

# MEMORY EFFICIENT BLOCK COORDINATE DESCENT METHOD FOR HESSIAN-INFORMED ZERO-ORDER OPTIMIZER

**Anonymous authors**

Paper under double-blind review

## ABSTRACT

Fine-tuning large language models (LLMs) for specific downstream tasks has traditionally relied on memory-intensive optimizers using classical backpropagation, which demands substantial memory to store model states for gradient computation, motivating the development of memory-efficient zeroth-order optimizers that operate in a forward-only manner. However, the slower convergence of the zeroth-order optimizer remains a challenge, which recent research addresses by incorporating Hessian information to accelerate training, although storing even the diagonal Hessian requires memory equivalent to that of the model weights, leading to significant memory usage. To mitigate this problem, we propose a novel approach that integrates the block coordinate descent (BCD) method with a Hessian-informed zeroth-order optimizer, allowing us to treat model layers as separate blocks and update only a subset of layers per training iteration, thereby reducing memory requirements and accelerating convergence. Specifically, at each iteration, an active block of layers is selected according to the chosen BCD rule, such as ascending order, and their weights are updated while the other layers remain fixed, with diagonal Hessian information stored and updated exclusively for the active layers. For fine-tuning foundation models of medium size (OPT-1.3B and LLaMA-2-7B), our method achieves up to 39% memory reduction compared to existing Hessian-informed zeroth-order methods, while preserving baseline accuracy and memory usage to zeroth-order methods across various tasks, offering a memory-efficient alternative method for LLMs fine-tuning, especially on memory-constrained devices.

## 1 INTRODUCTION

Fine-tuning transformer-based large models is an essential step in adapting pre-trained models to specific downstream tasks and further improving performance (Raffel et al., 2020). This process also allows the model to continue training on lower-end devices compared to those used for pre-training, thus improving accessibility and reducing the training cost. For this intent, parameter-efficient fine-tuning (PEFT) techniques, such as LoRA (Hu et al., 2021), have been proposed to enable fine-tuning on consumer-level GPUs or even edge devices, providing significant economic and practical benefits. Typically, fine-tuning employs traditional optimizers like SGD or Adam (Kingma, 2014), which use backpropagation to update model weights. This process requires storing parameters, gradients, activations, and possibly other optimizer states, significantly increasing memory requirements (Lv et al., 2023b;a; Rajbhandari et al., 2020). As model sizes have increased and larger batch sizes are employed for training, the memory demands of traditional optimizers have become a significant bottleneck for devices with limited memory resources, even when using existing PEFT methods (Cai et al., 2020). Our work aims to address this challenge by exploring memory-efficient techniques further to reduce the memory overhead during fine-tuning on low-end devices.

To tackle the memory inefficiency issue, recent advancements have explored the use of zeroth-order optimizers such as MeZO (Malladi et al., 2023) that estimate the gradients with only forward passes, which eliminates the need for backpropagation, thereby significantly reducing memory consumption by avoiding the storage of intermediate optimizer states. Though memory-efficient, the slower convergence rates of zeroth-order optimizers have limited their practical utility. To accelerate

054  
055  
056  
057  
058  
059  
060  
061  
062  
063  
064  
065  
066  
067  
068  
069  
070  
071  
072  
073  
074  
075  
076  
077  
078  
079  
080  
081  
082  
083  
084  
085  
086  
087  
088  
089  
090  
091  
092  
093  
094  
095  
096  
097  
098  
099  
100  
101  
102  
103  
104  
105  
106  
107

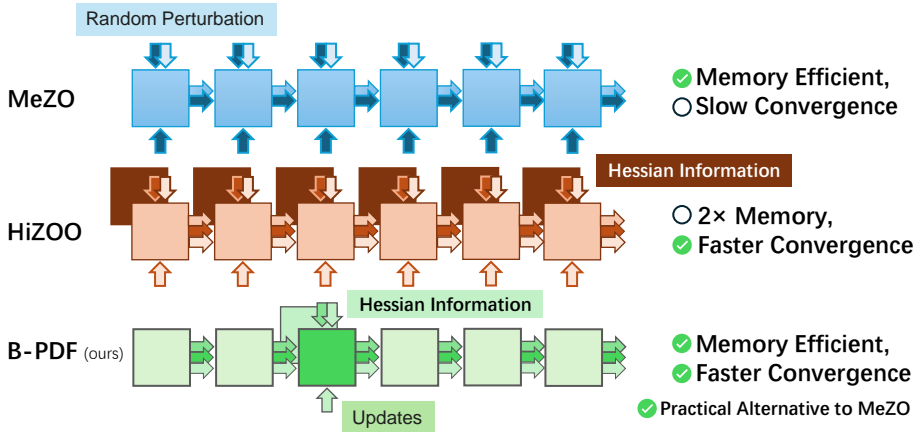


Figure 1: Illustration of MeZO, HiZOO and our proposed training pipeline.

convergence, researchers have incorporated second-order information, such as diagonal Hessian approximations as proposed in HiZOO (Zhao et al., 2024b), into the optimization process. However, this solution comes at the cost of memory overhead, as storing even diagonal Hessian values introduces substantial memory cost comparable to the storage required for the model weights themselves, thereby negating the original memory-saving intent of applying zeroth-order optimization.

To efficiently and effectively utilize second-order information, we consider the traditional block coordinate descent (BCD) method, which solves optimization problems successively along coordinate directions, and propose a block coordinate descent Newton method (BCD-Newton) to tackle our challenge. Inspired by recent advances in applying BCD with Adam (Kingma, 2014) and AdamW Loshchilov (2017) optimizations for training large language models (LLMs) (Pan et al., 2024; Luo et al., 2024), we introduce a layer-wise block coordinate descent scheme to optimize memory usage, treating model layers as independent blocks and selectively activates a subset of layers during each iteration. In practice, block selection is guided by various BCD rules, including ordered selection and block-wise importance sampling, to achieve optimal training performance. This approach significantly reduces the memory required to store Hessian information. In our optimization step, beyond the basic zeroth-order method with two forward passes, a three-step forward pass is employed to incorporate second-order updates, thereby facilitating faster convergence. The second-order term is stored as a diagonal Hessian estimate matrix, sized according to the active block of selected transformer layers, and is updated throughout the training process. Thus, we propose a novel optimizer that addresses the memory-convergence trade-off inherent in Hessian-informed zeroth-order optimization by integrating BCD techniques, while simultaneously improving convergence rates.

Through extensive experiments on a single RTX 4090 or RTX A6000 GPU, we demonstrate that our method enhances training efficiency and memory management while fine-tuning foundation models, including OPT-1.3B (Zhang et al., 2022b) and LLaMA-2-7B (Touvron et al., 2023). As a BCD zeroth-order Newton method, it empirically delivers superior convergence speed and accuracy compared to MeZO. Additionally, compared to the HiZOO baseline, our approach achieves approximately a 50% speedup and a 40% reduction in memory usage with comparable baseline accuracy across multiple GLUE (Wang, 2018) and SuperGLUE (Wang et al., 2019) tasks. These improvements make our method particularly well-suited for fine-tuning large models on devices with limited memory, expanding the accessibility of large language models in real-world applications.

In summary, our main contributions are three-fold:

- We propose a novel block coordinate descent fine-tuning pipeline that integrates the previous Hessian-informed zeroth-order optimizer, reducing the memory overhead to make the method a practical and convergence-enhanced alternative to MeZO.
- We design improved block coordinate descent schemes that reduce the compute and memory cost of the Hessian-informed forward-only optimizer. By adaptively updating weights across block coordinates of the model layers, this method manages blockwise updates efficiently and reduces memory and computational costs.

- We conduct experiments on fine-tuning OPT-1.3B and LLaMA2-7B, demonstrating that our method reduces training memory by over 39% without a loss in accuracy compared to the full diagonal Hessian baseline.

## 2 RELATED WORK

**First and second order Optimization for LLMs.** Traditional first-order optimizers, such as SGD, AdaGrad (Duchi et al., 2011), and RMSProp (Tieleman et al., 2012), are foundational tools in deep learning. Adam (Kingma, 2014), with its adaptive moment estimates for faster convergence, and its variant AdamW (Loshchilov, 2017), which modifies the weight decay term to improve generalization, have become the dominant optimizers for fine-tuning large language models (LLMs). Second-order optimization methods incorporating Hessian information, such as K-FAC (Martens & Grosse, 2015), EVA (Zhang et al., 2022a), Adahessian (Yao et al., 2021), and Sophia (Liu et al., 2023), have been explored to further accelerate convergence. However, estimating the Hessian is computationally and memory intensive, particularly with the growing size of LLMs, which makes these second-order methods less practical for fine-tuning on devices with limited memory resources.

