

# Flexora: Flexible Low-Rank Adaptation for Large Language Models

Chenxing Wei<sup>\*†§</sup>, Yao Shu<sup>\*‡</sup>, Ying Tiffany He<sup>†</sup>, Fei Richard Yu<sup>#†§‡</sup>

<sup>†</sup>College of Computer Science and Software Engineering, Shenzhen University, China

<sup>§</sup>Guangdong Lab of AI and Digital Economy (SZ), China

<sup>‡</sup>Hong Kong University of Science and Technology (Guangzhou), China

<sup>‡</sup>School of Information Technology, Carleton University, Canada

weichenxing2023@email.szu.edu.cn, yaoshu@hkust-gz.edu.cn

heyingszu.edu.cn, richard.yu@ieee.org

## Abstract

Large language models (LLMs) have revolutionized artificial intelligence, but their performance on specific tasks is often limited by knowledge boundaries. While fine-tuning techniques like low-rank adaptation (LoRA) aim to address this, they can suffer from overfitting. We propose *flexible low-rank adaptation* (Flexora), a novel method that automatically selects the most critical layers for fine-tuning to optimize performance across diverse downstream tasks. Flexora formulates layer selection as a hyperparameter optimization problem, employs unrolled differentiation for efficient solving, and identifies the most impactful layers based on optimized hyperparameters. Extensive experiments across various pre-trained models and natural language tasks demonstrate that Flexora consistently outperforms existing baselines. We provide theoretical insights and comprehensive ablation studies to elucidate the effectiveness of Flexora. Therefore, Flexora offers a robust solution to enhance LoRA fine-tuning for LLMs, potentially advancing the field of adaptive language model optimization.

## 1 Introduction

The advent of large language models (LLMs) (Zhao et al., 2023; Xu et al., 2023) has revolutionized artificial intelligence, offering unprecedented capabilities across various domains. However, this progress comes at a significant cost: LLMs demand substantial computational resources due to their vast parameter sets and complex functionalities (Wei et al., 2022; Touvron et al., 2023). This challenge has spurred the development of parameter-efficient fine-tuning (PEFT) methods (Li and Liang, 2021; Lester et al., 2021), with low-rank adaptation (LoRA) (Hu et al., 2021) emerging as a particularly promising approach. The innovation of LoRA lies in its ability to freeze pre-trained parameters

while introducing trainable auxiliary parameters ( $\Delta W$ ) at each layer, dramatically reducing training costs while maintaining impressive performance. However, despite its widespread adoption, LoRA is not without limitations. It can underperform on certain tasks, likely due to overfitting issues, as evidenced in benchmarks like GLUE (Wu et al., 2024b), summary tasks (Liu et al., 2024), and complex reasoning tasks (Zhang et al., 2024). Existing techniques to combat overfitting, such as dropout (Lin et al., 2024) and novel regularization strategies (Mao et al., 2024b), often yield performance comparable to or lower than vanilla LoRA and lack the flexibility to adapt across different tasks. Moreover, current methods typically require manual hyperparameter tuning, limiting their practical applicability in diverse scenarios. These challenges therefore underscore the urgent need for an algorithm that delivers superior performance, enables automatic hyperparameter tuning, and supports flexible training across various tasks.

To address these limitations, we introduce *flexible low-rank adaptation* (Flexora), a novel framework designed to flexibly fine-tune LLMs using an automated layer-level policy. Our approach is inspired by hyperparameter optimization (HPO) and offers several key innovations. First, we demonstrate that fine-tuning only the most critical layers can significantly reduce overfitting and enhance performance. Second, we frame the layer selection problem as an HPO task and employ unrolled differentiation (UD) to solve it efficiently. Third, we develop a three-stage process that automatically identifies and focuses on the most important layers for downstream tasks. As illustrated in Figure 1, Flexora operates through an initialization stage (Sec. 4.1) that injects defined hyperparameters into LoRA parameters, a flexible layer selection stage (Sec. 4.2) that optimizes these hyperparameters using UD, and a fine-tuning stage (Sec. 4.3) that selectively updates only the most

\* Equal contribution, # corresponding author.

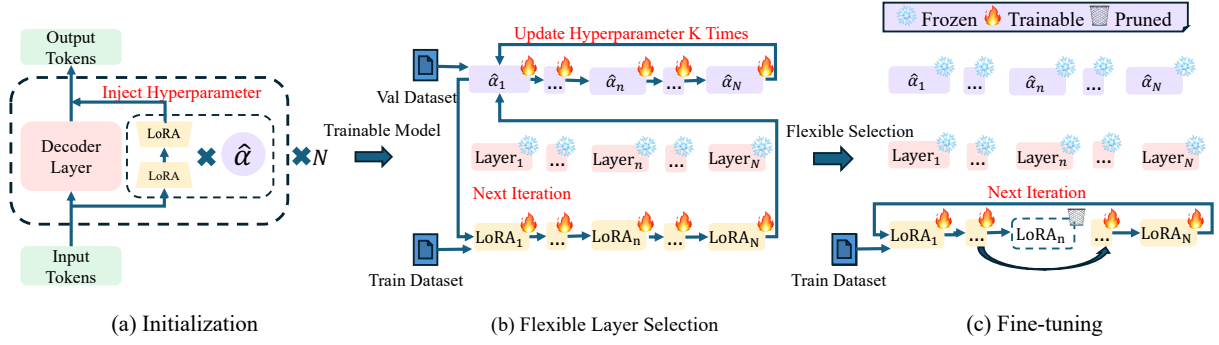


Figure 1: An overview of Flexora: (a) Initialization of hyperparameters  $\hat{\alpha}$  and their integration with LoRA parameters to produce the Trainable Model. (b) Simultaneous training of LoRA parameters and hyperparameters  $\hat{\alpha}$  using different datasets, minimizing empirical risk for both validation and training datasets. The hyperparameter vector  $\hat{\alpha}$  is then ranked based on magnitude. (c) Flexible selection of layers to be trained, where higher-ranked layers are activated for training while others remain frozen.

crucial layers, significantly reducing computational overhead. Our extensive empirical results (Sec. 5) demonstrate that Flexora effectively reduces unimportant LoRA parameters, mitigates overfitting, and enhances overall performance across a variety of tasks and model architectures.

In summary, our key contributions consist of: (a) the introduction of Flexora, a novel framework for automatic layer selection in LoRA fine-tuning; (b) a formulation of layer selection as an HPO task, efficiently solved using unrolled differentiation; (c) comprehensive validation through extensive experiments on various LLMs and downstream tasks; and (d) theoretical insights into the performance improvements achieved by Flexora, providing a deeper understanding of its effectiveness.

## 2 Related Work

**Low-Rank Adaptation (LoRA)** Low-Rank Adaptation (LoRA) methods are widely used to reduce training parameters when fine-tuning large language models (LLMs) for specific applications. However, LoRA often suffers from overfitting, which can degrade performance on downstream tasks. To mitigate this, various strategies have been proposed: LoRA-SP (Wu et al., 2024b) randomly freezes half of the LoRA parameters during fine-tuning to alleviate overfitting; LoRA-FA (Zhang et al., 2023a) freezes down-projection weights while updating only up-projection weights; VeRA (Kopiczko et al., 2024) introduces vector-based random matrix adaptation, significantly reducing trainable parameters compared to LoRA; LoRA-drop (Zhou et al., 2024) prunes less important parameters based on layer output analysis; AdaLoRA (Zhang et al., 2023b) dynamically allocates the

parameter budget across weight matrices based on importance scores; LoRAPrune (Zhang et al., 2024) jointly prunes parts of the LoRA matrix and LLM parameters based on gradients; and LoRAShear (Chen et al., 2023) employs knowledge-based structured pruning to reduce costs while enhancing generalization. Despite their benefits, these methods often (a) require significant design effort, (b) struggle to adapt across different tasks, and (c) can be overly complex for practical application. In contrast, we introduce Flexora, a framework designed for flexible LoRA fine-tuning across various tasks using a simple, automated layer-level policy.

**Hyperparameter Optimization (HPO)** HPO is widely applied across various domains. Specifically, in the domain of neural architecture search, DARTS (Liu et al., 2019) conceptualizes the coefficients defining the network architecture as hyperparameters. In the domains of feature learning, DS<sup>3</sup>L (Guo et al., 2020) considers feature extractors as hyperparameters. In the field of data science, TPOT (Olson et al., 2016) employs hyperparameters as weights to measure the importance of data. By minimizing the validation loss over these hyperparameters, the optimal variables, e.g., the architectures in Liu et al. (2019), the features in Guo et al. (2020), and the data in Olson et al. (2016), are identified, leading to superior performance in their respective domains. Drawing inspiration from these works, we initially formulated the layer selection in the LoRA method as an HPO problem. This involves optimizing hyperparameters to quantify the contributions of different layers, aiming to achieve optimal performance on downstream tasks and thereby select the most crucial layers for fine-

tuning. This formulation subsequently led to the development of our Flexora.

### 3 Preliminaries

In this section, we first provide empirical insights showing that layer selection is crucial for improving the performance of LLMs in Sec. 3.1, and then frame the layer selection problem as a well-defined HPO problem in Sec. 3.2.

#### 3.1 Empirical Insights

To study the impact of the number of LoRA fine-tuning layers on overall performance, we conducted a preliminary study using Llama3-8B (Meta, 2024) across a range of downstream tasks. Here, we randomly selected different subsets of layers, different ranks (e.g., 4, 8, 16, 32) for LoRA fine-tuning, and evaluated their performance on these tasks. The findings, shown in Figure 2 and Appendix C.9, reveal a clear trend: while increasing the number of fine-tuned layers generally improves model performance, there is a critical point beyond which fine-tuning more layers leads to potential overfitting and subsequent performance decline. This hence suggests that selecting an optimal subset of layers for LoRA fine-tuning is crucial for maximizing performance, which interestingly aligns with the previous empirical studies (Zhu et al., 2023; Zhou et al., 2024; Chen et al., 2023).

#### 3.2 Problem Formulation

Inspired by the empirical insights above, we aim to identify the most critical layers in LoRA fine-tuning to improve generalization performance across a variety of downstream tasks. Formally, we consider an  $N$ -layer LLM with LoRA fine-tuning parameters  $\theta \in \mathbb{R}^d$ , and let the hyperparameter  $\alpha \in \{0, 1\}^N$  denote the selection of fine-tuning layers, where a value of 1 indicates that a layer is selected for fine-tuning. Given the test data distribution  $\mathcal{D}_{\text{test}}$  and the training dataset  $S_{\text{train}}$ , we then define the expected test and training error as  $\mathcal{R}^{\text{test}}(\theta, \alpha) \triangleq \mathbb{E}_{x \sim \mathcal{D}_{\text{test}}} [\ell(x, \theta; \alpha)]$  and  $\mathcal{R}^{\text{train}}(\theta, \alpha) \triangleq \mathbb{E}_{x \sim S_{\text{train}}} [\ell(x, \theta; \alpha)]$ , respectively.

Hence, to select the optimal LoRA fine-tuning layers for maximized performance on downstream tasks, we aim to solve the following bilevel optimization problem:

$$\begin{aligned} \min_{\alpha \in \{0, 1\}^N} \quad & \mathcal{R}^{\text{test}}(\theta^*(\alpha), \alpha) \\ \text{s.t.} \quad & \theta^*(\alpha) = \arg \min_{\theta \in \mathbb{R}^d} \mathcal{R}^{\text{train}}(\theta, \alpha). \end{aligned} \quad (1)$$

This formulation follows a standard hyperparameter optimization (HPO) approach as demonstrated in Bao et al. (2021), where  $\alpha$  serves as the hyperparameter. Thus, the layer selection problem for LoRA fine-tuning in LLMs is framed as a well-defined HPO problem.

Unfortunately, it is typically infeasible to access the full test distributions, denoted by  $\mathcal{D}_{\text{test}}$ , for this optimization. This expected test error can typically be approximated by the empirical validation error based on the validation dataset  $S_{\text{val}}$ , which is defined as  $\hat{\mathcal{R}}^{\text{val}}(\theta, \alpha) \triangleq \mathbb{E}_{x \sim S_{\text{val}}} [\ell(x, \theta; \alpha)]$ . Therefore, (1) can be simplified as:

$$\begin{aligned} \min_{\alpha \in \{0, 1\}^N} \quad & \hat{\mathcal{R}}^{\text{val}}(\theta^*(\alpha), \alpha) \\ \text{s.t.} \quad & \theta^*(\alpha) = \arg \min_{\theta \in \mathbb{R}^d} \mathcal{R}^{\text{train}}(\theta, \alpha). \end{aligned} \quad (2)$$

### 4 The Flexora Framework

To address the layer selection problem defined above, we propose our *flexible low-rank adaptation for LLMs* (Flexora) framework in Figure 1. As illustrated in Figure 1, the Flexora framework consists of three key stages: an initial stage (detailed in Sec. 4.1), a flexible layer selection stage (detailed in Sec. 4.2), and a fine-tuning stage for the selected LoRA layers (detailed in Sec. 4.3).

#### 4.1 Initial Stage

We begin by introducing a special formulation of LoRA, which incorporates the layer selection hyperparameter  $\alpha = (\alpha^{(1)}, \dots, \alpha^{(N)}) \in \{0, 1\}^N$ , as follows:

$$h^{(i)} = Wz^{(i)} + \alpha^{(i)} B^{(i)} A^{(i)} z^{(i)}, \quad \text{s.t.} \quad \alpha_i \in \{0, 1\}. \quad (3)$$

Here,  $h^{(i)}$  is the output of the  $i$ -th layer, where  $W$  is the original weight matrix,  $z^{(i)}$  is the input, and  $B^{(i)}$  and  $A^{(i)}$  are the low-rank adaptation matrices of LoRA. The hyperparameter  $\alpha^{(i)}$  determines whether LoRA is applied for layer  $i$ . Specifically, if  $\alpha^{(i)} = 0$ , the equation simplifies to  $h^{(i)} = Wz^{(i)}$ , meaning the  $i$ -th layer reverts to standard computation without LoRA, implying that the additional complexity of LoRA is unnecessary for layer  $i$ . Conversely, when  $\alpha^{(i)} = 1$ , the equation becomes the standard LoRA form,  $h^{(i)} = Wz^{(i)} + B^{(i)} A^{(i)} z^{(i)}$ , indicating that LoRA significantly enhances the performance of layer  $i$  by allowing the low-rank matrices to better capture complex patterns. So, this dynamic adjustment allows the model to selectively apply LoRA when a

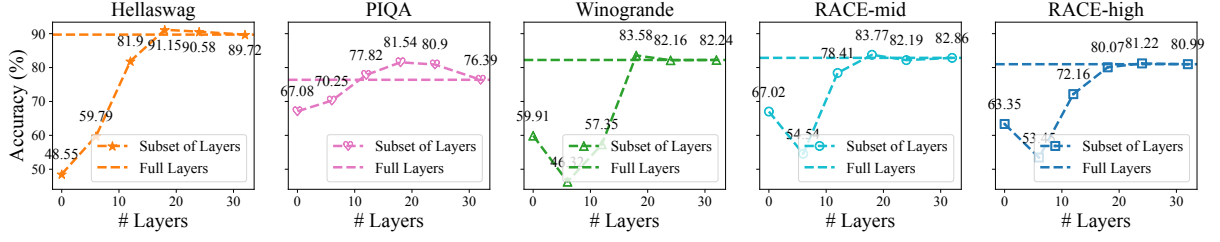


Figure 2: This figure depicts the relationship between the number of LoRA fine-tuning layers and model accuracy across four distinct datasets: HellaSwag, PIQA, Winogrande, and RACE, with the latter including two separate tasks, RACE-mid and RACE-high, which vary in difficulty. Results for LoRA rank 8 are shown here. The  $x$ -axis represents the number of fine-tuned layers, ranging from 0 to 32, where 0 corresponds to the base model without fine-tuning. Selected configurations include 6, 12, 18, 24, and 32 randomly fine-tuned layers. The full 32-layer configuration, representing the vanilla LoRA setup, is shown as a horizontal dashed line in the plots. The  $y$ -axis indicates model accuracy as a percentage.

specific layer is most beneficial, thereby optimizing the fine-tuning process and mitigating overfitting.

However, due to the inherent difficulty of directly optimizing the discrete layer selection hyperparameter  $\alpha$ , we adopt a continuous relaxation approach by replacing the  $\alpha$  in (3) with its continuous counterpart,  $\hat{\alpha} = (\hat{\alpha}^{(1)}, \dots, \hat{\alpha}^{(N)})$ :

$$h^{(i)} = W_0 z^{(i)} + \hat{\alpha}^{(i)} B^{(i)} A^{(i)} z^{(i)},$$

$$\text{s.t. } \hat{\alpha}^{(i)} = \frac{\exp(\alpha^{(i)})}{\sum_{i \in [N]} \exp(\alpha^{(i)})} N. \quad (4)$$

Notably,  $\alpha \in \mathbb{R}^N$  now and  $\alpha$  are typically initialized to zeros, providing a neutral starting point where no layer is initially excluded from LoRA fine-tuning. Meanwhile, the constant scale  $N$  ensures that when all layers are selected for fine-tuning, the scale of each selected layer for LoRA fine-tuning is preserved, resulting in  $\hat{\alpha}^{(i)} = 1$  for all layers, aligning with the vanilla LoRA scale as shown above.

## 4.2 Flexible Layer Selection Stage

**Optimization Strategy.** Given the continuous relaxation  $\hat{\alpha}$  defined above, we propose to solve the well-defined HPO problem in Equation 2 using the widely applied unrolled differentiation (UD) method (Franceschi et al., 2017, 2018; Fu et al., 2016; Maclaurin et al., 2015; Shaban et al., 2019). The UD method typically involves two alternating optimization processes: (a) the inner-level and (b) the outer-level optimization. In this paper, the outer-level optimization is defined as  $\arg \min_{\theta \in \mathbb{R}^d} \mathcal{R}^{\text{train}}(\theta, \alpha)$ , in which the layer selection hyperparameter  $\alpha$  is fixed, and the LoRA parameters  $\theta$  are updated using stochastic gradient methods (e.g., SGD (Sra et al., 2011)) on the train-

### Algorithm 1 The Flexora Framework

---

```

1: Input: Number of steps  $T$  and  $K$ ; Initialized LoRA parameters  $\theta_0$  and hyperparameter  $\alpha_0 = 0$ ; Learning rate  $\eta_\alpha$  and  $\eta_\theta$ 
2: for  $t = 0$  to  $T - 1$  do
3:   Sample a mini-batch  $B_{\text{train}} \sim S_{\text{train}}$ 
4:    $\theta_{t+1} \leftarrow \eta_\theta \nabla_\theta \left( \frac{1}{|B_{\text{train}}|} \sum_{x \in B_{\text{train}}} \ell(x, \theta; \alpha_t) \right) \big|_{\theta=\theta_t}$ 
5:    $\alpha_{t+1}^0 \leftarrow \alpha_t$ 
6:   for  $k = 0$  to  $K - 1$  do
7:     Sample a mini-batch  $B_{\text{val}} \sim S_{\text{val}}$ 
8:      $\alpha_{t+1}^{k+1} \leftarrow \eta_\alpha \nabla_\alpha \left( \frac{1}{|B_{\text{val}}|} \sum_{x \in B_{\text{val}}} \ell(x, \theta_{t+1}; \alpha) \right) \big|_{\alpha=\alpha_{t+1}^k}$ 
9:   end for
10:   $\alpha_{t+1} \leftarrow \alpha_{t+1}^K$ 
11: end for
12: return  $\alpha^* = \alpha_T$ 

```

---

ing dataset  $S_{\text{train}}$ . This step focuses on optimizing model performance by adjusting the parameters associated with the selected layers (line 3-4 in Algorithm 1). Meanwhile, the inner-level optimization is  $\arg \min_{\alpha \in \mathbb{R}^N} \hat{\mathcal{R}}^{\text{val}}(\theta, \alpha)$ , in which the layer selection hyperparameter  $\alpha$  is updated using stochastic gradient methods (e.g., SGD) based on the validation performance of the optimized LoRA parameters  $\theta$  from the inner-level process (lines 6-9 in Algorithm 1). This step intends to maximize the validation performance of LoRA fine-tuning based on a subset of selected layers. These two alternating processes therefore iteratively refine both the model parameters and the layer selection criteria, making LoRA layer selection more computationally efficient in practice. After  $T$  iterations of these alternating processes, the converged  $\alpha_T$  is output as the optimal layer selection denoted as  $\alpha^*$  (line 12 in Algorithm 1).



