

MoA: Heterogeneous Mixture of Adapters for Parameter-Efficient Fine-Tuning of Large Language Models

Jie Cao¹, Tianwei Lin¹, Hongyang He, Rolan Yan²,
Wenqiao Zhang^{*}¹, Juncheng Li¹, Dongping Zhang, Siliang Tang¹, Yueting Zhuang¹

¹ Zhejiang University, ² Tencent

Abstract

Recent studies integrate Low-Rank Adaptation (LoRA) and Mixture-of-Experts (MoE) to further enhance the performance of parameter-efficient fine-tuning (PEFT) methods in Large Language Model (LLM) applications. Existing methods employ *homogeneous* MoE-LoRA architectures composed of LoRA experts with either similar or identical structures and capacities. However, these approaches often suffer from representation collapse and expert load imbalance, which negatively impact the potential of LLMs. To address these challenges, we propose a *heterogeneous Mixture-of-Adapters (MoA)* approach. This method dynamically integrates PEFT adapter experts with diverse structures, leveraging their complementary representational capabilities to foster expert specialization, thereby enhancing the effective transfer of pre-trained knowledge to downstream tasks. MoA supports two variants: (i) *Soft MoA* achieves fine-grained integration by performing a weighted fusion of all expert outputs; (ii) *Sparse MoA* activates adapter experts sparsely based on their contribution, achieving this with negligible performance degradation. Experimental results demonstrate that heterogeneous MoA outperforms homogeneous MoE-LoRA methods in both performance and parameter efficiency. Our project is available at <https://github.com/DCDm1lm/MoA>.

1 Introduction

The rapid development of Large Language Models (LLMs) (Achiam et al., 2023; Yang et al., 2025; Grattafiori et al., 2024a,a) is reshaping the field of NLP, demonstrating cross-domain generalization capabilities across a wide range of tasks. However, as model sizes continue to grow, the computational, storage, and deployment costs of traditional full fine-tuning methods increasingly become bottlenecks in practical applications. As a

result, parameter-efficient fine-tuning (PEFT) methods (Houlsby et al., 2019; He et al., 2022; Li and Liang, 2021; Lester et al., 2021; Zhang et al., 2024; Han et al., 2024; Lin et al., 2025) are emerging as a key research direction for adapting LLMs. These methods introduce a small number of trainable, lightweight adapter modules on top of frozen pre-trained weights. By doing so, they effectively guide the activation and adjustment of the model’s existing knowledge representations while maintaining low overhead, enabling fast transfer and adaptation to specific downstream tasks.

The Low-Rank Adaptation (LoRA) method (Hu et al., 2021) reparameterizes weight matrices with low-rank decomposition, enabling the approximation of full fine-tuning using only a small number of trainable parameters. However, increasing the rank does not consistently improve performance, as the representational capacity of LoRA tends to saturate (Chen et al., 2022; Zhu et al., 2023). To overcome this limitation, recent methods integrate LoRA into the Mixture-of-Experts (MoE) architecture (Shazeer et al., 2017), giving rise to the *MoE-LoRA* framework (Zadouri et al., 2023; Zhu et al., 2023). In this setup, multiple structurally identical LoRA experts are dynamically routed at the token level via a learnable routing mechanism, thereby enhancing the model’s adaptability to diverse inputs and increasing its representational flexibility.

Although MoE-LoRA hybrids enhance representational capacity by introducing multiple LoRA experts, their homogeneous design causes the experts to tend toward learning similar representations during training, leading to **representation collapse** (Chi et al., 2022; Wang et al., 2024; Lin et al., 2024). This undermines expert diversity and specialization. Moreover, the dynamic routing mechanism is prone to **expert load imbalance** in homogeneous architectures, where a few initially well-performing experts consistently receive the majority of tokens, suppressing the participation

^{*}Corresponding author

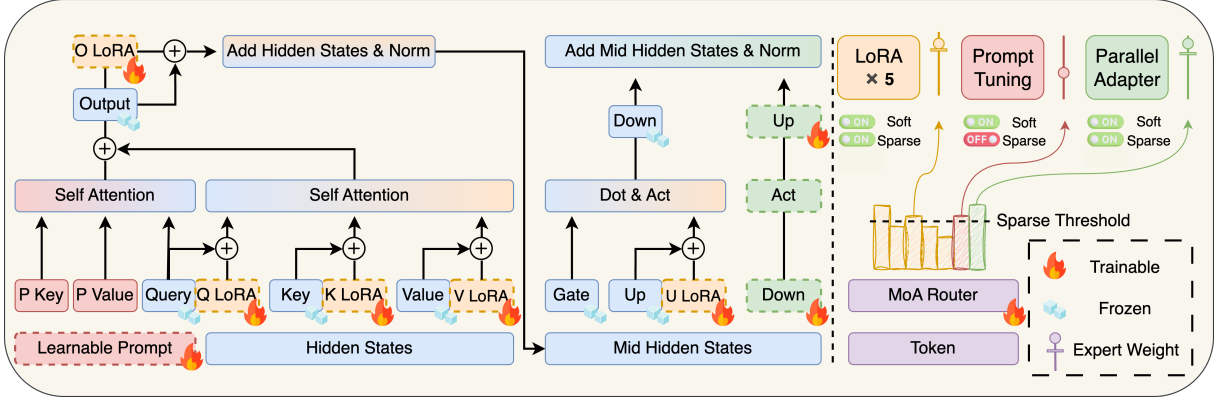


Figure 1: MoA architecture with heterogeneous PEFT adapters. It is worth noting that in Sparse MoA, the Prompt Tuning module is deactivated due to its non-token-level activation mechanism.

of other experts. These mechanisms ultimately result in redundant experts (Zhou et al., 2022), which limit the expressiveness of the model and reduce the efficiency of resource utilization.

Based on the above analysis, we posit that the representational convergence of homogeneous experts may limit the performance ceiling of PEFT methods. To address the issues of redundant experts and load imbalance in MoE-LoRA, we propose a **Mixture-of-Adapters (MoA)** approach with heterogeneous structures (Figure 1). MoA consists of PEFT adapter experts with diverse architectural forms and employs token-level dynamic routing to activate experts on demand. By leveraging the complementary representational capacities of different structures, MoA promotes expert specialization through functional differentiation among experts. This mechanism enables more refined utilization of pre-trained knowledge and enhances the ability of models to transfer and generalize across diverse downstream tasks. Our main contributions are as follows:

(i) We propose MoA based on heterogeneous experts, which constructs adapters with complementary representational capacities by integrating structurally diverse PEFT modules. This effectively enhances expert specialization, mitigating the issue of expert redundancy commonly observed in traditional MoE-LoRA architectures. In addition, MoA achieves efficient adaptation to downstream tasks using fewer trainable parameters, fully leveraging the knowledge representations embedded in the pre-trained model.

(ii) Building on this foundation, we further design two variants: **Soft MoA** and **Sparse MoA**. In Sparse MoA, a threshold function dynamically selects active experts that valuable contribute to the

current input. This reduces redundant computation while enabling the model to engage more experts for critical tokens, thereby improving representational capacity.

(iii) We systematically evaluate the architectural advantages and practical effectiveness of MoA across multiple tasks. Experimental results show that MoA significantly outperforms SoTA MoE-LoRA methods in GPU memory usage, training efficiency, and inference speed. Moreover, it achieves **higher parameter efficiency and knowledge transfer capability while maintaining or even improving downstream performance**.

2 Related Works

2.1 Parameter-Efficient MoE

In non-PEFT scenarios, standard sparse MoE (Shazeer et al., 2017; Fedus et al., 2022; Jiang et al., 2024) architectures must employ discrete top-k routing strategies to reduce computational overhead. In contrast, under PEFT settings, where PEFT experts occupy significantly less parameters than the pre-trained LLMs, it becomes feasible to compute a soft-weighted output by averaging experts’ predictions according to the router’s distribution. Thus, we categorize PEFT-based MoE into two paradigms: **sparse parameter efficient MoE** and **soft-weighted parameter efficient MoE**.

