SIMULTANEOUS COMPUTATION AND MEMORY EFFI CIENT ZEROTH-ORDER OPTIMIZER FOR FINE-TUNING LARGE LANGUAGE MODELS

Anonymous authors

006

008 009 010

011

013

014

015

016

017

018

019

021

023

025

026

027

028

029

031

032

Paper under double-blind review

ABSTRACT

Fine-tuning is powerful for adapting large language models to downstream tasks, but it often results in huge memory usages. A promising approach to mitigate this is using Zeroth-Order (ZO) optimization, which estimates gradients to replace First-Order (FO) gradient calculations, albeit with longer training time due to its stochastic nature. By revisiting the Memory-efficient ZO (MeZO) optimizer, we discover that the full-parameter perturbation and updating processes consume over 50% of its overall fine-tuning time cost. Based on these observations, we introduce a novel layer-wise sparse computation and memory efficient ZO optimizer, named LeZO. LeZO treats layers as fundamental units for sparsification and dynamically perturbs different parameter subsets in each step to achieve full-parameter finetuning. LeZO incorporates layer-wise parameter sparsity in the process of simultaneous perturbation stochastic approximation (SPSA) and ZO stochastic gradient descent (ZO-SGD). It achieves accelerated computation during perturbation and updating processes without additional memory overhead. We conduct extensive experiments with the OPT model family on the SuperGLUE benchmark and two generative tasks. The experiments show that LeZO accelerates training without compromising the performance of ZO optimization. Specifically, it achieves over $3 \times$ speedup compared to MeZO on the SST-2, BoolQ, and Copa tasks.

1 INTRODUCTION

033 Large language models (LLMs) have shown 034 remarkable capabilities in understanding and generating languages, and have been widely adopted in various applications (Brown et al., 2020; Wei et al., 2022; Rao et al., 2022; Wang 037 et al., 2023). To effectively adapt LLMs to downstream tasks, full-parameter fine-tuning has become crucial (Li & Liang, 2021; Lester 040 et al., 2021; Hu et al., 2022; Malladi et al., 041 2023). However, as the scale of LLMs con-042 tinues to increase, the memory usage and com-043 putational cost of fine-tuning also escalate (Ka-044 plan et al., 2020; Hoffmann et al., 2022), posing a significant challenge to the practical application of LLMs. 046



Figure 1: LeZO achieves a $3.4 \times$ speedup for finetuning the OPT-13b model in run-time compared with MeZO on the SST-2 dataset.

⁰⁴⁷ To address these challenges, researchers have

proposed efficient fine-tuning methods, which can be categorized into forward-backward (Li & Liang, 2021; Hu et al., 2022; Ding et al., 2022; Pu et al., 2023; Xu et al., 2023) and forward-only (Zhang et al., 2024; Liu et al., 2024; Guo et al., 2024) approaches. Forward-backward methods typically use First-Order (FO) optimizers (Robbins & Monro, 1951; Kingma & Ba, 2015), updating either partial model parameters or introducing small trainable modules to reduce memory usage. However, the performance of these methods can hardly match full-parameter fine-tuning. Dynamic update schemes (Brock et al., 2017; Liu et al., 2021; Pan et al., 2024) address this problem by dynam-

icly changing learnable parameter subsets, but they still require caching many activations. In contrast, forward-only methods utilize Zeroth-Order (ZO) optimizers, which estimate gradients without back-propagation, thereby reducing the memory overhead associated with forward-backward
methods. Malladi et al. (2023) first introduce a memory-efficient ZO optimizer (MeZO) for LLM
fine-tuning. ZO optimizers, however, exhibit higher stochasticity in gradient estimation compared
to FO optimizers, and this stochasticity increases as the scale of tuned parameters grows (Wang
et al., 2018; Balasubramanian & Ghadimi, 2018; Cai et al., 2022). Consequently, MeZO requires
significantly more computational cost than FO optimizers.

In this paper, we revisit MeZO and identify that its computational cost is primarily in the forward pass, perturbation, and updating stages. Specifically, perturbation and updating account for over 50% of the total time when fine-tuning the OPT-13B model (Zhang et al., 2022) on the SST-2 task.
Simplifying these steps can accelerate MeZO. One promising approach is to sparsify the tuned parameters, which has been shown to be effective for FO optimization (Pan et al., 2024).

067 Inspired by these insights, we propose a layer-wise sparse and efficient zeroth-order opti-068 mizer (LeZO). LeZO aims to improve computational efficiency without additional memory over-069 head. We use layers as the basic unit of sparsity, dynamically updating the set of layers to be tuned at different fine-tuning steps. We integrate this policy into Simultaneous Perturbation Stochastic 071 Approximation (SPSA) and Zeroth-Order Stochastic Gradient Descent (ZO-SGD), which enables sparse ZO optimization through localized parameter perturbation and updating. Over multiple steps, 072 LeZO can achieve full-parameter fine-tuning. LeZO accelerates computations because, during each 073 step, the weights of untuned layers remain dense. This allows us to skip perturbation and updating 074 processes for those layers, thereby reducing floating-point computations. Additionally, compared to 075 MeZO, LeZO does not introduce extra memory overhead. We employ a simple and efficient random 076 selection strategy, avoiding the need for new parameter modules. 077

We evaluate our approach by fine-tuning models from the OPT family on the SuperGLUE benchmark (Wang et al., 2019), and the SQuAD (Rajpurkar et al., 2016) and DROP (Dua et al., 2019) datasets. By sparsifying 75% of the layers and fine-tuning full-parameters, LeZO outperforms MeZO on most datasets across three model scales. For instance, on the OPT-13B model, LeZO surpasses MeZO by an average of 1% in accuracy or F1 score across seven datasets. When integrating the ZO optimizer with Parameter Efficient Fine Tuning (PEFT), LeZO exceeds MeZO on five datasets, highlighting the versatility and effectiveness of our sparsity scheme. Additionally, LeZO achieves over $2 \times$ computational acceleration across some tasks, with speedups increasing as layer sparsity exceeds 75%. Finally, LeZO achieves over $3 \times$ training acceleration for the OPT-13B model fine-tuning on the SST-2, BoolQ, and Copa datasets.

- In summary, the contributions of this paper are as follows:
 - We analyze the computational costs of MeZO during LLM fine-tuning and find that the perturbation and updating processes account for over 50% of the fine-tuning time cost.
 - We propose a dynamic layer-wise sparse and efficient zeroth-order optimizer that effectively reduces computational cost in fine-tuning LLMs, and achieve faster convergence than MeZO.
 - We evaluate LeZO on SuperGLUE and two generative tasks. This demonstrates its acceleration capability across different sparsity rates and showcases its performance improvement.
- 095 096

098

087

088

090

091

092

094

2 RELATED WORK

099 Forward-Backward Fine-Tuning Optimziers rely on the FO optimizer and compute gradients 100 through derivation (Robbins & Monro, 1951; Kingma & Ba, 2015), which results in high memory 101 usage. One approach to handle this problem is Parameter-Efficient Fine-Tuning (PEFT) (Xu et al., 102 2023). Relevant strategies involve adding trainable modules (Li & Liang, 2021; Hambardzumyan 103 et al., 2021) or learnable input embeddings (Houlsby et al., 2019) to the model. Moreover, the LoRA 104 series (Hu et al., 2022; Dettmers et al., 2024) adopt low-rank decomposition to reduce the size of 105 trainable modules. However, recent studies show that partial fine-tuning cannot match the performance of full-parameter fine-tuning (Ding et al., 2022; Pu et al., 2023). To handle this problem, 106 some methods enable full-parameter fine-tuning by adjusting parameters in an interleaved fashion. 107 For example, LISA (Pan et al., 2024) updates parameters of different layers in different iterations based on layer importance. Galore (Zhao et al., 2024) switches low-rank subspaces using the rank
 information while simultaneously modifying the learnable parameters. These techniques still require
 storing the intermediate states of FO optimizers, resulting in significant memory overhead.

111 Forward-Only Fine-tuning Optimziers do not require back-propagation for gradient computa-112 tion (Malladi et al., 2023; Zhang et al., 2024). Malladi et al. (2023) firstly introduce the ZO optimizer 113 for LLMs fine-tuning, known as MeZO, which significantly reduces memory costs. Because MeZO 114 employs SPSA to estimate gradients, it exhibits greater randomness compared to FO optimizers. 115 Consequently, it requires more computational cost to converge. To accelerate the convergence of 116 MeZO, Sparse-MeZO (Liu et al., 2024) only updates model parameters with small values to en-117 hance the effectiveness of ZO optimizer. Nevertheless, this approach necessitates ranking parameter 118 values and introduces a mask matrix; therefore, it results in additional memory and computational overhead. Guo et al. (2024) reduced the number of learnable parameters to one-thousandth. How-119 ever, it requires cloud computation to provide first-order gradient of the objective function, which is 120 unsuitable for more general scenarios. In summary, how to effectively accelerate the convergence 121 speed of ZO optimizers without incurring additional memory overhead still remains a challenge. 122

Our proposed LeZO adopts ZO optimization to update model parameters, reducing memory usage
 compared to forward-backward fine-tuning. By employing a dynamic layer-wise sparsity strategy,
 LeZO significantly reduces the number of trainable parameters and minimizes the computational
 costs during the forward-only fine-tuning. Thus, LeZO achieves high efficiency in both memory
 usage and computational cost.

128 129

130

3 EFFICIENCY ANALYSIS OF MEZO

In this section, we first revisit the classical ZO gradient estimator SPSA (Spall, 1992) and the ZO
 optimizer ZO-SGD (Malladi et al., 2023). Subsequently, we analyze the computational redundancy
 of MeZO (Malladi et al., 2023).

134 135

136

144 145

158

161

3.1 PRELIMINARIES

SPSA approximates gradients for multi-parameter systems by applying random perturbations, which
 eliminates the need of gradient back-propagation. It iteratively optimizes parameters by estimating
 gradients based on differences in function values around a given point. SPSA does not require the
 function to be smooth or convex, but it generally converges slowly.