Table 1: Comparison of accuracy across various common sense reasoning tasks using Llama3-8B. The baseline experimental configuration is detailed in Appendix B. Here, "Pre-trained" refers to using the base model for reasoning, "Full FT" indicates full parameter fine-tuning, and "Random (Greedy)" represents the best result from randomly selected layers. Unless otherwise specified, the results are based on the default LoRA Rank of 8.

Methods	Hellaswag	PIQA	Winogrande	RACE-mid	RACE-high	Average
Pre-trained	48.55	67.08	59.91	67.02	63.35	61.18
Full FT	90.53	79.32	81.16	81.92	79.36	82.46
LoRA( $r = 8$ )	89.72	76.39	82.24	82.86	80.99	83.04
LoRA( $r = 16$ )	89.99	78.47	82.77	81.63	79.68	82.51
LoRA( $r = 32$ )	90.01	79.56	84.36	82.36	80.99	83.46
LoRA-SP	89.37	78.97	83.67	83.27	79.01	82.86
LoRA-FA	89.16	75.97	82.16	82.79	79.03	81.83
VeRA	90.98	78.63	83.64	83.55	78.84	83.13
LoRAPrune (Ratio = 0.5)	88.42	77.12	81.23	82.96	80.42	82.03
AdaLoRA ( $r_0 = 4$ )	90.17	80.20	77.19	83.15	77.93	81.73
LoRA-drop	91.86	77.91	76.46	77.30	75.24	79.75
Random (Greedy)	91.15	81.54	83.58	83.77	81.22	84.25
<b>Flexora</b> ( $r = 8$ )	93.62	85.91	<b>85.79</b>	84.61	82.36	86.46
<b>Flexora</b> ( $r = 16$ )	93.71	85.26	84.99	<b>85.62</b>	<b>83.03</b>	<b>86.52</b>
<b>Flexora</b> ( $r = 32$ )	<b>93.87</b>	<b>86.02</b>	85.01	84.27	81.97	86.23

**Selection Strategy.** To begin with, we introduce the following proposition:

**Proposition 1.** *If  $\alpha$  is initialized to zeros, then for any  $T \geq 0$  and  $K \geq 0$  in Alg. 1,  $\sum_{i=1}^N \alpha^{(i)} = 0$ .*

The proof of this proposition is provided in Appendix A.1. This result highlights that the mean value of the hyperparameter  $\alpha$  remains 0, indicating that after the layer selection stage, the elements in the optimized hyperparameter  $\alpha^*$  can take on both positive and negative values. Therefore, we propose to determine the layers for LoRA fine-tuning by selecting layers with  $\alpha^{(i)} > 0$ , as these layers are expected to make positive contributions. In contrast, layers with  $\alpha^{(i)} \leq 0$  are believed to be less beneficial or even harmful to LoRA fine-tuning. As a result, this method not only facilitates automatic layer selection but also provides flexibility in adjusting the number and specific layers for LoRA fine-tuning, helping to mitigate the potential overfitting and improve overall performance.

### 4.3 Fine-Tuning Stage

During the fine-tuning stage, as illustrated in Figure 1c, we adopt a selective activation strategy. In this phase, we freeze the layers not selected for fine-tuning, keeping their parameters unchanged, and focus on retraining only the selected layers to enhance performance. This targeted approach concentrates computational resources on the most critical layers for the downstream task. By retraining the LoRA parameters from scratch in these

layers, the model adaptively learns optimal representations, reducing the risk of overfitting and improving performance, especially for simpler tasks. We will validate this approach with the empirical results presented below.

## 5 Empirical Results

In this section, we present comprehensive experiments to support the effectiveness of our Flexora framework with datasets and experimental setup detailed in Sec. 5.1, main results detailed in Sec. 5.2, and ablation studies detailed in Sec. 5.3.

### 5.1 Datasets and Setup

To evaluate the performance of our proposed Flexora method, we primarily focus on reasoning and reading comprehension tasks. Since Flexora is the first algorithm to select layers based on specific downstream tasks, we refer to and modify the dataset selection process of Hu et al. (2023). We selected the Winogrande(Sakaguchi et al., 2019), PIQA(Bisk et al., 2019), and Hellaswag(Zellers et al., 2019) reasoning benchmarks as recommended by Hu et al. (2023), and additionally included the reading comprehension benchmark RACE(Lai et al., 2017). Each of these datasets, measured by accuracy, has independent training, validation, and test sets. In all experiments, we use the training set to train the LoRA parameters, the validation set to tune the hyperparameters introduced by Flexora, and

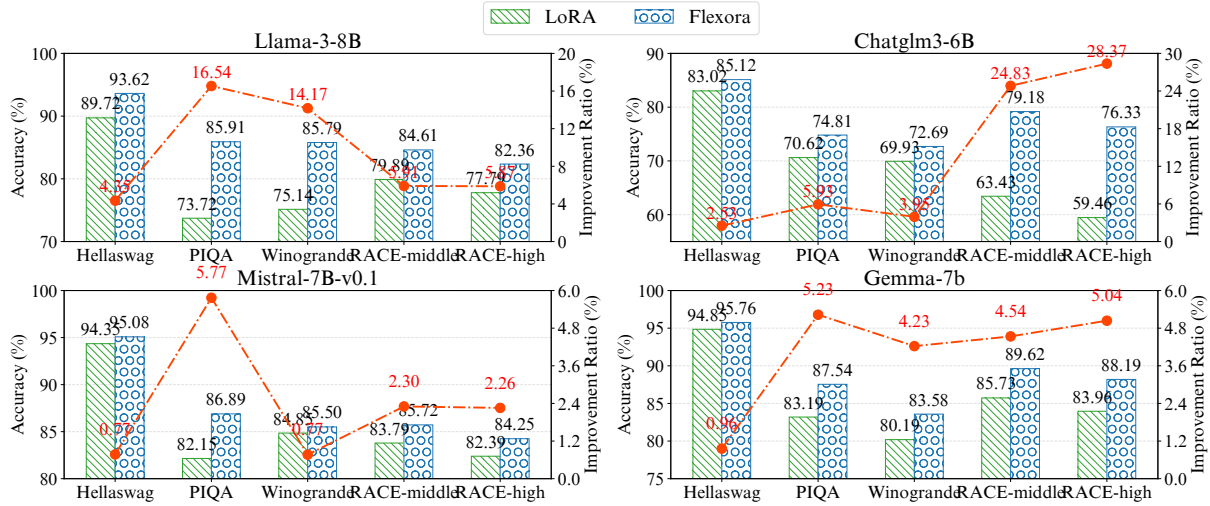


Figure 3: Comparison of the accuracy of various models (Llama-3-8B, ChatGLM3-6B, Mistral-7B-v0.1, and Gemma-7B) across different tasks. Bars with green diagonal stripes represent LoRA accuracy, while blue circles indicate Flexora accuracy, and the red dotted line represents the improvement ratio of Flexora over LoRA. Notably, Flexora generally outperforms LoRA in most tasks and models, demonstrating its effectiveness.

finally the test set for evaluation. It is important to emphasize that the test set remains unseen during the training phase. Our experimental setup includes 11 mainstream large-scale language models (LLMs), such as Llama3-8B (Meta, 2024) and others. Our Flexora method is implemented on the Llama-factory framework (Zheng et al., 2024) and evaluated using the Opencompass framework (Contributors, 2023). The benchmarks for comparison include pre-trained models, Full FT, LoRA, and various LoRA enhancement methods that reduce trainable parameters, such as LoRAPrune, AdaLoRA, LoRA-drop, and others. Detailed descriptions of the experimental setup are provided in Appendix B. All experiments are conducted on a single NVIDIA A100 GPU.

## 5.2 Main Results

In this section, we evaluate the performance improvement of Flexora on Llama3-8B, and the results are listed in Table 1. The loss metrics are discussed in Appendix D.1. The results show that Flexora outperforms all baseline methods. Specifically, compared with full fine-tuning and LoRA, Flexora fine-tunes 0.02% and 50% of its parameters, respectively, to achieve superior performance. This demonstrates that fine-tuning too many parameters can lead to overfitting, which not only fails to improve the performance of the model on downstream tasks but may also reduce the generalization ability of the model due to the overfitting effect. Therefore, it is crucial to select the layers most relevant

to the downstream tasks for optimization. The flexible layer selection stage of Flexora is able to consider the relationship between the pre-trained parameters of each LLM layer and the downstream task. This stage effectively identifies the most critical layers for various downstream tasks and minimizes the risk of model overfitting by focusing on training these layers, resulting in excellent performance. In Table 1, we also compare Flexora with other methods that attempt to enhance the model by reducing model parameters. Particularly, the experimental results using LoRAShear are detailed in Appendix C.4 and the experimental results using Flexora on full-parameter fine-tuning and on the instruction model are detailed in Appendix C.5 and C.6 respectively. The results show that Flexora can most accurately identify the most important parameters to achieve the largest performance improvement. Flexora introduces a flexible layer selection stage, but incurs no additional computational overhead (see Appendix C.1 for details). Flexora is a vertical method, and the discussion of integration with many LoRA enhancement methods is detailed in Appendix C.3. We also evaluate Flexora at different LoRA ranks, and the results show that changing the rank has a negligible impact on the performance of Flexora. The specific layers selected are listed in Table 18 in the Appendix. It is worth noting that the layers selected by Flexora are roughly consistent under different level conditions, which shows that Flexora effectively identifies the layers that are most suitable for downstream tasks.

Table 2: Comparison of the accuracy of different randomly selected fine-tuning layers with the same number of fine-tuning layers. We fixed the number of fine-tuning layers to match the number selected by Flexora, ensuring that the number of fine-tuning parameters remained constant while the layers were randomly selected for fine-tuning.

Methods	Hellaswag	PIQA	Winogrande	RACE-mid	RACE-high	Average
Random 1	92.97	82.91	80.98	83.98	81.10	84.39
Random 2	93.11	80.79	76.09	<b>85.45</b>	81.16	83.32
Random 3	92.52	80.47	83.50	84.54	81.93	84.59
Random (Avg.)	92.87	81.39	80.19	84.66	81.40	84.10
Flexora	<b>93.62</b>	<b>85.91</b>	<b>85.79</b>	84.61	<b>82.36</b>	<b>86.46</b>

Table 3: Comparison of the performance of models with and without a fine-tuning phase on various tasks.

Methods	Hellaswag	PIQA	Winogrande	RACE-mid	RACE-high	Average
Flexora (w/o Fine-Tuning Stage)	48.93	80.20	66.38	62.72	60.76	63.80
Flexora (w/ Fine-Tuning Stage)	<b>93.62</b>	<b>85.91</b>	<b>85.79</b>	<b>84.61</b>	<b>82.36</b>	<b>86.46</b>

We also discuss the impact of different search samples on the training time and final performance of the flexible layer selection stage, as detailed in Appendix C.7. In addition, Flexora shows strong generalization and scalability across different LLMs. As shown in Figure 3 and explained in detail in Appendix C.2, almost all LLMs can significantly improve performance with fewer fine-tuning parameters by leveraging Flexora. In Appendix G, we compare Flexora with LoRA using specific cases. The model fine-tuned with Flexora outperforms LoRA on challenging cases and provides correct explanations for answers not seen in the training set, demonstrating its strong learning and generalization capabilities. Finally, we discuss the impact of two hyperparameters  $K$  and  $T$  introduced by Algorithm 1 on the results. The results show that changes in  $K$  and  $T$  have little effect on the layer selection results and model performance. For a more detailed discussion, see Appendix C.8.

### 5.3 Ablation Studies

**Effective Layer Selection in Flexora.** In the first ablation experiment, we maintained the number of layers selected by Flexora unchanged but chose different layers for fine-tuning, aiming to verify whether Flexora selected the right layers. The experimental results are shown in Table 2. The result underscores two key points: First, Flexora can precisely determine the number of layers for fine-tuning. Even when the specific fine-tuning layers are chosen at random, the results continue to outperform LoRA. The theoretical explanation for this result can be found in Sec. 6. Secondly, Flexora enables adaptive layer selection for fine-tuning, optimizing performance and generalization by focus-

ing on crucial layers while mitigating local optima-induced performance degradation (see Appendix F for details). Appendix C.10 provides an analysis of the characteristics of the selected layers in Flexora, revealing a distinct layer selection pattern. The loss metrics are discussed in Appendix D.2.

**Flexible Layer Selection in Flexora.** In the second ablation experiment, we manually determine the number of fine-tuning layers and compare Flexora with random selection, highlighting the flexibility of Flexora. The results in Table 4 show that it can achieve the best performance regardless of the number of fine-tuning layers. The specific layers selected are shown in Table 19. The loss metrics are discussed in Appendix D.3. A noteworthy observation is that Flexora usually chooses the initial and final layers. An intuitive explanation is that the initial and final layers of the model have a significant impact on the data. The initial layers directly contact the original input, while the final layers are related to the model output, rendering them crucial. In addition, for the same downstream task, the input of the initial layer is consistent and closely coupled to the task, and the output of the final layer is also consistent. Focusing on optimizing these layers can improve learning efficiency. This conclusion has also been confirmed by other studies. LoRAShear(Chen et al., 2023) observed that the knowledge distribution in LLM is mainly concentrated in the initial and final layers. LASER(Sharma et al., 2023) revealed steep loss gradients in both initial and final layers, enhancing model training efficacy. LISA(Pan et al., 2024) found much higher weight norms in initial and final layers, indicating their increased importance.

Table 4: Comparison of the accuracy of fine-tuning a subset of layers. We standardized the number of layers to be fine-tuned and compared the performance of layers selected by Flexora against those selected randomly.

Methods	Hellaswag	PIQA	Winogrande	RACE-mid	RACE-high	Average
Random (6 Layers)	59.79	70.25	46.32	54.54	53.45	56.87
Flexora (First 6 Layers)	<b>60.04 (+0.25)</b>	<b>77.20 (+6.95)</b>	<b>57.54 (+11.22)</b>	<b>69.71 (+15.17)</b>	<b>58.35 (+4.90)</b>	<b>64.57 (+7.70)</b>
Random (12 Layers)	81.90	77.82	57.35	78.41	72.16	73.53
Flexora (First 12 Layers)	<b>88.85 (+6.95)</b>	<b>79.71 (+1.89)</b>	<b>65.82 (+8.47)</b>	<b>79.42 (+1.01)</b>	<b>72.33 (+0.17)</b>	<b>77.23 (+3.70)</b>
Random (18 Layers)	91.15	81.54	83.58	83.77	81.22	84.25
Flexora (First 18 Layers)	<b>91.31 (+0.16)</b>	<b>82.21 (+0.67)</b>	<b>84.69 (+1.11)</b>	<b>84.07 (+0.30)</b>	<b>81.53 (+0.31)</b>	<b>84.76 (+0.51)</b>
Random (24 Layers)	90.58	80.90	82.16	82.19	79.22	83.01
Flexora (First 24 Layers)	<b>91.01 (+0.43)</b>	<b>81.21 (+0.31)</b>	<b>82.87 (+0.71)</b>	<b>83.53 (+1.34)</b>	<b>80.22 (+1.00)</b>	<b>83.77 (+0.76)</b>

Table 5: Performance comparison under different training configurations

Methods(Epoch)	Hellaswag	PIQA	Winogrande	RACE-mid	RACE-high	Average
LoRA (3)	89.72	76.39	82.24	82.86	80.99	82.32
LoRA (4)	89.74	76.27	82.47	82.73	81.04	82.45
Flexora (1+3)	93.62 (+3.90)	85.91 (+9.52)	85.79 (+3.55)	84.61 (+1.75)	82.36 (+1.37)	86.46 (+4.14)

**Importance of the Fine-Tuning Stage** In the third ablation experiment, we investigated the significance of the Fine-Tuning Stage in the Flexora method by comparing model performance from the Flexible Layer Selection Stage and the Fine-Tuning Stage on the test set. Results in Table 3 show that omitting the Fine-Tuning Stage significantly degrades performance. This is because the layer selection stage outputs continuous values  $\hat{\alpha}^{(i)} \in [0, 1]^N$ , while we need discrete  $\alpha \in \{0, 1\}^N$  values. The discrepancy between continuous and discrete  $\alpha$  values leads to a performance gap. The Fine-Tuning Stage is crucial as it addresses this gap by refining the model to better approximate the discrete  $\alpha$  values, thereby mitigating the performance loss.

**Performance of Flexora from training framework** To determine whether the performance of Flexora enhancement primarily originates from its layer selection mechanism rather than an extended effective training duration (due to its two-stage process), we conducted a meticulous ablation study. A central question was whether exposing the model to the data set twice, once during layer selection (1 epoch) and again during fine-tuning (3 epochs) - artificially inflates the observed performance. We compared three distinct configurations: a) **LoRA (3 epochs)**: The standard LoRA baseline, trained for 3 epochs. b) **LoRA (4 epochs)**: An extended training control, designed to evaluate if merely increasing training epochs of LoRA to match Flexora’s total fine-tuning-equivalent duration yields comparable gains. c) **Flexora (1+3 epochs)**: The proposed Flexora method, which dedicates 1 epoch to layer selection followed by 3 epochs of fine-tuning on

the identified layers.

The results, detailed in Table 5, provide clear insights. Extending LoRA training from 3 to 4 epochs (LoRA (4 epochs)) yielded a marginal average improvement of +0.13%, suggesting that LoRA’s performance largely converges by the third epoch for these tasks. In contrast, Flexora (1+3 epochs) significantly outperformed both LoRA configurations, achieving an average score of 86.46%. This represents a substantial +4.14% improvement over the LoRA (3 epochs) baseline and +4.01% over LoRA (4 epochs). These findings robustly confirm that the primary driver of Flexora’s superior performance is its layer selection mechanism, not simply the cumulative number of training iterations. The framework’s capacity to identify and concentrate optimization efforts on task-critical layers results in significant performance enhancements that surpass those achievable by merely scaling up the training duration of standard LoRA.

## 6 Theoretical Insights

In this section, we provide theoretical explanations for *why Flexora (using only a subset of LoRA layers) can achieve excellent results*. We first introduce Theorem 1 below, and then derive our general Proposition 2, aiming to offer theoretical insights.

**Theorem 1** (Theorem 3.8 in (Hardt et al., 2016)). *Assume that  $f(\cdot; z) \in [0, 1]$  is an  $L$ -Lipschitz and  $\beta$ -smooth loss function for every sample  $z$ . Suppose that we run stochastic gradient method (e.g., SGD) for  $T$  steps with monotonically non-increasing step sizes  $\eta_t \leq c/t$  ( $t \in [T]$ ), and the number of samples is  $m$ . In particular, omitting*



constant factors that depend on  $\beta$ ,  $c$ , and  $L$ , we have  $\mathcal{R}^{\text{test}}(\theta, \eta) \leq \mathcal{R}^{\text{train}}(\theta, \eta) + \frac{T^{1-1/(\beta c+1)}}{m}$ .

Theorem 1 reveals that if all the conditions except for  $\beta$  in Theorem 1 remain the same, a smaller smoothness  $\beta$  will typically result in a smaller test error  $\mathcal{R}^{\text{test}}(\theta, \eta)$ , indicating a better generalization performance in practice. The specific definition of smoothness  $\beta$  can be found in Appendix A.2. To show how the number of LoRA layers is related to this  $\beta$ , we then follow the practice in (Shu et al., 2020) to prove our Proposition 2 below.

