

MoPE: A Massive Mixture of Passage-Level Experts for Knowledge Editing in Open-Domain Question Answering

Anonymous ACL submission

Abstract

As world knowledge continues to evolve, adapting LLMs to new knowledge is crucial, however, it poses significant challenges, as naively fine-tuning the entire model often leads to catastrophic forgetting and high computational costs. While RAG and model editing have been increasingly studied for knowledge adaptation, this paper moves beyond the ‘RAG vs. fine-tuning’ discussion to explore the ‘RAG vs. model editing’ issue, and propose a “Massive Mixture of Experts (MMoE)” approach for model editing, referred to as **MoPE**, i.e., Massive Mixture of Passage-Level Experts, which consists of the key components at training and inference stages: (1) **Massive passage-level editing with MMoE**, where a large set of passage-level experts is created using automatically generated question-answer pairs for each passage, and (2) **Retrieval-based routing with MMoE**, which employs dense retrieval to select the top- k passage-level experts without requiring additional training. Experimental results demonstrate that MoPE outperforms a naively designed variant of RAG, i.e., direct RAG, and when combined with direct RAG, it surpasses an advanced variant of RAG, significantly improving over LoRA-based parameter-efficient tuning methods. Our data and code will be available at <https://github.com/XXX/XXX>.

1 Introduction

Large language models (LLMs) have shown remarkable performances on various natural language processing (NLP) tasks, particularly effectively performing knowledge-intensive tasks by their enormously large pretrained knowledge (Zhao et al., 2023; Hadi et al., 2023). Given that world knowledge is continually updated, adapting LLMs to new information and knowledge is crucial, however, naively finetuning the entire model of LLMs meets the key challenging problems: 1) catastrophic forgetting, which may lead to disruption of

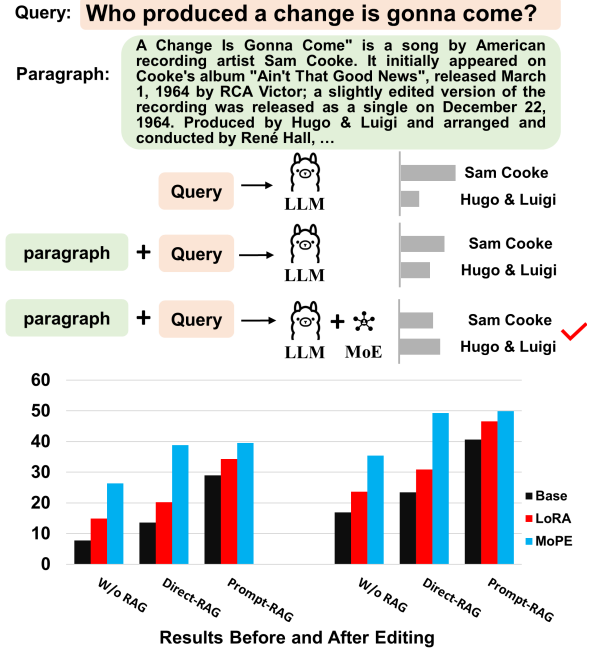


Figure 1: An illustration of context-aware attention: mitigating content focus limitations in RAG through paragraph editing, and a brief overview of the main results of MoPE comparing to the base model with our without RAG, on the NQ dataset. **Left:** EM metric, **Right:** F1 Score. MoPE outperforms direct RAG, which relies on a naively designed prompt for retrieval. When combined with direct RAG, MoPE surpasses prompt-RAG, which employs a more refined prompt variant.

the old knowledge and performance graduation on the pre-acquired tasks, and 2) nontrivial updating costs due to the large scale of parameters, which requires high computation and memory resource required for updating new knowledge. Previous adaptation approaches can be broadly categorized into two main methods: (1) **Retrieval-Augmented Generation (RAG)** (Lewis et al., 2020; Pan et al., 2024), which leverages the in-context learning (ICL) capabilities of large language models (LLMs) by augmenting input prompts with “retrieved” passages; and (2) **Knowledge Injection**, which can be

further divided into *parameter-efficient tuning* (Hu et al., 2021) – which adjusts a small set of trainable parameters while keeping the original ones intact – and *model editing* (De Cao et al., 2021; Meng et al., 2022a), which aims to balance the incorporation of new knowledge updates while preserving existing knowledge.

Among these approaches, unlike RAG and parameter-efficient tuning, which have demonstrated improvements across various NLP tasks, Model editing has been explored primarily in standard ‘editing’-specific tasks, rather than in practical knowledge-intensive NLP tasks, leaving its usefulness and impact in real-world applications unclear. Going beyond traditional editing-specific tasks, we aim to explore model editing within open-domain QA as a representative knowledge-intensive NLP task to assess its value and impact in achieving improvements over RAG, as open-domain QA has been widely used. Related to this issue, existing works on fine-tuning methods have explored the debate of ‘RAG vs. fine-tuning’ (Ovadia et al., 2023; Alghisi et al., 2024; Gupta et al., 2024), showing that parameter-efficient tuning and fine-tuning generally underperform compared to RAG, and their combination with RAG has been found to be advantageous, etc. Advancing beyond previous works, we extensively explore “model editing,” rather than fine-tuning or parameter-efficient tuning, so addressing ‘RAG vs. model editing’, leading to our key research question: “Are specific model editing methods effective in achieving improvements over RAG on standard open-domain QA tasks, compared to results obtained using fine-tuning or parameter-efficient tuning?”