Definition 1 (Simultaneous Perturbation Stochastic Approximation, SPSA (Spall, 1992)). *Given a group of parameters* $\theta \in \mathbb{R}^d$ *from a model and a loss function* \mathcal{L} *used for optimization, SPSA estimates the gradient on a mini-batch of data* \mathcal{B} *as follows:*

$$\widehat{\nabla}\mathcal{L}(\boldsymbol{\theta};\mathcal{B}) = \frac{\mathcal{L}(\boldsymbol{\theta} + \epsilon \boldsymbol{z};\mathcal{B}) - \mathcal{L}(\boldsymbol{\theta} - \epsilon \boldsymbol{z};\mathcal{B})}{2\epsilon} \boldsymbol{z} \approx \boldsymbol{z} \boldsymbol{z}^{\top} \nabla \mathcal{L}(\boldsymbol{\theta};\mathcal{B}),$$
(1)

where $z \in \mathbb{R}^d$ with $z \sim \mathcal{N}(0, I_d)$ and ϵ is a perturbation scale.

Malladi et al. (2023) applied SPSA (Definition 1) to optimize LLMs, using it to obtain approximate gradients of parameters. They modified the SGD algorithm (Robbins & Monro, 1951), creating ZO-SGD (Definition 2), which incorporates ZO differential gradients to update model parameters. MeZO employs both SPSA and ZO-SGD, and performs two forward passes in a single optimization step without caching the activation information of the parameters. Therefore, it reduces memory usage to the size occupied by the model parameters alone.

Definition 2 (Zeroth-Order Stochastic Gradient Descent, ZO-SGD (Malladi et al., 2023)). Given a learning rate η and a group of parameters $\theta \in \mathbb{R}^d$, the parameters at time t can be updated using the SPSA gradient estimate as follows:

$$\boldsymbol{\theta}_{t+1} = \boldsymbol{\theta}_t - \eta \widehat{\nabla} \mathcal{L}(\boldsymbol{\theta}; \mathcal{B}_t), \tag{2}$$

where \mathcal{B}_t is the mini-batch of data used at time t.

3.2 EFFICIENCY DILEMMA OF MEZO

Fine-tuning LLMs with MeZO is extremely time-consuming.
It often takes hundreds or even thousands of times longer than
using FO optimizers like SGD (Robbins & Monro, 1951) and
Adam (Kingma & Ba, 2015). We break down the MeZO computation process into three stages: forward pass, perturbation,
and updating. Algorithm 1 explains the computational operations for each of these stages. Figure 2 illustrates their relative
time cost in an optimization step.

170 Surprisingly, in classification tasks, the parameter perturbation 171 and updating stages account for more than one half of the over-172 all time. A common method to reduce the computational burden during optimization is to skip the calculations of certain 173 modules (Fan et al., 2019). However, this approach typically 174 introduces an inherent error during the forward pass (Evans & 175 Aamodt, 2021). This phenomenon is challenging to mitigate 176 and usually slows down the convergence rate. An alternative 177 strategy is to selectively optimize the parameters of specific 178 modules during the updating stage. Many studies have sub-179 stantiated the efficacy of this strategy in FO optimization (Liu 180



Figure 2: Proportion of computational time cost for each operation in a single step when fine-tuning the OPT-13b model on the SST-2 dataset using MeZO.

et al., 2021; Mi et al., 2022; Pan et al., 2024; Zhao et al., 2024). The above observations motivate us to implement sparse perturbation and updating for MeZO.

4 Methodology

183

184 185

187

188 189

190 191

192

193

194

195 196 197

199

200

201

207

209

210

211 212 In this section, we present the implementation details of LeZO. We explain how it achieves convergence and computation acceleration without additional memory consumption by updating only a subset of structured parameters during each optimization step.

4.1 DYNAMIC LAYER-WISE SPARSE ZEROTH-ORDER OPTIMIZATION

Building upon Equation (1), we apply structured pruning to the model parameter vector $\boldsymbol{\theta}$. In each iteration, we partition this vector into two distinct parts. We introduce a sparse function \mathcal{R} to extract sub-vectors from the input vector and stipulate this partitioning is structured by layer. Given a sparse rate ρ ($0 \le \rho \le 1$), a parameter vector $\boldsymbol{\theta}$, and a random number *s*, we obtain a new parameter vector $\boldsymbol{\theta}' \in \mathbb{R}^{\rho d}$, which is mathematically represented as:

$$\boldsymbol{\theta}' = \mathcal{R}(\boldsymbol{\theta}, \rho, s_t). \tag{3}$$

In each step, we maintain the number of elements in θ' . The random seed s_t is used to determine which layers' parameters are selected to form θ' .

Definition 3 (Layer-wise Sparse SPSA). During each step, the sparse parameter set is used as a condition for differential calculation in SPSA, which is defined as follows:

$$\widehat{\nabla}\mathcal{L}_{\mathcal{R}}(\boldsymbol{\theta};\mathcal{B}) = \frac{\mathcal{L}(\boldsymbol{\theta} + \epsilon \boldsymbol{z}';\mathcal{B}) - \mathcal{L}(\boldsymbol{\theta} - \epsilon \boldsymbol{z}';\mathcal{B})}{2\epsilon} \boldsymbol{z}', \qquad (4)$$

206 where $z' = \mathcal{R}(z, \rho, s_t)$.

208 Furthermore, incorporating layer-wise sparsity into ZO-SGD results in:

Definition 4 (LeZO-SGD). At time t, θ can be updated using the layer-wise sparse SPSA gradient estimation as:

$$\boldsymbol{\theta}_{t+1}^{'} = \boldsymbol{\theta}_{t}^{'} - \eta \widehat{\nabla} \mathcal{L}_{\mathcal{R}}(\boldsymbol{\theta}; \mathcal{B}_{t}).$$
(5)

Algorithm 1 presents the pseudocode for LeZO, employing a straightforward approach to implement the sparse function \mathcal{R} . We use **layers** as the fundamental unit for sparsity (Fan et al., 2020). By maintaining a subset *a* to store pruned layers, these layers are skipped during perturbation and parameter updating processes.

Algorithm 1: LeZO (Layer-wise Sparse and Efficient Zeroth-Order Opt	timiz	ation)
Require : parameters $\boldsymbol{\theta} \in \mathbb{R}^d$, loss $\mathcal{L} : \mathbb{R}^d \to \mathbb{R}$, step budget T , perturba B , learning rate schedule $\{\eta_t\}$, layer number N , and dropping layer number N .	ition imbe	scale μ , batch size or $n \in (0, N)$.
for $t = 1,, T$ do Sample batch $\mathcal{B} \subset \mathcal{D}$ and random seed s Randomly select <i>n</i> elements from $\{1,, N\}$ to generate subset <i>a</i> $\theta \leftarrow \texttt{PerturbParameters}(\theta, a, \mu, s)$ $\ell_+ \leftarrow \mathcal{L}(\theta; \mathcal{B})$		Perturbation Forward Pass
$ \begin{array}{c} \boldsymbol{\theta} \leftarrow \texttt{PerturbParameters} \left(\boldsymbol{\theta}, a, -2\mu, s \right) \\ \boldsymbol{\ell}_{-} \leftarrow \mathcal{L}(\boldsymbol{\theta}; \mathcal{B}) \\ \boldsymbol{\theta} \leftarrow \texttt{PerturbParameters} \left(\boldsymbol{\theta}, a, \mu, s \right) \end{array} $		Perturbation Forward Pass Perturbation
$\begin{array}{c c} \texttt{projected_grad} \leftarrow (\ell_+ - \ell)/(2\mu) \\ \texttt{Reset random number generator with seed } s \\ \texttt{for } \theta_i \in \pmb{\theta} \And i \notin a \ \texttt{do} \\ & z \sim \mathcal{N}(0,1) \\ & \theta_i \leftarrow \theta_i - \eta_t * \texttt{projected_grad} * z \\ \texttt{end} \end{array}$		⊳ Updating
end		
Subroutine PerturbParameters (θ , a , μ , s) Reset random number generator with seed s		
$\left \begin{array}{c} \text{for } \theta_i \in \boldsymbol{\theta} \And i \notin a \text{ do} \\ z \sim \mathcal{N}(0,1) \\ \theta_i \leftarrow \theta_i + \mu z \end{array}\right $		
end return θ		

Remark 1. As illustrated in Algorithm 1, it is evident that LeZO introduces a layer-wise selection operation \mathcal{R} compared to MeZO. Given the scale of parameters in the billions, this addition has a negligible computational impact. We have the following remarks. (1) During the Forward pass, LeZO's computation aligns with MeZO, incurring no additional memory or computational overhead. (2) During the perturbation process, we skip the parameter perturbation in the sparsified layers, thereby reducing the computational load. Both LeZO and MeZO utilize a random seed for fixed perturbation z. This approach eliminates the need for additional cache storage for perturbation information, thereby not increasing memory overhead. The same principle applies to the updating process. Therefore, it is evident that LeZO reduces computational overhead through the introduction of layer-wise parameter sparsification during the perturbation and updating processes and incurs no additional memory overhead.

4.2 CONVERGENCE ANALYSIS OF LEZO

When examining the individual optimization steps, LeZO effectively updates a sub-network, similar to Sparse-MeZO (Liu et al., 2024). We aim to demonstrate that as model parameters converge towards local optimal values, there exists a value σ^2 such that the speed at which the loss approaches a local optimum is directly proportional to the number of parameters. This implies that the parameter convergence rate is related to ρ .

Lemma 1 (Unbiased Estimation of Sparse Gradient). Let $\mathcal{L}_{\mathcal{R}} = \mathbb{E}_{\mathcal{R}}[\mathcal{L}(\boldsymbol{\theta} + \epsilon \boldsymbol{z}')]$. The relationship between the model's sparse gradient $\widehat{\nabla}_{\boldsymbol{\theta}'}\mathcal{L}_{\mathcal{R}}(\boldsymbol{\theta})$ and the estimated ZO sparse gradient $\widehat{\nabla}\mathcal{L}_{\mathcal{R}}(\boldsymbol{\theta})$ can be expressed as:

$$\widehat{\nabla}_{\boldsymbol{\theta}'} \mathcal{L}_{\mathcal{R}}(\boldsymbol{\theta}) = \mathcal{R}(\nabla \mathcal{L}_{\mathcal{R}}(\boldsymbol{\theta})) = \nabla_{\boldsymbol{\theta}'} \mathbb{E}_{\mathcal{R}}[\mathcal{L}(\boldsymbol{\theta} + \epsilon \boldsymbol{z}')] \\ = \mathbb{E}_{\mathcal{R}}[\frac{\mathcal{L}(\boldsymbol{\theta} + \epsilon \boldsymbol{z}') - \mathcal{L}(\boldsymbol{\theta} - \epsilon \boldsymbol{z}')}{2\epsilon} \boldsymbol{z}'] = \mathbb{E}_{\mathcal{R}}[\widehat{\nabla} \mathcal{L}_{\mathcal{R}}(\boldsymbol{\theta})].$$
(6)

From Lemma 1, we can see that the estimated gradient by sparse ZO is an unbiased estimation of the model's sparse parameter gradient.