**Zeroth-order (ZO) Optimization.** A classical zeroth-order optimization method is SPSA (Spall, 1992), along with its corresponding SGD variant, ZO-SGD, which estimates the gradient using two forward passes before and after parameter perturbation. Recently, MeZO (Malladi et al., 2023) adapted ZO-SGD by incorporating the random number generator, enabling an in-place implementation that significantly reduces memory usage for storing random vectors during training. Based on MeZO, recent work explores its variants like incorporated sparsity for memory efficiency (Guo et al., 2024). Additionally, Zhang et al. (2024) conducted a benchmark study to analyze and enhance zeroth-order fine-tuning methods. However, the convergence performance of ZO methods often falls behind that of first-order methods. To improve convergence, HiZOO (Zhao et al., 2024b) proposed to utilize Hessian information through diagonal Hessian estimation. Beyond these approaches, several other gradient-free methods have been proposed, such as using evolutionary algorithms for gradient-free optimization (Sun et al., 2022b;a).

**Memory-efficient Fine-tuning for LLMs.** Numerous algorithms have been developed to reduce memory costs for training LLMs. Based on backpropagation, practical techniques such as gradient checkpointing (Chen et al., 2016) recompute gradients, FlashAttention (Dao et al., 2022) employs tiling and recomputation to leverage cache for improved efficiency, and the ZeRO optimizers (Rajbhandari et al., 2020; Ren et al., 2021) enable offloading to manage memory usage effectively. Additionally, researchers have utilized compression and quantization methods to approximate gradients, activations, and other optimizer states, enhancing training performance (Jiang et al., 2022; Li et al., 2024). On another front, methods like LOMO (Lv et al., 2023b;a) fuse gradient updates to accelerate training. One notable approach to fine-tuning is parameter-efficient fine-tuning (PEFT) methods, which includes techniques such as Adapters (LoRA) (Hu et al., 2021; Housby et al., 2019), prompt tuning (Lester et al., 2021), and selective methods like bias-only fine-tuning (Zaken et al., 2021) and layer-wise freezing (Brock et al., 2017). In addition, Zhao et al. (2024a) recently introduced GaLore which reduces memory costs by projecting gradients into a low-rank compact space.

**Block Coordinate Descent (BCD) methods for LLM Optimization.** In BCD, the optimization objective is minimized successively along coordinate directions. When applied to LLM fine-tuning, this approach can be seen as a branch of selective methods in parameter-efficient fine-tuning. The recently proposed BAdam (Luo et al., 2024) showcases the effectiveness of combining block coordinate descent with Adam. Similarly, LiSA (Pan et al., 2024) improves performance by selectively updating transformer layers with AdamW optimizer, outperforming LoRA across tasks on LLaMA-2-70B.

Overall, our method offers a complementary optimizer-based solution that can be combined with techniques like compression and system-level approaches to improve memory efficiency. Amid the rapid advancements in efficient training for LLMs and other foundation models, the most closely related works to ours are HiZOO and BAdam. However, our approach distinguishes itself by addressing the memory overhead of these methods in two key ways: first, by eliminating the need for backpropagation through zeroth-order optimization, and second, by reducing the memory cost of Hessian-informed methods through block coordinate descent. Furthermore, unlike PEFT methods, our approach enables full parameter fine-tuning, which has been demonstrated to yield superior performance in various tasks (Ding et al., 2022).

### 3 REVISITING MEMORY CONSUMPTION FROM BCD PERSPECTIVE

In this section, we provide a brief overview of how zeroth-order (ZO) and Hessian-informed ZO optimizer methods work by introducing the core concepts of MeZO (Malladi et al., 2023) and HiZOO (Zhao et al., 2024b). Next, we introduce block coordinate descent (BCD) methods such as BAdam (Luo et al., 2024). To ensure consistency, we have adapted the definitions from these works. Finally, we reconsider the memory consumption of these methods, which leads us to propose our BCD-integrated Newton method optimizer.

#### 3.1 PRELIMINARIES OF ZEROTH-ORDER OPTIMIZERS

##### 3.1.1 SPSA, ZO-SGD, AND MEZO

Let  $\mathcal{L}(\theta; \mathcal{B})$  represent the loss function for training the model with parameters  $\theta \in \mathbb{R}^d$  on the minibatch  $\mathcal{B}$ , omitting the  $\mathcal{B}$  for simplicity. The SPSA algorithm (Spall, 1992) perturbs the model using  $z \in \mathbb{R}^d$ , sampled from  $\mathcal{N}(0, \mathbf{I}_d)$ , and estimates the gradient on the minibatch as follows:

$$\hat{\nabla} \mathcal{L}(\theta) = \frac{\mathcal{L}(\theta + \mu z) - \mathcal{L}(\theta - \mu z)}{2\mu} z \approx z z^\top \nabla \mathcal{L}(\theta) \quad (1)$$

where  $\mu$  is the perturbation scale.

The corresponding SPSA optimizer, ZO-SGD, employs two forward passes to estimate the gradients. With learning rate  $\eta$ , ZO-SGD updates the parameters as  $\theta_{t+1} = \theta_t - \eta \hat{\nabla} \mathcal{L}(\theta; \mathcal{B}_t)$ . In this vanilla algorithm, the sampled vector  $z$  requires memory equivalent to that of the perturbed weights, resulting in a memory cost that is twice that of inference.

In contrast, MeZO (Malladi et al., 2023) introduces an in-place implementation using a random number generator. Only a random seed  $s$  needs to be sampled and stored at each step, allowing the generator to be reset by  $s$  to regenerate the vector  $z$ . This approach eliminates the need to save the vector, reducing the memory cost to match that of inference.

##### 3.1.2 HiZOO

To harness second-order information through MeZO for enhanced convergence rates, Zhao et al. (2024b) introduce HiZOO, utilizing a diagonal Hessian-based preconditioner that adjusts the update sizes of parameters based on their curvature. By estimating and storing only the diagonal Hessian, HiZOO requires  $\mathcal{O}(d)$  memory, significantly less than the  $\mathcal{O}(d^2)$  needed for the full Hessian matrix.

Let  $\Sigma$  denote the estimated inverse Hessian matrix, approximating the diagonal Hessian as a positive definite matrix, with  $\Sigma^{-1} \approx \nabla^2 \mathcal{L}(\theta)$ . Define  $\Sigma_t$  as the matrix at training step  $t$ , initialized as  $\Sigma_0 = \mathbf{I}_d$ . Storing  $\Sigma_t$  incurs a memory cost of  $\mathcal{O}(d)$ , and it is updated at each step. In addition, to mitigate noise in the computation, an exponential moving average (EMA) is employed, leading to the following update rule for the diagonal Hessian estimate:

$$\Sigma_{t+1}^{-1} = (1 - \alpha_t) \Sigma_t^{-1} + \alpha_t \left| \text{diag}(\hat{\Sigma}_t) \right|, \quad (2)$$

where  $\alpha_t$  is a smooth scale, and  $|\text{diag}(\hat{\Sigma}_t)|$  ensures that all entries of  $\Sigma_t$  remain non-negative.

