

GRAPHMOE: Amplifying Cognitive Depth of Mixture-of-Experts Network via Introducing Self-Rethinking Mechanism

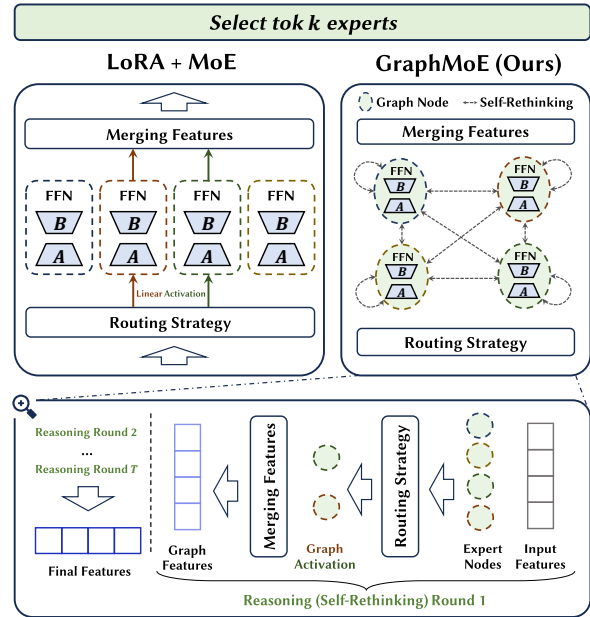
Anonymous ACL submission

Abstract

Traditional Mixture-of-Experts (MoE) networks benefit from utilizing multiple smaller expert models as opposed to a single large network. However, these experts typically operate independently, leaving a question open about whether interconnecting these models could enhance the performance of MoE networks. In response, we introduce GRAPHMOE, aiming at augmenting the cognitive depth of language models via a self-rethinking mechanism constructed on Pseudo Graph MoE networks. GRAPHMOE employs a recurrent routing strategy, allowing the model to iteratively revisit and refine intermediate reasoning states. Simultaneously, by treating experts as interconnected nodes in a pseudo-graph, it facilitates iterative information exchange among experts rather than treating them as isolated modules. We implement the GRAPHMOE architecture using Low-Rank Adaptation techniques (LoRA) and conduct extensive experiments on various benchmark datasets. The experimental results reveal that GRAPHMOE outperforms other LoRA based models, achieving state-of-the-art performance.

1 Introduction

Large Language Models (LLMs), exemplified by the GPT series, have dramatically advanced the field of Natural Language Processing (NLP), demonstrating exceptional performance across a variety of tasks including commonsense reasoning (Yang et al., 2024), creative generation (Loakman et al., 2023), dialogue generation (Lv et al., 2024; Zhao et al., 2023, 2024b), and summarization (Goldsack et al., 2023; Zhao et al., 2024a). In recent years, there has been a growing interest in refining LLM architectures to deliver superior performance with reduced parameters and lower GPU memory consumption (Tang et al., 2023). A prominent strategy in this endeavor is the Mixture-of-Experts (MoE) architecture (Cai et al., 2024), which pre-trains multiple Feed-Forward Networks (FFNs) as specialized experts, activating a select few during inference. This strategy allows experts to



Selected expert nodes transmit messages to the nodes selected for next reasoning round

Figure 1: Comparison between LoRA+MoE and GRAPHMOE architectures.

operate optimally within their domains of expertise, often providing enhanced performance compared to LLMs with a similar number of activated parameters.

Despite their effectiveness, existing MoE architectures typically employ experts independently under a linear routing strategy. Since individual experts are optimized for distinct input distributions, such isolated utilization may limit their collective problem-solving potential. We argue that fostering collaboration among experts, interconnected nodes in a graph, could further enhance MoE reasoning capacity. Inspired by knowledge aggregation mechanisms in pseudo-graphs (Tang et al., 2023; Zequn Sun, 2020), we consider a routing paradigm in which expert models are treated as graph nodes that interact cyclically. We hypothesize that, as graphical message passing among expert nodes increases, MoE outputs can progressively narrow the semantic gap with respect to human references.

061 Based on this motivation, we introduce GRAPH-
062 MOE (Graph Mixture of Experts), an MoE framework
063 that enhances the cognitive capabilities of language
064 models by integrating a self-rethinking mechanism
065 into a pseudo-graph-structured expert network.
066 GRAPHMOE emulates human-like iterative interac-
067 tion by allowing the model to continuously revisit
068 and refine its comprehension of input data through
069 multiple reasoning cycles (Ding et al., 2019; Yan et al.,
070 2023). This is accomplished by implementing a re-
071 current routing strategy on a pseudo-graph formed by
072 expert models. These models transmit their outputs as
073 signals for selecting subsequent expert batches and as
074 inputs encapsulating aggregated graph features from
075 prior reasoning rounds—a mechanism we term the
076 “self-rethinking mechanism”. Given the challenges
077 inherent in pre-training a MoE LLM from scratch,
078 we chose to implement the GRAPHMOE architecture
079 utilizing Low-Rank Adaptation (LoRA) techniques.

080 In Figure 1, we provide an illustration of
081 GRAPHMOE(LoRA) and compare it with the
082 conventional LoRA+MoE framework. In standard
083 LoRA+MoE approaches, the features of activated
084 expert models are integrated directly. This process
085 renders other inactive expert models non-contributory
086 to problem-solving, particularly when the model
087 tends to consistently select a specific set of expert
088 models, as evidenced by low workload balancing. In
089 contrast, GRAPHMOE(LoRA) conceptualizes expert
090 nodes as graph nodes, enabling message transmission
091 across these interconnected nodes through graph
092 activation over multiple reasoning rounds. This
093 approach maintains the reduced memory usage
094 characteristic of conventional LoRA+MoE models
095 (with same activation number of expert models in
096 each cycle) while increasing computational iterations.
097 Experimental results demonstrate the effectiveness of
098 GRAPHMOE, showing consistent improvements over
099 LoRA-based baselines and achieving state-of-the-art
100 performance across evaluated benchmarks.

101 The main contributions of this work are as follows:

- 102 • We introduce a self-rethinking mechanism that
103 enables MoE models to iteratively revisit and
104 refine intermediate representations, facilitating
105 human-like iterative reasoning.
- 106 • We present a pseudo-graph-based MoE architec-
107 ture that models experts as interconnected nodes,
108 enabling iterative information exchange across
109 experts rather than treating them as independent
110 modules.

- Extensive experiments on a diverse range of
benchmarks have been conducted to validate the
effectiveness of GRAPHMOE, providing empir-
ical insights into expert collaboration in LLMs.

2 Related Work

The concept of Mixture of Experts (MoE) was first
introduced by Jacobs et al. (1991), who proposed
training multiple networks (experts) on different sub-
sets of data and aggregating their outputs. Recently,
as LLMs have become a focal point of research, MoE
layers have been integrated into Transformer-based
architectures. Specifically, researchers have replaced
standard Feed-Forward Networks (FFNs) with sparse
MoE layers, employing novel routing strategies (Zuo
et al., 2021; Zhong et al., 2024; Wu et al., 2024a;
Muqeeth et al., 2023; Fu et al., 2024) and advanced
expert segmentation techniques (Dai et al., 2024; He,
2024; Jiang et al., 2024; Xiao et al., 2024).

Numerous studies have investigated the application
of Parameter-Efficient Fine-Tuning (PEFT) methods
to introduce additional trainable parameters for
implementing pseudo MoE structures within LLMs
(Dou et al., 2024; Luo et al., 2024; Gao et al., 2024; Li
et al., 2024; Gou et al., 2023). These models expand
the conventional single feed-forward network (FFN)
architecture and its corresponding representation
space into multiple subspaces, thereby effectively
emulating the behavior of MoE architectures.

Prominent examples include MoLA (Gao et al.,
2024), LoRAMoE (Dou et al., 2024), and MixLoRA
(Li et al., 2024), all of which have demonstrated
state-of-the-art performance compared to other
PEFT-based methods across various benchmarks.
Consequently, we integrate these three distinct LoRA
MoE methodologies into our GRAPHMOE frame-
work. The primary differences between these models
are as follows: LoRAMoE incorporates vanilla atten-
tion LoRA layers and MLP plugged with LoRA-MoE
layers; MoLA integrates plugged LoRA-MoE layers
in both attention and MLP layers; and MixLoRA
combines fused LoRA-MoE modules in attention and
MLP layers. Further analysis of these LoRA-MoE
implementations can be found in Li et al. (2024).

In this study, we augment MoE architectures by
enhancing their reasoning depth and demonstrate the
efficacy of our novel architecture when integrated
with PEFT-based LLM baselines. For a detailed
review of related works, including general PEFT,
Transformer architectures, and standard MoE, please
refer to the Appendix A.