Unlike editing-specific tasks, applying model editing to open-domain QA involves handling a set of *passage-level edits* within a collection, thereby requiring the editor to effectively inject the required passages into LLMs, treating each passage as an edit request. To address this passage-level editing problem, inspired by the remarkable adaptability of the mixture of experts (MoE) on knowledge-intensive and editing tasks (Wang and Li, 2024b,a), we propose the use of a massive MoE (MMoE) for model editing to inject all required passages into a set of experts, referred to as MoPE, i.e., massive mixture of passage-level experts, by assigning a dedicated passage-level expert to each individual passage and integrating these passage-level experts into LLMs, using the “retrieval”-based router, without modifying the original parameters. MoPE

consists of the key components for editing and inference stages, as follows:

1. Massive passage-level editing with MMoE.

In the editing stage, all target passages are considered as massive edit requests. Similar to most MoE approaches that utilize feed-forward network (FFN) layers (Fedus et al., 2022; Du et al., 2022), we augment an FFN layer in the Transformer with an additional expertized FFN module, referred to as a passage-level expert, specialized for each individual passage; Given a passage, we first generate a set of question-answer pairs using a simple prompting method, and then train the passage-level expert based on the loss function used in the machine reading comprehension (MRC) task. During the editing stage, it should be noted that each passage-level expert is trained individually, i.e., only a single passage-level expert is augmented with the original FFN layers in the Transformer, without incorporating all other experts. As a result, we will have a massive mixture of N static passage-level experts, given N number of passages in a collection.

2. Retrieval-based Router with MMoE.

In the inference stage, since we have passage-level experts, routing mechanism is effectively managed through retrieval. We extensively utilize dense retrieval to design a router, ensuring that only the top- k experts corresponding to the top- k retrieved passages are selected to form a sparse MoE layer, resulting in consisting $k + 1$ FFNs at a specific layer. Thus, no separate additional training is required for designing a router.

Experimental results on standard open-domain QA datasets – natural question (NQ) (Kwiatkowski et al., 2019) and Trivia QA (Joshi et al., 2017) show that MoPE improves the baseline performance of RAG and the results using the LoRA-based parameter efficient tuning, as briefly shown in Figure 1: MoPE outperforms “direct-RAG,” which relies on a naively designed prompt for retrieval; MoPE combined with direct-RAG surpasses prompt-RAG, which utilizes a more refined prompt variant. For details on the input template, see Appendix A.

Our main contributions are as follows: 1) We investigate the ‘RAG vs. model editing’ paradigm, expanding beyond previous studies that focused on

‘RAG vs. fine-tuning’ (Alghisi et al., 2024). 2) We propose MoPE, a novel approach for model editing in open-domain QA. 3) We offer new experimental findings, showing that MoPE outperforms LoRA-based tuning and direct RAG, and the combination of MoPE and direct RAG surpasses even the more advanced variant of RAG, i.e., prompt RAG.

2 Related Work

2.1 Retrieval Augmented Generation

In open-domain question answering, Retrieval-Augmented Generation (RAG) (Lewis et al., 2020) methods effectively enhance the accuracy of model responses by incorporating external knowledge. Recent studies have made significant progress in retrieval quality control and timing optimization. Shi et al. (2023) and Dong et al. (2024) emphasized the importance of enhancing the relevance of retrieved content to user query and the role of knowledge alignment in this process. Lin et al. (2023) proposed an approach that optimizes the query encoder by combining supervised and unsupervised tasks. Additionally, the Fusion-in-Decoder (FiD) method (Izacard and Grave, 2020; Hofstätter et al., 2023) significantly improved generation performance by integrating information from multiple documents during the decoding stage.

In structural model optimization, Self-RAG (Asai et al., 2023) enhances retrieval accuracy and robustness through fine-grained self-reflection, while Yan et al. (2024) and Jiang et al. (2023b) introduced an error-correction strategy to handle inaccurate information. SAIL (Luo et al., 2023) improves instruction tracking in complex contexts by integrating internal and external search engines.

For post-generation output enhancement, KNN-LM (Khandelwal et al., 2019) and RETOMATON (Alon et al., 2022) significantly improve generation quality through nearest-neighbor memory retrieval and weighted finite-state machines. Furthermore, end-to-end training (Lewis et al., 2020; Guu et al., 2020) has emerged as a crucial direction to optimize RAG architectures, allowing more efficient knowledge integration and response generation by minimizing manual intervention and iteratively refining the model.

2.2 Knowledge Injection

Knowledge injection has emerged as a crucial technique for enhancing the performance of language

models by efficiently integrating external knowledge. Two primary approaches have gained significant attention: parameter-efficient tuning and model editing.

Parameter-efficient tuning enhances large pre-trained models by introducing small trainable modules while keeping most parameters frozen. Key approaches include adapter (Houlsby et al., 2019), which adds neural network bottlenecks to transformer blocks; Prompt Tuning (Li and Liang, 2021), which optimizes the appended prompts for task adaptation; and LoRA (Hu et al., 2021), which updates rank decomposition matrices. Recent advances, such as DyLoRA (Valipour et al., 2022), improve efficiency by selectively updating partial parameters. Building on DyLoRA, MELO (Yu et al., 2024) introduces a neuron-indexed dynamic LoRA mechanism.

