

MoSLD: An Extremely Parameter-Efficient Mixture-of-Shared LoRAs for Multi-Task Learning

Lulu Zhao^{1*}, Weihao Zeng², Xiaofeng Shi¹, Hua Zhou¹

¹Beijing Academy of Artificial Intelligence (BAAI)

²School of Artificial Intelligence, Beijing University of Posts and Telecommunications
llzhao@baai.ac.cn

Abstract

Recently, LoRA has emerged as a crucial technique for fine-tuning large pre-trained models, yet its performance in multi-task learning scenarios often falls short. In contrast, the MoE architecture presents a natural solution to this issue. However, it introduces challenges such as mutual interference of data across multiple domains and knowledge forgetting of various tasks. Additionally, MoE significantly increases the number of parameters, posing a computational cost challenge. Therefore, in this paper, we propose MoSLD, a mixture-of-shared-LoRAs model with a dropout strategy. MoSLD addresses these challenges by sharing the upper projection matrix in LoRA among different experts, encouraging the model to learn general knowledge across tasks, while still allowing the lower projection matrix to focus on the unique features of each task. The application of dropout alleviates the imbalanced update of parameter matrix and mitigates parameter overfitting in LoRA. Extensive experiments demonstrate that our model exhibits excellent performance in both single-task and multi-task scenarios, with robust out-of-domain generalization capabilities.

1 Introduction

The emergence of Large Language Models (LLMs) has significantly advanced Natural Language Processing (NLP) technology, serving as a robust foundation with broad applicability (Touvron et al., 2023a,b; Ouyang et al., 2022). However, as the parameter scale increases, the process of full parameter fine-tuning (FP-tuning) demands substantial computational and memory resources. To strike a balance between resource requirements and effectiveness, the research community is increasingly turning to parameter-efficient fine-tuning (PEFT) methods (Zhao et al., 2022a; Zeng et al., 2023),

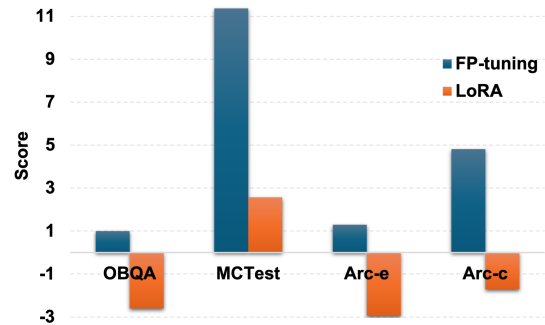


Figure 1: The increase between mixture setting and single setting for FP-tuning and LoRA on four datasets. The vertical axis is Score (mixture)-Score (single).

with LoRA emerging as the most prevalent and effective choice. Nevertheless, training an LLM via LoRA with multi-faceted capabilities faces significant challenges due to the differences and diversity inherent in various tasks. Figure 1 illustrates that while FP-tuning demonstrates competitive performance in a multi-task mixed training data setting, plain LoRA exhibits a drop. This decline underscores the challenge posed by the heterogeneity and imbalance in training data, resulting in interference between data from different tasks and consequently degrading the performance of plain LoRA on in-domain tasks. In essence, plain LoRA proves highly sensitive to the configuration of training data.

As we all know, MoE (Shazeer et al., 2017) has demonstrated remarkable advantages in amalgamating multiple capabilities. Particularly, the integration of MoE and LoRA (Hu et al., 2022) stands out as a promising approach to leveraging MoE in a parameter-efficient manner. This method preserves domain knowledge while significantly reducing training costs by introducing a limited number of domain-specific parameters (Dou et al., 2024; Luo et al., 2024; Liu et al., 2023). Presently, several works are devoted to applying MoE to LoRA. Some directly combine trained LoRAs linearly (Zhang et al., 2023; Huang et al., 2024), while others ap-

* Lulu Zhao is the corresponding author.

ply combinations of MoE and LoRA to different backbones (Chen et al., 2024; Dou et al., 2024). Another approach involves training a LoRA module for each distinct task type and employing a routing mechanism to integrate the LoRA modules under a shared LLM (Feng et al., 2024). However, we contend that these methods inadequately address the issue of data conflicts across different domains during LoRA training. Three primary challenges emerge: (1) The MoE architecture emphasizes the unique attributes of each LoRA and overlooks the transfer of general knowledge between different LoRAs, thereby impeding cross-task generalization in LLMs; (2) Requires a large number of trained LoRA modules (for each task); (3) Multiple LoRAs escalate the number of parameters and computational costs.

To solve these issues, in this paper, we propose a parameter-sharing method applied to the mixture-of-LoRAs, called MoSLD. The plain LoRA module comprises the upper projection matrix (A) and the lower projection matrix (B), which can be viewed as naturally decoupled general-feature and specific-feature matrices, respectively. Building upon the classic MoE architecture, we enable all experts at each layer to share a general-feature matrix while retaining the specific-feature matrix of each expert. This approach compels the model to capture shared general knowledge across various tasks to the fullest extent. The shared operation notably reduces the parameters of the MoE architecture, aligning with findings indicating parameter redundancy among experts (Fedus et al., 2022b; Kim et al., 2021). Despite the majority of parameters in the LoRA module being shared, differences can still be learned in each expert’s specific-feature matrix due to the tight coupling between the general and specific features. We posit that this mechanism can adaptively generalize to any new task. Furthermore, recognizing that the general-feature matrix is updated more frequently than the specific-feature matrix during fine-tuning, and overfitting tends to occur in LoRA (Wang et al., 2024), we apply the dropout strategy to the general-feature matrix, that is some weight values are randomly set to zero during training. This approach helps balance the updates between the general-feature and specific-feature matrices. Consequently, it not only facilitates a more balanced information exchange between different experts but also mitigates issues related to parameter redundancy and optimization imbalance.

In summary, our contributions are as follows: (1) We introduce a parameter-efficient MoSLD approach that disentangles domain knowledge and captures general knowledge by sharing a general-feature matrix, thus mitigating interference between heterogeneous datasets. (2) We implement a dropout strategy on the general-feature matrix to effectively mitigate overfitting and address the imbalance in directly optimizing MoE. (3) We conduct extensive experiments on various benchmarks to validate the effectiveness of our methods. Additionally, our approach demonstrates superior generalization to out-of-domain data.

2 Related Work

2.1 Mixture-of-Expert

The Mixture of Experts (MoE) functions as an ensemble method, conceptualized as a collection of sub-modules or experts, each tailored to process distinct types of input data. Guided by a router, each expert is selectively activated based on the input data type. This technique has garnered increasing attention and demonstrated remarkable performance across various domains, including computer vision, speech recognition, and multimodal applications (Fedus et al., 2022a). Evolution of MoE techniques spans from early sample-level approaches (Jacobs et al., 1991) to contemporary token-level implementations (Shazeer et al., 2017; Riquelme et al., 2021), which have now become mainstream. Concurrently, some researchers (Zhou et al., 2022; Chi et al., 2022) are delving into the router selection problem within MoE. Notably, the majority of these endeavors aim to scale up model parameters while mitigating computational costs.