**Proposition 2.** *For an  $N$ -layer linear multi-layer perceptron (MLP):  $y^{(N)} \triangleq \prod_{j=1}^N W^{(j)} \mathbf{x}$  with MSE function  $\ell \triangleq (y^{(N)} - y)^2/2$  where  $y$  denotes the true label, let  $\lambda^{(i)} = \|W^{(i)}\|$  for any  $i \in [N]$ , we then have  $\left\| \frac{\partial \ell}{\partial W_1^{(i)}} - \frac{\partial \ell}{\partial W_2^{(i)}} \right\| \leq \left( \prod_{j=1, j \neq i}^N \lambda^{(j)} \right)^2 \|\mathbf{x}\|^2 \|W_1^{(i)} - W_2^{(i)}\|$ .*

The proof of Proposition 2 is in Appendix A.3. Given Proposition 2, the block-wise smoothness  $\beta_i^{(N)}$  on layer  $i \in [N]$  of an  $N$ -th layer MLP can be bounded by:  $\beta_i^{(N)} \leq \left( \prod_{j=1, j \neq i}^N \lambda^{(j)} \right)^2 \|\mathbf{x}\|^2$ . From this bound, we can see that as the number of layers  $N$  increases, the upper bound of  $\beta_i^{(N)}$  will also be increasing as  $\lambda^{(i)} > 1$  for  $i \in [N]$ . Thus, shallow MLP of fewer layers are more likely to have smaller overall smoothness  $\beta$ . Thanks to this smaller overall smoothness  $\beta$ , shallow MLP of fewer layers are more likely to achieve a smaller generalization gap (i.e., the second term on the right-hand side of Theorem 1) than deep MLP with more layers. When the training error  $\mathcal{R}^{\text{train}}(\theta, \eta)$  is the same, that is, both shallow and deep MLPs are fully trained to converge, the shallower MLP may have a lower test error  $\mathcal{R}^{\text{test}}(\theta, \eta)$  and thus may exhibit better performance on downstream tasks. To demonstrate that Proposition 2 is also applicable to the Transformer model and LoRA method, we present theoretical insights and experiments in Appendix E. These experiments demonstrate that, the smoothness of the Transformer model also increases exponentially with the number of layers.

We can now answer the question posed earlier. Flexora employs LoRA adapters to a subset of LLM layers, effectively reducing the smoothness of the network. When sufficiently trained to convergence, the aforementioned theory suggests that networks with less smoothness are more likely to better generalization and performance on downstream

tasks. In summary, the reason Flexora achieves excellent results is that it makes the model more suitable for downstream tasks.

## 7 Conclusion

We introduce Flexora, a method to improve the efficiency and effectiveness of fine-tuning in LLMs by automatically selecting key layers, by formulating layer selection as an HPO problem, and using UD. Experiments show Flexora decreases parameters, mitigates overfitting, is scalable and outperforms baselines.

## Limitations

In this section, we aim to highlight some potential considerations that may lead to suboptimal performance of Flexora. The layer selection strategy in Flexora is primarily based on the magnitude of the optimized hyperparameters. If the validation set used for optimizing these hyperparameters is too small, especially when the downstream task is complex, it may result in the optimization process converging to a hyperparameter gap that is too narrow. In such cases, the layer selection strategy may fail, leading to the incorrect choice of layers for subsequent optimization stages, ultimately resulting in poor performance. To address the issue of having a minimal validation set for different datasets, we conducted additional experiments on search samples, as detailed in Appendix C.7. These experiments demonstrate that an insufficient number of samples can indeed lead to poor performance. However, this issue can be mitigated by increasing the number of search samples. Furthermore, although Flexora is a vertical method and can theoretically be combined with all LoRA methods, there are certain methods for which fine-tuning only specific layers significantly impacts the model’s fine-tuning effectiveness. In such cases, these methods may not be compatible with Flexora.

## Ethics Statement

We have manually reevaluated the dataset we created to ensure it is free of any potential for discrimination, human rights violations, bias, exploitation, and any other ethical concerns.

## References

- Fan Bao, Guoqiang Wu, Chongxuan Li, Jun Zhu, and Bo Zhang. 2021. [Stability and generalization of bilevel programming in hyperparameter optimization](#).
- Yonatan Bisk, Rowan Zellers, Ronan Le Bras, Jianfeng Gao, and Yejin Choi. 2019. [Piqa: Reasoning about physical commonsense in natural language](#).
- Charles Blair. 1985. Problem complexity and method efficiency in optimization (as nemirovsky and dubudn). *Siam Review*, 27(2):264.
- Tianyi Chen, Tianyu Ding, Badal Yadav, Ilya Zharkov, and Luming Liang. 2023. [Lorashear: Efficient large language model structured pruning and knowledge recovery](#).
- OpenCompass Contributors. 2023. Opencompass: A universal evaluation platform for foundation models. <https://github.com/open-compass/opencompass>.
- Luca Franceschi, Michele Donini, Paolo Frasconi, and Massimiliano Pontil. 2017. [Forward and reverse gradient-based hyperparameter optimization](#).
- Luca Franceschi, Paolo Frasconi, Saverio Salzo, Riccardo Grazi, and Massimiliano Pontil. 2018. [Bilevel programming for hyperparameter optimization and meta-learning](#).
- Jie Fu, Hongyin Luo, Jiashi Feng, Kian Hsiang Low, and Tat-Seng Chua. 2016. [Drmad: Distilling reverse-mode automatic differentiation for optimizing hyperparameters of deep neural networks](#).
- Lan-Zhe Guo, Zhen-Yu Zhang, Yuan Jiang, Yu-Feng Li, and Zhi-Hua Zhou. 2020. Safe deep semi-supervised learning for unseen-class unlabeled data. In *International conference on machine learning*, pages 3897–3906. PMLR.
- Moritz Hardt, Ben Recht, and Yoram Singer. 2016. Train faster, generalize better: Stability of stochastic gradient descent. In *International conference on machine learning*, pages 1225–1234. PMLR.
- Edward J. Hu, Yelong Shen, Phillip Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang, and Weizhu Chen. 2021. [Lora: Low-rank adaptation of large language models](#).
- Zhiqiang Hu, Lei Wang, Yihuai Lan, Wanyu Xu, Ee-Peng Lim, Lidong Bing, Xing Xu, Soujanya Poria, and Roy Lee. 2023. [LLM-adapters: An adapter family for parameter-efficient fine-tuning of large language models](#). In *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*, pages 5254–5276, Singapore. Association for Computational Linguistics.
- Damjan Kalajdzievski. 2023. [A rank stabilization scaling factor for fine-tuning with lora](#).

- Dawid Jan Kopiczko, Tijmen Blankevoort, and Yuki M Asano. 2024. [VeRA: Vector-based random matrix adaptation](#). In *The Twelfth International Conference on Learning Representations*.
- Guokun Lai, Qizhe Xie, Hanxiao Liu, Yiming Yang, and Eduard Hovy. 2017. [Race: Large-scale reading comprehension dataset from examinations](#).
- Brian Lester, Rami Al-Rfou, and Noah Constant. 2021. [The power of scale for parameter-efficient prompt tuning](#).
- Xiang Lisa Li and Percy Liang. 2021. [Prefix-tuning: Optimizing continuous prompts for generation](#).
- Yang Lin, Xinyu Ma, Xu Chu, Yujie Jin, Zhibang Yang, Yasha Wang, and Hong Mei. 2024. [Lora dropout as a sparsity regularizer for overfitting control](#).
- Hanxiao Liu, Karen Simonyan, and Yiming Yang. 2019. [Darts: Differentiable architecture search](#).
- Zeyu Liu, Souvik Kundu, Anni Li, Junrui Wan, Lianghao Jiang, and Peter Anthony Bearel. 2024. [Aflora: Adaptive freezing of low rank adaptation in parameter efficient fine-tuning of large models](#).
- Dougal Maclaurin, David Duvenaud, and Ryan P. Adams. 2015. [Gradient-based hyperparameter optimization through reversible learning](#).
- Yulong Mao, Kaiyu Huang, Changhao Guan, Ganglin Bao, Fengran Mo, and Jinan Xu. 2024a. [Dora: Enhancing parameter-efficient fine-tuning with dynamic rank distribution](#).
- Yuzhu Mao, Siqi Ping, Zihao Zhao, Yang Liu, and Wenbo Ding. 2024b. [Enhancing parameter efficiency and generalization in large-scale models: A regularized and masked low-rank adaptation approach](#).
- Meta. 2024. Introducing meta llama 3: The most capable openly available LLM to date. *Meta Blog*.
- Randal S. Olson, Nathan Bartley, Ryan J. Urbanowicz, and Jason H. Moore. 2016. [Evaluation of a tree-based pipeline optimization tool for automating data science](#).
- Rui Pan, Xiang Liu, Shizhe Diao, Renjie Pi, Jipeng Zhang, Chi Han, and Tong Zhang. 2024. [Lisa: Layer-wise importance sampling for memory-efficient large language model fine-tuning](#).
- Keisuke Sakaguchi, Ronan Le Bras, Chandra Bhagavatula, and Yejin Choi. 2019. [Winogrande: An adversarial winograd schema challenge at scale](#).
- Amirreza Shaban, Ching-An Cheng, Nathan Hatch, and Byron Boots. 2019. [Truncated back-propagation for bilevel optimization](#).
- Pratyusha Sharma, Jordan T. Ash, and Dipendra Misra. 2023. [The truth is in there: Improving reasoning in language models with layer-selective rank reduction](#).
- Yao Shu, Wei Wang, and Shaofeng Cai. 2020. [Understanding architectures learnt by cell-based neural architecture search](#).
- Suvrit Sra, Sebastian Nowozin, and Stephen J Wright. 2011. *Optimization for machine learning*, page 351–368. Mit Press.
- Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier Martinet, Marie-Anne Lachaux, Timothée Lacroix, Baptiste Rozière, Naman Goyal, Eric Hambro, Faisal Azhar, Aurelien Rodriguez, Armand Joulin, Edouard Grave, and Guillaume Lample. 2023. [Llama: Open and efficient foundation language models](#).
- Jason Wei, Yi Tay, Rishi Bommasani, Colin Raffel, Barret Zoph, Sebastian Borgeaud, Dani Yogatama, Maarten Bosma, Denny Zhou, Donald Metzler, Ed H. Chi, Tatsunori Hashimoto, Oriol Vinyals, Percy Liang, Jeff Dean, and William Fedus. 2022. [Emergent abilities of large language models](#).
- Taiqiang Wu, Jiahao Wang, Zhe Zhao, and Ngai Wong. 2024a. [Mixture-of-subspaces in low-rank adaptation](#).
- Yichao Wu, Yafei Xiang, Shuning Huo, Yulu Gong, and Penghao Liang. 2024b. [Lora-sp: Streamlined partial parameter adaptation for resource-efficient fine-tuning of large language models](#).
- Zhengxuan Wu, Aryaman Arora, Zheng Wang, Atticus Geiger, Dan Jurafsky, Christopher D. Manning, and Christopher Potts. 2024c. [Reft: Representation finetuning for language models](#).
- Lingling Xu, Haoran Xie, Si-Zhao Joe Qin, Xiaohui Tao, and Fu Lee Wang. 2023. [Parameter-efficient fine-tuning methods for pretrained language models: A critical review and assessment](#).
- Rowan Zellers, Ari Holtzman, Yonatan Bisk, Ali Farhadi, and Yejin Choi. 2019. [Hellaswag: Can a machine really finish your sentence?](#)
- Longteng Zhang, Lin Zhang, Shaohuai Shi, Xiaowen Chu, and Bo Li. 2023a. [Lora-fa: Memory-efficient low-rank adaptation for large language models fine-tuning](#).
- Mingyang Zhang, Hao Chen, Chunhua Shen, Zhen Yang, Linlin Ou, Xinyi Yu, and Bohan Zhuang. 2024. [Loraprune: Pruning meets low-rank parameter-efficient fine-tuning](#).
- Qingru Zhang, Minshuo Chen, Alexander Bukharin, Nikos Karampatziakis, Pengcheng He, Yu Cheng, Weizhu Chen, and Tuo Zhao. 2023b. [Adalora: Adaptive budget allocation for parameter-efficient fine-tuning](#).
- Wayne Xin Zhao, Kun Zhou, Junyi Li, Tianyi Tang, Xiaolei Wang, Yupeng Hou, Yingqian Min, Beichen Zhang, Junjie Zhang, Zican Dong, Yifan Du, Chen Yang, Yushuo Chen, Zhipeng Chen, Jinhao Jiang, Ruiyang Ren, Yifan Li, Xinyu Tang, Zikang Liu, Peiyu Liu, Jian-Yun Nie, and Ji-Rong Wen. 2023. [A survey of large language models](#).

Yaowei Zheng, Richong Zhang, Junhao Zhang, Yanhan Ye, Zheyang Luo, Zhangchi Feng, and Yongqiang Ma. 2024. [Llamafactory: Unified efficient fine-tuning of 100+ language models](#). In *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 3: System Demonstrations)*, Bangkok, Thailand. Association for Computational Linguistics.

Hongyun Zhou, Xiangyu Lu, Wang Xu, Conghui Zhu, Tiejun Zhao, and Muyun Yang. 2024. [Lora-drop: Efficient lora parameter pruning based on output evaluation](#).

Yunqi Zhu, Xuebing Yang, Yuanyuan Wu, and Wensheng Zhang. 2023. [Parameter-efficient fine-tuning with layer pruning on free-text sequence-to-sequence modeling](#).



## A Theorems and proofs

We first prove Proposition 1, then introduce the theorems proposed by (Blair, 1985) and (Hardt et al., 2016), which reveal the properties of  $\beta$ -smooth, a necessary theoretical basis for proving Proposition 2. Finally, we prove Proposition 2.

### A.1 Proof of proposition 1

The proof of Proposition 1 is expressed as follows:

*Proof.* It is easy to verify that

$$\frac{\partial \hat{\alpha}^{(j)}}{\partial \alpha^{(i)}} = \begin{cases} \hat{\alpha}^{(j)}(1 - \frac{1}{n}\hat{\alpha}^{(j)}), & \text{if } j = i \\ -\frac{1}{n}\hat{\alpha}^{(j)}\hat{\alpha}^{(i)}, & \text{if } j \neq i \end{cases}.$$

Therefore, given that  $\sum_{i=1}^n \hat{\alpha}^{(i)} = n$

$$\begin{aligned} \sum_{i=1}^n \frac{\partial \hat{\mathcal{R}}^{\text{val}}}{\partial \hat{\alpha}^{(j)}} \frac{\partial \hat{\alpha}^{(j)}}{\partial \alpha^{(i)}} &= \frac{\partial \hat{\mathcal{R}}^{\text{val}}}{\partial \hat{\alpha}^{(j)}} \left( \hat{\alpha}^{(j)} - \frac{1}{n} (\hat{\alpha}^{(j)})^2 - \frac{1}{n} \sum_{i=1, i \neq j}^n \hat{\alpha}^{(j)} \hat{\alpha}^{(i)} \right) \\ &= \frac{\partial \hat{\mathcal{R}}^{\text{val}}}{\partial \hat{\alpha}^{(j)}} \left( \hat{\alpha}^{(j)} - \frac{\hat{\alpha}^{(j)}}{n} \sum_{i=1}^n \hat{\alpha}^{(i)} \right) \\ &= 0. \end{aligned}$$

When applying SGD to update  $\alpha$ , we have

$$\sum_{i=1}^n \alpha^{(i)} - \eta \sum_{i=1}^n \sum_{j=1}^n \frac{\partial \hat{\mathcal{R}}^{\text{val}}}{\partial \hat{\alpha}^{(j)}} \frac{\partial \hat{\alpha}^{(j)}}{\partial \alpha^{(i)}} = \sum_{i=1}^n \alpha^{(i)}.$$

That is, the updated  $\alpha$  shares the same summation as the one before the updates, which therefore concludes our proof.  $\square$

### A.2 Definition of $\beta$ -Smooth

**Definition 1.**  $\beta$ -smooth refers to the Lipschitz continuity of the gradient of the loss function, that is, for all  $w$  and  $w'$ :

$$\|\nabla f(w; z) - \nabla f(w'; z)\| \leq \beta \|w - w'\|$$

where  $\|\cdot\|$  denotes the norm of the vector, and  $f(w; z)$  is the loss function with parameter  $w$  for sample  $z$ .

Let  $f_{\text{deep}}(w)$  and  $f_{\text{shallow}}(w)$  be the loss functions for deep and shallow architectures, respectively. According to Definition 1, the relationship between  $\beta_{\text{deep}}$  and  $\beta_{\text{shallow}}$  illustrates the relationship between the generalization and performance of deep and shallow networks.

### A.3 Proof of proposition 2

**Abstract LLM into a layered network:**(Shu et al., 2020) As shown in Figure 4, we abstract LLM into a hierarchical network, and the weight of each layer is represented by  $W^{(i)}$ . Figure 4 represents the general case. The output of the  $i$ -th layer network is:

$$\mathbf{y} = \prod_{j=1}^n W^{(j)} \mathbf{x}. \quad (5)$$

**Gradient analysis:** For the abstract network, represented in Equation 5. The gradient of the loss function  $\ell$  with respect to the weight  $W^{(i)}$  is:

$$\frac{\partial \ell}{\partial W^{(i)}} = \left( \prod_{j=i+1}^n W^{(j)} \right) \frac{\partial \ell}{\partial \mathbf{y}^{(i)}} \mathbf{x} \left( \prod_{j=1}^{i-1} W^{(j)} \right). \quad (6)$$

The proof of Proposition 2 is expressed as follows:

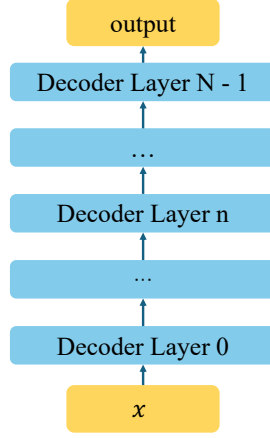


Figure 4: We present LLM as a hierarchical network. In this context, all parameters of a Decoder layer are represented as a weight matrix  $W$  for subsequent analysis.