Model editing focuses on maintaining the reliability of edited knowledge, ensuring that the changes successfully address the target queries. Additionally, it emphasizes enhancing the generality, allowing the edited model to generalize the new knowledge to related queries effectively. Furthermore, it seeks to preserve locality, ensuring that the modifications do not interfere with the retention of unrelated original knowledge. It is categorized into three types: Meta-learning editors (De Cao et al., 2021; Mitchell et al., 2021; Tan et al., 2023), which use hyper-networks to adjust gradients; Locate-then-edit editors (Meng et al., 2022a,b; Li et al., 2024), which identify and update relevant parameters; and Memory-based editors, where (Zheng et al., 2023; Zhong et al., 2023; Gu et al., 2023; Cheng et al., 2023) update knowledge from prompts using in-context learning without gradient updates or parameter modifications, where other approaches like T-Patcher (Huang et al.) and MEMoE (Wang and Li, 2024b) store the memory of edited facts using additional parameters.

2.3 Mixture of Experts

In transformer-based Large Language Models (LLMs), Mixture-of-Experts (MoE) layers utilize a set of expert networks and a gating mechanism to route inputs to the most suitable experts (Shazeer et al., 2017; Antoniak et al., 2023). These layers are strategically positioned after the self-attention sub-layer to optimize feed-forward network (FFN) selection, significantly reducing computational overhead in large models like PaLM (Chowdhery et al., 2023), where FFN layers account for the majority

of parameters.

Dense MoE approaches activate all available experts simultaneously, which enhances predictive accuracy but demands substantial computational resources. Early implementations (Jacobs et al., 1991; Rasmussen and Ghahramani, 2001; Aljundi et al., 2017) demonstrated this effectiveness, and more recent methods like EvoMoE (Nie et al., 2021), MoLE (Wu et al., 2024), LoRAMoE (Dou et al., 2023), and DSMoE (Pan et al., 2024) have refined the dense MoE structure to balance performance and efficiency.

Sparse MoE improves computational efficiency by selecting only the top- k experts for each input, thereby maintaining accuracy while reducing processing demands (Shazeer et al., 2017). However, this selective activation can cause load imbalances, where certain experts are overused while others are underutilized. To counter this, auxiliary loss functions are introduced to distribute tokens more evenly across experts, as seen in (Lepikhin et al., 2020; Jiang et al., 2024; Du et al., 2022; Fedus et al., 2022). This strategy allows sparse MoE models to scale effectively by expanding parameter capacity without a corresponding increase in computational cost.

3 Methodology

Figure 2 presents an overview of MoPE, highlighting its training and inference stages: (1) *Passage-level editing with MMoe* and (2) Inference using a *retrieval-based router with MMoe*. In the passage-level editing stage, a set of question-answer (QA) pairs is automatically generated for each passage using named entity recognition (NER)-based answer extraction and a simple prompting method. Each passage-level expert is trained individually using the generated QA pairs, following an MRC-like objective function. Given N passages, we construct N corresponding passage-level experts. During inference, given these N passage-level experts and a test question, we first perform retrieval-based routing. Specifically, we use dense retrieval with a reranking method to identify the top- k retrieved passages. The corresponding k passage-level experts are then selected from the N experts and integrated with the original FFN layer. This selected sparse MoE is used to generate an answer to the test question.

3.1 Passage-level Editing with MMoe

Suppose that there is N number of passages in a collection, formulated as $\mathcal{C} = \{p_i\}_{i=1}^N$ where \mathcal{C} is a collection, and p_i is i -th passage¹. The passage-level editing process creates N passage-level experts, comprising two key steps: (1) the *construction of QA pairs* and (2) the *training passage-level experts*, as detailed below.

3.1.1 Construction of QA pairs

Given a passage, we create its set of QA pairs by using the automatic data augmentation method (Yu et al., 2022). To this end, we first separately train a question generator G_θ based on the T5-xl model to learn how to generate questions from textual segments, based on an additional training dataset $D_{QG} = \{(p_{\psi(j)}, q_j, a_j)\}_{j=1}^M$, where $p_{\psi(j)}$ is the $\psi(j)$ -th passage in the training set², q_j is the j -th question for $p_{\psi(j)}$, a_j is the j -answer for q_j on $p_{\psi(j)}$. It is important to note that D_{train} is separate from the test dataset in the final evaluation. In our case, we use only the training samples from the NQ dataset. We also have the instruction template function \mathcal{T} is defined as follows: "Generate a question with <answer> as ' $\{a\}$ ' from the <context>: ' $\{p\}$ '.". Thus, $\mathcal{T}(p, a)$ indicates the prompted input that takes an answer a for a passage p . The generator G_θ is trained by minimizing the negative log-likelihood (NLL) loss as follows:

$$\mathcal{L}(\theta) = - \sum_{(p,q,a) \in \mathcal{D}_{QG}} \log P(q|\mathcal{T}(p, a); \theta) \quad (1)$$

where $P(y|x; \theta)$ represents the conditional probability of generating question y given input x .

Once G_θ is trained, given a i -th passage $p_i \in \mathcal{C}$, we first extract named entities in p_i and answers using a NER method, and generate corresponding questions to the extracted answers. The resulting set of QA pairs for p_i is denoted as \mathcal{D}_{p_i} . The generation process please refer to Appendix B.

3.1.2 Training Passage-Level Experts

Once \mathcal{D}_{p_i} is obtained, the editing step uses \mathcal{D}_{p_i} as training set to create the corresponding passage-level expert E_i under the MoE framework. To

¹While we need to create experts for all passages in a collection, evaluation only requires the experts corresponding to the top-retrieved passages for test questions. Therefore, instead of generating a full set of experts, we construct only the minimum necessary number of experts, ensuring that the subset is sufficient for evaluation.