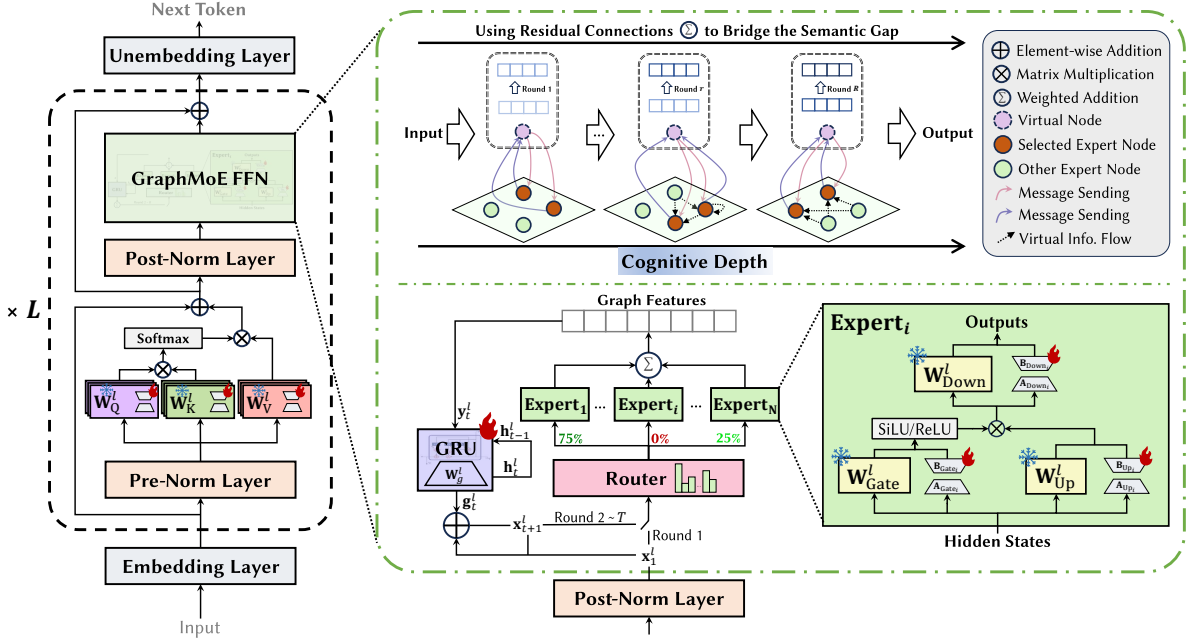


Figure 2: Overview of GRAPHMOE architecture. In this figure, the original Feed-Forward Network (FFN) layer in each transformer block is modified.

3 Method

We propose GRAPHMOE, a pseudo-graph-based MoE framework that connects experts as graph nodes to enable iterative information exchange, inspired by the modular reasoning patterns of the human brain. The overall approach includes three components: transforming existing LoRA-based MoE architectures, constructing a pseudo reasoning graph among experts, and training the model with a recurrent routing mechanism. As shown in Figure 2, the self-rethinking mechanism operates within the GraphMoE Feed-Forward Network (FFN) layer, which replaces the original FFN layer found in standard transformer-based LLMs. While MixLoRA is used as the base example, the framework is compatible with MoLA and LoRAMoE. Across all variants, a unified self-rethinking module based on GRU enables reasoning across expert layers.

3.1 MoE Transformation

Incorporating multiple expert models into LoRA layers, inspired by the work of Li et al. (2024) and Dou et al. (2024), the standard Feed-Forward Networks (FFNs) can be converted into MoE structures comprising n experts, denoted as $\{E_i\}_{i=1}^n$. As illustrated in the bottom-right corner of Figure 2, for the ℓ -th layer of the LLM ($1 \leq \ell \leq L$), we freeze the original FFN parameters and train three pairs of low-rank matrices for each expert: $\mathbf{A}_{\text{Down}_i}^\ell \in \mathbb{R}^{r \times d_1}$, $\mathbf{B}_{\text{Down}_i}^\ell \in \mathbb{R}^{d_2 \times r}$, $\mathbf{A}_{\text{Gate}_i}^\ell$, $\mathbf{B}_{\text{Gate}_i}^\ell$, $\mathbf{A}_{\text{Up}_i}^\ell$, and $\mathbf{B}_{\text{Up}_i}^\ell$. Here, i represents the i -th expert, and d_1 , d_2 are the dimensions of the cor-

responding weight matrices $\mathbf{W}_{\text{Down}_i}^\ell \in \mathbb{R}^{d_2 \times d_1}$, with other dimensions defined similarly. Each expert $E_i(\cdot)$ is constructed as a fusion of the original FFN and trainable low-rank matrices, as shown in Equations 1 and 2, where $\sigma(\cdot)$ denotes activation functions (e.g., SiLU, ReLU), and \odot denotes element-wise multiplication.

$$E_i^\ell(\mathbf{x}) = \tilde{\mathbf{W}}_{\text{Down}_i}^\ell \left(\sigma(\tilde{\mathbf{W}}_{\text{Gate}_i}^\ell \mathbf{x}) \odot (\tilde{\mathbf{W}}_{\text{Up}_i}^\ell \mathbf{x}) \right) \quad (1)$$

$$\begin{aligned} \tilde{\mathbf{W}}_{\text{Gate}_i}^\ell &= \mathbf{W}_{\text{Gate}_i}^\ell + \mathbf{B}_{\text{Gate}_i}^\ell \cdot \mathbf{A}_{\text{Gate}_i}^\ell \\ \tilde{\mathbf{W}}_{\text{Down}_i}^\ell &= \mathbf{W}_{\text{Down}_i}^\ell + \mathbf{B}_{\text{Down}_i}^\ell \cdot \mathbf{A}_{\text{Down}_i}^\ell \\ \tilde{\mathbf{W}}_{\text{Up}_i}^\ell &= \mathbf{W}_{\text{Up}_i}^\ell + \mathbf{B}_{\text{Up}_i}^\ell \cdot \mathbf{A}_{\text{Up}_i}^\ell \end{aligned} \quad (2)$$

Additionally, we introduce a trainable linear layer $\mathbf{R}^\ell \in \mathbb{R}^{n \times d}$, referred to as the Router, where n denotes the total number of experts, and d represents the dimension of the hidden state. The Router determines the importance of each expert during forward propagation by computing a relevance score for every expert. Based on the Router's output, we select the top k experts and normalize their corresponding weights to ensure effective routing, as described in Equations 3 and 4. In subsequent experiments, unless stated otherwise, our setup activates 2 out of 8 experts, thus setting k to 2.

$$\hat{\mathbf{s}}^\ell = \text{Softmax}(\mathbf{R}^\ell \cdot \mathbf{x}) \quad (3)$$

$$\begin{aligned} \mathbf{s}^\ell &= \text{Top-}k(\hat{\mathbf{s}}^\ell) \\ &= \begin{cases} \hat{\mathbf{s}}^\ell[i], & \text{if } \hat{\mathbf{s}}^\ell[i] \text{ is in the top } k, \\ 0, & \text{otherwise.} \end{cases} \end{aligned} \quad (4)$$

Finally, as shown in Equation 5, the output of the MoE module is computed as a weighted sum of the selected experts. This result replaces the original output of the vanilla FFN module.

$$\mathbf{y}^\ell = \text{MoE}(\mathbf{x}) = \sum_{i=1}^n \mathbf{s}^\ell[i] \cdot E_i^\ell(\mathbf{x}) \quad (5)$$

3.2 Pseudo Reasoning Graph Construction

The development of the pseudo reasoning graph represents a pivotal step in implementing our proposed self-rethinking mechanism, which is designed to enhance the internal consistency of LLMs (Liang et al., 2024). The pseudocode outlining the inference process of GRAPHMOE is detailed in Appendix B. For clarity and conciseness, the terms Gate, Up, and Down in the subscripts are abbreviated as G , U , and D , respectively.

Initially, the pseudo graph is constructed as illustrated in the top right corner of Figure 2. Within this framework, expert models are represented as graph nodes, each functioning according to the procedure described in §3.1. A distinctive virtual node is simultaneously established alongside these expert model graph nodes. Unlike other nodes, the virtual node is characterized by a graph feature vector and a GRU module. This node serves as a recurrent router, acting as a conduit for expert node connectivity and a layer for forward feature aggregation.

In contrast to conventional knowledge graphs, there is no prior knowledge that defines the edges between nodes. Instead, the edges connecting graph nodes are dynamically generated at each reasoning round, denoted as t , with the virtual node serving as the intermediary, as indicated by the black dashed arrows in Figure 2. This process facilitates the dynamical activation of temporal edges between nodes, enabling feature transformation across nodes and the newly created edges. Consequently, expert nodes at each reasoning round can communicate and collaborate to address the given task effectively.

Therefore, the core of this method lies in implementing the functionality of the virtual node. In this work, we employ a Low-Rank GRU to serve as the virtual node. Specifically, we reduce the GRU’s hidden size to $\bar{d} \ll d$ and incorporate a low-rank

linear layer to compute its output \mathbf{g}_t^ℓ . The virtual node integrates features from previous reasoning round and processes the expert results \mathbf{y}^ℓ , subsequently selecting k experts for the next round.

The reasoning (self-thinking) process involves a total of T rounds of MoE computations. The first round is processed as described in Equations 1-5. After obtaining the representation \mathbf{y}_t^ℓ from the MoE, it is fed into the Low-Rank GRU along with the hidden state \mathbf{h}_{t-1}^ℓ from the previous timestep (here, “round” and “timestep” are equivalent since there is only one step per round). A residual connection is applied by adding the result to the input of the previous round, \mathbf{x}_t^ℓ , to produce the input for the next round, \mathbf{x}_{t+1}^ℓ . This iterative process extracts meta-information for subsequent reasoning rounds, as shown in Equations 6-7.

$$\begin{aligned} \mathbf{z}_t^\ell &= \sigma(\mathbf{W}_z^\ell \cdot [\mathbf{h}_{t-1}^\ell, \mathbf{y}_t^\ell]), \mathbf{W}_z^\ell \in \mathbb{R}^{\bar{d} \times (\bar{d}+d)} \\ \mathbf{r}_t^\ell &= \sigma(\mathbf{W}_r^\ell \cdot [\mathbf{h}_{t-1}^\ell, \mathbf{y}_t^\ell]), \mathbf{W}_r^\ell \in \mathbb{R}^{\bar{d} \times (\bar{d}+d)} \\ \hat{\mathbf{h}}_t^\ell &= \sigma(\mathbf{W}_o^\ell \cdot [\mathbf{r}_t^\ell \odot \mathbf{h}_{t-1}^\ell, \mathbf{y}_t^\ell] + \mathbf{b}_o^\ell) \\ \mathbf{h}_t^\ell &= (\mathbf{1} - \mathbf{z}_t^\ell) \odot \mathbf{h}_{t-1}^\ell + \mathbf{z}_t^\ell \odot \hat{\mathbf{h}}_t^\ell \end{aligned} \quad (6)$$

$$\begin{aligned} \mathbf{g}_t^\ell &= \mathbf{W}_g^\ell \cdot \mathbf{h}_t^\ell, \mathbf{W}_g^\ell \in \mathbb{R}^{d \times \bar{d}} \\ \mathbf{x}_{t+1}^\ell &= \mathbf{x}_t^\ell + \mathbf{g}_t^\ell \end{aligned} \quad (7)$$

In a nutshell, from a model architecture standpoint, the GRU and the Router together form the aforementioned “virtual node”. Specifically, the GRU functions to aggregate information from multiple experts E_i , whereas the Router is responsible for determining which experts should receive the integrated information for further inference, as depicted in Figure 1 in Appendix A.