273	Task	SST-2	RTE	CB	BoolQ	WSC	WIC	Copa	SQuAD	AVG
274	Task type			—- classif	ication —			multiple choice	generation	<i>n</i> vo.
214	Zero-Shot	58.8	59.6	46.4	59.0	38.5	55.0	80.0	46.2	55.4
275	ICL	87.0	62.1	57.1	66.9	39.4	50.5	87.0	75.9	65.7
276	FT (12× memory)	92.0	70.8	83.9	77.1	63.5	70.1	79.0	84.9	77.7
277	MeZO (Malladi et al., 2023)	91.4	66.1	67.9	67.6	63.5	61.1	88.0	84.7	73.8
278	MeZO (Reproduce)	91.1 ± 0.1	70.8±1.4	68.8±1.8	69.1±0.6	63.5±1.2	58.7±0.9	87.0±1.6	84.2±1.0	74.2±1.1
279	SSZO	93.0±0.3	71.4±2.1	69.2±1.7	74.3±0.8	62.7±0.4	60.7±1.5	87.2±0.8	84.3±0.7	75.4±1.1

270 Table 1: Experiments on the OPT-13b model. SSZO sparsifies 75% of the layers (30 layers out of 271 40). Two A100-40G GPUs are used for BoolQ and SQuAD tasks.

Task Task type	SST-2	RTE	СВ	BoolQ classficatio	WSC	WIC	MultiRC	Copa —— mult	ReCoRD iple choice ——	SQuAD — gener	DROP ation —-
zero-shot	53.6	53.4	39.3	61.1	43.3	57.5	45.4	75.0	70.6	27.2	11.2
ICL	67.7	53.1	44.6	67.2	53.8	56.0	44.7	70.0	69.9	59.0	20.3
MeZO	89.7±1.3	65.4±2.5	69.3±3.0	64.1±0.8	63.5±0.0	58.7±0.5	60.4±1.5	76.0±1.9	71.6±0.5	76.9±0.8	23.1±0.3
LeZO	91.9±0.3	64.5±2.4	69.6±1.8	65.3±1.2	63.3±1.1	59.4±1.3	61.4±1.2	76.8±2.3	71.7±0.4	76.8±0.7	24.1±0.8

Assumption 1 (Lipschitz Continuous). Let $\nabla \mathcal{L}(\theta; x)$ denotes the first-order gradient of \mathcal{L} with respect to $\boldsymbol{\theta}$ at x. \mathcal{L} satisfies Lipschitz continuity, then

$$\|\nabla \mathcal{L}(\boldsymbol{\theta}; \boldsymbol{x}) - \nabla \mathcal{L}(\boldsymbol{\theta}_t; \boldsymbol{x})\| \le \frac{L(l)}{2} \|\boldsymbol{\theta} - \boldsymbol{\theta}_t\|^2,$$
(7)

where L(l) is a constant ensuring \mathcal{L} satisfies Lipschitz continuity.

Assuming that the loss function $\mathcal{L}(\theta; x)$ is Lipschitz continuous (Assumption 1), we can calculate the norm distance between the sparse ZO estimated gradient and the true sparse gradient $\nabla \mathcal{L}_{\mathcal{R}}(\theta)$. This distance, denoted as $\|\widehat{\nabla}_{\theta'}\mathcal{L}_{\mathcal{R}}(\theta) - \nabla \mathcal{L}_{\mathcal{R}}(\theta)\|$, helps establish their relationship using Lemma 1. **Lemma 2** (Relationship between Sparse Gradient and Estimate Value). Let the loss function $\mathcal{L}(\theta; x)$ to be Lipschitz continuous. According to Minkowski inequality, we have

> $\|\nabla \mathcal{L}_{\mathcal{R}}(\boldsymbol{\theta})\|^{2} \leq 2\|\widehat{\nabla}_{\boldsymbol{\theta}'}\mathcal{L}_{\mathcal{R}}(\boldsymbol{\theta})\|^{2} + \frac{\epsilon^{2}L^{2}(l)}{2}(\rho d + 4)^{3},$ (8)

where $\nabla \mathcal{L}_{\mathcal{R}}(\boldsymbol{\theta}) = \mathcal{R}(\nabla \mathcal{L}(\boldsymbol{\theta})).$

304 Finally, we can obtain the convergence rate of LeZO.

305 **Lemma 3** (Convergence of LeZO). Assuming a sequence of generated parameters $\{\theta_t\}_{t>0}$ in LeZO. 306 We have $[||\nabla a(a)||^2] < 2$ (9)

307

308

309

272

281

283 284 285

287 288

289

290 291

292

293

295

296

297

298

299

300

301 302

303

$$\mathbb{E}_{\mathcal{R},x}[\|\nabla_{\boldsymbol{\theta}}\mathcal{L}_{\mathcal{R}}(\boldsymbol{\theta}_{T})\|^{2}] \leq \sigma^{2}$$

for any $T = \mathcal{O}(\frac{\rho dL}{\sigma^2})$. t represents the step of fine-tuning and $L(l) \leq L$ for all $\mathcal{L}(\boldsymbol{\theta}_t)$.

310 From Lemma 3, it is apparent that as the sparsity ratio ρ decreases, the upper bound on the time re-311 quired to converge to the expected value also decreases. Hence, parameter sparsity and the smoothness of the objective function can enhance the convergence speed. 312

- 5 EXPERIMENTS
- 315 316 317

313 314

In this section, we present extensive experiments to validate the effectiveness of the LeZO algorithm.

318 5.1 EXPERIMENTAL SETTING 319

320 Models and Datasets. To evaluate the efficacy of LeZO, we follow the validation methodology 321 by Malladi et al. (2023). We conduct experiments on the OPT model family (Zhang et al., 2022) at various scales: 1.3 billion, 13 billion, and 30 billion parameters, with 32, 40, and 48 Trans-322 former blocks, respectively. These models utilize the core components of the Transformer archi-323 tecture (Vaswani et al., 2017), including self-attention mechanisms, feed-forward neural networks, Table 4: Fine-tuning performance of the ZO optimizer with PEFT. †indicates the results reported in
 the paper of Malladi et al. (2023). LeZO (LoRA/prefix) outperforms MeZO (LoRA/prefix) on 4 out
 of 5 tasks. LeZO (LoRA) sparsifies 50% of layers, while LeZO (prefix) sparsifies 75%.

327						
328	Task	SST-2	CB	BoolQ	Copa	SQuAD
329	MeZO (LoRA)†	89.6	66.1	73.8	84.0	83.8
330	MeZO (LoRA)	90.6±1.6	70.4±1.6	71.8±1.2	86.0±1.2	81.1±0.8
331	MeZO (prefix)†	90.7	69.6	73.1	87.0	84.2
332	MeZO (prefix)	90.7±0.9	67.9±2.5	73.2±0.4	87.4±1.7	83.3±1.0
333	LeZO (LoRA)	92.3±0.5	71.1±0.8	73.7±0.8	86.4±1.5	81.5±1.1
334	LeZO (prefix)	92.1±0.7	69.7±2.5	74.5±0.6	87.6±1.1	83.1±0.4

335

residual connections, and layer normalization within each block. For downstream tasks, we use the
 SuperGLUE benchmark (Wang et al., 2019), which includes classification tasks, multiple-choice
 tasks, and two question-answering tasks.

339 Methods. We utilize zero-shot and 32-shot In-Content 340 Learning (ICL) (Brown et al., 2020) alongside fine-341 tuning with Adamw (FT) as comparison. To demonstrate LeZO's ability to accelerate model convergence and re-342 duce time cost per training step without additional mem-343 ory costs, we compared it with MeZO (Malladi et al., 344 2023). In addition to full-parameter fine-tuning, we in-345 tegrate LeZO with two PEFT methods: LoRA (Hu et al., 346 2022) and prefix tuning (Li & Liang, 2021), to further 347 enhance fine-tuning efficiency. 348

Hyper-parameters. We follow most of the experimental settings used in MeZO, including configurations for LoRA and prefix tuning, and methods for selecting training and testing data. We conduct a grid search for the learning rate and perturbation scale to identify the optimate of the perturbation scale

Table 3: Experiments on the OPT-30b model. LeZO sparsifies 75% of the layers (36 layers out of 48). We report the best results of MeZO and MeZO with prefix from the paper of Malladi et al. (2023). Two A100-40G GPUs are used for SST2, and four for BoolQ.

Task	SST-2	BoolQ
zero-shot	56.7	39.1
ICL	81.9	66.2
MeZO/MeZO (prefix)	90.6	73.5
MeZO (reproduce)	90.3±0.5	69.5±0.4
LeZO	92.8±0.6	73.7±0.9

mal parameters for each experimental group. We use the experimental code for ZO optimization 354 provided by Zhang et al. (2024) and set the validation steps to 2000. Experiments are conducted on 355 an A100-40G GPU. However, due to its memory limitations compare to the A100-80G GPU used 356 by Malladi et al. (2023), we employ a multi-card parallel strategy for fine-tuning the OPT-13b and 357 30b model on certain downstream tasks to ensure parameter consistency. We select the final model 358 based on the checkpoint with the lowest validation loss. All numerical results report in the experiments are reproduced. We conduct parameter fine-tuning experiments using five different random 359 seeds, reporting the average and standard deviation of these results. Detailed settings are provided in 360 Appendix A. Unless otherwise specified, all baseline experiments are conducted using the OPT-13b 361 model with 75% sparsity (30 out of 40 layers). 362

363 364

5.2 MAIN RESULT

Comparison between LeZO and Other Methods. LeZO outperforms both non-training methods 366 and MeZO, and surpasses the FO optimizer on some tasks. Table 1 presents a comparison between 367 LeZO and the baseline methods such as zero-shot, the training-free method (ICL), Fine-Tuning 368 (FT) with AdamW for FO optimization, and MeZO across various downstream tasks. We obtain the 369 following observations. 1) LeZO significantly outperforms non-training methods such as zero-shot 370 and ICL on all tasks. This demonstrates LeZO's superior convergence capabilities during training. 371 2) LeZO surpasses MeZO in 7 out of 8 tasks, with an average improvement exceeding 1% across 372 all tasks. Notable enhancements are observed on the SST-2, BoolQ, and WIC datasets, with average 373 accuracy increases over 2% across five experiments. This indicates that layer-wise sparsity schemes 374 can effectively accelerate convergence. However, there is a slight performance degradation on the 375 WSC task compared to MeZO, likely due to the inherent instability of ZO optimization methods and MeZO's sensitivity to prompts. It is speculated that LeZO may inherit these drawbacks, as 376 the fundamental differential gradient computation process in SPSA remains unchanged. 3) LeZO 377 marginally outperforms FT on the SST-2, RTE, and Copa datasets, while demonstrating comparable





390 Figure 3: Impact of learning rate and sparsity ratio on fine-tuning models using LeZO on the SST-2 task. Results the OPT-13b model using LeZO on the 391 with accuracy exceeding 90% are displayed. Experiments 392 were conducted with a single random seed. "Dropout Number" indicates the number of sparse layers in the OPT 394 13B model. To enhance clarity, the learning rates are logarithmically scaled. The curves mapping onto the straight 396 lines in the lower plane are used to illustrate the range 397 of learning rates that lead to improved performance after sparsifying different numbers of layers in the model. The 399 performance of MeZO performance is the red line. 400

Figure 4: Correlation between the sparsity ratio and runtime in fine-tuning SST-2 task. This figure presents the optimal experimental results at various sparsity levels from Figure 3, annotated with different colored pentagrams. The purple line indicates the total time required for fine-tuning in the corresponding experiments.