HiZOO approximates the diagonal Hessian using three forward passes to compute  $\mathcal{L}(\theta + \mu \Sigma^{1/2} z)$ ,  $\mathcal{L}(\theta - \mu \Sigma^{1/2} z)$ , and  $\mathcal{L}(\theta)$ . By applying Taylor’s expansion, we obtain:

$$\mathcal{L}(\theta \pm \mu \Sigma^{1/2} z) = \mathcal{L}(\theta) \pm \mu \langle \mathcal{L}(\theta), \Sigma^{1/2} z \rangle + \frac{\mu^2}{2} z^\top \Sigma^{1/2} \nabla^2 \mathcal{L}(\theta) \Sigma^{1/2} z + \mathcal{O}(\mu^3), \quad (3)$$

the difference  $\Delta \mathcal{L}$  is then calculated as:

$$\begin{aligned} \Delta \mathcal{L} &= \mathcal{L}(\theta + \mu \Sigma^{1/2} z) + \mathcal{L}(\theta - \mu \Sigma^{1/2} z) - 2\mathcal{L}(\theta) \\ &= \mu^2 z^\top \Sigma^{1/2} \nabla^2 \mathcal{L}(\theta) \Sigma^{1/2} z + \mathcal{O}(\mu^3). \end{aligned}$$

Based on the lemma from Ye (2023):

$$\frac{1}{2} \cdot \mathbb{E}_{\mathbf{z}}(\mathbf{z}^\top \Sigma^{1/2} \nabla^2 \mathcal{L}(\boldsymbol{\theta}) \Sigma^{1/2} \mathbf{z} \cdot (\Sigma^{-1/2} \mathbf{z} \mathbf{z}^\top \Sigma^{-1/2} - \Sigma^{-1})) = \nabla^2 \mathcal{L}(\boldsymbol{\theta}), \quad (4)$$

substitute  $\Delta \mathcal{L}$  into the left side, and we obtain:

$$\frac{1}{2} \mathbb{E} \left[ \frac{\Delta \mathcal{L}}{\mu^2} \cdot \left( \Sigma^{-1/2} \mathbf{z} \mathbf{z}^\top \Sigma^{-1/2} - \Sigma^{-1} \right) \right] = \nabla^2 \mathcal{L}(\boldsymbol{\theta}) + \mathcal{O}(\mu).$$

Consequently, the estimation of the diagonal Hessian  $\nabla^2 \mathcal{L}(\boldsymbol{\theta})$  at  $\boldsymbol{\theta}$  is:

$$\hat{\Sigma}_t(\boldsymbol{\theta}) = \frac{\Delta \mathcal{L}}{2\mu^2} \left( \Sigma_t^{-1/2} \mathbf{z}_i \mathbf{z}_i^\top \Sigma_t^{-1/2} - \Sigma_t^{-1} \right). \quad (5)$$

In this manner, HiZOO approximates the diagonal entries of  $\nabla^2 \mathcal{L}(\boldsymbol{\theta})$  by  $\text{diag}(\Sigma'_t(\boldsymbol{\theta}))$ , requiring one more forward pass per step compared with MeZO.

### 3.1.3 BLOCK COORDINATE DESCENT

At each iteration, block coordinate descent (BCD) fixes all other parameters and optimizes the objective function over the selected coordinates, resulting in an optimization problem with reduced dimension. For large language models, a natural block partition is to organize transformer layers in ascending order. Formally, an ordered block partition  $\pi = \{\pi_1, \dots, \pi_i, \dots, \pi_D\}$  divides the entire model parameters  $\boldsymbol{\theta} \in \mathbb{R}^d$  into  $D$  blocks, such that  $\boldsymbol{\theta} = \{\boldsymbol{\theta}_{\pi_1}, \dots, \boldsymbol{\theta}_{\pi_i}, \dots, \boldsymbol{\theta}_{\pi_D}\}$  with  $\boldsymbol{\theta}_{\pi_i} \in \mathbb{R}^{d_i}$  and  $\sum_{i=1}^D d_i = d$ . Based on the main idea of BCD, BAdam (Luo et al., 2024) propose to incorporate Adam updates as its inner solver and optimize over only one active block  $\boldsymbol{\theta}_{\pi_i}$  at a time while keeping the other inactive blocks fixed. Mathematically, BAdam solves the following subproblem at the  $t$ -th block-epoch for  $i = 1, \dots, D$  to update the active block  $\boldsymbol{\theta}_{\pi_i}$ :

$$\boldsymbol{\theta}_{\pi_i}^{t+1} \in \arg \min_{\boldsymbol{\theta}_{\pi_i} \in \mathbb{R}^{d_i}} \mathcal{L}(\boldsymbol{\theta}_{\pi_1}^{t+1}, \dots, \boldsymbol{\theta}_{\pi_{i-1}}^{t+1}, \boldsymbol{\theta}_{\pi_i}, \boldsymbol{\theta}_{\pi_{i+1}}^t, \dots, \boldsymbol{\theta}_{\pi_D}^t). \quad (6)$$

This subproblem Equation 6 keeps inactive blocks fixed at their latest values, leading to a significantly lower-dimensional optimization problem compared to  $\min_{\boldsymbol{\theta}} \mathcal{L}(\boldsymbol{\theta})$ .

## 3.2 REVISITING MEMORY CONSUMPTION FROM BCD PERSPECTIVE

**Who consumed my memory?** Second-order methods incorporate full or diagonal Hessian matrix, or its estimation, as a preconditioner to accelerate convergence, but this introduces a significant memory cost of  $\mathcal{O}(d)$ . For large models such as LLaMA-2-7B (Touvron et al., 2023) with  $d = 7$  billion parameters, this requires  $2d$  memory in FP16 precision, resulting in approximately over 14GB of memory storage. When combined with the memory required for model parameters, this easily exceeds the capacity of consumer-level devices, undermining MeZO’s original goal of achieving memory efficiency. Our experiments further demonstrate that directly applying Hessian-based optimization steps significantly increases memory usage, as shown in Table 1. Even though approaches such as HiZOO offer performance improvements, the considerable memory overhead from storing Hessian information becomes a bottleneck, particularly when fine-tuning large models. This dilemma leads to a situation where the benefits of second-order methods are outweighed by their heavy memory consumption, limiting their practicality in memory-constrained environments. Furthermore, the memory consumption increases with batch size for both first-order and Hessian-based methods, intensifying the memory overhead, as illustrated in Figure 2.

**How to reduce Hessian memory consumption?** To address this memory-convergence trade-off, we propose integrating block coordinate descent (BCD) into the zeroth-order Newton optimization. BCD allows us to partition the model into blocks, optimizing only a subset of layers at each iteration while keeping the rest fixed. This approach dramatically reduces the memory required for storing Hessian information, as it is only computed for the active blocks. For instance, by partitioning the aforementioned LLaMA-2-7B model into  $D = 32$  blocks, corresponding to its 32 transformer layers, we reduce the additional memory cost associated with Hessian storage to  $\frac{2d}{D}$ , bringing it to under 1GB of memory. This significantly improves memory efficiency while preserving the advantages of second-order optimization.



Moreover, we further optimize memory usage by applying MeZO to update the embedding and language modeling head layers, avoiding the instability and overhead often associated with second-order methods. Our integration of BCD not only achieves comparable memory usage to MeZO but also leverages the improved convergence rates of Hessian-informed updates.

Table 1: Experiments of actual GPU memory consumption for various algorithms.

DEVICE	MODEL	SGD	BCD	LoRA	MEZO	HiZOO	OURS B-PDF
RTX 4090	OPT-1.3B	23G	21G	11G	4.4G	7.5G	4.6G
RTX A6000	LLAMA-2-7B	> 48G	46G	40G	31G	> 48G	32G
THEORETICAL AVG MEMORY IN UNITS PARAM+ACTIVATION+GRAD+HESSIAN		SGD $3d$	BCD OR LoRA $d \sim 3d$	MEZO $d$	HiZOO $2d$	OURS B-PDF $d + d/D$	

To validate our analysis, we conducted preliminary experiments (detailed in Section 5) measuring the GPU memory consumption of various optimizers during the fine-tuning of medium-sized language models, specifically OPT-1.3B (Zhang et al., 2022b) on an RTX 4090 (24GB) and LLaMA-2-7B on an RTX A6000 (48GB). As Table 1 and Figure 2 briefly illustrate, HiZOO’s incorporation of second-order information increases memory demand by over 70%. Notably, the actual allocated memory includes residual state memory such as temporary buffers and fragments (Rajbhandari et al., 2020), which means the overall memory requirement exceeds that of the parameters alone, resulting in the overall increase short of a full 100%. In contrast, our BCD-integrated method significantly reduces memory consumption, bringing it in line with MeZO while maintaining comparable performance. As we will further demonstrate in Section 5, our proposed B-PDF method achieves comparable accuracy, and offers a practical, memory-efficient alternative to MeZO with the extra benefit of incorporating second-order information.

**Flexibility in BCD Block Selection.** Beyond the natural block partitioning of model layers in ascending order, BCD can be adapted with various strategies such as descending order, random reshuffling, or importance sampling (Luo et al., 2024; Pan et al., 2024). For instance, LiSA (Pan et al., 2024) proposes a layer-wise importance sampling approach, which updates selected layers while keeping others frozen, utilizing AdamW as the optimizer. In this approach, layers are randomly selected based on predefined probability values. In Section 4.1, we will present several BCD methods for block selection. This flexibility allows BCD to be adapted to different optimization scenarios, enhancing the overall training process while maintaining memory efficiency.

