

Balancing Performance and Efficiency in Large Language Model Fine-Tuning through Hierarchical Freezing

Jinxu Guo

Dartmouth College, Hanover, USA

jinxuguo2c@gmail.com

Abstract: This study investigates the efficiency of fine-tuning in large language models and proposes an optimization method based on hierarchical parameter freezing. The method divides model parameters into three levels: lower, middle, and upper. It adopts freezing, partial updating, and full updating strategies for these levels, respectively, to balance stability and adaptability. The lower-level parameters preserve general syntax and semantic knowledge. The middle-level parameters are flexibly controlled according to task complexity. The upper-level parameters focus on task-specific semantic modeling. In this way, the method reduces computational and storage costs significantly while maintaining performance. To verify the effectiveness of the method, systematic comparison experiments were conducted. Multiple metrics, including ROUGE, BLEU, and EM, were evaluated. The results show that the proposed method achieves a better balance between performance and efficiency. In addition, sensitivity experiments were carried out in three dimensions: hyperparameters, environmental settings, and data conditions. The analysis covered learning rate, sequence length, training data ratio, and text noise. The findings further demonstrate the robustness of the method under diverse conditions. By combining hierarchical freezing mechanisms with parameter updating strategies, this study provides a new approach for efficient use of large language models in resource-constrained environments. It also confirms the broad applicability of the method in real-world tasks.

Keywords: Hierarchical parameter freezing; efficient fine-tuning; large language model; sensitivity analysis

1. Introduction

In the accelerating wave of intelligent development, large language models have become a core infrastructure in artificial intelligence due to their strong capabilities in knowledge representation and reasoning. However, as the scale of these models continues to expand, the computational resources, storage costs, and energy consumption required for training and deployment have increased significantly[1,2]. This not only limits their broad application in academia and industry but also raises challenges for green computing and sustainable development. Therefore, how to reduce resource consumption while maintaining performance has become a pressing problem. In this context, parameter-efficient fine-tuning methods have gained attention. By enabling knowledge transfer and adaptation with a limited number of parameters, they demonstrate both practical value and theoretical significance[3].

Traditional full-parameter fine-tuning requires updating the entire set of large parameter matrices inside the model. This leads to significant storage and computational overhead, and also makes the model prone to overfitting and forgetting in new tasks. In contrast, parameter freezing strategies fix part of the parameters

during training while updating only specific layers or additional modules. This reduces computational demand and improves training stability. Yet, a single-level freezing strategy is difficult to balance efficiency and performance. Freezing too many parameters may harm expressive power, while freezing too few cannot effectively reduce resource consumption. Designing a hierarchical and controllable freezing mechanism that allows efficient knowledge transfer at different levels has become a key breakthrough for improving fine-tuning efficiency[4].

The hierarchical parameter freezing method was proposed to address this challenge. Its core idea lies in exploiting the semantic differences across layers of large language models. Lower layers capture general features and syntactic knowledge. Middle layers focus on contextual dependencies and structural modeling. Higher layers represent task-specific semantics. By applying differentiated freezing and updating strategies to different layers, it is possible to stabilize general knowledge while maximizing task adaptation. This layered mechanism reduces redundant updates and unnecessary computation, while also ensuring both performance and efficiency in resource-constrained environments.

Beyond optimizing resource use, the significance of hierarchical parameter freezing also lies in its ability to improve the universality and adaptability of artificial intelligence across diverse application scenarios. Large language models are increasingly applied in education, healthcare, finance, and public services. These tasks have very different requirements, and full-parameter fine-tuning is often insufficient. Hierarchical freezing enables flexible control of training depth and parameter scale, allowing personalized adjustments according to task complexity and resource availability. This approach promotes efficient transfer across domains and tasks. It not only improves the operability of models in real-world deployment but also provides a technical path for broader and fairer access to AI applications.

