

# Sensitivity-LoRA : Low-Load Sensitivity-Based Fine-Tuning for Large Language Models

Anonymous ACL submission

## Abstract

Large Language Models (LLMs) have revolutionized various domains, including everyday applications and scientific research. However, adapting LLMs from general-purpose models to specialized tasks remains challenging, particularly in resource-constrained environments. Parameter-Efficient Fine-Tuning (PEFT), especially Low-Rank Adaptation (LoRA), has emerged as a promising approach to fine-tuning LLMs while reducing computational and memory overhead. Despite its advantages, existing LoRA techniques still struggle with computational inefficiency, complexity, and instability, limiting their practical applicability. To address these limitations, we propose Sensitivity-LoRA, an efficient fine-tuning method that dynamically allocates ranks to weight matrices based on both their global and local sensitivities. It leverages the second-order derivatives (Hessian Matrix) of the loss function to effectively capture weight sensitivity, enabling optimal rank allocation with minimal computational overhead. Our experimental results have demonstrated robust effectiveness, efficiency and stability of Sensitivity-LoRA across diverse tasks and benchmarks.

## 1 Introduction

Large language models (LLMs) have emerged as revolutionary tools, delivering state-of-the-art (SOTA) performance across a wide spectrum of tasks and applications (Ding et al., 2022; Qin et al., 2023; Zhu et al., 2023b,a; Li et al., 2023; Zhang et al., 2023a; Huang et al., 2023; Wang et al., 2023). Despite these advancements, fine-tuning remains a critical technique for effectively adapting LLMs from general-purpose models to specialized applications, especially in resource-constrained environments. However, full-parameter fine-tuning can be prohibitively resource-intensive, requiring significant computational power and GPU capacity. To address this limitation, the research community

introduced parameter-efficient fine-tuning (PEFT) (Houlsby et al., 2019; Lester et al., 2021; Li and Liang, 2021; Zaken et al., 2022; Hu et al., 2022a), which aims to balance accuracy and efficiency by selectively updating a subset of model parameters. A decent amount of PEFT methods have been validated to be effective across a variety of models and tasks, often achieving results comparable to those of full-parameter fine-tuning (Lester et al., 2021; Zaken et al., 2022; Hu et al., 2022a).

Among these PEFT methods, the Low-Rank Adaptation (LoRA) method based on reparameterization (Hu et al., 2022b) is considered one of the most efficient and effective approaches. It utilizes low-rank decomposition to approximate the updates of model weights. This strategy draws from prior research demonstrating that the learned over-parametrized models in fact reside on a low intrinsic dimension (Li et al., 2018; Aghajanyan et al., 2020). During the training process, the update of the weight matrix ( $\Delta W$ ), can be approximated as the product of two smaller matrices  $B$  and  $A$ , that is:

$$\Delta W \approx B \cdot A \quad (1)$$

where  $\Delta W \in R^{d_1 \times d_2}$ ,  $A \in R^{r \times d_2}$  and  $B \in R^{d_1 \times r}$  with  $r \ll \{d_1, d_2\}$ . Through this method, it can approximate the update of the weight matrix with fewer parameters, thereby improving the parameter efficiency of the model. However, the full potential of LoRA remains constrained by its inherent design limitations. Specifically, it assumes a uniform rank  $r$  for each incremental matrix, not accounting for the varying significance of weight matrices across different modules and layers (Hu et al., 2023; Zhang et al., 2023b).

To address this limitation, dynamic rank allocation has emerged as a key solution by allocating the rank  $r$  to each different module or layer according to its specific requirements. Existing methods achieve this through three main approaches:

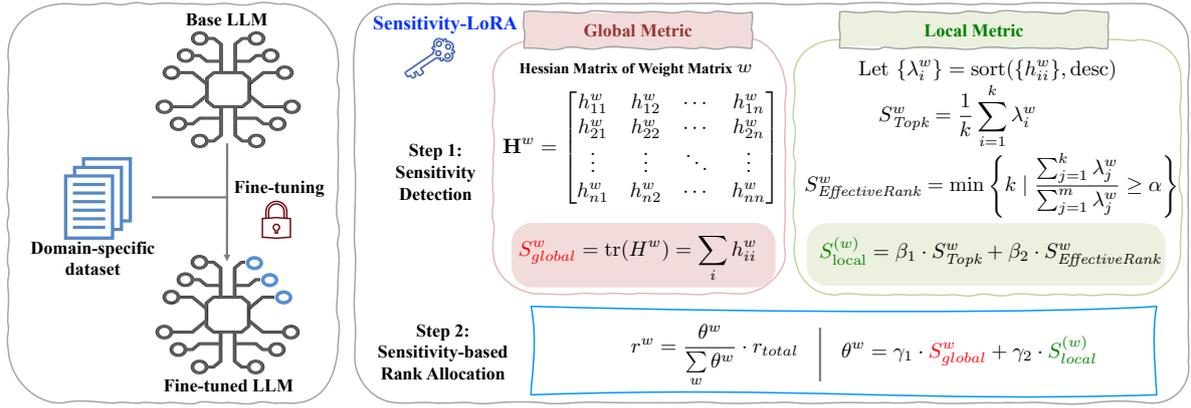


Figure 1: Pipeline of the Sensitivity-LoRA Method: **Step 1 - Sensitivity detection** via Hessian-based metrics, including global and local sensitivity measures. ( $h_{ij}^w = \frac{\partial^2 E^w}{\partial w_i \partial w_j}$ ,  $h_{ii} = \frac{\partial^2 E^w}{\partial w_i^2}$ , where  $E^w$  denotes the change in loss function regarding weight matrix  $w$ ;  $\text{tr}(H^w)$  denotes the trace of  $H^w$ .) **Step 2 - Dynamic rank allocation** based on global and local sensitivity. ( $r^w$  denotes the allocated rank of weight matrix  $w$ ,  $r_{total}$  denotes the total rank of all weight matrices in the model.)