*Proof.* For the abstract network, we begin with Definition 1:

$$\begin{aligned} & \left\| \frac{\partial \ell}{\partial W_1^{(i)}} - \frac{\partial \ell}{\partial W_2^{(i)}} \right\| \\ &= \left\| \left( \prod_{j=i+1}^n W^{(j)} \right) \left( \frac{\partial \ell}{\partial \mathbf{y}_1^{(i)}} - \frac{\partial \ell}{\partial \mathbf{y}_2^{(i)}} \right) \mathbf{x} \left( \prod_{j=1}^{i-1} W^{(j)} \right) \right\|. \end{aligned} \quad (7)$$

Taking MSE Loss as an example, for one predictions  $\mathbf{y}^{(N)}$  and their corresponding true values  $\mathbf{y}$ :

$$\ell \triangleq (\mathbf{y}^{(N)} - \mathbf{y})^2 / 2, \quad (8)$$

therefore:

$$\frac{\partial \ell}{\partial \mathbf{y}^{(i)}} = (\mathbf{y}^{(N)} - \mathbf{y}) \prod_{j=i+1}^N W^{(j)}. \quad (9)$$

We select the MSE loss function and calculate the  $i$ -th layer of  $N$  layers network, Substituting Equation 9 into Equation 7:

$$\begin{aligned} \left\| \frac{\partial \ell}{\partial W_1^{(i)}} - \frac{\partial \ell}{\partial W_2^{(i)}} \right\| &= \left\| \left( \prod_{j=i+1}^N W^{(j)} \right)^2 (\mathbf{y}_1^{(N)} - \mathbf{y} - \mathbf{y}_2^{(N)} + \mathbf{y}) \mathbf{x} \left( \prod_{j=1}^{i-1} W^{(j)} \right) \right\| \\ &\leq \left( \frac{1}{\lambda^{(i)}} \left( \prod_{j=1}^N \lambda^{(j)} \right) \right) \left( \prod_{j=i+1}^N \lambda^{(j)} \right) \left\| \prod_{j=1}^{i-1} W^{(j)} (W_1^{(i)} - W_2^{(i)}) \mathbf{x}^2 \right\| \\ &\leq \left( \prod_{j=1, j \neq i}^N \lambda^{(j)} \right) \|\mathbf{x}^2\| \left( \prod_{j=i+1}^N \lambda^{(j)} \right) \left( \prod_{j=1}^{i-1} \lambda^{(j)} \right) \| (W_1^{(i)} - W_2^{(i)}) \| \\ &\leq \left( \prod_{j=1, j \neq i}^N \lambda^{(j)} \right)^2 \|\mathbf{x}^2\| \|W_1^{(i)} - W_2^{(i)}\|, \end{aligned} \quad (10)$$

which therefore concludes our proof.  $\square$

## B Experimental setting

In the main experiment, we compared Flexora with the baseline. The datasets and experimental parameters are as follows:

### B.1 Dataset

In this section, we introduce the statistics of the dataset and the additional processing performed on the dataset. The statistics of the dataset are shown in Table 6. In addition, We added new templates to the original dataset to ensure the model could complete the required tasks and output formats. It is important to note that the added templates did not alter the original dataset, and special processing was performed for different LLMs. The specific examples are as follows:

#### Dataset Format of Hellaswag

```
dataset: Hellaswag
dataset format:
{
  "instruction": "{Article}\n
Question: {Question}\n
A. {Option A}\n
B. {Option B}\n
C. {Option C}\n
D. {Option D}\n
You may choose from 'A', 'B', 'C', 'D'.\n Answer:",
  "output": "{Answer}"
}
example:
{
  "instruction": "A man is sitting on a roof. He\n
Question: Which ending makes the most sense?\n
A. is using wrap to wrap a pair of skis.\n
B. is ripping level tiles off.\n
C. is holding a Rubik's cube.\n
D. starts pulling up roofing on a roof.\n
You may choose from 'A', 'B', 'C', 'D'.\n Answer:",
  "output": "D"
}
```

#### Dataset Format of PIQA

```
dataset: PIQA
dataset format:
{
  "instruction": "There is a single choice question.
Answer the question by replying A or B.\n
Question: {Question}\n
A. {Option A}\n
B. {Option B}\n
Answer:",
  "output": "{Answer}"
}
example:
{
  "instruction": "There is a single choice question.
Answer the question by replying A or B.\n
Question: When boiling butter, when it's ready, you can\n
A. Pour it onto a plate\n
B. Pour it into a jar\n
Answer:",
  "output": "B"
}
```

## Dataset Format of Winogrande

```
dataset: Winogrande
dataset format:
{
  "instruction": "There is a single choice question,
you need to choose the correct option to fill in the blank.
Answer the question by replying A or B."\n
Question: {Question}\n
A. {Option A}\n
B. {Option B}\n
Answer:",
  "output": "{Answer}"
}
example:
{
  "instruction": "There is a single choice question,
you need to choose the correct option to fill in the blank.
Answer the question by replying A or B."\n
Question: Sarah was a much better surgeon than Maria so _ always got the
easier cases.\n
A. Sarah\n
B. Maria\n
Answer:",
  "output": "B"
}
```

## Dataset Format of RACE

```
dataset: RACE
dataset format:
{
  "instruction": "{Article}
{Question}\n
[ {Option A}, {Option B}\, {Option C}, {Option D}]",
  "output": "{Answer}"
}
example:
{
  "instruction": "I am a psychologist. I first met Timothy, a quiet,
overweight eleven-year-old boy, when his mother brought him to me to discuss
his declining grades. A few minutes with Timothy were enough to confirm that
his self-esteem and general happiness were falling right along with _ .
I asked about Timothy's typical day. He awoke every morning at six thirty
so he could reach his school by eight and arrived home around four thirty each
afternoon. He then had a quick snack, followed by either a piano lesson
or a lesson with his math tutor. He finished dinner at 7 pm, and then he sat
down to do homework for two to three hours. Quickly doing the math in my
head, I found that Timothy spent an average of thirteen hours a day
at a writing desk.\n
What if Timothy spent thirteen hours a day at a sewing machine instead of
a desk? We would immediately be shocked, because that would be called
children being horribly mistreated. Timothy was far from being mistreated,
but the mountain of homework he faced daily resulted in a similar consequence
--he was being robbed of his childhood. In fact, Timothy had no time
to do anything he truly enjoyed, such as playing video games, watching
movies, or playing board games with his friends.\n
Play, however, is a crucial part of healthy child development.
It affects children's creativity, their social skills, and even their brain
development. The absence of play, physical exercise, and freefrom social
interaction takes a serious toll on many children. It can also cause
significant health problems like childhood obesity, sleep problems
and depression.\nExperts in the field recommend the minutes children
spend on their homework should be no more than ten times the number
of their grade level.\nWhat did the writer think of Timothy after
learning about his typical day?\n
['Timothy was very hardworking.',
'Timothy was being mistreated.',
'Timothy had a heavy burden.',
'Timothy was enjoying his childhood.'],
  "output": "C"
}
```

## B.2 Specific experimental parameters

Based on the Llama3-8B model configuration, several adjustments were made to optimize model performance. In the baseline model experiment, generation parameters were adjusted to ensure the correct



Table 6: Number of samples in the train, validation, and test datasets for various dataset.

Number of samples	train dataset	validation dataset	test dataset
Hellaswag	39900	10000	10000
PIQA	16000	2000	3000
Winogrande	40398	1267	1767
RACE	87866	4887	4934

Table 7: Detailed experimental parameters. This table lists the specific parameters we used in the experiments for various methods. These parameters include the target module of LoRA (Lora Target), the maximum sequence length (Max Length), the number of samples for supervised fine-tuning (SFT Samples), the learning rate (LR), the number of search samples (Search Samples), the initial rank (Init Rank), the target rank (Target Rank), and the ratio of pruning (Ratio). Epoch represents the epoch of training. In particular, the epochs of Flexora in the Flexora Layer Selection stage and the Fine-tuning stage are different. In the table, the former is the epoch of the Flexible Layer Selection stage and the latter is the epoch of the Fine-tuning stage. All other parameters not listed here remain consistent across all experiments.

Methods	LoRA Target	Max Length	SFT Samples	LR	Search Samples	Init Rank	Target Rank	Ratio	Epoch
LoRA	q & v Proj	1024	20000	0.0001	-	-	-	-	3
Flexora	q & v Proj	1024	20000	0.0001	20000	-	-	-	1/3
AdaLoRA	q & v Proj	1024	20000	0.0001	-	4	8	-	3
LoRA-drop	q & v Proj	1024	20000	0.0001	20000	-	-	-	3
LoRAShear	q & v Proj	1024	20000	0.0001	20000	-	-	0.5	3
Dora	q & v Proj	1024	20000	0.0001	20000	-	-	-	3
rsLoRA	q & v Proj	1024	20000	0.0001	20000	-	-	-	3
LoRAPrune	q & v Proj	1024	20000	0.0001	20000	-	-	0.5	3

output. In the LoRA experiment, adjustments to the generation parameters were retained, and LoRA-related parameters were adjusted. In the Flexora experiment, the size of the validation set was adjusted to control the time required to search for the optimal layer. In the AdaLoRA experiment, the initial rank size was modified to ensure that the fine-tuning parameters are consistent with Flexora. In the LoRA-drop experiment, the number of fine-tuning layers was set to be consistent with Flexora to ensure that the fine-tuning parameters are consistent. In the LoRAShear experiment, the pruning ratio was modified, where the parameter amount with a pruning ratio of 50% is consistent with Flexora. For specific experimental parameters, see the table 7.

### B.3 Other LLMs experimental parameters

In order to explore the versatility and scalability of Flexora, we conducted experiments on multiple different LLMs. The specific training parameters are shown in Table 8.

## C More results

### C.1 Computational Overhead of Flexora

This section analyzes the computational overhead of Flexora, focusing on the flexible layer selection stage and comparing the overall cost with LoRA.

**Flexible Layer Selection Cost** Flexora operates in two phases: (1) flexible layer selection and (2) fine-tuning. The layer selection phase identifies the optimal layer combination, and its computational cost scales with the number of samples used for the search. Table 15 shows the average search time increases with the number of samples, but remains manageable. For example, searching with 1,000 samples takes 0.08 hours, while searching with 10,000 samples takes 0.8 hours. This demonstrates the efficiency of the search process, even for larger datasets. It is important to note that the cost associated with the flexible layer selection phase is relatively low, especially when considering the significant improvements

Table 8: Detailed LLM experiment parameters. This table provides a comprehensive overview of the specific parameters used for different large language models (LLMs) in our experiments. These parameters include the LoRA alpha value (LoRA Alpha), the dropout rate of LoRA (LoRA Dropout), the rank used in LoRA (LoRA Rank), and the target module of LoRA (LoRA Target). In addition, the table lists the specific templates used for each LLM, which are derived from Llama-factory (Template). For experiments involving different downstream tasks using the same model, all other parameters are kept consistent to ensure fair comparison and best performance.

LLM	LoRA Alpha	LoRA Dropout	LoRA Rank	LoRA Target	Template (From Llama-factory)
Llama3	16	0	8	q & v Proj	llama3
Llama	16	0	8	q & v Proj	default
Llama2	16	0	8	q & v Proj	llama2
chatglm3	16	0	8	query_key_value	chatglm3
Mistral-v0.1	16	0	8	q & v Proj	mistral
gemma	16	0	8	q & v Proj	gemma
zephyr	16	0	8	q & v Proj	zephyr
vicuna	16	0	8	q & v Proj	vicuna
xuanyuan	16	0	8	q & v Proj	xuanyuan
qwen1.5	16	0	8	q & v Proj	qwen
yi	16	0	8	q & v Proj	yi

in accuracy it yields. While the search time does increase with the number of samples, the overall cost remains acceptable, particularly in tasks where accuracy is of paramount importance. In such scenarios, the benefits of identifying the optimal layer combination far outweigh the modest computational expense incurred during the search process. Moreover, the efficiency of this search process allows for exploration of a wider range of potential layer combinations, increasing the likelihood of discovering highly effective architectures that contribute to improved model performance and generalization. This efficient layer selection strategy is a crucial component of Flexora, enabling it to achieve state-of-the-art results without incurring prohibitive computational costs.

**Comparison with LoRA** During fine-tuning, Flexora significantly reduces both training time and the number of trainable parameters compared to LoRA, as shown in Table 9. Flexora reduces training time by 4.0% to 22.6% and the number of trainable parameters by 41.2% to 50.0% across various datasets. Importantly, the total computational cost of Flexora (search plus fine-tuning) is comparable to LoRA’s fine-tuning cost alone. For instance, on Hellaswag, LoRA fine-tuning requires 5.30 hours, while Flexora takes  $4.71 + 0.08 = 4.79$  hours. On Winogrande, LoRA requires 4.96 hours, and Flexora takes  $3.84 + 0.08 = 3.92$  hours. This shows Flexora doesn’t introduce significant additional overhead compared to LoRA, while achieving better performance and efficiency.

**Resource-Constrained Scenarios** Flexora’s efficiency is particularly beneficial in resource-constrained settings. The layer search can use a small number of samples (e.g., 1,000), requiring minimal resources (e.g., 0.08 hours). The reduction in trainable parameters and training time further makes Flexora suitable for deployment in resource-limited environments. In conclusion, Flexora exhibits minimal computational overhead, comparable to LoRA, while substantially improving efficiency and performance, making it a practical approach for fine-tuning large language models.

## C.2 The results of other LLMs experiment

**Wide Applicability of Flexora.** According to the parameter settings in Table 8, the verification results for various LLMs are presented in Table 10. The selected LLMs include Llama3-8B, Llama-7B, Llama2-7B, ChatGLM3-6B, Mistral-7B-v0.1, Gemma-7B, Zephyr-7B-beta, Vicuna-7B-v1.5, XuanYuan-6B, Qwen1.5-7B, and Yi-6B. These models demonstrate unique characteristics in terms of training data, architecture design, and optimized training. First, the models utilize varied training data, leading to differences in data distribution. Additionally, some models have enhanced attention mechanisms: Mistral-7B-v0.1 employs grouped query attention (GQA) and sliding window attention (SWA), while ChatGLM3-

Table 9: Comparison of training time and parameters, with the green font indicating the reduction ratio, is conducted on a single NVIDIA A100 GPU using Llama3-8B. The time metric reflects the wallclock time for the fine-tuning phase of LoRA and Flexora, excluding the layer selection phase.

Metrics	Method	Hellaswag	PIQA	Winogrande	RACE
Time (h)	LoRA	5.30	4.03	4.96	8.37
	Flexora	<b>4.71 (11.1%)</b>	<b>3.87 (4.0%)</b>	<b>3.84 (22.6%)</b>	<b>7.46 (10.9%)</b>
# Params (M)	LoRA	3.4	3.4	3.4	3.4
	Flexora	<b>2.00 (41.2%)</b>	<b>1.70 (50.0%)</b>	<b>1.70 (50.0%)</b>	<b>1.70 (50.0%)</b>

6B features a special attention design to support tool calling and code execution capabilities. Activation functions vary across these models. Llama3-8B uses the SwiGLU activation function, inspired by the PaLM model, to improve performance and convergence speed, while ChatGLM3-6B uses the Swish activation function. Furthermore, differences in reasoning optimization and multilingual capabilities contribute to varied reasoning abilities across fields. The experimental result of each model is shown in Table 10, which presents the scores of each model on different downstream tasks after LoRA and Flexora fine-tuning. It should be noted that all models fine-tuned using LoRA will have a certain degree of overfitting, while Flexora can effectively identify and analyze unnecessary layers in specific downstream tasks and prune them to reduce model overfitting. After optimization by Flexora, these LLMs showed significant performance improvements on downstream tasks. In particular, models that originally performed poorly on some tasks, such as ChatGLM3-6B, experienced significant improvements, achieving more than a 15% increase on the RACE-mid and RACE-high tasks. This improvement is attributable to the key layer selection by Flexora and efficient model learning. In summary, Flexora is applicable across Transformer models of various structures, excels in diverse tasks, and effectively enhances areas where model capabilities are lacking.

### C.3 The results of other LoRAs experiment

**Strong Scalability of Flexora.** Recently, as highlighted in the introduction, many LoRA improvement methods have been proposed and have achieved excellent performance in specific fine-tuning tasks. In this section, we explore the potential of combining our algorithm with other emerging LoRA algorithms. Four promising LoRA variants are selected from different methods, each demonstrating impressive performance. Specifically, DoRA (Decomposed Low Rank Adaptation by Weight) (Mao et al., 2024a) achieves low-rank adaptation through weight decomposition, and rsLoRA (Rank-Stabilized LoRA) (Kalajdzievski, 2023) addresses the slow training speed of traditional LoRA by introducing a rank-stable scaling factor when increasing the rank. These methods primarily solve the parameter overfitting problem within the LoRA parameters but overlook the overall overfitting problem. By innovatively combining these methods with our algorithm, we first address the overall overfitting problem and then tackle the overfitting issue of the remaining LoRA parameters, thereby significantly improving performance. Additionally, we attempt to integrate with other methods to enhance the representation ability of LoRA. For instance, MoSLoRA (Mixture-of-Subspaces in Low-Rank Adaptation) (Wu et al., 2024a) decomposes LoRA into subspaces via structural re-parameterization, employing a learnable mixer to fuse more subspaces more flexibly. LoReFT (Low-rank Linear Subspace ReFT) (Wu et al., 2024c) is a parameter-efficient finetuning method that operates on a frozen base model, learning task-specific interventions on hidden representations. The specific experimental results are shown in Table 11. The results indicate that Flexora can be effectively integrated with DoRA and rsLoRA, alleviating the overfitting problem of LLM and improving performance with less than half of the parameters. Notably, the integration of Flexora and LoReFT can further enhance performance. Flexora helps LoReFT identify the most suitable layer for fine-tuning, avoiding performance loss caused by manually selecting the fine-tuning layer. However, MoSLoRA is not suitable for integration with Flexora because MoSLoRA combines the A and B matrices of all LoRA layers. Deleting a layer would cause significant changes and degrade performance. The specific implementation requires replacing LoRA with DoRA, rsLoRA, MoSLoRA, or LoReFT for inner layer optimization during the flexible

Table 10: Detailed comparison of the accuracy of different LLMs. This table presents a comprehensive comparison of the accuracy results obtained by fine-tuning various mainstream Large Language Models (LLMs) using Flexora and LoRA methods. The accuracy metrics are reported across multiple benchmark datasets, including HellaSwag, PIQA, Winogrande, RACE-mid, and RACE-high. The average accuracy across all datasets is also provided. The exact values of accuracy improvements for each method, highlighted in red, indicate the performance gains achieved.

Methods	Hellaswag	PIQA	Winogrande	RACE-mid	RACE-high	Average
Llama3-8B-LoRA	89.72	76.39	82.24	82.86	80.99	83.04
Llama3-8B-Flexora	93.62 (+3.90)	85.91 (+9.52)	85.79 (+3.55)	84.61 (+1.75)	82.36 (+1.37)	86.46 (+3.42)
Llama-7B-LoRA	76.10	69.80	67.01	75.69	70.81	71.88
Llama-7B-Flexora	85.28 (+9.18)	71.93 (+2.13)	74.11 (+7.10)	81.62 (+5.93)	78.62 (+7.81)	78.31 (+6.43)
Llama2-7B-LoRA	79.60	75.90	78.60	79.32	75.07	77.70
Llama2-7B-Flexora	90.89 (+11.29)	81.72 (+5.82)	82.85 (+4.25)	84.89 (+5.57)	83.19 (+8.12)	84.71 (+7.01)
Chatglm3-6B-LoRA	83.02	70.62	69.93	63.43	59.46	69.29
Chatglm3-6B-Flexora	85.12 (+2.10)	74.81 (+4.19)	72.69 (+2.76)	79.18 (+15.75)	76.33 (+16.87)	77.63 (+8.33)
Mistral-7B-v0.1-LoRA	94.35	82.15	84.85	83.79	82.39	85.51
Mistral-7B-v0.1-Flexora	95.08 (+0.73)	86.89 (+4.74)	85.50 (+0.65)	85.72 (+1.93)	84.25 (+1.86)	87.49 (+1.98)
Gemma-7B-LoRA	94.85	83.19	80.19	85.73	83.96	85.58
Gemma-7B-Flexora	95.76 (+0.91)	87.54 (+4.35)	83.58 (+3.39)	89.62 (+3.89)	88.19 (+4.23)	88.94 (+3.35)
Zephyr-7B-beta-LoRA	93.77	75.03	78.37	83.45	82.25	82.57
Zephyr-7B-beta-Flexora	95.05 (+1.28)	85.58 (+10.55)	84.95 (+6.58)	86.19 (+2.74)	84.30 (+2.05)	87.21 (+4.64)
Vicuna-7B-v1.5-LoRA	87.64	69.48	63.85	67.30	73.90	72.43
Vicuna-7B-v1.5-Flexora	90.43 (+2.79)	79.49 (+10.01)	76.06 (+12.21)	82.94 (+15.64)	81.90 (+8.00)	82.16 (+9.73)
XuanYuan-6B-LoRA	82.38	74.16	65.27	78.04	72.11	74.39
XuanYuan-6B-Flexora	88.41 (+6.03)	79.43 (+5.27)	73.40 (+8.13)	84.89 (+6.85)	80.70 (+8.59)	81.37 (+6.97)
Qwen1.5-7B-LoRA	91.75	75.03	78.14	87.59	81.36	82.77
Qwen1.5-7B-Flexora	91.96 (+0.21)	84.33 (+9.30)	80.69 (+2.55)	89.90 (+2.31)	87.08 (+5.72)	86.79 (+4.02)
Yi-6B-LoRA	89.46	78.29	76.01	80.02	85.13	81.78
Yi-6B-Flexora	92.24 (+2.78)	84.82 (+6.53)	84.96 (+8.95)	88.72 (+8.70)	86.91 (+1.78)	87.53 (+5.75)
Qwen2.5-32B-LoRA	90.35	78.94	84.59	84.17	81.03	83.82
Qwen2.5-32B-Flexora	94.01 (+3.66)	86.72 (+7.78)	87.36 (+2.77)	87.16 (+2.99)	84.27 (+3.24)	87.90 (+4.08)



Table 11: Detailed comparison of the accuracy of the combination of Flexora and different LoRA algorithms on Llama3-8B. This table presents a detailed comparison of the accuracy results obtained by integrating Flexora with various improved LoRA algorithms, including DoRA, rsLoRA, MoSLoRA and LoReFT, while maintaining other experimental settings constant. The accuracy metrics are reported across multiple benchmark datasets, including HellaSwag, PIQA, Winogrande, RACE-mid, and RACE-high, with the average accuracy across all datasets also provided. The results are compared against those obtained from direct fine-tuning without Flexora. The experimental findings indicate that the application of Flexora can significantly reduce model overfitting and enhance overall performance.