3.3 GRAPHMOE Training

Trainable Components. As shown in Figure 2, the modules marked with a flame icon represent the trainable components in GRAPHMOE. Specifically, we need to train $3 \cdot L \cdot n$ pairs of low-rank adapters applied to the MoE modules, as well as a single linear layer serving as the Router. Additionally, in each decoder layer, the Low-Rank GRU module requires training 4 linear layers: \mathbf{W}_z^ℓ , \mathbf{W}_r^ℓ , \mathbf{W}_o^ℓ , and \mathbf{W}_g^ℓ .

Inspired by Li et al. (2024), we also enhance the attention module by adding 4 pairs of trainable low-rank matrices to each decoder layer. These matrices are applied to the attention components \mathbf{Q}_h^ℓ , \mathbf{K}_h^ℓ , $\mathbf{V}_h^\ell \in \mathbb{R}^{d \times \frac{d}{H}}$, and $\mathbf{O}_h^\ell \in \mathbb{R}^{\frac{d}{H} \times d}$, where h and H represent the h -th attention head and the total number of heads,

respectively. This enhances the adaptability of certain attention heads to specific tasks (Zheng et al., 2024).

In summary, as shown in Table 1, GRAPHMOE training involves a total of $(3n + 4) \cdot L$ pairs of low-rank adapters and $4 \cdot L$ linear layers.

Component	Attention Head	MoE	GRU	Total
Low-rank adapters	$4 \cdot L$	$3 \cdot L \cdot n$	0	$(3n+4) \cdot L$
Linear layers	0	0	$4 \cdot L$	$4 \cdot L$

Table 1: Summary of components to be trained.

Loss Design. Similar to standard continual learning (Wu et al., 2024c) and supervised fine-tuning (Patil and Gudivada, 2024), the primary objective during training is to minimize the cross-entropy loss \mathcal{L}_{CE} between the predicted token logits and the target tokens.

However, the Router is a unique and critical component of the MoE architecture, as it determines which experts are activated during inference. Prior research has shown that unconstrained Router training can lead to over-reliance on a few specific experts (Zoph et al., 2022; Fedus et al., 2022), resulting in poor load balancing. Inspired by the work of Shen et al. (2024) and Muennighoff et al. (2024), we introduce a **Load Balancing Loss** \mathcal{L}_{LB} (as defined in Equation 8) as an auxiliary objective. This loss penalizes the Router for disproportionately selecting the same experts across the sequence set \mathcal{M} , which consists of several sequences and the corresponding next tokens to be predicted.

$$\mathcal{L}_{LB} = n \cdot \sum_{i=1}^n f_i \cdot p_i,$$

$$\text{where } f_i = \frac{1}{|\mathcal{M}|} \sum_{(\mathbf{x}, \mathbf{y}) \in \mathcal{M}} \mathbb{1}\{\mathbf{s}_{\mathbf{x}}[i] > 0\}, \quad (8)$$

$$p_i = \frac{1}{|\mathcal{M}|} \sum_{(\mathbf{x}, \mathbf{y}) \in \mathcal{M}} \mathbf{s}_{\mathbf{x}}[i]$$

Setting a hyperparameter λ , the total loss can be expressed as follows:

$$\mathcal{L}_{\text{total}} = \mathcal{L}_{CE} + \lambda \cdot \mathcal{L}_{LB} \quad (9)$$

4 Experiment

4.1 Experimental Setup

Datasets. We evaluate our method on a diverse range of benchmarks: ARC (Clark et al., 2018), OpenBookQA (Mihaylov et al., 2018), PIQA (Bisk et al., 2020), and SocialIQA (Sap et al., 2019) for QA tasks; BoolQ (Clark et al., 2019) for classification tasks; Hellaswag (Zellers et al., 2019) for science

completion; and Winogrande (Sakaguchi et al., 2021) for fill-in-the-blank tasks.

Baselines. We adopt LLaMA-3-8B as the backbone for all models. To benchmark GRAPHMOE, we compare it with standard LoRA (Hu et al., 2021a) and three recent state-of-the-art LoRA+MoE methods: MoLA (Gao et al., 2024), LoRAMoE (Dou et al., 2024), and MixLoRA (Li et al., 2024). As Weight-Decomposed LoRA (DoRA) (Liu et al., 2024) has shown to improve LoRA’s performance, we also implement DoRA-enhanced variants of each baseline.

Evaluation Metrics. We evaluate all methods using accuracy across all datasets. For Multiple-Choice QA tasks, we adopt the first-token evaluation strategy, which measures the probability of each option letter (Yu et al., 2024). To analyze expert utilization, we track the selection frequency of each expert during inference and compute their selection ratios. Additionally, we report several indicators: GRU hidden size scale (\bar{d}/d), parameter ratio (Param), and normalized inference time (Infer Time). GRU hidden size scale is defined as the ratio of GRU hidden size to model dimension; Param denotes the proportion of trainable parameters; Infer Time is normalized by dividing each model’s inference time by the minimum across all models.

Implementation Details. All methods use LLaMA-3-8B¹ as the backbone. Training settings follow Li et al. (2024). The key hyperparameters for all methods are as follows: the cutoff length is set to 512, the learning rate is $2e-4$, and the optimizer used is AdamW. We use a batch size of 16 with 8 accumulation steps, and the dropout rate is 0.05. Training is conducted for 2 epochs. For LoRA and DoRA methods, the rank r is set to 80, and the LoRA alpha α is set to 160. For MixLoRA, MixDoRA and GRAPHMOE, the LoRA rank and α is set to 16 and 32, respectively. Each MoE layer consists of 8 experts, with top-2 experts activated during inference. For GRAPHMOE, (1) reasoning round $T=3$ (the first round disables GRU-based feature transmission); (2) a 1-layer GRU with a hidden size scaled to 0.1×4096 ; (3) a load-balance loss weight $\lambda=0.01$. All experiments are conducted on a single A800 GPU using BF16 precision².

4.2 Experimental Result

Effect of GRAPHMOE. As shown in Table 2, integrating GRAPHMOE, a pseudo graph-based

¹<https://huggingface.co/meta-llama/Meta-Llama-3-8B-Instruct>

²Additional implementation details are provided in Appendix C

Method	ARC-E	ARC-C	BoolQ	OBQA	PIQA	SIQA	HellaS	WinoG	AVG.
MoLA	86.4	77.9	74.0	84.4	86.7	76.4	93.9	83.3	82.9
GRAPHMOE (MoLA)	89.7	81.0	75.8	87.6	87.8	79.7	95.5	83.2	85.0
Δ	3.3	3.1	1.8	3.2	1.1	3.3	1.6	-0.1	2.1
LoRAMoE	87.8	79.5	72.4	85.0	87.1	74.8	94.8	83.4	83.1
GRAPHMOE (LoRAMoE)	87.7	79.9	75.9	88.8	87.8	79.4	94.8	83.4	84.7
Δ	-0.1	0.4	3.5	3.8	0.7	4.6	0.0	0.0	1.6
MixLoRA	86.9	77.0	74.0	84.4	86.0	75.5	93.7	83.3	82.6
GRAPHMOE (MixLoRA)	89.2	80.9	75.8	89.6	87.8	77.0	95.4	83.2	84.9
Δ	2.3	3.9	1.8	5.2	1.8	1.5	1.7	-0.1	2.3

Table 2: Comparison of the performance between baseline LoRA+MoE models and GRAPHMOE (LoRA+MoE) models. GRAPHMOE (·) indicates that the conventional MoE module in the baseline has been substituted with the self-rethinking mechanism proposed in § 3.2. Scores highlighted in **bold** denote superior performance for each respective method. **Green** and **red** indicate increase and decrease in performance following the application of GRAPHMOE respectively.