In summary, the study of efficient fine-tuning for large language models based on hierarchical parameter freezing responds to the current challenges of computational bottlenecks and resource pressure. At the same time, it offers a new solution for lightweight and sustainable model development. This line of research has both theoretical importance and practical value. On one hand, it alleviates the difficulties faced by institutions with limited resources. On the other hand, it helps extend large language models into broader fields, allowing them to play a more enduring and profound role in building an intelligent society.

2. Related work

In recent years, the rapid development of large language models has given rise to various parameter-efficient fine-tuning methods[5]. These approaches aim to address the high cost of computation, storage, and energy consumption associated with traditional full-parameter fine-tuning. Researchers have proposed techniques such as adapter insertion, low-rank decomposition, and prompt tuning. By updating only a small portion of additional parameters or imposing structured constraints on existing ones, these methods greatly reduce the resources needed for fine-tuning[6]. At the same time, they maintain model performance and enable flexible knowledge transfer, providing new feasibility for cross-task and cross-domain applications. However, these methods often rely on additional structures or specialized parameterization forms. This reliance may lead to high implementation complexity and strong coupling with the underlying architecture, which limits their broader deployment.

To balance full-parameter freezing and partial updating, parameter freezing strategies have been considered as an intuitive and effective solution. By fixing most of the parameters during fine-tuning and updating only a few critical or task-related layers, such strategies can significantly reduce computational requirements and mitigate catastrophic forgetting. Yet, traditional freezing approaches usually adopt a single-level scheme without recognizing the differences in knowledge representation across layers. This coarse-grained design often struggles to balance efficiency and performance. Excessive freezing reduces adaptability to new tasks, while insufficient freezing diminishes the benefits of cost reduction. Thus, refining the layer-wise design of parameter freezing has become a critical bottleneck in this field[7].

Against this backdrop, hierarchical parameter freezing has emerged as a new research direction. This method highlights the functional differences of model layers in semantic representation. It applies differentiated freezing and updating strategies to lower, middle, and upper layers, achieving layered management of general and task-specific knowledge. Lower layers usually capture syntax and general representations. Middle layers emphasize contextual dependencies and structural modeling. Upper layers align more closely with task-related semantics. Hierarchical parameter freezing stabilizes general knowledge while enhancing adaptability to new tasks. It avoids redundant computation from full updates and overcomes the limitations of coarse-grained freezing. This approach improves resource efficiency in training and provides a solid foundation for deploying large models in resource-constrained environments.

At the same time, the demand for parameter efficiency has become more urgent with the widespread adoption of large language models across domains. Hierarchical parameter freezing, as a scalable and flexible mechanism, offers a new paradigm for rapid transfer across multiple tasks and fields. Compared with other parameter-efficient fine-tuning methods, it reduces resource costs without relying on additional structures. It shows stronger generality and practicality. Future trends suggest that hierarchical freezing can be combined with existing techniques, such as adapters or low-rank decomposition, to further improve fine-tuning efficiency. It may also become a key foundation for building green intelligent systems and promoting more accessible artificial intelligence.

3. Proposed Approach

In this study, the core idea of the method is to achieve efficient fine-tuning of large language models through hierarchical parameter freezing. The overall model architecture is shown in Figure 1.

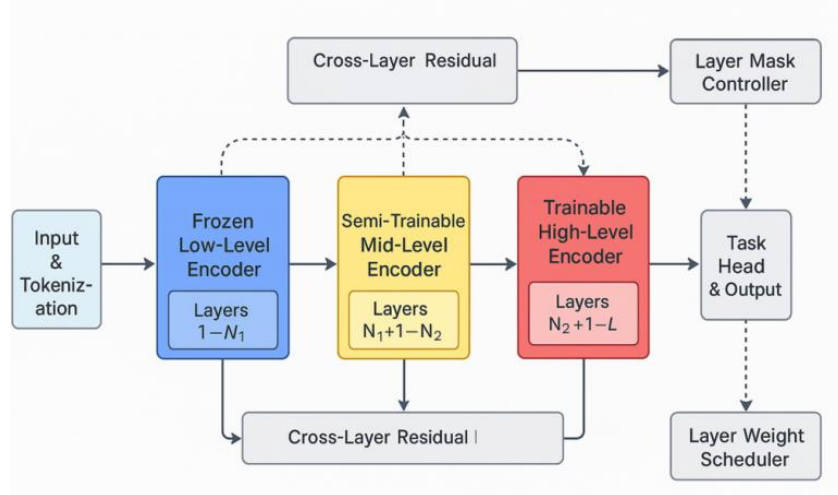


Figure 1. Overall model architecture

Suppose the parameter set of the original model is:

$$\theta = \{\theta_1, \theta_2, \dots, \theta_L\}$$

Where L represents the number of layers of the model and θ_l is the parameter of the l th layer. The goal of hierarchical parameter freezing is to divide these layers into different subsets according to semantic functions, such as the underlying parameter set θ_{low} , the middle-level parameter set θ_{mid} , and the high-level parameter set θ_{high} , to realize differentiated update strategies during the training process. Specifically, the underlying layer focuses on maintaining stability, while the middle and high layers perform selective updates according to task needs to achieve efficient knowledge migration and adaptation.

In the design of the update mechanism, assuming the training data is $D = \{(x_i, y_i)\}_{i=1}^N$, its objective function can be expressed as:

$$L(\theta) = \frac{1}{N} \sum_{i=1}^N l(f(x_i; \theta); y_i)$$

Where $l(\cdot)$ is the loss function and $f(x_i; \theta)$ represents the prediction result of the model. To achieve hierarchical parameter freezing, we introduce a mask function $m \in \{0,1\}$ to control whether the l th layer parameters participate in the update. The parameter update rule can then be expressed as:

$$\theta_l^{(t+1)} = \theta_l^{(t)} - \eta m_l \nabla_{\theta_l} L(\theta)$$

Where η is the learning rate. When $m_l = 0$, the parameters of this layer remain frozen and do not participate in backpropagation update.

To further balance stability and adaptability, we introduce a hierarchical weight coefficient λ_l , which is used to explicitly adjust the update importance of different levels in the objective function. The adjusted loss function is:

$$L_{layered}(\theta) = \frac{1}{N} \sum_{i=1}^N \sum_{l=1}^L \lambda_l l(f_l(x_i; \theta_l), y_i)$$

Where $f_l(x_i; \theta_l)$ represents the intermediate representation generated by the l -th parameter. By rationally setting λ_l , hierarchical dynamic control can be achieved during the training process, so that the model can maintain the stability of general representation in different task scenarios, and strengthen the task-related expression ability.

In addition, considering that the freezing layer may lead to limited information transfer, we construct a residual parameter interaction mechanism inside the model. Specifically, the cross-layer residual connection is defined as:

$$h_l = f_l(h_{l-1}; \theta_l) + \alpha h_{l-k}$$

Where h_l represents the output of the l th layer, α is the residual balance factor, and k is the cross-layer distance. This design can ensure sufficient flow of gradients and information while keeping some layers frozen, thereby improving the stability and convergence of the overall training. To sum up, this method provides an efficient and flexible fine-tuning path for large language models through the organic combination of hierarchical parameter freezing, mask control, weight adjustment, and residual interaction.

4. Performance Evaluation

4.1 Dataset

The dataset used in this study is The Pile, a large-scale open corpus designed for the training and evaluation of large language models. It contains heterogeneous text from multiple sources, including academic articles, news reports, web content, code snippets, and books. The corpus covers a wide range of semantic domains and writing styles. Its size reaches several hundred gigabytes. The sources are diverse and of high quality, providing the model with rich linguistic knowledge and contextual information. This makes it suitable for evaluating fine-tuning methods in terms of adaptability and generalization across diverse contexts.

The construction of the dataset followed a unified cleaning and deduplication process to ensure consistency and completeness. All text underwent preprocessing, including character encoding normalization, removal

of noisy content, and formatting standardization. These steps improved the overall quality and usability of the data. In addition, the dataset was carefully designed to cover both general-purpose texts and domain-specific professional materials. This guarantees adaptability to different task environments.