²Here, we use ψ as an existential Skolem-like function that maps the j -th training example to its corresponding passage index, for notational convenience.

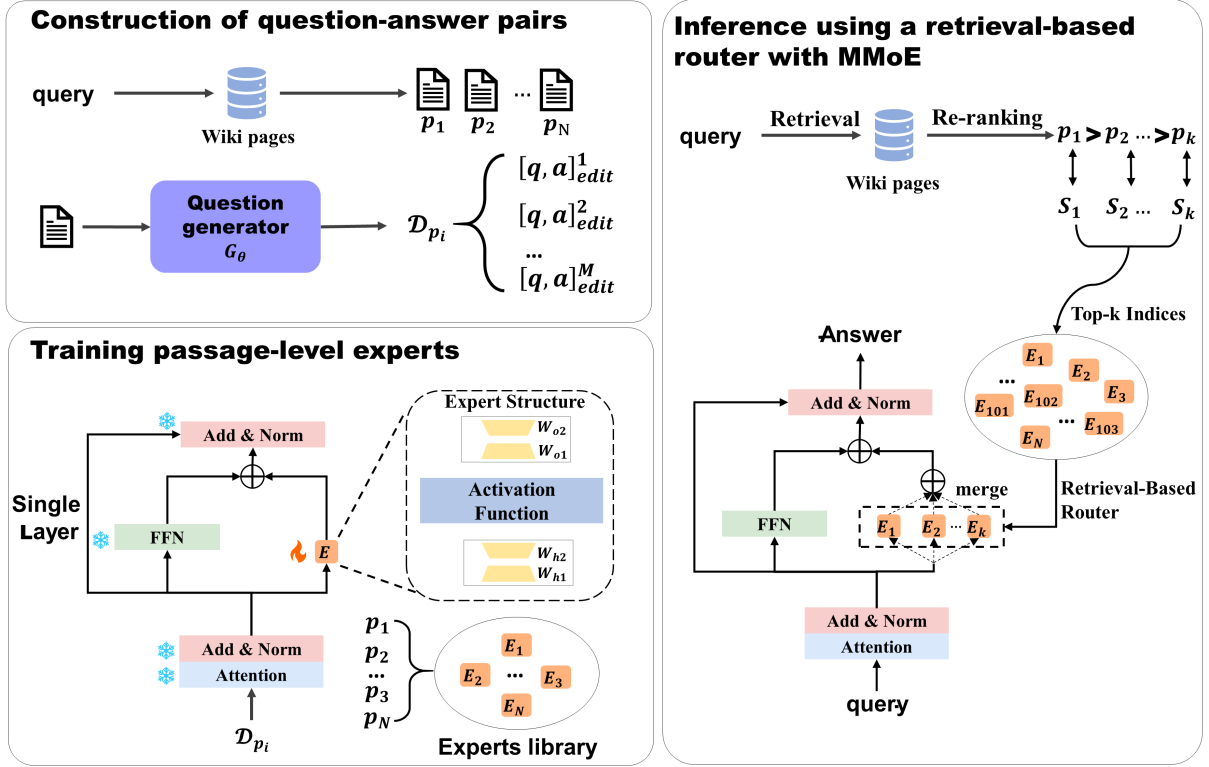


Figure 2: An overview of the proposed MoPE framework: 1) Passage-Level Editing with MMoE: For a given passage, its QA pairs are automatically constructed using the fine-tuned question generator (see Section 3.1.1). The passage-specific expert is then trained within the MoE layer using these QA pairs. 2) Inference with Retrieval-Based Router: MMoE is deployed, where dense retrieval and a fine-tuned reranker are applied to select the top-retrieved relevant passages (see Section 3.2.1), and the corresponding passage-level experts are then incorporated to form a sparse MoE. The relevance vector of the router is computed based on either the top-1 or top- k passages, using Eq. (5) and Eq. (12), respectively.

formally describe the MoE setting, suppose that at each token position t , $\mathbf{x}_t^l \in \mathbb{R}^{d_m}$ is given input at layer l . The original FNN block, denoted as $\text{FFN}^l(\mathbf{x}_t^l)$, is formulated as follows:

$$\text{FFN}^l(\mathbf{x}_t^l) = \text{relu}(\mathbf{x}_t^l \mathbf{W}_K^l) \mathbf{W}_V^l \quad (2)$$

where $\mathbf{W}_K^l \in \mathbb{R}^{d_m \times d}$ and $\mathbf{W}_V^l \in \mathbb{R}^{d \times d_m}$ represent the key and value projection matrices of the original FNN layer, respectively. For a i -th passage, we now have the passage-level expert, $E_i^l(\mathbf{x}_t^l)$ is formulated, based on another FNN block, as follows:

$$E_i^l(\mathbf{x}_t^l) = \text{relu}(\mathbf{x}_t^l \mathbf{W}_{i,K}^l) \mathbf{W}_{i,V}^l \quad (3)$$

where $\mathbf{W}_{i,K}^l \in \mathbb{R}^{d_m \times d}$ and $\mathbf{W}_{i,V}^l \in \mathbb{R}^{d \times d_m}$ represent the key and value projection matrices of the expert FNN layer for i -th passage, respectively. Given N passages, we have N passage-level experts, thereby forming the following MMoE:

$$\mathbf{y}_t^l = \text{FFN}^l(\mathbf{x}_t^l) + \lambda_t^l \sum_{j=1}^N \mathbf{r}[j] \cdot E_j^l(\mathbf{x}_t^l) \quad (4)$$

where λ_t^l is the mixing parameter, which is set to 1 for a specific layer l and 0 otherwise depending on token index t , and $\mathbf{r} \in \mathbb{R}^N$ is the “relevance” vector to a given query q , obtained by the “retrieval-based router,” which consists of N relevance values, defined as follows:

$$\mathbf{r}[j] = \begin{cases} 1, & \text{if } p_j \text{ in top-}k(q, \mathcal{C}) \\ 0, & \text{else} \end{cases} \quad (5)$$

where $\text{top-}k(q, \mathcal{C})$ is the set of the top- k retrieved passages, in general. During the editing, when $q \in \mathcal{D}_{p_i}$, $\text{top-}k(q, \mathcal{C}) = \{p_i\}$, because it is clear that p_i is the gold relevant passage for q . Thus, during editing, given \mathcal{D}_{p_i} , Eq. (4) is simplified into:

$$\mathbf{y}_t^l = \text{FFN}^l(\mathbf{x}_t^l) + \lambda_t^l E_i^l(\mathbf{x}_t^l) \quad (6)$$

Under Eq. (6), given \mathcal{D}_{p_i} , we use an autoregressive generative loss function to train the i -th passage-level expert E_i , independently from other experts.

Furthermore, to effectively maintain parameters of E_i , instead of using the full parameters of FNN

as in Eq. (3), we perform low-rank matrix factorization on its standard linear transformation layers to reduce computational overhead and parameter redundancy, thereby decomposing the fully connected weight matrices into low-rank components.

$$E_i^\ell(\mathbf{x}_i^\ell) = \text{relu} \left(\mathbf{x}_i^\ell \mathbf{W}_{i,K2}^l \mathbf{W}_{i,K1}^l \right) \mathbf{W}_{i,V1}^l \mathbf{W}_{i,V2}^l \quad (7)$$

Where $\mathbf{W}_{i,K2}^l \in \mathbb{R}^{d_m \times k}$, $\mathbf{W}_{i,K1}^l \in \mathbb{R}^{k \times d_m}$, $\mathbf{W}_{i,V1}^l \in \mathbb{R}^{d_m \times k}$ and $\mathbf{W}_{i,V2}^l \in \mathbb{R}^{k \times d_m}$ represent the corresponding low-rank key and value projection matrices of the i -th passage-level expert.

3.2 Inference with a Retrieval-Based Router with MMoE

Once MMoE is trained during the passage-level editing stage, inference begins by retrieving relevant passages for a given test query, based on dense retrieval, followed by a reranking step to refine the results. Then, the retrieval-based router is applied to determine the relevance vector \mathbf{r} in Eq. (4), based on similarity scores obtained from retrieval. This process depends on two different variants of using MMoE: 1) *Single-expert*: Uses only the top-1 expert, where the relevance vector remains the same as in Eq. (5) with $k = 1$ used during the editing stage. 2) *Multi-expert*: Uses the top- k passage-level experts, where the relevance vector is a soft version computed using a softmax-based probability distribution after top- k truncation.

3.2.1 Retrieval-Reranker

Given a query q , DPR uses both the query vector $\text{Emd}(q)$ and the passage vector $\text{Emd}(p)_i$ for p_i , and computes inner product similarity between them:

$$\text{score}(q, p_i) = \text{Emd}(q)^\top \text{Emd}(p)_i \quad (8)$$

For the reranking, we incorporate *BAAI/bge-reranker-v2-gemma*³ (Xiao et al., 2024) as a reranker model and fine-tune it to improve passage relevance. For the finetuning the reranker, we additionally construct a training dataset passages retrieved by DPR. To optimize the reranker, we adopt a contrastive loss function (Sohn, 2016), ensuring that relevant passages receive higher scores than irrelevant ones. For a given question q in the training set, let p_+ be a positive passage, and let $(p_-^{(1)}, \dots, p_-^{(L)})$ be L negative passages for q . The

loss function is defined as:

$$\mathcal{L}_{\text{reranker}} = -\log \frac{\exp(\text{score}_{\text{reranker}}(q, p_+))}{\sum_{i=1}^{L+1} \exp(\text{score}_{\text{reranker}}(q, p^{(i)}))} \quad (9)$$

The re-ranker results can be found in Appendix C.

3.2.2 Retrieval-Based Router

The retrieval-based router determines the relevance vector \mathbf{r} in Eq. (4), depending on we use the top-1 and top- k experts.

Single-Expert: In the single-expert case ($k = 1$), the relevance vector is just one-hot vector, the router-applied MMoE layer is degenerated into Eq. (6).

Multi-Expert: In the multi-experts case, the relevance vector based on the dense retrieval is computed as follows:

$$\mathbf{r} = \text{softmax} \left(\text{top}_k \left(\text{Emd}(q)^\top W_C \right) \right) \quad (10)$$

where $\text{Emd}(q) \in \mathbb{R}^{d_{\text{emd}}}$ is the query embedding vector, d_{emd} is the dimension of embedding vectors for query and passages, and $W_C \in \mathbb{R}^{d_{\text{emd}} \times N}$ is the matrix consisting of all passage embeddings in the collection \mathcal{C} , defined as follows:

$$W_C = [\text{Emd}(p_1), \dots, \text{Emd}(p_N)] \quad (11)$$

When reranking is applied, we then use reranking-based scores over all passages, resulting in the relevance vector, with a slight abuse of notation, as follows:

$$\mathbf{r} = \text{softmax} \left(\text{top}_k \left([\text{score}_{\text{reranker}}(q, p_i)]_{i=1}^N \right) \right) \quad (12)$$

4 Experiments

In this part, we explain the experimental setup and present the key results of our findings.