Sparse parameter efficient MoE Building on the demonstrated superiority of standard sparse MoE architectures, recent works has integrated sparse MoE with parameter-efficient fine-tuning (PEFT) methods to enhance resource-constrained adaptation performance. Most of these approaches are based on LoRA, leveraging its superior efficacy in parameter-efficient adaptation as a foun-

dational component. SiRA (Zhu et al., 2023) enforces top- k experts routing with per-expert capacity constraints, limiting the maximum tokens each expert processes, and uses dropout to the gating network to mitigate overfitting. MoELoRA (Luo et al., 2024) employs contrastive learning to encourage experts to learn different features, thus mitigating the random routing phenomenon observed in MoE. To reduce redundancy, MoLA (Gao et al., 2024) explores layer-wise expert redundancy by assigning a different number of LoRA experts to each Transformer layer, revealing that the lower layers exhibit higher redundancy in the sparse MoE-LoRA frameworks. Advancing this line, AdaMoLE (Liu and Luo, 2024) dynamically adjusts the activation of LoRA experts for each input token via a threshold network integrated into the top- k routing.

Soft-weighted parameter efficient MoE MoLoRA (Zadouri et al., 2023) proposes a soft-weighted MoE-LoRA method with a token-level routing strategy. LoRAMoE (Dou et al., 2024) integrates LoRA experts using a linear router to alleviate the forgetting of the world knowledge previously stored in LLMs. MOLE (Wu et al., 2024) proposes to efficiently compose multiple trained LoRAs while preserving all their characteristics by training only the router within each layer. HydraLoRA (Tian et al., 2024) introduces an asymmetric structure with a shared matrix, further enhancing parameter efficiency compared to other MoE-LoRA methods. Diverging from methods that ensemble the outputs of individual experts, certain approaches such as AdaMix (Wang et al., 2022), SMEAR (Muqeeth et al., 2024), and MoV (Zadouri et al., 2023) adopt a strategy of first merging expert parameters via routing, followed by computation using the merged parameters. This aims to reduce computational overhead.

2.2 Heterogeneous MoE

AutoMoE (Jawahar et al., 2023) employs the neural architecture search to obtain heterogeneous MoE subnetworks from original FFN experts on a small-scale MoE model that utilizes the top-1 routing strategy. HMoE (Wang et al., 2024) constructs heterogeneous experts by varying the dimensions of FFN experts, thereby endowing them with diverse capacities. UNIPEFT (Mao et al., 2022) employs distinct gated networks for various PEFT methods in instance-level to adapt to different examples.

3 Method

Homogeneous MoE methods inherently suffer from the representation collapse and load imbalance issues, fundamentally stemming from their experts’ identical structures and capacities, which lead to insufficient specialization of individual experts. To address this, we construct inherently heterogeneous MoE experts leveraging structurally diverse and positionally varied PEFT adapters for each Transformer block.

LoRA, Parallel Adapter, and Prompt Tuning approaches can be simplified into a unified view (He et al., 2022):

$$\mathbf{h} = F(\mathbf{x}) + E(\mathbf{x}), \quad (1)$$

where $F(\mathbf{x})$ denotes the transformer module to which the corresponding PEFT Adapter E is applied. Our selection of heterogeneous PEFT experts comprises LoRA (Hu et al., 2021), Parallel Adapters (He et al., 2022), and the zero-initialized Prompt Tuning (Zhang et al., 2023), the details of these methods and the unified view are shown in Appendix A.

3.1 Soft MoA

In each transformer layer, our Soft MoA comprises a soft-weighted router R and a set of experts $\{E_1, E_2, \dots, E_n\}$ that are composed of PEFT adapters of different types and placements. Our router network commonly consists of a dense layer with trainable weights $\mathbf{W}_r \in \mathbb{R}^{d \times n}$ followed by a sigmoid function which takes an intermediate token representation \mathbf{x} as input and output weights $R(\mathbf{x})$ for experts:

$$R(\mathbf{x}) = \text{Sigmoid}(\mathbf{W}_r^T \mathbf{x}). \quad (2)$$

Then, each expert’s weight $R(\mathbf{x})_i$ is applied to its output $E_i(\mathbf{x})$, controlling the influence of expert E_i in that transformer layer:

$$\mathbf{h} = F_i(\mathbf{x}) + R(\mathbf{x})_i E_i(\mathbf{x}). \quad (3)$$

Unlike previous MoE methods (Zadouri et al., 2023; Dou et al., 2024), our MoA leverages heterogeneous experts of various types within the Transformer module to process each token. Due to the diverse characteristics and capacities of these heterogeneous experts, MoA neither suffer from the representation collapse problem nor need a load-balancing loss. Moreover, MoA employs a sigmoid function in the router, rather than the traditional softmax activation used in conventional

MoE (Shazeer et al., 2017; Fedus et al., 2022), encouraging cooperation among experts instead of competition and fully capitalizing on their unique specializations.

3.2 Sparse MoA

Methods based on soft-weighted routing, including MoA and other MoE approaches (Zadouri et al., 2023; Dou et al., 2024; Muqeeth et al., 2024), require computing experts E_i even when their weights $R(x)_i \approx 0$, despite negligible contributions to outputs, resulting in unnecessary computational overhead and time consumption. To address this issue, we propose Sparse MoA, which dynamically computes only the necessary experts, thereby avoiding the computation of low-contributing ones.

To implement sparse MoA, a straightforward approach is to apply a fixed threshold Γ to the experts' weights. An expert E_i is activated for computation if and only if its weight satisfies $R(x)_i > \Gamma$, where unselected experts are entirely excluded from the computational graph, thus eliminating redundant computations. The detailed computation is defined as follows:

$$h = \begin{cases} F_i(x) + R(x)_i E_i(x) & \text{if } R(x)_i > \Gamma, \\ F_i(x) & \text{otherwise.} \end{cases} \quad (4)$$

Note that when the expert E_i is selected, its weight $R(x)_i$ is still used to control its influence. The threshold-based approach can activate varying numbers of heterogeneous experts for each token, dynamically adapting to the current task and data.

However, due to the varying semantic importance of tokens in contexts (Xia et al., 2025), assigning token-specific thresholds can more effectively harness multi-expert collaboration: lower thresholds for critical tokens promote broader expert engagement, while higher thresholds for less critical tokens restrict processing to fewer specialized experts. Inspired by AdaMoLE (Liu and Luo, 2024), we replace fixed thresholds with a learnable threshold function. Specifically, the threshold for each token is calculated via a linear layer followed by a sigmoid activation:

$$\Gamma = \Gamma_{max} \text{Sigmoid}(\mathbf{W}_{\Gamma}^T \mathbf{x} + \mathbf{b}_{\Gamma}), \quad (5)$$

where $\mathbf{W}_{\Gamma}^T \in \mathbb{R}^{d \times 1}$ and $\mathbf{b}_{\Gamma} \in \mathbb{R}^1$ are trainable parameters, Γ_{max} is a hyperparameter that defines the upper bound of threshold value.

As the core objective of PEFT methods is to fully leverage LLMs' pre-trained knowledge and capa-

bilities via adapters for downstream tasks (Ghosh et al., 2024), Sparse MoA enhances this efficiency by adaptively activating diverse combinations of heterogeneous experts. Moreover, Sparse MoA dynamically adjusts the number of participating experts based on the contextual importance of tokens, thereby mitigating redundant expert computations. For instance, in extreme cases where certain tokens require no expert intervention and can be processed solely by the frozen LLM, Sparse MoA deactivates all experts, thus maximizing computational efficiency.

Notably, due to its computational mechanism, Prompt Tuning does not facilitate token-level expert allocation within a batch. Consequently, Prompt Tuning experts are not included in our Sparse MoA framework.

4 Experiments

4.1 Datasets

We conduct extensive experiments and analyses on both mathematical and commonsense reasoning tasks. For mathematical reasoning, we evaluate performance on six test datasets: GSM8K (Cobbe et al., 2021), SVAMP (Patel et al., 2021), Multi-Arith (Roy and Roth, 2016), AddSub (Hosseini et al., 2014), AQuA (Ling et al., 2017), and SingleEq (Koncel-Kedziorski et al., 2015). The training dataset is Math14K constructed by (Hu et al., 2023), which includes the training sets of GSM8K and AQuA, augmented with rationales generated by ChatGPT and GPT-4. For commonsense reasoning, we test on eight datasets: BoolQ (Clark et al., 2019), PIQA (Bisk et al., 2020), SIQA (Sap et al., 2019), HellaSwag, WinoGrande (Sakaguchi et al., 2021), ARC Challenge (ARC-c), ARC Easy (ARC-e) (Chollet, 2019), and OBQA (Mihaylov et al., 2018). The training dataset for commonsense reasoning, Commonsense15K, is sampled from Commonsense170K (Hu et al., 2023) which comprises eight training sets.