In this study, the diversity and scale of The Pile provide a strong foundation for validating the hierarchical parameter freezing method. Its content supports stable modeling of syntax and general semantics at lower layers, while also offering abundant material for semantic adaptation at higher layers. By using this dataset, it is possible to better explore the potential and advantages of different freezing strategies in large-scale language knowledge transfer.

4.2 Experimental Results

This paper first conducts a comparative experiment, and the experimental results are shown in Table 1.

Table 1: Comparative experimental results

Method	ROUGE-1	ROUGE-L	EM	BLEU
SplitLoRA[8]	42.8	39.5	36.2	25.7
Lisa[9]	44.1	40.3	37.5	26.9
Mixlora[10]	45.6	41.7	38.8	28.1
Flora[11]	46.3	42.5	39.6	28.8
Ours	49.7	45.2	42.4	31.5

From the results in Table 1, it can be seen that the hierarchical parameter freezing method outperforms existing approaches across all evaluation metrics. Compared with traditional parameter-efficient fine-tuning methods, it achieves significant improvements on ROUGE-1 and ROUGE-L, reaching 49.7 and 45.2, respectively. This indicates that in text generation tasks, the proposed method captures semantic information and contextual dependencies more effectively, leading to outputs that are more consistent with the reference answers. Through a well-designed hierarchical freezing strategy, the model preserves stability in low-level syntax and general knowledge while enhancing task-specific semantic modeling at higher layers.

For the EM metric, the proposed method reaches 42.4, which shows a clear advantage over other approaches. This demonstrates that hierarchical parameter freezing not only improves the fluency and semantic coherence of generated text but also enhances accuracy and consistency with the gold standard. Traditional approaches often face a trade-off when freezing too many or too few layers. In contrast, the hierarchical strategy achieves a better balance through fine-grained parameter control, which significantly improves the precision of model outputs.

The improvement in BLEU further confirms the effectiveness of this method in capturing diverse expressions and maintaining translation consistency. The hierarchical freezing mechanism enables differentiated roles across layers, allowing the model to generate semantically equivalent yet diverse outputs. This reduces redundancy and repetitive patterns. Such capability is especially important for real-world applications of large language models, as different tasks and scenarios often require greater flexibility in expression and stronger semantic adaptability.

Overall, the experimental results demonstrate the potential and advantages of hierarchical parameter freezing for efficient fine-tuning. It not only reduces resource consumption effectively but also achieves comprehensive performance gains over comparison methods. This mechanism provides a practical path to address the computational bottlenecks and generalization challenges of large models, while also laying a solid foundation for lightweight and universal applications of large language models.

This paper also gives experimental results for the sensitivity of learning rate to experimental results, and the experimental results are shown in Figure 2.

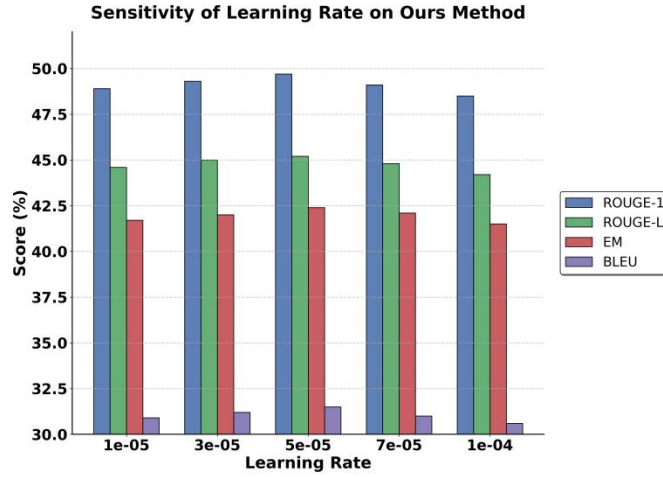


Figure 2. The sensitivity of the learning rate to experimental results

From the figure, it can be observed that under different learning rate settings, the model shows some fluctuations across metrics, but the overall trend remains relatively stable. In particular, around 5×10^{-5} , the values of ROUGE-1, ROUGE-L, EM, and BLEU reach higher levels. This indicates that the hierarchical parameter freezing method achieves a better balance between adaptability and stability at this learning rate. Compared with very low or very high learning rates, a medium learning rate balances parameter update speed and gradient convergence, resulting in better overall performance.