4.1 Experimental Setup

Dataset: We utilize the Natural Questions (NQ) dataset (Kwiatkowski et al., 2019) for training to enhance the model’s performance on open-domain question answering. To further evaluate the model’s generalization ability, we also test it on the TriviaQA (TQA) (Joshi et al., 2017) dataset, which presents diverse and challenging questions from trivia domains.

4.2 RAG settings

We explore two variants of RAGs – **Direct-RAG**, and **Prompt-RAG** – using prompting templates to evaluate the RAG capabilities of the models, with the input format shown in Table 4.

³<https://huggingface.co/BAAI/bge-reranker-v2-gemma>

4.3 MoPE settings

General: For a given test query q_{test} , we first perform retrieval to identify top- k experts and apply the retrieval-based router by determining relevance vectors depending on single-expert and multi-expert cases, as in Section 3.2.2. Suppose that \mathcal{M}_{MoPE} is the model using MMoE based on Eq. (4) at λ_l^t is 1 only for $l = 12$. The final answer a is determined as follows.

$$a = \arg \max_a \mathcal{M}_{MoPE}(a|q_{test}) \quad (13)$$

RAG on MoPE We apply RAG under MoPE but use a mixed approach that combines the original pre-edited base model and the MoPE-based model. More specifically, we use the original model until the passage knowledge is encoded, and then switch to MoPE when processing the question part.

Thus, given the specific layer l where new experts are additionally located, λ_l^t is defined as:

$$\lambda_l^t = \begin{cases} 0, & \text{if } t \leq L_p \text{ or } l \neq l_{moe} \\ 1, & \text{otherwise} \end{cases} \quad (14)$$

where L_p is the last token position where passage knowledge is injected into a prompted RAG, and l_{moe} is the layer into which experts are newly added.

Backbone Language Models: We select Llama2-7B (Touvron et al., 2023) as the base model and integrate expert modules into it to implement our framework. We also compare our approach with several advanced models (Yang et al., 2023; Team, 2024; Jiang et al., 2023a) built on Llama2-7B, as well as the larger Llama2-13B model (Touvron et al., 2023) with more parameters.

Metric: We use Exact Match (EM) and F1-score (F1), which are standard metrics (Rajpurkar, 2016) in question answering for evaluating answer accuracy and completeness.

Implementation Details: All experiments, including data construction, knowledge injection, and evaluation, were conducted on workstations with 8xNVIDIA RTX A6000 GPUs. For training the experts, we used the AdamW optimizer for 10 epochs, with the learning rate decaying from 1e-4 to 1e-6 using cosine annealing.

4.4 Why $l=12$ was Chosen?

(Wang et al., 2023) suggests that the information captured at different layers can vary in its contribution to the final output, especially for inputs with

Methods	NQ		TQA	
	EM	F1	EM	F1
<i>Without-RAG</i>				
Llama2 _{7B}	7.80	16.89	49.47	59.38
Llama2 _{13B}	10.40	20.62	50.87	61.21
Baichuan2 _{7B}	8.60	16.84	43.73	53.29
QWen2.5 _{7B}	1.20	6.24	31.53	40.19
Mistral _{7B}	3.27	12.11	52.60	61.71
PEFT-Lora _{7B}	14.87	23.66	51.73	61.37
MoPE _{7B}	26.40	35.40	56.13	64.08
<i>Direct-RAG</i>				
Llama2 _{7B}	13.60	23.46	54.40	64.49
Llama2 _{13B}	15.20	25.79	55.53	65.77
Baichuan2 _{7B}	11.67	21.32	50.40	60.11
QWen2.5 _{7B}	0.07	7.47	4.27	19.09
Mistral _{7B}	11.20	21.76	53.60	63.75
PEFT-Lora _{7B}	20.20	30.86	62.21	71.13
MoPE _{7B}	38.87	49.27	68.27	75.64
<i>Prompt-RAG</i>				
Llama2 _{7B}	29.00	40.61	60.80	69.39
Llama2 _{13B}	29.73	40.57	63.00	71.71
Baichuan2 _{7B}	28.87	38.14	57.20	65.11
QWen2.5 _{7B}	29.80	38.80	55.27	61.95
Mistral _{7B}	32.07	44.43	59.80	68.71
PEFT-Lora _{7B}	34.26	46.54	65.34	74.28
MoPE _{7B}	39.53	49.86	68.33	75.48

Table 1: Results on the NQ dataset, where inference is performed with a single expert.

contextual information. We randomly selected a small subset of sample questions for an expert insertion experiment across different layers. As show in figure 3, the results indicate that the 12th layer has the most significant impact on model performance, outperforming other layers. Although performance differences across layers are generally minor, the middle layers, particularly the 12th layer, demonstrate better performance in integrating injected knowledge compared to the initial and final layers.

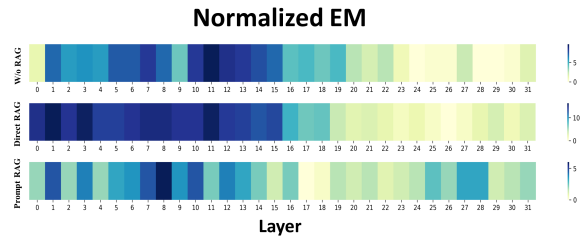


Figure 3: Normalized EM metric for expert insertion across different layers.