## 4 METHODOLOGY

### 4.1 BCD-INTERGRATED ZEROth-ORDER NEWTON METHOD OPTIMIZER

Motivated by the revisiting of the second-order Hessian memory consumption, we identified a significant bottleneck caused by the storage of diagonal Hessian estimation, which introduced substantial memory overhead, particularly for large models. Ultimately, to address this memory-convergence trade-off, we propose a new method that integrates block coordinate descent (BCD) with a zeroth-order Newton optimizer, termed **Block-wise diagonal-Hessian Preconditioned Coordinate Descent Forward-only optimizer (B-PDF)**. Recognizing the layerwise structure of the transformer model, we treat each layer as a block for Hessian-informed zeroth-order optimization. By partitioning the model into blocks and updating only a subset of layers at each iteration, we reduce the Hessian storage requirement while maintaining the convergence benefits of second-order methods. Additionally, we update the embedding and language model (LM) head layers solely through ZO optimization to mitigate the instability and overhead typically associated with second-order methods, resulting in a more memory-efficient and scalable approach for fine-tuning large models.

The block partitioning is naturally arranged in ascending order, and various BCD algorithms can be employed to determine the active block  $\theta_{\pi_i}$ . Possible strategies include using ascending order, a layerwise importance sampling scheme based on the mean weight norms of each block  $\pi_i$ , the Gaussian-Southwell-Diagonal rule (Nutini et al., 2017), or dynamically updated probabilistic lists employing a bandit method. Formally, for the current step  $T$ , the parameter block  $\theta_{\pi_b}$  to update can be selected using several types of BCD algorithms:

$$\theta_{\pi_b} = \begin{cases} \theta_{\pi_i}, i \leftarrow [1, \dots, D], & \text{ascending / descending / random order,} \\ \arg \max_{\theta_{\pi_i}} \frac{1}{T} \sum_{t=1}^T \|\theta_{\pi_i}^t\|_2, & \text{mean weight norms,} \\ \arg \max_{\theta_{\pi_i}} \frac{|\Delta \mathcal{L}(\theta_{\pi_i})|^2}{\Sigma_{\pi_i}}, & \text{Gauss-Southwell-Diagonal rule,} \\ \theta_{\pi_z}, z \sim p_z & \text{importance sampling / bandit method} \\ \dots & \dots \end{cases} \quad (7)$$

We note that while different methods can be effective in their original settings, they vary significantly in terms of memory and computational costs. In practice, the substantial computation and storage required for updates by importance-score-based methods led us to select the more efficient default ascending order rule, and our experiments empirically demonstrate its performance. Now we present the pseudocode for the proposed algorithm in Algorithm 1.

---

**Algorithm 1** Training Pipeline of the Propose B-PDF Method.

---

```

1: Input: parameters  $\theta \in \mathbb{R}^d$ , loss function  $\mathcal{L}$ , perturbation scale  $\mu$ , learning rate  $\eta$ , smooth scale  $\alpha$ 
2: for  $t = 1, \dots, T$  do
3:   Select block  $\theta_{\pi_b} \in \theta$  according to the BCD rule
4:   if a new block is selected then
5:      $\Sigma \leftarrow \mathbf{I}_{|\theta_{\pi_b}|}$  ▷ Diagonal Hessian initialization
6:   end if
7:   Freeze other layers
8:   Sample a random seed  $s$  ▷ First-time sampling random vector
9:   for  $\mu_i = 0, +\mu, -2\mu$  do
10:    for  $\theta_i \in \theta_{\pi_b}$  do
11:      Sample  $z \sim \mathcal{N}_s(0, \mathbf{I}_{|\theta_i|})$ 
12:       $\theta_i \leftarrow \theta_i + \mu_i \Sigma_t^{1/2} z$  ▷ In-place perturbation
13:    end for
14:     $\ell_{\text{sign}(\mu_i)} = \mathcal{L}(\theta)$ 
15:  end for
16:   $\hat{\Sigma}_t \leftarrow \frac{1}{2\mu^2} (\ell_+ + \ell_- - 2\ell_{\text{original}}) (\Sigma_{t-1}^{-1/2} z_i z_i^\top \Sigma_{t-1}^{-1/2})$  ▷ Update diagonal Hessian
17:   $\Sigma_t^{-1} \leftarrow (1 - \alpha_t) \Sigma_{t-1}^{-1} + \alpha_t \left| \text{diag}(\hat{\Sigma}_t) \right|$ 
18:  projected_grad  $\leftarrow (\ell_+ - \ell_-) \Sigma_t^{1/2} / 2\mu$ 
19:  Reset random number generator with seed  $s$  ▷ For regenerating the random vector
20:  for  $\theta_i \in \theta_{\pi_b}$  do
21:    Sample  $z \sim \mathcal{N}_s(0, \mathbf{I}_{|\theta_i|})$ 
22:     $\theta_i \leftarrow \theta_i - \eta_t * \text{projected\_grad} * z$  ▷ Update weights
23:  end for
24: end for

```

---

**Remark 1.** The optimization objective is to minimize the loss function  $\mathcal{L}(\theta)$  by leveraging diagonal Hessian preconditioning within the memory-efficient framework. During each iteration, after selecting the active blocks for updates, zeroth-order optimization with a diagonal Hessian preconditioner is performed for the chosen layers. The diagonal Hessian estimate will be reinitialized for a newly selected block, and updates for that block will occur over several subsequent iterations. The algorithm applies in-place perturbations to the parameters in three steps with the perturbation scale corresponding to  $\mu_i = 0, +\mu, -2\mu$ , sampling a normally distributed random vector  $z$  to perturb the selected block  $\theta_{\pi_b}$ . For each perturbation, the loss function  $\mathcal{L}(\theta)$  is computed to estimate the gradient information. Afterwards, the diagonal Hessian is updated based on the difference in the computed losses from the perturbed parameters. The gradient for the selected block is then projected using the updated Hessian, and the weights of the active block are updated accordingly.

**Remark 2.** The proposed algorithm efficiently combines BCD with a zeroth-order Newton method by updating only a subset of model layers per iteration. This approach reduces memory usage by eliminating backpropagation and utilizing block-wise gradient updates, while maintaining convergence speed through the use of diagonal Hessian approximations. The consistent use of random vectors and selective parameter perturbation further enhance the method’s memory efficiency.

## 4.2 CONVERGENCE ANALYSIS

As a block coordinate descent variant of HiZOO (Zhao et al., 2024b), our proposed B-PDF preserves HiZOO’s convergence properties. Since our focus is primarily on the practical analysis and implementation of memory efficiency, this work emphasizes practical solutions over theoretical exploration. Nevertheless, we provide a brief summary adapted from their convergence analysis.

Adopt the classical assumptions and update rule  $\theta_{t+1} = \theta_t - \eta_t \hat{\nabla} \mathcal{L}_\mu(\theta_t)$  as detailed by Zhao et al. (2024b), with iteration number  $T$  and a suitable step size  $\eta_t$ , we have:

$$\mathbb{E}\left[\frac{1}{T} \sum_{t=1}^T \|\nabla \mathcal{L}(\theta_t)\|^2\right] \leq \mathcal{O}\left(\frac{1}{\sqrt{T}}\right)(\mathcal{L}(\theta_0) - \mathcal{L}^*) + \mathcal{O}(\mu^2), \quad (8)$$

where  $\mathcal{L}^*$  denotes the minimization of the function  $\mathcal{L}(\theta; \mathcal{B})$ . As training progresses, the first term on the right-hand side of the equation gradually diminishes to zero, while the second term remains bounded by the perturbation scale. This establishes that our proposed method converges to a bounded neighbourhood around a stationary point. Moreover, as  $T \rightarrow \infty$ , the method converges to the optimal point, as demonstrated by the equation above.

*Proof.* A brief proof adapted from HiZOO (Zhao et al., 2024b) is provided in the Appendix B. For further theoretical details, we refer readers to their original work.  $\square$

## 5 EXPERIMENTS

In this section, we build on the experimental settings of MeZO (Malladi et al., 2023) and HiZOO (Zhao et al., 2024b) to evaluate our proposed B-PDF method in terms of memory consumption, runtime, and convergence. Our experimental code builds on their open-source repositories, with the block coordinate descent method integrated. Detailed implementations are provided in Appendix A. To facilitate implementation and reduce resource requirements, we focus on performance across several GLUE and SuperGLUE tasks, following their approach. All experiments are conducted on either a single RTX 4090 (24GB) or RTX A6000 (48GB) GPU. Specific details regarding the hyperparameter grids can be found in the Appendix C.