When the learning rate is set to a lower value, such as 1×10^{-5} , the metrics remain at reasonable levels but are slightly below the optimal. This reflects that a too small learning rate causes parameter updates to be too slow, preventing the model from fully adapting to task features within limited iterations. Although low-level parameters remain stable due to freezing, task-related parameters at higher levels are not sufficiently optimized, leading to some performance gaps. This shows that relying on an extremely small learning rate cannot maximize the potential of the hierarchical freezing strategy.

When the learning rate increases to 1×10^{-4} , all four metrics decline to some extent, with BLEU and ROUGE-L showing more obvious drops. This suggests that an excessively high learning rate causes the updates of trainable parameters at higher layers to be too aggressive, which disrupts the semantic consistency of the middle and upper layers. Since the hierarchical freezing method emphasizes both stability and adaptability, a too large learning rate breaks this balance, leading to reduced fluency and accuracy in the generated text.

Overall, these results confirm the learning rate sensitivity pattern of the hierarchical parameter freezing method. A moderate learning rate allows the model to maintain stability in lower layers while effectively enhancing adaptability in higher layers. As a result, the model achieves improvements across multiple metrics. This finding highlights not only the importance of reasonable hyperparameter selection but also the robustness and practical value of the proposed method in efficient fine-tuning.

This paper also gives the impact of the scale ratio of the training data on the experimental results, and the experimental results are shown in Figure 3.

From the figure, it can be seen that as the proportion of training data increases, the model shows a stable upward trend across all metrics. The improvements in ROUGE-1 and ROUGE-L are particularly significant. When the data proportion increases from 10% to 100%, both metrics improve markedly. This indicates that

the hierarchical parameter freezing method can make full use of larger training data to enhance semantic modeling ability, thereby producing outputs closer to the reference texts in generation tasks.

For the EM metric, the exact match rate gradually increases as the data proportion grows. This shows that more data supports better optimization of task-related parameters at higher layers. Since the lower and middle layers are kept frozen or partially updated, model stability is maintained. At the same time, the higher layers are continuously optimized with large-scale samples, improving consistency between predictions and gold answers. This trend also reflects that the hierarchical freezing strategy can effectively sustain performance gains when facing data expansion.