083 singular value decomposition (SVD), single-rank  
084 decomposition (SRD), and rank sampling. SVD-  
085 based methods (Zhang et al., 2023c; Hu et al., 2023;  
086 Zhang et al., 2023b) decompose matrix BA into  
087 an SVD form and selectively truncate the singular  
088 values in order to allocate the matrix rank. How-  
089 ever, this process is computationally expensive and  
090 requires additional memory to store singular val-  
091 ues and vectors. SRD-based methods (Mao et al.,  
092 2024; Zhang et al., 2024; Liu et al., 2024) decom-  
093 pose matrix BA into single-rank components and  
094 allocate the ranks by selecting the proper compo-  
095 nents. However, optimizing single-rank compo-  
096 nents and the pruning process can increase compu-  
097 tational complexity, potentially offsetting efficiency  
098 gains. Rank sampling-based methods (Valipour  
099 et al., 2022) allocate ranks directly by random sam-  
100 pling. However, the randomness introduced by  
101 sampling could increase potential instability in the  
102 training.

103 In view of these challenges and opportunities,  
104 we propose Sensitivity-LoRA, which can rapidly  
105 allocate rank to the weight matrix based on the  
106 sensitivity of the parameters, without incurring a  
107 significant computational load. Specifically, we uti-  
108 lize the second derivatives of the loss function with  
109 respect to the parameters (Hessian matrix) to as-  
110 certain the sensitivity of each parameter within the  
111 weight matrix. To comprehensively evaluate the  
112 sensitivity of the parameter matrix, we employ met-  
113 rics such as the trace of the Hessian matrix, *Topk*  
114 and *Effective Rank* to measure its global and local  
115 sensitivities respectively. By integrating various

116 metrics, we determine the rank allocation weights  
117 corresponding to the weight matrices to achieve  
118 rank allocation. The efficiency, stability, and gener-  
119 ality of our approach have been validated through  
120 extensive experiments on various tasks, such as sen-  
121 timent analysis, natural language inference, ques-  
122 tion answering, and text generation.

123 In summary, the main contributions of our paper  
124 are listed as follows:

- 125 • We introduce the second derivatives of the loss  
126 function with respect to the weight matrix to  
127 measure their sensitivity.
- 128 • We achieve rank allocation by taking into ac-  
129 count both the global and local sensitivity of  
130 the weight matrix.
- 131 • Extensive experiments demonstrate the ef-  
132 fectiveness, stability, and efficiency of our  
133 method.

## 134 2 Related Work

135 Existing PEFT approaches can be classified into  
136 four main types in terms of memory efficiency, stor-  
137 age efficiency, and inference overhead, as follows:

### 138 2.1 Additive PEFT

139 Additive PEFT introduces lightweight modules into  
140 the model architecture, such as adapters and soft  
141 prompts, while keeping the pre-trained backbone  
142 frozen. Adapters add small networks with down-  
143 projection and up-projection matrices, enabling  
144 task-specific learning with minimal parameter up-  
145 dates (Houlsby et al., 2019; Lester et al., 2021).

Soft prompts prepend learnable embeddings to the input sequence, allowing fine-tuning by modifying input activations only (Li and Liang, 2021; Zaken et al., 2022). These methods typically require updating less than 1% of the total parameters, significantly reducing computation and memory costs, making them ideal for resource-constrained environments (Hu et al., 2022a).

## 2.2 Selective PEFT

Selective PEFT fine-tunes a subset of the existing parameters in a pre-trained model, rather than adding new modules. It employs binary masks to identify and update only the most important parameters while keeping the majority frozen. Techniques like Diff pruning and FishMask leverage Fisher information or parameter sensitivity analysis to select critical parameters for fine-tuning (Zaken et al., 2022; Li and Liang, 2021). This approach avoids increasing model complexity and is particularly suited for scenarios where only a small fraction of the model contributes significantly to performance.

## 2.3 Reparameterized PEFT

Reparameterized PEFT utilizes low-rank parameterization techniques to represent model weights in a reduced form during training. LoRA (Low-Rank Adaptation) is a prominent example, introducing low-rank matrices to fine-tune specific weights while maintaining high inference efficiency (Hu et al., 2022a). Other methods, such as Compacter, use the Kronecker product for parameter reparameterization, further reducing memory requirements and computational costs (Houlsby et al., 2019). Reparameterized PEFT is highly effective for large-scale models where resource constraints are critical.

## 2.4 Hybrid PEFT

Hybrid PEFT combines the strengths of Additive, Selective, and Reparameterized PEFT methods into a unified framework. For example, UniPELT integrates LoRA, adapters, and soft prompts, allowing dynamic selection of the most suitable module for specific tasks through gating mechanisms (Zaken et al., 2022). This hybrid approach enhances adaptability and task performance by leveraging the complementary advantages of different PEFT strategies (Li and Liang, 2021; Hu et al., 2022a).

## 3 Methodology

In this section, we firstly introduce the concept of weight sensitivity, followed by a formal definition of global and local sensitivity metrics of weight matrices. Next, we propose effective allocation strategies to optimize the dynamic rank allocation process based on these sensitivity metrics. The pipeline of our method is presented in figure 1.