Methods	Hellaswag	PIQA	Winogrande	RACE-mid	RACE-high	Average
LoRA	89.72	76.39	82.24	85.86	80.99	83.04
Flexora (w/ LoRA)	<b>93.62</b>	<b>85.91</b>	<b>85.79</b>	<b>84.61</b>	<b>82.36</b>	<b>86.46</b>
rsLoRA	94.33	87.21	85.32	87.60	84.36	87.76
Flexora (w/ rsLoRA)	<b>94.83</b>	<b>87.58</b>	<b>86.69</b>	<b>88.21</b>	<b>85.46</b>	<b>88.55</b>
DoRA	93.62	85.75	84.77	86.77	83.39	86.86
Flexora (w/ DoRA)	<b>94.10</b>	<b>86.05</b>	<b>86.32</b>	<b>87.12</b>	<b>84.45</b>	<b>87.61</b>
LoReFT	96.31	90.24	<b>87.48</b>	88.21	<b>85.33</b>	89.51
Flexora (w/ LoReFT)	<b>96.47</b>	<b>91.06</b>	87.23	<b>88.36</b>	84.97	<b>89.62</b>
MoSLoRA	93.53	85.97	84.26	<b>86.13</b>	<b>83.75</b>	<b>86.73</b>
Flexora (w/ MoSLoRA)	<b>93.76</b>	<b>86.43</b>	<b>85.36</b>	85.09	82.07	86.54

Table 12: Detailed comparison of commonsense reasoning task accuracy. This table provides a comprehensive comparison of the accuracy results for various methods applied to common sense reasoning tasks, conducted on the Llama-7B model. The methods compared include the pre-trained model, LoRA, LoRAShear with different pruning ratios (0.5), and Flexora. The accuracy metrics are reported across multiple benchmark datasets, including BoolQ, PIQA, HellaSwag, Winogrande, ARC-e, ARC-c, and OBQA. The average accuracy across all datasets is also provided. The "Ratio" column represents the ratio of parameter pruning in LoRAShear.

Methods	BoolQ	PIQA	HellaSwag	WinoGrande	ARC-e	ARC-c	OBQA	Average
Pre-trained	57.98	60.94	34.35	52.25	31.82	27.30	35.80	42.92
LoRA	67.76	69.80	76.10	67.01	67.21	35.23	38.60	60.24
LoRAShear (Ratio = 0.5)	63.40	<b>72.15</b>	49.83	56.40	49.45	34.31	35.86	51.63
<b>Flexora</b>	<b>73.54</b>	71.93	<b>85.28</b>	<b>74.11</b>	<b>71.22</b>	<b>45.64</b>	<b>39.86</b>	<b>65.94</b>

layer selection stage, while the outer layer optimization remains unchanged. These adjustments can be achieved through direct modification. The results demonstrate that Flexora exhibits strong scalability when combined with algorithms for enhancing LoRA parameters, highlighting its great potential.

#### C.4 Comparison with LoRAShear

**Better Performance of Flexora.** In this section, the accuracy of Flexora is compared with that of LoRAShear across various datasets, with specific results presented in Table 12. Since LoRAShear is not open source and poses challenges for direct experimentation, the comparison relies on the experimental configurations and results reported in the LoRAShear paper. Notably, Flexora can freely adjust the selected layers according to the dataset, achieving an average pruning parameter rate of 50%. Consequently, under the same pruning rate, Flexora outperforms by 14% (Ratio = 0.5). Experiments have shown that under the same pruning rate, Flexora can achieve better performance. Note that during the test, BoolQ only has training and test sets. We still keep the test set unchanged for testing the model, use 80% of the training set data to train LoRA parameters, and use the other 20% of the data to train the hyperparameters introduced by Flexora. In addition, the reason for the poor performance on the ARC and OBQA datasets is that the number of validation sets is small, and the layer selection may not be accurate enough. For a discussion on the number of validation sets and the accuracy, see section C.7.

#### C.5 Flexora in full-parameter fine-tuning

**Flexora can improve the performance of full parameter fine-tuning.** In this section, we evaluate the performance of Flexora in the context of full-parameter fine-tuning. Specifically, while maintaining

Table 13: Performance comparison of Full Fine-Tuning (Full FT), LoRA, and their Flexora-enhanced variants across multiple reasoning and reading comprehension tasks, including Hellaswag, PIQA, Winogrande, RACE-mid, and RACE-high. Results are reported as accuracy, with improvements over baseline methods indicated in parentheses. Flexora significantly enhances both Full FT and LoRA, with the largest gains observed for LoRA, achieving an average accuracy of 86.46% (+3.42% over LoRA).

Method	Hellaswag	PIQA	Winogrande	RACE-mid	RACE-high	Average
Full FT	90.53	79.32	81.16	81.92	79.36	82.46
Flexora(w/ Full FT)	91.32(+0.79)	83.21(+3.89)	81.73(+0.57)	83.13(+1.21)	80.37(+1.01)	83.95(+1.49)
LoRA	89.72	76.39	82.24	82.86	80.99	83.04
Flexora(w/ LoRA)	<b>93.62 (+3.90)</b>	<b>85.91 (+9.52)</b>	<b>85.79 (+3.55)</b>	<b>84.61 (+1.75)</b>	<b>82.36 (+1.37)</b>	<b>86.46 (+3.42)</b>

Table 14: Comparison of Flexora with base models and instruction-tuned models on various datasets. Numbers in parentheses indicate the relative improvement over the corresponding base model or instruction-tuned model. Flexora demonstrates consistent improvements across all datasets and settings.

Method	Hellaswag	PIQA	Winogrande	RACE-mid	RACE-high	Average
Base Model	48.55	67.08	59.91	67.02	63.35	61.18
Flexora(w/ Base Model)	<b>93.62 (+45.07)</b>	<b>85.91 (+18.83)</b>	<b>85.79 (+25.88)</b>	<b>84.61 (+17.59)</b>	<b>82.36 (+19.01)</b>	<b>86.46 (+25.28)</b>
Instruct Model	89.38	80.36	81.35	81.36	79.48	82.39
Flexora(w/ Instruct Model)	93.53 (+4.15)	85.76 (+5.4)	85.67 (+4.32)	84.58 (+3.22)	82.19 (+2.71)	86.35 (+3.96)

the inner loop optimization steps unchanged, we modify the trainable parameters in the outer loop from the LoRA adapter to the full set of parameters of the LLM itself. The experimental results, presented in Table 13, demonstrate that although Flexora enhances performance in full-parameter fine-tuning compared to baseline methods, it still falls short of the performance achieved by LoRA + Flexora. This suggests that full-parameter fine-tuning, which involves adjusting all layers of the model, is more susceptible to overfitting, even when employing the layer selection mechanism of Flexora. These findings underscore the significance of parameter-efficient approaches like LoRA, particularly in scenarios where overfitting is a critical concern, such as in customization or personalization tasks with limited disk storage. The capability of Flexora to identify and prioritize critical layers proves especially advantageous in these contexts, offering a balanced trade-off between model adaptability and resource efficiency.

## C.6 Flexora in instruct model

**Effectiveness of Flexora on instruct models.** To demonstrate the effectiveness of Flexora on the instruct model, we conducted experiments on Meta-Llama-3-8B-Instruct. The experimental results are shown in Table 14, which show that Flexora always maintains excellent performance on both the base model and the instruct model. This confirms that Flexora is effective in different fine-tuning scenarios (including instruct model adaptation).

## C.7 Different search sample

**Flexibility of Flexora in search sample .** In Flexora, search time is managed by adjusting the maximum number of search samples (corresponding to the size of the validation dataset) to align with the requirements of the downstream task. In Table 15, we explore the relationship between different numbers of search samples, downstream task performance, and search time. For simpler datasets like Hellaswag and PIQA, a 10-minute search with 1,000 samples significantly improves performance. For more challenging tasks, at least 1 hour of search time is required for 5,000 samples. In more difficult tasks, using too few samples can prevent validation loss from converging. To optimize performance, it is recommended to dynamically adjust the number of search samples based on the convergence of the validation loss. In summary, for simpler downstream tasks, Flexora can be rapidly applied to reduce model overfitting significantly and enhance performance. For more challenging downstream tasks, Flexora balances performance and training resources by adjusting the number of search samples.

Table 15: Detailed analysis of the impact of different numbers of search samples on the Flexora accuracy of Llama3-8B. This table investigates how varying the number of search samples, i.e., different validation dataset sizes, affects the performance of Flexora. The accuracy metrics are reported across multiple benchmark datasets, including HellaSwag, PIQA, Winogrande, RACE-mid, and RACE-high, with the average accuracy across all datasets also provided. The number of search samples tested includes 1000, 2000, 5000, 10000, and 200000. All experimental conditions remain unchanged except for the size of the validation set, allowing for a focused analysis on the impact of search sample size on model performance.

# Samples	Hellaswag	PIQA	Winogrande	RACE-mid	RACE-high	Average Time(h)
1000	93.00	80.52	83.04	76.74	72.93	0.08
1267	-	-	<b>85.79</b>	-	-	0.12
2000	92.29	<b>85.91</b>	-	80.15	78.82	0.17
4887	-	-	-	<b>84.82</b>	<b>82.36</b>	0.4
5000	93.17	-	-	-	-	0.42
10000	<b>93.62</b>	-	-	-	-	0.8

Table 16: Performance comparison of different hyperparameter settings  $K$  and  $T$  on various datasets. The rows represent different combinations of hyperparameters  $K$  and  $T$ . The columns represent the accuracy results on different datasets: HellaSwag, PIQA, Winogrande, RACE-mid, and RACE-high. The last column shows the average accuracy across all datasets.

# $K$ and $T$	Hellaswag	PIQA	Winogrande	RACE-mid	RACE-high	Average
$K = 4, T = 1$	93.57	85.76	85.72	84.62	<b>82.46</b>	86.43
$K = 8, T = 1$	93.62	85.91	<b>85.79</b>	84.61	82.36	86.46
$K = 4, T = 2$	93.78	85.37	84.16	83.96	83.17	86.09
$K = 8, T = 2$	93.07	85.16	85.01	84.57	82.06	85.97
$K = 4, T = 4$	92.97	85.72	85.56	<b>85.07</b>	82.11	86.29
$K = 8, T = 4$	<b>93.89</b>	<b>86.01</b>	<b>85.79</b>	84.99	<b>82.46</b>	<b>86.63</b>

## C.8 Ablation experiments on training settings

**The choice of  $K$  and  $T$  is not important** The  $K$  and  $T$  parameters of Flexora are inherited from the UD algorithm (Bao et al., 2021). As discussed in Section 2, some applications of the UD algorithm are highly sensitive to the choice of  $K$  and  $T$  (Liu et al., 2019). Therefore, we conducted ablation experiments to determine the optimal values for  $K$  and  $T$ . The specific experimental results are presented in Table 16. The results demonstrate that Flexora is highly robust to variations in  $K$  and  $T$ . This robustness may be attributed to the significant variability in the contribution of LLM layers to downstream tasks. For a detailed discussion, see Section C.10. Regardless of the settings for  $K$  and  $T$ , Flexora consistently identifies the layers that contribute the most to downstream tasks.

## C.9 More results for preliminary study

This section provides additional experimental results that are not shown in Section 3.1. In these experiments, we kept the randomly selected layers unchanged and only varied the LoRA rank. The specific experimental results are shown in Table 17. The results indicate that regardless of the selected rank, the model’s performance improves with an increasing number of LoRA fine-tuned layers up to a certain threshold. Beyond this threshold, further increasing the number of fine-tuned layers may lead to a decline in model performance. This intriguing phenomenon motivates our research.

## C.10 Selection of layers

For different LLMs and datasets, the layers chosen by Flexora vary due to the different parameters learned in the pre-training stage and the diversity of downstream tasks. In Table 18, Table 19, Table 20, Table 21, and Table 22, we show the layers chosen by Flexora in all experiments and the corresponding training parameters. In this section, the preferences of the layers chosen by Flexora are analyzed in detail, providing layer-wise insights for LLMs.

**The Effectiveness of Flexora Comes from Reducing Overfitting.** In Table 18, the layers and parameter amounts selected by different LoRA methods are presented. A comparison between LoRA-drop and Flexora reveals that Flexora is more effective. LoRA-drop tends to select the later layers, as these outputs exhibit a larger two-norm, aligning with Proposition 2. This result suggests that layers selected during fine-tuning should not concentrate in a specific range but rather be distributed across various ranges, fully utilizing the extensive knowledge system of LLMs. Comparing LoRA with DoRA and rsLoRA shows that LoRA selects more layers, requiring more training parameters but yielding worse performance. This suggests a higher degree of overfitting when Flexora is applied to LoRA compared to the other two methods. Therefore, using more advanced LoRA improvement algorithms can significantly reduce overfitting and enhance performance, underscoring the importance of the fine-tuning approach. Interestingly, certain layers are consistently fine-tuned in the same downstream task, regardless of whether LoRA, DoRA, or rsLoRA is used. For example, in Hellaswag, layers [0, 1, 2, 4, 14, 15, 19, 20, 21, 23, 26, 27, 28, 29, 31] are consistently selected, suggesting these layers are crucial for this task or represent general knowledge layers (see the next two paragraphs for details), closely related to the LLM itself.

**General Knowledge Layers.** In Table 19, the layers and parameters selected in the second ablation study are shown. Observing the "Select first 6 layers by Flexora" row reveals that certain layers, such as [27, 28], are crucial for any downstream task. These layers may store general knowledge, suggesting that their fine-tuning could enhance the performance across most downstream tasks.

**Downstream task-specific layers.** Table 20 displays the layers and parameter amounts selected by various LLMs for different downstream tasks. As evident from the table, the same model utilizes the aforementioned general knowledge layers across different tasks. Additionally, unique layers for each downstream task, termed downstream task-specific layers, are predominantly found in the first and last layers. The distinction between general knowledge layers and downstream task-specific layers can be attributed to the self-attention mechanism, which effectively differentiates these layers. In the self-attention mechanism, similar knowledge is aggregated, leading to this layer differentiation. Furthermore, concerning downstream task-specific layers, two conclusions are drawn: (a) Fewer layers are selected for simpler datasets to minimize overfitting. (b) Typically, the initial and final layers are selected for a given dataset. This selection pattern may stem from the initial layer processing the original input and the final layer generating the model’s output representation. Given the consistent and predefined input and output, learning these parameters is deemed effective.

**Poor Effects with No Critical Layers** Tables 21 and 22 serve as evidence for the existence of downstream task-specific and general knowledge layers. Failure to select these layers, due to reasons like random selection or lack of convergence, leads to poor performance.

Table 17: Performance of the model on various datasets (Hellaswag, PIQA, Winogrande, RACE-mid, RACE-high) under different LoRA ranks and varying numbers of LoRA fine-tuned layers.

Rank	Layers	Hellaswag	PIQA	Winogrande	RACE-mid	RACE-high	Average
$r = 4$	6 layers	58.36	68.23	45.71	53.35	52.99	55.73
	12 layers	78.23	76.53	54.78	79.04	54.99	68.71
	18 layers	<b>89.01</b>	<b>80.57</b>	82.79	82.37	<b>80.96</b>	<b>83.14</b>
	24 layers	88.21	79.36	<b>82.97</b>	<b>82.39</b>	80.12	82.61
	32 layers	87.68	74.36	81.74	81.10	79.63	80.90
$r = 8$	6 layers	59.79	70.25	46.32	54.54	53.45	56.87
	12 layers	81.9	77.82	57.35	78.41	72.16	73.53
	18 layers	<b>91.15</b>	<b>81.54</b>	<b>83.58</b>	<b>83.77</b>	<b>81.22</b>	<b>84.25</b>
	24 layers	90.58	80.9	82.16	82.19	79.22	83.01
	32 layers	89.72	76.39	82.24	82.86	80.99	82.44
$r = 16$	6 layers	60.98	71.36	47.12	55.78	54.26	57.90
	12 layers	80.23	78.01	62.69	79.55	75.62	75.22
	18 layers	<b>91.63</b>	<b>81.69</b>	<b>85.06</b>	<b>84.27</b>	<b>83.69</b>	<b>85.27</b>
	24 layers	90.11	79.60	83.57	82.13	78.39	82.76
	32 layers	89.99	78.47	82.77	81.63	79.68	82.51
$r = 32$	6 layers	60.45	71.46	50.36	57.36	55.13	58.95
	12 layers	82.4	79.07	63.17	80.13	78.63	76.68
	18 layers	<b>92.08</b>	<b>82.14</b>	<b>86.07</b>	85.35	83.04	<b>85.74</b>
	24 layers	91.55	81.37	85.13	<b>85.75</b>	<b>83.17</b>	82.76
	32 layers	90.01	79.56	84.36	82.36	80.99	83.46



In summary, it is evident that almost all LLMs feature downstream task-specific layers and general knowledge layers. Fine-tuning these layers effectively mitigates model overfitting and enhances both generalization and performance. Fortunately, Flexora accurately and efficiently identifies both the downstream task-specific layers and the general knowledge layers.

Table 18: Comprehensive overview of layer selection strategies in main experiments. This table presents a detailed breakdown of the layer selection strategies used in different experiments involving the Llama3-8B model and its variants (Flexora, LoRA-drop, DoRA + Flexora, and rsLoRA + Flexora). For each model, the specific datasets utilized (HellaSwag, PIQA, RACE, and Winogrande) are listed along with the corresponding layers selected for each dataset. The ‘‘Layer selection’’ column provides the indices of the layers chosen for each experiment, indicating the specific layers of the model that were fine-tuned or modified. Additionally, the ‘‘Parameter(M)’’ column indicates the total number of parameters (in millions) used in each configuration. This detailed breakdown allows for a clear understanding of the experimental setup, the layer selection process, and the parameter allocation across different models and datasets, facilitating a deeper analysis of the impact of these strategies on model performance. Unless otherwise specified, the results are based on the default LoRA Rank of 8.