2.2 Mixture-of-LoRA

As LoRA gradually becomes the most common parameter-efficient fine-tuning method, researchers pay more attention to combining MoE and LoRA for more efficient and effective model tuning. Huang et al. (2024) and Feng et al. (2024) pioneer the approach of training several LoRA weights on upstream tasks and then integrating the LoRA modules into a shared LLM using a routing mechanism. However, these methods necessitate the training of numerous pre-defined LoRA modules. Chen et al. (2024) initially engage in instruction fine-tuning through sparse mixing of LoRA experts in the multi-modal domain, while Dou et al. (2024) split the LoRA experts into two groups to explicitly

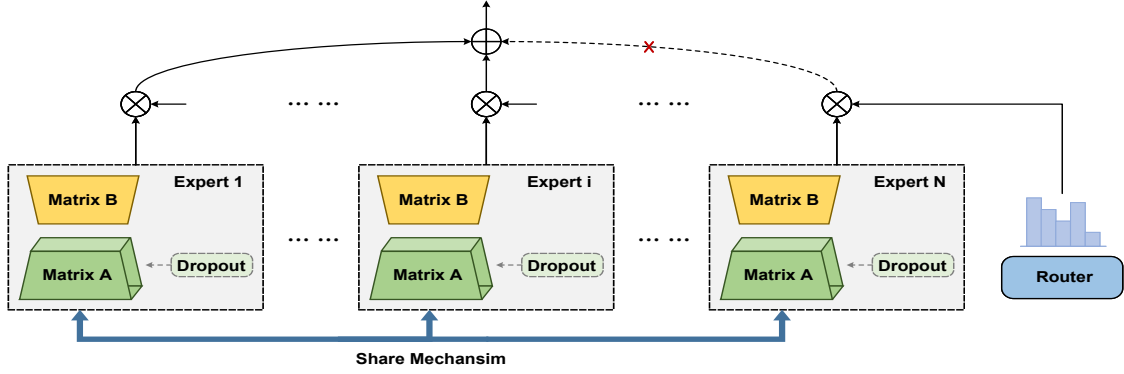


Figure 2: Overview of the share mechanism and dropout strategy in our MoSLD. Noted that the matrix A is shared among all experts in each layer.

learn different capabilities for each group. These mixture-of-LoRA methods typically involve predefined hyperparameters that require careful selection, and they densely mix LoRA experts, significantly increasing computational costs. To tackle overfitting resulting from an excessive number of experts, Gao et al. (2024) allocate a varying number of experts to each layer. Wu et al. (2024) propose MOLE, treating each layer of trained LoRAs as a distinct expert and implementing hierarchical weight control through a learnable gating function within each layer to tailor composition weights specific to a given domain’s objectives. However, these approaches overlook the issue of data conflicts across different datasets during LoRA training. As our concurrent work, MixLoRA (Li et al., 2024) also focuses on multi-task learning, which fuses multiple LoRAs with the shared FFN layer and employ a plain LoRA on the self-attention layer. We believe this method will introduce a large number of additional trainable parameters and incur a huge computational cost. In our study, we conduct extensive experimental analysis for both single and mixture data settings in a more lightweight way.

3 Methodology

In this section, we describe our MoSLD from the sharing mechanism, dropout strategy and optimization details, as shown in Figure 2.

3.1 Sharing Mechanism of LoRAs

In the area of parameter-efficient fine-tuning, LoRA introduces the concept of training only two low-rank matrices as an alternative to dense layer updates. In other words, it reformulates the parameter fine-tuning process in LLMs as a low-rank decomposition. Specifically, the equation $W_0 + \Delta W = W_0 + BA$ captures this decom-

position. Here, $W_0 \in \mathcal{R}^{d_{in} \times d_{out}}$ represents the parameter matrix of the pre-trained LLM, while $\Delta W \in \mathcal{R}^{d_{in} \times d_{out}}$ denotes the matrix updated during fine-tuning. The matrices $B \in \mathcal{R}^{d_{in} \times r}$ and $A \in \mathcal{R}^{r \times d_{out}}$ are low-rank and trainable.

In order to achieve the transfer of general features between different tasks and capture the shared general knowledge, we design a novel sharing mechanism. Specifically, given a Transformer model with L layers, we allocate N_l experts for layer l and create N_l pairs of low-rank matrices $\{A_{i,l}, B_{i,l}\}_{i=1}^{N_l}$, where $A_{i,l}$ is initialized from a random Gaussian distribution and each $B_{i,l}$ is set to zero. It is worth noting that the matrix $A_{i,l}$ is shared among all experts in each layer, i.e., $A_{1,l} = A_{2,l} \dots = A_{N_l,l}$ ($l \in L$). In other words, the core idea is to share the matrix A as the general-feature matrix and keep matrix B as specific-feature matrix. In this way, we can only keep L central general-feature matrices for a L -layer MoE architecture, which significantly reduces the parameters of the MoE architecture. A router with a trainable weight matrix $W_l \in \mathcal{R}^{d_{in} \times N_l}$ is used to specify different experts for the input x . As in the original MoE, MoSLD selects the top K experts for computation, and the gate score S_l^k is calculated as follows:

$$S_l^k(x) = \frac{\text{TopK}(\text{softmax}(W_l x), K)_k}{\sum_{k=1}^K \text{TopK}(\text{softmax}(W_l x), K)_k} \quad (1)$$

3.2 Dropout Strategy

In order to alleviate the imbalance and over-fitting problems caused by frequent general-feature matrix updates, we propose to apply the dropout strategy on the general-feature parameter matrix A_l . Dropout involves randomly ignoring a proportion of updates to the parameter matrix during each it-

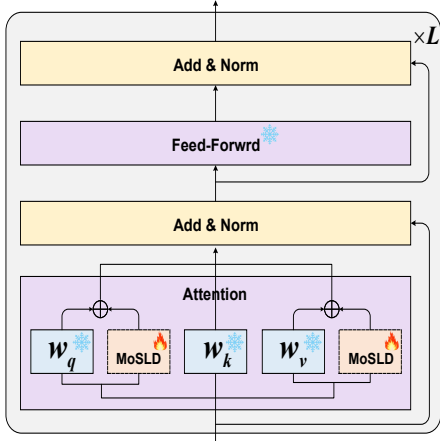


Figure 3: The overview of our proposed Mixture-of-Shared-LoRA with dropout strategy applied on W_q and W_v .

eration of training. This technique helps prevent over-reliance on specific parameters and encourages robust learning by introducing noise. That is, at each iteration, we take a certain probability p to discard the update in the general-feature matrix. Specifically, we generate a binary mask matrix drawn from Bernoulli distribution with a mask probability p , where each element in the general-feature matrix independently takes a value of 1 (keeping the parameter) with probability $1 - p$ or 0 (dropping the parameter) with probability p . The general-feature matrix is updated as follows:

$$\begin{aligned} \text{Mask} &\sim \text{Bernoulli}(p) \\ \mathbf{A}'_l &= \text{Mask} \odot \mathbf{A}_l \\ \widetilde{\mathbf{A}}'_l &= \mathbf{A}'_l / (1 - p) \end{aligned} \quad (2)$$

3.3 The Overall Procedure

Our method is a combination of shared LoRA modules and MoE framework, as shown in Figure 3. Here, we apply our MoSLD on the matrix Q and matrix V of the self-attention layer:

$$h_l = W_0 x + \frac{\alpha}{r} \sum_{k=1}^K S_l^k(x) A_{k,l} B_{k,l} x \quad (3)$$

where $W_0 \in \{W_q, W_v\}$ and h_l is the output embedding. Besides, similar to previous sparse MoE works, the load balancing loss L_b is also applied on each MoE layer, which is formulated as:

$$\begin{aligned} L_b &= \sum_{k=1}^K c_k \cdot s^k \\ p_k &= \sum_{x \in X} \frac{e^{S^k(x)}}{\sum_{k=1}^K e^{S^k(x)}} \end{aligned} \quad (4)$$

where c_k is the number of tokens assigned to the k -th expert.

