivations.
 326
 327

328 By combining Proposition 1 and Lemma 2 within the framework proposed in (Yu et al., 2024),
 329 Theorem 1 demonstrates that our SVD-0 achieves a convergence rate of $\mathcal{O}\left(\sqrt{\frac{d}{T}}\right)$, matching the rate
 330 derived in (Yu et al., 2024). For a more detailed explanation, please see Appendix D.
 331

332 6 EXPERIMENTS

333 To evaluate the effectiveness of our approach, we implemented SVD-0 using the PyTorch framework
 334 (version 20.10). All experiments were conducted on a Linux workstation equipped with CentOS,
 335 featuring two NVIDIA A100-40GB GPUs, dual Intel Xeon 6240R CPUs, and 384GB of RAM. The
 336 following presents the dataset settings and ZO baselines used in the experiments. Please refer to
 337 Appendix A for our detailed model settings.
 338

339 **Dataset Settings.** For OPT models, we experimented with the SuperGLUE benchmark (Wang et al.,
 340 2019), which consists of various types of tasks, including classification tasks (e.g., SST-2 (Socher
 341 et al., 2013), RTE (Bar Haim et al., 2006; Bentivogli et al., 2009; Dagan et al., 2005; Giampiccolo
 342 et al., 2007), CB (de Marneffe et al., 2019), BoolQ (Clark et al., 2019), WSC (Levesque et al.,
 343 2012), WIC (Pilehvar & Camacho-Collados, 2019), and MultiRC (Khashabi et al., 2018)), multiple
 344 choice tasks (e.g., COPA (Roemmele et al., 2011) and ReCoRD (Zhang et al., 2018)), and generation
 345 tasks (e.g., SQuAD (Rajpurkar et al., 2016) and DROP (Dua et al., 2019)). Here, for each task,
 346 we randomly selected 1000 samples for training, 500 samples for validation, and 1000 samples for
 347 testing. For the RoBERTa-large model, in addition to the task SST-2, we investigated three more
 348 tasks, i.e., SST-5 (Socher et al., 2013), SNLI (Bowman et al., 2015), and MNLI (Williams et al.,
 349 2018). In this case, we fixed the parameter k at 512 throughout the training and validation phases,
 350 indicating that 512 samples are allocated for each category. For the testing phase, we randomly chose
 351 a total of 1000 samples.
 352

353 **ZO Baselines.** Our SVD-0 method was evaluated against six latest ZO optimization algorithms,
 354 i.e., MeZO (Malladi et al., 2023), ZO-AdaMU (Jiang et al., 2024), S-MeZO (Liu et al., 2024),
 355 SubZero (Yu et al., 2024), LOZO (Chen et al., 2025), and HiZOO (Zhao et al., 2024b). Meanwhile,
 356 we examined three memory-efficient inference-only approaches, i.e., zero-shot evaluation, in-context
 357 learning (ICL) (Brown et al., 2020), and linear probing (LP) (Kumar et al., 2022).

358 We designed our experiments to explore the following research questions (RQs).

359 **RQ1 (Superiority of SVD-0):** To what extent does SVD-0 outperform SOTA methods in accuracy?

360 **RQ2 (Impact of Hyperparameters):** What are the impacts of critical hyperparameters (e.g., learning
 361 rate, subspace rank, subspace update frequency) on SVD-0-based fine-tuning?

362 **RQ3 (Applicability of SVD-0):** How does SVD-0 perform when fine-tuning models of varying sizes
 363 or architectures (e.g., masked or causal language models)?

364 6.1 COMPARISON WITH STATE-OF-THE-ARTS (RQ1)

365 We compared our proposed SVD-0 method with the SOTA ZO optimizers. The experiments were
 366 conducted on the SuperGLUE benchmark employing both the OPT-13B and OPT-1.3B language
 367 models of different sizes. Note that in each experiment, we applied the adopted stochastic gradient
 368 descent (SGD) or ZO method to all model parameters.
 369

370 Table 3 compares the fine-tuning performance of the OPT-13B model on SuperGLUE benchmark
 371 tasks. Here, we considered three types of fine-tuning methods: i) the traditional fine-tuning method
 372 (i.e., SGD) with backpropagation; ii) inference-only methods (i.e., Zero-shot, ICL, and LP) without
 373 fine-tuning; and iii) memory-efficient ZO-based methods. To enable a fair comparison between ZO-
 374 based methods, we used the MeZO method here as a reference. We evaluated the overall performance
 375 across each classification task category and denoted the improvement in performance compared to
 376 the baseline (i.e., MeZO) in the sub-column labeled “Total”. For example, the total performance
 377 on multiple choice tasks with MeZO and SVD-0 is 169.0 and 171.2, respectively. In this case,
 SVD-0 improves inference performance by 1.30% compared to MeZO. From the results provided in

Table 3: Comparison of OPT-1.3B fine-tuning performance (%) on SuperGLUE, where the best results are presented in **bold** and the second-best results are highlighted with underlines.

Method	Classification Task							Multiple Choice Task		Generation Task		All Task			
	SST-2	RTE	CB	BoolQ	WSC	WIC	MultiRC	Total	COPA	ReCoRD	Total	SQuAD	DROP	Total	
SGD	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 (Brown et al., 2020)	87.0	62.1	57.1	66.9	39.4	50.5	53.1	-	87.0	82.5	-	75.9	29.6	-	-
LP (Kumar et al., 2022)	93.4	68.6	67.9	59.3	63.5	60.2	63.5	-	55.0	27.1	-	3.7	11.1	-	-
MeZO (Malladi et al., 2023)	92.1	71.5	71.4	74.4	61.5	60.0	60.1	0%	87.0	82.0	0%	84.2	31.2	0%	0%
ZO-AdaMU (Jiang et al., 2024)	92.1	72.9	67.9	73.0	61.5	60.7	63.0	0.02%	89.0	83.0	1.78%	82.4	32.0	-0.87%	0.27%
S-MeZO (Liu et al., 2024)	92.3	76.9	75.0	76.5	61.1	58.2	63.3	2.51%	87.0	71.2	-6.39%	77.9	31.9	-4.85%	-0.53%
HiZOO (Zhao et al., 2024b)	91.3	69.3	69.4	67.3	<u>63.5</u>	59.4	55.5	-3.12%	88.0	81.4	0.24%	81.9	31.3	-1.91%	-2.21%
LOZO (Chen et al., 2025)	91.7	70.4	69.6	71.9	<u>63.5</u>	60.8	63.0	-0.02%	89.0	81.3	0.77%	84.9	30.7	0.17%	0.18%
SubZero (Yu et al., 2024)	92.1	74.0	<u>73.2</u>	<u>75.3</u>	65.4	60.8	61.0	2.20%	88.0	82.3	0.77%	84.5	32.0	0.95%	<u>1.70%</u>
SVD-0	93.6	<u>75.5</u>	71.4	75.2	<u>63.5</u>	65.4	60.6	2.89%	89.0	82.2	<u>1.30%</u>	85.1	30.9	<u>0.52%</u>	2.19%

the ‘‘Total’’ sub-columns, we can find that SVD-0 can always achieve top-2 inference performance. Furthermore, we used the final column to show the relative performance improvement for all tasks. From this column, we can find that SVD-0 achieves the best overall performance. Interestingly, while S-MeZO matches SVD-0 in the number of tasks where it excels, its overall performance, shown in the final column, is noticeably inferior to SVD-0 and even falls short of the reference (i.e., MeZO).

6.2 IMPACTS OF HYPERPARAMETERS (RQ2)

In this experiment, we investigate three key hyperparameters (i.e., subspace update frequency, rank, and learning rate) to evaluate their impacts on fine-tuning performance.

Table 4: Impact of subspace update frequency, where the best results are highlighted in **bold**.

Frequency	SST-2	RTE	CB	BoolQ	WSC	WIC	MultiRC	COPA	ReCoRD	SQuAD	DROP
50	90.5	57.0	64.3	65.0	63.5	55.6	57.5	72.0	72.4	74.2	23.0
500	89.5	55.6	69.6	64.1	63.5	53.9	58.1	73.0	72.2	74.3	22.9
1000	90.6	58.5	71.4	65.2	63.5	56.4	58.2	73.0	72.1	73.7	24.0
2000	89.2	56.7	73.2	64.5	62.5	57.4	58.1	73.0	71.7	72.6	23.8
20000	89.8	56.3	71.4	65.3	62.5	57.5	58.2	72.0	72.1	72.6	22.6

For the subspace update frequency F , our goal is to evaluate the impact of varying this frequency on model performance across different tasks. We conducted experiments based on the OPT-1.3B model, with a fixed rank of $r = 24$ and a learning rate of $1e-7$. In this analysis, we evaluated five frequencies at varying magnitudes, specifically selected from the set $\{50, 500, 1000, 2000, 20000\}$. Table 4 provides the experimental results. From this table, we can find that when the frequency is set to 1000 (i.e., the subspace is updated every 1000 steps), SVD-0 achieves the best performance in six of the eleven tasks. Note that SVD-0-based fine-tuning is not sensitive to the hyperparameter F . Therefore, we suggest setting F to 1000 by default for fine-tuning.