### 3.1 Weight Sensitivity

Consider a neural network whose dynamics is driven by a collection of parameters  $w$  and a loss function  $E$ , which guides its learning dynamics. When a small perturbation  $\delta w$  is introduced to the parameters, the resulting change in the loss function can be expressed using a Taylor series expansion up to the second-order term, with higher-order terms captured by  $O(\|\delta w\|^3)$  as follows:

$$E(w+\delta w) = E(w) + g^T \delta w + \frac{1}{2} \delta w^T H \delta w + O(\|\delta w\|^3) \quad (2)$$

where  $g$  denotes the gradient vector of the loss function  $E$  with respect to the parameters  $w$ , indicating the rate of change of the loss function in the direction of each parameter.  $H$  represents the Hessian matrix of the loss function  $E$ , which is a matrix of second-order partial derivatives and contains information about the curvature of the loss function at the current parameter point. The change in the loss function  $\Delta E$  can be represented by the following expression:

$$\Delta E = g^T \delta w + \frac{1}{2} \delta w^T H \delta w + O(\|\delta w\|^3) \quad (3)$$

By expanding the components of  $\Delta E$ , we have:

$$\Delta E = \sum_i g_i \delta w_i + \frac{1}{2} \sum_{i,j} h_{ij} \delta w_i \delta w_j + O(\|\delta w\|^3) \quad (4)$$

where  $g_i$  and  $h_{ij}$  are the gradient and Hessian elements, respectively.

For a well-trained neural network, when the parameter  $w$  is located at a local minimum of the loss function, the gradient  $g$  becomes zero. Then, the above equation can be simplified to

$$\Delta E = \frac{1}{2} \sum_{i,j} h_{ij} \delta w_i \delta w_j + O(\|\delta w\|^3) \quad (5)$$

Additionally, the parameters  $w$  are nearly independent at the minimum, the Hessian matrix  $H$  tends to be diagonal-dominant, meaning that the

off-diagonal elements, which represent interactions between different parameters, become negligible. Then, the above equation can be simplified to

$$\Delta E \approx \frac{1}{2} \sum_i h_{ii} \delta_{w_i}^2 + O(\|\delta_w\|^3) \quad (6)$$

Given that the perturbation in the weights ( $\delta w$ ) is sufficiently small, the higher-order term becomes negligible compared to the quadratic term, and therefore can be disregarded. Consequently, the above formula can be further simplified to

$$\Delta E \approx \frac{1}{2} \sum_i h_{ii} \delta_{w_i}^2 \quad (7)$$

Consequently, the diagonal elements of the Hessian matrix serve as a reliable indicator of weight sensitivity. Specifically, the diagonal elements of the Hessian matrix represent the local curvature of the loss function in the direction of each parameter. Larger diagonal elements indicate more dramatic changes in the loss function in that parameter direction, thereby showing that the parameter has a more significant impact on model performance. Furthermore, the methods for approximating the Hessian matrix have become increasingly sophisticated, ensuring that they do not place a substantial load on computational resources.

## 3.2 Rank Allocation Metric

### 3.2.1 Global Metric

The global sensitivity measurement aims to evaluate the overall impact of an entire parameter (or weight) matrix on model output. It quantifies how variations in this weight matrix affect the loss function. To capture this dynamics, the Hessian matrix, which consists of the second-order partial derivatives of the loss function with respect to the weight matrix, is used. Given that the Hessian matrix tends to be diagonal-dominant at the minimum, its trace can serve as an effective global sensitivity indicator. Formally, the global sensitivity  $S_{global}^w$  of weight matrix  $w$  can be defined as:

$$S_{global}^w = \text{tr}(H^w) = \sum_i h_{ii}^w \quad (8)$$

where  $h_{ii}^w$  is the  $i$ -th diagonal element of the Hessian matrix  $H^w$ , and  $\text{tr}(H^w)$  denotes the trace of  $H^w$ . Since the diagonal elements reflect the impact of individual parameter changes on the loss function, a larger trace value indicates that the model is more sensitive to its changes. This suggests more

parameters make significant contributions to the changes in the loss function, emphasizing their role in model performance.

### 3.2.2 Local Metric

While certain weight matrices might have low overall sensitivities, specific weight elements within these matrices can still have high sensitivity, significantly impacting model performance. As such, it is essential to account for local sensitivity to capture finer-grained variations in parameter influence on the loss function. To address this, we introduce two metrics: *Topk* and *Efficient Rank*.

The *Topk* metric approximates local sensitivity of a weight matrix by averaging its largest  $k$  diagonal elements of Hessian matrix, based on the assumption that most of the matrix's energy or sensitivity is concentrated in these large values. By focusing on *Top k* diagonal elements, the *Topk* metric can guide us to prioritize these critical weights during weight pruning or optimization processes. It reduces computational complexity while preserving the most impactful weights for model performance. The computation formula for the *Topk* metric of weight matrix  $w$  is as follows:

$$S_{Topk}^w = \frac{1}{k} \sum_{i=1}^k \lambda_i^w \quad (9)$$

where  $\lambda_i^w$  represents the diagonal elements of Hessian matrix  $H^w$  sorted in descending order, and  $k$  denotes the number of diagonal elements selected.