	Methods	1 expert		2 expert		3 expert		4 expert	
		EM	F1	EM	F1	EM	F1	EM	F1
NQ	W/o-RAG	26.40	35.40	28.20	37.15	28.93	37.41	29.07	37.65
	D-RAG	38.87	49.27	39.87	50.36	40.40	50.77	40.53	50.99
	P-RAG	39.53	49.86	40.20	50.20	40.40	50.44	40.06	50.59
TQA	W/o-RAG	56.13	64.08	57.07	64.47	58.27	65.74	57.73	65.39
	D-RAG	68.27	75.64	69.27	76.12	69.47	76.25	69.20	75.92
	P-RAG	68.33	75.48	69.53	76.08	69.07	75.69	68.47	75.38

Table 2: The results of inference with varying numbers of experts.

4.5 Main Results

Table 1 presents a comparison of MoPE with other baseline methods under the W/o-RAG, D-RAG, and P-RAG settings. As shown in the table, when expert modules are injected with paragraph knowledge, the performance of MoPE on the LLaMa2-7B model significantly surpasses that of the original model in answering questions. On the other hand, by injecting question-style knowledge, MoPE’s performance under D-RAG is comparable to the results obtained from providing prompts in P-RAG, suggesting that the method can achieve similar outcomes without explicit prompts. This opens up new possibilities for tasks such as simplifying or distilling in-context learning.

We also tested the performance after varying the number of experts involved in the query, as shown in the table 2. The results indicate that as the number of experts increases, the model’s performance, measured by both EM and F1 scores, gradually improves, particularly under the D-RAG and P-RAG settings, where the performance boost is more pronounced. This suggests that knowledge injection through expert modules effectively enhances the model’s performance, and with the increase in the number of experts, the model is able to better leverage different paragraph of knowledge. Under the W/o-RAG setting, despite the absence of external prompts, the model’s performance remains relatively stable, yet still improves as the number of experts increases, demonstrating the effectiveness of expert injection.

4.6 Efficiency Analysis of Expert Insertion

Table 3 shows the throughput comparison between the base LLaMa2-7B model and the MoPE framework. Despite the additional steps required for loading expert parameters, and performing inference, the throughput of MoPE (0.598 samples/s) remains close to that of the base model (0.628 sam-

ples/s). This indicates that the expert selection and integration process introduces minimal time overhead. The marginal reduction in throughput is compensated by a significant improvement in accuracy, demonstrating that MoPE effectively balances efficiency and performance. This makes it a practical solution for enhancing model accuracy without compromising too much on processing speed.

Method	Throughput
Base	0.628
MoPE	0.598

Table 3: Throughput comparison between the base model and MoPE.

5 Conclusion

In this study, we proposed the MoPE framework, which applies knowledge injection via MMoE to address the issue of LLMs’ insufficient attention to external prompts in RAG. Given a query, relevant paragraphs are transformed into edited knowledge that the model can learn and store in the designed expert modules. Experiments on NQ and TQA demonstrated that this expert memory-based approach mitigates the inadequate attention to external prompts, significantly enhancing performance in both W/o RAG and W/ RAG modes.

In future work, we aim to explore the concept of ‘meta-learning experts,’ where trained offline experts do not require additional storage. Instead, we propose dynamically storing expert parameters within a small-scale supernet, using it to generate expert parameters based on the query. This approach reduces storage costs while enabling more flexible expert selection and knowledge injection. Furthermore, we plan to investigate stable continual learning mechanisms for updating the supernet, allowing it to adapt to new textual data and improve the model’s long-term adaptability.

597 Limitations

598 In the current setup, our framework requires training
599 an offline expert repository, which necessitates
600 a certain amount of storage space to maintain these
601 expert modules. Future research will explore the
602 use of compact hypernetworks to dynamically generate
603 expert parameters, potentially reducing the
604 storage requirements. Additionally, more advanced
605 expert merging techniques need to be investigated
606 to better leverage complementary knowledge from
607 different experts.

608 Furthermore, the quality of the generated QA
609 pairs for editing significantly influences the performance
610 of the experts. It is crucial to avoid injecting
611 incorrect knowledge into the experts. However,
612 as shown in Appendix B, the quality of QA pairs
613 generated by T5-xl as the question generator is not
614 always satisfactory and sometimes contains errors.
615 Future work should explore or propose more advanced
616 components and methodologies for converting
617 paragraphs into correct editing facts without
618 relying on structured data. This approach could
619 also overcome limitations of existing editing methods,
620 such as MEMIT (Meng et al., 2022b), which
621 depend on fixed editing formats like triple-based
622 data.

623 References

624 Simone Alghisi, Massimo Rizzoli, Gabriel Roccabruna,
625 Seyed Mahed Mousavi, and Giuseppe Riccardi. 2024.
626 Should we fine-tune or rag? evaluating different techniques
627 to adapt llms for dialogue. *arXiv preprint arXiv:2406.06399*.
628

629 Rahaf Aljundi, Punarjay Chakravarty, and Tinne Tuytelaars. 2017. Expert gate: Lifelong learning with a
630 network of experts. In *Proceedings of the IEEE conference on computer vision and pattern recognition*,
631 pages 3366–3375.
632

633 Uri Alon, Frank Xu, Junxian He, Sudipta Sengupta, Dan Roth, and Graham Neubig. 2022. [Neuro-symbolic language modeling with automaton-augmented retrieval](#). In *Proceedings of the 39th International Conference on Machine Learning*, volume 162 of *Proceedings of Machine Learning Research*, pages 468–485. PMLR.