We investigated the rank of hyperspace (i.e., r) and the learning rate in tandem. Table 5 presents the fine-tuning performance under various combinations of these two hyperparameters, where the rank is selected from $\{2, 24, 48, 64, 128\}$ and the learning rate is selected from $\{1e-7, 5e-7, 1e-6\}$. All the experimental results are collected based on the SST-2 task using the OPT-1.3B model, with a fixed subspace update frequency of 1000. This table shows that the fine-tuning performance is weak when the rank is low (i.e., $r = 2$). While elevating the rank can enhance fine-tuning performance, the extent of this enhancement becomes negligible once the rank surpasses 24. At low ranks, the performance can vary significantly with different learning rates. In contrast, increasing rank tends to reduce this variability in performance. Moreover, we observe a similar trend for the learning rate hyperparameter, where setting the learning rate to $5e-7$ achieves the best performance for most rank settings. However, increasing the learning rates can lead to a decline in inference performance.

Table 5: Impacts of rank and learning rate on inference.

Rank \ LR	1e-7	5e-7	1e-6
2	87.7	91.2	86.7
24	90.6	92.2	90.3
48	89.5	91.6	90.1
64	89.9	90.4	91.6
128	90.0	91.3	90.6

6.3 IMPACT OF MODEL SIZES AND ARCHITECTURES (RQ3)

In Table 3, we have evaluated the adaptability of SVD-0 to large-scale LLMs. To further validate the generalizability of our approach, we extended our evaluation to the OPT-1.3B model, using

representative tasks of different types. These tasks include SST-2 and WIC, which are classification tasks, ReCoRD, a multiple-choice task, and SQuAD, a generation task. Table 6 presents the results of the comparison between four ZO-based fine-tuning methods, where the last column shows the average fine-tuning performance of the four tasks. According to this table, we can see that SVD-0 is also well-suited for fine-tuning on small-scale LLMs. Although LOZO delivers the highest performance in this experiment, the difference in average fine-tuning performance between SVD-0 and LOZO is minimal (i.e., only 0.2%). Note that SVD-0 achieves better performance than MeZO, the reference method, while SubZero fails to beat MeZO. Moreover, SVD-0 can consistently outperform its counterpart (i.e., SubZero) by an average of 0.7%. All these observations substantiate the efficiency of our method in enhancing subspaces for optimizing LLMs.

Table 6: Fine-tuning performance (%) comparison for OPT-1.3B, where the top-2 results are marked in bold and with underlines, respectively.

Method	SST-2	WIC	ReCoRD	SQuAD	AVG.
MeZO (Malladi et al., 2023)	91.7	<u>61.1</u>	<u>72.2</u>	77.4	75.6
LOZO (Chen et al., 2025)	93.2	62.4	71.9	78.1	76.4
SubZero (Yu et al., 2024)	91.9	60.7	72.0	<u>77.6</u>	75.5
SVD-0 (Ours)	<u>93.0</u>	<u>61.1</u>	73.0	<u>77.6</u>	<u>76.2</u>

Table 7: Fine-tuning performance (%) comparison for RoBERTa-large, where the top-2 results are marked in bold and with underlines, respectively.

Method	SST-2	SST-5	SNLI	MNLI
Zero-shot	79.0	35.5	50.2	48.8
MeZO (Malladi et al., 2023)	93.7 (0.4)	53.9 (1.9)	84.8 (1.1)	76.6 (0.8)
LOZO (Chen et al., 2025)	<u>94.1 (0.7)</u>	<u>53.0 (0.4)</u>	85.4 (0.8)	80.4 (1.0)
SVD-0 (Ours)	<u>94.4 (0.7)</u>	<u>54.4 (0.7)</u>	<u>85.4 (1.3)</u>	<u>80.4 (1.5)</u>

We investigated the fine-tuning performance of different optimization methods on RoBERTa-large, where we considered four downstream tasks, including two sentiment classification tasks (i.e., SST-2 and SST-5) and two natural language inference tasks (i.e., SNLI and MNLI). For a fair comparison, like the work in (Chen et al., 2025), we performed fine-tuning on each task five times using different random seeds. Table 7 presents the experimental results, reflecting both the average inference performance and its standard deviation (indicated in parentheses) for each combination of fine-tuning methods and tasks. From this table, we can see that SVD-0 performs the best compared to SOTA ZO optimization methods, demonstrating the adaptability of our approach to various model architectures.

Table 8: Fine-tuning performance (%) comparison for Qwen-1.8B, where the top-2 results are marked in bold and with underlines, respectively.

Method	SST-2	WIC	ReCoRD	Total
MeZO (Malladi et al., 2023)	78.3	55.6	64.8	0%
LOZO (Chen et al., 2025)	81.7	55.2	64.8	1.51%
SubZero (Yu et al., 2024)	80.8	<u>56.7</u>	<u>65.2</u>	<u>2.01%</u>
SVD-0	82.2	57.2	65.3	3.02%

Table 9: Fine-tuning performance (%) comparison for OPT-1.3B on financial datasets, where the top-2 results are marked in bold and with underlines, respectively.

Method	FPB	FIQA-SA	TFNS	NWGI	Total
MeZO (Malladi et al., 2023)	65.3	81.4	74.7	48.5	0%
LOZO (Chen et al., 2025)	61.3	85.1	71.6	53.7	0.67%
SubZero (Yu et al., 2024)	<u>66.4</u>	<u>84.0</u>	78.4	49.7	<u>3.19%</u>
SVD-0	74.1	<u>84.0</u>	<u>76.3</u>	<u>52.8</u>	6.41%

To further validate the generalization ability of our method on cutting-edge models, we conducted experiments based on the Qwen-1.8B model. We exclusively compared our method against the baseline (i.e., MeZO (Malladi et al., 2023)) and the two most recent ZO baseline techniques (i.e., LOZO (Chen et al., 2025) and SubZero (Yu et al., 2024)). Table 8 shows that our approach still achieves the best performance, indicating the adaptability and generalizability of our method in cutting-edge models. Moreover, we assessed SVD-0 on datasets derived from four financial sentiment analysis benchmarks, including FPB (Malo et al., 2014), FIQA-SA (Maia et al., 2018), TFNS (Magic, 2022), and NWGI (Yang, 2023). As shown in Table 9, SVD-0 achieves the highest total performance, demonstrating the method’s reliability and efficiency across various domains and task types.

6.4 DISCUSSION

Limitations. While the SVD-0 technique improves the ZO subspace fine-tuning approach, the accuracy of the subspace projection matrices is significantly influenced by the precision of the ZO gradients. In smaller models, such as the OPT-1.3B, the ZO gradients may have a greater approximation error, which can result in decreased precision in obtaining the projection matrices.

Border Impacts. In this paper, we introduced a new approach to derive more precise projection matrices, which can be used to improve the effectiveness of ZO subspace fine-tuning techniques for LLMs. Our method utilizes SVD on ZO gradients to extract projection matrices, eliminating the need for the memory-demanding FO gradients. Our theoretical convergence analysis, in conjunction with the experimental findings, demonstrates that our research makes a positive contribution to the advancement of memory-efficient fine-tuning methods for LLMs.

486 7 CONCLUSION

487
488 Although various zeroth-order (ZO) optimization methods have been proposed to enable memory-
489 efficient fine-tuning for large language models (LLMs), due to the use of random subspaces, most
490 of them suffer from inaccurate gradient estimation, resulting in inferior training performance. To
491 address this problem, this paper presents a novel ZO subspace fine-tuning method named SVD-0. By
492 precisely capturing fine-tuning subspaces, SVD-0 enables the construction of projection matrices
493 with higher accuracy, thereby achieving more accurate gradient estimation and improving the LLM
494 fine-tuning performance. Extensive experimental findings demonstrate the efficacy of SVD-0 in
495 dealing with complex language modeling tasks. In the future, we plan to integrate our SVD-0 method
496 with parameter quantization techniques to reduce the memory requirements of LLM fine-tuning.