4 Experimental Setup

4.1 Datasets

To evaluate the effectiveness of MoSLD, we conduct experiments on six commonsense reasoning datasets, including commonsense QA task (OBQA (Mihaylov et al., 2018), CSQA (Talmor et al., 2019)), reading comprehension task (Race (Lai et al., 2017), MCTest (Richardson et al., 2013)), and subject knowledge QA task (Arc-e (Clark et al., 2018), and Arc-c (Clark et al., 2018)). We denote the six datasets as $\{D_1, D_2, \dots, D_6\}$, and we also create a mixed dataset D_{mix} , corresponding to the single setting and the mixture setting respectively. The dataset sizes are as follows for training and testing: 5,457/500, 10,962/1140, 10,083/4934, 1,330/147, 2,821/2,376, and 1,418/1,172. We allocate 10% of the training set for validation. For all datasets, we use answer accuracy as the evaluation metric.

4.2 Baselines

We compare MoSLD with five parameter-efficient fine-tuning methods: Prefix-tuning (Li and Liang, 2021; Zhao et al., 2022b), LoRA (Hu et al., 2022), MoLoRA (Zadouri et al., 2024), SiRA (Zhu et al., 2023), MoLA (Gao et al., 2024), MixLoRA (Li et al., 2024). Additionally, we evaluate full-parameter fine-tuning. The details can be seen in Appendix A.

4.3 Training Details

We take LLaMA2-7B (Touvron et al., 2023b) which contains 32 layers as our base model. For plain LoRA and its variants, the r is set to 8 and α is 16. Besides, the LoRA modules are used in matrix Q and matrix V in attention layers. Our MoSLD also follows the same settings. We allocate 8 experts to each layer for 1-8 layers, 6 experts to each layer for 9-16 layers, 4 experts to each layer for 17-24 layers, and 2 experts to each layer for the last 8 layers. The K of the selected experts is 2. For training details, we finetune models with 10 epochs

Model		OBQA	CSQA	Race	MCTest	Arc-e	Arc-c	Avg
FP-tuning	single	75.00	75.74	80.62	39.05	72.39	60.63	67.24
	mixture	76.00	75.27	81.46	50.42	73.69	65.45	70.38
Prefix-tuning	single	47.76	42.65	53.77	25.19	45.65	35.50	41.70
	mixture	46.51	44.98	49.88	22.46	47.92	35.30	41.18
LoRA	single	75.40	76.33	76.06	53.10	73.82	62.71	69.57
	mixture	72.80	76.30	78.23	55.67	70.87	61.00	69.15
MoLoRA	single	74.71	76.65	74.26	49.08	74.14	61.38	68.37
	mixture	75.04	75.27	73.88	55.37	75.25	62.86	69.61
SiRA	single	73.99	76.26	75.63	48.28	74.02	62.86	68.51
	mixture	74.34	76.22	75.04	52.33	74.98	63.16	69.35
MoLA	single	74.60	77.23	75.29	44.90	72.73	60.80	67.59
	mixture	76.60	73.46	75.25	54.42	76.34	63.91	70.00
MixLoRA	single	75.60	74.83	75.47	50.88	74.51	60.10	68.57
	mixture	75.80	76.81	74.79	54.26	74.41	63.62	69.95
MoSL (our)	single	76.30	77.56	74.63	49.66	76.30	60.48	69.16
	mixture	76.80 (+0.50)	75.02 (-2.54)	74.69 (+0.06)	58.50 (+8.84)	76.09 (-0.21)	64.16 (+3.68)	70.88 (+1.72)
MoSLD (our)	single	78.40	75.84	76.08	53.06	76.35	61.49	70.20
	mixture	78.80 (+0.40)	76.43 (+0.59)	76.96 (+0.88)	54.42 (+1.36)	76.60 (+0.25)	66.13 (+4.64)	71.56 (+1.36)

Table 1: Results of different methods on the in-domain test sets of six commonsense reasoning datasets. We also report the increase of mixture setting compared to single setting. Results are averaged over three random runs. ($p < 0.01$ under t-test)

and a peak of $3e-4$ learning rate. The drop ratio applied to matrix A is set to 0.1. The batch size during model tuning is 128. The experiments are run on 16 NVIDIA A100 40GB GPUs.

4.4 Main Results

Table 1 presents the experimental outcomes of various baselines under both single and mixture settings across different datasets. Initially, we report the performance of models trained on individual datasets. LoRA notably outperforms other baselines, exhibiting improvements of 2.33% and 27.87% over FP-tuning (single) and Prefix-tuning (single), respectively. MoLoRA, SiRA, MoLA, and MixLoRA trail behind LoRA by 1.20%, 1.06%, 1.98%, and 1.00%, indicating that simply combining LoRA and MoE does not confer an advantage in single in-domain datasets. After establishing a robust baseline in the single setting, we proceed to report results for the mixture setting. Here, we observe a decline in LoRA’s performance, trailing 1.23 points behind FP-tuning (70.38%). Conversely, applying the MoE framework to LoRA, i.e., MoLoRA, SiRA, MoLA, and MixLoRA, achieves scores of 69.61%, 69.35%, 70.00%, and 69.95%, demonstrating MoE’s suitability for multi-task scenarios and MoLA is the best performing baseline in the mixture setting. Further comparison between single and mixture settings reveals that FP-tuning and MoLA improve by 3.14% and 2.41%, respectively, in the mixture setting compared to the single setting. However, LoRA’s performance decreases by 0.42% in the mixture setting compared to the single setting, indicating conflicts between multi-

task data and the mixture strategy’s detrimental impact on performance.