4.2 Baseline Models

To validate the efficacy of our method, we compare it against: (i) the standard LoRA (Hu et al., 2021); (ii) SoTA soft-weighted parameter-efficient MoE methods: MoLoRA (Zadouri et al., 2023) and HydraLoRA (Tian et al., 2024); and (iii) SoTA sparse parameter-efficient MoE approaches: MoLA (Gao et al., 2024) and AdaMoLE (Liu and Luo, 2024).

Models	AddSub	AQuA	GSM8k	MultiArith	SingleEq	SVAMP	Average	Param
LoRA	87.68 \pm 1.14	36.75 \pm 0.99	76.14 \pm 0.88	96.28 \pm 0.79	96.98 \pm 0.30	81.27 \pm 0.45	79.18 \pm 0.26	23.06M
MoLoRA	91.73 \pm 1.02	38.06 \pm 1.86	77.46 \pm 0.68	99.33 \pm 0.17	97.18 \pm 0.41	81.77 \pm 0.68	80.92 \pm 0.19	100.14M
HydraLoRA	91.22 \pm 0.15	39.11 \pm 2.17	77.51 \pm 0.59	99.17 \pm 0.33	96.85 \pm 0.34	81.87 \pm 0.25	80.95 \pm 0.28	45.09M
MoLA	90.46 \pm 0.53	39.37 \pm 1.80	77.61 \pm 0.61	98.94 \pm 0.19	97.05 \pm 0.20	82.27 \pm 0.60	80.95 \pm 0.14	100.14M
AdaMoLE	90.97 \pm 0.53	40.03 \pm 2.37	77.51 \pm 0.27	98.33 \pm 2.05	97.31 \pm 0.30	82.40 \pm 0.50	81.09 \pm 0.78	101.12M
Soft MoA	92.24 \pm 0.89	38.98 \pm 1.57	78.01 \pm 0.23	98.94 \pm 0.54	97.44 \pm 0.71	83.43 \pm 0.55	81.51 \pm 0.35	24.52M
Sparse MoA	91.90 \pm 0.44	38.06 \pm 2.27	78.09 \pm 0.95	99.17 \pm 0.17	97.31 \pm 0.23	82.70 \pm 0.00	81.20 \pm 0.47	22.29M

Table 1: Experiment results on mathematical reasoning benchmarks. The evaluation metric is accuracy. All methods are run with three random seeds; mean and standard deviation are reported. “Param” indicates trainable parameters.

Models	BoolQ	PIQA	SIQA	HellaSwag	WinoGrande	ARC-C	ARC-E	OBQA	Average	Param
LoRA	72.94 \pm 0.29	86.22 \pm 0.64	78.78 \pm 0.64	86.34 \pm 1.37	80.01 \pm 1.87	82.08 \pm 0.51	92.63 \pm 0.98	86.93 \pm 0.81	83.24 \pm 0.58	23.06M
MoLoRA	73.89 \pm 0.77	87.32 \pm 0.73	78.98 \pm 0.96	88.52 \pm 0.58	82.11 \pm 0.41	84.84 \pm 1.04	93.45 \pm 0.45	87.33 \pm 0.99	84.56 \pm 0.58	100.14M
HydraLoRA	73.52 \pm 0.25	86.65 \pm 0.65	79.17 \pm 0.42	87.46 \pm 1.00	82.30 \pm 0.39	83.87 \pm 0.45	93.28 \pm 0.49	86.73 \pm 1.50	84.09 \pm 0.31	45.09M
MoLA	73.19 \pm 0.42	87.01 \pm 0.44	79.32 \pm 0.75	87.60 \pm 0.67	82.53 \pm 0.90	84.95 \pm 0.42	93.35 \pm 0.33	87.60 \pm 0.92	84.45 \pm 0.22	100.14M
AdaMoLE	73.60 \pm 0.37	88.08 \pm 0.61	79.72 \pm 0.48	88.70 \pm 0.65	81.29 \pm 1.14	85.04 \pm 0.81	93.27 \pm 0.19	89.13 \pm 0.81	84.85 \pm 0.21	101.12M
Soft MoA	72.86 \pm 0.40	88.08 \pm 0.55	79.32 \pm 0.72	87.82 \pm 0.44	83.03 \pm 0.70	85.84 \pm 0.31	93.49 \pm 0.46	89.27 \pm 0.50	84.96 \pm 0.26	24.52M
Sparse MoA	73.69 \pm 0.31	87.36 \pm 0.27	78.51 \pm 0.23	87.97 \pm 0.24	81.66 \pm 0.91	85.55 \pm 0.13	93.53 \pm 0.09	88.67 \pm 0.58	84.62 \pm 0.09	22.29M

Table 2: Experiment results on commonsense reasoning benchmarks. The evaluation metric is accuracy.

4.3 Experiment Settings

All experiments use the LLaMA-3.1 8B model (Grattafiori et al., 2024b). To ensure fair comparisons, we apply identical configurations to all LoRA-based methods: LoRA is inserted into four weight matrices in the self-attention module (W_q , W_k , W_v , W_o) and one in the FFN module (W_{up}), with rank 8 and alpha 8. For all MoE-LoRA baselines, we set the number of experts to 8. For the original LoRA baseline, we increase the rank to 16 to match the number of trainable parameters across methods. MoA integrates 7 heterogeneous experts, including five LoRA modules, parallel adapters in the FFN layer, and a zero-initialized prompt-tuning module. The bottleneck dimension of the parallel adapter is set to 16, and the prompt length to 10. We use the AdamW optimizer with a learning rate of 6e-3. For commonsense reasoning, input sequences are truncated to 200 tokens; for mathematical reasoning, to 300 tokens. All models are trained with a batch size of 32 on 2xA6000 GPUs.

4.4 Experiment Results

As shown in Tables 1 and 2, Soft MoA and Sparse MoA consistently outperform homogeneous MoE-LoRA baselines on both mathematical and commonsense reasoning tasks. Soft MoA achieves the highest accuracy on math benchmarks (81.51%) with only 24.52M trainable parameters—nearly 4 \times fewer than AdaMoLE and MoLoRA. Sparse MoA delivers strong performance (81.20% math, 84.62% commonsense) with the smallest parameter count

(22.29M), surpassing all other methods in mathematical accuracy. These results demonstrate that heterogeneous experts in **MoA enhance both performance and efficiency, effectively addressing redundancy in existing MoE-LoRA designs.**

5 Ablation Study

5.1 Ablation Study on Soft MoA Components

The token-level soft-weighted router effectively leverages the distinct characteristics of heterogeneous experts. Table 3 shows the performance of individual PEFT methods and their naive combination on two reasoning tasks. Among the three standalone PEFT methods, LoRA achieves the highest accuracy. Notably, the naive composition of these three PEFT methods degrades performance, failing to effectively leverage the functionalities of diverse experts. In contrast, our Soft MoA integrates seven heterogeneous experts via a token-level router, achieving significant improvements over standalone LoRA (+2.33 on mathematical reasoning, +1.72 on commonsense reasoning). This demonstrates that Soft MoA’s router dynamically assigns experts based on token-specific contextual demands, fully exploiting the complementary strengths of these structurally diverse modules.

The sigmoid activation function outperforms softmax in the MoA router. Conventional MoE methods with homogeneous experts typically use routers with softmax activation, which enforces weight trade-offs over experts ($\sum_i R(x)_i = 1$). For comparison, we replaced the sigmoid func-

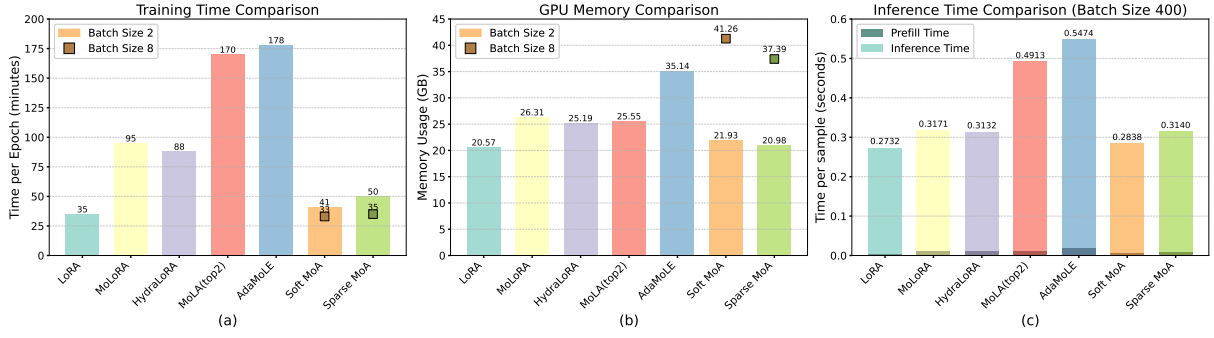


Figure 2: Comparison of different models under identical batch sizes (1* 48G GPU): (a) Training time per epoch, (b) GPU memory consumption during training, and (c) Average inference time per sample.