401 performance on the other datasets. Notably, LeZO and MeZO have the same memory consumption, and FT incurs a memory cost twelve times higher than LeZO. 402

403 Performance on Multi-Size Models. LeZO achieves a notable performance improvement over 404 MeZO across models of varying scales. Tables 2 and 3 display the performance of LeZO and 405 other methods on the OPT-1.3b and OPT-30b models, respectively. When fine-tuning OPT-1.3b on 406 downstream tasks, LeZO achieves superior results in 9 out of 11 tasks. It achieves comparable results on the WSC task and its performance drop is less than 1% on RTE. When fine-tuning OPT-30b on 407 the SST-2 and BoolQ datasets, LeZO demonstrates significant improvements compared to MeZO. 408 The experiments across Tables 1, 2, and 3 reveal that LeZO achieves more pronounced performance 409 gains on larger LLMs. This phenomenon may relate to the stability of model convergence, with 410 larger models being more stable. Additionally, LeZO demonstrates greater stability compared to 411 MeZO. 412

Performance of Combining ZO optimizer 413 and PEFT. LeZO can be effectively integrated 414 with PEFT methods. We combine ZO op-415 timizers with PEFT by updating the train-416 able parameters introduced by PEFT across 417 different methods. Table 4 showcases the 418 fine-tuning results of LeZO and MeZO com-419 bined with the PEFT across different datasets. 420 LeZO (LoRA) outperforms MeZO (LoRA) on 421 all five datasets, while LeZO (prefix) outper-422 forms MeZO (LoRA) on 4 out of 5 datasets. 423 Additionally, the combination of LeZO with PEFT results in performance improvements ex-424 ceeding 1% on the SST-2, CB, and BoolQ tasks 425 compared to MeZO. This indicates that LeZO 426 enhances fine-tuning effectiveness even when 427 combined with PEFT. 428



Figure 5: Comparison in computation efficiency between LeZO and MeZO on various tasks.

- **Convergence and Computation Speedup by** 429 LeZO. LeZO is more efficient than MeZO. Fig-430
- ure 1 illustrates the variation in accuracy on the test set when fine-tuning OPT-13b using LeZO on 431 the SST-2 datasets. LeZO achieves $2 \times$ computation speedup and $1.7 \times$ convergence speedup. In
 - 8

432 Figure 5, we present the convergence and computation speedups obtained when using LeZO to fine-433 tune OPT-13b on eight downstream tasks. The experimental results show that with the layer-wise 434 sparsity policy, LeZO achieves effective convergence and computation speedups. In conclusion, 435 LeZO achieves more efficient LLM fine-tuning than MeZO.

436 437

438

439

5.3 HYPERPARAMETER ANALYSIS

440 Influence of Learning Rate and Sparsity Rate. The layer-wise sparsity scheme we use shows 441 how the number of sparse layers affects LeZO's performance. Our experiments show that with an increase in sparsity rate (more layers being sparsified), the ZO optimization process needs a higher 442 learning rate. Figure 3 illustrates the relationship between the sparsity rate (Dropout Number), 443 learning rate, and model accuracy on downstream tasks. From the figure, we obtain several key 444 observations: 1) LeZO demonstrates significantly superior optimization effectiveness compared to 445 MeZO (Dropout Number=0). Under full-parameter fine-tuning, LeZO consistently outperforms 446 MeZO regardless of the number of sparsified layers. 2) LeZO enhances the robustness of the ZO 447 optimization process. Comparing the lengths of the dashed lines, it can be observed that as the 448 sparsity ratio increases, the range of learning rates that achieve over 90% accuracy on the SST-2 task 449 expands. This indicates that the robustness of LeZO improves with an increased sparsity ratio. 3) 450 The behavior of LeZO markedly differs from fine-tuning only a subset of parameters. Performance 451 experiences a substantial decline when the sparsity rate reaches $\rho = 1$ (sparse 40 layers in OPT-13b, 452 while fine-tuning solely the embedding and linear layers).

453 Impact of Downstream Tasks on Compu-454 tational Speedup. LeZO achieves varying 455 speedup ratios across different tasks, as shown 456 in Figure 5, which depicts the computation and 457 convergence speedup ratios obtained by LeZO across all tasks. The differences in computa-458 tional speedup ratios come from the varying 459 proportions of the forward propagation process 460 within the total computation. The speed of for-461 ward propagation depends on the average token 462 length of different tasks within the same model 463 framework. Figure 6 illustrates the relationship 464 between the average input token length of dif-465 ferent datasets and the computational speedup 466 achieved by LeZO. When input token length in-467 creases, computation speedup decreases. With 468 a constant sparsity ratio, LeZO's floating-point 469 computational savings per downstream task re-



Figure 6: The relationship between the length of input tokens and the computational speedup achieved by LeZO.

main fixed. Consequently, LeZO's computation acceleration is more suited for relatively straight-470 forward training samples. However, a similar trend is not observed in the convergence acceleration 471 of LeZO across different tasks. This disparity may be attributed to the varying difficulty levels of 472 tasks, their convergence characteristics, and the distributions of training and extracted testing sets. 473

- 474 475
- CONCLUSION 6

477

476

478 In this study, we introduce LeZO, a layer-wise sparse and efficient ZO optimizer for fine-tuning 479 LLMs. We integrate the layer-wise pruning scheme into the ZO optimization process, which 480 achieves both computation and convergence speedups in LLMs training without additional mem-481 ory costs. Moreover, LeZO combines effectively with existing PEFT methods to further enhance 482 acceleration. We conduct fine-tuning experiments on models from the OPT family using the Super-483 GLUE benchmark and two question-answering datasets. Empirical results show that LeZO significantly reduces training time and improves performance. Overall, LeZO provides an efficient and 484 effective approach to ZO fine-tuning. It reduces the computational load caused by full-parameter 485 perturbations in zeroth-order optimization without introducing additional memory overhead.

486 REFERENCES

- Krishnakumar Balasubramanian and Saeed Ghadimi. Zeroth-order (non)-convex stochastic opti mization via conditional gradient and gradient updates. In *Advances in Neural Information Processing Systems*, volume 31, 2018.
- Andrew Brock, Theodore Lim, James Millar Ritchie, and Nicholas J Weston. Freezeout: Accelerate
 training by progressively freezing layers. In *NIPS 2017 Workshop on Optimization: 10th NIPS* Workshop on Optimization for Machine Learning, 2017.
- Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal,
 Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, et al. Language models are
 few-shot learners. *Advances in neural information processing systems*, 33:1877–1901, 2020.
- HanQin Cai, Daniel Mckenzie, Wotao Yin, and Zhenliang Zhang. Zeroth-order regularized opti mization (zoro): Approximately sparse gradients and adaptive sampling. *SIAM Journal on Opti- mization*, 32(2):687–714, 2022.
- Tim Dettmers, Artidoro Pagnoni, Ari Holtzman, and Luke Zettlemoyer. Qlora: Efficient finetuning of quantized llms. *Advances in Neural Information Processing Systems*, 36, 2024.
- Ning Ding, Yujia Qin, Guang Yang, Fuchao Wei, Zonghan Yang, Yusheng Su, Shengding Hu, Yulin
 Chen, Chi-Min Chan, Weize Chen, et al. Delta tuning: A comprehensive study of parameter
 efficient methods for pre-trained language models. *arXiv preprint arXiv:2203.06904*, 2022.
- 507 Dheeru Dua, Yizhong Wang, Pradeep Dasigi, Gabriel Stanovsky, Sameer Singh, and Matt Gardner.
 508 Drop: A reading comprehension benchmark requiring discrete reasoning over paragraphs. In
 509 Proceedings of the 2019 Conference of the North American Chapter of the Association for Com 510 putational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers), pp.
 511 2368–2378, 2019.
- R David Evans and Tor Aamodt. Ac-gc: Lossy activation compression with guaranteed convergence. *Advances in Neural Information Processing Systems*, 34:27434–27448, 2021.
- Angela Fan, Edouard Grave, and Armand Joulin. Reducing transformer depth on demand with
 structured dropout. *arXiv preprint arXiv:1909.11556*, 2019.
- Angela Fan, Edouard Grave, and Armand Joulin. Reducing transformer depth on demand with structured dropout. In *International Conference on Learning Representations*, 2020.
- Wentao Guo, Jikai Long, Yimeng Zeng, Zirui Liu, Xinyu Yang, Yide Ran, Jacob R Gardner, Osbert Bastani, Christopher De Sa, Xiaodong Yu, et al. Zeroth-order fine-tuning of llms with extreme sparsity. *arXiv preprint arXiv:2406.02913*, 2024.
- Karen Hambardzumyan, Hrant Khachatrian, and Jonathan May. Warp: Word-level adversarial re programming. In *Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume Long Papers)*, pp. 4921–4933, 2021.
- Jordan Hoffmann, Sebastian Borgeaud, Arthur Mensch, Elena Buchatskaya, Trevor Cai, Eliza Rutherford, Diego de las Casas, Lisa Anne Hendricks, Johannes Welbl, Aidan Clark, Tom Hennigan, Eric Noland, Katherine Millican, George van den Driessche, Bogdan Damoc, Aurelia Guy, Simon Osindero, Karen Simonyan, Erich Elsen, Oriol Vinyals, Jack William Rae, and Laurent Sifre. An empirical analysis of compute-optimal large language model training. In Alice H. Oh, Alekh Agarwal, Danielle Belgrave, and Kyunghyun Cho (eds.), *Advances in Neural Information Processing Systems*, 2022.
- Neil Houlsby, Andrei Giurgiu, Stanislaw Jastrzebski, Bruna Morrone, Quentin De Laroussilhe, Andrea Gesmundo, Mona Attariyan, and Sylvain Gelly. Parameter-efficient transfer learning for NLP. In *Proceedings of the 36th International Conference on Machine Learning*, 2019.
- Edward J Hu, Yelong Shen, Phillip Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang,
 and Weizhu Chen. LoRA: Low-rank adaptation of large language models. In *International Conference on Learning Representations*, 2022.