The *Effective Rank* metric determines the minimum rank that captures most of the energy of a weight matrix based on the cumulative contribution of the diagonal elements of its Hessian matrix. By establishing a threshold for the cumulative contribution rate (such as 0.9 or 0.95), the *Effective Rank* metric identifies the minimum number of eigenvalues needed to achieve this threshold, thereby appropriately ranking the weight matrix. The key benefit of this metric is ensuring the stability of the rank allocation process. The formula for *Effective Rank* of weight matrix  $w$  is as follows:

$$S_{EffectiveRank}^w = \min \left\{ k \mid \frac{\sum_{j=1}^k \lambda_j^w}{\sum_{j=1}^m \lambda_j^w} \geq \alpha \right\} \quad (10)$$

where  $\lambda_j^w$  is the  $j$ -th diagonal element of  $H^w$  in non-increasing order,  $m$  is the total number of diagonal elements, and  $k$  is the minimum number of diagonal elements required for the cumulative contribution rate to reach the threshold  $\alpha$ .

Method	MNLI	SST-2	MRPC	CoLA	QNLI	QQP	RTE	STS-B	Avg.
LoRA	87.26	93.46	87.01	58.83	92.95	90.50	79.42	91.03	85.06
$S_g$ -LoRA	87.15	94.61	88.48	60.83	93.08	90.54	79.78	91.03	85.69
$S_T$ -LoRA	86.94	93.81	88.48	60.82	92.75	90.61	79.06	91.20	85.46
$S_E$ -LoRA	87.54	94.04	89.46	63.07	92.77	90.58	80.51	91.16	86.14
Ours	<b>87.84</b>	<b>95.11</b>	<b>89.53</b>	<b>63.21</b>	<b>93.08</b>	<b>90.72</b>	<b>80.64</b>	<b>91.31</b>	<b>86.43</b>

Table 1: The comparative results of Sensitivity-LoRA and baseline methods for fine-tuning RoBERTa-base on the eight datasets of the GLUE benchmark. The bolded values represent the optimal outcomes. We reported the overall accuracy for MNLI (both matched and mismatched), the Matthew’s correlation for CoLA, the Combine Score for STS-B, and the accuracy for the other tasks. ( $S_g$ -LoRA,  $S_T$ -LoRA and  $S_E$ -LoRA represent LoRA fine-tuning using the global metric, *Topk* metric, and *Efficient Rank* metric for rank allocation, respectively.)

To ensure the effectiveness and stability, we integrate *Topk* and *Efficient Rank* metrics together to define the local sensitivity metric  $S_{local}^w$  of weight matrix  $w$  as follows:

$$S_{local}^w = \beta_1 \cdot S_{Topk}^w + \beta_2 \cdot S_{EfficientRank}^w \quad (11)$$

where  $\beta_1 + \beta_2 = 1$ .

### 3.3 Rank Allocation Strategy

Taking into account both global and local metrics, we define a refined rank allocation strategy to determine the rank allocation weights  $\theta^w$  of weight matrix  $w$  by integrating global and local sensitivities:

$$\theta^w = \gamma_1 \cdot S_{global}^w + \gamma_2 \cdot S_{local}^w \quad (12)$$

subject to the constraints  $\gamma_1 + \gamma_2 = 1$ . Hence, we can derive the formula for rank allocation as follows:

$$r^w = \frac{\theta^w}{\sum_w \theta^w} \cdot r_{total} \quad (13)$$

where  $r^w$  denotes the rank allocated to weight matrix  $w$ , and  $r_{total}$  represents the total rank of all weight matrices in the model.

## 4 Experiments

### 4.1 Experimental Setup

**Benchmarks** We evaluate the performance of Sensitivity-LoRA across diverse tasks, including NLG (Natural Language Generation) and NLU (Natural Language Understanding). Specifically, for the NLU tasks, we select RoBERTa-base (Liu, 2019) as the base model and evaluate its performance on eight sub-tasks of the GLUE (General Language Understanding Evaluation) benchmark (Wang, 2018): MNLI (Williams et al., 2017), SST-2 (Socher et al., 2013), MRPC (Dolan and Brockett, 2005), CoLA (Warstadt, 2019), QNLI (Rajpurkar

et al., 2018), QQP<sup>1</sup>, RTE (Wang, 2018), and STS-B (Cer et al., 2017). For the NLG tasks, we assess the performance of the proposed method with GPT-2 Large (Radford et al., 2019) on E2E benchmark (Novikova et al., 2017).

**Baselines** The LoRA (Hu et al., 2022b) with a uniform rank allocation is adopted as the baseline for comparison with the proposed Sensitivity-LoRA. We maintain the same hyperparameter settings as those used in LoRA, varying only the rank allocation method. We further compare the fine-tuning methods that allocate ranks based solely on global or local sensitivity metrics with the proposed Sensitivity-LoRA to highlight its effectiveness. Specifically,  $S_g$ -LoRA,  $S_T$ -LoRA and  $S_E$ -LoRA represent LoRA fine-tuning methods that allocate ranks using the global metric, *Topk* metric, and *Efficient Rank*, respectively.

**Implementation Details** Our code is implemented using the PyTorch (Paszke et al., 2019) framework and Transformers (Wolf, 2020) libraries, with all experiments conducted on four NVIDIA A100 GPUs. We fine-tune GPT-2 Large (36 layers) and RoBERTa-base (12 layers) using Sensitivity-LoRA. For a fair comparison, we adhere to the same total rank  $r_{total}$  as specified in LoRA. Specifically, for GPT-2 Large, we set  $r_{total} = 144$ , while for RoBERTa-base, we set  $r_{total} = 96$ . We designate the local metric  $S_{Topk}$  with  $k$  set to half of the total number of diagonal elements, and set the parameter  $\alpha$  in the *Efficient Rank* metric to 0.8. Considering the scale characteristics of the different metrics, we set the parameter  $\beta_1$  to 0.45,  $\beta_2$  to 0.55,  $\gamma_1$  to 0.005, and  $\gamma_2$  to 0.095.