### 5.1 EXPERIMENTS ON OPT-1.3B WITH A SINGLE RTX 4090

**Settings.** First, we conduct experiments by fine-tuning the OPT-1.3B model on a single RTX 4090 GPU. Following the settings of previous work, we select several GLUE and SuperGLUE tasks to evaluate the performance of our proposed B-PDF method. These NLP tasks include sentence classification and text generation. We note that MeZO highlights the significance of incorporating prompts for optimal performance and is structured accordingly. Therefore, we maintained MeZO’s original setup and refrained from introducing additional baselines in our experiments. For the first-order baselines, we include SGD, BCD-based SGD (referred to as BCD in the tables), and LoRA with SGD. For the zeroth-order methods, we compare MeZO, HiZOO, and our proposed B-PDF. The batch size is set to 8 for zeroth-order methods and 2 for first-order methods to prevent memory exhaustion. Our primary goal is to demonstrate that B-PDF reduces HiZOO’s memory consumption while maintaining speed and accuracy.

Table 2: Experiments on OPT-1.3B on SST-2 dataset.

Method		Accuracy	Runtime (20,000 steps)		Average Memory	
First-Order	Forward+Backward	SGD	94.3	4min 05s	22688 MB	
		BCD	92.4	3min 09s	20510 MB	
		LoRA	92.0	0min 55s	10624 MB	
Zeroth-Order	2×Forward	MeZO	91.7	54min 55s	baseline	<b>4368 MB</b> baseline
	3×Forward	HiZOO	91.7	99min 44s	+ 81.61%	7524 MB + 72.25%
	3×Forward	<b>(ours)</b>	<b>91.9</b>	<b>51min 38s</b>	- 6.98%	4589 MB + 5.06%



432  
433  
434  
435  
436  
437  
438  
439  
440  
441  
442  
443  
444  
445  
446  
447  
448  
449  
450  
451  
452  
453  
454  
455  
456  
457  
458  
459  
460  
461  
462  
463  
464  
465  
466  
467  
468  
469  
470  
471  
472  
473  
474  
475  
476  
477  
478  
479  
480  
481  
482  
483  
484  
485

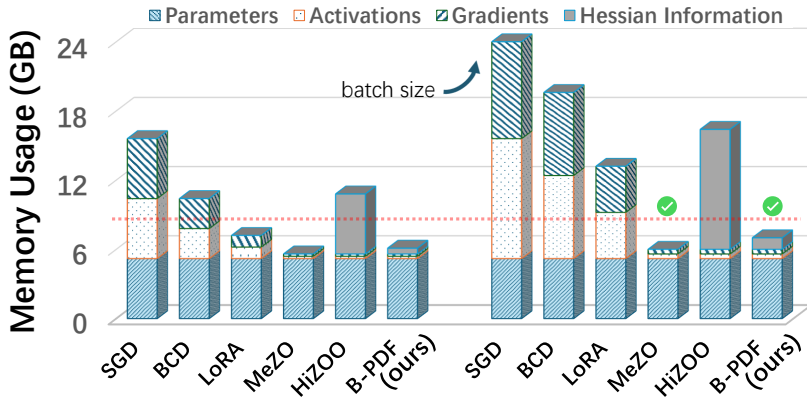


Figure 2: Illustration of average GPU memory consumption for fine-tuning the OPT-1.3B model using different methods, with  $batchsize = \{1, 2\}$ . As the batch size increases, our proposed B-PDF and MeZO maintain low memory usage, while other methods easily surpass the memory threshold of devices (red dashed line represents an 8GB memory limit on low-end devices).

**Memory Efficiency.** For memory efficiency, B-PDF significantly reduces the memory overhead of incorporating Hessian information while maintaining accuracy. As shown in Tables 1 and 2, our report on average GPU memory usage during experiments demonstrates that B-PDF has a comparable memory cost to MeZO, while offering substantial savings in memory consumption compared to HiZOO and first-order methods such as SGD, BCD-SGD, and LoRA ( $rank=8$ ). This notable improvement ensures the practical adoption of the proposed method on low-end devices, where memory is a primary bottleneck for training, which is also the original reason why the forward-only approach was developed to save memory down to inference-level requirements. This makes our method a suitable solution for low-memory training environments. In contrast, HiZOO incurs a significant 72% higher memory cost than MeZO, indicating an impractical convergence-memory tradeoff in memory-limited scenarios. Additionally, first-order methods consume even more memory due to the overhead introduced by backpropagation. For instance, BCD-SGD still requires nearly full fine-tuning memory to store activations and gradients for backpropagation. Consequently, due to their substantial memory demands, they are rendered impractical in low-end environments, making faster convergence irrelevant. This further highlights the advantages and rationale of our approach.

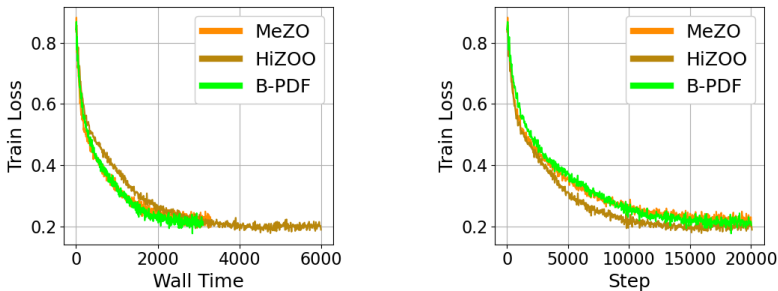


Figure 3: Convergence curves of MeZO, HiZOO and proposed B-PDF on SST-2 training OPT-1.3B.

**Convergence Study.** Regarding convergence rate, we present the convergence curve relative to wall-clock time or steps for training on the SST-2 dataset, as illustrated in Figure 3. The results show that while HiZOO converges more effectively than MeZO for 20,000 steps, it requires nearly double the completion time. Conversely, our proposed B-PDF achieves convergence between the results of MeZO and HiZOO while finishing first among the three zeroth-order methods. This speedup is due to the BCD strategy, which activates only a subset of layers (in the experiment, two layers per iteration), thereby reducing computational demands. Consequently, our method with three forward passes still finishes faster than MeZO, which requires only two forward passes.

Table 3: Experiments on OPT-1.3B across different datasets.

Task Type	Task	SST-2	RTE	CB	BoolQ	WSC	WIC	SQuAD	Average
		classification						generation	
First-Order	SGD	94.3	68.6	71.4	70.0	63.5	61.4	81.6	73.0
	BCD	92.4	69.7	69.6	63.2	63.5	61.6	78.8	71.3
	LoRA	92.4	66.4	69.6	66.8	63.5	58.5	80.5	71.1
Zeroth-Order	MeZO	91.7	64.3	69.6	65.5	63.5	57.7	77.9	70.0
	HiZOO	91.7	64.6	71.4	65.5	63.5	58.5	78.7	70.6
	<b>(ours)</b>	91.9	65.3	69.6	65.2	63.5	57.7	77.9	70.2

Furthermore, the experimental crossover results, presented in Table 3 and visualized in Figure 4 in Appendix D, demonstrate that our method achieves comparable accuracy to baseline methods across benchmarks. While first-order methods deliver superior results, their actual memory consumption is several times higher than that of zeroth-order methods, rendering them impractical for low-end environments. In contrast, our method improves memory efficiency while enhancing convergence, outperforming MeZO baseline and offering a practical, efficient solution for low-end settings. These findings position B-PDF as a memory-efficient optimizer and an effective alternative to MeZO.

Table 4: Experiments of fine-tuning LLaMA-2-7B on SST-2 dataset on an RTX A6000 (48 GB).

Zeroth-Order Method	Accuracy	Average Memory	First-Order Method	Accuracy	Average Memory
MeZO	85.2	<b>31GB</b> baseline	LoRA	94.8	41GB +32.3%
B-PDF	<b>90.6</b>	32GB + 3.23%	OOM for SGD, BCD, and HiZOO.		

## 5.2 EXPERIMENTS ON LLAMA-2-7B WITH A SINGLE RTX A6000

To further evaluate our proposed method on larger models, we fine-tuned a LLaMA-2-7B model in FP16 precision on an RTX A6000, using the aforementioned optimization algorithms, as shown in Table 4. Due to the increased model size, both the first-order method and HiZOO encountered out-of-memory (OOM) errors, and B-PDF required longer completion times because of the higher computational cost associated with the larger Hessian estimation. We compared the performance of three methods: LoRA ( $\text{rank}=8$ ), MeZO, and our proposed B-PDF, on the SST-2 dataset with  $\text{batchsize}=1$ . The remaining settings were kept consistent with those in Section 5.1. Despite the limited batch size and hardware constraints, which caused an accuracy drop from incomplete convergence, B-PDF still demonstrated performance gains while maintaining comparable memory consumption as a Hessian-informed method, unlike first-order methods draining GPU memory, underscoring its potential in memory-constrained environments.