Method	Param	ARC-E	ARC-C	BoolQ	OBQA	PIQA	SIQA	HellaS	WinoG	AVG.
Base (ICL-3)	0%	90.6	76.7	45.9	73.8	61.6	71.5	64.1	59.2	65.6
Base (ICL-3, Self-cons-5)	0%	91.1	79.4	50.2	73.7	57.0	70.3	66.9	62.1	66.3
Base (ICL-3, Self-cons-10)	0%	91.6	78.1	42.8	75.3	59.8	71.3	69.9	61.4	66.2
Base (ICL-3, Self-cons-15)	0%	91.9	78.0	45.6	72.0	67.9	71.7	63.3	63.2	66.4
LoRA*	2.6%	89.0	75.7	67.2	85.0	80.7	78.3	74.2	75.3	78.2
DoRA*	2.6%	88.1	76.4	61.7	80.6	82.3	76.2	78.8	83.7	78.5
MixLoRA*	3.0%	86.5	79.9	75.0	84.8	87.6	78.8	93.3	82.1	83.5
MixDoRA*	3.0%	87.7	78.9	76.8	86.9	83.4	80.1	94.6	84.2	84.1
LoRA	2.6%	86.2	76.7	70.2	78.2	82.9	74.2	80.1	50.7	74.9
DoRA	2.6%	87.6	76.8	69.3	80.8	84.2	76.8	85.3	50.4	76.4
MoLA	2.7%	86.4	77.9	74.0	84.4	86.7	76.4	93.9	83.3	82.9
MoDA	2.7%	86.4	78.0	74.0	84.8	87.4	76.4	95.5	84.1	83.3
LoRAMoE	3.2%	87.8	79.5	72.4	85.0	87.1	74.8	94.8	83.4	83.1
DoRAMoE	3.2%	87.9	80.1	74.6	85.8	86.3	74.5	94.0	83.1	83.3
MixLoRA	3.0%	86.9	77.0	74.0	84.4	86.0	75.5	93.7	83.3	82.6
MixDoRA	3.0%	86.7	76.7	75.5	84.8	85.2	76.9	93.0	83.1	82.7
GRAPHMOE (MixLoRA)	5.9%	89.2	80.9	75.8	89.6	87.8	77.0	95.4	83.2	84.9
GRAPHMOE (MixDoRA)	5.9%	90.3	80.6	75.9	88.2	88.8	79.4	95.3	83.7	85.3

Table 3: Performance of LoRA-based SFT methods across benchmarks. ICL- n ” indicates n in-context examples; Self-cons-5” denotes 5 samples for self-consistency. Param” is the number of trainable parameters. Results marked with * are from Li et al. (2024) and use full precision (FP32); all others use BF16. DoRA-based variants are denoted by replacing L” with “D” in model names.

MoE with recurrent routing, consistently improves performance over all SOTA LoRA+MoE baselines. This demonstrates GRAPHMOE’s ability to enhance expert collaboration by simply replacing the original MoE module, enabling more effective feature exchange among experts. Performance gains vary across base models due to differences in how LoRA and MoE are combined. As noted by Li et al. (2024), MoLA and LoRAMoE use a “plugged-in” design, injecting LoRA modules into attention or MLP outputs, while MixLoRA adopts a “fused” strategy that embeds LoRA experts directly into the network layers—yielding stronger integration and performance boost when enhanced by GRAPHMOE. Across all tasks, GRAPHMOE consistently improves or at least maintains accuracy. Notably, on SocialQA, applying

GRAPHMOE to LoRAMoE improves accuracy from 74.8 to 79.4 (+4.6, over 6%). These observations underscores the effectiveness of the GRAPHMOE approach within this architectural framework.

Main Comparison Results. The experimental results are presented in Table 3. We compare the performance of four base models (LLaMA-3-8B) incorporating in-context learning (ICL- n) (Dong et al., 2024) and self-consistency (Self-cons- n) (Wang et al., 2023). While increasing the number of samples for self-consistency to 5, 10, and 15 shows minimal variation in performance (scores of 66.3, 66.2, and 66.4, respectively), these models still fall short compared to those using LoRA-based SFT methods. Due to computational constraints, we use BF16 precision for training and inference. However, we include full

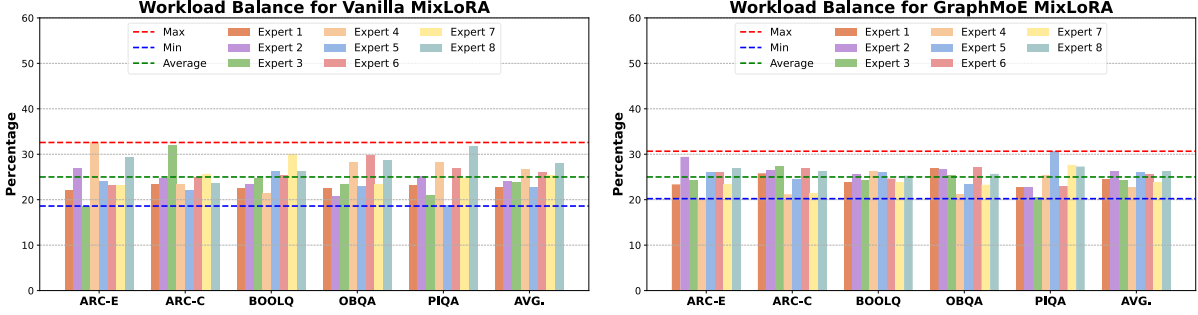


Figure 3: Normalized expert workloads across routing steps. Red and blue dashed lines mark the max and min workloads; green indicates the average. The workload std is 0.0313 for MixLoRA and 0.0215 for GRAPHMOE.

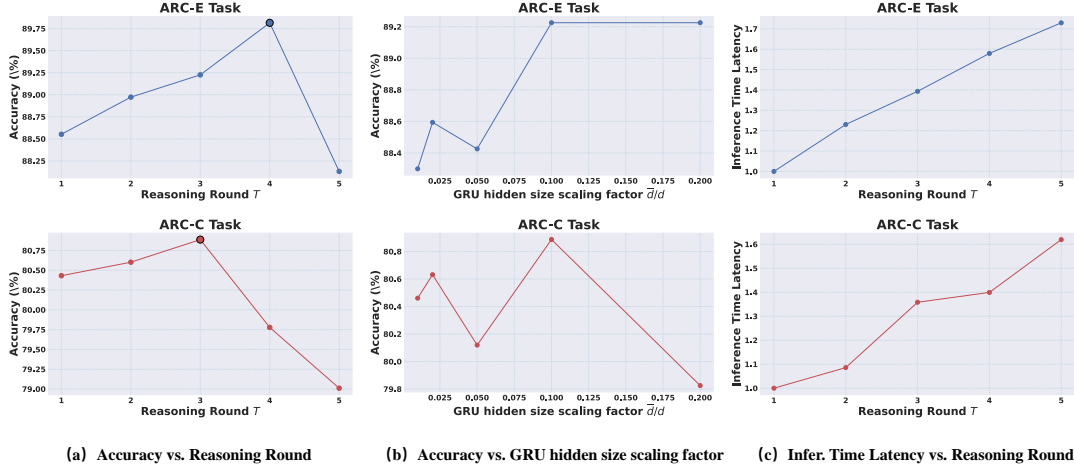


Figure 4: Sensitivity analysis of GRAPHMOE hyperparameters. Subfigure (c) shows how inference time scales with reasoning round (T), where the first round’s time is set as the unit. “Infer. Time” denotes the relative time per round.

412 precision MixLoRA results for comparison. Despite
 413 a significant performance drop in all models when
 414 using reduced precision, our GRAPHMOE models
 415 with MixLoRA and MixDoRA still outperform full
 416 precision models, demonstrating the efficacy of our
 417 architecture. Moreover, MixLoRA does not represent
 418 the optimal base model with the LoRA+MoE frame-
 419 work, as MoLA and LoRAMoE have better accuracy
 420 scores. This indicates that the performance gains
 421 achieved by GRAPHMOE surpass those of existing
 422 methodologies applied within the LoRA+MoE
 423 framework. The self-rethinking mechanism enhances
 424 the MoE approach, evident in both LoRA and
 425 DoRA settings. This suggests that the improvements
 426 carried by the PRG are not merely a result of simple
 427 adjustments on additional trainable parameters, but
 428 rather arise from the effective collaboration among
 429 diverse expert models facilitated by our approach.

430 **Analysis of Workload Balance.** The main difference
 431 between MixLoRA and GRAPHMOE (MixLoRA)
 432 is the self-rethinking mechanism enabled by the re-
 433 current router. We analyze the impact of expert model

434 selection at each routing stage. Figure 3 compares
 435 the workload distribution across expert models in
 436 MixLoRA and GRAPHMOE (MixLoRA) for various
 437 tasks. Results show that GRAPHMOE achieves more
 438 balanced expert selection, with standard deviations of
 439 0.0313 for MixLoRA and 0.0249 for GRAPHMOE
 440 (MixLoRA), a 20% reduction. This indicates that the
 441 reasoning process in GRAPHMOE activates a broader
 442 range of experts, enhancing MoE’s representation
 443 learning. In each reasoning step, a distinct set of
 444 expert models is more likely to be activated to incre-
 445 mentally construct the representations within the MoE.
 446 This approach allows more expert models to operate
 447 in their optimal conditions, sequentially addressing
 448 complex problems. It is as if a chain-of-thoughts
 449 (CoT) is implicitly formed, using different expert
 450 models to fit features at different stages.

451 **Sensitivity Analysis and Overhead.** We further
 452 conduct sensitivity analysis on ARC-E and ARC-C to
 453 provide insights into the efficacy of our self-rethinking
 454 mechanism. In subfigure (a) of Figure 4, increasing
 455 the Reasoning Round T improves accuracy for both

Method	Param	ARC-E	ARC-C	BoolQ	OBQA	PIQA	SIQA	HellaS	WinoG	AVG
LoRA	1.9%	87.5	78.9	69.8	84.4	84.4	77.1	88.4	49.1	77.5
LoRA	2.6%	86.2	76.7	70.2	78.2	82.9	74.2	80.1	50.7	74.9
LoRA	5.9%	84.7	78.6	63.6	82.2	82.1	32.9	84.1	50.4	69.8
LoRA	7.9%	83.9	73.1	62.2	84.4	82.5	74.6	92.7	49.5	75.4
LoRA	9.9%	85.9	73.3	70.0	82.0	85.7	32.1	89.1	49.6	71.0
MixLoRA	3.0%	86.9	77.0	74.0	84.4	86.0	75.5	93.7	83.3	82.6
MixLoRA	6.0%	87.3	79.0	74.9	88.4	88.1	76.7	93.6	82.5	83.8
MixLoRA	9.0%	86.7	77.0	71.7	86.2	80.5	74.0	89.5	80.4	80.8
GRAPHMOE	5.9%	89.2	80.9	75.8	89.6	87.8	77.0	95.4	83.2	84.9

Table 4: Performance Comparison of Parameter Variants. In the table, GRAPHMOE refers to the GRAPHMOE (MixLoRA) model.