<sup>1</sup><https://quoradata.quora.com/First-Quora-Dataset-Release-Question-Pairs>

Method	BLEU	NIST	MET	ROUGE-L	CIDEr
LoRA	67.15	8.58	46.27	69.02	2.39
$S_g$ -LoRA	67.61	8.61	46.24	69.34	2.38
$S_T$ -LoRA	67.37	8.60	46.14	69.47	2.41
$S_E$ -LoRA	67.94	8.64	46.14	69.40	2.40
Ours	<b>67.99</b>	<b>8.65</b>	<b>46.58</b>	<b>69.64</b>	<b>2.41</b>

Table 2: The results of GPT-2 Large on the E2E dataset for the baseline and Proposed method. For all metrics, higher values indicate better performance, and the bold values represent the best results. ( $S_g$ -LoRA,  $S_T$ -LoRA and  $S_E$ -LoRA represent LoRA fine-tuning using the global metric, *Topk* metric, and *Efficient Rank* metric for rank allocation, respectively.)

## 4.2 Main Result

At first, we assess the performance of Sensitivity-LoRA on NLU tasks by fine-tuning the RoBERTa-base model on the eight tasks of the GLUE benchmark. The overall accuracy for MNLI, the Matthew’s correlation for CoLA, the Combined Score for STS-B, and the accuracy for the other tasks are reported in table 1. Sensitivity-LoRA demonstrates outstanding performance in a variety of natural language understanding tasks. Compared to baseline methods, our approach achieve an average score of 86.43 over the eight datasets, representing a gain of 1.37 over the LoRA. The evaluation metrics of each dataset also significantly surpass those of other methods. Sensitivity-LoRA leverages the second-order derivatives of the loss function to extract weight-wise importance metrics, incorporating both local and global sensitivity. Based on these metrics, it dynamically determines the optimal rank allocation, thereby achieving exceptional performance.

Next, we fine-tune the GPT-2 Large model on the E2E (Novikova et al., 2017) datasets for NLG task evaluation. The E2E dataset is commonly employed for end-to-end data-to-text generation tasks. We train the model on 42,000 data-to-text pairs from the restaurant domain. To comprehensively evaluate both the baseline methods and Sensitivity-LoRA, we utilized five evaluation metrics: BLEU (Papineni et al., 2002), NIST (Lin and Och, 2004), METEOR (Banerjee and Lavie, 2005), ROUGE-L (Lin, 2004), and CIDEr (Vedantam et al., 2015). Table 2 presents the experimental results of Sensitivity-LoRA on the E2E dataset. We evaluate the impact of global and local metrics for rank allocation, alongside a combined rank allocation strategy incorporating both types of information. The experimental results indicate that our proposed rank allocation strategy, achieve su-

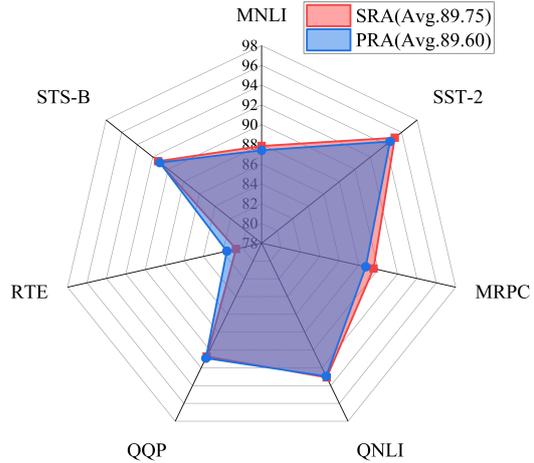


Figure 2: Comparison of the evaluation metrics for RoBERTa-base fine-tuned on seven datasets from the GLUE benchmark, using both PRA and SRA rank assignment methods.

perior performance across all evaluation metrics. Although relying solely on global or local information for rank allocation can improve model performance to some extent, its performance on several metrics (METEOR and CIDEr) still falls slightly short of the LoRA model. This may be because an overemphasis on either global or local information diminishes the importance of certain critical weights, thereby affecting the overall performance of the model.

## 4.3 Comparison of Rank Assignment Methods

In this section, we compare two rank allocation methods for model weights, based on global or local sensitivity metrics. The Progressive Rank Assignment (PRA) method first sorts the global or local sensitivity metrics in descending order, subsequently allocating ranks progressively within a specified range. Weights with higher sensitivity are assigned higher ranks, with rank decrementing in 20% intervals. The rank ranges for GPT-2