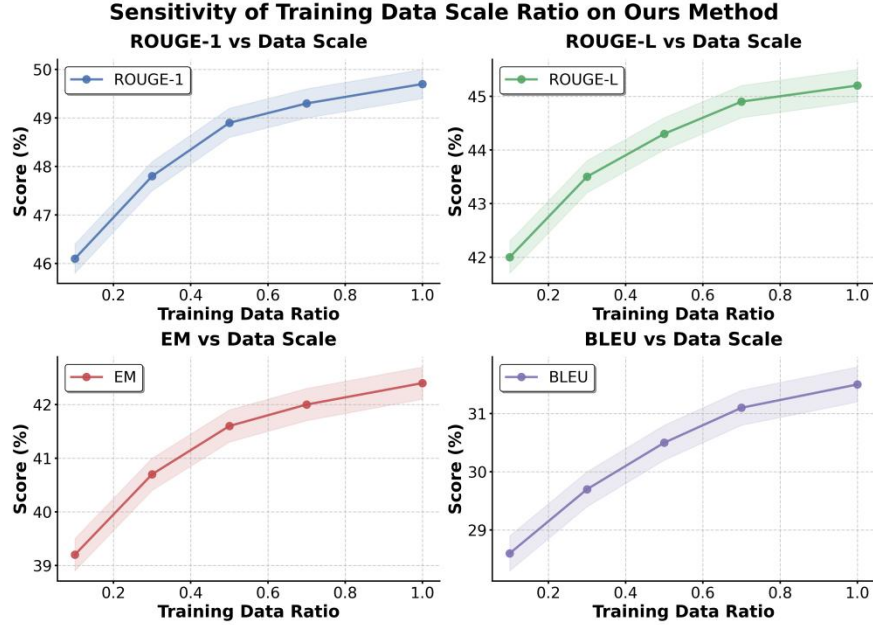


Figure 3. Effect of training data scale ratio on experimental results

The growth of the BLEU metric highlights the role of data scale in improving diversity of expression and translation consistency. With a larger proportion of training data, the model is exposed to richer structures and expression styles. This allows it to generate more natural text while maintaining semantic equivalence. The hierarchical parameter freezing mechanism secures low-level syntactic knowledge in this process, enabling more flexible adaptation at higher levels. As a result, the quality and diversity of generated outputs are improved.

Overall, these experimental results clearly demonstrate the importance of data scale for the hierarchical parameter freezing method. Increasing the proportion of training data not only improves overall metric performance but also highlights the robustness and scalability of the method in large-scale corpus settings. This shows that the strategy has strong potential for practical applications. It can leverage accumulated data resources to further enhance performance and provide reliable support for efficient fine-tuning and broad deployment of large language models.

This paper also gives the effect of the maximum length of the sequence on the experimental results, and the experimental results are shown in Figure 4.

From the figure, it can be seen that as the maximum sequence length increases, the model performance across all metrics first improves and then stabilizes. Around the length of 512, ROUGE-1 and ROUGE-L both reach higher values. This indicates that under this configuration, the model captures global semantic dependencies more effectively and shows stronger semantic consistency in generation tasks. Compared with shorter

sequences, longer contexts help the model, under the hierarchical freezing strategy, to make better use of contextual information, thereby enhancing overall modeling ability.

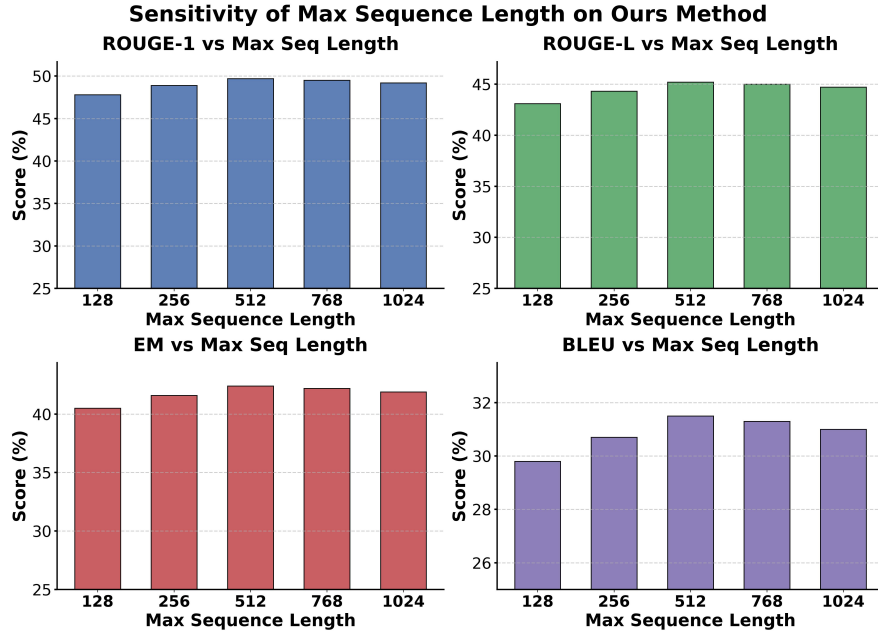


Figure 4. Effect of maximum sequence length on experimental results

For the EM metric, as the maximum sequence length increases from 128 to 512, the exact match rate improves continuously. This shows that task-related parameters are more fully optimized with longer sequence inputs. This phenomenon is closely related to the design of hierarchical parameter freezing. Low-level parameters remain stable, while middle and upper layers learn more detailed semantic features in the extended input space, improving the accuracy of predictions. After 512, the gains in EM begin to level off, suggesting that excessively long inputs do not provide significant additional benefits.