555

559

560

561

565

573

577

578

579

580

585

586

587 588

589

- Jared Kaplan, Sam McCandlish, Tom Henighan, Tom B Brown, Benjamin Chess, Rewon Child,
 Scott Gray, Alec Radford, Jeffrey Wu, and Dario Amodei. Scaling laws for neural language
 models. arXiv preprint arXiv:2001.08361, 2020.
- Diederik P. Kingma and Jimmy Lei Ba. Adam: A method for stochastic optimization. In *Interna- tional Conference on Learning Representations*, pp. 1–13, 2015.
- 546
 547
 548
 549
 549
 549
 540
 541
 541
 541
 542
 544
 545
 544
 545
 546
 547
 548
 549
 549
 549
 549
 549
 549
 549
 549
 540
 540
 541
 541
 541
 542
 542
 544
 545
 545
 546
 547
 546
 547
 548
 549
 549
 549
 549
 549
 549
 549
 549
 540
 540
 541
 541
 541
 542
 542
 542
 543
 544
 544
 545
 545
 546
 547
 547
 547
 548
 549
 549
 549
 549
 549
 549
 549
 540
 540
 541
 541
 541
 542
 542
 542
 544
 545
 545
 546
 547
 547
 547
 548
 549
 549
 549
 549
 549
 549
 549
 540
 540
 541
 541
 541
 542
 542
 542
 542
 542
 542
 542
 542
 542
 542
 542
 542
 542
 542
 542
 542
 542
 542
 542
 542
 542
 542
 542
 542
 542
 542
 542
 542
 542
 542
 542
 542
 542
 542
 542
 542
 542
 542
 542
 542
 542
 542
 542
 542
 542
 542
 542
 542
 542
 542
 544
 544
 544
 544
 544
- Xiang Lisa Li and Percy Liang. Prefix-tuning: Optimizing continuous prompts for generation. In Chengqing Zong, Fei Xia, Wenjie Li, and Roberto Navigli (eds.), *Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume 1: Long Papers)*, pp. 4582–4597, Online, August 2021. Association for Computational Linguistics. doi: 10.18653/v1/2021.acl-long.353.
- Yong Liu, Zirui Zhu, Chaoyu Gong, Minhao Cheng, Cho-Jui Hsieh, and Yang You. Sparse
 mezo: Less parameters for better performance in zeroth-order llm fine-tuning. *arXiv preprint arXiv:2402.15751*, 2024.
 - Yuhan Liu, Saurabh Agarwal, and Shivaram Venkataraman. Autofreeze: Automatically freezing model blocks to accelerate fine-tuning. *arXiv preprint arXiv:2102.01386*, 2021.
- Sadhika Malladi, Tianyu Gao, Eshaan Nichani, Alex Damian, Jason D. Lee, Danqi Chen, and Sanjeev Arora. Fine-tuning language models with just forward passes. In *Thirty-seventh Conference on Neural Information Processing Systems*, 2023.
- Peng Mi, Li Shen, Tianhe Ren, Yiyi Zhou, Xiaoshuai Sun, Rongrong Ji, and Dacheng Tao. Make
 sharpness-aware minimization stronger: A sparsified perturbation approach. In Alice H. Oh,
 Alekh Agarwal, Danielle Belgrave, and Kyunghyun Cho (eds.), Advances in Neural Information *Processing Systems*, 2022.
- ⁵⁷⁰ Rui Pan, Xiang Liu, Shizhe Diao, Renjie Pi, Jipeng Zhang, Chi Han, and Tong Zhang. Lisa: Layer⁵⁷¹ wise importance sampling for memory-efficient large language model fine-tuning. *arXiv preprint*⁵⁷² *arXiv:2403.17919*, 2024.
- George Pu, Anirudh Jain, Jihan Yin, and Russell Kaplan. Empirical analysis of the strengths and
 weaknesses of PEFT techniques for LLMs. In *ICLR 2023 Workshop on Mathematical and Empirical Understanding of Foundation Models*, 2023.
 - Pranav Rajpurkar, Jian Zhang, Konstantin Lopyrev, and Percy Liang. Squad: 100,000+ questions for machine comprehension of text. In *Proceedings of the 2016 Conference on Empirical Methods in Natural Language Processing*, pp. 2383–2392, 2016.
- Jun Rao, Fei Wang, Liang Ding, Shuhan Qi, Yibing Zhan, Weifeng Liu, and Dacheng Tao. Where
 does the performance improvement come from? -a reproducibility concern about image-text retrieval. In *Proceedings of the 45th international ACM SIGIR conference on research and devel-*opment in information retrieval, pp. 2727–2737, 2022.
 - Herbert Robbins and Sutton Monro. A stochastic approximation method. *The annals of mathematical statistics*, pp. 400–407, 1951.
 - J.C. Spall. Multivariate stochastic approximation using a simultaneous perturbation gradient approximation. *IEEE Transactions on Automatic Control*, 37(3):332–341, 1992.
- Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Ł ukasz Kaiser, and Illia Polosukhin. Attention is all you need. In I. Guyon, U. Von Luxburg, S. Bengio, H. Wallach, R. Fergus, S. Vishwanathan, and R. Garnett (eds.), *Advances in Neural Information Processing Systems*, volume 30. Curran Associates, Inc., 2017.

- Alex Wang, Yada Pruksachatkun, Nikita Nangia, Amanpreet Singh, Julian Michael, Felix Hill, Omer
 Levy, and Samuel Bowman. Superglue: A stickier benchmark for general-purpose language
 understanding systems. In H. Wallach, H. Larochelle, A. Beygelzimer, F. d'Alché-Buc, E. Fox,
 and R. Garnett (eds.), Advances in Neural Information Processing Systems, volume 32. Curran
 Associates, Inc., 2019.
- Fei Wang, Liang Ding, Jun Rao, Ye Liu, Li Shen, and Changxing Ding. Can linguistic knowledge improve multimodal alignment in vision-language pretraining? *arXiv preprint arXiv:2308.12898*, 2023.
- Yining Wang, Simon Du, Sivaraman Balakrishnan, and Aarti Singh. Stochastic zeroth-order op timization in high dimensions. In *Proceedings of the Twenty-First International Conference on Artificial Intelligence and Statistics*, volume 84 of *Proceedings of Machine Learning Research*,
 pp. 1356–1365. PMLR, 2018.
- Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, brian ichter, Fei Xia, Ed H. Chi, Quoc V Le, and Denny Zhou. Chain of thought prompting elicits reasoning in large language models. In Alice H. Oh, Alekh Agarwal, Danielle Belgrave, and Kyunghyun Cho (eds.), *Advances in Neural Information Processing Systems*, 2022.
- Lingling Xu, Haoran Xie, Si-Zhao Joe Qin, Xiaohui Tao, and Fu Lee Wang. Parameter-efficient fine-tuning methods for pretrained language models: A critical review and assessment. *arXiv preprint arXiv:2312.12148*, 2023.
- Susan Zhang, Stephen Roller, Naman Goyal, Mikel Artetxe, Moya Chen, Shuohui Chen, Christopher Dewan, Mona Diab, Xian Li, Xi Victoria Lin, et al. Opt: Open pre-trained transformer
 language models. *arXiv preprint arXiv:2205.01068*, 2022.
- Yihua Zhang, Pingzhi Li, Junyuan Hong, Jiaxiang Li, Yimeng Zhang, Wenqing Zheng, Pin-Yu Chen, Jason D Lee, Wotao Yin, Mingyi Hong, et al. Revisiting zeroth-order optimization for memory-efficient llm fine-tuning: A benchmark. *arXiv preprint arXiv:2402.11592*, 2024.
- Jiawei Zhao, Zhenyu Zhang, Beidi Chen, Zhangyang Wang, Anima Anandkumar, and Yuandong
 Tian. Galore: Memory-efficient LLM training by gradient low-rank projection. In *Forty-first International Conference on Machine Learning*, 2024.

648 DETAILED EXPERIMENTAL SETTINGS А 649

650 We aligned the hyperparameter configurations for our primary experiments with the standards set 651 by MeZO, as described by Malladi et al. (2023). Our experimental framework is based on the work 652 of Zhang et al. (2024). Our approach incorporates a sparsity rate in LeZO. Our experiments indicate 653 that this requires larger learning rates to enhance model convergence. Consequently, the grid search 654 ranges for learning rates differ between LeZO and MeZO in the experiments presented in Table 1, 2, 3, and 4. We strictly adhered to the settings outlined by Malladi et al. (2023) when replicating 655 MeZO experiments. The hyperparameter search ranges for grid search are detailed in Table 5. We 656 conducted five trials with fixed random seeds for each set of experiments to compute average results 657 and standard deviations. This facilitated a robust comparison of the efficiency of different methods. 658 In the experiments depicted in Figure 3, due to the extensive number of trials, we conducted a 659 single experiment with one random seed to obtain results under various hyperparameter settings. To 660 streamline the evaluation process without compromising the representation of convergence across 661 diverse models and datasets, we set the testing interval at 2000 steps. The convergence behavior of 662 LeZO is contingent upon that of MeZO and is influenced by prompts; therefore, we did not conduct 663 prompt removal experiments similar to those in MeZO. The prompts used in the experiments are 664 identical to those in MeZO.