Modules	Router	Commonsense	Math
Prompt-tuning	✗	78.03	78.06
Parallel Adapter	✗	81.62	79.20
LoRA	✗	83.24	79.18
LoRA	Sigmoid R	84.23	80.50
Composition	✗	81.58	78.28
Composition	Softmax R	84.16	80.21
Composition	Sigmoid R	84.96	81.51

Table 3: Ablation study of each type of PEFT expert and the activation function of the Router in MoA.

Max Gamma	Threshold Function	Math
0.2	✗	81.22
0.5	✗	80.92
0.8	✗	78.58
0.2	✓	81.24
0.5	✓	81.51
0.8	✓	81.36
1	✓	80.61

Table 4: Comparative performance of fixed-threshold baseline and threshold function in Sparse MoA under varying hyperparameters Γ_{\max} .

tion with softmax in our framework. As shown in Table 3, while the softmax-based router improves upon standalone LoRA and the naive composition, it underperforms MoA equipped with a sigmoid router. This indicates a fundamental operational difference: homogeneous MoE experts function within a competitive paradigm, whereas MoA’s heterogeneous experts leverage collaborative computation (sigmoid permits non-exclusive activation). The sigmoid router dynamically scales expert contributions - enabling full-capacity collaboration ($R(x)_i \approx 1$ for all experts) or complete deactivation ($R(x)_i \approx 0$)—thereby maximizing the utility of architectural diversity.

The five LoRA modules contribute the most, while the Parallel Adapter and Prompt Tuning further improve performance. The superior performance of LoRA over Prompt-tuning and Parallel Adapter, combined with its inherent inclusion of positionally diverse expert modules, motivated our extension of LoRA’s multi-module experts with a token-level router. While this adaptation yielded significant improvements over standalone LoRA (e.g., +1.32 on mathematical reasoning), it underperformed the full MoA (-1.01), highlighting two

critical insights: 1. MoA’s primary strength derives from leveraging multiple LoRA modules as distinct experts. 2. The distinct mechanisms of Prompt Tuning and Parallel Adapter further complement these LoRA-based experts.

5.2 Effect of Threshold Function

To validate the effect of the threshold function in Sparse MoA, we compare it against fixed-threshold baselines under varying hyperparameters $\Gamma_{\max} \in \{0.2, 0.5, 0.8\}$, with an additional setting of $\Gamma_{\max} = 1$ exclusive to the threshold function. As shown in Table 4, the threshold function achieves higher performance ceilings and consistently outperforms fixed-threshold methods across all Γ_{\max} configurations. This demonstrates that tokens exhibit varying contextual importance: the dynamic threshold function assigns lower thresholds to critical tokens (promoting multi-expert collaboration) and higher thresholds to trivial ones (reducing redundant computations), thereby better leveraging the LLM’s capabilities while maintaining efficiency.

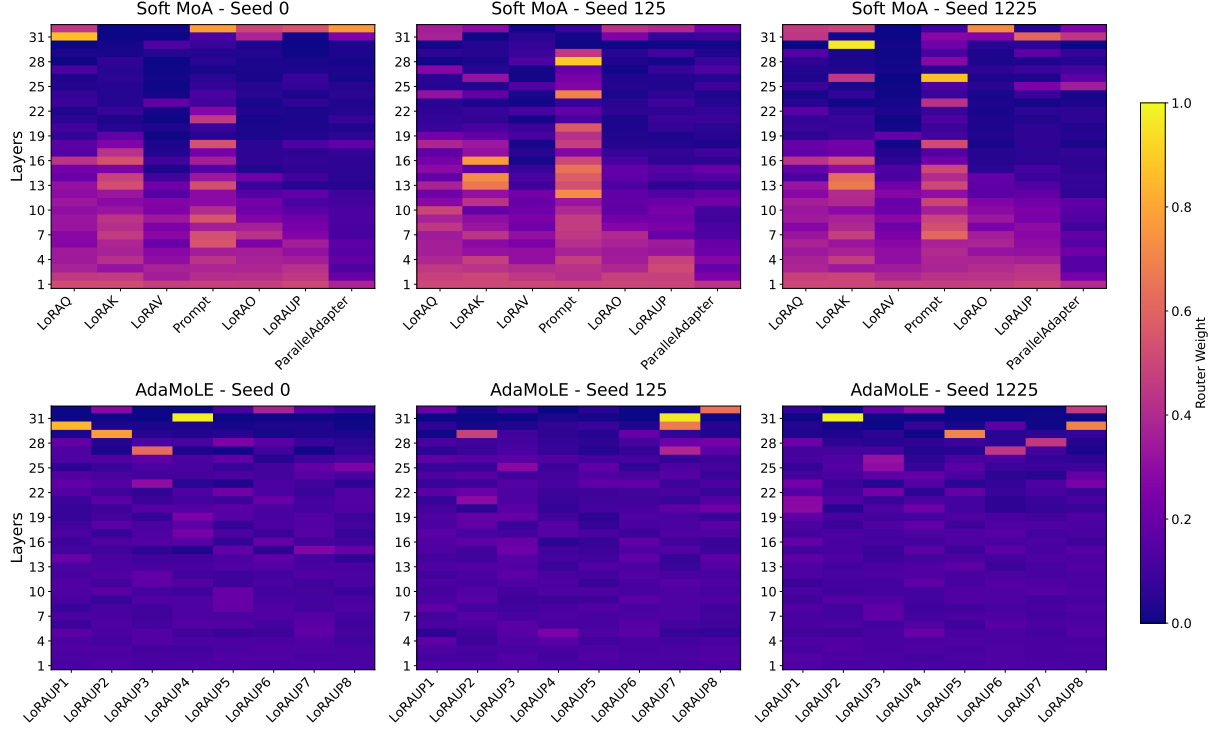


Figure 3: Comparison of router weight distributions between Soft MoA and AdaMoLE under different random seeds. The MoA method exhibits strong consistency, whereas AdaMoLE does not.

6 In-depth Analysis

6.1 Efficiency Comparison

To systematically compare the differences between Soft MoA, Sparse MoA, and baseline methods, Figure 2 presents the training time per epoch, GPU memory requirements, and average inference latency per sample under identical configurations (single GPU setup, fixed batch size (BS)).

As shown in Figure 2(a), under identical batch size conditions (BS=2), both Soft MoA and Sparse MoA require significantly less training time compared to other MoE methods. Notably, Soft MoA achieves less than half the training time of MoLoRA and HydraLoRA. This efficiency advantage stems not only from reduced parameter counts but also from architectural optimizations - for instance, while HydraLoRA contains half the parameters of MoLoRA, their training times remain comparable. Sparse MoE methods (MoLA and AdaMoLE) exhibit nearly double the training time of MoLoRA despite having the same parameter counts. This discrepancy occurs because, unlike traditional MoE approaches, PEFT-based MoE methods feature experts with significantly fewer parameters and computations relative to the frozen LLM backbone; besides, the sparse token-to-expert

routing process within each batch brings additional computational overhead. This phenomenon is also observable between Soft MoA and Sparse MoA. However, as the batch size increases (e.g., BS=8 in Figure 2(a)), two key effects emerge: (i) the sparse per-batch routing overhead becomes relatively less significant, and (ii) the reduction strategy of redundant expert computation in Sparse MoA becomes more effective, ultimately causing its training time to converge with Soft MoA's.

As shown in Figure 2(b), both Soft MoA and Sparse MoA maintain memory requirements comparable to LoRA, with Sparse MoA demonstrating additional memory efficiency through its dynamic threshold function that selectively activates experts per token. This advantage becomes more pronounced at larger batch sizes - for instance, at BS=8, Sparse MoA reduces memory consumption by 3.86GB compared to Soft MoA. Figure 2(c) reveals that the inference latency of both Soft MoA and Sparse MoA approaches that of standalone LoRA, while maintaining significant advantages over other MoE methods.