497 REFERENCES

- 498
499 Roy Bar Haim, Ido Dagan, Bill Dolan, Lisa Ferro, Danilo Giampiccolo, Bernardo Magnini, and Idan
500 Szpektor. The second PASCAL recognising textual entailment challenge. In *Proceedings of the*
501 *Second PASCAL Challenges Workshop on Recognising Textual Entailment*, 2006.
- 502
503 Luisa Bentivogli, Peter Clark, Ido Dagan, and Danilo Giampiccolo. The fifth PASCAL recognizing
504 textual entailment challenge. In *Proceedings of the Second Text Analysis Conference*, 2009.
- 505
506 Samuel R. Bowman, Gabor Angeli, Christopher Potts, and Christopher D. Manning. A large annotated
507 corpus for learning natural language inference. In *Proceedings of the Conference on Empirical*
508 *Methods in Natural Language Processing (EMNLP)*, pp. 632–642, 2015.
- 509
510 Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal,
511 Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, et al. Language models
512 are few-shot learners. *Proceedings of the Advances on Neural Information Processing Systems*
513 *(NeurIPS)*, 33:1877–1901, 2020.
- 514
515 Yiming Chen, yuan zhang, Liyuan Cao, Kun Yuan, and Zaiwen Wen. Enhancing zeroth-order fine-
516 tuning for language models with low-rank structures. In *Proceedings of International Conference*
on Learning Representations (ICLR), 2025.
- 517
518 Christopher Clark, Kenton Lee, Ming-Wei Chang, Tom Kwiatkowski, Michael Collins, and Kristina
519 Toutanova. BoolQ: Exploring the surprising difficulty of natural yes/no questions. In *Proceedings*
520 *of the Conference of the North American Chapter of the Association for Computational Linguistics:*
521 *Human Language Technologies (NAACL)*, 2019.
- 522
523 Ido Dagan, Oren Glickman, and Bernardo Magnini. The PASCAL recognising textual entailment
524 challenge. In *Proceedings of the International Conference on Machine Learning Challenges: Eval-*
525 *uating Predictive Uncertainty Visual Object Classification, and Recognizing Textual Entailment*,
2005.
- 526
527 Marie-Catherine de Marneffe, Mandy Simons, and Judith Tonhauser. The commitmentbank: Inves-
528 tigating projection in naturally occurring discourse. In *Proceedings of Sinn und Bedeutung 23*,
2019.
- 529
530 Tim Dettmers, Mike Lewis, Younes Belkada, and Luke Zettlemoyer. Llm. int8 () 8-bit matrix
531 multiplication for transformers at scale. In *Proceedings of the Advances on Neural Information*
532 *Processing Systems (NeurIPS)*, pp. 30318–30332, 2022.
- 533
534 Tim Dettmers, Artidoro Pagnoni, Ari Holtzman, and Luke Zettlemoyer. Qlora: Efficient finetuning of
535 quantized llms. *Proceedings of the Advances in Neural Information Processing Systems (NeurIPS)*,
36:10088–10115, 2023.
- 536
537 Dheeru Dua, Yizhong Wang, Pradeep Dasigi, Gabriel Stanovsky, Sameer Singh, and Matt Gardner.
538 DROP: A reading comprehension benchmark requiring discrete reasoning over paragraphs. In
539 *Proceedings of the Conference of the North American Chapter of the Association for Computational*
Linguistics: Human Language Technologies (NAACL), pp. 2368–2378, 2019.

- 540 Elias Frantar, Saleh Ashkboos, Torsten Hoefler, and Dan Alistarh. Gptq: Accurate post-training
541 quantization for generative pre-trained transformers. *arXiv preprint arXiv:2210.17323*, 2022.
542
- 543 Tanmay Gautam, Youngsuk Park, Hao Zhou, Parameswaran Raman, and Wooseok Ha. Variance-
544 reduced zeroth-order methods for fine-tuning language models. In *Proceedings of the International
545 Conference on Machine Learning (ICML)*, pp. 15180–15208, 2024.
- 546 Danilo Giampiccolo, Bernardo Magnini, Ido Dagan, and Bill Dolan. The third PASCAL recognizing
547 textual entailment challenge. In *Proceedings of the ACL-PASCAL Workshop on Textual Entailment
548 and Paraphrasing*, 2007.
- 549 Daya Guo, Dejian Yang, Haowei Zhang, Junxiao Song, Ruoyu Zhang, Runxin Xu, Qihao Zhu,
550 Shirong Ma, Peiyi Wang, Xiao Bi, et al. Deepseek-r1: Incentivizing reasoning capability in llms
551 via reinforcement learning. *arXiv preprint arXiv:2501.12948*, 2025.
552
- 553 Zeyu Han, Chao Gao, Jinyang Liu, Jeff Zhang, and Sai Qian Zhang. Parameter-efficient fine-tuning
554 for large models: A comprehensive survey. *Transactions on Machine Learning Research (TMLR)*,
555 2024. ISSN 2835-8856.
- 556 Yongchang Hao, Yanshuai Cao, and Lili Mou. Flora: low-rank adapters are secretly gradient
557 compressors. In *Proceedings of the International Conference on Machine Learning (ICML)*, pp.
558 17554–17571, 2024.
- 559 Edward J. Hu, Yelong Shen, Phillip Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang,
560 and Weizhu Chen. Lora: Low-rank adaptation of large language models. In *Proceedings of the
561 International Conference on Learning Representations (ICLR)*, 2022.
562
- 563 Jia-Hong Huang, Yixian Shen, Hongyi Zhu, Stevan Rudinac, and Evangelos Kanoulas. Gradient
564 weight-normalized low-rank projection for efficient llm training. In *Proceedings of the AAAI
565 Conference on Artificial Intelligence (AAAI)*, volume 39, pp. 24123–24131, 2025.
- 566 Shuoran Jiang, Qingcai Chen, Youcheng Pan, Yang Xiang, Yukang Lin, Xiangping Wu, Chuanyi Liu,
567 and Xiaobao Song. Zo-adamu optimizer: Adapting perturbation by the momentum and uncertainty
568 in zeroth-order optimization. In *Proceedings of the AAAI Conference on Artificial Intelligence
569 (AAAI)*, volume 38, pp. 18363–18371, 2024.
570
- 571 Daniel Khashabi, Snigdha Chaturvedi, Michael Roth, Shyam Upadhyay, and Dan Roth. Looking
572 beyond the surface: A challenge set for reading comprehension over multiple sentences. In
573 *Proceedings of the Conference of the North American Chapter of the Association for Computational
574 Linguistics: Human Language Technologies*, pp. 252–262, 2018.
- 575 Diederik P Kingma and Jimmy Ba. Adam: A method for stochastic optimization. *arXiv preprint
576 arXiv:1412.6980*, 2014.
- 577 Ananya Kumar, Aditi Raghunathan, Robbie Jones, Tengyu Ma, and Percy Liang. Fine-tuning can
578 distort pretrained features and underperform out-of-distribution. *arXiv preprint arXiv:2202.10054*,
579 2022.
580
- 581 Hector Levesque, Ernest Davis, and Leora Morgenstern. The winograd schema challenge. In
582 *Proceedings of the International Conference on the Principles of Knowledge Representation and
583 Reasoning*, 2012.
- 584 Ji Lin, Jiaming Tang, Haotian Tang, Shang Yang, Wei-Ming Chen, Wei-Chen Wang, Guangxuan
585 Xiao, Xingyu Dang, Chuang Gan, and Song Han. Awq: Activation-aware weight quantization for
586 on-device llm compression and acceleration. *Proceedings of the Machine Learning and Systems
587 (MLSys)*, 6:87–100, 2024.
- 588 Xiao Liu, Kaixuan Ji, Yicheng Fu, Weng Tam, Zhengxiao Du, Zhilin Yang, and Jie Tang. P-tuning:
589 Prompt tuning can be comparable to fine-tuning across scales and tasks. In *Proceedings of the
590 Annual Meeting of the Association for Computational Linguistics (ACL)*, pp. 61–68, 2022.
591
- 592 Yinhan Liu, Myle Ott, Naman Goyal, Jingfei Du, Mandar Joshi, Danqi Chen, Omer Levy, Mike Lewis,
593 Luke Zettlemoyer, and Veselin Stoyanov. RoBERTa: A robustly optimized BERT pretraining
approach. *arXiv:1907.11692*, 2019.