Methods	Dataset	Layer selection	Parameter(M)
Llama3-8B + Flexora( $r = 8$ )	Hellaswag	[0, 1, 2, 3, 4, 5, 6, 14, 15, 19, 20, 21, 23, 24, 26, 27, 28, 29, 31]	2.0
	PIQA	[1, 2, 3, 4, 5, 7, 8, 9, 14, 20, 25, 26, 27, 28, 29, 30]	1.7
	RACE	[0, 1, 2, 3, 4, 7, 8, 9, 12, 14, 25, 26, 27, 28, 29, 31]	1.7
	Winogrande	[0, 1, 2, 3, 4, 16, 20, 22, 23, 24, 25, 26, 27, 28, 29, 31]	1.7
Llama3-8B + Flexora( $r = 16$ )	Hellaswag	[0, 1, 2, 3, 4, 5, 10, 14, 18, 19, 20, 21, 23, 24, 26, 27, 28, 29, 31]	2.0
	PIQA	[1, 2, 3, 4, 5, 7, 8, 10, 14, 20, 25, 26, 27, 28, 29, 30]	1.7
	RACE	[0, 1, 2, 3, 4, 7, 8, 9, 11, 14, 25, 26, 27, 28, 29, 31]	1.7
	Winogrande	[0, 1, 2, 3, 4, 18, 19, 22, 23, 24, 25, 26, 27, 28, 29, 31]	1.7
Llama3-8B + Flexora( $r = 32$ )	Hellaswag	[0, 1, 2, 3, 4, 5, 6, 11, 15, 18, 20, 21, 23, 24, 26, 27, 28, 29, 31]	2.0
	PIQA	[1, 2, 3, 4, 5, 7, 10, 12, 14, 20, 25, 26, 27, 28, 29, 30]	1.7
	RACE	[0, 1, 2, 3, 4, 7, 8, 10, 12, 14, 24, 26, 27, 28, 29, 31]	1.7
	Winogrande	[0, 1, 2, 3, 4, 16, 19, 20, 23, 24, 25, 26, 27, 28, 29, 31]	1.7
Llama3-8B + LoRA-drop	Hellaswag	[13, 14, 15, 16, 17, 18, 19, 20, 21, 22, 23, 24, 25, 26, 27, 28, 29, 30, 31]	2.0
	PIQA	[16, 17, 18, 19, 20, 21, 22, 23, 24, 25, 26, 27, 28, 29, 30, 31]	1.7
	RACE	[16, 17, 18, 19, 20, 21, 22, 23, 24, 25, 26, 27, 28, 29, 30, 31]	1.7
	Winogrande	[16, 17, 18, 19, 20, 21, 22, 23, 24, 25, 26, 27, 28, 29, 30, 31]	1.7
Llama3-8B + DoRA + Flexora	Hellaswag	[0, 1, 2, 4, 5, 14, 15, 19, 20, 21, 23, 26, 27, 28, 29, 31]	1.8
	PIQA	[0, 1, 2, 4, 7, 23, 24, 25, 26, 27, 28, 29, 31]	1.5
	RACE	[1, 3, 4, 7, 9, 12, 14, 23, 25, 27, 28, 29, 31]	1.3
	Winogrande	[0, 1, 2, 3, 20, 21, 22, 23, 24, 25, 26, 27, 28, 29, 31]	1.7
Llama3-8B + rsLoRA + Flexora	Hellaswag	[0, 1, 2, 4, 6, 14, 15, 19, 20, 21, 23, 25, 26, 27, 28, 29, 31]	1.8
	PIQA	[0, 1, 2, 3, 15, 20, 21, 25, 26, 27, 28, 29, 31]	1.3
	RACE	[0, 1, 2, 3, 7, 8, 12, 13, 25, 26, 27, 28, 29, 31]	1.5
	Winogrande	[1, 2, 3, 6, 14, 15, 18, 20, 21, 22, 23, 24, 25, 26, 27, 28, 29, 31]	1.9
Llama3-8B + LoReFT + Flexora	Hellaswag	[0, 1, 2, 3, 4, 5, 6, 19, 20, 21, 23, 26, 27, 28, 29, 31]	1.8
	PIQA	[0, 1, 2, 4, 7, 22, 24, 25, 26, 27, 28, 29, 30, 31]	1.5
	RACE	[0, 1, 3, 4, 7, 9, 14, 23, 25, 27, 28, 29, 31]	1.3
	Winogrande	[0, 1, 2, 3, 4, 5, 22, 23, 24, 25, 26, 27, 28, 29, 31]	1.7
Llama3-8B + MoSLoRA + Flexora	Hellaswag	[0, 1, 2, 4, 5, 14, 16, 19, 20, 21, 23, 25, 26, 27, 28, 29, 31]	1.8
	PIQA	[0, 1, 2, 3, 4, 20, 21, 25, 26, 27, 28, 29, 31]	1.3
	RACE	[0, 1, 2, 3, 7, 9, 11, 13, 25, 26, 27, 28, 29, 31]	1.5
	Winogrande	[1, 2, 3, 8, 10, 15, 18, 20, 21, 22, 23, 24, 25, 26, 27, 28, 29, 31]	1.9

## D Loss

This section presents the training, evaluation, and validation loss during the Flexora flexible layer selection and fine-tuning stages, accompanied by intuitive explanations.

### D.1 Effectiveness of Flexora.

Figure 5 plots the training and validation loss curves for Llama-8B during the flexible layer selection stage across four different datasets over one epoch. Both inner and outer layer optimizations are observed to converge well during the flexible layer selection stage, demonstrating the effectiveness of Flexora.

### D.2 Flexora can Correctly Identify Critical Layers.

Figures 6, 7, 8, and 9 depict the training and evaluation loss from the first ablation study. In all experiments, the training loss converges effectively, demonstrating robust training performance. However, variations in evaluation loss underscore the model’s generalization capabilities. Flexora generally surpasses methods that randomly select an equivalent number of layers, demonstrating its ability to accurately identify critical layers for more effective improvements.

Table 19: Detailed display of selected layers in the second ablation study. In the second ablation experiment, we manually determined the number of fine-tuning layers and contrasted the performance of Flexora with random layer selection strategies. This table presents the results of this experiment, showcasing different configurations where a specific number of layers (6, 12, 18, and 24) were selected for fine-tuning. For each configuration, the table compares the layers selected by Flexora with those selected randomly. The datasets used in this experiment include HellaSwag, PIQA, RACE, and Winogrande. The “Layer selection” column lists the indices of the layers chosen for fine-tuning in each dataset, while the “Parameter(M)” column indicates the total number of parameters (in millions) used in each configuration. This detailed breakdown provides insights into how different layer selection strategies, with a manually determined number of fine-tuning layers, impact the performance of model across different datasets, facilitating a comprehensive comparison between Flexora and random selection methods.

Methods	Dataset	Layer selection	Parameter(M)
Select first 6 layers by Flexora	Hellaswag	[0, 26, 27, 28, 29, 31]	0.6
	PIQA	[2, 4, 26, 27, 28, 29]	0.6
	RACE	[0, 7, 12, 27, 28, 29]	0.6
	Winogrande	[22, 23, 24, 26, 27, 28]	0.6
Random selection 6 layers	Hellaswag	[2, 4, 11, 19, 23, 25]	0.6
	PIQA	[2, 4, 11, 19, 23, 25]	0.6
	RACE	[2, 4, 11, 19, 23, 25]	0.6
	Winogrande	[2, 4, 11, 19, 23, 25]	0.6
Select first 12 layers by Flexora	Hellaswag	[0, 2, 3, 14, 15, 21, 23, 26, 27, 28, 29, 31]	1.3
	PIQA	[1, 2, 3, 4, 7, 20, 25, 26, 27, 28, 29, 30]	1.3
	RACE	[0, 1, 3, 7, 8, 12, 13, 25, 27, 28, 29, 31]	1.3
	Winogrande	[0, 3, 20, 22, 23, 24, 25, 26, 27, 28, 29, 31]	1.3
Random selection 12 layers	Hellaswag	[1, 3, 4, 12, 14, 18, 20, 21, 22, 27, 29, 31]	1.3
	PIQA	[1, 3, 4, 12, 14, 18, 20, 21, 22, 27, 29, 31]	1.3
	RACE	[1, 3, 4, 12, 14, 18, 20, 21, 22, 27, 29, 31]	1.3
	Winogrande	[1, 3, 4, 12, 14, 18, 20, 21, 22, 27, 29, 31]	1.3
Select first 18 layers by Flexora	Hellaswag	[0, 1, 2, 3, 4, 5, 6, 14, 15, 19, 21, 23, 26, 27, 28, 29, 30, 31]	1.9
	PIQA	[0, 1, 2, 3, 4, 5, 7, 8, 19, 20, 23, 25, 26, 27, 28, 29, 30, 31]	1.9
	RACE	[0, 1, 2, 3, 4, 7, 8, 9, 10, 12, 13, 15, 25, 27, 28, 29, 30, 31]	1.9
	Winogrande	[0, 1, 3, 5, 7, 9, 15, 20, 21, 22, 23, 24, 25, 26, 27, 28, 29, 31]	1.9
Random selection 18 layers	Hellaswag	[1, 2, 5, 8, 9, 10, 12, 13, 17, 18, 20, 21, 22, 23, 24, 25, 26, 30]	1.9
	PIQA	[1, 2, 5, 8, 9, 10, 12, 13, 17, 18, 20, 21, 22, 23, 24, 25, 26, 30]	1.9
	RACE	[1, 2, 5, 8, 9, 10, 12, 13, 17, 18, 20, 21, 22, 23, 24, 25, 26, 30]	1.9
	Winogrande	[1, 2, 5, 8, 9, 10, 12, 13, 17, 18, 20, 21, 22, 23, 24, 25, 26, 30]	1.9
Select first 24 layers by Flexora	Hellaswag	[0, 1, 2, 3, 4, 5, 6, 11, 12, 13, 14, 15, 18, 19, 20, 21, 23, 24, 26, 27, 28, 29, 30, 31]	2.6
	PIQA	[0, 1, 2, 3, 4, 5, 6, 7, 8, 9, 15, 18, 19, 20, 21, 23, 24, 25, 26, 27, 28, 29, 30, 31]	2.6
	RACE	[0, 1, 2, 3, 4, 5, 6, 7, 8, 9, 10, 11, 12, 13, 14, 15, 23, 24, 25, 27, 28, 29, 30, 31]	2.6
	Winogrande	[0, 1, 2, 3, 4, 5, 7, 8, 9, 10, 15, 19, 20, 21, 22, 23, 24, 25, 26, 27, 28, 29, 30, 31]	2.6
Random selection 24 layers	Hellaswag	[0, 1, 4, 6, 7, 8, 9, 10, 12, 13, 14, 15, 16, 17, 18, 19, 21, 23, 25, 26, 27, 28, 30, 31]	2.6
	PIQA	[0, 1, 4, 6, 7, 8, 9, 10, 12, 13, 14, 15, 16, 17, 18, 19, 21, 23, 25, 26, 27, 28, 30, 31]	2.6
	RACE	[0, 1, 4, 6, 7, 8, 9, 10, 12, 13, 14, 15, 16, 17, 18, 19, 21, 23, 25, 26, 27, 28, 30, 31]	2.6
	Winogrande	[0, 1, 4, 6, 7, 8, 9, 10, 12, 13, 14, 15, 16, 17, 18, 19, 21, 23, 25, 26, 27, 28, 30, 31]	2.6

### D.3 Flexora can Reduce Overfitting.

Figures 10, 11, 12, and 13 present the training and evaluation loss from the sencond ablation study. Consistent with previous experiments, the training loss converges, indicating a strong training effect on the training set. Notably, the 24-layer (red) model consistently shows the lowest training loss, suggesting optimal learning, whereas the 6-layer (blue) model consistently records the highest, indicating poorer training performance. However, differences in evaluation loss reveal variations in model generalization across different layers. The 18-layer (green) model consistently exhibits the lowest evaluation loss, indicating superior generalization and downstream task performance, corroborated by actual results. The 24-layer (red) model’s evaluation loss consistently exceeds that of the 18-layer (green) model, suggesting significant overfitting. Similarly, the 6-layer (blue) model consistently records the highest evaluation loss, indicative of underfitting.

In summary, too few training layers can lead to underfitting and poor performance, as seen in the 6-layer (blue) model. Conversely, too many layers can also result in overfitting, as evidenced by the 24-layer (red) model’s performance. However following the selection strategy of Flexora, choosing the right number of layers can minimize overfitting and improve performance

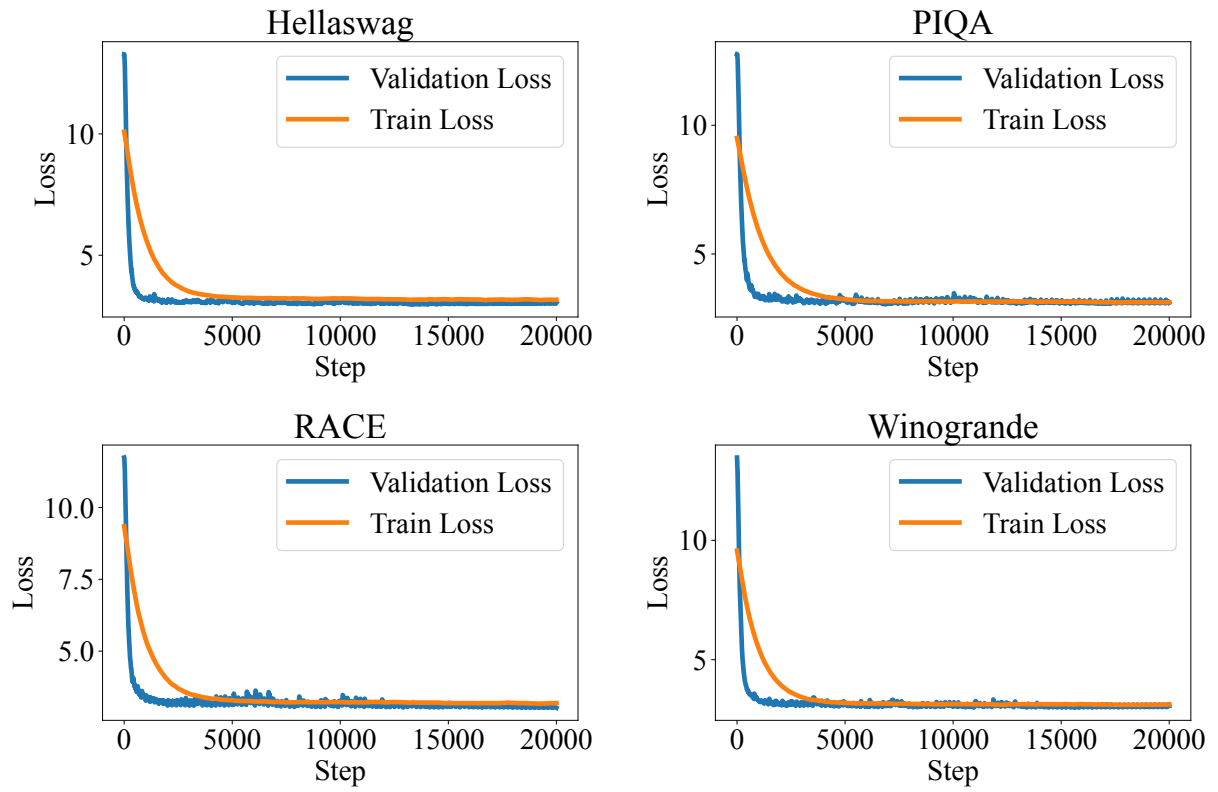


Figure 5: Training and validation loss during the flexible layer selection phase. The figure shows the training and validation loss over 20,000 steps for four different datasets (Hellaswag, PIQA, RACE, and Winogrande), where the batch size at each step is 1. The blue line shows the validation loss and the orange line shows the training loss. These plots visually compare how the performance of the models changes during the flexible layer selection phase, highlighting the convergence behavior.

Table 20: Comprehensive overview of layer selection strategies and parameter allocation in various experiments. This table provides an in-depth breakdown of the layer selection strategies employed across different models and datasets in the experiments. The models tested include Llama3-8B, Chatglm3-6B, Mistral-7B-v0.1 and others, all combined with Flexora. For each model, the specific datasets used (Hellaswag, PIQA, RACE, and Winogrande) are listed along with the corresponding layers selected for each dataset. The “Layer selection” column details the indices of the layers chosen for each experiment, indicating the specific layers of the model that were fine-tuned or modified. Additionally, the “Parameter(M)” column indicates the total number of parameters (in millions) used in each configuration. This detailed breakdown allows for a clear understanding of the experimental setup, the layer selection process, and the parameter allocation across different models and datasets, facilitating a deeper analysis of the impact of these strategies on model performance.

Methods	Dataset	Layer selection	Parameter(M)
Llama3-8B + Flexora	Hellaswag	[0, 1, 2, 3, 4, 5, 6, 14, 15, 19, 20, 21, 23, 24, 26, 27, 28, 29, 31]	2.0
	PIQA	[1, 2, 3, 4, 5, 7, 8, 9, 14, 20, 25, 26, 27, 28, 29, 30]	1.7
	RACE	[0, 1, 2, 3, 4, 7, 8, 9, 12, 14, 25, 26, 27, 28, 29, 31]	1.7
	Winogrande	[0, 1, 2, 3, 4, 16, 20, 22, 23, 24, 25, 26, 27, 28, 29, 31]	1.7
Chatglm3-6B + Flexora	Hellaswag	[1, 2, 3, 4, 5, 6, 7, 10, 12, 13, 16, 18, 20]	0.9
	PIQA	[0, 1, 2, 3, 5, 6, 7, 8, 9, 19, 21, 23, 25, 27]	1.0
	RACE	[2, 6, 8, 9, 10, 11, 14, 15, 16, 17, 18, 20, 23, 26]	1.0
	Winogrande	[0, 2, 6, 8, 9, 11, 12, 13, 16, 17, 18, 20, 25, 26]	1.0
Mistral-7B-v0.1 + Flexora	Hellaswag	[0, 1, 2, 3, 4, 5, 6, 7, 14, 22, 26, 27, 30]	1.5
	PIQA	[6, 8, 14, 17, 18, 22, 23, 24, 25, 26, 27, 28, 29, 30]	1.7
	RACE	[0, 2, 3, 4, 5, 6, 7, 8, 9, 10, 11, 13, 14, 17, 30, 31]	1.7
	Winogrande	[0, 1, 2, 3, 4, 5, 6, 7, 22, 23, 24, 25, 26, 27, 28, 29, 30, 31]	1.9
Gemma-7B + Flexora	Hellaswag	[0, 1, 2, 3, 4, 5, 6, 7, 8, 9, 14, 15, 16, 18, 20, 23, 27]	1.9
	PIQA	[0, 1, 8, 9, 10, 12, 15, 16, 17, 20, 21, 22, 23, 24, 25, 26, 27]	1.9
	RACE	[0, 1, 2, 3, 4, 5, 6, 7, 8, 9, 15, 16]	1.4
	Winogrande	[0, 1, 2, 3, 4, 5, 6, 7, 8, 18, 19, 20, 21, 22, 23, 24, 25, 26, 27]	2.1
Vicuna-7B-v1.5 + Flexora	Hellaswag	[0, 1, 2, 3, 4, 6, 7, 8, 9, 10, 11, 12]	1.6
	PIQA	[1, 2, 3, 5, 7, 8, 11, 12, 13, 14, 21, 31]	1.6
	RACE	[2, 3, 4, 5, 6, 7, 8, 9, 10, 11, 12, 13]	1.6
	Winogrande	[0, 2, 3, 4, 6, 8, 9, 12, 20, 21, 22, 23, 24, 25, 26, 27, 28, 29, 30, 31]	2.6
Zephyr-7B-beta + Flexora	Hellaswag	[1, 13, 15, 17, 18, 22, 23, 24, 25, 26, 27, 28, 30, 31]	1.5
	PIQA	[2, 3, 6, 7, 14, 15, 16, 17, 22, 26, 27, 28]	1.4
	RACE	[1, 2, 4, 6, 7, 9, 11, 13, 14, 17, 26, 30, 31]	1.4
	Winogrande	[1, 3, 5, 6, 8, 13, 27, 28, 29, 30, 31]	1.2
Yi-6B + Flexora	Hellaswag	[0, 1, 2, 3, 4, 6, 8, 9, 10, 19, 20, 21, 22]	1.3
	PIQA	[1, 2, 3, 5, 6, 7, 8, 9, 12, 13, 15, 16, 17, 18, 20, 23]	1.6
	RACE	[1, 3, 5, 6, 7, 9, 11, 12, 13, 14, 17, 21]	1.2
	Winogrande	[0, 1, 2, 3, 5, 6, 7, 11, 23, 26, 27, 30, 31]	1.3
Llama-7B + Flexora	Hellaswag	[0, 1, 2, 4, 5, 6, 8, 12, 16, 30, 31]	1.4
	PIQA	[2, 12, 14, 15, 16, 21, 22, 23, 24, 25, 26, 27, 28, 29, 30, 31]	2.1
	RACE	[4, 5, 6, 7, 8, 10, 11, 23, 30, 31]	1.3
	Winogrande	[0, 2, 3, 6, 7, 8, 10, 11, 13, 16, 23, 28, 29, 30, 31]	2.0
Llama2-7B + Flexora	Hellaswag	[0, 1, 2, 3, 4, 5, 6, 7, 8, 12]	1.3
	PIQA	[0, 1, 2, 3, 7, 8, 11, 13, 14, 21, 24, 29, 30, 31]	1.8
	RACE	[0, 1, 2, 4, 5, 6, 7, 8, 9, 10, 11, 12, 14, 16]	1.8
	Winogrande	[0, 1, 3, 4, 8, 14, 15, 16, 17, 20, 21, 22, 23, 24, 25, 26, 28, 29, 30]	2.5
XuanYuan-6B + Flexora	Hellaswag	[1, 2, 3, 4, 5, 6, 7, 8, 9, 12, 13, 14, 17]	1.7
	PIQA	[3, 4, 7, 8, 12, 14, 16, 17, 19, 21, 23, 25, 28, 29]	1.8
	RACE	[0, 4, 5, 6, 7, 8, 9, 10, 11, 12, 14, 16, 17, 20, 21, 22, 25, 28, 29]	2.5
	Winogrande	[2, 3, 4, 8, 9, 10, 14, 15, 19, 20, 21, 22, 23, 24, 25, 26, 27, 28, 29]	2.5
Qwen1.5-7B + Flexora	Hellaswag	[0, 1, 2, 3, 4, 5, 6, 7, 9, 17]	1.3
	PIQA	[0, 1, 2, 3, 4, 5, 6, 7, 8, 11, 13, 14, 15, 17]	1.8
	RACE	[0, 1, 2, 3, 4, 5, 6, 7, 8, 9, 10, 11, 12]	1.7
	Winogrande	[0, 1, 2, 3, 4, 5, 6, 7, 8, 21, 24, 25, 27, 28, 30]	2.0

Table 21: Detailed display of selected layers in the first ablation study. This table presents the results of the first ablation experiment, where the number of layers selected by Flexora was kept constant, but different layers were chosen for fine-tuning. The table includes three different random layer selection strategies (Random1, Random2, and Random3) applied to various datasets (HellaSwag, PIQA, RACE, and Winogrande). For each random selection method, the “Layer selection” column lists the indices of the layers chosen for fine-tuning in each dataset. The “Parameter(M)” column indicates the total number of parameters (in millions) used in each configuration. This detailed breakdown allows for a clear understanding of how different layer selection strategies impact the performance of model across different datasets while maintaining a consistent number of layers for fine-tuning.