Method	Param	ARC-E	ARC-C	BoolQ	OBQA	PIQA	SIQA	HellaS	WinoG	AVG
MixLoRA†	6.0%	87.4	79.2	73.9	85.2	85.3	72.9	95.0	80.4	82.4
MixLoRA	6.0%	87.3	79.0	74.9	88.4	88.1	76.7	93.6	82.5	83.8
Δ	-	-0.1	-0.2	1	3.2	2.8	3.8	-1.4	2.1	1.4
GRAPHMOE†	5.9%	88.1	79.0	74.1	87.8	85.6	74.3	93.2	85.1	83.4
GRAPHMOE	5.9%	89.2	80.9	75.8	89.6	87.8	77.0	95.4	83.2	84.9
Δ	-	1.1	1.9	1.7	1.8	2.2	2.7	2.2	-1.9	1.5

Table 5: Performance Comparison of GRAPHMOE (MixLoRA) with different activated experts. GRAPHMOE† indicates the GRAPHMOE (MixLoRA) variant with 4 out of 8 experts activated, whereas the other LoRA MoE models in above experiments have 2 out of 8 experts activated.

tasks, demonstrating the self-rethinking mechanism’s effectiveness in enhancing MoE. However, accuracy drops significantly beyond a certain threshold ($T=4$ for ARC-E and $T=3$ for ARC-C), likely due to overfitting from "overthinking" the problem. Subfigure (b) of Figure 4 further illustrates that increasing the GRU hidden size does not always improve performance, highlighting the importance of the self-rethinking mechanism over the added GRU parameters. Finally, subfigure (c) shows that, compared to conventional LoRA+MoE models (MixLoRA), which employ a single reasoning round, the increase in inference time for GRAPHMOE is only approximately 0.16 times for each additional reasoning round. These findings shows the promising impact GRAPHMOE on enhancing MoE’s cognitive depth while maintaining acceptable computational overhead.

Analysis of PEFT Parameters. We analyze the effect of the parameter size on model performance. From Table 4, GRAPHMOE (MixLoRA) gains the best overall performance, which uses 5.9% of parameters yet achieves the highest average score of 84.9. In addition, it can be observed that simply increasing parameter size does not lead to proportional improvements in SFT performance. For instance, MixLoRA with 3.0% of parameters achieves the best average score of 82.6%, while increasing the parameter size to 9.0% causes fluctuations below that

score. Similarly, in Figure 4(b), performance does not improve as the GRU hidden size increases.

Analysis of Active Expert Number. To further investigate MoE utilization, we tested the effect of activating additional experts in the MoE layers, as shown in Table 5. Surprisingly, increasing the number of active experts led to a performance decline. This highlights that simply adding parameters does not guarantee better reasoning performance. Instead, the effectiveness of the neural architectural, as demonstrated in Table 4, is a more critical factor.

5 Conclusion

In this work, we introduce GRAPHMOE, a novel approach that enhances the cognitive depth of language models through the integration of a self-rethinking mechanism in pseudo-graph Mixture-of-Experts (MoE) networks. Unlike traditional MoE models that operate independently, GRAPHMOE facilitates communication among expert nodes, allowing for iterative refinement and enhanced reasoning capabilities. Implemented using Low-Rank Adaptation techniques (LoRA), GRAPHMOE demonstrates significant performance improvements across various benchmark datasets, outperforming existing LoRA & LoRA+MoE baseline models and achieving state-of-the-art results.

510 Limitations

511 Despite the notable performance improvements
512 achieved by the GRAPHMOE framework, there are
513 several limitations to be acknowledged.

- 514 • While the GRAPHMOE architecture consistently leads to performance gains across
515 diverse tasks and base models, the degree of
516 improvement is variable and can depend on
517 the interaction between the MoE and LoRA
518 components. This variability underscores the
519 need for comprehensive exploration into the
520 integration modalities and parameter space to
521 ensure maximum efficiency across different
522 LoRA+MoE configurations.
523
- 524 • The workload distribution analysis indicates a
525 marked improvement in workload balance with
526 GRAPHMOE. However, the potential impact of
527 workload imbalance on model performance over
528 a broader range of tasks remains underexplored.
529 Further investigation into balancing expert
530 model selection and activation across diverse
531 scenarios could unveil additional pathways for
532 performance optimization.
- 533 • The sensitivity analysis indicated potential
534 overfitting issues when increasing the reasoning
535 rounds beyond a particular threshold, revealing
536 the necessity for careful hyperparameter tuning
537 to mitigate issues of over-complexity and model
538 overthinking.

539 In summary, while the GRAPHMOE framework
540 shows promise in enhancing cognitive depth and
541 performance of MoE architectures, future research
542 should address computational precision limits,
543 explore broader integration strategies, and focus on
544 hyperparameter optimization to further refine and
545 substantiate these initial findings.

546 References

- 547 Yonatan Bisk, Rowan Zellers, Jianfeng Gao, Yejin Choi,
548 and 1 others. 2020. Piqa: Reasoning about physical
549 commonsense in natural language. In *Proceedings
550 of the AAAI conference on artificial intelligence*,
551 volume 34, pages 7432–7439.
- 552 Weilin Cai, Juyong Jiang, Fan Wang, Jing Tang, Sunghun
553 Kim, and Jiayi Huang. 2024. A survey on mixture of
554 experts. *arXiv preprint arXiv:2407.06204*.
- 555 Yupeng Chang, Xu Wang, Jindong Wang, Yuan Wu, Linyi
556 Yang, Kaijie Zhu, Hao Chen, Xiaoyuan Yi, Cunxiang

- Wang, Yidong Wang, and 1 others. 2024. A survey on
557 evaluation of large language models. *ACM Transactions
558 on Intelligent Systems and Technology*, 15(3):1–45.
559
- Junyoung Chung, Caglar Gulcehre, KyungHyun Cho, and
560 Yoshua Bengio. 2014. Empirical evaluation of gated
561 recurrent neural networks on sequence modeling. *arXiv
562 preprint arXiv:1412.3555*.
563
- Christopher Clark, Kenton Lee, Ming-Wei Chang, Tom
564 Kwiatkowski, Michael Collins, and Kristina Toutanova.
565 2019. Boolq: Exploring the surprising difficulty of natu-
566 ral yes/no questions. *arXiv preprint arXiv:1905.10044*.
567
- Peter Clark, Isaac Cowhey, Oren Etzioni, Tushar Khot,
568 Ashish Sabharwal, Carissa Schoenick, and Oyvind
569 Tafjord. 2018. Think you have solved question
570 answering? try arc, the ai2 reasoning challenge. *arXiv
571 preprint arXiv:1803.05457*.
572
- Damai Dai, Chengqi Deng, Chenggang Zhao, RX Xu,
573 Huazuo Gao, Deli Chen, Jiashi Li, Wangding Zeng,
574 Xingkai Yu, Y Wu, and 1 others. 2024. Deepseek-
575 moe: Towards ultimate expert specialization in
576 mixture-of-experts language models. *arXiv preprint
577 arXiv:2401.06066*.
578
- Ming Ding, Chang Zhou, Qibin Chen, Hongxia Yang, and
579 Jie Tang. 2019. Cognitive graph for multi-hop reading
580 comprehension at scale. In *Proceedings of the 57th
581 Annual Meeting of the Association for Computational
582 Linguistics*.
583
- Qingxiu Dong, Lei Li
584 a@miscdong2024surveyincontextlearning, title=A
585 Survey on In-context Learning, author=Qingxiu Dong
586 and Lei Li and Damai Dai and Ce Zheng and Jingyuan
587 Ma and Rui Li and Heming Xia and Jingjing Xu
588 and Zhiyong Wu and Tianyu Liu and Baobao Chang
589 and Xu Sun and Lei Li and Zhifang Sui, year=2024,
590 eprint=2301.00234, archivePrefix=arXiv, primary
591 Class=cs.CL, url=https://arxiv.org/abs/2301.00234, nd
592 Damai Dai, Ce Zheng, Jingyuan Ma, Rui Li, Heming
593 Xia, Jingjing Xu, Zhiyong Wu, Tianyu Liu, Baobao
594 Chang, Xu Sun, Lei Li, and Zhifang Sui. 2024. [A survey
595 on in-context learning](#). *Preprint*, arXiv:2301.00234.
596
- Shihan Dou, Enyu Zhou, Yan Liu, Songyang Gao,
597 Wei Shen, Limao Xiong, Yuhao Zhou, Xiao Wang,
598 Zhiheng Xi, Xiaoran Fan, Shiliang Pu, Jiang Zhu,
599 Rui Zheng, Tao Gui, Qi Zhang, and Xuanjing Huang.
600 2024. [LoRAMoE: Alleviating world knowledge
601 forgetting in large language models via MoE-style
602 plugin](#). In *Proceedings of the 62nd Annual Meeting of
603 the Association for Computational Linguistics (Volume
604 1: Long Papers)*, pages 1932–1945, Bangkok, Thailand.
605 Association for Computational Linguistics.
606
- William Fedus, Barret Zoph, and Noam Shazeer. 2022.
607 Switch transformers: Scaling to trillion parameter
608 models with simple and efficient sparsity. *Journal of
609 Machine Learning Research*, 23(120):1–39.
610
- Tianyu Fu, Haofeng Huang, Xuefei Ning, Genghan Zhang,
611 Boju Chen, Tianqi Wu, Hongyi Wang, Zixiao Huang,
612 Shiyao Li, Shengen Yan, and 1 others. 2024. Moa:
613