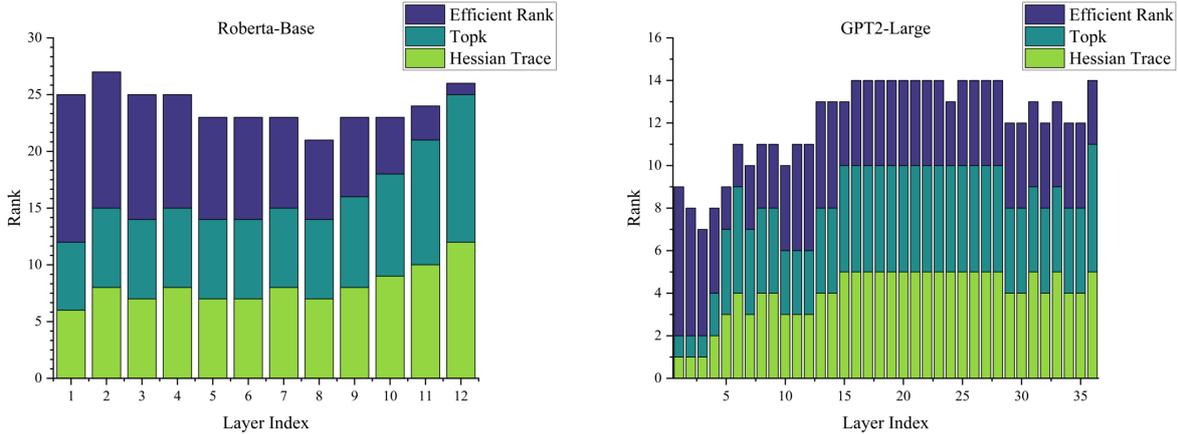


Figure 3: The rank allocation for each layer of GPT-2 Large and RoBERTa-base under different rank allocation strategies. Different colors represent different allocation strategies, and the height of each bar in the histogram corresponds to the rank allocated to that layer by the respective strategy.

Large and RoBERTa-base are 2~6 and 6~10, respectively. The Scaled Rank Assignment (SRA) method (mentioned in section 3.3) allocates ranks according to the proportion of each weight’s global or local sensitivity relative to the model’s total sensitivity. To visually compare the effectiveness of these two allocation methods, we apply both strategies to the global and local sensitivity metrics and subsequently combine them using the corresponding rank allocation strategy. We then fine-tune RoBERTa-base on seven datasets from the GLUE benchmark. As depicted in figure 2, the evaluation metrics for both allocation methods are compared across datasets. SRA and PRA achieve average scores of 89.75 and 89.60, respectively. Although SRA perform marginally worse than PRA on the RTE dataset, it surpass PRA on all other datasets. This highlights the effectiveness of the SRA allocation strategy. Consequently, we adapt the SRA method for rank allocation in this paper.

#### 4.4 Optimized Rank Allocation

As shown in figure 3, we visualize the global and local rank allocation results for GPT-2 Large and RoBERTa-base under the SRA rank assignment method described in section 3.3. From the figure 3, it is clear that the global sensitivity metric, Hessian Trace, allocates a larger rank budget to the intermediate and deeper layers of the model, with less emphasis on the initial layers. The local sensitivity metric, *Topk*, primarily allocates rank to the middle layers of the model. In contrast, the *Efficient Rank* approach assigns higher ranks to the initial layers, while exhibiting a decreasing trend

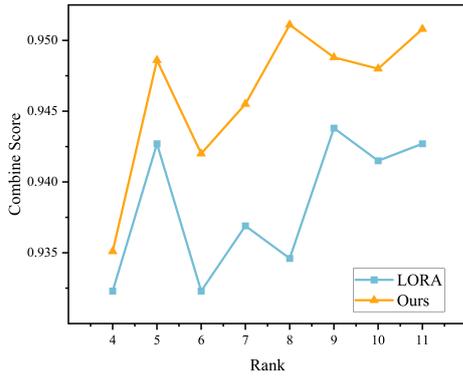
in rank allocation for subsequent layers. Each of these three sensitivity metrics emphasizes different aspects, indicating that relying on a single source of information for decision-making is insufficient. This underscores the necessity of Sensitivity-LoRA in integrating these diverse information sources for dynamic rank allocation.

#### 4.5 Analysis of Different Rank Budgets

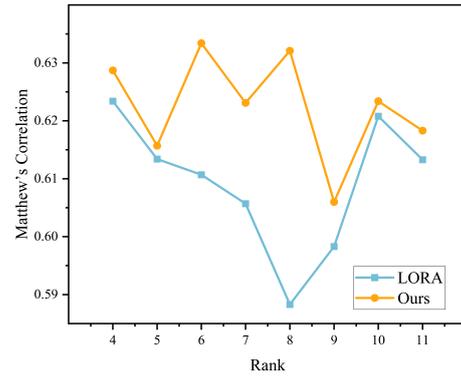
We further analyze the performance of Sensitivity-LoRA and the baseline method, LoRA, under various rank budget scenarios. Figure 4 illustrates the changes in evaluation metrics for RoBERTa-base on the SST-2 and CoLA datasets for different rank budgets. We evaluate the model’s performance under eight rank budgets, ranging from 4 to 11. For Sensitivity-LoRA, the ranks in the figure 4 represent the average rank of each model layer. As seen in the figure 4, the proposed Sensitivity-LoRA consistently outperforms LoRA across different rank budget scenarios on both datasets. On the SST-2 dataset, as the rank budget increases, the model performance exhibits an overall upward trend, with a notable improvement at  $r = 5$  and a peak at  $r = 8$ . On the CoLA dataset, however, as the rank budget increases from 4, the model performance for the baseline method LoRA decreases after fine-tuning, while Sensitivity-LoRA initially declines slightly before increasing again. Despite some fluctuations, Sensitivity-LoRA consistently outperforms LoRA.

#### 4.6 Case study

Figure 5 illustrates the performance of the GPT-2 Large and RoBERTa-base models, fine-tuned using



(a) SST-2



(b) CoLA

Figure 4: The evaluation metric changes for RoBERTa-base fine-tuned on the SST-2 and CoLA datasets under different rank budget scenarios for Sensitivity-LoRA and LoRA.