## 6 CONCLUSION

In this paper, we propose a novel memory-efficient zeroth-order Newton method that integrates block coordinate descent (BCD) with a diagonal Hessian-preconditioned zeroth-order optimizer for fine-tuning large language models (LLMs). Our approach effectively mitigates the substantial memory overhead commonly associated with second-order methods by employing selective block-wise updates. By combining the BCD technique with the Hessian preconditioner, we achieve significant reductions in memory consumption while preserving competitive accuracy and convergence speed performance. Our extensive experiments on OPT-1.3B and LLaMA-2-7B demonstrate that our method can reduce memory usage by up to 39% compared to existing second-order optimizers while maintaining baseline accuracy across various downstream tasks. Furthermore, our approach exhibits faster wall-clock convergence than conventional zeroth-order methods, making it a practical and scalable solution for fine-tuning large models on resource-constrained devices. Future work will aim to extend this methodology to larger models and more complex tasks, as well as refine the block selection strategies to further enhance both efficiency and performance. In summary, our method provides a promising direction for memory-efficient fine-tuning of LLMs, offering practical advantages, particularly in memory-limited environments.

## REFERENCES

- 540 Andrew Brock, Theodore Lim, James M Ritchie, and Nick Weston. Freezeout: Accelerate training  
541 by progressively freezing layers. [arXiv preprint arXiv:1706.04983](#), 2017.
- 542 Han Cai, Chuang Gan, Ligeng Zhu, and Song Han. Tinytl: Reduce activations, not trainable  
543 parameters for efficient on-device learning. [arXiv preprint arXiv:2007.11622](#), 2020.
- 544 Tianqi Chen, Bing Xu, Chiyuan Zhang, and Carlos Guestrin. Training deep nets with sublinear  
545 memory cost. [arXiv preprint arXiv:1604.06174](#), 2016.
- 546 Tri Dao, Dan Fu, Stefano Ermon, Atri Rudra, and Christopher Ré. Flashattention: Fast and memory-  
547 efficient exact attention with io-awareness. [Advances in Neural Information Processing Systems](#),  
548 35:16344–16359, 2022.
- 549 Ning Ding, Yujia Qin, Guang Yang, Fuchao Wei, Zonghan Yang, Yusheng Su, Shengding Hu, Yulin  
550 Chen, Chi-Min Chan, Weize Chen, et al. Delta tuning: A comprehensive study of parameter  
551 efficient methods for pre-trained language models. [arXiv preprint arXiv:2203.06904](#), 2022.
- 552 John Duchi, Elad Hazan, and Yoram Singer. Adaptive subgradient methods for online learning and  
553 stochastic optimization. [Journal of machine learning research](#), 12(7), 2011.
- 554 Wentao Guo, Jikai Long, Yimeng Zeng, Zirui Liu, Xinyu Yang, Yide Ran, Jacob R Gardner, Osbert  
555 Bastani, Christopher De Sa, Xiaodong Yu, et al. Zeroth-order fine-tuning of llms with extreme  
556 sparsity. [arXiv preprint arXiv:2406.02913](#), 2024.
- 557 Neil Houlsby, Andrei Giurgiu, Stanislaw Jastrzebski, Bruna Morrone, Quentin De Laroussilhe,  
558 Andrea Gesmundo, Mona Attariyan, and Sylvain Gelly. Parameter-efficient transfer learning for  
559 nlp. In [International Conference on Machine Learning](#), pp. 2790–2799. PMLR, 2019.
- 560 Edward J Hu, Yelong Shen, Phillip Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang,  
561 and Weizhu Chen. Lora: Low-rank adaptation of large language models. [arXiv preprint  
arXiv:2106.09685](#), 2021.
- 562 Ziyu Jiang, Xuxi Chen, Xueqin Huang, Xianzhi Du, Denny Zhou, and Zhangyang Wang. Back  
563 razor: Memory-efficient transfer learning by self-sparsified backpropagation. [Advances in neural  
information processing systems](#), 35:29248–29261, 2022.
- 564 Diederik P Kingma. Adam: A method for stochastic optimization. [arXiv preprint arXiv:1412.6980](#),  
565 2014.
- 566 Brian Lester, Rami Al-Rfou, and Noah Constant. The power of scale for parameter-efficient prompt  
567 tuning. [arXiv preprint arXiv:2104.08691](#), 2021.
- 568 Bingrui Li, Jianfei Chen, and Jun Zhu. Memory efficient optimizers with 4-bit states. [Advances in  
Neural Information Processing Systems](#), 36, 2024.
- 569 Hong Liu, Zhiyuan Li, David Hall, Percy Liang, and Tengyu Ma. Sophia: A scalable stochastic  
570 second-order optimizer for language model pre-training. [arXiv preprint arXiv:2305.14342](#), 2023.
- 571 I Loshchilov. Decoupled weight decay regularization. [arXiv preprint arXiv:1711.05101](#), 2017.
- 572 Qijun Luo, Hengxu Yu, and Xiao Li. Badam: A memory efficient full parameter training method for  
573 large language models. [arXiv preprint arXiv:2404.02827](#), 2024.
- 574 Kai Lv, Hang Yan, Qipeng Guo, Haijun Lv, and Xipeng Qiu. Adalomo: Low-memory optimization  
575 with adaptive learning rate. [arXiv preprint arXiv:2310.10195](#), 2023a.
- 576 Kai Lv, Yuqing Yang, Tengxiao Liu, Qinghui Gao, Qipeng Guo, and Xipeng Qiu. Full parameter  
577 fine-tuning for large language models with limited resources. [arXiv preprint arXiv:2306.09782](#),  
578 2023b.
- 579 Sadhika Malladi, Tianyu Gao, Eshaan Nichani, Alex Damian, Jason D Lee, Danqi Chen, and Sanjeev  
580 Arora. Fine-tuning language models with just forward passes. [Advances in Neural Information  
Processing Systems](#), 36:53038–53075, 2023.
- 581  
582  
583  
584  
585  
586  
587  
588  
589  
590  
591  
592  
593