In summary, both MoA variants demonstrate notable advantages over MoE-LoRA methods across training time, GPU memory usage, and inference latency, nearing the performance of LoRA. Specifi-

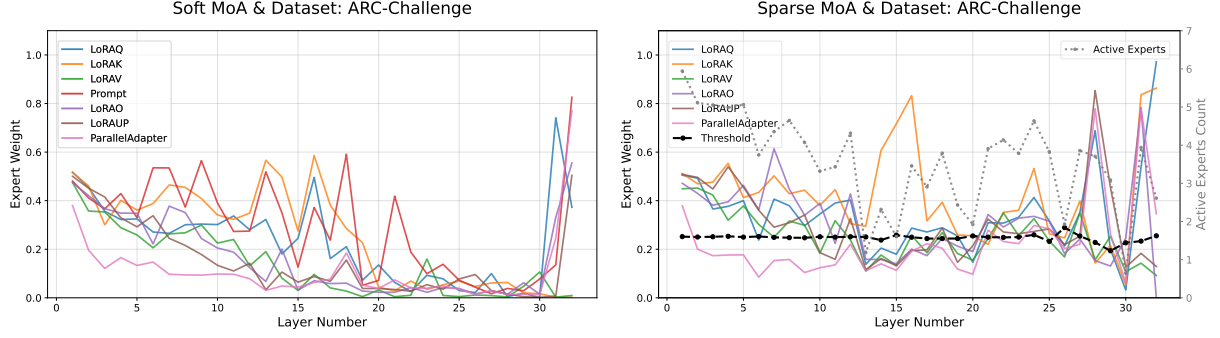


Figure 4: Visualization of average router weights per layer for Soft MoA (left) and Sparse MoA (right) on ARC-Challenge, averaged over tokens within 50 samples. Sparse MoA also includes the average per-layer threshold and the average count of activated experts. The average count of activated experts across layers in Sparse MoA is 3.55.

cally: (i) Sparse MoA exhibits superior memory efficiency compared to Soft MoA through its sparse token-to-expert routing strategy. (ii) The memory and computational benefits of Sparse MoA scale favorably with increasing batch size, making it particularly well-suited for resource-intensive applications.

6.2 Difference between MoA and MoE

The heterogeneous experts in the MoA method ensure expert specialization. As shown in Figure 3, we compare the expert weight distributions across 32 Transformer layers on 50 test samples for Soft MoA and AdaMoLE, trained with different random seeds. Soft MoA exhibits highly consistent expert weight distributions across seeds, as illustrated in the first row. Experts in the lower 16 layers tend to be more active than those in the upper layers. Among the experts, LoRA(Q, K) and Prompt experts show higher activation than LoRA(O, V, UP) and the Parallel Adapter. This indicates that MoA is capable of leveraging the distinct adaptation capabilities of heterogeneous experts according to their structural characteristics and positions within the model. In contrast, as shown in the second row, AdaMoLE fails to exhibit such consistency under different seeds, aligning with findings from prior work such as HMoE (Wang et al., 2024), which identifies under-specialization as a common issue in homogeneous expert settings.

6.3 Expert Activation Comparison

Sparse MoA significantly reduces expert computations with nearly no accuracy loss relative to Soft MoA. Figures 4 and 5 (see Appendix) present a comparison of expert activation weights between Soft MoA and Sparse MoA on ARC-Challenge and GSM8K datasets. The

distribution patterns of router weights for experts exhibit notable similarities between Soft MoA and Sparse MoA. However, Soft MoA activates a fixed total number of experts per token, whereas Sparse MoA employs a thresholding function for selective expert activation. As illustrated in Figure 4 (right), the average number of activated experts in Sparse MoA drops to merely 1-2 in certain intermediate layers. The average number of activated experts across all layers is reduced from 6 to 3.55, reducing 40% expert computations with a negligible accuracy drop of 0.3%.

6.4 Instance- and Token-level Routing

In contrast to token-level routing approaches, instance-level routing, wherein the model makes a single routing decision for the entire input example, significantly reduces the computational overhead associated with the router. This brings the overall computational cost closer to that of standard, non-MoE architectures. We conducted a comparison between token-level MoA and instance-level variant of Soft MoA and the UniPEFT approach. The details are shown in Appendix B.1.

7 Conclusion

We present MoA, a novel heterogeneous parameter-efficient fine-tuning (PEFT) method that adapts LLMs to downstream tasks using PEFT adapters with diverse architectures and a minimal number of trainable parameters. Comprehensive experiments are conducted on commonsense and mathematical reasoning tasks, demonstrating that heterogeneous MoA outperforms the state-of-the-art homogeneous MoE-LoRA approaches, exhibiting superior performance in training time, inference latency, and GPU memory consumption.

Limitations

While our proposed Sparse MoA method reduces GPU memory consumption compared to Soft MoA, its training and inference times are longer at a small batch size due to the additional computational overhead of sparse token-to-expert routing, achieving comparable or reduced time consumption only at larger batch sizes. Furthermore, sparse routing methods, including Sparse MoA, necessitate token-to-expert assignment within a batch, thus limiting support to PEFT Adapters where input tokens within a sample can be computed independently. Consequently, methods like zero-initialized Prompt Tuning, where token computations within a sample are interdependent, are not compatible with Sparse MoA.

References

- Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, Florencia Leoni Aleman, Diogo Almeida, Janko Altenschmidt, Sam Altman, Shyamal Anadkat, and 1 others. 2023. Gpt-4 technical report. *arXiv preprint arXiv:2303.08774*.
- Yonatan Bisk, Rowan Zellers, Jianfeng Gao, Yejin Choi, and 1 others. 2020. Piqa: Reasoning about physical commonsense in natural language. In *Proceedings of the AAAI conference on artificial intelligence*, volume 34, pages 7432–7439.
- Guanzheng Chen, Fangyu Liu, Zaiqiao Meng, and Shangsong Liang. 2022. Revisiting parameter-efficient tuning: Are we really there yet? In *Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing*, pages 2612–2626.
- Zewen Chi, Li Dong, Shaohan Huang, Damai Dai, Shuming Ma, Barun Patra, Saksham Singhal, Payal Bajaj, Xia Song, Xian-Ling Mao, and 1 others. 2022. On the representation collapse of sparse mixture of experts. *Advances in Neural Information Processing Systems*, 35:34600–34613.
- François Chollet. 2019. On the measure of intelligence. *arXiv preprint arXiv:1911.01547*.
- Christopher Clark, Kenton Lee, Ming-Wei Chang, Tom Kwiatkowski, Michael Collins, and Kristina Toutanova. 2019. Boolq: Exploring the surprising difficulty of natural yes/no questions. *arXiv preprint arXiv:1905.10044*.
- Karl Cobbe, Vineet Kosaraju, Mohammad Bavarian, Mark Chen, Heewoo Jun, Lukasz Kaiser, Matthias Plappert, Jerry Tworek, Jacob Hilton, Reiichiro Nakano, and 1 others. 2021. Training verifiers to solve math word problems. *arXiv preprint arXiv:2110.14168*.
- 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](#). *arXiv preprint*. ArXiv:2312.09979 [cs].
- William Fedus, Barret Zoph, and Noam Shazeer. 2022. [Switch Transformers: Scaling to Trillion Parameter Models with Simple and Efficient Sparsity](#). *arXiv preprint*. ArXiv:2101.03961 [cs].
- Chongyang Gao, Kezhen Chen, Jinmeng Rao, Baochen Sun, Ruibo Liu, Daiyi Peng, Yawen Zhang, Xiaoyuan Guo, Jie Yang, and V. S. Subrahmanian. 2024. [Higher Layers Need More LoRA Experts](#). *arXiv preprint*. ArXiv:2402.08562 [cs].
- Sreyan Ghosh, Chandra Kiran Reddy Evuru, Sonal Kumar, Ramaneswaran S, Deepali Aneja, Zeyu Jin, Ramani Duraiswami, and Dinesh Manocha. 2024. [A Closer Look at the Limitations of Instruction Tuning](#). *arXiv preprint*. ArXiv:2402.05119 [cs].
- Aaron Grattafiori, Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad Al-Dahle, Aiesha Letman, Akhil Mathur, Alan Schelten, Alex Vaughan, and 1 others. 2024a. The llama 3 herd of models. *arXiv preprint arXiv:2407.21783*.
- Aaron Grattafiori, Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad Al-Dahle, Aiesha Letman, Akhil Mathur, Alan Schelten, Alex Vaughan, and 1 others. 2024b. The llama 3 herd of models. *arXiv preprint arXiv:2407.21783*.
- Zeyu Han, Chao Gao, Jinyang Liu, Sai Qian Zhang, and 1 others. 2024. Parameter-efficient fine-tuning for large models: A comprehensive survey. *arXiv preprint arXiv:2403.14608*.
- Junxian He, Chunting Zhou, Xuezhe Ma, Taylor Berg-Kirkpatrick, and Graham Neubig. 2022. [Towards a Unified View of Parameter-Efficient Transfer Learning](#). *arXiv preprint*. ArXiv:2110.04366 [cs].
- Mohammad Javad Hosseini, Hannaneh Hajishirzi, Oren Etzioni, and Nate Kushman. 2014. Learning to solve arithmetic word problems with verb categorization. In *Proceedings of the 2014 conference on empirical methods in natural language processing (EMNLP)*, pages 523–533.
- N. Houlsby, A. Giurgiu, Stanislaw Jastrzebski, Bruna Morrone, Quentin de Laroussilhe, Andrea Gesmundo, Mona Attariyan, and S. Gelly. 2019. [Parameter-Efficient Transfer Learning for NLP](#). *ArXiv*.
- Edward J. Hu, Yelong Shen, Phillip Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang, and Weizhu Chen. 2021. [LoRA: Low-Rank Adaptation of Large Language Models](#). *arXiv preprint*. ArXiv:2106.09685 [cs].