Methods	Dataset	Layer selection	Parameter(M)
Random1	Hellaswag	[0, 1, 2, 3, 4, 6, 7, 8, 10, 11, 14, 18, 19, 20, 21, 25, 26, 27, 28]	2.0
	PIQA	[0, 2, 4, 10, 12, 16, 17, 18, 23, 24, 25, 26, 27, 28, 29, 30]	1.7
	RACE	[1, 2, 4, 7, 9, 11, 12, 14, 15, 18, 20, 23, 24, 26, 28, 30]	1.7
	Winogrande	[1, 2, 4, 5, 9, 10, 11, 13, 15, 17, 20, 21, 24, 26, 30, 31]	1.7
Random2	Hellaswag	[0, 2, 3, 4, 5, 6, 10, 12, 13, 15, 17, 20, 21, 22, 23, 24, 28, 29, 30]	2.0
	PIQA	[0, 1, 3, 4, 8, 13, 14, 18, 19, 22, 24, 26, 28, 29, 30, 31]	1.7
	RACE	[5, 6, 7, 8, 9, 11, 12, 13, 15, 19, 20, 21, 25, 27, 28, 30]	1.7
	Winogrande	[2, 5, 6, 7, 8, 10, 11, 13, 14, 17, 18, 22, 25, 26, 28, 30]	1.7
Random3	Hellaswag	[0, 1, 3, 4, 6, 9, 12, 13, 17, 18, 19, 20, 21, 22, 25, 26, 27, 28, 29]	2.0
	PIQA	[0, 3, 4, 9, 12, 13, 14, 15, 16, 24, 25, 26, 27, 28, 30, 31]	1.7
	RACE	[0, 1, 2, 9, 11, 12, 14, 18, 19, 20, 21, 23, 25, 26, 29, 30]	1.7
	Winogrande	[2, 4, 6, 8, 10, 12, 14, 16, 17, 18, 20, 22, 23, 29, 30, 31]	1.7

Table 22: Detailed display of layer selection with varying numbers of searching samples. This table presents the results of an experiment where different numbers of searching samples (1000, 2000, 5000, and 10000) were used to determine the layers for Flexora. The datasets involved in this experiment include HellaSwag, PIQA, RACE, and Winogrande. For each number of searching samples, the “Layer selection” column lists the indices of the layers chosen for fine-tuning in each dataset. The “Parameter(M)” column indicates the total number of parameters (in millions) used in each configuration. This detailed breakdown provides insights into how the number of searching samples impacts the layer selection process and the performance of model across different datasets.

Methods	Dataset	Layer selection	Parameter(M)
1000 searching samples	Hellaswag	[0, 2, 4, 5, 6, 8, 10, 16, 21, 26, 27, 28, 30, 31]	1.5
	PIQA	[0, 1, 2, 3, 4, 16, 25, 26, 27, 28, 29, 30, 31]	1.4
	RACE	[0, 1, 2, 3, 4, 16, 21, 28, 29, 30, 31]	1.2
	Winogrande	[0, 1, 2, 3, 4, 16, 20, 25, 26, 27, 28, 29, 30, 31]	1.5
2000 searching samples	Hellaswag	[1, 2, 3, 4, 8, 10, 11, 16, 30, 31]	1.0
	PIQA	[0, 1, 2, 20, 22, 23, 24, 25, 26, 27, 28, 29, 30, 31]	1.5
	RACE	[0, 1, 2, 3, 4, 10, 20, 23, 27, 28, 29, 30, 31]	1.4
	Winogrande	[0, 1, 2, 3, 4, 20, 25, 27, 30, 31]	1.0
5000 searching samples	Hellaswag	[0, 1, 2, 3, 4, 8, 31]	0.7
	PIQA	[0, 2, 3, 4, 20, 22, 23, 24, 25, 26, 27, 28, 29, 30, 31]	1.6
	RACE	[1, 3, 4, 6, 9, 10, 11, 12, 14, 27, 28, 29, 30, 31]	1.5
	Winogrande	[1, 2, 3, 4, 6, 7, 8, 9, 26, 27, 30, 31]	1.3
10000 searching samples	Hellaswag	[0, 1, 4, 10, 12, 14, 21, 24, 26, 27, 28, 29, 30, 31]	1.5
	PIQA	[0, 1, 3, 4, 7, 20, 21, 22, 24, 25, 26, 27, 28, 29, 30, 31]	1.7
	RACE	[1, 2, 7, 13, 14, 23, 25, 26, 27, 28, 29, 31]	1.3
	Winogrande	[6, 7, 9, 10, 15, 19, 20, 22, 26, 27, 30, 31]	1.3



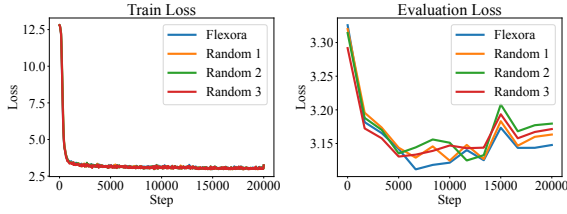


Figure 6: Comparison of train loss and evaluation loss in the Hellaswag dataset during the first ablation study. This figure presents the train loss (left) and evaluation loss (right) over 20,000 steps for the Hellaswag dataset, where the batch size at each step is 1. The performance of the Flexora method is compared against three different random layer selection strategies (Random 1, Random 2, and Random 3). The train loss graph shows how the training performance of model evolves, while the evaluation loss graph highlights the generalization capability of model on the validation set. This detailed comparison provides insights into the effectiveness of Flexora relative to random selection methods in terms of both training and evaluation metrics.

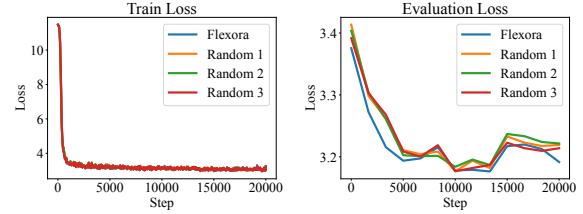


Figure 8: Comparison of train loss and evaluation loss in the RACE dataset during the first ablation study. This figure presents the train loss (left) and evaluation loss (right) over 20,000 steps for the RACE dataset, where the batch size at each step is 1. The performance of the Flexora method is compared against three different random layer selection strategies (Random 1, Random 2, and Random 3). The train loss graph shows how the training performance of model evolves, while the evaluation loss graph highlights the generalization capability of model on the validation set. This detailed comparison provides insights into the effectiveness of Flexora relative to random selection methods in terms of both training and evaluation metrics.

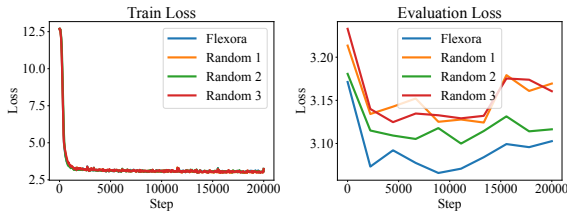


Figure 7: Comparison of train loss and evaluation loss in the PIQA dataset during the first ablation study. This figure presents the train loss (left) and evaluation loss (right) over 20,000 steps for the PIQA dataset, where the batch size at each step is 1. The performance of the Flexora method is compared against three different random layer selection strategies (Random 1, Random 2, and Random 3). The train loss graph shows how the training performance of model evolves, while the evaluation loss graph highlights the generalization capability of model on the validation set. This detailed comparison provides insights into the effectiveness of Flexora relative to random selection methods in terms of both training and evaluation metrics.

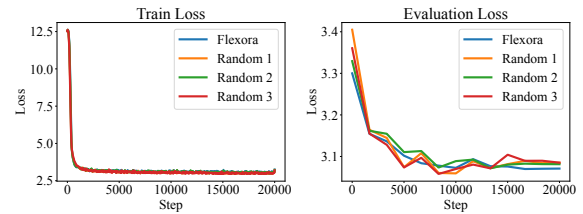


Figure 9: Comparison of train loss and evaluation loss in the Winogrande dataset during the first ablation study. This figure presents the train loss (left) and evaluation loss (right) over 20,000 steps for the Winogrande dataset, where the batch size at each step is 1. The performance of the Flexora method is compared against three different random layer selection strategies (Random 1, Random 2, and Random 3). The train loss graph shows how the training performance of model evolves, while the evaluation loss graph highlights the generalization capability of model on the validation set. This detailed comparison provides insights into the effectiveness of Flexora relative to random selection methods in terms of both training and evaluation metrics.

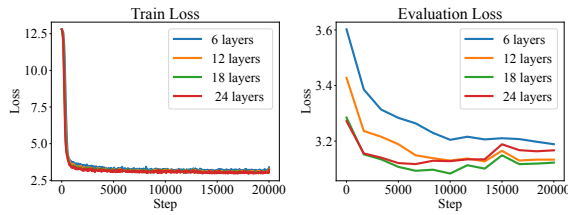


Figure 10: Training loss and evaluation loss during fine-tuning of different numbers of layers in the Flexora on the Hellaswag dataset. This figure presents the training loss (left) and evaluation loss (right) over 20,000 steps for the Hellaswag dataset. The performance is compared across four different configurations where the first 6, 12, 18, and 24 layers of the Flexora model are fine-tuned. The training loss graph shows that the model with 24 layers (red) achieves the lowest training loss, indicating it fits the training data very well. However, the evaluation loss graph reveals that the model with 18 layers (green) achieves the lowest evaluation loss, suggesting better generalization to unseen data. This discrepancy highlights the overfitting issue, where the model with 24 layers performs well on the training data but does not generalize as effectively as the model with 18 layers.

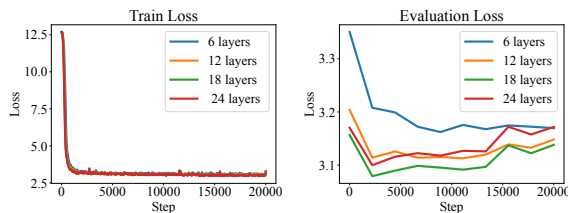


Figure 11: Training loss and evaluation loss during fine-tuning of different numbers of layers in the Flexora on the PIQA dataset. This figure presents the training loss (left) and evaluation loss (right) over 20,000 steps for the PIQA dataset. The performance is compared across four different configurations where the first 6, 12, 18, and 24 layers of the Flexora model are fine-tuned. The training loss graph shows that the model with 24 layers (red) achieves the lowest training loss, indicating it fits the training data very well. However, the evaluation loss graph reveals that the model with 18 layers (green) achieves the lowest evaluation loss, suggesting better generalization to unseen data. This discrepancy highlights the overfitting issue, where the model with 24 layers performs well on the training data but does not generalize as effectively as the model with 18 layers.

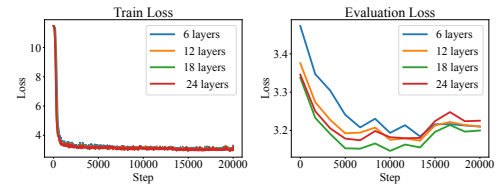


Figure 12: Training loss and evaluation loss during fine-tuning of different numbers of layers in the Flexora on the RACE dataset. This figure presents the training loss (left) and evaluation loss (right) over 20,000 steps for the RACE dataset. The performance is compared across four different configurations where the first 6, 12, 18, and 24 layers of the Flexora model are fine-tuned. The training loss graph shows that the model with 24 layers (red) achieves the lowest training loss, indicating it fits the training data very well. However, the evaluation loss graph reveals that the model with 18 layers (green) achieves the lowest evaluation loss, suggesting better generalization to unseen data. This discrepancy highlights the overfitting issue, where the model with 24 layers performs well on the training data but does not generalize as effectively as the model with 18 layers.

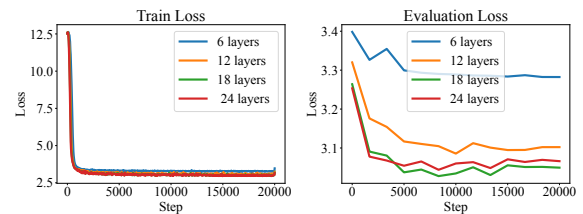


Figure 13: Training loss and evaluation loss during fine-tuning of different numbers of layers in the Flexora on the Winogrande dataset. This figure presents the training loss (left) and evaluation loss (right) over 20,000 steps for the Winogrande dataset. The performance is compared across four different configurations where the first 6, 12, 18, and 24 layers of the Flexora model are fine-tuned. The training loss graph shows that the model with 24 layers (red) achieves the lowest training loss, indicating it fits the training data very well. However, the evaluation loss graph reveals that the model with 18 layers (green) achieves the lowest evaluation loss, suggesting better generalization to unseen data. This discrepancy highlights the overfitting issue, where the model with 24 layers performs well on the training data but does not generalize as effectively as the model with 18 layers.

## E Theoretical insights and numerical experiments on the smoothness constant of Llama3-8B

### E.1 Theoretical insights

Consider a neural network with weight matrices decomposed as:

$$W^{(i)} = W_1^{(i)}(\text{frozen}) + W_2^{(i)}(\text{trainable}),$$

similar to LoRA's low-rank adaptation. We analyze the gradient difference for the trainable component  $W_2^{(i)}$ .

The network output is given by:

$$y^{(N)} = \left( \prod_{j=1}^N (W_1^{(j)} + W_2^{(j)}) \right) x,$$

Let  $\ell$  be the MSE loss:

$$\ell = \frac{1}{2} \left( y^{(N)} - y \right)^2,$$

where  $y^{(N)}$  is the network output. The gradient of  $\ell$  w.r.t.  $W_2^{(i)}$  is:

$$\frac{\partial \ell}{\partial W_2^{(i)}} = \left( \prod_{j=i+1}^N (W_1^{(j)} + W_2^{(j)}) \right) \cdot \frac{\partial \ell}{\partial y^{(i)}} \cdot x \left( \prod_{j=1}^{i-1} (W_1^{(j)} + W_2^{(j)}) \right). \quad (11)$$

For two networks differing only in  $W_2^{(i)}$ , the gradient difference is:

$$\left\| \frac{\partial \ell}{\partial W_{2,1}^{(i)}} - \frac{\partial \ell}{\partial W_{2,2}^{(i)}} \right\| = \left\| \left( \prod_{j=i+1}^N (W_1^{(j)} + W_2^{(j)}) \right) \left( \frac{\partial \ell}{\partial y_1^{(i)}} - \frac{\partial \ell}{\partial y_2^{(i)}} \right) x \left( \prod_{j=1}^{i-1} (W_1^{(j)} + W_2^{(j)}) \right) \right\|. \quad (12)$$

For MSE loss, the gradient at layer  $i$  is:

$$\frac{\partial \ell}{\partial y^{(i)}} = (y^{(N)} - y) \prod_{j=i+1}^N (W_1^{(j)} + W_2^{(j)}). \quad (13)$$

The difference between the two network outputs is:

$$y_1^{(N)} - y_2^{(N)} = \left( \left( \prod_{j=1}^i (W_{1,1}^{(j)} + W_{1,2}^{(j)}) \right) - \left( \prod_{j=1}^i (W_{2,1}^{(j)} + W_{2,2}^{(j)}) \right) \right) x \quad (14)$$

$$= \left( \prod_{j=1}^{i-1} (W_1^{(j)} + W_2^{(j)}) \right) (W_{1,1}^{(i)} + W_{2,1}^{(i)} - W_{1,2}^{(i)} - W_{2,2}^{(i)}) x \quad (15)$$

$$= \left( \prod_{j=1}^{i-1} (W_1^{(j)} + W_2^{(j)}) \right) (W_{2,1}^{(i)} - W_{2,2}^{(i)}) x. \quad (16)$$

The difference between the differentials of the two network loss functions with respect to the output is:

$$\frac{\partial \ell}{\partial y_1^{(i)}} - \frac{\partial \ell}{\partial y_2^{(i)}} = \left( y_1^{(N)} - y_2^{(N)} \right) \prod_{j=i+1}^N (W_1^{(j)} + W_2^{(j)}) \quad (17)$$

$$= \left( \prod_{j=1}^{i-1} (W_1^{(j)} + W_2^{(j)}) \right) (W_{2,1}^{(i)} - W_{2,2}^{(i)}) x \prod_{j=i+1}^N (W_1^{(j)} + W_2^{(j)}) \quad (18)$$

$$= \left( \prod_{\substack{j=1 \\ j \neq i}}^N (W_1^{(j)} + W_2^{(j)}) \right) (W_{2,1}^{(i)} - W_{2,2}^{(i)}) x. \quad (19)$$

Let  $\lambda^{(j)} = \|W_1^{(j)} + W_2^{(j)}\|$  be the spectral norm.

$$\left\| \frac{\partial \ell}{\partial W_{2,1}^{(i)}} - \frac{\partial \ell}{\partial W_{2,2}^{(i)}} \right\| = \left\| \left( \prod_{j=i+1}^N (W_1^{(j)} + W_2^{(j)}) \right) \left( \frac{\partial \ell}{\partial y_1^{(i)}} - \frac{\partial \ell}{\partial y_2^{(i)}} \right) x \left( \prod_{j=1}^{i-1} (W_1^{(j)} + W_2^{(j)}) \right) \right\| \quad (20)$$

$$= \left\| \left( \prod_{\substack{j=1 \\ j \neq i}}^N (W_1^{(j)} + W_2^{(j)}) \right)^2 (W_{2,1}^{(i)} - W_{2,2}^{(i)}) x^2 \right\| \quad (21)$$

$$\leq \left( \prod_{j=1, j \neq i}^N \lambda^{(j)} \right)^2 \|\mathbf{x}\|^2 \|W_{2,1}^{(i)} - W_{2,2}^{(i)}\|. \quad (22)$$

Therefore the block-wise smoothness  $\beta_i^{(N)}$  on layer  $i \in [N]$  of an  $N$ -th layer MLP can be bounded by:

$$\beta_i^{(N)} \leq \left( \prod_{j=1, j \neq i}^N \lambda^{(j)} \right)^2 \|\mathbf{x}\|^2.$$

Let the weights of the  $i$ -th layer be decomposed as  $W^{(i)} = W_1^{(i)} + W_2^{(i)}$ , with the spectral norm given by:

$$\lambda^{(i)} = \|W_1^{(i)} + W_2^{(i)}\|,$$

where  $W_1^{(i)}$  is frozen, and  $W_2^{(i)}$  is trainable. We aim to decompose the following expression:

$$\left( \prod_{j=1, j \neq i}^N \lambda^{(j)} \right)^2.$$

Using the triangle inequality:

$$\lambda^{(j)} = \|W_1^{(j)} + W_2^{(j)}\| \leq \|W_1^{(j)}\| + \|W_2^{(j)}\| \triangleq \lambda_1^{(j)} + \lambda_2^{(j)},$$

where  $\lambda_1^{(j)} = \|W_1^{(j)}\|$  and  $\lambda_2^{(j)} = \|W_2^{(j)}\|$ . Extract  $\lambda_1^{(j)}$  to facilitate product operations:

$$\lambda^{(j)} \leq \lambda_1^{(j)} \left( 1 + \frac{\lambda_2^{(j)}}{\lambda_1^{(j)}} \right).$$

Substituting the spectral norms of each layer into the product and expanding:

$$\prod_{j=1, j \neq i}^N \lambda^{(j)} \leq \prod_{j=1, j \neq i}^N \lambda_1^{(j)} \cdot \prod_{j=1, j \neq i}^N \left( 1 + \frac{\lambda_2^{(j)}}{\lambda_1^{(j)}} \right).$$

Squaring the product and expanding:

$$\left( \prod_{j=1, j \neq i}^N \lambda^{(j)} \right)^2 \leq \left( \prod_{j=1, j \neq i}^N \lambda_1^{(j)} \cdot \prod_{j=1, j \neq i}^N 1 + \frac{\lambda_2^{(j)}}{\lambda_1^{(j)}} \right)^2.$$

The final decomposition result is:

$$\left( \prod_{j=1, j \neq i}^N \lambda^{(j)} \right)^2 \leq \left( \prod_{j=1, j \neq i}^N \lambda_1^{(j)} \right)^2 \cdot \prod_{j=1, j \neq i}^N \left( 1 + \frac{\lambda_2^{(j)}}{\lambda_1^{(j)}} \right)^2$$

Finally, we can get the upper bound of the block-wise smoothness  $\beta_i^{(N)}$  as:

$$\beta_i^{(N)} \leq \left( \prod_{j=1, j \neq i}^N \lambda^{(j)} \right)^2 \|\mathbf{x}\|^2 \leq \left( \prod_{j=1, j \neq i}^N \lambda_1^{(j)} \right)^2 \cdot \prod_{j=1, j \neq i}^N \left( 1 + \frac{\lambda_2^{(j)}}{\lambda_1^{(j)}} \right)^2 \|\mathbf{x}\|^2 \quad (23)$$

Therefore, for the same LLM backbone,  $N$  and  $\lambda_1$  are the same. If we do not add a LoRA adapter to a certain layer, then  $\lambda_2$  of that layer is 0. The fewer LoRA adapters we add to the LLM backbone, the smaller the second term, and the lower the upper bound of the block-wise smoothness  $\beta_i^{(N)}$ .