The performance of BLEU further illustrates that an appropriate sequence length helps the model generate fluent and diverse expressions while maintaining semantic consistency. When the sequence length is too short, the model tends to ignore long-range dependencies, leading to fragmented outputs. When the sequence length is too long, redundant information may dilute useful context, causing dispersed attention and limited improvement. The range between 512 and 768 shows relatively balanced performance, reflecting a reasonable trade-off between context modeling and redundancy control.

Overall, these results show that the maximum sequence length has a significant impact on the hierarchical parameter freezing method. A proper configuration allows full use of the complementarity between frozen and trainable layers. It ensures the stability of general representations while strengthening task-specific semantic modeling. In practical applications, choosing an appropriate sequence length not only improves performance but also avoids unnecessary computational costs, further highlighting the practical value of this method for efficient fine-tuning.

This paper further gives the effect of text noise on experimental results, and the experimental results are shown in Figure 5.

From the figure, it can be seen that as the proportion of text noise increases, the model performance on all metrics shows a clear downward trend. In particular, for ROUGE-1 and ROUGE-L, when the noise ratio increases from 0 to 0.3, the scores drop by about 5 points. This indicates that semantic and contextual consistency are significantly affected by noise interference. The result shows that although hierarchical

parameter freezing can maintain stability at lower layers to some extent, excessive noise still weakens the task adaptation ability of higher-level parameters.

For the EM metric, the increase in noise ratio also leads to a gradual performance decline. Exact match, as a strict standard measuring complete consistency between model outputs and reference answers, is highly sensitive to input disturbances. When text contains a high proportion of spelling errors or character perturbations, the model struggles to maintain stable exact predictions, resulting in a continuous drop in EM scores. This suggests that under noisy conditions, task-related parameters at higher layers fail to extract effective features from a limited context, and model robustness is restricted.

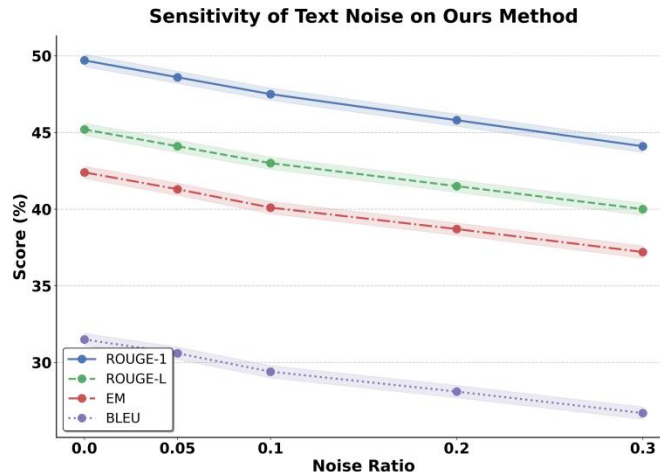


Figure 5. Effect of text noise on experimental results

The decline in BLEU further highlights the impact of text noise on fluency and diversity in generation. As the noise ratio increases, the model’s ability to maintain semantically equivalent expressions weakens. Generated sentences become more prone to local incoherence or redundancy. This also reflects that while hierarchical freezing can preserve some robustness under perturbations, additional regularization or denoising strategies are needed to effectively mitigate the negative effects of noise on higher-level modeling.

Overall, these results clearly reveal the negative impact of text noise on the hierarchical parameter freezing method. A moderate level of noise may still be tolerated by the model, but when the proportion becomes too high, both semantic modeling and task-related representations are weakened, leading to overall performance decline. This indicates that in practical applications, data quality control and preprocessing mechanisms are essential to ensure that the method can achieve efficient and robust performance in low-noise environments.