- Zhiqiang Hu, Lei Wang, Yihuai Lan, Wanyu Xu, Ee-Peng Lim, Lidong Bing, Xing Xu, Soujanya Poria, and Roy Ka-Wei Lee. 2023. [LLM-Adapters: An Adapter Family for Parameter-Efficient Fine-Tuning of Large Language Models](#). *arXiv preprint*. ArXiv:2304.01933 [cs].
- Ganesh Jawahar, Subhabrata Mukherjee, Xiaodong Liu, Young Jin Kim, Muhammad Abdul-Mageed, Laks V. S. Lakshmanan, Ahmed Hassan Awadallah, Sebastian Bubeck, and Jianfeng Gao. 2023. [AutoMoE: Heterogeneous Mixture-of-Experts with Adaptive Computation for Efficient Neural Machine Translation](#). *arXiv preprint*. ArXiv:2210.07535 [cs].
- Albert Q Jiang, Alexandre Sablayrolles, Antoine Roux, Arthur Mensch, Blanche Savary, Chris Bamford, Devendra Singh Chaplot, Diego de las Casas, Emma Bou Hanna, Florian Bressand, and 1 others. 2024. Mixture of experts. *arXiv preprint arXiv:2401.04088*.
- Rik Koncel-Kedziorski, Hannaneh Hajishirzi, Ashish Sabharwal, Oren Etzioni, and Siena Dumas Ang. 2015. Parsing algebraic word problems into equations. *Transactions of the Association for Computational Linguistics*, 3:585–597.
- Brian Lester, Rami Al-Rfou, and Noah Constant. 2021. [The Power of Scale for Parameter-Efficient Prompt Tuning](#). In *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing*, pages 3045–3059, Online and Punta Cana, Dominican Republic. Association for Computational Linguistics.
- 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.
- Tianwei Lin, Jiang Liu, Wenqiao Zhang, Zhaocheng Li, Yang Dai, Haoyuan Li, Zhelun Yu, Wanggui He, Juncheng Li, Hao Jiang, and 1 others. 2024. Teamlora: Boosting low-rank adaptation with expert collaboration and competition. *arXiv preprint arXiv:2408.09856*.
- Tianwei Lin, Wenqiao Zhang, Sijing Li, Yuqian Yuan, Binhe Yu, Haoyuan Li, Wanggui He, Hao Jiang, Mengze Li, Xiaohui Song, and 1 others. 2025. Healthgpt: A medical large vision-language model for unifying comprehension and generation via heterogeneous knowledge adaptation. *arXiv preprint arXiv:2502.09838*.
- Wang Ling, Dani Yogatama, Chris Dyer, and Phil Blunsom. 2017. Program induction by rationale generation: Learning to solve and explain algebraic word problems. *arXiv preprint arXiv:1705.04146*.
- Zefang Liu and Jiahua Luo. 2024. [AdaMoLE: Fine-Tuning Large Language Models with Adaptive Mixture of Low-Rank Adaptation Experts](#). *arXiv preprint*. ArXiv:2405.00361 [cs].
- 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](#). *arXiv preprint*. ArXiv:2402.12851 [cs].
- Yuning Mao, Lambert Mathias, Rui Hou, Amjad Almahairi, Hao Ma, Jiawei Han, Wen-tau Yih, and Madian Khabisa. 2022. [UniPELT: A Unified Framework for Parameter-Efficient Language Model Tuning](#). *arXiv preprint*. ArXiv:2110.07577 [cs].
- 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. *arXiv preprint arXiv:1809.02789*.
- Mohammed Muqeeth, Haokun Liu, and Colin Raffel. 2024. [Soft Merging of Experts with Adaptive Routing](#). *arXiv preprint*. ArXiv:2306.03745 [cs].
- Arkil Patel, Satwik Bhattamishra, and Navin Goyal. 2021. Are nlp models really able to solve simple math word problems? *arXiv preprint arXiv:2103.07191*.
- Subhro Roy and Dan Roth. 2016. Solving general arithmetic word problems. *arXiv preprint arXiv:1608.01413*.
- Keisuke Sakaguchi, Ronan Le Bras, Chandra Bhagavatula, and Yejin Choi. 2021. Winogrande: An adversarial winograd schema challenge at scale. *Communications of the ACM*, 64(9):99–106.
- Maarten Sap, Hannah Rashkin, Derek Chen, Ronan LeBras, and Yejin Choi. 2019. Socialliqa: Commonsense reasoning about social interactions. *arXiv preprint arXiv:1904.09728*.
- Noam Shazeer, Azalia Mirhoseini, Krzysztof Maziarz, Andy Davis, Quoc Le, Geoffrey Hinton, and Jeff Dean. 2017. Outrageously large neural networks: The sparsely-gated mixture-of-experts layer. *arXiv preprint arXiv:1701.06538*.
- Chunlin Tian, Zhan Shi, Zhijiang Guo, Li Li, and Chengzhong Xu. 2024. [HydraLoRA: An Asymmetric LoRA Architecture for Efficient Fine-Tuning](#). *arXiv preprint*. ArXiv:2404.19245 [cs].
- An Wang, Xingwu Sun, Ruobing Xie, Shuaipeng Li, Jiaqi Zhu, Zhen Yang, Pinxue Zhao, J. N. Han, Zhanhui Kang, Di Wang, Naoaki Okazaki, and Chengzhong Xu. 2024. [HMoE: Heterogeneous Mixture of Experts for Language Modeling](#). *arXiv preprint*. ArXiv:2408.10681 [cs].
- Yaqing Wang, Sahaj Agarwal, Subhabrata Mukherjee, Xiaodong Liu, Jing Gao, Ahmed Hassan Awadallah, and Jianfeng Gao. 2022. [AdaMix: Mixture-of-Adaptations for Parameter-efficient Model Tuning](#).

In *Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing*, pages 5744–5760, Abu Dhabi, United Arab Emirates. Association for Computational Linguistics.

Xun Wu, Shaohan Huang, and Furu Wei. 2024. [Mixture of LoRA Experts](#). *arXiv preprint*. ArXiv:2404.13628 [cs].

Heming Xia, Yongqi Li, Chak Tou Leong, Wenjie Wang, and Wenjie Li. 2025. Tokenskip: Controllable chain-of-thought compression in llms. *arXiv preprint* arXiv:2502.12067.

An Yang, Anfeng Li, Baosong Yang, Beichen Zhang, Binyuan Hui, Bo Zheng, Bowen Yu, Chang Gao, Chengen Huang, Chenxu Lv, and 1 others. 2025. Qwen3 technical report. *arXiv preprint* arXiv:2505.09388.

Ted Zadori, Ahmet Üstün, Arash Ahmadian, Beyza Ermiş, Acyr Locatelli, and Sara Hooker. 2023. [Pushing Mixture of Experts to the Limit: Extremely Parameter Efficient MoE for Instruction Tuning](#). *arXiv preprint*. ArXiv:2309.05444 [cs].