614	Mixture of sparse attention for automatic large language model compression. <i>arXiv preprint arXiv:2406.14909</i> .	669
615		670
616	Chongyang Gao, Kezhen Chen, Jinmeng Rao, Baochen Sun, Ruibo Liu, Daiyi Peng, Yawen Zhang, Xiaoyuan Guo, Jie Yang, and VS Subrahmanian. 2024. Higher layers need more lora experts. <i>arXiv preprint arXiv:2402.08562</i> .	671
617		672
618		673
619		674
620		675
621	Tomas Goldsack, Zhihao Zhang, Chen Tang, Carolina Scarton, and Chenghua Lin. 2023. Enhancing biomedical lay summarisation with external knowledge graphs. In <i>Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing</i> , pages 8016–8032.	676
622		677
623		678
624		679
625		680
626	Yunhao Gou, Zhili Liu, Kai Chen, Lanqing Hong, Hang Xu, Aoxue Li, Dit-Yan Yeung, James T Kwok, and Yu Zhang. 2023. Mixture of cluster-conditional lora experts for vision-language instruction tuning. <i>arXiv preprint arXiv:2312.12379</i> .	681
627		682
628		683
629		684
630		685
631	Zishan Guo, Renren Jin, Chuang Liu, Yufei Huang, Dan Shi, Linhao Yu, Yan Liu, Jiaxuan Li, Bojian Xiong, Deyi Xiong, and 1 others. 2023. Evaluating large language models: A comprehensive survey. <i>arXiv preprint arXiv:2310.19736</i> .	686
632		687
633		688
634		689
635		690
636	Xu Owen He. 2024. Mixture of a million experts. <i>arXiv preprint arXiv:2407.04153</i> .	691
637		692
638		693
639	S Hochreiter. 1997. Long short-term memory. <i>Neural Computation MIT-Press</i> .	694
640	Edward J Hu, Yelong Shen, Phillip Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang, and Weizhu Chen. 2021a. Lora: Low-rank adaptation of large language models. <i>arXiv preprint arXiv:2106.09685</i> .	695
641		696
642		697
643		698
644	Edward J Hu, Yelong Shen, Phillip Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang, and Weizhu Chen. 2021b. Lora: Low-rank adaptation of large language models. <i>arXiv preprint arXiv:2106.09685</i> .	699
645		700
646		701
647		702
648	Robert A Jacobs, Michael I Jordan, Steven J Nowlan, and Geoffrey E Hinton. 1991. Adaptive mixtures of local experts. <i>Neural computation</i> , 3(1):79–87.	703
649		704
650		705
651	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. Mixtral of experts. <i>arXiv preprint arXiv:2401.04088</i> .	706
652		707
653		708
654		709
655		710
656	Dawid J Kopiczko, Tijmen Blankevoort, and Yuki M Asano. 2023. Vera: Vector-based random matrix adaptation. <i>arXiv preprint arXiv:2310.11454</i> .	711
657		712
658		713
659	Dengchun Li, Yingzi Ma, Naizheng Wang, Zhiyuan Cheng, Lei Duan, Jie Zuo, Cal Yang, and Mingjie Tang. 2024. Mixlora: Enhancing large language models fine-tuning with lora based mixture of experts. <i>arXiv preprint arXiv:2404.15159</i> .	714
660		715
661		716
662		717
663		718
664	Xun Liang, Shichao Song, Zifan Zheng, Hanyu Wang, Qingchen Yu, Xunkai Li, Rong-Hua Li, Yi Wang, Zhonghao Wang, Feiyu Xiong, and 1 others. 2024. Internal consistency and self-feedback in large language models: A survey. <i>arXiv preprint arXiv:2407.14507</i> .	719
665		720
666		721
667		722
668		723
	Shih-Yang Liu, Chien-Yi Wang, Hongxu Yin, Pavlo Molchanov, Yu-Chiang Frank Wang, Kwang-Ting Cheng, and Min-Hung Chen. 2024. Dora: Weight-decomposed low-rank adaptation. <i>arXiv preprint arXiv:2402.09353</i> .	
	Tyler Loakman, Chen Tang, and Chenghua Lin. 2023. TwistList: Resources and baselines for tongue twister generation. In <i>Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 2: Short Papers)</i> .	
	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. <i>arXiv preprint arXiv:2402.12851</i> .	
	Bo Lv, Chen Tang, Yanan Zhang, Xin Liu, Ping Luo, and Yue Yu. 2024. URG: A unified ranking and generation method for ensembling language models. In <i>Findings of the Association for Computational Linguistics ACL 2024</i> , pages 4421–4434, Bangkok, Thailand and virtual meeting. Association for Computational Linguistics.	
	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. <i>arXiv preprint arXiv:1809.02789</i> .	
	Niklas Muennighoff, Luca Soldaini, Dirk Groeneveld, Kyle Lo, Jacob Morrison, Sewon Min, Weijia Shi, Pete Walsh, Oyvind Tafjord, Nathan Lambert, and 1 others. 2024. Olmoe: Open mixture-of-experts language models. <i>arXiv preprint arXiv:2409.02060</i> .	
	Mohammed Muqeeth, Haokun Liu, and Colin Raffel. 2023. Soft merging of experts with adaptive routing. <i>arXiv preprint arXiv:2306.03745</i> .	
	Rajvardhan Patil and Venkat Gudivada. 2024. A review of current trends, techniques, and challenges in large language models (llms). <i>Applied Sciences</i> , 14(5):2074.	
	Bo Peng, Eric Alcaide, Quentin Anthony, Alon Albalak, Samuel Arcadinho, Stella Biderman, Huanqi Cao, Xin Cheng, Michael Chung, Matteo Grella, and 1 others. 2023. Rwkv: Reinventing rnns for the transformer era. <i>arXiv preprint arXiv:2305.13048</i> .	
	Subhojeet Pramanik, Esraa Elelimy, Marlos C Machado, and Adam White. 2023. Recurrent linear transformers. <i>arXiv preprint arXiv:2310.15719</i> .	
	Keisuke Sakaguchi, Ronan Le Bras, Chandra Bhagavatula, and Yejin Choi. 2021. Winogrande: An adversarial winograd schema challenge at scale. <i>Communications of the ACM</i> , 64(9):99–106.	
	Maarten Sap, Hannah Rashkin, Derek Chen, Ronan Le-Bras, and Yejin Choi. 2019. Socialqa: Commonsense reasoning about social interactions. <i>arXiv preprint arXiv:1904.09728</i> .	
	Yikang Shen, Zhen Guo, Tianle Cai, and Zengyi Qin. 2024. Jetmoe: Reaching llama2 performance with 0.1 m dollars. <i>arXiv preprint arXiv:2404.07413</i> .	