- 594 Yong Liu, Zirui Zhu, Chaoyu Gong, Minhao Cheng, Cho-Jui Hsieh, and Yang You. Sparse mezo: Less
595 parameters for better performance in zeroth-order llm fine-tuning. *arXiv preprint arXiv:2402.15751*,
596 2024.
- 597 Shaocong Ma and Heng Huang. Revisiting zeroth-order optimization: Minimum-variance two-point
598 estimators and directionally aligned perturbations. In *Proceedings of The Thirteenth International*
599 *Conference on Learning Representations (ICLR)*, 2025.
- 601 Neural Magic. Twitter financial news sentiment. [https://huggingface.co/datasets/zeroshot/twitter-](https://huggingface.co/datasets/zeroshot/twitter-financial-news-sentiment)
602 *financial-news-sentiment*, 2022.
- 603 Macedo Maia, Siegfried Handschuh, André Freitas, Brian Davis, Ross McDermott, Manel Zarrouk,
604 and Alexandra Balahur. Wwv'18 open challenge: Financial opinion mining and question
605 answering. *Companion Proceedings of the The Web Conference 2018*, 2018. URL <https://api.semanticscholar.org/CorpusID:13866508>.
- 606 Sadhika Malladi, Tianyu Gao, Eshaan Nichani, Alex Damian, Jason D Lee, Danqi Chen, and Sanjeev
607 Arora. Fine-tuning language models with just forward passes. In *Proceedings of Advances in*
608 *Neural Information Processing Systems (NeurIPS)*, volume 36, pp. 53038–53075, 2023.
- 609 P. Malo, A. Sinha, P. Korhonen, J. Wallenius, and P. Takala. Good debt or bad debt: Detecting
610 semantic orientations in economic texts. *Journal of the Association for Information Science and*
611 *Technology*, 65, 2014.
- 612 Yurii Nesterov and Vladimir Spokoiny. Random gradient-free minimization of convex functions.
613 *Foundations of Computational Mathematics*, 17(2):527–566, 2017.
- 614 Ryota Nozawa, Pierre-Louis Poirion, and Akiko Takeda. Zeroth-order random subspace algorithm
615 for non-smooth convex optimization. *Journal of Optimization Theory and Applications*, 204(3):53,
616 2025.
- 617 Rui Pan, Xiang Liu, Shizhe Diao, Renjie Pi, Jipeng Zhang, Chi Han, and Tong Zhang. Lisa: layerwise
618 importance sampling for memory-efficient large language model fine-tuning. *Proceedings of the*
619 *Advances in Neural Information Processing Systems (NeurIPS)*, 37:57018–57049, 2024.
- 620 Sihwan Park, Jihun Yun, SungYub Kim, Souvik Kundu, and Eunho Yang. Unraveling zeroth-
621 order optimization through the lens of low-dimensional structured perturbations. *arXiv preprint*
622 *arXiv:2501.19099*, 2025.
- 623 Mohammad Taher Pilehvar and Jose Camacho-Collados. WiC: the word-in-context dataset for
624 evaluating context-sensitive meaning representations. In *Proceedings of the Conference of the*
625 *North American Chapter of the Association for Computational Linguistics: Human Language*
626 *Technologies (NAACL)*, pp. 1267–1273, 2019.
- 627 Pranan Rajpurkar, Jian Zhang, Konstantin Lopyrev, and Percy Liang. SQuAD: 100,000+ questions
628 for machine comprehension of text. In *Proceedings of the Conference on Empirical Methods in*
629 *Natural Language Processing (EMNLP)*, pp. 2383–2392, 2016.
- 630 Melissa Roemmele, Cosmin Adrian Bejan, and Andrew S Gordon. Choice of plausible alternatives:
631 An evaluation of commonsense causal reasoning. 2011.
- 632 Murray Shanahan and Catherine Clarke. Evaluating large language model creativity from a literary
633 perspective. *arXiv preprint arXiv:2312.03746*, 2023.
- 634 Richard Socher, Alex Perelygin, Jean Wu, Jason Chuang, Christopher D. Manning, Andrew Ng,
635 and Christopher Potts. Recursive deep models for semantic compositionality over a sentiment
636 treebank. In *Proceedings of the Conference on Empirical Methods in Natural Language Processing*
637 *(EMNLP)*, 2013.
- 638 Yan Sun, Tiansheng Huang, Liang Ding, Li Shen, and Dacheng Tao. Tezo: Empowering the low-
639 rankness on the temporal dimension in the zeroth-order optimization for fine-tuning llms. *arXiv*
640 *preprint arXiv:2501.19057*, 2025.

- 648 Qitao Tan, Jun Liu, Zheng Zhan, Caiwei Ding, Yanzhi Wang, Jin Lu, and Geng Yuan. Harmony
649 in divergence: Towards fast, accurate, and memory-efficient zeroth-order llm fine-tuning. *arXiv*
650 *preprint arXiv:2502.03304*, 2025.
- 651 Alex Wang, Yada Pruksachatkun, Nikita Nangia, Amanpreet Singh, Julian Michael, Felix Hill, Omer
652 Levy, and Samuel Bowman. Superglue: A stickier benchmark for general-purpose language
653 understanding systems. In *Proceedings of Advances in Neural Information Processing Systems*
654 *(NeurIPS)*, volume 32, 2019.
- 655 Adina Williams, Nikita Nangia, and Samuel Bowman. A broad-coverage challenge corpus for
656 sentence understanding through inference. In *Proceedings of the Conference of the North American*
657 *Chapter of the Association for Computational Linguistics: Human Language Technologies*
658 *(NAACL)*, 2018.
- 660 Hongyang Yang. Data-centric fingpt. open-source for open finance. [https://github.com/
661 AI4Finance-Foundation/FinGPT](https://github.com/AI4Finance-Foundation/FinGPT), 2023.
- 662 Yifan Yang, Kai Zhen, Ershad Banijamali, Athanasios Mouchtaris, and Zheng Zhang. Adazeta:
663 Adaptive zeroth-order tensor-train adaption for memory-efficient large language models fine-
664 tuning. In *Proceedings of the Conference on Empirical Methods in Natural Language Processing*
665 *(EMNLP)*, pp. 977–995, 2024.
- 666 Ziming Yu, Pan Zhou, Sike Wang, Jia Li, and Hua Huang. Subzero: Random subspace zeroth-order
667 optimization for memory-efficient llm fine-tuning. *arXiv preprint arXiv:2410.08989*, 2024.
- 668 Sheng Zhang, Xiaodong Liu, Jingjing Liu, Jianfeng Gao, Kevin Duh, and Benjamin Van Durme.
669 ReCoRD: Bridging the gap between human and machine commonsense reading comprehension.
670 *arXiv preprint 1810.12885*, 2018.
- 671 Susan Zhang, Stephen Roller, Naman Goyal, Mikel Artetxe, Moya Chen, Shuohui Chen, Christopher
672 Dewan, Mona Diab, Xian Li, Xi Victoria Lin, et al. OPT: Open pre-trained transformer language
673 models. *arXiv:2205.01068*, 2022.
- 674 Yihua Zhang, Pingzhi Li, Junyuan Hong, Jiayang Li, Yimeng Zhang, Wenqing Zheng, Pin-Yu
675 Chen, Jason D Lee, Wotao Yin, Mingyi Hong, et al. Revisiting zeroth-order optimization for
676 memory-efficient llm fine-tuning: a benchmark. In *Proceedings of the International Conference*
677 *on Machine Learning (ICML)*, pp. 59173–59190, 2024.
- 678 Jiawei Zhao, Zhenyu Zhang, Beidi Chen, Zhangyang Wang, Anima Anandkumar, and Yuandong
679 Tian. Galore: Memory-efficient llm training by gradient low-rank projection. In *Proceedings of*
680 *the International Conference on Machine Learning (ICML)*, pp. 61121–61143. PMLR, 2024a.
- 681 Yanjun Zhao, Sizhe Dang, Haishan Ye, Guang Dai, Yi Qian, and Ivor W Tsang. Second-order
682 fine-tuning without pain for llms: A hessian informed zeroth-order optimizer. *arXiv preprint*
683 *arXiv:2402.15173*, 2024b.