Upon closer examination, our proposed MoSLD demonstrates performance enhancements of 2.61% and 1.56% over MoLA in single and mixture settings, respectively. This emphasizes the effectiveness of the sharing mechanism and dropout strategy in alleviating data conflicts and retaining shared knowledge between various tasks. Furthermore, conducting ablation experiments by removing the dropout strategy, MoSL experiences performance decreases of 1.04% and 0.68%, respectively, compared to MoSLD. This highlights the crucial role of the dropout strategy in mitigating training overfitting and optimization imbalance. Nevertheless, MoSL still achieves competitive results of 69.16% and 70.88%. We also found that our model not only achieves good results in the mixture setting, but also achieves good results in the single setting, which overcomes the disadvantage of MoLA’s poor performance in the single setting. However, we find that our models, especially MoSL, do not have much advantage over plain LoRA, which is consistent with the performance of all baselines combining MoE with LoRA. This is because the complexity of the model ensemble causes overfitting on a single simple task, resulting in little improvement. In conclusion, our approach exhibits significant advantages under both single and mixture settings, particularly in alleviating data conflicts across multiple tasks and addressing knowledge forgetting issues in multi-task learning. In addition, we also pay attention to the efficiency of training. Due to the introduction of multiple LoRAs, the trainable

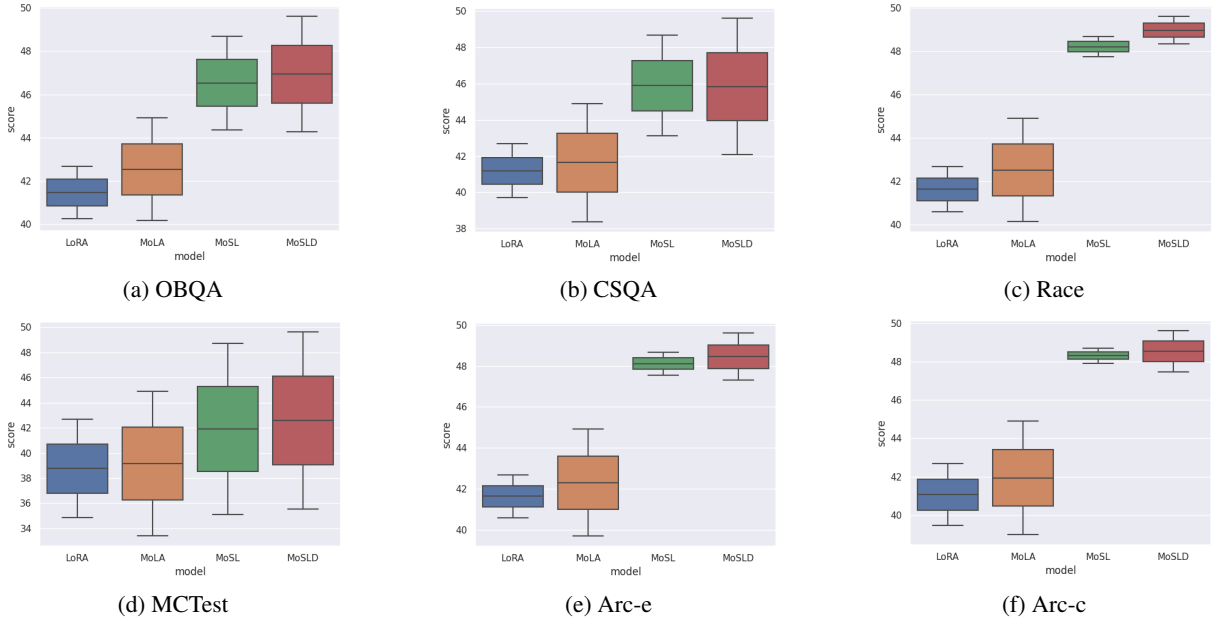


Figure 4: A comparison of performance for LoRA, MoLA, MoSL, and MoSLD on single and mixture settings for MMLU test set.

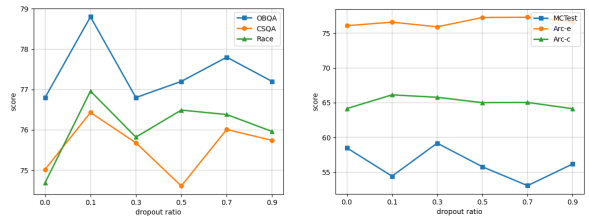
parameters of MoLA are higher than those of plain LoRA. However, although our MoSLD expands LoRA several times through the MoE architecture, it does not introduce a large number of additional parameters and also enables the LoRA training to have multiple capabilities. Details can be seen in Section 5.5.

5 Qualitative Analysis

5.1 Out-of-domain Test

To assess the generalization capability of our proposed model, we conducted out-of-domain experiments using the test set of MMLU. Figure 4 presents a boxplot, where the top and bottom horizontal lines represent the mixture and single settings, respectively. Our models, MoSL and MoSLD, consistently outperform others in both settings, exhibiting significant improvements, particularly on Race, Arc-e, and Arc-c datasets. This highlights the effectiveness of our models in disentangling domain knowledge and transferring general features across diverse datasets. OBQA and CSQA exhibit similar trends in the boxplot, indicating similar data distributions between the two datasets. Conversely, for MCTest, while improvements are observed in the mixture settings, the single settings remain relatively unchanged. This divergence may stem from the substantial differences between the MCTest and MMLU test sets, suggesting that introducing data from other domains

or tasks could inspire general domain knowledge. In summary, our model demonstrates strong generalization capabilities, particularly in multi-task scenarios.



(a) OBQA&CSQA&Race (b) MCTest&Arc-e&Arc-c

Figure 5: Results of six datasets under different dropout ratios. Here, we are based on the mixture setting.

5.2 Effect of Model Parameters

In this section, we conduct parameter search experiments.

Dropout Location As shown in Table 2, we show the results of applying our methods on matrix A and matrix B. We found that in the single setting, MoSLD (matrix B) does not achieve much improvement, 0.94 points lower than the ordinary LoRA and 1.04 points higher than MoLA. The mixture setting still achieves good results. However, the results of applying our method on matrix B are lower than those of applying it on matrix A in both the single and mixture settings. This also shows that matrix A is more used to extract general features.

Dropout Ratio In Figure 5, we depict the per-

Model		OBQA	CSQA	Race	MCTest	Arc-e	Arc-c	Avg
LoRA	single	75.40	76.33	76.06	53.10	73.82	62.71	69.57
	mixture	72.80	76.30	78.23	55.67	70.87	61.00	69.15
MoLA	single	74.60	77.23	75.29	44.90	72.73	60.80	67.59
	mixture	76.60	73.46	75.25	54.42	76.34	63.91	70.00
MoSLD (matrix A)	single	78.40	75.84	76.08	53.06	76.35	61.49	70.20
	mixture	78.80	76.43	76.96	54.42	76.60	66.13	71.56
MoSLD (matrix B)	single	77.60	75.76	74.58	46.94	76.09	60.83	68.63
	mixture	76.40	74.11	75.25	56.46	77.15	65.02	70.73

Table 2: The results for applying our methods on matrix A and matrix B.