Experiment	Hyperparameters	Values
LeZO	Batch size	16
	Learning rate	$\{1e-6, 7e-7\}$ for OPT-13b/30b, $\{3e-6, 1e-6\}$ for OPT-1.3b
	ϵ	1e-3 0.75
	Sparse Kate	0.75
MeZO	Batch size	
	Learning rate	$\{1e-6, 1e-7\}$ or $\{1e-6, 5e-7, 1e-7\}$ for RTE and SQUAD
	e	16-2
LeZO (prefix)	Batch size	16
	Learning rate	$\{3e-2, 1e-2\}$
	Enorse Dete	1e-1 0.75
	# prefix tokens	5
		10
MeZO (prefix)	Batch size	$\begin{bmatrix} 10 \\ 10 \\ 2 \\ 10 \\ 2 \\ 10 \\ 2 \\ 10 \\ 2 \\ 10 \\ 2 \\ 50 \\ 2 \\ 10 \\ 2 \\ 50 \\ 2 \\ 10 \\ 2 \\ 50 \\ 2 \\ 10 \\ 10 \\ 10 \\ 10 \\ 10 \\ 10 \\ 10$
	Leanning rate	1e-2, $1e-3$ or $1e-2$, $1e-2$, $5e-3$ for SQUAD
	# prefix tokens	5
LeZO (LoRA)	Batch size	16
	Learning rate	$\{3e-5, 5e-5, 7e-5\}$
	ϵ	1e-2
	Sparse Rate	0.50
	(r, α)	(8, 16)
MeZO (LoRA)	Batch size	16
	Learning rate	$\{1e-4, 5e-5\}$ or $\{1e-4, 5e-5, 1e-5\}$ for SQuAD
	ϵ	1e-2
	(r, α)	(8,16)
FT with Adam	Batch size	8
	Learning Rate	$\{1e-5, 5e-5, 8e-5\}$

Table 5: The hyperparameter grids used for the experiments. Weight decay is set to 0. FT uses 5 epochs and linear scheduled learning rates. ZO optimizers use 20K steps and constant learning rates. We check validation performance and save the best checkpoint every 1/10 total training steps (2K).

В Proof

665

694

695

696 697 698

699 700

From a single-step perspective, LeZO also involves updating a sub-network, which aligns with the 701 theory of Sparse-MeZO (Liu et al., 2024). Therefore, the following proof draws inspiration from it. ~

Proof of Lemma 1. Let $\mathcal{L}_{\boldsymbol{z}}(\boldsymbol{\theta})$ be the expectation of $\mathcal{L}(\boldsymbol{\theta} + \epsilon \boldsymbol{z})$:

$$\mathcal{L}_{\mathcal{R}}(\boldsymbol{\theta}) := \mathbb{E}_{z}[\mathcal{L}(\boldsymbol{\theta} + \epsilon \mathcal{R}(\boldsymbol{z}))] = \mathbb{E}_{\mathcal{R}}[\mathcal{L}(\boldsymbol{\theta} + \epsilon \boldsymbol{z}')].$$
(10)

,

Then,

$$\begin{aligned} \widehat{\nabla}_{\theta'} \mathcal{L}_{\mathcal{R}}(\theta) &= \widehat{\nabla}_{\theta'} \mathbb{E}_{\mathcal{R}} [\mathcal{L}(\theta + \epsilon z')] \\ &= \widehat{\nabla}_{\theta'} \int_{z'} \mathrm{pdf}_{z'}(z) \mathcal{L}(\theta + \epsilon z) dz \\ &= \mathcal{R}(\nabla \int_{z'} \mathrm{pdf}_{z'}(z) \mathcal{L}(\theta + \epsilon z) dz) \\ &= \mathcal{R}(\int_{z'} \nabla \mathrm{pdf}_{z'}(z) \mathcal{L}(\theta + \epsilon z) dz) \\ &= \frac{1}{k} \mathcal{R}(\int_{z'} \nabla \mathrm{pdf}_{z'}(z) \mathcal{L}(\theta + \epsilon z) dz) \\ &= \frac{1}{k} \mathcal{R}(\int_{y'} \nabla \mathrm{e}^{-\frac{1}{2} || \frac{y - \theta}{\theta} ||^2} \mathcal{L}(y) \frac{1}{\epsilon^n} dy) \\ &= \frac{1}{k} \mathcal{R}(\int_{y'} \frac{y - \theta}{\epsilon^2} \mathrm{e}^{-\frac{1}{2\epsilon^2} || y - \theta ||^2} \mathcal{L}(y) \frac{1}{\epsilon^n} dy) \\ &= \frac{1}{k} \mathcal{R}(\int_{z'} \frac{z}{\epsilon} \mathrm{e}^{-\frac{1}{2} || x ||^2} \mathcal{L}(\theta + \epsilon z) dz) \\ &= \frac{1}{k} \mathcal{R}(\int_{z'} \frac{z}{\epsilon} \mathrm{e}^{-\frac{1}{2} || x ||^2} \mathcal{L}(\theta + \epsilon z) dz) \\ &= \mathcal{R}(\int_{z'} \mathrm{pdf}_{z'}(z) \mathcal{L}(\theta + \epsilon z) \frac{z}{\epsilon} dz) \\ &= \mathbb{E}_{\mathcal{R}}[\mathcal{R}(\frac{\mathcal{L}(\theta + \epsilon z')}{\epsilon} z')] \\ &= \mathbb{E}_{\mathcal{R}}[\frac{\mathcal{L}(\theta + \epsilon z')}{\epsilon} z'], \end{aligned}$$

where $y = \boldsymbol{\theta} + \epsilon \boldsymbol{z}, \boldsymbol{y}' = \boldsymbol{\theta} + \epsilon \boldsymbol{z}'$, and $k = \sqrt{(2\pi)^{\rho d}}$. Next,

$$\mathbb{E}_{\mathcal{R}}\left[\frac{\mathcal{L}(\boldsymbol{\theta}-\epsilon\boldsymbol{z}')}{\epsilon}\boldsymbol{z}'\right] = \frac{1}{k} \int_{-\boldsymbol{z}'} \frac{\mathcal{L}(\boldsymbol{\theta}+\epsilon(-\boldsymbol{z}))}{\epsilon} - \boldsymbol{z}e^{-\frac{1}{2}\|-\boldsymbol{z}\|^{2}}d(-\boldsymbol{z})$$
$$= \frac{1}{k} \int_{\hat{\boldsymbol{z}}} \frac{\mathcal{L}(\boldsymbol{\theta}+\epsilon\boldsymbol{z})}{\epsilon} \boldsymbol{z}e^{-\frac{1}{2}\|\boldsymbol{z}\|^{2}}d\boldsymbol{z}$$
$$= \mathbb{E}_{\mathcal{R}}\left[\frac{\mathcal{L}(\boldsymbol{\theta}+\epsilon\boldsymbol{z}')}{\epsilon}\boldsymbol{z}'\right].$$
(12)

 $= \frac{1}{2} \left(\mathbb{E}_{\mathcal{R}} \left[\frac{\mathcal{L}(\boldsymbol{\theta} + \epsilon \boldsymbol{z}^{'})}{\epsilon} \boldsymbol{z}^{'} \right] - \mathbb{E}_{\mathcal{R}} \left[\frac{\mathcal{L}(\boldsymbol{\theta} - \epsilon \boldsymbol{z}^{'})}{\epsilon} \boldsymbol{z}^{'} \right] \right)$ $= \mathbb{E}_{\mathcal{R}} \left[\frac{\mathcal{L}(\boldsymbol{\theta} + \epsilon \boldsymbol{z}^{'}) - \mathcal{L}(\boldsymbol{\theta} - \epsilon \boldsymbol{z}^{'})}{2\epsilon} \boldsymbol{z}^{'} \right]$

Therefore,

Q.E.D.

(13)

 $\widehat{\nabla}_{\boldsymbol{\theta}^{'}}\mathcal{L}_{\mathcal{R}}(\boldsymbol{\theta}) = \mathbb{E}_{\mathcal{R}}[\frac{\mathcal{L}(\boldsymbol{\theta} + \epsilon \boldsymbol{z}^{'})}{\epsilon} \boldsymbol{z}^{'}]$

 $= \mathbb{E}_{\mathcal{R}}[\widehat{\nabla}\mathcal{L}_{\mathcal{R}}(\boldsymbol{\theta})].$

Proof of Lemma 2. Firstly, compute the norm distance between the sparse ZO estimated gradient $\nabla_{\boldsymbol{\theta}'} \mathcal{L}_{\mathcal{R}}(\boldsymbol{\theta})$ and the sparse FO gradient $\nabla \mathcal{L}_{\mathcal{R}}(\boldsymbol{\theta})$:

`

$$\begin{split} \|\widehat{\nabla}_{\boldsymbol{\theta}'}\mathcal{L}_{\mathcal{R}}(\boldsymbol{\theta}) - \nabla\mathcal{L}_{\mathcal{R}}(\boldsymbol{\theta})\| &= \|\frac{1}{k} \int_{\boldsymbol{z}} (\frac{\mathcal{L}(\boldsymbol{\theta} + \epsilon \boldsymbol{z}) - \mathcal{L}(\boldsymbol{\theta} - \epsilon \boldsymbol{z})}{2\epsilon} - \langle \nabla\mathcal{L}_{\mathcal{R}}(\boldsymbol{\theta}), \boldsymbol{z} \rangle) \boldsymbol{z} e^{-\frac{1}{2} \|\boldsymbol{z}\|^{2}} d\boldsymbol{z}' \| \\ &= \|\frac{1}{k} \int_{\boldsymbol{z}} (\frac{\mathcal{L}(\boldsymbol{\theta} + \epsilon \boldsymbol{z}) - \mathcal{L}(\boldsymbol{\theta})}{\epsilon} - \langle \mathcal{R}(\nabla\mathcal{L}(\boldsymbol{\theta})), \boldsymbol{z} \rangle) \boldsymbol{z} e^{-\frac{1}{2} \|\boldsymbol{z}\|^{2}} d\boldsymbol{z}' \| \\ &\leq \frac{1}{k\epsilon} \int_{\boldsymbol{z}} |\mathcal{L}(\boldsymbol{\theta} + \epsilon \boldsymbol{z}) - \mathcal{L}(\boldsymbol{\theta}) - \epsilon \langle \nabla\mathcal{L}(\boldsymbol{\theta}), \epsilon \rangle \| \|\mathcal{R}(\boldsymbol{z})\| e^{-\frac{1}{2} \|\boldsymbol{z}\|^{2}} d\boldsymbol{z}' \\ &\leq \frac{\epsilon L(l)}{2k} \int_{\boldsymbol{\varepsilon}} \|\boldsymbol{z}\|^{2} \|\mathcal{R}(\boldsymbol{z})\| e^{-\frac{1}{2} \|\boldsymbol{z}\|^{2}} d\boldsymbol{z}' \\ &= \frac{\epsilon L(l)}{2} \mathbb{E}_{\mathcal{R}}[\|\boldsymbol{z}'\|^{3}] \\ &\leq \frac{\epsilon L(l)}{2} (\rho d + 3)^{\frac{3}{2}}. \end{split}$$
(14)