688 A DETAILED EXPERIMENTAL SETTINGS

689 A.1 MODEL SETTINGS

690 In our experiments, we considered both large-scale autoregressive language models (i.e., OPT-1.3B
691 and OPT-13B (Zhang et al., 2022)) and a masked language model (i.e., RoBERTa-large (Liu et al.,
692 2019)). In the experiments, all ZO methods used a batch size of 16, except where specified, since
693 larger batches help minimize the gradient approximation variance. We chose MeZO as the main
694 baseline because it is the first widely adopted ZO optimizer for LLMs, and included the first-order
695 SGD as a reference for optimization. In line with previous research (Malladi et al., 2023; Zhang
696 et al., 2024), our experiments utilized standardized prompt templates, which are crucial in influencing
697 the performance of ZO methods. Moreover, to ensure a fair comparison, we considered multiple
698 values for each key hyperparameter. For example, we investigated the following hyperparameter
699 configurations for OPT-13B: a learning rate in $\{1e-7, 2e-7, 5e-7, 1e-6\}$, $\epsilon = 1e-3$, a batch size
700 of 16 (except for MultiRC and DROP which have a batch size of 8), a rank in $\{24, 32, 48, 64, 128\}$,
701

and a subspace update frequency in $\{500, 1000, 2000\}$. Please refer to Appendix A for detailed configurations of other models. Similar to the work in (Yu et al., 2024), we conducted an exhaustive grid search over hyperparameters for each pairing of ZO methods and LLMs, using the best results for an equitable comparison.

A.2 DATASET SETTINGS.

For OPT models, we experimented with the SuperGLUE benchmark (Wang et al., 2019), which consists of various types of tasks, including classification tasks (e.g., SST-2 (Socher et al., 2013), RTE (Bar Haim et al., 2006; Bentivogli et al., 2009; Dagan et al., 2005; Giampiccolo et al., 2007), CB (de Marneffe et al., 2019), BoolQ (Clark et al., 2019), WSC (Levesque et al., 2012), and WIC (Pilehvar & Camacho-Collados, 2019)), multiple choice tasks (e.g., COPA (Roemmele et al., 2011) and ReCoRD (Zhang et al., 2018)), and generation tasks (e.g., SQuAD (Rajpurkar et al., 2016) and DROP (Dua et al., 2019)). Here, for each task, we randomly selected 1000 samples for training, 500 samples for validation, and 1000 samples for testing. For the RoBERTa-large model, in addition to the task SST-2, we investigated three more tasks, i.e., SST-5 (Socher et al., 2013), SNLI (Bowman et al., 2015), and MNLI (Williams et al., 2018). In this case, we fixed the parameter k at 512 throughout the training and validation phases, indicating that 512 samples are allocated for each category. For the testing phase, we randomly chose a total of 1000 samples.

A.3 HYPERPARAMETER SETTINGS

This section provides a detailed overview of the hyperparameters employed in our grid search across the experiments, as depicted in Tables 10 and 12. For the OPT model, we carried out 20,000 steps for each method. Both the SGD and ZO methodologies were implemented for an identical number of steps. In the remaining RoBERTa experiments, ZO optimization strategies were applied over 100,000 training steps. For both models, we evaluated the validation loss every 1,000 training steps to determine the optimal model checkpoint. In the S-MeZO strategy, the sparsity rate is set to 0.75.

Table 10: The hyperparameter grids used for OPT-13B experiments.

Method	Hyperparameters				
	Batch Size	Learning Rate	ϵ	Rank	Update Interval
SGD	16	$\{1e-4, 1e-3, 5e-3\}$	-	-	-
MeZO (Malladi et al., 2023)	16	$\{1e-7, 2e-7, 5e-7, 1e-6\}$	$1e-3$	-	-
S-MeZO (Liu et al., 2024)	16	$\{1e-6, 5e-6\}$	$1e-3$	-	-
LOZO (Chen et al., 2025)	16	$\{1e-7, 1e-6\}$	$\{1e-3, 1e-4\}$	$\{1, 2, 4\}$	$\{50, 100\}$
SubZero (Yu et al., 2024)	16	$\{1e-7, 2e-7, 5e-7, 1e-6\}$	$1e-3$	$\{32, 64, 128, 256\}$	$\{500, 1000, 2000\}$
SVD-0	16	$\{1e-7, 2e-7, 5e-7, 1e-6\}$	$1e-3$	$\{24, 32, 48, 64, 128\}$	$\{500, 1000, 2000\}$

Table 11: The hyperparameter grids used for OPT-1.3B experiments.

Method	Hyperparameters				
	Batch Size	Learning Rate	ϵ	Rank	Update Interval
MeZO (Malladi et al., 2023)	16	$\{1e-7, 5e-7, 1e-6\}$	$1e-3$	-	-
LOZO (Chen et al., 2025)	16	$\{1e-7, 1e-6\}$	$\{1e-3, 1e-4\}$	$\{1, 2, 4\}$	$\{50, 100\}$
SubZero (Yu et al., 2024)	16	$\{1e-7, 5e-7, 1e-6\}$	$1e-3$	$\{24, 48\}$	1000
SVD-0	16	$\{1e-7, 5e-7, 1e-6\}$	$1e-3$	$\{8, 24, 48\}$	$\{50, 500, 1000\}$

For all previously mentioned ZO methods, we utilized a consistent learning rate schedule and set the weight decay to zero. Typically, we chose a batch size of 16 for the OPT-1.3B and OPT-13B models across various tasks. Nonetheless, due to limited GPU resources, we reduced the batch size to 8 for the DROP, MultiRC, and SQuAD evaluations.

Table 12: The hyperparameter grids used for RoBERTa-large experiments.

Method	Hyperparameters				
	Batch Size	Learning Rate	ϵ	Rank	Update Interval
MeZO (Malladi et al., 2023)	64	{1e-7, 1e-6, 1e-5}	1e-3	-	-
LOZO (Chen et al., 2025)	64	2e-7	1e-3	{4, 8}	{50, 100}
SVD-0	64	1e-6	1e-3	{8, 16, 24}	1000

B DETAILED PRESTUDY RESULTS

Table 13: Impact of rank and batch size.

Rank\Batch Size	1	2	4	8	16
8	0.87	0.86	0.50	0.72	0.67
24	0.85	0.84	0.45	0.70	0.63
48	0.85	0.84	0.45	0.70	0.63

For SVD, we perform truncation by restricting the rank to 24 for the singular vectors. In Figure 1, the x-axis indicates the number of training steps completed during fine-tuning, while the y-axis plots the cosine similarity between the ZO gradient at the current step and the FO gradient, both after applying SVD at each step. It is important to note that this cosine similarity is calculated using only a subset of the parameters, rather than as an average across the entirety of the parameters.

We further explored how similarity changes under various rank and batch size configurations. To do this, we examined the ranks of 8, 24, and 48, paired with batch sizes of 1, 2, 4, 8, and 16. Table 13 summarizes these findings. Our results show that reducing the rank leads to a minor decrease in similarity, though the effect is minimal, and overall similarity remains stable. In contrast, modifying batch size results in more noticeable fluctuations in cosine similarity, but there is no clear trend of increase or decrease. In summary, the similarity remains consistently high in most settings.

Table 14: Comparison between SVD and other dimensionality reduction techniques.

Method	Mean	Std	Min	Max	Median
Origin	0.0000	0.0005	-0.0013	0.0012	0.0000
After SVD (ours)	0.4249	0.0257	0.3276	0.4774	0.4278
PCA	-0.0003	0.0044	-0.0141	0.0101	-0.0004
NMF	0.2464	0.0442	0.1097	0.3594	0.2456
Factor Analysis	0.0001	0.0045	-0.0124	0.0136	-0.0002
Random Proj	-0.0002	0.0045	-0.0112	0.0177	-0.0001
t-SNE	-0.0009	0.0157	-0.0437	0.0567	-0.0012

We conducted experiments to investigate why SVD outperforms other methods for dimensionality reduction. Table 14 presents a comparison between our dimensionality reduction strategy ("After SVD") and other techniques (PCA, NMF, Factor Analysis, Random Projection, and t-SNE) based on the calculated gradients. As shown in Table 14, our method yields the highest mean value (0.4249) among all dimensionality reduction techniques, accompanied by a lower standard deviation (0.0257), which highlights its superior and consistent performance.

C ADDITIONAL RESULTS

Figure 3 presents the loss curves of SVD-0 and SubZero during training, where the x-axis represents the wall-clock time. Although SVD-0 takes a bit longer than SubZero for each iteration, it consistently reaches lower training loss from the beginning to the end. This indicates that, overall, SVD-0 converges faster in terms of total training time.