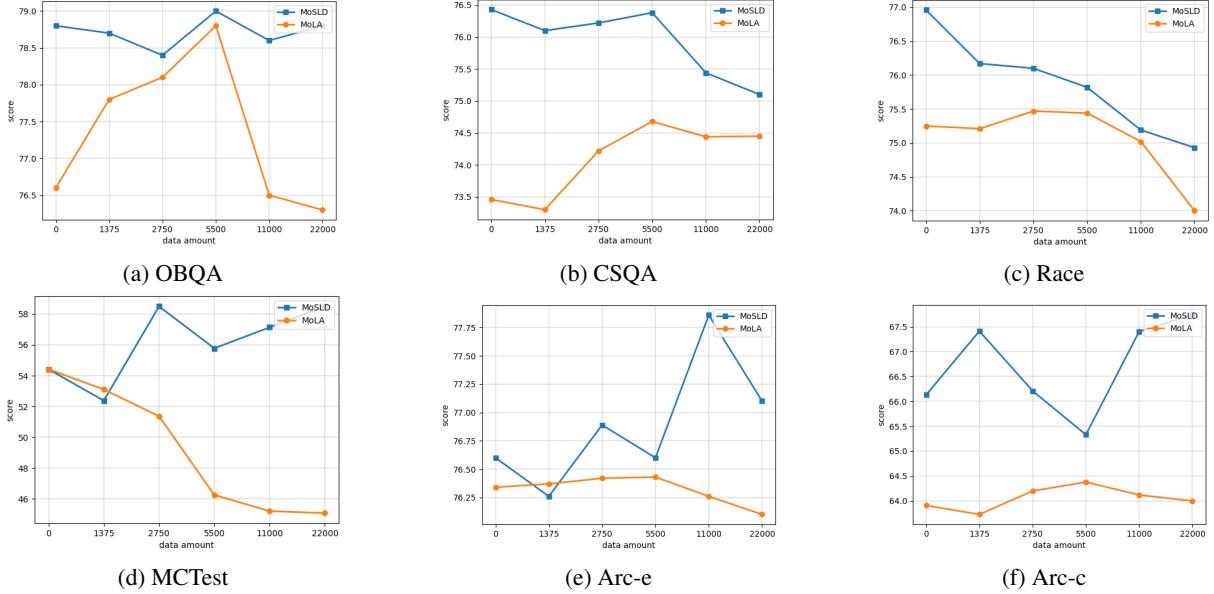


Figure 6: Different data amount of OpenOrca between MoSLD and MoLA on six datasets. Here, we use the mixture setting.

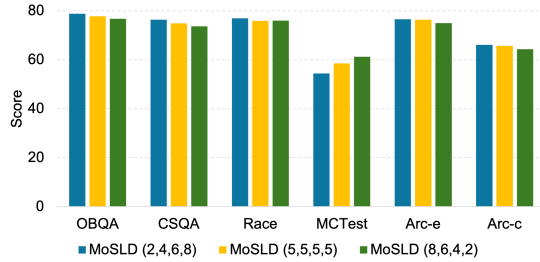


Figure 7: Different allocation strategies for the number of experts at different layers of the model. Here, we use the mixture setting.

formance of six datasets under the mixture setting with varying dropout ratios. We observe a general downward trend in most results as the dropout ratio increases. This phenomenon occurs because while dropout can mitigate overfitting to some extent, excessively high dropout rates may diminish the model’s capabilities. Therefore, careful selection of the dropout ratio parameter is necessary. Interestingly, the curves for the Arc-e and Arc-c datasets remain relatively stable across different dropout ratios. We attribute this stability to the simplicity of

these two datasets, where model sparsification has minimal impact on the results.

Expert Number Considering the redundancy among experts, following (Gao et al., 2024), we set different numbers of experts at different layers in Figure 7. Keeping the total number of experts constant, we choose three settings, i.e., (2,4,6,8), (5,5,5,5), (8,6,4,2). It is observed that assigning more experts at higher layers and fewer experts at lower layers, i.e., (2,4,6,8), works better. This is consistent with people’s intuition: the lower layers of the model mainly extract general knowledge, which can be well learned by a small number of experts. While the higher layers of the model focus more on acquiring specific features of different tasks, and a larger number of experts can better capture multi-aspect capabilities.

5.3 Mix with Other Data

Mathematical Reasoning Data We construct a new multi-task setting, including commonsense QA task (OBQA), reading comprehension task

Model		OBQA	MCTest	Arc-c	GSM8K	Avg
LoRA	single	75.40	53.10	62.71	23.12	53.58
	mixture	73.20	55.10	64.08	17.51	52.47
MoSLD	single	78.40	53.06	61.49	22.06	53.75
	mixture	79.80	53.90	63.29	22.73	53.93

Table 3: The results of the mixture setting of tasks with little commonality.

Model		OBQA	CSQA	Race	MCTest	Arc-e	Arc-c	Avg
LLaMA2-7B	single	78.40	75.84	76.08	53.06	76.35	61.49	70.20
	mixture	78.80	76.43	76.96	54.42	76.60	66.13	71.56
LLaMA2-13B	single	81.4	77.95	78.01	57.86	78.93	65.05	73.20
	mixture	82.2	78.46	79.87	58.50	79.67	70.14	74.81
LLaMA-33B	single	83.93	81.49	83.27	65.99	85.10	68.52	78.05
	mixture	84.55	83.26	84.90	66.73	85.95	74.36	79.96

Table 4: The results of six datasets in single and mixture settings based on LLaMA2-7B, LLaMA2-13B and LLaMA-33B.

(MCTest), subject knowledge QA task (Arc-c), and mathematical reasoning task (GSM8K). As shown in Table 3, we found that for plain LoRA, the mixture setting was 1.11 points lower than the single setting on average, especially for GSM8K, it is reduced by 5.61%, which shows that it is very challenging for plain LoRA to train multiple tasks with little commonality. However, for our MoSLD, the mixture setting is 1.18 points higher than the single setting on average. For the GSM8K with the largest difference, it is also improved by 0.67%. This shows that MoSLD is also effective for tasks with little commonality. This is because for tasks with little commonality, although the role of the shared general-feature matrix becomes smaller, the unique-feature matrix still captures the knowledge of each task, which further shows that our MoSLD can effectively alleviate the data conflict problem in multi-task learning.

Mix with General Data In Figure 6, we illustrate the impact of adding varying amounts of randomly filtered data from OpenOrca¹ to the mixed dataset D_{mix} . The data amount from OpenOrca ranges from 1,375 to 22,000. We observed that for MoLA, as the amount of general data increases, performance initially improves before eventually declining. This suggests that mixing a large amount of general data can lead to data conflicts and domain knowledge forgetting. In contrast, MoSLD demonstrates an upward trend in performance with the increase in data amount for OBQA, MCTest, Arc-e, and Arc-c. However, performance on CSQA and Race experiences a decline. We attribute this to

significant distribution differences between these datasets and the general data. Overall, our model consistently outperforms MoLA when mixing various amounts of generic data. This underscores our model’s ability to effectively leverage general knowledge across different tasks.

5.4 Scaling of Model Size

Table 4 shows the results of our model for the six datasets both in single and mixture settings as the model size scalings. We find that the performance of our model increases with the size of the model, whether in single or mixture settings, which is in line with our expectations. In addition, it is observed that the results improve by 1.36%, 1.61%, and 1.91% from single to mixture for LLaMA2-7B, LLaMA2-13B, and LLaMA-33B, respectively. The experimental results show that our method has achieved good performance on models of different sizes, and has a certain scaling ability. We also give the model size scaling results of other LoRA-based baselines, which can be seen in the Appendix C.