724	Chen Tang. 2024. <i>Knowledge Enhanced Natural Language Generation</i> . Ph.D. thesis, University of Surrey.	777
725		778
726	Chen Tang, Hongbo Zhang, Tyler Loakman, Chenghua Lin, and Frank Guerin. 2023. Enhancing dialogue generation via dynamic graph knowledge aggregation. In <i>Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)</i> .	779
727		780
728		781
729		782
730		783
731	A Vaswani. 2017. Attention is all you need. <i>Advances in Neural Information Processing Systems</i> .	784
732		785
733	Xuezhi Wang, Jason Wei, Dale Schuurmans, Quoc Le, Ed Chi, Sharan Narang, Aakanksha Chowdhery, and Denny Zhou. 2023. Self-consistency improves chain of thought reasoning in language models . <i>Preprint</i> , arXiv:2203.11171.	786
734		787
735		788
736		789
737		790
738	Shaohua Wu, Jiangang Luo, Xi Chen, Lingjun Li, Xudong Zhao, Tong Yu, Chao Wang, Yue Wang, Fei Wang, Weixu Qiao, and 1 others. 2024a. Yuan 2.0-m32: Mixture of experts with attention router. <i>arXiv preprint arXiv:2405.17976</i> .	791
739		792
740		793
741		794
742		795
743	Taiqiang Wu, Jiahao Wang, Zhe Zhao, and Ngai Wong. 2024b. Mixture-of-subspaces in low-rank adaptation. <i>arXiv preprint arXiv:2406.11909</i> .	796
744		797
745		798
746	Tongtong Wu, Linhao Luo, Yuan-Fang Li, Shirui Pan, Thuy-Trang Vu, and Gholamreza Haffari. 2024c. Continual learning for large language models: A survey. <i>arXiv preprint arXiv:2402.01364</i> .	799
747		800
748		801
749		802
750	Da Xiao, Qingye Meng, Shengping Li, and Xingyuan Yuan. 2024. Improving transformers with dynamically composable multi-head attention. <i>arXiv preprint arXiv:2405.08553</i> .	803
751		804
752		805
753		806
754	Junbing Yan, Chengyu Wang, Taolin Zhang, Xiaofeng He, Jun Huang, and Wei Zhang. 2023. From complex to simple: Unraveling the cognitive tree for reasoning with small language models. <i>arXiv preprint arXiv:2311.06754</i> .	807
755		808
756		809
757		810
758		811
759	Bohao Yang, Chen Tang, Kun Zhao, Chenghao Xiao, and Chenghua Lin. 2024. Effective distillation of table-based reasoning ability from LLMs . In <i>Proceedings of the 2024 Joint International Conference on Computational Linguistics, Language Resources and Evaluation (LREC-COLING 2024)</i> , pages 5538–5550, Torino, Italia. ELRA and ICCL.	812
760		813
761		814
762		815
763		816
764		817
765		818
766	Qingchen Yu, Zifan Zheng, Shichao Song, Zhiyu Li, Feiyu Xiong, Bo Tang, and Ding Chen. 2024. xfinder: Robust and pinpoint answer extraction for large language models . <i>arXiv preprint arXiv:2405.11874</i> .	819
767		820
768		821
769		822
770	Rowan Zellers, Ari Holtzman, Yonatan Bisk, Ali Farhadi, and Yejin Choi. 2019. Hellaswag: Can a machine really finish your sentence? <i>arXiv preprint arXiv:1905.07830</i> .	823
771		824
772		825
773	Wei Hu Muhao Chen Jian Dai Wei Zhang Yuzhong Qu Zequn Sun, Chengming Wang. 2020. Knowledge graph alignment network with gated multi-hop neighborhood aggregation. In <i>AAAI</i> , pages 222–229.	826
774		827
775		828
776		
	Qingru Zhang, Minshuo Chen, Alexander Bukharin, Pengcheng He, Yu Cheng, Weizhu Chen, and Tuo Zhao. 2023. Adaptive budget allocation for parameter-efficient fine-tuning. In <i>The Eleventh International Conference on Learning Representations</i> .	
	Zhaoyang Zhang, Wenqi Shao, Yixiao Ge, Xiaogang Wang, Jinwei Gu, and Ping Luo. 2024. Cached transformers: Improving transformers with differentiable memory cache. In <i>Proceedings of the AAAI Conference on Artificial Intelligence</i> , volume 38, pages 16935–16943.	
	Kun Zhao, Chenghao Xiao, Chen Tang, Bohao Yang, Kai Ye, Noura Al Moubayed, Liang Zhan, and Chenghua Lin. 2024a. X-ray made simple: Radiology report generation and evaluation with layman’s terms. <i>arXiv preprint arXiv:2406.17911</i> .	
	Kun Zhao, Bohao Yang, Chenghua Lin, Wenge Rong, Aline Villavicencio, and Xiaohui Cui. 2023. Evaluating open-domain dialogues in latent space with next sentence prediction and mutual information. In <i>Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)</i> , pages 562–574.	
	Kun Zhao, Bohao Yang, Chen Tang, Chenghua Lin, and Liang Zhan. 2024b. Slide: A framework integrating small and large language models for open-domain dialogues evaluation. <i>arXiv preprint arXiv:2405.15924</i> .	
	Zifan Zheng, Yezhaohui Wang, Yuxin Huang, Shichao Song, Mingchuan Yang, Bo Tang, Feiyu Xiong, and Zhiyu Li. 2024. Attention heads of large language models: A survey. <i>arXiv preprint arXiv:2409.03752</i> .	
	Zexuan Zhong, Mengzhou Xia, Danqi Chen, and Mike Lewis. 2024. Lory: Fully differentiable mixture-of-experts for autoregressive language model pre-training. <i>arXiv preprint arXiv:2405.03133</i> .	
	Barret Zoph, Irwan Bello, Sameer Kumar, Nan Du, Yanping Huang, Jeff Dean, Noam Shazeer, and William Fedus. 2022. St-moe: Designing stable and transferable sparse expert models. <i>arXiv preprint arXiv:2202.08906</i> .	
	Simiao Zuo, Xiaodong Liu, Jian Jiao, Young Jin Kim, Hany Hassan, Ruofei Zhang, Tuo Zhao, and Jianfeng Gao. 2021. Taming sparsely activated transformer with stochastic experts. <i>arXiv preprint arXiv:2110.04260</i> .	
	A Related Work	
	A.1 Parameter-Efficient Fine-Tuning (PEFT)	
	LLMs have demonstrated remarkable improvements in general natural language understanding (NLU) and natural language generation (NLG) tasks (Guo et al., 2023; Chang et al., 2024; Tang, 2024). However, their performance often suffers from limited generalization and effectiveness in domain-specific applications. To address this, numerous studies have proposed Parameter-Efficient Fine-Tuning (PEFT) methods,	

829 which optimize only a small subset of LLM param- 879
830 eters, significantly reducing computational overhead. 880

831 One of the most popular PEFT methods is Low- 881
832 Rank Adaptation (LoRA) (Hu et al., 2021b), which 882
833 introduces two low-rank matrices to each weight 883
834 matrix \mathbf{W} in LLMs. The product of these matrices 884
835 represents the weight adjustment $\Delta\mathbf{W}$. Building 885
836 upon LoRA, subsequent works have introduced more 886
837 efficient variants, such as VeRa (Kopiczko et al., 887
838 2023), AdaLoRA (Zhang et al., 2023), DoRA (Liu 888
839 et al., 2024), and MoSLoRA (Wu et al., 2024b). These 889
840 methods aim to better capture task-specific features 890
841 and integrate diverse feature subspaces effectively. 891

842 A.2 Transformer & Recurrent Models 892

843 The self-attention mechanism, as a fundamental algo- 893
844 rithm within Transformers (Vaswani, 2017), facilitates 894
845 parallel sequence processing and enhances model 895
846 capacity through increased width. However, it lacks 896
847 the temporal reasoning capabilities inherent in recur- 897
848 rent architectures, which contribute to model depth, 898
849 such as Long Short-Term Memory (LSTM) networks 899
850 (Hochreiter, 1997) and Gated Recurrent Units (GRUs) 900
851 (Chung et al., 2014). To address this limitation, recent 901
852 research has investigated integrating Transformer archi- 902
853 tectures with recurrent structures (Peng et al., 2023; 903
854 Pramanik et al., 2023; Zhang et al., 2024; Tang, 2024). 904
855 These efforts aim to blend recurrent neural network 905
856 capabilities with the ability to capture long-distance 906
857 dependencies. In this study, we likewise endeavor to 907
858 incorporate a recurrent mechanism to model multiple 908
859 expert systems. However, rather than focusing on 909
860 capturing long-distance dependencies, our primary 910
861 objective is to emulate the stepwise cognitive pro- 911
862 cesses characteristic of human cognition. To achieve 912
863 this, we employ GRUs to augment reasoning depth 913
864 by aggregating hidden representations from attention 914
865 features at each stage of the recurrent routing process. 915
866 This integration is designed to enhance the model’s 916
867 proficiency in apprehending complex dependencies 917
868 and thereby improve its overall reasoning capabilities. 918

869 A.3 Mixture of Experts (MoE) 919

870 The concept of Mixture of Experts (MoE) was first 920
871 introduced by Jacobs et al. (1991), who proposed 921
872 training multiple networks (experts) on different sub- 922
873 sets of data and aggregating their outputs. Recently, 923
874 as LLMs have become a focal point of research, MoE 924
875 layers have been integrated into Transformer-based 925
876 architectures. Specifically, researchers have replaced 926
877 standard Feed-Forward Networks (FFNs) with sparse 927
878 MoE layers, employing novel routing strategies (Zuo

879 et al., 2021; Zhong et al., 2024; Wu et al., 2024a; 879
880 Muqeeth et al., 2023; Fu et al., 2024) and advanced 880
881 expert segmentation techniques (Dai et al., 2024; He, 881
882 2024; Jiang et al., 2024; Xiao et al., 2024). 882

883 Numerous studies have investigated the application 883
884 of Parameter-Efficient Fine-Tuning (PEFT) methods 884
885 to introduce additional trainable parameters for 885
886 implementing pseudo MoE structures within LLMs 886
887 (Dou et al., 2024; Luo et al., 2024; Gao et al., 2024; Li 887
888 et al., 2024; Gou et al., 2023). These models expand 888
889 the conventional single feed-forward network (FFN) 889
890 architecture and its corresponding representation 890
891 space into multiple subspaces, thereby effectively 891
892 emulating the behavior of MoE architectures. 892

893 Prominent examples include MoLA (Gao et al., 893
894 2024), LoRAMoE (Dou et al., 2024), and MixLoRA 894
895 (Li et al., 2024), all of which have demonstrated 895
896 state-of-the-art performance compared to other 896
897 PEFT-based methods across various benchmarks. 897
898 Consequently, we integrate these three distinct LoRA 898
899 MoE methodologies into our GRAPHMOE frame- 899
900 work. The primary differences between these models 900
901 are as follows: LoRAMoE incorporates vanilla atten- 901
902 tion LoRA layers and MLP plugged with LoRA-MoE 902
903 layers; MoLA integrates plugged LoRA-MoE layers 903
904 in both attention and MLP layers; and MixLoRA 904
905 combines fused LoRA-MoE modules in attention and 905
906 MLP layers. Further analysis of these LoRA-MoE 906
907 implementations can be found in Li et al. (2024). 907

908 In this study, we augment MoE architectures by 908
909 enhancing their reasoning depth and demonstrate the 909
910 efficacy of our novel architecture when integrated 910
911 with PEFT-based LLM baselines. 911

912 B Pseudo Reasoning 912 913 Graph Construction Algorithm 913

914 To enhance clarity in illustrating the long processing 914
915 pipeline, we summarize the entire methodology in 915
916 Algorithm 1. The process involves T reasoning iter- 916
917 ations, referred to as self-rethinking. Initially, features 917
918 from the FFN layer of transformer blocks are derived 918
919 via a LoRA MoE forward pass. Subsequently, these 919
920 features are transmitted to the subsequent Low-Rank 920
921 GRU module, which computes both the residual 921
922 connection within the pseudo graph and the merged 922
923 feature for the input of the succeeding iteration. 923