Renrui Zhang, Jiaming Han, Chris Liu, Peng Gao, Aojun Zhou, Xiangfei Hu, Shilin Yan, Pan Lu, Hongsheng Li, and Yu Qiao. 2023. [LLaMA-Adapter: Efficient Fine-tuning of Language Models with Zero-init Attention](#). *arXiv preprint*. ArXiv:2303.16199 [cs].

Wenqiao Zhang, Tianwei Lin, Jiang Liu, Fangxun Shu, Haoyuan Li, Lei Zhang, He Wanggui, Hao Zhou, Zheqi Lv, Hao Jiang, and 1 others. 2024. Hyperllava: Dynamic visual and language expert tuning for multimodal large language models. *arXiv preprint* arXiv:2403.13447.

Yanqi Zhou, Tao Lei, Hanxiao Liu, Nan Du, Yanping Huang, Vincent Zhao, Andrew M Dai, Quoc V Le, James Laudon, and 1 others. 2022. Mixture-of-experts with expert choice routing. *Advances in Neural Information Processing Systems*, 35:7103–7114.

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](#). *arXiv preprint*. ArXiv:2311.09179 [cs].

A PEFT Methods And The Unified Form

LoRA (Hu et al., 2021): LoRA injects trainable low-rank matrices into transformer layers to approximate the weight updates. For a pre-trained weight matrix $\mathbf{W} \in \mathbb{R}^{d \times k}$, LoRA represents its update with a low-rank decomposition $\mathbf{W} + \Delta\mathbf{W} = \mathbf{W} + \mathbf{W}_{down}\mathbf{W}_{up}$, where $\mathbf{W}_{down} \in \mathbb{R}^{d \times r}$, $\mathbf{W}_{up} \in \mathbb{R}^{r \times k}$ are tunable parameters. For a specific input \mathbf{x} to the linear projection in the transformer layer, the computation of LoRA is as follows:

$$\mathbf{h} = \mathbf{x}\mathbf{W} + \alpha \cdot \mathbf{x}\mathbf{W}_{down}\mathbf{W}_{up}, \quad (6)$$

where α is a scalar hyperparameter.

Parallel Adapter (He et al., 2022): Parallel Adapters involve adding small adapter modules in parallel to the Feed-Forward Network (FFN) or Attention blocks within Transformer layers. The adapter module is normally moulded by a two-layer feed-forward neural network with a bottleneck: a down-projection with $\mathbf{W}_{down} \in \mathbb{R}^{d \times r}$ to project the input \mathbf{h} to a lower-dimensional space specified by bottleneck dimension r , followed by a nonlinear activation function $f(\cdot)$, an up-projection with $\mathbf{W}_{up} \in \mathbb{R}^{r \times d}$ to project back to the input size. The parallel adapter can be defined as:

$$\mathbf{h} = F(\mathbf{x}) + f(\mathbf{x}\mathbf{W}_{down})\mathbf{W}_{up}, \quad (7)$$

where $F(\cdot)$ represents FFN or Attention block here.

Prompt tuning (Zhang et al., 2023): An advanced zero-initialized prompt tuning method adaptively incorporates instructional signals while preserving the pre-trained knowledge in LLMs. In each transformer layer, let $\mathbf{P} \in \mathbb{R}^{K \times d}$ denote the learnable prompts, the zero-initialized prompt tuning in the attention block computed as:

$$\mathbf{Q} = \mathbf{x}\mathbf{W}_q, \quad (8)$$

$$\mathbf{h} = \text{Attn}(\mathbf{Q}, \mathbf{x}\mathbf{W}_k, \mathbf{x}\mathbf{W}_v) + g \cdot \text{Attn}(\mathbf{Q}, \mathbf{P}\mathbf{W}_k, \mathbf{P}\mathbf{W}_v), \quad (9)$$

where g is a learnable scalar for each head initialized with zero, $\mathbf{W}_q, \mathbf{W}_k, \mathbf{W}_v$ are frozen pre-trained matrix. Note that the first term in Equation 9 is the original attention without prompts; only the second term includes trainable parameters.

The unified form: As demonstrated above, LoRA, Parallel Adapter, and Prompt Tuning approaches can be simplified into a unified form:

$$\mathbf{h} = F(\mathbf{x}) + E(\mathbf{x}), \quad (10)$$

where $F(\mathbf{x})$ denotes the transformer module to which the corresponding PEFT Adapter E is applied.

B Additional Experiments and Analyses

B.1 Comparison between Instance-level and Token-level Routing Methods

We conducted comparative experiments on LLaMA-3.1 8B, evaluating an instance-level variant of our Soft MoA method and the UniPEFT approach. UniPEFT utilizes a separate gated network for each type of adapter to control their influence on the input sample. As presented in Table 5, the results indicate that both instance-level approaches yield performance inferior to token-level routing methods. Furthermore, their performance is also surpassed by employing a standard LoRA adaptation alone.

This performance deficit can be attributed to key limitations of instance-level routing in PEFT of LLM. Firstly, these methods typically rely on the mean pooling of input text hidden states as the router’s input, which fails to capture fine-grained characteristics of the sample adequately. Secondly, the constraint of making a single routing choice for the entire sequence precludes the possibility of dynamically adjusting the contributions of different experts in response to the diverse contextual nuances within the input.

Model	Routing level	Commonsense15K	Math14K
UniPEFT	instance	78.29	74.65
Soft MoA	instance	73.18	76.11
Soft MoA	token	84.96	81.51
Sparse MoA	token	84.62	81.20

Table 5: Comparison between the token-level MoA approaches and instance-level methods. The instance-level methods include an instance-level variant of Soft MoA and the UniPEFT method.

B.2 Case Study of Router Weight

Figures 6 and 7 illustrate the router weights for a specific sample in a particular layer for Soft MoA and Sparse MoA models, respectively.

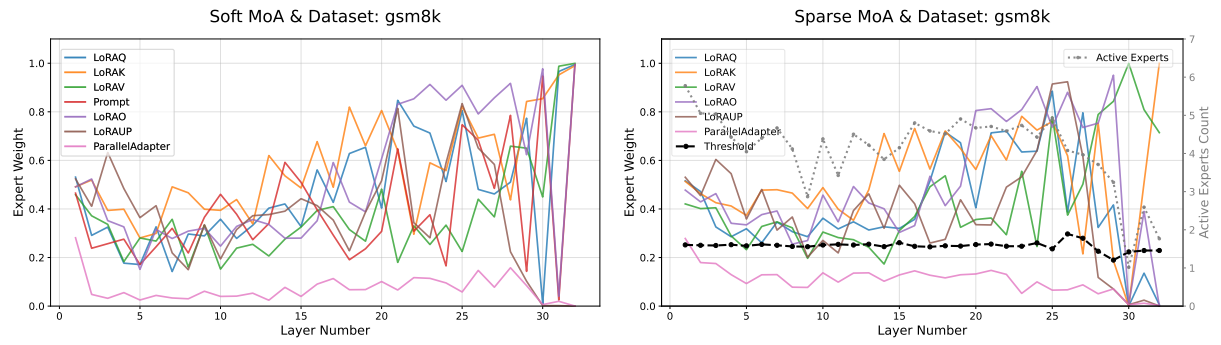


Figure 5: Visualization of average router weights per layer for Soft MoA (left) and Sparse MoA (right) on gsm8k, averaged over tokens within 50 samples. The Sparse MoA plot (right) also includes the average per-layer threshold and the average count of activated experts. The average count of activated experts across layers in Sparse MoA is 4.13.