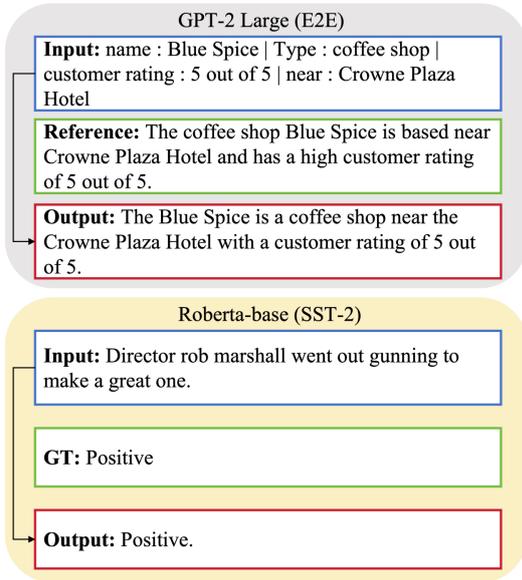


Figure 5: The Case Study of the GPT-2 Large and RoBERTa-base models. The blue boxes represent the input test data, the green boxes indicate the reference text or ground truth output, and the red boxes highlight the model's actual output.

the dynamic rank allocation method (Sensitivity-LoRA) on the E2E and SST-2 datasets. For the E2E dataset, the output of the GPT-2 Large model retains the input information while generating fluent and grammatically correct natural language text that closely aligns with the reference text. This suggests that the model effectively processes structured input data and excels in generating accurate and fluent natural language descriptions. For the SST-2 dataset, the RoBERTa-base model demonstrates strong performance in sentiment classification

tasks, accurately classifying input text as "Positive". This outcome validates the efficacy of the Sensitivity-LoRA method in improving model performance on text classification tasks.

## 5 Conclusion

In this work, we introduce Sensitivity-LoRA, a method that efficiently allocates ranks to weight matrices based on their sensitivity, without a significant computational burden. Sensitivity-LoRA first performs sensitivity detection by analyzing both global and local sensitivities. It utilizes the second-order derivatives (Hessian matrix) of the loss function to accurately capture parameter sensitivity. Next, it optimizes rank allocation by aggregating global and local sensitivities, ensuring a comprehensive and fair evaluation metric. Extensive experiments consistently demonstrate the efficiency, effectiveness and stability of our method.

## 6 Limitations

In this study, we opt for predefined parameters ( $\gamma_1, \gamma_2, etc.$ ) when integrating global metrics and local metrics. While our approach effectively guides the rank allocation process, it does not guarantee an optimal solution for fine-tuning. Although dynamically adjusting parameters through heuristic learning could be a potential solution, it may significantly extend training time and increase resource consumption. Given that our core design principle is to maintain a low computational load, we ultimately decide to use predefined parameters to ensure the efficiency and simplicity of the system. In the future, we aim to explore an effective

558	solution for dynamically modifying predefined parameters.	C. Li et al. 2023. Fine-tuning techniques for efficient model adaptation. <i>AI Research Journal</i> .	609
559			610
560	<b>References</b>	Chunyuan Li, Heerad Farkhoor, Rosanne Liu, and Jason Yosinski. 2018. Measuring the intrinsic dimension of objective landscapes. <i>arXiv preprint arXiv:1804.08838</i> .	611
561	Armen Aghajanyan, Luke Zettlemoyer, and Sonal Gupta. 2020. Intrinsic dimensionality explains the effectiveness of language model fine-tuning. <i>arXiv preprint arXiv:2012.13255</i> .		612
562			613
563		X. Li and P. Liang. 2021. Prefix-tuning: Optimizing continuous prompts for generation tasks. <i>Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics (ACL)</i> , 2021:4582–4597.	614
564			615
565	Satanjeev Banerjee and Alon Lavie. 2005. Meteor: An automatic metric for mt evaluation with improved correlation with human judgments. In <i>Proceedings of the acl workshop on intrinsic and extrinsic evaluation measures for machine translation and/or summarization</i> , pages 65–72.		616
566			617
567		Chin-Yew Lin. 2004. Rouge: A package for automatic evaluation of summaries. In <i>Text summarization branches out</i> , pages 74–81.	618
568			619
569		Chin-Yew Lin and Franz Josef Och. 2004. Automatic evaluation of machine translation quality using longest common subsequence and skip-bigram statistics. In <i>Proceedings of the 42nd annual meeting of the association for computational linguistics (ACL-04)</i> , pages 605–612.	620
570			621
571	Daniel Cer, Mona Diab, Eneko Agirre, Inigo Lopez-Gazpio, and Lucia Specia. 2017. Semeval-2017 task 1: Semantic textual similarity-multilingual and cross-lingual focused evaluation. <i>arXiv preprint arXiv:1708.00055</i> .		622
572			623
573			624
574			625
575			626
576	G. Ding et al. 2022. Efficient fine-tuning for resource-constrained systems. <i>Proceedings of the Machine Learning Conference</i> .		627
577			628
578			629
579	Bill Dolan and Chris Brockett. 2005. Automatically constructing a corpus of sentential paraphrases. In <i>Third international workshop on paraphrasing (IWP2005)</i> .		630
580			631
581			632
582			633
583	N. Houlisby, A. Giurgiu, S. Jastrzebski, et al. 2019. Parameter-efficient transfer learning for nlp. <i>Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing (EMNLP)</i> , 2019:279–285.		634
584			635
585			636
586			637
587			638
588	E. J. Hu, Y. Shen, P. Wallis, et al. 2022a. Lora: Low-rank adaptation of large language models. <i>International Conference on Learning Representations (ICLR)</i> .		639
589			640
590			641
591			642
592	Edward J Hu, Yelong Shen, Phillip Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang, and Weizhu Chen. 2022b. LoRA: Low-rank adaptation of large language models. In <i>International Conference on Learning Representations</i> .		643
593			644
594			645
595			646
596			647
597	Yahao Hu, Yifei Xie, Tianfeng Wang, Man Chen, and Zhisong Pan. 2023. Structure-aware low-rank adaptation for parameter-efficient fine-tuning. <i>Mathematics</i> , 11(20):4317.		648
598			649
599			650
600			651
601	E. Huang et al. 2023. Evaluating large language models in complex scenarios. <i>Journal of Computational Linguistics</i> .		652
602			653
603			654
604	B. Lester, R. Al-Rfou, and N. Constant. 2021. The power of scale for parameter-efficient prompt tuning. <i>Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing (EMNLP)</i> , 2021:3045–3061.		655
605			656
606			657
607			658
608			659
			660
			661
			662
			663
			664
			665
			666
			667
			668
			669
			670
			671
			672
			673
			674
			675
			676
			677
			678
			679
			680
			681
			682
			683
			684
			685
			686
			687
			688
			689
			690
			691
			692
			693
			694
			695
			696
			697
			698
			699
			700