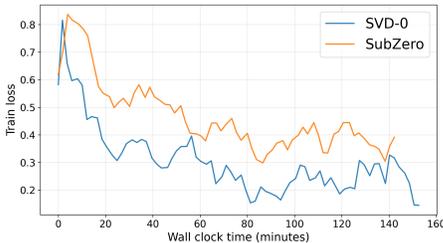


Figure 3: Loss curves plotted against wall-clock time.

We examined the cosine similarity between gradients obtained from different ZO methods and those from FO. Throughout the fine-tuning process of

the OPT-1.3B model on the SST2 (Socher et al., 2013) dataset, we periodically sampled the actual gradients for a chosen subset of parameters, along with the ZO gradients. The cosine similarity between these ZO and FO gradients was then calculated. The resulting values were averaged across the entire fine-tuning process and normalized to the $[0, 1]$ range to facilitate comparison.

Table 15: Cosine similarity between ZO gradients and the FO gradient.

Method	Cosine Similarity
MeZO (Malladi et al., 2023)	0
SubZero (Yu et al., 2024)	0.73
SVD-0	1

The results presented in Table 15 indicate that SVD-0 achieves the highest level of similarity. This suggests that using SVD to extract the shared features between ZO and FO gradients results in a closer match to the FO gradient. These observations further reinforce the conclusions drawn in our preliminary study.

Table 16: Variance of the estimated gradients.

Method	Variance
MeZO (Malladi et al., 2023)	5.5
SubZero (Yu et al., 2024)	5.6
SVD-0	5.0

We also conducted experiments to evaluate the variance of the estimated gradients. For each approach, we performed tests on the SST-2 dataset with OPT-1.3B using identical hyperparameters. We collected 1,000 samples per method and calculated the variance for the resulting estimated gradients. A summary of the results is provided in Table 16. The findings indicate that SVD-0 yields the lowest variance, further supporting the superiority of SVD-0.

We performed experiments to explore the behavior of SVD-0 when subjected to significantly noisy zero-order (ZO) gradients. Gaussian noise with varying variances was introduced into the ZO gradients for these tests. In particular, we injected zero-mean Gaussian noise into the ZO gradients, with its variance scaled to K times that of the projected gradients, simulating noisy conditions. We then fine-tuned the OPT-1.3B model on the SST-2 (Socher et al., 2013), RTE (Bar Haim et al., 2006; Bentivogli et al., 2009; Dagan et al., 2005; Giampiccolo et al., 2007), and WIC (Pilehvar & Camacho-Collados, 2019) datasets. The results for these experiments are shown in Table 17.

Table 17: Impact of different noise conditions.

K	SST-2	RTE	WIC
No noise	92.2	68.2	56.3
0.5	87.5	68.6	56.6
1	88.9	66.1	56.4
2	87.5	66.4	55.0

In general, increasing the noise variance leads to a modest drop in model performance, yet SVD-0 continues to show resilience in noisy scenarios. Furthermore, the best results are achieved in the absence of noise or when noise levels are minimal, suggesting that SVD-0 effectively leverages the valuable aspects of the ZO gradients to enhance fine-tuning.

To demonstrate that SVD-0 consistently delivers superior performance across various random seeds on OPT, we conducted an experiment where we fine-tuned the OPT-1.3B model on the SST-2 (Socher et al., 2013) dataset. As shown in Table 18, SVD-0 consistently outperforms MeZO (Malladi et al., 2023) across all random seeds.

Table 18: Impact of random seed.

Method\Seed	0	42	123
MeZO (Malladi et al., 2023)	90.9	88	89.2
SVD-0	93.0	90.1	90.4

As presented in Table 19, we fine-tuned the OPT-1.3B model on various datasets using the Adam optimizer, applying the same hyperparameters as in the SGD setup. Our results demonstrate that SVD-0, in combination with Adam, consistently surpasses the MeZO baseline. This suggests that SVD-0 is especially effective for real-world fine-tuning scenarios involving the Adam optimizer.

Table 19: Performance on Adam optimizer.

Method	SST-2	WIC	RTE
MeZO(Adam)	75.6	55.8	54.2
SVD-0(Adam)	92.9	62.4	75.5

Given that vanilla LoRA uses Adam for fine-tuning, we set up experiments to compare SVD-0 with SGD in fine-tuning scenarios, alongside vanilla LoRA with Adam. All experiments were conducted using the pretrained OPT-1.3B model and evaluated on the SST-2 dataset. Table 20 shows that both SVD-0 (SGD) and LoRA (Adam) deliver similar results on SST-2. However, SVD-0 (SGD) stands out by using much less memory than LoRA (Adam), underscoring its key advantage as a zeroth-order method: substantial memory savings without a loss in performance.

Table 20: Comparison with LoRA.

Method	Accuracy(%)	Memory(GB)
LoRA(Adam)	93.2	12.926
SVD-0(SGD)	93.0	4.891

We also integrated SVD-0 with LoRA and performed experiments using OPT-13B. We set the learning rates for MeZO (LoRA) and SVD-0 (LoRA) to $\{1.5e-5, 3e-5, 5e-5\}$, the rank of SVD-0 (LoRA) to $\{4, 8, 16\}$, and the subspace change frequency to $\{500, 1000, 2000\}$.

Table 21: Performance of SVD-0 combined with LoRA.

Method	SST-2	RTE	CB	BoolQ	WSC	WIC	MultiRC
MeZO(LoRA)	92.2	74.4	69.6	75.2	64.4	59.7	58.2
SVD-0(LoRA)	93.8	76.9	71.4	76.0	65.4	61.4	59.3

As shown in Table 21, SVD-0 (LoRA) consistently achieves superior results compared to MeZO (LoRA) across all datasets, suggesting that integrating SVD-0 with LoRA yields a substantial performance boost.

Table 22: Performance of SVD-0 on various ZO estimators.

Estimators	MeZO	SPSA	NES
Accuracy(%)	92.2	90.5	90.3

As presented in Table 22, we replaced MeZO with alternative zeroth-order (ZO) gradient estimators to provide SVD-0 with the necessary gradient information. Specifically, we utilized SPSA and NES (Natural Evolution Strategies) as gradient estimators and conducted experiments on the SST-2 dataset using OPT-1.3B, maintaining consistent hyperparameter settings. The performance outcomes are summarized below. Our results show that, although SPSA and NES perform slightly below MeZO, both methods still deliver robust results. This suggests that SVD-0 is a generally applicable subspace technique that is not dependent on MeZO in particular.

D PROOFS

Here, we introduce some definitions and lemmas for continuous proofs. The following two lemmas illustrate that the low-rank perturbation matrix for each layer can be represented as a layer-scale projection matrix that is orthogonal across its columns.

Lemma 3. Let $\tilde{\mathbf{Z}} = \mathbf{U}\mathbf{Z}\mathbf{V}^\top$, 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 $\mathbf{U}^\top \mathbf{U} = \mathbf{V}^\top \mathbf{V} = \mathbf{I}_r$. Then we have $\text{vec}(\tilde{\mathbf{Z}}) = \mathbf{P}\text{vec}(\mathbf{Z})$ and $\mathbf{P}^\top \mathbf{P} = \mathbf{I}_{r^2}$, where $\mathbf{P} = \mathbf{V} \otimes \mathbf{U}$.

Proof. Since $\text{vec}(\mathbf{U}\mathbf{Z}\mathbf{V}^\top) = (\mathbf{V} \otimes \mathbf{U})\text{vec}(\mathbf{Z})$, 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 show that the low-rank perturbation matrices across all layers can be represented as a model-scale projection matrix.

Lemma 4. 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 $\mathbf{P}_i^\top \mathbf{P}_i = \mathbf{I}_{r_i}$ 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 $\mathbf{P}^\top \mathbf{P} = \mathbf{I}_{lr}$.

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

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

The proof is completed. \square

According to Lemma 4 and Proposition 1, the perturbation vector of SVD-0 is given by $\tilde{\mathbf{z}} = \mathbf{P} \mathbf{z}$. This is similar to existing random subspace methods, but SVD-0's projection matrix is block diagonal and orthogonal by columns.

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

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

Lemma 5. Let $\mathbf{z} \sim \mathcal{N}(\mathbf{0}, \mathbf{I}_n)$. For any orthogonal $(n \times n)$ matrix \mathbf{Q} and any continuous function f , we have $\mathbb{E}_{\mathbf{z}}[f(\mathbf{z})] = \mathbb{E}_{\mathbf{z}}[f(\mathbf{Q}\mathbf{z})]$.