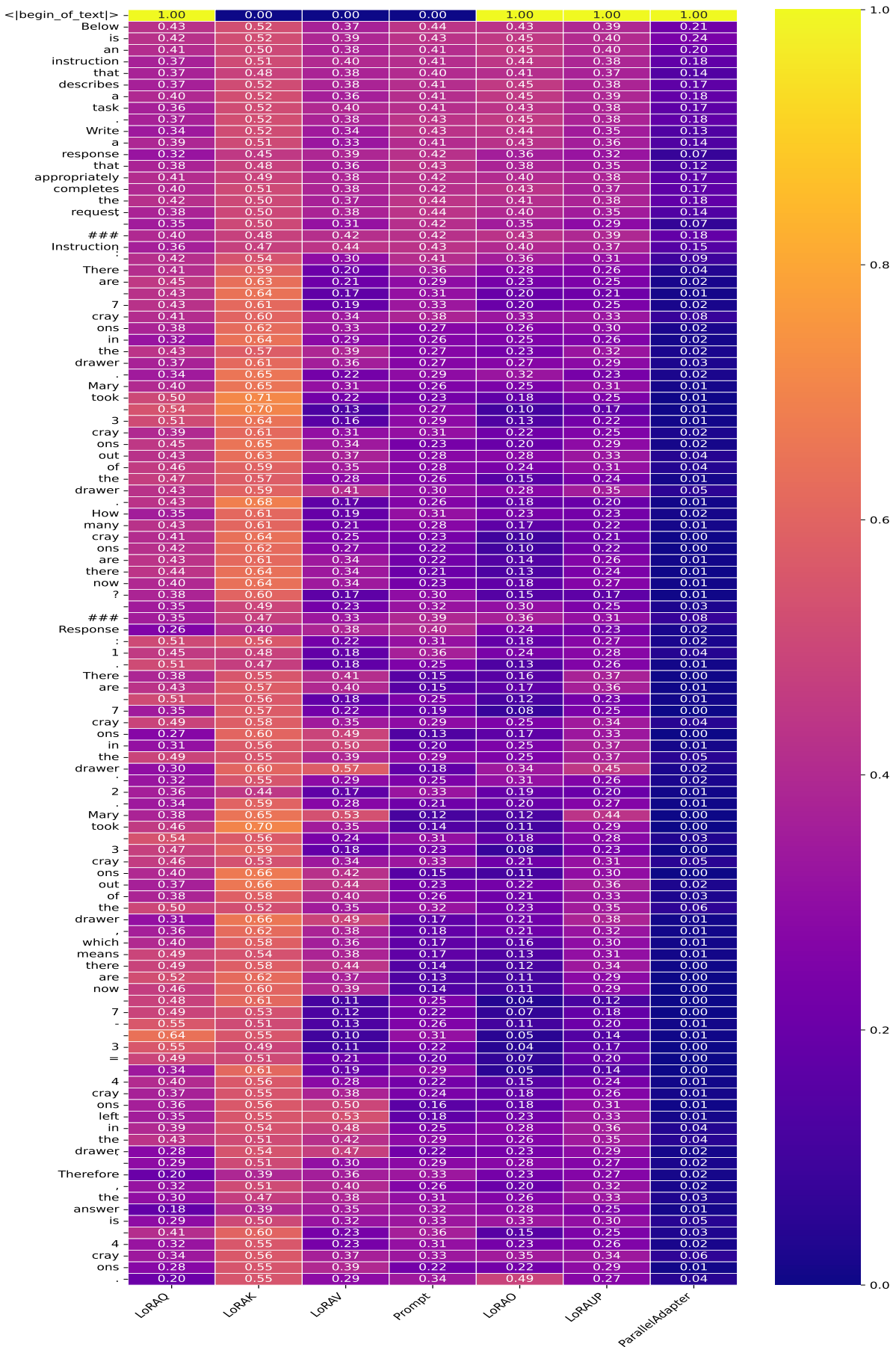


Figure 6: Router Weights of an example of at Layer 14 in Soft MoA.

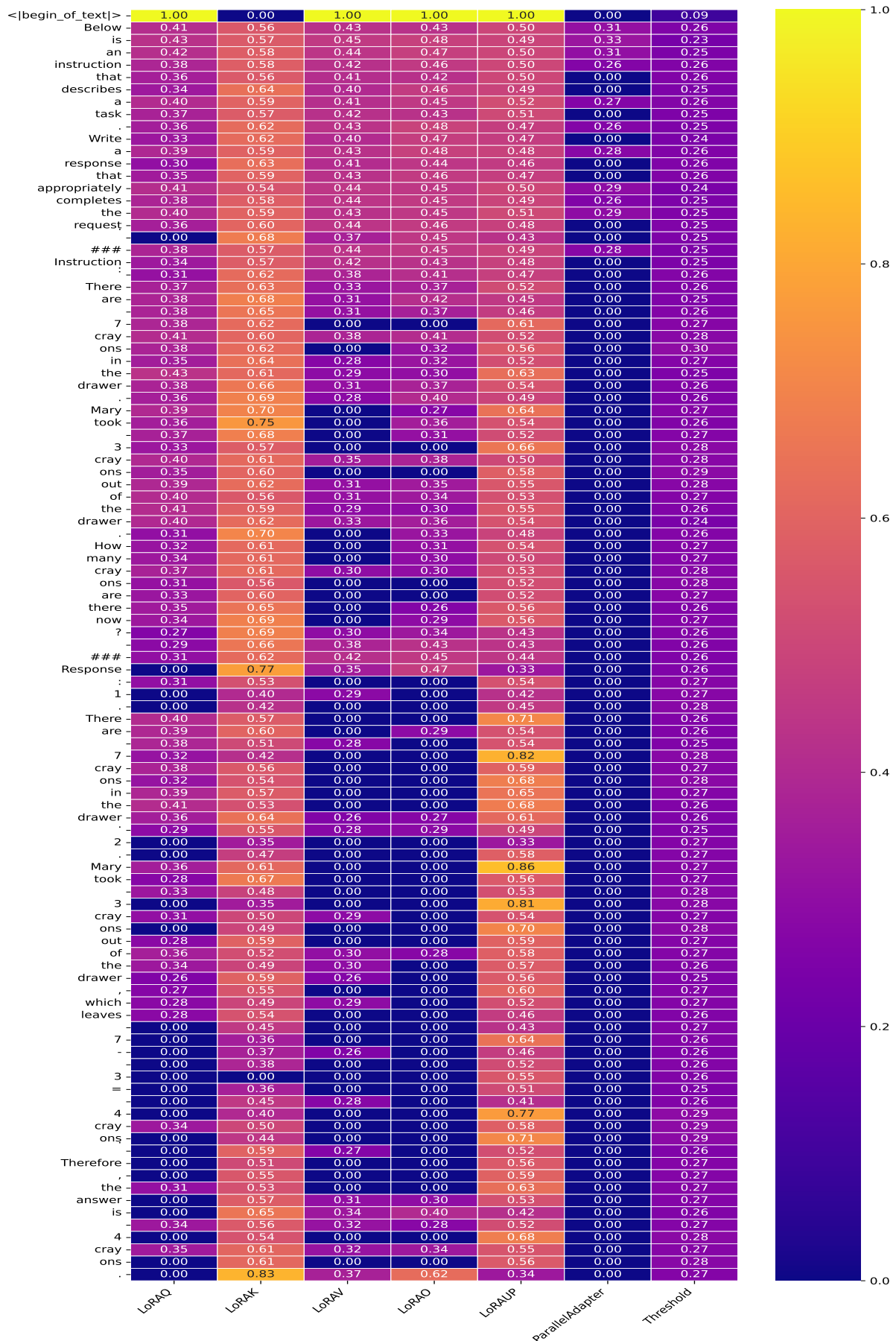


Figure 7: Router Weights of an example at Layer 14 in Sparse MoA. "Threshold" denotes the dynamic value of the threshold function for Sparse MoA. A router weight of "0.00" for an expert indicates that the expert's computation is skipped. Sparse MoA avoids unnecessary computations from the low-contribution expert for each token.

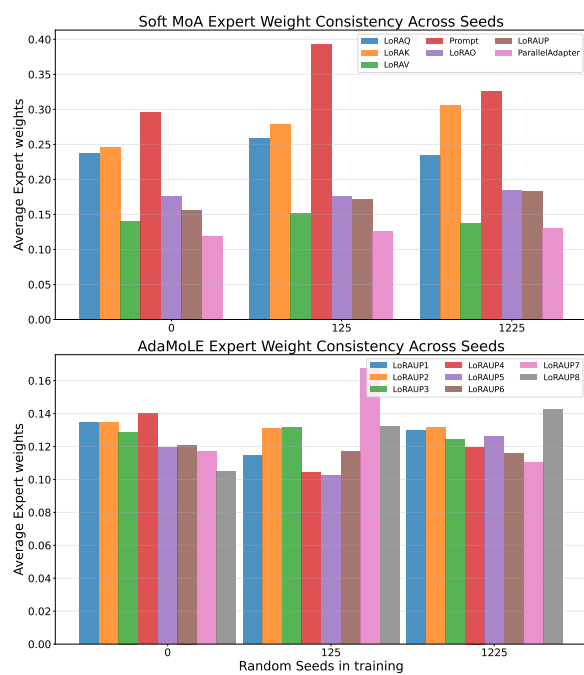


Figure 8: Average router weights across layers. Soft MoA exhibits strong consistency under different random seeds in training, whereas AdaMoLE does not.