5. Conclusion

This study addresses the challenges of efficient fine-tuning in large language models and proposes an optimization method based on hierarchical parameter freezing. By managing parameters at different layers in a differentiated way, the lower layers maintain stability of general representations, while the middle and upper layers are flexibly adjusted according to task requirements. This effectively alleviates the limitations of full-parameter fine-tuning in terms of computational cost and storage burden. Experimental results show that the method outperforms existing approaches across multiple evaluation metrics while significantly reducing resource consumption. It provides a new path for adapting and deploying large language models in resource-constrained environments.

The study further verifies the robustness of the hierarchical freezing strategy in handling sensitivity to hyperparameters, environmental constraints, and data perturbations. Under different settings of learning rate, sequence length, and training data scale, the method maintains relatively balanced performance and avoids the fluctuations often caused by instability in traditional approaches. In addition, when faced with uncertainty

such as text noise, the hierarchical freezing mechanism still demonstrates strong resistance to interference. This highlights its potential to ensure stable outputs in complex application scenarios. Such robustness enhances the theoretical value of the method and provides a reliable foundation for practical implementation.

From the application perspective, the proposed fine-tuning strategy has important significance. The use of large language models in education, healthcare, finance, and government is expanding rapidly. However, limited computational and storage resources make traditional full-parameter fine-tuning impractical. Hierarchical parameter freezing reduces resource consumption significantly, enabling more small and medium-sized institutions to apply large models under restricted hardware conditions. This supports the wider adoption of intelligent services. It not only improves accessibility but also promotes fairness in diverse tasks, helping to address the imbalance caused by concentrated computing power.

Overall, this study presents an efficient and scalable solution at the methodological level and demonstrates practical value for the democratization of artificial intelligence. The fine-tuning method based on hierarchical parameter freezing balances performance and efficiency, offering a new direction for lightweight and sustainable development of large language models. In the future, as demand for large model capabilities continues to grow, this approach is expected to have a deeper impact on cross-domain knowledge transfer, multi-task collaborative modeling, and large-scale deployment. It will provide new momentum for the sustained development of artificial intelligence.

References

- [1] Hu E J, Shen Y, Wallis P, et al. Lora: Low-rank adaptation of large language models[J]. ICLR, 2022, 1(2): 3.
- [2] Liu X, Ji K, Fu Y, et al. P-tuning v2: Prompt tuning can be comparable to fine-tuning universally across scales and tasks[J]. arXiv preprint arXiv:2110.07602, 2021.
- [3] Zhang R, Han J, Liu C, et al. LLaMA-adapter: Efficient fine-tuning of large language models with zero-initialized attention[C]//The Twelfth International Conference on Learning Representations. 2024.
- [4] Dettmers T, Pagnoni A, Holtzman A, et al. Qlora: Efficient finetuning of quantized llms[J]. Advances in neural information processing systems, 2023, 36: 10088-10115.
- [5] Zhang Q, Chen M, Bukharin A, et al. Adalora: Adaptive budget allocation for parameter-efficient fine-tuning[J]. arXiv preprint arXiv:2303.10512, 2023.
- [6] Liu S Y, Wang C Y, Yin H, et al. Dora: Weight-decomposed low-rank adaptation[C]//Forty-first International Conference on Machine Learning. 2024.
- [7] Hayou S, Ghosh N, Yu B. Lora+: Efficient low rank adaptation of large models[J]. arXiv preprint arXiv:2402.12354, 2024.
- [8] Lin Z, Hu X, Zhang Y, et al. Splitlora: A split parameter-efficient fine-tuning framework for large language models[J]. arXiv preprint arXiv:2407.00952, 2024.
- [9] Pan R, Liu X, Diao S, et al. Lisa: Layerwise importance sampling for memory-efficient large language model fine-tuning[J]. Advances in Neural Information Processing Systems, 2024, 37: 57018-57049.
- [10] Li D, Ma Y, Wang N, et al. Mixlora: Enhancing large language models fine-tuning with lora-based mixture of experts[J]. arXiv preprint arXiv:2404.15159, 2024.
- [11] Honovich, O., Scialom, T., Levy, O. and Schick, T., "Unnatural Instructions: Tuning Language Models with (Almost) No Human Labor", arXiv preprint arXiv:2212.09689, 2022.