5.5 Analysis of Computation Efficiency

In Table 5, we further show the computational efficiency of our model. We first analyze the number of new LoRA modules inserted in ordinary LoRA, MoLA, and MoSLD. Since MoLA introduces the MoE framework, the trainable parameters become 5 times that of ordinary LoRA, and its results are improved by 0.43 points from 69.57 to 70.00. We believe that despite the introduction of a large number of trainable parameters, the change in results is not very large, which is a method of sacrificing efficiency for effect. In addition, we also found that

¹<https://huggingface.co/datasets/Open-Orca/OpenOrca>

Model	LoRA number	Forward param	Trainable param	Avg_score
FP-tuning	/	6.738B	6.738B	70.38
LoRA	(1A+1B)*32	6.743B	0.419B	69.57
MoLA	(5A+5B)*32	6.761B	2.228B	70.00
MoSLD	(1A+5B)*32	6.572B	1.389B	71.56

Table 5: The number of LoRA matrices, forward parameters, and trainable parameters for FP-tuning, LoRA, MoLA, and our MoSLD during training. Here, "A" is matrix A, "B" is matrix B, and "5" is the average number of experts per layer. We also report the average results across 6 datasets under the mixture setting.

although our method reduces 128 matrix A compared to MoLA, it is still 1.56% higher than MoLA and 1.99% higher than LoRA. This shows that although our MoSLD introduces multiple LoRAs through the MoE framework, the expert sharing mechanism greatly reduces the additional parameters and achieves a balance between effect and efficiency. We also compare FP-tuning. Although our trainable parameters are 20.6% of FP-tuning, but it still achieves a 1.18 point improvement. This also proves that our MoSLD is indeed an extremely efficient-parameter fine-tuning method.

6 Conclusion

In this paper, we propose MoSLD, which is a mixture-of-shared-LoRAs model with dropout strategy. Unlike traditional LoRA-MoE approaches, we design a sharing mechanism for matrix A, which aims to capture the general-feature among various tasks. A dropout strategy is also applied to the matrix A, solving the overfitting caused by parameter redundancy to a certain extent. Evaluations show that MoSLD outperforms the baseline in both single-task and multi-task scenarios. Especially in multi-task scenarios, where it can effectively alleviate knowledge conflict and forgetting problems. In general, our model is extremely parameter-efficient for fine-tuning.

Acknowledgement

We thank all anonymous reviewers. This work was supported by National Science and Technology Major Project No.2022ZD0116314.

Limitations

Although MoSLD achieves significant improvements over existing baselines, there are still avenues worth exploring in future research. (1) This paper focuses on applying MoSLD on the matrix Q and V of the attention layer. We hope to extend this method to the FFN layer. (2) This paper explores the multi-task setting of directly mixing

multiple datasets and compares with the performance of a single task. We plan to study the impact of multi-task data ratio on MoSLD. (3) This paper emphasizes the extraction of general and unique features by the upper and lower projection matrices in LoRA, and intends to visualize this phenomenon in the future.

Ethics Statement

LoRA has emerged as a pivotal technique for refining extensive pre-trained models. Nevertheless, its efficacy tends to fail in multi-task learning. Conversely, the MoE architecture offers a promising remedy to this setback. However, it introduces hurdles such as the interference of data across diverse domains and the risk of forgetting knowledge from various tasks. Furthermore, MoE substantially inflates parameter counts, presenting computational challenges. In light of these considerations, we present MoSLD in this paper, a model that integrates the strengths of both approaches. MoSLD, a mixture-of-shared-LoRAs model with a dropout strategy, addresses these obstacles ingeniously. By sharing the upper projection matrix in LoRA among different experts, MoSLD fosters the acquisition of broad knowledge across tasks while allowing the lower projection matrix to concentrate on task-specific features. Additionally, the application of dropout mitigates parameter overfitting in LoRA. The experimental results prove the effectiveness of our model and evaluation framework. Besides, there is no huge biased content in the datasets and the models. If the knowledge base is further used, the biased content will be brought into the generated responses, just like biased content posted by content creator on the Web which is promoted by a search engine. To prevent the technology from being abused for disinformation, we look forward to more research effort being paid to fake/biased/offensive content detection and encourage developers to carefully choose the proper dataset and content to build the knowledge base.

References

- Shaoxiang Chen, Zequn Jie, and Lin Ma. 2024. [Llava-mole: Sparse mixture of lora experts for mitigating data conflicts in instruction finetuning mllms](#). *Preprint*, arXiv:2401.16160.
- Zewen Chi, Li Dong, Shaohan Huang, Damai Dai, Shuming Ma, Barun Patra, Saksham Singhal, Payal Bajaj, Xia Song, Xian-Ling Mao, Heyan Huang, and Furu Wei. 2022. [On the representation collapse of sparse mixture of experts](#). In *Advances in Neural Information Processing Systems*.
- Peter Clark, Isaac Cowhey, Oren Etzioni, Tushar Khot, Ashish Sabharwal, Carissa Schoenick, and Oyvind Tafjord. 2018. [Think you have solved question answering? try arc, the ai2 reasoning challenge](#). *arXiv:1803.05457v1*.
- Shihan Dou, Enyu Zhou, Yan Liu, Songyang Gao, Jun Zhao, Wei Shen, Yuhao Zhou, Zhiheng Xi, Xiao Wang, Xiaoran Fan, Shiliang Pu, Jiang Zhu, Rui Zheng, Tao Gui, Qi Zhang, and Xuanjing Huang. 2024. [Loramoe: Alleviate world knowledge forgetting in large language models via moe-style plugin](#). *Preprint*, arXiv:2312.09979.
- William Fedus, Jeff Dean, and Barret Zoph. 2022a. [A review of sparse expert models in deep learning](#). *Preprint*, arXiv:2209.01667.
- William Fedus, Barret Zoph, and Noam Shazeer. 2022b. [Switch transformers: scaling to trillion parameter models with simple and efficient sparsity](#). *J. Mach. Learn. Res.*, 23(1).
- Wenfeng Feng, Chuzhan Hao, Yuewei Zhang, Yu Han, and Hao Wang. 2024. [Mixture-of-loras: An efficient multitask tuning for large language models](#). *Preprint*, arXiv:2403.03432.
- Chongyang Gao, Kezhen Chen, Jinmeng Rao, Baochen Sun, Ruiyao Liu, Daiyi Peng, Yawen Zhang, Xiaoyuan Guo, Jie Yang, and VS Subrahmanian. 2024. [Higher layers need more lora experts](#). *Preprint*, arXiv:2402.08562.
- Edward J Hu, yelong shen, Phillip Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang, and Weizhu Chen. 2022. [LoRA: Low-rank adaptation of large language models](#). In *International Conference on Learning Representations*.
- Chengsong Huang, Qian Liu, Bill Yuchen Lin, Chao Du, Tianyu Pang, and Min Lin. 2024. [Lorahub: Efficient cross-task generalization via dynamic loRA composition](#).
- Robert A. Jacobs, Michael I. Jordan, Steven J. Nowlan, and Geoffrey E. Hinton. 1991. [Adaptive mixtures of local experts](#). *Neural Computation*, 3(1):79–87.
- Young Jin Kim, Ammar Ahmad Awan, Alexandre Muzio, Andres Felipe Cruz Salinas, Liyang Lu, Amr Hendy, Samyam Rajbhandari, Yuxiong He, and Hany Hassan Awadalla. 2021. [Scalable and efficient moe training for multitask multilingual models](#). *Preprint*, arXiv:2109.10465.
- Guokun Lai, Qizhe Xie, Hanxiao Liu, Yiming Yang, and Eduard Hovy. 2017. [RACE: Large-scale ReAding comprehension dataset from examinations](#). In *Proceedings of the 2017 Conference on Empirical Methods in Natural Language Processing*, pages 785–794, Copenhagen, Denmark. Association for Computational Linguistics.
- Dengchun Li, Yingzi Ma, Naizheng Wang, Zhengmao Ye, Zhiyuan Cheng, Yinghao Tang, Yan Zhang, Lei Duan, Jie Zuo, Cal Yang, and Mingjie Tang. 2024. [Mixlora: Enhancing large language models fine-tuning with lora-based mixture of experts](#). *Preprint*, arXiv:2404.15159.
- Xiang Lisa Li and Percy Liang. 2021. [Prefix-tuning: Optimizing continuous prompts for generation](#). In *Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume 1: Long Papers)*, pages 4582–4597, Online. Association for Computational Linguistics.
- Qidong Liu, Xian Wu, Xiangyu Zhao, Yuanshao Zhu, Derong Xu, Feng Tian, and Yefeng Zheng. 2023. [Moelora: An moe-based parameter efficient fine-tuning method for multi-task medical applications](#). *Preprint*, arXiv:2310.18339.
- Tongxu Luo, Jiahe Lei, Fangyu Lei, Weihao Liu, Shizhu He, Jun Zhao, and Kang Liu. 2024. [Moelora: Contrastive learning guided mixture of experts on parameter-efficient fine-tuning for large language models](#). *Preprint*, arXiv:2402.12851.
- Todor Mihaylov, Peter Clark, Tushar Khot, and Ashish Sabharwal. 2018. [Can a suit of armor conduct electricity? a new dataset for open book question answering](#). In *EMNLP*.
- Long Ouyang, Jeffrey Wu, Xu Jiang, Diogo Almeida, Carroll Wainwright, Pamela Mishkin, Chong Zhang, Sandhini Agarwal, Katarina Slama, Alex Ray, John Schulman, Jacob Hilton, Fraser Kelton, Luke Miller, Maddie Simens, Amanda Askell, Peter Welinder, Paul F Christiano, Jan Leike, and Ryan Lowe. 2022. [Training language models to follow instructions with human feedback](#). In *Advances in Neural Information Processing Systems*, volume 35, pages 27730–27744. Curran Associates, Inc.
- Matthew Richardson, Christopher J.C. Burges, and Erin Renshaw. 2013. [MCTest: A challenge dataset for the open-domain machine comprehension of text](#). In *Proceedings of the 2013 Conference on Empirical Methods in Natural Language Processing*, pages 193–203, Seattle, Washington, USA. Association for Computational Linguistics.