- 594 James Martens and Roger Grosse. Optimizing neural networks with kronecker-factored approximate  
595 curvature. In International conference on machine learning, pp. 2408–2417. PMLR, 2015.
- 596
- 597 Julie Nutini, Issam Laradji, and Mark Schmidt. Let’s make block coordinate descent converge faster:  
598 Faster greedy rules, message-passing, active-set complexity, and superlinear convergence. arXiv  
599 preprint arXiv:1712.08859, 2017.
- 600 Rui Pan, Xiang Liu, Shizhe Diao, Renjie Pi, Jipeng Zhang, Chi Han, and Tong Zhang. Lisa:  
601 Layerwise importance sampling for memory-efficient large language model fine-tuning. arXiv  
602 preprint arXiv:2403.17919, 2024.
- 603
- 604 Colin Raffel, Noam Shazeer, Adam Roberts, Katherine Lee, Sharan Narang, Michael Matena, Yanqi  
605 Zhou, Wei Li, and Peter J Liu. Exploring the limits of transfer learning with a unified text-to-text  
606 transformer. Journal of machine learning research, 21(140):1–67, 2020.
- 607 Samyam Rajbhandari, Jeff Rasley, Olatunji Ruwase, and Yuxiong He. Zero: Memory optimiza-  
608 tions toward training trillion parameter models. In SC20: International Conference for High  
609 Performance Computing, Networking, Storage and Analysis, pp. 1–16. IEEE, 2020.
- 610 Jie Ren, Samyam Rajbhandari, Reza Yazdani Aminabadi, Olatunji Ruwase, Shuangyan Yang, Minjia  
611 Zhang, Dong Li, and Yuxiong He. {ZeRO-Offload}: Democratizing {Billion-Scale} model  
612 training. In 2021 USENIX Annual Technical Conference (USENIX ATC 21), pp. 551–564, 2021.
- 613
- 614 James C Spall. Multivariate stochastic approximation using a simultaneous perturbation gradient  
615 approximation. IEEE transactions on automatic control, 37(3):332–341, 1992.
- 616 Tianxiang Sun, Zhengfu He, Hong Qian, Yunhua Zhou, Xuanjing Huang, and Xipeng Qiu. Bbtv2:  
617 Towards a gradient-free future with large language models. arXiv preprint arXiv:2205.11200,  
618 2022a.
- 619 Tianxiang Sun, Yunfan Shao, Hong Qian, Xuanjing Huang, and Xipeng Qiu. Black-box tuning  
620 for language-model-as-a-service. In International Conference on Machine Learning, pp. 20841–  
621 20855. PMLR, 2022b.
- 622
- 623 Tijmen Tieleman, Geoffrey Hinton, et al. Lecture 6.5-rmsprop: Divide the gradient by a running  
624 average of its recent magnitude. COURSERA: Neural networks for machine learning, 4(2):26–31,  
625 2012.
- 626 Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Nikolay  
627 Bashlykov, Soumya Batra, Prajjwal Bhargava, Shruti Bhosale, et al. Llama 2: Open foundation  
628 and fine-tuned chat models. arXiv preprint arXiv:2307.09288, 2023.
- 629
- 630 Alex Wang. Glue: A multi-task benchmark and analysis platform for natural language understanding.  
631 arXiv preprint arXiv:1804.07461, 2018.
- 632
- 633 Alex Wang, Yada Pruksachatkun, Nikita Nangia, Amanpreet Singh, Julian Michael, Felix Hill, Omer  
634 Levy, and Samuel Bowman. Superglue: A stickier benchmark for general-purpose language  
understanding systems. Advances in neural information processing systems, 32, 2019.
- 635
- 636 Zhewei Yao, Amir Gholami, Sheng Shen, Mustafa Mustafa, Kurt Keutzer, and Michael Mahoney.  
637 Adahessian: An adaptive second order optimizer for machine learning. In proceedings of the  
638 AAAI conference on artificial intelligence, volume 35, pp. 10665–10673, 2021.
- 639
- 640 Haishan Ye. Mirror natural evolution strategies. arXiv preprint arXiv:2308.00469, 2023.
- 641
- 642 Elad Ben Zaken, Shauli Ravfogel, and Yoav Goldberg. Bitfit: Simple parameter-efficient fine-tuning  
643 for transformer-based masked language-models. arXiv preprint arXiv:2106.10199, 2021.
- 644
- 645 Lin Zhang, Shaohuai Shi, and Bo Li. Eva: Practical second-order optimization with kronecker-  
646 vectorized approximation. In The Eleventh International Conference on Learning Representations,  
647 2022a.
- 648
- 649 Susan Zhang, Stephen Roller, Naman Goyal, Mikel Artetxe, Moya Chen, Shuohui Chen, Christopher  
650 Dewan, Mona Diab, Xian Li, Xi Victoria Lin, et al. Opt: Open pre-trained transformer language  
651 models. arXiv preprint arXiv:2205.01068, 2022b.

648 Yihua Zhang, Pingzhi Li, Junyuan Hong, Jiaxiang Li, Yimeng Zhang, Wenqing Zheng, Pin-Yu  
649 Chen, Jason D Lee, Wotao Yin, Mingyi Hong, et al. Revisiting zeroth-order optimization for  
650 memory-efficient llm fine-tuning: A benchmark. [arXiv preprint arXiv:2402.11592](#), 2024.  
651

652 Jiawei Zhao, Zhenyu Zhang, Beidi Chen, Zhangyang Wang, Anima Anandkumar, and Yuandong  
653 Tian. Galore: Memory-efficient llm training by gradient low-rank projection. [arXiv preprint](#)  
654 [arXiv:2403.03507](#), 2024a.

655 Yanjun Zhao, Sizhe Dang, Haishan Ye, Guang Dai, Yi Qian, and Ivor W Tsang. Second-order  
656 fine-tuning without pain for llms: A hessian informed zeroth-order optimizer. [arXiv preprint](#)  
657 [arXiv:2402.15173](#), 2024b.  
658  
659  
660  
661  
662  
663  
664  
665  
666  
667  
668  
669  
670  
671  
672  
673  
674  
675  
676  
677  
678  
679  
680  
681  
682  
683  
684  
685  
686  
687  
688  
689  
690  
691  
692  
693  
694  
695  
696  
697  
698  
699  
700  
701



## A IMPLEMENTATION DETAILS

The implementation of the B-PDF function is designed to enhance zero-order optimization (ZOO) by incorporating selective layer-wise updates based on Hessian-informed perturbations.

In the `zo_Hessian_step_update` function (Zhao et al., 2024b), the Hessian matrix is initialized if it does not already exist. This matrix is created by iterating over each trainable parameter of the model and initializing a tensor of ones with the same dimensions as the respective parameter. The estimate Hessian matrix serves as a second-order approximation that is updated during the optimization process.

In our framework, we introduce a layer-specific update mechanism within the optimization function `hizoo_step_update`, which implements a periodic selection of layers, referred to as "hizoo layers." These layers are chosen iteratively every fixed number of steps, with the cycle determined by a step counter.

The layers can be selected either sequentially, in an ordered manner, or via other rules such as Gauss-Southwell quadratic diagonal selection (GSQ) (Nutini et al., 2017), which prioritizes layers based on previous scores.

Following MeZO (Malladi et al., 2023) and HiZOO (Zhao et al., 2024b), we apply a noise-based perturbation to the selected Hessian-informed layers during each iteration using Gaussian noise. The noise is scaled by the square root of the corresponding Hessian matrix and a random vector sampled from a normal distribution. This approach allows the optimization process to focus on specific layers while updating their parameters iteratively.

By controlling the frequency and scope of these updates, we distribute optimization efforts across different parts of the network over time. This can ensure that updates are not applied uniformly but are instead targeted based on layer importance, thereby improving the overall efficiency of the training process.

Additionally, memory management is considered throughout the implementation, as the Hessian matrix is periodically cleared, and GPU memory is freed using `torch.cuda.empty_cache()`, ensuring that the training process remains efficient, even in memory-constrained environments.

In addition, we use `torch.clamp` API to clamp the intermediate results to meet the precision requirements and reduce the instability of second-order methods.

Overall, the B-PDF implementation introduces a structured and targeted optimization approach that leverages layer-wise perturbations to enhance the zero-order optimization process effectively.

## B CONVERGENCE ANALYSIS

As a block coordinate descent variant of HiZOO (Zhao et al., 2024b), our proposed B-PDF retains the convergence properties of HiZOO. Since our focus is on practical memory reduction rather than theoretical analysis, we offer a brief convergence analysis of our method, adapted from Zhao et al. (2024b), with adjustments made primarily for consistency. For more in-depth theoretical details, we direct readers to their original work.

We adopt several classical assumptions:

**Assumptions.** 1. The objective function  $\mathcal{L}(\boldsymbol{\theta}; \mathcal{B})$  is  $L_d$ -smooth with respect to  $\boldsymbol{\theta}_d$ , and  $L_\infty = \max_d L_d$ ; 2. The stochastic gradient  $\nabla \mathcal{L}(\boldsymbol{\theta}; \mathcal{B})$  has  $\sigma^2$  variance, i.e.  $\mathbb{E} [\|\nabla \mathcal{L}(\boldsymbol{\theta}; \mathcal{B}) - \nabla \mathcal{L}(\boldsymbol{\theta})\|^2] \leq \sigma^2$ ; 3. Each entry of  $\boldsymbol{\Sigma}_t$  lies in the range  $[\beta_\ell, \beta_u]$  with  $0 < \beta_\ell \leq \beta_u$ .

The the descent direction  $\hat{\nabla} \mathcal{L}_\mu(\boldsymbol{\theta}_t)$  defined as:

$$\hat{\nabla} \mathcal{L}_\mu = \sum_{i=1}^{\pi_b} \frac{\mathcal{L}(\boldsymbol{\theta}_t + \mu \boldsymbol{\Sigma}_t^{1/2} \mathbf{z}_i; \mathcal{B}) - \mathcal{L}(\boldsymbol{\theta}_t - \mu \boldsymbol{\Sigma}_t^{1/2} \mathbf{z}_i; \mathcal{B})}{2b\mu} \boldsymbol{\Sigma}_t^{1/2} \mathbf{z}_i. \quad (9)$$

and update rule is  $\boldsymbol{\theta}_{t+1} = \boldsymbol{\theta}_t - \eta_t \hat{\nabla} \mathcal{L}_\mu(\boldsymbol{\theta}_t)$ .