924 C Additional Implementation Details 924

925 Our hyperparameters strictly follow the configurations 925
926 detailed in (Li et al., 2024) including the LoRA and 926
927 DoRA methods and their three MoE derivatives, e.g. 927

Algorithm 1: GRAPHMOE at Layer ℓ

Input: $\mathbf{x}_1^\ell: (\mathbf{B}, \mathbf{P}, \mathbf{D})$ after the Post-Norm.
Output: $\mathbf{y}_T^\ell: (\mathbf{B}, \mathbf{P}, \mathbf{D})$ the hidden states.

```

1 /*  $\mathbf{B}$  for Batch Num,  $\mathbf{P}$  for Seq Len,  $\mathbf{D}$  for Dim */
2 /*  $\mathbf{D}'$  for Upper Dim,  $\bar{\mathbf{D}}$  for GRU Hidden Dim */
3 for  $t \leftarrow 1$  to  $T$  do
4   /* MoE Forward */
5    $\tilde{\mathbf{s}}_t^\ell: (\mathbf{B}, \mathbf{P}, \mathbf{N}) \leftarrow \text{Softmax}(\mathbf{R}^\ell * \mathbf{x}_t^\ell)$ ;
6    $\mathbf{s}_t^\ell: (\mathbf{B}, \mathbf{P}, \mathbf{N}) \leftarrow \text{Top-}k(\tilde{\mathbf{s}}_t^\ell)$ ;
7   for  $i$  in  $\{\text{selected experts } E\}$  do
8      $\tilde{\mathbf{W}}_{G_i}^\ell: (\mathbf{D}', \mathbf{D}) \leftarrow \mathbf{W}_{G_i}^\ell + \mathbf{B}_{G_i}^\ell * \mathbf{A}_{G_i}^\ell$ ;
9     Similarly, we can get  $\tilde{\mathbf{W}}_{U_i}^\ell$  and  $\tilde{\mathbf{W}}_{D_i}^\ell$ ;
10     $\hat{\mathbf{c}}_i^G: (\mathbf{B}, \mathbf{P}, \mathbf{D}') \leftarrow \tilde{\mathbf{W}}_{G_i}^\ell * \mathbf{x}_t^\ell$ ;
11     $\hat{\mathbf{c}}_i^U: (\mathbf{B}, \mathbf{P}, \mathbf{D}') \leftarrow \tilde{\mathbf{W}}_{U_i}^\ell * \mathbf{x}_t^\ell$ ;
12     $\hat{\mathbf{c}}_i: (\mathbf{B}, \mathbf{P}, \mathbf{D}') \leftarrow \sigma(\hat{\mathbf{c}}_i^G) \odot \hat{\mathbf{c}}_i^U$ ;
13     $\mathbf{c}_i: (\mathbf{B}, \mathbf{P}, \mathbf{D}) \leftarrow \tilde{\mathbf{W}}_{D_i}^\ell * \hat{\mathbf{c}}_i$ ;
14     $\mathbf{y}_t^\ell: (\mathbf{B}, \mathbf{P}, \mathbf{D}) \leftarrow \mathbf{y}_t^\ell + \mathbf{s}_t^\ell[:, i] \odot \mathbf{c}_i$ ;
15  /* Low-Rank GRU */
16  if  $t \neq T$  then
17     $\mathbf{z}_t^\ell: (\mathbf{B}, \mathbf{P}, \bar{\mathbf{D}}) \leftarrow \sigma(\mathbf{W}_z^\ell * [\mathbf{h}_{t-1}^\ell, \mathbf{y}_t^\ell])$ ;
18     $\mathbf{r}_t^\ell: (\mathbf{B}, \mathbf{P}, \bar{\mathbf{D}}) \leftarrow \sigma(\mathbf{W}_r^\ell * [\mathbf{h}_{t-1}^\ell, \mathbf{y}_t^\ell])$ ;
19     $\hat{\mathbf{h}}_t^\ell: (\mathbf{B}, \mathbf{P}, \bar{\mathbf{D}} + \mathbf{D}) \leftarrow [\mathbf{r}_t^\ell \odot \mathbf{h}_{t-1}^\ell, \mathbf{x}_t^\ell]$ ;
20     $\mathbf{h}_t^\ell: (\mathbf{B}, \mathbf{P}, \bar{\mathbf{D}}) \leftarrow \sigma(\mathbf{W}_o^\ell * \hat{\mathbf{h}}_t^\ell + \mathbf{b}_o^\ell)$ ;
21     $\mathbf{h}_t^\ell \leftarrow (\mathbf{1} - \mathbf{z}_t^\ell) \odot \mathbf{h}_{t-1}^\ell + \mathbf{z}_t^\ell \odot \mathbf{h}_t^\ell$ ;
22     $\mathbf{g}_t^\ell: (\mathbf{B}, \mathbf{P}, \mathbf{D}) \leftarrow \mathbf{W}_g^\ell \mathbf{h}_t^\ell$ ;
23    /* Residual Connection */
24     $\mathbf{x}_{t+1}^\ell: (\mathbf{B}, \mathbf{P}, \mathbf{D}) \leftarrow \mathbf{x}_t^\ell + \mathbf{g}_t^\ell$ ;
25 return  $\mathbf{y}_T^\ell: (\mathbf{B}, \mathbf{P}, \mathbf{D})$ ;

```

MixLoRA or MixDoRA. We summarize the settings for LoRA/DoRA and their derivatives of MoE-based methods, as shown in Table 6. During evaluation, the batch size is set to 8.

Both techniques utilized a cutoff length of 512, a learning rate of $2e - 4$, the AdamW optimizer, a batch size of 16, accumulation steps of 8, a dropout rate of 0.05, and were trained over 2 epochs. In LoRA/DoRA, the methods were applied to the Q, K, V, O, Up, Down, and Gate parameters with a LoRA rank r of 80 and a LoRA alpha α of 160. Conversely, MixLoRA/MixDoRA targeted the same parameters with a LoRA rank r of 16 and a LoRA alpha α of 32, incorporating 8 experts with a Top- k value of 2.

The experiments described in this paper required a single A800 GPU for at least one month. Including debugging and downtime, the entire process extended to around two months. The base model, LLaMA-8B,

Setting	LoRA/DoRA	MixLoRA/MixDoRA	GRAPHMOE
Cutoff Length		512	
Learning Rate		$2e-4$	
Optimizer		AdamW	
Batch Size		16	
Accumulation Steps		8	
Dropout Rate		0.05	
Training Epochs		2	
LoRA Rank (r)	80		16
LoRA Alpha (α)	160		32
Expert Number	-		8
Top- k	-		2
Targeted Parameters	$\mathbf{Q}_{\text{Gate}}^\ell, \mathbf{K}_{\text{Up}}^\ell, \mathbf{V}_{\text{Up}}^\ell, \mathbf{O}_{\text{Down}}^\ell, \mathbf{W}_{\text{Gate}}^\ell, \mathbf{W}_{\text{Up}}^\ell, \mathbf{W}_{\text{Down}}^\ell$		$\mathbf{Q}_{h'}^\ell, \mathbf{K}_{h'}^\ell, \mathbf{V}_{h'}^\ell, \mathbf{O}_{h'}^\ell, \mathbf{W}_{\text{Gate}}^\ell, \mathbf{W}_{\text{Up}}^\ell, \mathbf{W}_{\text{Down}}^\ell, \mathbf{W}_z^\ell, \mathbf{W}_r^\ell, \mathbf{W}_o^\ell, \mathbf{W}_g^\ell$

Table 6: Training Settings for baselines and GRAPHMOE.

was obtained from the Huggingface repository “meta-llama/Meta-Llama-3-8B-Instruct”³.

For experiments utilizing the GRAPHMOE framework, the key hyperparameters were configured as follows: (1) The number of GRU gate iterations within the recurrent routing—referred to as the reasoning (or self-rethinking) round T —was set to a default of 3. In the initial round, expert outputs are obtained without feature transmission via the GRU gate. (2) The GRU network consists of a single layer. (3) The GRU hidden size was scaled to 0.1×4096^4 , where 4096 corresponds to the model dimension of LLaMA-3-8B. (4) The hyperparameter λ for the auxiliary load balancing loss \mathcal{L}_{LB} was set to 0.01. (5) Unless otherwise specified, the MoE architecture employed a total of 8 experts per module, with 2 experts activated per pass.

We address several key questions regarding our implementation and design choices in the following discussion:

- **What is the value of k in the top- k selection? How is it selected?** Currently, k is treated as a hyperparameter, with two main considerations. First, aligning with other MoE architectures, the number of selected experts typically maintains fixed and sparse, with 2 out of 8 being a common configuration. Second, our framework facilitates expert collaboration, so the top k is not a critical parameter for fully leveraging expert potential.
- **Is the top- k selection in the routing operation (Equation 4) not differentiable?** Regarding the differentiability issue of the top- k in Equation 4, it is important to note that while the top- k operation itself is non-differentiable, it does not directly impact the gradient calculations. Rather,

³<https://huggingface.co/meta-llama/Meta-Llama-3-8B-Instruct>

⁴The actual hidden size is rounded to the nearest integer using the Python function “int()”.

980 the fixed number of experts chosen through the
981 top-k operation contribute to the subsequent
982 gradient computations.

983 • **Are some experts learning to do the initial**
984 **processing, while others more frequently do**
985 **the later processing?** In this study, the MoE
986 is applied within each transformer block, so
987 the experts do not possess specific physical
988 significance. Therefore, we have not conducted
989 an analysis of individual expert usage patterns
990 or preferences within the reasoning process.