- Carlos Riquelme, Joan Puigcerver, Basil Mustafa, Maxim Neumann, Rodolphe Jenatton, André Susano Pinto, Daniel Keysers, and Neil Houlsby. 2021. [Scaling vision with sparse mixture of experts](#). In *Advances in Neural Information Processing Systems*, volume 34, pages 8583–8595. Curran Associates, Inc.
- Noam Shazeer, *Azalia Mirhoseini, *Krzysztof Maziarczyk, Andy Davis, Quoc Le, Geoffrey Hinton, and Jeff Dean. 2017. [Outrageously large neural networks: The sparsely-gated mixture-of-experts layer](#). In *International Conference on Learning Representations*.
- Alon Talmor, Jonathan Herzig, Nicholas Lourie, and Jonathan Berant. 2019. [CommonsenseQA: A question answering challenge targeting commonsense knowledge](#). In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)*, pages 4149–4158, Minneapolis, Minnesota. Association for Computational Linguistics.
- 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. 2023a. [Llama: Open and efficient foundation language models](#). *Preprint*, arXiv:2302.13971.
- Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Nikolay Bashlykov, Soumya Batra, Prajwal Bhargava, Shruti Bhosale, Dan Bikel, Lukas Blecher, Cristian Canton Ferrer, Moya Chen, Guillem Cucurull, David Esiobu, Jude Fernandes, Jeremy Fu, Wenyin Fu, Brian Fuller, Cynthia Gao, Vedanuj Goswami, Naman Goyal, Anthony Hartshorn, Saghar Hosseini, Rui Hou, Hakan Inan, Marcin Kardas, Viktor Kerkez, Madian Khabsa, Isabel Kloumann, Artem Korenev, Punit Singh Koura, Marie-Anne Lachaux, Thibaut Lavril, Jenya Lee, Diana Liskovich, Yinghai Lu, Yuning Mao, Xavier Martinet, Todor Mihaylov, Pushkar Mishra, Igor Molybog, Yixin Nie, Andrew Poulton, Jeremy Reizenstein, Rashi Rungta, Kalyan Saladi, Alan Schelten, Ruan Silva, Eric Michael Smith, Ranjan Subramanian, Xiaoqing Ellen Tan, Binh Tang, Ross Taylor, Adina Williams, Jian Xiang Kuan, Puxin Xu, Zheng Yan, Iliyan Zarov, Yuchen Zhang, Angela Fan, Melanie Kambadur, Sharan Narang, Aurelien Rodriguez, Robert Stojnic, Sergey Edunov, and Thomas Scialom. 2023b. [Llama 2: Open foundation and fine-tuned chat models](#). *Preprint*, arXiv:2307.09288.
- Sheng Wang, Liheng Chen, Jiyue Jiang, Boyang Xue, Lingpeng Kong, and Chuan Wu. 2024. [Lora meets dropout under a unified framework](#). *Preprint*, arXiv:2403.00812.
- Xun Wu, Shaohan Huang, and Furu Wei. 2024. [Mixture of loRA experts](#). In *The Twelfth International Conference on Learning Representations*.
- Ted Zadouri, Ahmet Üstün, Arash Ahmadian, Beyza Ermis, Acyr Locatelli, and Sara Hooker. 2024. [Pushing mixture of experts to the limit: Extremely parameter efficient moe for instruction tuning](#). In *The Twelfth International Conference on Learning Representations*.
- Weihao Zeng, Lulu Zhao, Keqing He, Ruotong Geng, Jingang Wang, Wei Wu, and Weiran Xu. 2023. [Seen to unseen: Exploring compositional generalization of multi-attribute controllable dialogue generation](#). In *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 14179–14196, Toronto, Canada. Association for Computational Linguistics.
- Jinghan Zhang, Shiqi Chen, Junteng Liu, and Junxian He. 2023. [Composing parameter-efficient modules with arithmetic operation](#). In *Thirty-seventh Conference on Neural Information Processing Systems*.
- Lulu Zhao, Fujia Zheng, Weihao Zeng, Keqing He, Ruotong Geng, Huixing Jiang, Wei Wu, and Weiran Xu. 2022a. [Adpl: Adversarial prompt-based domain adaptation for dialogue summarization with knowledge disentanglement](#). In *Proceedings of the 45th International ACM SIGIR Conference on Research and Development in Information Retrieval, SIGIR '22*, page 245–255, New York, NY, USA. Association for Computing Machinery.
- Lulu Zhao, Fujia Zheng, Weihao Zeng, Keqing He, Weiran Xu, Huixing Jiang, Wei Wu, and Yanan Wu. 2022b. [Domain-oriented prefix-tuning: Towards efficient and generalizable fine-tuning for zero-shot dialogue summarization](#). In *Proceedings of the 2022 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 4848–4862, Seattle, United States. Association for Computational Linguistics.
- Yanqi Zhou, Tao Lei, Hanxiao Liu, Nan Du, Yanping Huang, Vincent Zhao, Andrew M Dai, zhifeng Chen, Quoc V Le, and James Laudon. 2022. [Mixture-of-experts with expert choice routing](#). In *Advances in Neural Information Processing Systems*, volume 35, pages 7103–7114. Curran Associates, Inc.
- Yun Zhu, Nevan Wichers, Chu-Cheng Lin, Xinyi Wang, Tianlong Chen, Lei Shu, Han Lu, Canoe Liu, Liangchen Luo, Jindong Chen, and Lei Meng. 2023. [Sira: Sparse mixture of low rank adaptation](#). *Preprint*, arXiv:2311.09179.