*Proof.* By the update rule of  $\boldsymbol{\theta}_t$  and above assumptions, we have

$$\begin{aligned} & \mathbb{E} [\mathcal{L}(\boldsymbol{\theta}_{t+1})] - \mathbb{E} [\mathcal{L}(\boldsymbol{\theta}_t)] \\ & \leq -\eta_t \mathbb{E} [\langle \nabla \mathcal{L}(\boldsymbol{\theta}_t), \hat{\nabla} \mathcal{L}_\mu(\boldsymbol{\theta}_t) \rangle] + \frac{L_\infty \eta_t^2}{2} \mathbb{E} [\|\hat{\nabla} \mathcal{L}_\mu(\boldsymbol{\theta}_t)\|^2] \\ & \leq -\eta_t \|\nabla \mathcal{L}(\boldsymbol{\theta}_t)\|_{\boldsymbol{\Sigma}_t}^2 + \eta_t \mathcal{O}(\mu \|\nabla \mathcal{L}(\boldsymbol{\theta}_t)\|) \\ & \quad + 2\eta_t^2 L_\infty (\text{tr}(\boldsymbol{\Sigma}_t) + \beta_u) \|\nabla \mathcal{L}(\boldsymbol{\theta}_t)\|_{\boldsymbol{\Sigma}_t}^2 \\ & \quad + 2\eta_t^2 L_\infty (\text{tr}(\boldsymbol{\Sigma}_t) + \beta_u) \sigma^2 + \mathcal{O}(\mu^2) \\ & \leq -\frac{\eta_t}{2} \|\nabla \mathcal{L}(\boldsymbol{\theta}_t)\|_{\boldsymbol{\Sigma}_t}^2 + 2\eta_t^2 L_\infty (\text{tr}(\boldsymbol{\Sigma}_t) + \beta_u) \|\nabla \mathcal{L}(\boldsymbol{\theta}_t)\|_{\boldsymbol{\Sigma}_t}^2 \\ & \quad + 2\eta_t^2 L_\infty (\text{tr}(\boldsymbol{\Sigma}_t) + \beta_u) \sigma^2 + \mathcal{O}(\mu^2) \\ & = -\frac{\eta_t}{2} (1 - 4\eta_t L (\text{tr}(\boldsymbol{\Sigma}_t) + \beta_u)) \|\nabla \mathcal{L}(\boldsymbol{\theta}_t)\|_{\boldsymbol{\Sigma}_t}^2 \\ & \quad + 2\eta_t^2 L_\infty (\text{tr}(\boldsymbol{\Sigma}_t) + \beta_u) \sigma^2 + \mathcal{O}(\mu^2) \\ & \leq -\frac{\eta_t}{4} \|\nabla \mathcal{L}(\boldsymbol{\theta}_t)\|_{\boldsymbol{\Sigma}_t}^2 + 2\eta_t^2 L_\infty (\text{tr}(\boldsymbol{\Sigma}_t) + \beta_u) \sigma^2 + \mathcal{O}(\mu^2), \end{aligned}$$

where the second inequality is derived from the following lemma (Zhao et al., 2024b):

$$\begin{aligned} \mathbb{E} [\hat{\nabla} \mathcal{L}_\mu(\boldsymbol{\theta}_t)] & = \boldsymbol{\Sigma}_t \nabla \mathcal{L}(\boldsymbol{\theta}_t) + \mathcal{O}(\mu) \\ \mathbb{E} [\|\hat{\nabla} \mathcal{L}_\mu(\boldsymbol{\theta}_t)\|^2] & \leq 4(\text{tr}(\boldsymbol{\Sigma}_t) + \beta_u) \|\nabla \mathcal{L}(\boldsymbol{\theta}_t)\|_{\boldsymbol{\Sigma}_t}^2 + 4\beta_u (\text{tr}(\boldsymbol{\Sigma}_t) + \beta_u) \boldsymbol{\Sigma}_t^2 + \mathcal{O}(\mu^2). \end{aligned}$$

By rearranging and summing over  $T$  iterations, we have:

$$\begin{aligned} \mathbb{E} \left[ \frac{1}{T} \sum_{t=1}^T \|\nabla \mathcal{L}(\boldsymbol{\theta}_t)\|^2 \right] & \leq \frac{1}{T\beta_\ell} \sum_{t=1}^T \|\nabla \mathcal{L}(\boldsymbol{\theta}_t)\|_{\boldsymbol{\Sigma}_t}^2 \\ & \leq \frac{4(\mathcal{L}(\boldsymbol{\theta}_1; \mathcal{B}) - \mathcal{L}(\boldsymbol{\theta}_*; \mathcal{B}))}{T\beta_\ell \eta} + \frac{8\eta L_\infty (\text{tr}(\boldsymbol{\Sigma}_t) + \beta_u)}{T\beta_\ell} \sigma^2 + \mathcal{O}(\mu^2) \\ & = \frac{32L_\infty (\text{tr}(\boldsymbol{\Sigma}_t) + \beta_u) (\mathcal{L}(\boldsymbol{\theta}_1; \mathcal{B}) - \mathcal{L}(\boldsymbol{\theta}_*; \mathcal{B}))}{\sqrt{T}\beta_\ell} + \frac{\sigma^2}{T^{3/2}\beta_\ell} + \mathcal{O}(\mu^2), \end{aligned}$$

where the first inequality is based on the assumption 3, and  $\eta$  selected as  $\frac{1}{8\sqrt{T}L_\infty (\max_t(\text{tr}(\boldsymbol{\Sigma}_t) + \beta_u))}$ .

756  
757  
758  
759  
760  
761  
762  
763  
764  
765  
766  
767  
768  
769  
770  
771  
772  
773  
774  
775  
776  
777  
778  
779  
780  
781  
782  
783  
784  
785  
786  
787  
788  
789  
790  
791  
792  
793  
794  
795  
796  
797  
798  
799  
800  
801  
802  
803  
804  
805  
806  
807  
808  
809

□

## C HYPERPARAMETER SEARCH

Here, we present the detailed hyperparameter grids used in our experiments. Empirically, we found that the optimal learning rate for B-PDF is an order of magnitude higher than that for MeZO. Some outlier values in the results may stem from insufficient parameter search or incomplete convergence, likely caused by limited training steps and small batch sizes due to hardware memory constraints.

Model	Method	Hyperparameters	Values
General Settings in Common		Learning rate schedule	Linear decay
		Steps	20000
		LoRA rank	8
OPT-1.3B	First-order	Batch size	{1, 2}
		Learning rate	$\{1, 3\}$ or $\{5, 7\} \times \{1e-6, 1e-7\}$
		$\mu$	$1e-3$
		Weight Decay	0
OPT-1.3B	Zeroth-order	Batch size	{1, 2, 8}
		Learning rate	$\{1, 3\}$ or $\{5, 7\} \times \{1e-5, 1e-6\}$
		$\mu$	$1e-3$
		Weight Decay	0
		Hessian Smooth Type	Constant $1e-9$
		BCD-Hessian Smooth Type	Constant $1e-5$
		BCD-Update Interval	{5, 10}
BCD-selected layers	{1, 2}		
LLaMA-2-7B	First or Zeroth-order	Batch size	{1}
		Learning rate	$\{3\} \times \{1e-6, 1e-7\}$
		$\mu$	$1e-3$
		Weight Decay	0

Table 5: The hyperparameter grids used for OPT-1.3B and LLaMA-2-7B experiments.

## D ADDITIONAL VISULIZATION RESULTS

Here, we present the bar chart illustrating the test results of OPT-1.3B.

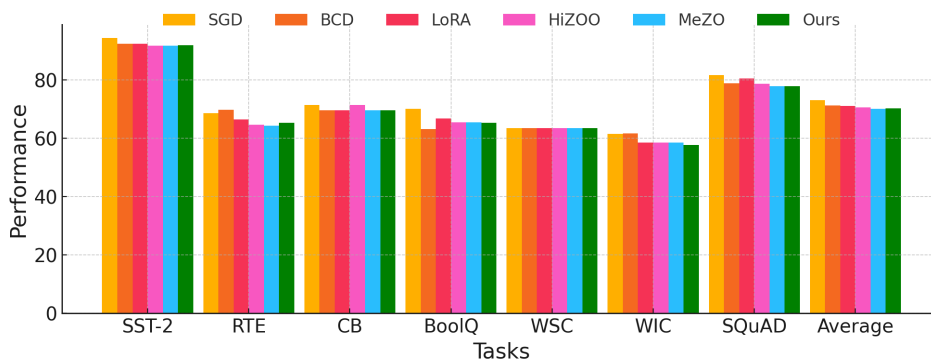


Figure 4: Bar chart illustrating the results of training OPT-1.3B with different methods across various benchmarks.