Subsequently, employing the Minkowski inequality for norms, which can be further generalized to $\|\boldsymbol{a} + \boldsymbol{b}\|^2 \leq 2\|\boldsymbol{a}\|^2 + 2\|\boldsymbol{b}\|^2$, by letting $\boldsymbol{a} = \nabla \mathcal{L}_{\mathcal{R}}(\boldsymbol{\theta}) - \widehat{\nabla}_{\boldsymbol{\theta}'}\mathcal{L}_{\mathcal{R}}(\boldsymbol{\theta})$ and $\boldsymbol{b} = \widehat{\nabla}_{\boldsymbol{\theta}'}\mathcal{L}_{\mathcal{R}}(\boldsymbol{\theta})$, we obtain:

$$\|\nabla \mathcal{L}_{\mathcal{R}}(\boldsymbol{\theta})\|^{2} \leq 2 \|\nabla \mathcal{L}_{\mathcal{R}}(\boldsymbol{\theta}) - \nabla_{\boldsymbol{\theta}'} \mathcal{L}_{\mathcal{R}}(\boldsymbol{\theta})\|^{2} + 2 \|\nabla_{\boldsymbol{\theta}'} \mathcal{L}_{\mathcal{R}}(\boldsymbol{\theta})\|^{2}$$

$$= 2 \|\widehat{\nabla}_{\boldsymbol{\theta}'} \mathcal{L}_{\mathcal{R}}(\boldsymbol{\theta}) - \nabla \mathcal{L}_{\mathcal{R}}(\boldsymbol{\theta})\|^{2} + 2 \|\widehat{\nabla}_{\boldsymbol{\theta}'} \mathcal{L}_{\mathcal{R}}(\boldsymbol{\theta})\|^{2}$$

$$= 2 \|\widehat{\nabla}_{\boldsymbol{\theta}'} \mathcal{L}_{\mathcal{R}}(\boldsymbol{\theta}) - \nabla \mathcal{L}_{\mathcal{R}}(\boldsymbol{\theta})\|^{2} + 2 \|\widehat{\nabla}_{\boldsymbol{\theta}'} \mathcal{L}_{\mathcal{R}}(\boldsymbol{\theta})\|^{2}$$

$$\leq \frac{\epsilon^{2} L^{2}(l)}{2} (\rho d + 3)^{3} + 2 \|\widehat{\nabla} \mathcal{L}_{\mathcal{R}}(\boldsymbol{\theta})\|^{2}$$

$$\leq \frac{\epsilon^{2} L^{2}(l)}{2} (\rho d + 4)^{3} + 2 \|\widehat{\nabla} \mathcal{L}_{\mathcal{R}}(\boldsymbol{\theta})\|^{2}.$$

$$(15)$$

Q.E.D.

Proof of Lemma 3.

$$\mathcal{L}_{\mathcal{R}}(\boldsymbol{\theta}) - \mathcal{L}(\boldsymbol{\theta}) = \mathbb{E}_{\mathcal{R}}[\mathcal{L}(\boldsymbol{\theta} + \epsilon \boldsymbol{z}') - \mathcal{L}(\boldsymbol{\theta})]$$

$$= \mathbb{E}_{\mathcal{R}}[\mathcal{L}(\boldsymbol{\theta} + \epsilon \boldsymbol{z}') - \mathcal{L}(\boldsymbol{\theta}) - \epsilon \langle \nabla \mathcal{L}(\boldsymbol{\theta}), \boldsymbol{z}' \rangle]$$

$$= \frac{1}{k} \int_{\boldsymbol{z}'} [\mathcal{L}(\boldsymbol{\theta} + \epsilon \boldsymbol{z}) - \mathcal{L}(\boldsymbol{\theta}) - \epsilon \langle \nabla \mathcal{L}(\boldsymbol{\theta}), \boldsymbol{z} \rangle] e^{-\frac{1}{2} \|\boldsymbol{z}\|^{2}} d\boldsymbol{z}$$

$$\leq \frac{1}{k} \int_{\boldsymbol{z}'} \frac{\epsilon^{2} L(l)}{2} \|\boldsymbol{z}\|^{2} e^{-\frac{1}{2} \|\boldsymbol{z}\|^{2}} d\boldsymbol{z}$$

$$= \frac{\epsilon^{2} L(l)}{2} \mathbb{E}_{\mathcal{R}}[\|\boldsymbol{z}'\|^{2}]$$

$$\leq \frac{\epsilon^{2} L(l)}{2} \rho d,$$
(16)

(17)

where $\theta_t = \theta + \epsilon z$. Given that \mathcal{L} satisfies Lipschitz continuity, we can derive $|\mathcal{L}(\theta') - \mathcal{L}(\theta) - \mathcal{L}(\theta)|$ $\langle \nabla \mathcal{L}(\boldsymbol{\theta}), \boldsymbol{\theta}_t - \boldsymbol{\theta} \rangle | \leq \frac{L(l)}{2} \| \boldsymbol{\theta}_t - \boldsymbol{\theta} \|^2$, establishing the validity of the first inequality.

$$[(\mathcal{L}_{\mathcal{R}}(\boldsymbol{\theta}) - \mathcal{L}(\boldsymbol{\theta})) - (\mathcal{L}_{\mathcal{R}}(\boldsymbol{\theta} + \epsilon \boldsymbol{z}') - \mathcal{L}(\boldsymbol{\theta} + \epsilon \boldsymbol{z}'))]^{2}$$

$$\leq 2[\mathcal{L}_{\mathcal{R}}(\boldsymbol{\theta}) - \mathcal{L}(\boldsymbol{\theta})]^{2} + 2[\mathcal{L}_{\mathcal{R}}(\boldsymbol{\theta} + \epsilon \boldsymbol{z}') - \mathcal{L}(\boldsymbol{\theta} + \epsilon \boldsymbol{z}')]^{2}$$

807
808
809

$$\leq 2[2\mathcal{R}(0) - \mathcal{L}(0)] + 2[\mathcal{L}(0 + \mathcal{L}) - \mathcal{L}$$

808
809
$$\leq \frac{\epsilon^{2}L^{2}(l)}{2}\rho^{2}d^{2} + \epsilon^{4}L^{2}(l)\rho^{2}d^{2}.$$

$$(\mathcal{L}_{\mathcal{R}}(\boldsymbol{\theta} + \epsilon \boldsymbol{z}') - \mathcal{L}_{\mathcal{R}}(\boldsymbol{\theta}))^{2} \leq 2(\mathcal{L}_{\mathcal{R}}(\boldsymbol{\theta} + \epsilon \boldsymbol{z}') - \mathcal{L}_{\mathcal{R}}(\boldsymbol{\theta}) - \epsilon \langle \widehat{\nabla} \mathcal{L}_{\mathcal{R}}(\boldsymbol{\theta}), \boldsymbol{z}' \rangle)^{2} + 2(\epsilon \langle \widehat{\nabla} \mathcal{L}_{\mathcal{R}}(\boldsymbol{\theta}), \boldsymbol{z}' \rangle)^{2} \leq \frac{\epsilon^{4} L^{2}(l)}{2} \|\boldsymbol{z}'\|^{4} + 2\epsilon^{2} \langle \widehat{\nabla} \mathcal{L}_{\mathcal{R}}(\boldsymbol{\theta}), \boldsymbol{z}' \rangle^{2} \leq \frac{\epsilon^{4} L^{2}(l)}{2} \|\boldsymbol{z}'\|^{4} + 2\epsilon^{2} \|\widehat{\nabla} \mathcal{L}_{\mathcal{R}}(\boldsymbol{\theta})\|^{2} \|\boldsymbol{z}'\|^{2}.$$
(18)

$$(\mathcal{L}(\boldsymbol{\theta} + \epsilon \boldsymbol{z}') - \mathcal{L}(\boldsymbol{\theta}))^{2} \leq 2((\mathcal{L}_{\mathcal{R}}(\boldsymbol{\theta}) - \mathcal{L}(\boldsymbol{\theta})) - (\mathcal{L}_{\mathcal{R}}(\boldsymbol{\theta} + \epsilon \boldsymbol{z}') - \mathcal{L}(\boldsymbol{\theta} + \epsilon \boldsymbol{z}')))^{2} + 2(\mathcal{L}_{\mathcal{R}}(\boldsymbol{\theta} + \epsilon \boldsymbol{z}') - \mathcal{L}_{\mathcal{R}}(\boldsymbol{\theta}))^{2} \leq 2\epsilon^{4}L^{2}(l)\rho^{2}d^{2} + \epsilon^{4}L^{2}(l)\|\boldsymbol{z}'\|^{4} + 4\epsilon^{2}\|\widehat{\nabla}\mathcal{L}_{\mathcal{R}}(\boldsymbol{\theta})\|^{2}\|\boldsymbol{z}'\|^{2}$$

$$(19)$$

$$\mathbb{E}_{\boldsymbol{z},\boldsymbol{x}}[\|\widehat{\nabla}\mathcal{L}_{\mathcal{R}}(\boldsymbol{\theta})\|^{2}] = \mathbb{E}_{\mathcal{R}}[\|\frac{\mathcal{L}(\boldsymbol{\theta} + \epsilon \boldsymbol{z}') - \mathcal{L}(\boldsymbol{\theta} - \epsilon \boldsymbol{z}')}{2\epsilon} \boldsymbol{z}'\|^{2}]$$

$$= \mathbb{E}_{\mathcal{R}}[\|\frac{\mathcal{L}(\boldsymbol{\theta} + \epsilon \boldsymbol{z}') - \mathcal{L}(\boldsymbol{\theta})}{2\epsilon} \boldsymbol{z}' + \frac{\mathcal{L}(\boldsymbol{\theta}) - \mathcal{L}(\boldsymbol{\theta} - \epsilon \boldsymbol{z}')}{2\epsilon} \boldsymbol{z}'\|^{2}]$$

$$\leq \mathbb{E}_{\mathcal{R}}[2\|\frac{\mathcal{L}(\boldsymbol{\theta} + \epsilon \boldsymbol{z}') - \mathcal{L}(\boldsymbol{\theta})}{2\epsilon} \boldsymbol{z}'\|^{2} + 2\|\frac{\mathcal{L}(\boldsymbol{\theta}) - \mathcal{L}(\boldsymbol{\theta} - \epsilon \boldsymbol{z}')}{2\epsilon} \boldsymbol{z}'\|^{2}]$$