A Baselines

In this section, we introduce the baselines in detail.

Prefix-tuning (Li and Liang, 2021; Zhao et al., 2022b): This method involves incorporating soft prompts into each attention layer of the Large Language Model (LLM). These soft prompts are a series of virtual tokens pre-appended to the text.

During fine-tuning, the LLM remains frozen, and only the virtual tokens are optimized.

LoRA (Hu et al., 2022): A popular parameter-efficient tuning approach widely used in LLM fine-tuning, LoRA leverages low-rank matrix decomposition of pre-trained weight matrices to significantly reduce the number of training parameters.

MoLoRA (Zadouri et al., 2024): A method which is a parameter-efficient MoE by uniquely combining MoE architecture with lightweight experts.

SiRA (Zhu et al., 2023): A method leverages the Sparse Mixture of Expert (SMoE) and enforces the top k experts routing with a capacity limit restricting the maximum number of tokens each expert can process.

MoLA (Gao et al., 2024): A LoRA variant with layer-wise expert allocation, MoLA flexibly assigns a different number of LoRA experts to each Transformer layer.

MixLoRA (Li et al., 2024): It inserts multiple LoRA-based experts within the feed-forward network block of a frozen pre-trained dense model and employs a commonly used top-k router.

B Effect on Rank

In this section, we add experiments on the effect of rank for our MoSLD, with r ranging from 2 to 32. Overall, the results of the six datasets did not fluctuate much, and the best value was obtained at 8 or 16. From the perspective of efficiency, 8 is indeed a suitable hyperparameter, which is also in line with the change law of LoRA's rank. The results are as shown in Table 6:

C Scaling of Model Size

In this section, We add model scaling experiments on LoRA-based baselines, such as LoRA, MoLoRA, SiRA, and MoLA. We find that for each baseline, the results improve as the model size increases, among which our model MoSLD scales even better. The results are shown in Table 7 :

Dataset	r=2	r=4	r=8	r=16	r=32
OBQA	76.19	76.34	78.80	75.53	74.27
CSQA	74.35	75.16	76.43	77.39	76.62
Race	75.22	76.01	76.96	76.74	74.73
MCTest	52.28	54.17	54.42	54.16	53.52
Arc-e	75.51	76.98	76.70	75.88	75.63
Arc-c	63.28	64.06	66.13	66.10	65.87

Table 6: The performance of our MoSLD as different rank values.

Model			OBQA	CSQA	Race	MCTest	Arc-e	Arc-c	Avg
LoRA	7B	single	75.40	76.33	76.06	53.10	73.82	62.71	69.57
		mixture	72.80	76.30	78.23	55.67	70.87	61.00	69.15
	13B	single	77.21	79.84	77.34	58.29	74.99	63.89	71.93
		mixture	77.98	78.32	77.83	55.74	74.05	64.11	71.34
	33B	single	79.06	80.97	81.78	59.54	77.36	64.79	73.92
		mixture	79.05	80.02	82.95	58.27	75.33	64.88	73.42
MoLoRA	7B	single	75.40	76.33	76.06	53.10	73.82	62.71	69.57
		mixture	72.80	76.30	78.23	55.67	70.87	61.00	69.15
	13B	single	77.46	81.26	75.33	51.79	75.83	64.27	70.99
		mixture	77.95	82.44	80.25	54.73	74.21	62.65	72.04
	33B	single	78.23	83.18	79.59	59.41	82.11	65.28	74.63
		mixture	77.54	81.35	81.78	61.62	82.07	64.35	74.79
SiRA	7B	single	73.99	76.26	75.63	48.28	74.02	62.86	68.51
		mixture	74.34	76.22	75.04	52.33	74.98	63.16	69.35
	13B	single	75.15	77.93	78.28	50.78	73.85	62.03	69.67
		mixture	75.01	76.45	78.11	50.24	74.52	61.74	69.35
	33B	single	78.99	81.34	80.03	53.59	75.78	64.55	72.38
		mixture	79.46	82.02	80.00	56.84	75.81	66.75	73.48
MoLA	7B	single	74.60	77.23	75.29	44.90	72.73	60.80	67.59
		mixture	76.60	73.46	75.25	54.42	76.34	63.91	70.00
	13B	single	76.82	80.55	76.87	48.35	74.84	63.66	70.18
		mixture	77.61	77.59	77.04	60.83	76.71	65.27	72.51
	33B	single	80.36	82.94	79.06	50.88	76.00	67.06	72.72
		mixture	81.79	85.03	79.82	57.35	76.48	68.82	74.88
MixLoRA	7B	single	75.60	74.83	75.47	50.88	74.51	60.10	68.57
		mixture	75.80	76.81	74.79	54.26	74.41	63.62	69.95
	13B	single	77.33	78.34	76.82	53.12	77.39	64.53	71.26
		mixture	76.98	78.05	77.31	56.88	78.00	66.92	72.36
	33B	single	80.57	81.04	78.99	55.62	81.25	67.45	74.15
		mixture	80.03	82.87	79.45	58.98	79.73	70.87	75.32
MoSLD	7B	single	78.40	75.84	76.08	53.06	76.35	61.49	70.20
		mixture	78.80	76.43	76.96	54.42	76.60	66.13	71.56
	13B	single	81.40	77.95	78.01	57.86	78.93	65.05	73.20
		mixture	82.20	78.46	79.87	58.50	79.67	70.14	74.81
	33B	single	83.93	81.94	83.27	65.99	85.10	68.52	78.05
		mixture	84.55	83.26	84.90	66.73	85.95	74.36	79.96

Table 7: The model scaling results about LLaMA2-7B, LLaMA2-13B, and LLaMA-33B of six datasets in single and mixture settings for LoRA, MoLoRA, SiRA, MoLA, MixLoRA, and our MoSLD.