Lemma 6. If $x \sim \chi^2(n)$, then we have:

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

Lemma 7. (Nesterov & Spokoiny, 2017) Let $f \in C_{L_2}^{2,2}(\mathbb{R}^n)$. 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 8. (Nesterov & Spokoiny, 2017) 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 9. 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 5, 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 6, 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. \square

Here we provide the proof of Lemma 2.

Proof. **a)** The conclusion is clearly supported by Lemma 3 and Lemma 4.

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 7 implies that:

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

Note that:

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

Therefore, in accordance with Lemma 8, we have:

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

The proof is completed. \square

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:

$$\mathbb{E}_z[\widehat{\nabla} f(\boldsymbol{\theta})] = \mathbf{P}\mathbf{P}^\top \nabla f(\mathbf{x}), \quad (6)$$

$$\mathbb{E}_z[\|\widehat{\nabla} f(\boldsymbol{\theta})\|^2] = (q+2)\|\mathbf{P}^\top \nabla f(\mathbf{x})\|^2, \quad (7)$$

$$\mathbb{E}_z \left[\frac{\langle \nabla f(\mathbf{x}), \widehat{\nabla} f(\boldsymbol{\theta}) \rangle^2}{\|\mathbf{P}^\top \nabla f(\mathbf{x})\|^2 \|\widehat{\nabla} f(\boldsymbol{\theta})\|^2} \right] = \frac{1}{q}. \quad (8)$$

Proof. It is straightforward to verify that

$$\widehat{\nabla} f(\boldsymbol{\theta}) = \mathbf{P}\langle \mathbf{P}^\top \nabla f(\mathbf{x}), \mathbf{z} \rangle \mathbf{z}.$$

Therefore, we can express the expected value as

$$\mathbb{E}_z[\widehat{\nabla} f(\boldsymbol{\theta})] = \mathbf{P}\mathbf{P}^\top \nabla f(\mathbf{x}).$$

Combining this with Lemma 9, we find that

$$\mathbb{E}_z[\|\widehat{\nabla} f(\boldsymbol{\theta})\|^2] = (q+2)\|\mathbf{P}^\top \nabla f(\mathbf{x})\|^2.$$

Additionally, it is important to note that for any orthogonal $(q \times q)$ matrix \mathbf{Q} , we have:

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

In accordance with Lemma 5, we can set $\mathbf{P}^\top \nabla f(\mathbf{x}) = [1, 0, \dots, 0]^\top$. Thus, we have:

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

The proof is completed. \square

To demonstrate the convergence of SVD-0 with SGD, we can structure our analysis into two main segments. The first segment examines the convergence behavior of the SVD-0 solution process while keeping the projection matrix \mathbf{P} constant. The second segment evaluates the impact of performing lazy updates to \mathbf{P} . Through these assessments, we aim to establish the global convergence of SVD-0, specifically in the context of a single layer.

In the initial stage, while \mathbf{P} remains constant, we can reformulate the original SVD-0 problem as an optimization problem constrained within that subspace. We define $h(\mathbf{y}) = f(\mathbf{x} + \mathbf{P}\mathbf{y})$, $h_\varepsilon(\mathbf{y}) = \mathbb{E}_{\mathbf{z}}[h(\mathbf{y} + \varepsilon\mathbf{z})]$, and $g_\varepsilon(\mathbf{y}) = \frac{h(\mathbf{y} + \varepsilon\mathbf{z}) - f(\mathbf{y})}{\varepsilon}\mathbf{z}$. According to Lemma 10, if f demonstrates first L_1 -smoothness, then h will also exhibit first L_1 -smoothness.

Lemma 10. *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 demonstrates 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

We then analyze the convergence of SVD-0 while maintaining a fixed subspace.

Lemma 11. (Nesterov & Spokoiny, 2017) *Let $f \in C_{L_1}^{1,1}(\mathbb{R})$. Then, for any $\mathbf{x} \in \mathbb{R}$, we have:*

$$E_{\mathbf{z}}[\|g_\varepsilon(\mathbf{x})\|^2] = E_{\mathbf{z}}\left[\left\|\frac{f(\mathbf{x} + \varepsilon\mathbf{z}) - f(\mathbf{x})}{\varepsilon}\right\|^2\right] \leq 4(n+4)\|\nabla f_\varepsilon(\mathbf{x})\|^2 + 3\varepsilon^2 L_1^2(f)(n+4)^3, \quad (9)$$

and

$$\|\nabla f(\mathbf{x})\|^2 \leq 2\|\nabla f_\varepsilon(\mathbf{x})\|^2 + \frac{\varepsilon^2}{2} L_1^2(f)(n+6)^3, \quad (10)$$

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

Lemma 12. *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 $\mathbf{z}_k \sim \mathcal{N}(\mathbf{0}, \mathbf{I}_q)$ 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 Proposition 1, which is fixed, we have:*

$$\phi_{k+1} - \phi_k \leq -\frac{1}{4}\eta \mathbb{E}_{\mathcal{E}_k}[\|\nabla h(\mathbf{y}_k)\|^2] + \frac{\varepsilon^2(q+6)^3}{8} L_1^2 + \frac{3\varepsilon^2(q+4)}{32} L_1 \quad (11)$$

Proof. If we have a fixed subspace $\mathbf{P} \in \mathbb{R}^{d \times q}$, we can reformulate the optimization objective as follows:

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

Let \mathbf{y}_0 represent the initial point, and let $\{\eta_k\}_{k \geq 0}$ be a sequence of positive real numbers. We will consider the randomized gradient search algorithm $\mathcal{RG}_\varepsilon(\varepsilon > 0)$:

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 $\mathbf{y}_{k+1} = \mathbf{y}_k - \eta_k g_\varepsilon(\mathbf{y}_k)$.

Our goal is to estimate the evolution of the function h_ε after one iteration of this algorithm.

Given that h is L_1 -Lipschitz continuous for its first derivative, and h_ε is L_ε -Lipschitz continuous for its first derivative (where $L_\varepsilon \leq L_1$)(Nesterov & Spokoiny, 2017), 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 have:

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

Since $h \in C^{1,1}(\mathbb{R}^q)$, from Lemma 11, 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 11, 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:

$$\phi_{k+1} - \phi_k \leq -\frac{1}{4} \hat{\eta} \mathbb{E}_{\mathcal{E}_k} [\|\nabla h(\mathbf{y}_k)\|^2] + \frac{\varepsilon^2(q+6)^3}{8} L_1^2 + \frac{3\varepsilon^2(q+4)}{32} L_1. \quad (12)$$

The proof is completed. \square

Next, we need to assess the randomness of our random subspace. According to Lemma 16, if we obtain the projection matrix using Algorithm 1, the expected value can be expressed as $\mathbb{E}[\mathbf{P}\mathbf{P}^T] = \frac{q}{d}\mathbf{I}$. In this equation, q represents the dimension of the subspace, d indicates the dimension of the original space, and \mathbf{P} is defined as $\mathbf{V} \otimes \mathbf{U}$.

Lemma 13. *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) $\mathbb{E} \left[\frac{[\mathbf{u}_k]_i [\mathbf{u}_k]_j}{\langle \mathbf{u}_k, \mathbf{u}_k \rangle} \right] = 0$, where $i \neq j$.

Proof. In real analysis, a matrix \mathbf{A} usually possesses full rank when it follows a Gaussian distribution, and it is common for both \mathbf{u}_k and \mathbf{e}_k to be non-zero.

(1) Given that \mathbf{a}_k is independently and identically distributed, this condition clearly applies.

(2) For the base case $k = 1$, it is obviously true. Assume that 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

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

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

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

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

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

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

Lemma 14. 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 $\mathbf{u}_k = \mathbf{a}_k - \sum_{s=1}^{k-1} \langle \mathbf{a}_k, \mathbf{e}_s \rangle \mathbf{e}_s$.

$$\text{Thus, } \mathbf{e}_k \mathbf{e}_k^T = \frac{\mathbf{u}_k \mathbf{u}_k^T}{\langle \mathbf{u}_k, \mathbf{u}_k \rangle}.$$

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 13(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 13(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 15. 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 the matrix \mathbf{A} is expressed as $\mathbf{A} = \mathbf{Q}\mathbf{R}$, where \mathbf{Q} is an orthogonal matrix with columns denoted as $\mathbf{e}_1, \mathbf{e}_2, \dots, \mathbf{e}_r$, and \mathbf{R} is an upper triangular matrix. Since \mathbf{Q} is orthogonal, it satisfies the condition $\mathbf{Q}\mathbf{Q}^T = \mathbf{I}_r$, where \mathbf{I}_r represents the $r \times r$ identity matrix. Our objective is to compute $\mathbb{E}[\mathbf{Q}\mathbf{Q}^T]$. By leveraging the linearity of expectation and the fact that the columns of \mathbf{Q} are orthonormal, we find:

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

From Lemma 14, 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. \square

Lemma 16. 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 $\mathbf{P} = \mathbf{Q}_2 \otimes \mathbf{Q}_1$ produces a matrix $\mathbf{P} \in \mathbb{R}^{mn \times r^2}$. According to Lemma 15, we know that $\mathbf{Q}_1 \mathbf{Q}_1^T = \frac{r}{m} \mathbf{I}$ and $\mathbf{Q}_2 \mathbf{Q}_2^T = \frac{r}{n} \mathbf{I}$. Our goal is to compute $\mathbb{E}[\mathbf{P}\mathbf{P}^T]$. By utilizing the properties of the Kronecker product, we can proceed with our computation:

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

The proof is completed. \square

Now we can assess the impact of the lazy updates to \mathbf{P} . Here we provide the proof of Theorem 1.