662	Pranav Rajpurkar, Robin Jia, and Percy Liang. 2018.	B. Zhu et al. 2023a. Expanding frontiers in large lan-	717
663	Know what you don't know: Unanswerable questions	guage models. <i>AI Frontier Research</i> .	718
664	for squad. <i>arXiv preprint arXiv:1806.03822</i> .		
665	Richard Socher, Alex Perelygin, Jean Wu, Jason	B. Zhu et al. 2023b. Large language models: Progress	719
666	Chuang, Christopher D Manning, Andrew Y Ng, and	and applications. <i>Advances in NLP</i> .	720
667	Christopher Potts. 2013. Recursive deep models for		
668	semantic compositionality over a sentiment treebank.		
669	In <i>Proceedings of the 2013 conference on empirical</i>		
670	<i>methods in natural language processing</i> , pages		
671	1631–1642.		
672	Mojtaba Valipour, Mehdi Rezagholizadeh, Ivan		
673	Kobyzev, and Ali Ghodsi. 2022. Dylora: Parameter		
674	efficient tuning of pre-trained models using dynamic		
675	search-free low-rank adaptation. <i>arXiv preprint</i>		
676	<i>arXiv:2210.07558</i> .		
677	Ramakrishna Vedantam, C Lawrence Zitnick, and Devi		
678	Parikh. 2015. Cider: Consensus-based image de-		
679	scription evaluation. In <i>Proceedings of the IEEE</i>		
680	<i>conference on computer vision and pattern recogni-</i>		
681	<i>tion</i> , pages 4566–4575.		
682	Alex Wang. 2018. Glue: A multi-task benchmark and		
683	analysis platform for natural language understanding.		
684	<i>arXiv preprint arXiv:1804.07461</i> .		
685	F. Wang et al. 2023. Practical applications of llms in		
686	specialized domains. <i>Specialized AI Applications</i> .		
687	A Warstadt. 2019. Neural network acceptability judg-		
688	ments. <i>arXiv preprint arXiv:1805.12471</i> .		
689	Adina Williams, Nikita Nangia, and Samuel R Bow-		
690	man. 2017. A broad-coverage challenge corpus for		
691	sentence understanding through inference. <i>arXiv</i>		
692	<i>preprint arXiv:1704.05426</i> .		
693	Thomas Wolf. 2020. Transformers: State-of-the-		
694	art natural language processing. <i>arXiv preprint</i>		
695	<i>arXiv:1910.03771</i> .		
696	E. Zaken, Y. Goldberg, and S. Ravfogel. 2022. Bitfit:		
697	Simple parameter-efficient fine-tuning for transfor-		
698	mers. <i>Transactions of the Association for Computa-</i>		
699	<i>tional Linguistics (TACL)</i> , 10:1–16.		
700	D. Zhang et al. 2023a. Parameter-efficient fine-tuning		
701	methods for llms. <i>Journal of Machine Learning Re-</i>		
702	<i>search</i> .		
703	Feiyu Zhang, Liangzhi Li, Junhao Chen, Zhouqiang		
704	Jiang, Bowen Wang, and Yiming Qian. 2023b. In-		
705	crelora: Incremental parameter allocation method		
706	for parameter-efficient fine-tuning. <i>arXiv preprint</i>		
707	<i>arXiv:2308.12043</i> .		
708	Qingru Zhang, Minshuo Chen, Alexander Bukharin,		
709	Nikos Karampatziakis, Pengcheng He, Yu Cheng,		
710	Weizhu Chen, and Tuo Zhao. 2023c. Adalora: Adap-		
711	tive budget allocation for parameter-efficient fine-		
712	tuning. <i>arXiv preprint arXiv:2303.10512</i> .		
713	Ruiyi Zhang, Rushi Qiang, Sai Ashish Somayajula, and		
714	Pengtao Xie. 2024. Autolora: Automatically tuning		
715	matrix ranks in low-rank adaptation based on meta		
716	learning. <i>arXiv preprint arXiv:2403.09113</i> .		