$$= \mathbb{E}_{\mathcal{R}}[\frac{1}{2\epsilon^{2}}[\mathcal{L}(\boldsymbol{\theta} + \epsilon \boldsymbol{z}') - \mathcal{L}(\boldsymbol{\theta})]^{2} \cdot \|\boldsymbol{z}'\|^{2} + \frac{1}{2\epsilon^{2}}[\mathcal{L}(\boldsymbol{\theta}) - \mathcal{L}(\boldsymbol{\theta} - \epsilon \boldsymbol{z}')]^{2} \cdot \|\boldsymbol{z}'\|^{2}]$$

$$\leq \mathbb{E}_{\mathcal{R}}[2\epsilon^{2}L^{2}(l)\rho^{2}d^{2}\|\boldsymbol{z}'\|^{2} + \epsilon^{2}L^{2}(l)\|\boldsymbol{z}'\|^{6} + 4\|\widehat{\nabla}\mathcal{L}_{\mathcal{R}}(\boldsymbol{\theta})\|^{2}\|\boldsymbol{z}'\|^{4}]$$

$$\leq 2\epsilon^{2}L^{2}(l)\rho^{3}d^{3} + \epsilon^{2}L^{2}(l)(\rho d + 6)^{3} + 4(\rho d + 4)^{2}\|\widehat{\nabla}\mathcal{L}_{\mathcal{R}}(\boldsymbol{\theta})\|^{2}$$

$$\leq 3\epsilon^{2}L^{2}(l)(\rho d + 4)^{3} + 4(\rho d + 4)^{2}\|\widehat{\nabla}\mathcal{L}_{\mathcal{R}}(\boldsymbol{\theta})\|^{2}.$$
(20)

As $\mathbb{E}_{\mathcal{R}}[\|\boldsymbol{z}\|^p] \leq (\rho d + p)^{\frac{r}{2}}$ for $p \geq 2$, the third inequality holds. Additionally, since $2\rho^3 d^3 + (\rho d + 6)^3 \leq 3(\rho d + 4)^3$, the fourth inequality is valid. Again, given that \mathcal{L} is Lipschitz continuous, we have $|\mathcal{L}(\boldsymbol{\theta}_{t+1}) - \mathcal{L}(\boldsymbol{\theta}_t) - \langle \nabla \mathcal{L}(\boldsymbol{\theta}_t), \boldsymbol{\theta}_{t+1} - \boldsymbol{\theta}_t \rangle| \leq 1$

Again, given that \mathcal{L} is Lipschitz continuous, we have $|\mathcal{L}(\boldsymbol{\theta}_{t+1}) - \mathcal{L}(\boldsymbol{\theta}_t) - \langle \nabla \mathcal{L}(\boldsymbol{\theta}_t), \boldsymbol{\theta}_{t+1} - \boldsymbol{\theta}_t \rangle| \le \frac{L(l)}{2} \|\boldsymbol{\theta}_{t+1} - \boldsymbol{\theta}_t\|^2$. Therefore, we can derive:

$$\mathcal{L}_{\mathcal{R}}(\boldsymbol{\theta}_{t+1}) - \mathcal{L}_{\mathcal{R}}(\boldsymbol{\theta}_{t}) - \langle \widehat{\nabla} \mathcal{L}_{\mathcal{R}}(\boldsymbol{\theta}_{t}), \boldsymbol{\theta}_{t+1} - \boldsymbol{\theta}_{t} \rangle \\ \leq |\mathcal{L}_{\mathcal{R}}(\boldsymbol{\theta}_{t+1}) - \mathcal{L}_{\mathcal{R}}(\boldsymbol{\theta}_{t}) - \langle \widehat{\nabla} \mathcal{L}_{\mathcal{R}}(\boldsymbol{\theta}_{t}), \boldsymbol{\theta}_{t+1} - \boldsymbol{\theta}_{t} \rangle| \\ \leq \frac{L(l)}{2} \|\boldsymbol{\theta}_{t+1} - \boldsymbol{\theta}_{t}\|^{2}.$$
(21)

By further examining the iterative process of Structured Sparse ZO-SGD as outlined in Equation (5), we can obtain:

$$\mathcal{L}_{\mathcal{R}}(\boldsymbol{\theta}_{t+1}) \leq \mathcal{L}_{\mathcal{R}}(\boldsymbol{\theta}_{t}) + \langle \widehat{\nabla} \mathcal{L}_{\mathcal{R}}(\boldsymbol{\theta}_{t}), \boldsymbol{\theta}_{t+1} - \boldsymbol{\theta}_{t} \rangle + \frac{L(l)}{2} \|\boldsymbol{\theta}_{t} - \boldsymbol{\theta}_{t+1}\|^{2} = \mathcal{L}_{\mathcal{R}}(\boldsymbol{\theta}_{t}) - \eta_{t} \langle \widehat{\nabla} \mathcal{L}_{\mathcal{R}}(\boldsymbol{\theta}_{t}), \widehat{\nabla} \mathcal{L}_{\mathcal{R}}(\boldsymbol{\theta}_{t}) \rangle + \frac{(\eta_{t})^{2} L(l)}{2} \|\widehat{\nabla} \mathcal{L}_{\mathcal{R}}(\boldsymbol{\theta}_{t})\|^{2},$$
(22)

where η_t represents the learning rate at step t.

Subsequently, we can derive the expected loss function of the structured sparse model at step t + 1 as:

$$\mathbb{E}_{\boldsymbol{z}',\boldsymbol{x}}[\mathcal{L}_{\mathcal{R}}(\boldsymbol{\theta}_{t+1})] \leq \mathbb{E}_{\boldsymbol{z}',\boldsymbol{x}}[\mathcal{L}_{\mathcal{R}}(\boldsymbol{\theta}_{t})] - \eta_{t}\mathbb{E}_{\boldsymbol{z}',\boldsymbol{x}}[\|\widehat{\nabla}\mathcal{L}_{\mathcal{R}}(\boldsymbol{\theta}_{t})\|^{2}] + \frac{(\eta_{t})^{2}L(l_{\boldsymbol{z}})}{2}\mathbb{E}_{\boldsymbol{z}',\boldsymbol{x}}[\|\widehat{\nabla}\mathcal{L}(\boldsymbol{\theta}_{t})\|^{2}] \leq \mathbb{E}_{\boldsymbol{z}',\boldsymbol{x}}[\mathcal{L}_{\mathcal{R}}(\boldsymbol{\theta}_{t})] - \eta_{t}\mathbb{E}_{\boldsymbol{z}',\boldsymbol{x}}[\|\widehat{\nabla}\mathcal{L}_{\mathcal{R}}(\boldsymbol{\theta}_{t})\|^{2}]$$
(23)

+
$$\frac{(\eta_t)^2 L(l)}{2} (4(\rho d + 4) \mathbb{E}_{\mathbf{z}', x}[\|\widehat{\nabla} \mathcal{L}_{\mathcal{R}}(\boldsymbol{\theta}_t)\|^2] + 3\epsilon^2 L^2(l)(\rho d + 4)^3).$$

Then, let learning rate be $\eta_t = \frac{1}{4(\rho d + 4)L(l)}$ and obtain:

$$\mathbb{E}_{\boldsymbol{z}',\boldsymbol{x}}[\mathcal{L}_{\mathcal{R}}(\boldsymbol{\theta}_{t+1})] \leq \mathbb{E}_{\boldsymbol{z}',\boldsymbol{x}}[\mathcal{L}_{\mathcal{R}}(\boldsymbol{\theta}_{t})] - \frac{1}{8(\rho d+4)L(l)} \mathbb{E}_{\boldsymbol{z}',\boldsymbol{x}}[\|\widehat{\nabla}\mathcal{L}_{\mathcal{R}}(\boldsymbol{\theta}_{t})\|^{2}] + \frac{3\epsilon^{2}}{32}L(l)(\rho d+4).$$
(24)

Summing Equation (24) from 0 to T + 1, where T denotes a sufficiently large number of training steps, yields:

$$\mathbb{E}_{\boldsymbol{z}',\boldsymbol{x}}[\|\widehat{\nabla}\mathcal{L}_{\mathcal{R}}(\boldsymbol{\theta}_{T})\|^{2}] \leq 8(\rho d+4)L[\frac{\mathcal{L}_{\mathcal{R}}(\boldsymbol{\theta}_{0})-\mathcal{L}_{\mathcal{R}}^{*}}{T+1}+\frac{3\epsilon^{2}}{32}L(\rho d+4)],$$
(25)

where $L(l) \leq L$ for all $\mathcal{L}(\boldsymbol{\theta}_t)$. Thus, based on Lemma 2, we can have:

$$\mathbb{E}_{\boldsymbol{z}',\boldsymbol{x}}[\|\nabla\mathcal{L}_{m}(\boldsymbol{\theta}_{T})\|^{2}] \leq \frac{\epsilon^{2}L^{2}}{2}(\rho d+4)^{3} + 2\mathbb{E}_{\boldsymbol{z}',\boldsymbol{x}}[\|\widehat{\nabla}\mathcal{L}_{\mathcal{R}}(\boldsymbol{\theta}_{T})\|^{2}]$$

$$\leq 16(\rho d+4)L\frac{\mathcal{L}_{\mathcal{R}}(\boldsymbol{\theta}_{0}) - \mathcal{L}_{\mathcal{R}}^{*}}{T+1} + \frac{\epsilon^{2}L^{2}}{2}(\rho d+4)^{2}(\rho d+\frac{11}{2}).$$
(26)

To obtain σ -accurate solution $\mathbb{E}_{\boldsymbol{z}',x}[\|\nabla \mathcal{L}_m(\boldsymbol{\theta}_T)\|^2] \leq \sigma^2$, we can define $\epsilon = \Omega(\frac{\sigma}{\rho^{\frac{3}{2}}d^{\frac{3}{2}}L})$.

$$16(\rho d+4)L\frac{\mathcal{L}_{\mathcal{R}}(\boldsymbol{\theta}_{0}) - \mathcal{L}_{\mathcal{R}}^{*}}{T+1} + \mathcal{O}(\epsilon^{2}L^{2}\rho^{3}d^{3})$$

=
$$16(\rho d+4)L\frac{\mathcal{L}_{\mathcal{R}}(\boldsymbol{\theta}_{0} - \mathcal{L}_{\mathcal{R}}^{*})}{T+1} + \mathcal{O}(\sigma^{2}).$$
 (27)

From the above, we can get:

$$T = \mathcal{O}(\frac{\rho dL}{\sigma^2}). \tag{28}$$

Q.E.D.