## E.2 numerical experiments

We introduced  $\beta$ -smoothness in definition 1, which refers to the Lipschitz continuity of the gradient of the loss function. As shown in Appendix E.1, Proposition 2 is a general proposition that can be well extended to the LoRA method. In this section, we use the results of numerical experiments to prove that our theory is reasonable in the LoRA method when discussing the relationship between the number of network layers and  $\beta$ -smoothness.

Assuming that the function we are discussing is continuous and differentiable, we introduce a very small perturbation  $\epsilon = 1e - 5$ . Then, Definition 1 can be simplified to:

$$\|\nabla f(w; z) - \nabla f(w + \epsilon; z)\| \leq \beta \|\epsilon\|. \quad (24)$$

Therefore, the estimation formula for  $\beta$ -smoothness can be obtained:

$$\frac{\|\nabla f(w; z) - \nabla f(w + \epsilon; z)\|}{\|\epsilon\|} \leq \beta. \quad (25)$$

According to Equation 25, we use Llama3-8B as an example to calculate the  $\beta$ -smoothness of different layers of the model across various datasets (Hellaswag, PIQA, RACE, Winogrande). This is done to verify the relationship between the  $\beta$ -smoothness of model and the number of layers. The specific experimental steps are as follows: we selected a model fine-tuned with LoRA for each dataset, perturbed its trainable LoRA parameters, randomly sampled 10 data points from the corresponding dataset as input, and calculated the average  $\beta$ -smoothness of these ten data points. We use this average  $\beta$ -smoothness to represent the  $\beta$ -smoothness of network. The experimental results are shown in Figure 14. The vertical axis uses a logarithmic scale, and it can be seen that across different datasets, the  $\beta$ -smoothness of the Llama3-8B network (i.e., the Smoothness constant in Figure 14) increases exponentially with the number of network layers. In summary, large language models represented by Llama3-8B exhibit properties similar to those of MLP networks as described in Proposition 2, where  $\beta$ -smoothness increases exponentially with the number of network layers.



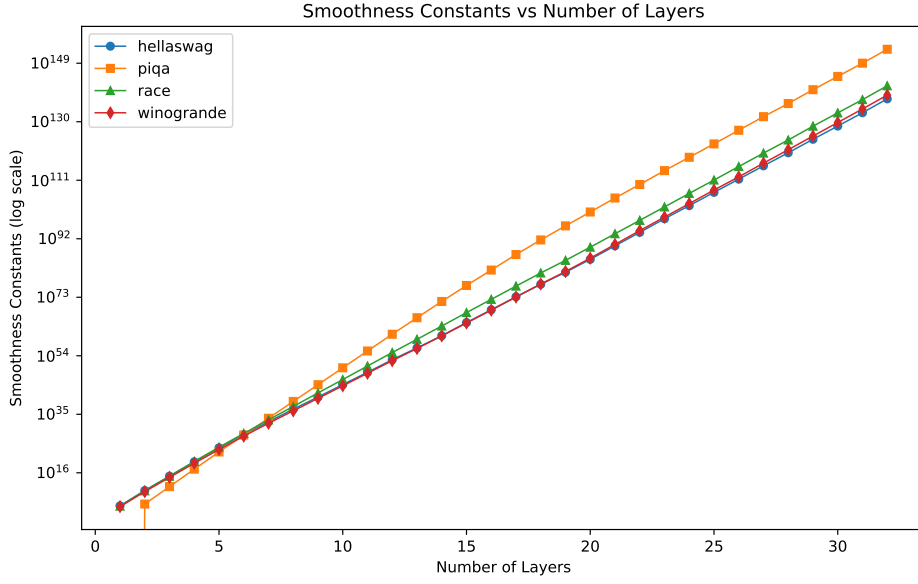


Figure 14: The  $\beta$ -smoothness constants of the Llama3-8B model across different datasets (Hellaswag, PIQA, RACE, Winogrande) as a function of the number of layers. The vertical axis is on a logarithmic scale, demonstrating that the  $\beta$ -smoothness increases exponentially with the number of layers for all datasets.

## F Flexora can identify the best single solution

The literature extensively documents the prevalence of multiple local optima in hyperparameter optimization (HPO) problems (Bao et al., 2021; Franceschi et al., 2017). This phenomenon is particularly relevant in layer selection for LoRA fine-tuning, where varying layer combinations yield divergent performance outcomes due to intricate inter-layer interactions and task-specific characteristics. Here, we systematically analyze the emergence of multiple solutions and demonstrate how Flexora effectively identifies the optimal configuration.

### F.1 Why Multiple Solutions Emerge?

The multiplicity of solutions arises from two primary factors.

**Local vs. Global Optima:** In layer selection problems, multiple local optima typically exist due to varying performance of different layer combinations across tasks. Certain combinations may excel on training data but underperform on validation sets (indicating overfitting), while others demonstrate superior validation performance despite marginally weaker training results.

**Optimization Trajectory:** As shown in Table 15, the evolutionary path of model during optimization leads to distinct local optima at different stages. Initial phases often favor simpler layer configurations, while subsequent optimization may uncover more sophisticated combinations that better capture task-specific characteristics.

### F.2 How Flexora Finds the Optimal Single Solution

To improve the performance of Flexora and prevent convergence to suboptimal local solutions, we implemented a dual-strategy approach. On the training strategy, we established validation set performance as our primary optimization criterion, systematically evaluating various layer combinations to determine the most general configuration. In addition, we adopted an early stopping mechanism that monitors validation performance and terminates optimization when a steady state is reached, thereby preventing overfitting and selecting the best layer combination. On the optimization target preconditioning, we apply continuous relaxation to the target hyperparameters using formula 4 in Section 4.1, which significantly

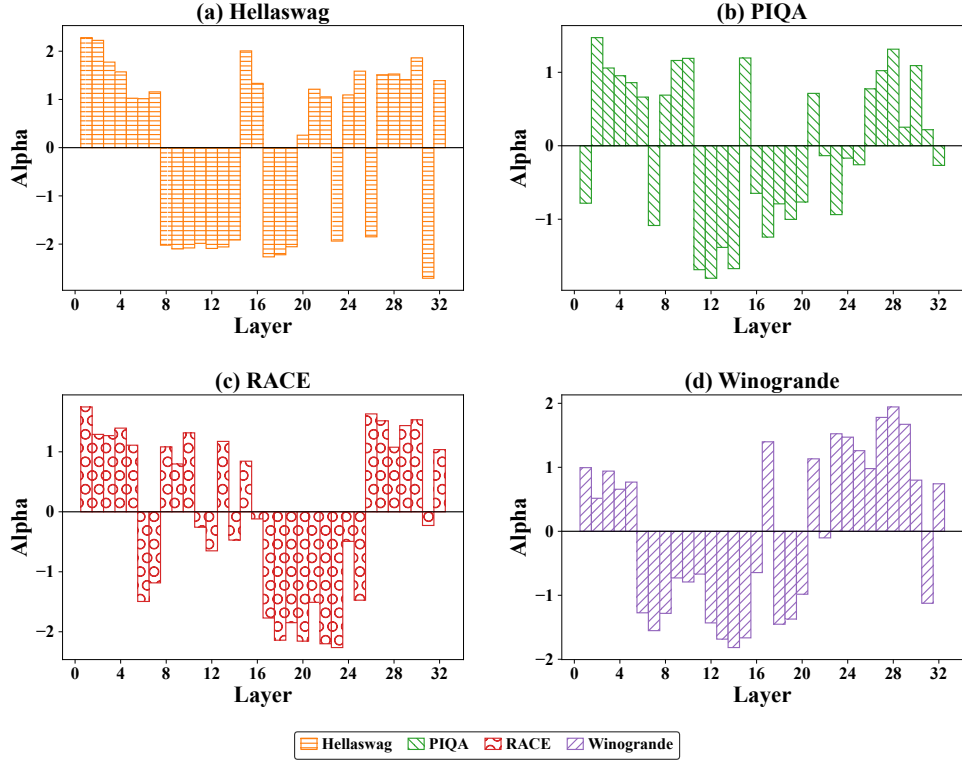


Figure 15: Distribution of  $\alpha$  values across different layers for the Llama3-8B model using LoRA with a rank of 8. The panels (a) through (d) correspond to the Hellaswag, PIQA, RACE, and Winogrande datasets, respectively. Each bar represents the  $\alpha$  value for a specific layer, with the x-axis indicating the layer number and the y-axis showing the  $\alpha$  value.

improves optimization stability and efficiency. As shown in Figure 15, our strategy produces well-converged hyperparameters with significant absolute value differences ( $\alpha > 0$  vs.  $\alpha \leq 0$ ), reflecting significant changes in layer importance. This differential importance ultimately enables Flexora to consistently identify a single optimal solution.

## G Special cases

This section details the performance of Flexora and LoRA across four distinct datasets. The results indicate that Flexora demonstrates superior comprehension and judgment on more challenging questions within the test dataset, compared to LoRA. In certain instances, Flexora successfully explains problems not previously encountered during training, showcasing its robust learning and generalization capabilities.

### Special cases of Hellaswag

```
dataset: Hellaswag
"1": {
  "origin_prompt": "A lady walks to a barbell. She bends down and grabs
the pole. The lady\n
Question: Which ending makes the most sense?\n
A. swings and lands in her arms.\n
B. pulls the barbell forward.\n
C. pulls a rope attached to the barbell.\n
D. stands and lifts the weight over her head.\n
You may choose from 'A', 'B', 'C', 'D'.\n
Answer:",
  "Flexora prediciton": " D",
  "LoRA prediciton" : "B",
  "gold": "D"
},
"2": {
  "origin_prompt": "Two women in a child are shown in a canoe while a man
pulls the canoe while standing in the water, with other individuals
visible in the background. The child and a different man\n
Question: Which ending makes the most sense?\n
A. are then shown paddling down a river in a boat while a woman talks.\n
B. are driving the canoe, they go down the river flowing side to side.\n
C. sit in a canoe while the man paddles.\n
D. walking go down the rapids, while the man in his helicopter almost
falls and goes out of canoehood.\n
You may choose from 'A', 'B', 'C', 'D'.\n
Answer:",
  "Flexora prediciton": " C",
  "LoRA prediciton" : "B",
  "gold": "C"
},
"3": {
  "origin_prompt": "The boy lifts his body above the height of a pole.
The boy lands on his back on to a red mat. The boy\n
Question: Which ending makes the most sense?\n
A. turns his body around on the mat.\n
B. gets up from the mat.\n
C. continues to lift his body over the pole.\n
D. wiggles out of the mat.\n
You may choose from 'A', 'B', 'C', 'D'.\n
Answer:",
  "Flexora prediciton": " B",
  "LoRA prediciton" : "B",
  "gold": "B"
},
"4": {
  "origin_prompt": "We see a person holding face wash then putting it on
their face. They rinse the face and add the face wash with a brush. We\n
Question: Which ending makes the most sense?\n
A. see a closing title screen.\n
B. see a black screen with the credits.\n
C. see an illustration on how to add the wash using a brush.\n
D. then see a replay then the person putting the face wash on.\n
You may choose from 'A', 'B', 'C', 'D'.\n
Answer:",
  "Flexora prediction": " C",
  "LoRA prediciton" : "A",
  "gold": "C"
},
},
```

## Special cases of PIQA

```
dataset: PIQA
"1": {
  "origin_prompt": "ice box\n
  A. will turn into a cooler if you add water to it\n
  B. will turn into a cooler if you add soda to it\n
  Answer:",
  "Flexora prediciton": "A",
  "LoRA prediciton" : "A",
  "gold": "A"
},
"2": {
  "origin_prompt": "how do you put eyelashes on?\n
  A. glue them on with mascara.\n
  B. put eyelash glue on the fake eyelashes and then let it get tacky.\n
  then, place it on top of your actual eyelashes and let it dry on.\n
  Answer:",
  "Flexora prediciton": "A",
  "LoRA prediciton" : "B",
  "gold": "B"
},
"3": {
  "origin_prompt": "How do I fill holes and tiny gaps in the concrete when\n
  making a concrete countertop?\n
  A. Use a concrete slurry\n
  B. Use a concrete brush\n
  Answer:",
  "Flexora prediciton": "A",
  "LoRA prediciton" : "B",
  "gold": "A"
}
"4": {
  "origin_prompt": "When I'm deep frying a turkey, how defrosted\n
  should it be?\n
  A. It should be completely defrosted, otherwise the oil may\n
  boil over and start a fire.\n
  B. It should be completely frozen, otherwise the oil may boil\n
  over and start a fire.\n
  Answer:",
  "\"ours{} prediction\": \" A\\nExplanation: A turkey should be completely\n
  defrosted before deep frying. If it is not, the oil may boil over\n
  and start a fire.\",
  "Flexora prediciton" : "A",
  "LoRA prediciton" : "B",
  "gold": "A"
},
"5": {
  "origin_prompt": "How do you properly prepare a steak.\n
  A. Take the steak out of warm storage and let come to room temperature,\n
  generously add salt and pepper to both sides and let sit\n
  for 10 minutes.\n
  B. Take the steak out of cold storage and let come to room temperature,\n
  generously add salt and pepper to both sides and let sit\n
  for 10 minutes.\n
  Answer:",
  "Flexora prediction": \" B\\nExplanation: B. Taking the steak\n
  out of cold storage and letting it come to room temperature is\n
  the correct answer. The steak should be at room temperature before\n
  cooking. The steak should be generously salted and peppered on both\n
  sides and let sit for 10 minutes.\",
  "LoRA prediciton" : "B",
  "gold": "B"
},
"6": {
  "origin_prompt": "To cream butter and sugar together, you can\n
  A. Place it in a bowl and use a hand warmer\n
  B. Place in a bowl and use a hand mixer\n
  Answer:",
  "Flexora prediction": \" B\\nExplanation: B. Place in a bowl and\n
  use a hand mixer\\nExplanation: To cream butter and sugar together,\n
  you can place it in a bowl and use a hand mixer.\",
  "LoRA prediciton" : "B",
  "gold": "B"
},
}
```

## Special cases of RACE

```
dataset: RACE
  "1": {
    "origin_prompt": "Read the article, and answer the question by replying A,
    B, C or D.\n\n
    Article:\n
    The rain had continued for a week and the flood
    had created a big river which were running by Nancy Brown's
    farm. As she tried to gather her cows to a higher ground,
    she slipped and hit her head on a fallen tree trunk.
    The fall made her unconscious for a moment or two. When she came to,
    Lizzie, one of her oldest and favorite cows, was licking her face. \n
    At that time, the water level on the farm was still rising.
    Nancy gathered all her strength to get up and began walking
    slowly with Lizzie. The rain had become much heavier,
    and the water in the field was now waist high. Nancy's pace
    got slower and slower because she felt a great pain in her head.
    Finally, all she could do was to throw her arm around Lizzie's
    neck and try to hang on. About 20 minutes later, Lizzie managed
    to pull herself and Nancy out of the rising water and onto
    a bit of high land, which seemed like a small island in
    the middle of a lake of white water. \n
    Even though it was about noon, the sky was so dark and the rain
    and lightning was so bad that it took rescuers more than
    two hours to discover Nancy. A man from a helicopter
    lowered a rope, but Nancy couldn't catch it. A moment later,
    two men landed on the small island from a ladder in the helicopter.
    They raised her into the helicopter and took her to the school gym,
    where the Red Cross had set up an emergency shelter.
    \n
    When the flood disappeared two days later, Nancy immediately
    went back to the \island.\" Lizzie was gone. She was one of
    19 cows that Nancy had lost in the flood. \n
    I owe my life to
    her,\" said Nancy with tears.\n\n
    Q: What did Nancy try to do before she fell over?\n\n
    A. Measure the depth of the river\n
    B. Look for a fallen tree trunk\n
    C. Protect her cows from being drowned\n
    D. Run away from the flooded farm\n
    ",
    "Flexora prediciton": "D",
    "LoRA prediciton" : "B",
    "gold": "D"
  }
```



## Special cases of Winogrande

```
dataset: Winogrande
"1": {
  "origin_prompt": "Question: Sarah was a much better surgeon
than Maria so _ always got the easier cases.\n
A. Sarah\n
B. Maria\n
Answer:",
  "Flexora prediciton": "B",
  "LoRA prediciton" : "B",
  "gold": "B"
},
"2": {
  "origin_prompt": "Question: Sarah was a much better surgeon
than Maria so _ always got the harder cases.\n
A. Sarah\n
B. Maria\n
Answer:",
  "Flexora prediciton": "B",
  "LoRA prediciton" : "B",
  "gold": "A"
},
"3": {
  "origin_prompt": "Question: They were worried the wine would ruin
the bed and the blanket, but the _ was't ruined.\n
A. blanket\n
B. bed\n
Answer:",
  "Flexora prediciton": "B",
  "LoRA prediciton" : "A",
  "gold": "B"
},
"4": {
  "origin_prompt": "Question: Terry tried to bake the eggplant
in the toaster oven but the _ was too big.\n
A. eggplant\n
B. toaster\n
Answer:",
  "Flexora prediction": " A\nExplanation: The eggplant was
too big to fit in the toaster oven.",
  "LoRA prediciton" : "B",
  "gold": "A"
},
"5": {
  "origin_prompt": "Question: At night, Jeffrey always stays up
later than Hunter to watch TV because _ wakes up late.\n
A. Jeffrey\n
B. Hunter\n
Answer:",
  "Flexora prediciton": "A",
  "LoRA prediciton" : "B",
  "gold": "A"
},
"6": {
  "origin_prompt": "Question: The cat of Sarah has some mouth problems,
so she takes it to see Maria. _ is a responsible cat owner.\n
A. Sarah\n
B. Maria\n
Answer:",
  "Flexora prediction": " A\nExplanation: A is a responsible cat
owner because she takes her cat to see a veterinarian.",
  "LoRA prediciton" : "B",
  "gold": "A"
},
"7": {
  "origin_prompt": "Question: Benjamin was chosen instead of Brett to
be the makeup artist for the play because _ was less experienced.\n
A. Benjamin\n
B. Brett\n
Answer:",
  "Flexora prediction": " B",
  "LoRA prediciton" : "A",
  "gold": "B"
},
}
```