Proof. Let $\mathcal{P}_j = (\mathbf{P}_0, \mathbf{P}_1, \dots, \mathbf{P}_j)$, where \mathbf{P}_j is the sequence generated by Proposition 1 and $j \leq K$. According to Lemma 10 and Lemma 12, when the subspace is fixed, we can transform the original problem $f \in C_{L_1}^{1,1}(\mathbb{R}^d)$ into $h \in C_{L_1}^{1,1}(\mathbb{R}^q)$ using the transformation $h(\mathbf{y}) = f(\mathbf{x} + \mathbf{P}\mathbf{y})$. Consider the following update rule:

$$\mathbf{y}_{j,0} = 0, h_j(\mathbf{y}) = f(\mathbf{x}_{jF} + \mathbf{P}_j \mathbf{y}), \forall j \in 0, 1, \dots, K-1 \quad (13)$$

$$\mathbf{y}_{j,k} = \mathbf{y}_{j,k-1} - \eta \widehat{\nabla} h_j(\mathbf{y}_{j,k-1}), \forall k \in 0, 1, \dots, F \quad (14)$$

$$\mathbf{x}_{jF+k} = \mathbf{x}_{jF} + \mathbf{P}_j \mathbf{y}_k, \quad (15)$$

In the j -th subspace, the projection matrix \mathbf{P}_j is constant, allowing us to accumulate the changes of ϕ within this subspace. By applying Lemma 12, we obtain:

$$\phi_{(j+1)F} - \phi_{jF} \leq -\frac{1}{4} \widehat{\eta} \sum_{i=0}^{K-1} \mathbb{E}_{\mathcal{E}_{jF+i}} [\|\nabla h_j(\mathbf{y}_{j,i})\|^2] + \frac{\varepsilon^2(q+6)^3}{8} KL_1^2 + \frac{3\varepsilon^2(q+4)}{32} KL_1 \quad (16)$$

$$\leq -\frac{1}{4} \widehat{\eta} \mathbb{E}_{\mathcal{E}_{jF}} [\|\nabla h_j(\mathbf{y}_{j,0})\|^2] + \frac{\varepsilon^2(q+6)^3}{8} KL_1^2 + \frac{3\varepsilon^2(q+4)}{32} KL_1. \quad (17)$$

Furthermore, we note that $\nabla h_j(\mathbf{y}_{j,0}) = (\mathbf{P}_j)^T \nabla f(\mathbf{x}_{jF})$. By taking expectations over the overall historical projection matrix \mathcal{P}_j and applying Lemma 16, we find that $\mathbb{E}[\mathbf{P}_j (\mathbf{P}_j)^T] = \frac{q}{d} \mathbf{I}$, with \mathbf{P}_j

1242 being independent of \mathbf{x}_{jF} . Thus, we obtain:

$$1243 \mathbb{E}_{\mathcal{P}_{j+1}}[\phi_{(j+1)F}] - \mathbb{E}_{\mathcal{P}_j}[\phi_{jF}] \leq -\frac{1}{4}\hat{\eta}\mathbb{E}_{\mathcal{E}_{jF},\mathcal{P}_j}[\|(\mathbf{P}_j)^\top \nabla f(\mathbf{x}_{jF})\|^2] + \frac{\varepsilon^2(q+6)^3}{8}KL_1^2 + \frac{3\varepsilon^2(q+4)}{32}KL_1$$

$$1244 = -\frac{q}{4d}\hat{\eta}\mathbb{E}_{\mathcal{E}_{jF},\mathcal{P}_j}[\|\nabla f(\mathbf{x}_{jF})\|^2] + \frac{\varepsilon^2(q+6)^3}{8}KL_1^2 + \frac{3\varepsilon^2(q+4)}{32}KL_1. \quad (18)$$

1245 Assuming $f(\mathbf{x}) \geq f^*$ holds for all $\mathbf{x} \in \mathbb{R}^d$, and letting $T = KF$, summing the inequality yields:

$$1246 \mathbb{E}_{\mathcal{P}_{K-1}}[\phi_T] \leq \mathbb{E}_{\mathcal{P}_0}[\phi_0] - \frac{q}{4d}\hat{\eta} \sum_{j=0}^{K-1} \mathbb{E}_{\mathcal{E}_{jF},\mathcal{P}_j}[\|\nabla f(\mathbf{x}_{jF})\|^2] + T \frac{\varepsilon^2(q+6)^3}{8}L_1^2 + T \frac{3\varepsilon^2(q+4)}{32}L_1. \quad (19)$$

1247 Since $\mathbb{E}_{\mathcal{P}_{K-1}}[\phi_T] \geq f^*$, we have:

$$1248 f^* \leq \mathbb{E}_{\mathcal{P}_0}[\phi_0] - \frac{q}{4d}\hat{\eta} \sum_{j=0}^{K-1} \mathbb{E}_{\mathcal{E}_{jF},\mathcal{P}_j}[\|\nabla f(\mathbf{x}_{jF})\|^2] + T \frac{\varepsilon^2(q+6)^3}{8}L_1^2 + T \frac{3\varepsilon^2(q+4)}{32}L_1. \quad (20)$$

1249 Rearranging the inequality, we get:

$$1250 \frac{q}{4d}\hat{\eta} \sum_{j=0}^{K-1} \mathbb{E}_{\mathcal{E}_{jF},\mathcal{P}_j}[\|\nabla f(\mathbf{x}_{jF})\|^2] \leq \mathbb{E}_{\mathcal{P}_0}[\phi_0] - f^* + T \frac{\varepsilon^2(q+6)^3}{8}L_1^2 + T \frac{3\varepsilon^2(q+4)}{32}L_1. \quad (21)$$

1251 Substituting $\hat{\eta} = \frac{1}{4(q+4)L_1}$, we obtain:

$$1252 \frac{q}{16d(q+4)L_1} \sum_{j=0}^{K-1} \mathbb{E}_{\mathcal{E}_{jF},\mathcal{P}_j}[\|\nabla f(\mathbf{x}_{jF})\|^2] \leq \mathbb{E}_{\mathcal{P}_0}[\phi_0] - f^* + T \frac{\varepsilon^2(q+6)^3}{8}L_1^2 + T \frac{3\varepsilon^2(q+4)}{32}L_1. \quad (22)$$

1253 Thus, we have:

$$1254 \frac{1}{T} \sum_{k=0}^{T-1} \mathbb{E}_{\mathcal{E}_k, \mathcal{P}_{\lfloor k/F \rfloor}}[\|\nabla f(\mathbf{x}_k)\|^2] \leq \frac{16(q+4)dL_1(\mathbb{E}_{\mathcal{P}_0}[\phi_0] - f^*)}{qT} + \frac{2\varepsilon^2(q+6)^3(q+4)d}{q}L_1^3 + \frac{3\varepsilon^2(q+4)^2d}{2q}L_1^2. \quad (23)$$

1255 To ensure $\frac{1}{T} \sum_{k=0}^{T-1} \mathbb{E}_{\mathcal{E}_k, \mathcal{P}_{\lfloor k/F \rfloor}}[\|\nabla f(\mathbf{x}_k)\|^2] \leq \epsilon$, we can choose:

$$1256 \epsilon \leq \mathcal{O}\left(\frac{1}{q^{3/2}d^{1/2}L_1^{3/2}}\right).$$

1257 As a result, the convergence rate is $\mathcal{O}(\sqrt{\frac{d}{T}})$. The proof is completed. \square