634 Szymon Antoniak, Sebastian Jaszczur, Michał Krutul, Maciej Pióro, Jakub Krajewski, Jan Ludziejewski, Tomasz Odrzygóźdź, and Marek Cygan. 2023. Mixture of tokens: Efficient llms through cross-example aggregation. *arXiv preprint arXiv:2310.15961*.

635 Akari Asai, Zeqiu Wu, Yizhong Wang, Avirup Sil, and Hannaneh Hajishirzi. 2023. Self-rag: Learning to

648 retrieve, generate, and critique through self-reflection.
649 *arXiv preprint arXiv:2310.11511*.

650 Siyuan Cheng, Bozhong Tian, Qingbin Liu, Xi Chen, Yongheng Wang, Huajun Chen, and Ningyu Zhang. 2023. [Can we edit multimodal large language models?](#) In *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*, pages 13877–13888, Singapore. Association for Computational Linguistics.

657 Aakanksha Chowdhery, Sharan Narang, Jacob Devlin, Maarten Bosma, Gaurav Mishra, Adam Roberts, Paul Barham, Hyung Won Chung, Charles Sutton, Sebastian Gehrmann, et al. 2023. Palm: Scaling language modeling with pathways. *Journal of Machine Learning Research*, 24(240):1–113.

663 Nicola De Cao, Wilker Aziz, and Ivan Titov. 2021. Editing factual knowledge in language models. *arXiv preprint arXiv:2104.08164*.

666 Guanting Dong, Yutao Zhu, Chenghao Zhang, Ze Chen Wang, Zhicheng Dou, and Ji-Rong Wen. 2024. Understand what llm needs: Dual preference alignment for retrieval-augmented generation. *arXiv preprint arXiv:2406.18676*.

671 Shihan Dou, Enyu Zhou, Yan Liu, Songyang Gao, Jun Zhao, Wei Shen, Yuhao Zhou, Zhiheng Xi, Xiao Wang, Xiaoran Fan, et al. 2023. Lora-moe: Revolutionizing mixture of experts for maintaining world knowledge in language model alignment. *arXiv preprint arXiv:2312.09979*, 4(7).

677 Nan Du, Yanping Huang, Andrew M Dai, Simon Tong, Dmitry Lepikhin, Yuanzhong Xu, Maxim Krikun, Yanqi Zhou, Adams Wei Yu, Orhan Firat, et al. 2022. Glam: Efficient scaling of language models with mixture-of-experts. In *International Conference on Machine Learning*, pages 5547–5569. PMLR.

683 William Fedus, Barret Zoph, and Noam Shazeer. 2022. Switch transformers: Scaling to trillion parameter models with simple and efficient sparsity. *Journal of Machine Learning Research*, 23(120):1–39.

687 Hengrui Gu, Kaixiong Zhou, Xiaotian Han, Ninghao Liu, Ruobing Wang, and Xin Wang. 2023. Pokemqa: Programmable knowledge editing for multi-hop question answering. *arXiv preprint arXiv:2312.15194*.

691 Aman Gupta, Anup Shirgaonkar, Angels de Luis Balaguer, Bruno Silva, Daniel Holstein, Dawei Li, Jennifer Marsman, Leonardo O Nunes, Mahsa Rouzbahman, Morris Sharp, et al. 2024. Rag vs fine-tuning: Pipelines, tradeoffs, and a case study on agriculture. *arXiv preprint arXiv:2401.08406*.

697 Kelvin Guu, Kenton Lee, Zora Tung, Panupong Pasupat, and Mingwei Chang. 2020. Retrieval augmented language model pre-training. In *International conference on machine learning*, pages 3929–3938. PMLR.

701	Muhammad Usman Hadi, Rizwan Qureshi, Abbas Shah,	Vladimir Karpukhin, Barlas Oguz, Sewon Min, Patrick	755
702	Muhammad Irfan, Anas Zafar, Muhammad Bilal	Lewis, Ledell Wu, Sergey Edunov, Danqi Chen, and	756
703	Shaikh, Naveed Akhtar, Jia Wu, Seyedali Mirjalili,	Wen-tau Yih. 2020. Dense passage retrieval for open-	757
704	et al. 2023. A survey on large language models:	domain question answering . In <i>Proceedings of the</i>	758
705	Applications, challenges, limitations, and practical	<i>2020 Conference on Empirical Methods in Natural</i>	759
706	usage. <i>Authorea Preprints</i> .	<i>Language Processing (EMNLP)</i> , pages 6769–6781,	760
		Online. Association for Computational Linguistics.	761
707	Sebastian Hofstätter, Jiecao Chen, Karthik Raman, and	Urvashi Khandelwal, Omer Levy, Dan Jurafsky, Luke	762
708	Hamed Zamani. 2023. Fid-light: Efficient and effective	Zettlemoyer, and Mike Lewis. 2019. Generalization	763
709	retrieval-augmented text generation. In <i>Proceed-</i>	through memorization: Nearest neighbor language	764
710	<i>ings of the 46th International ACM SIGIR Confer-</i>	models. <i>arXiv preprint arXiv:1911.00172</i> .	765
711	<i>ence on Research and Development in Information</i>		
712	<i>Retrieval</i> , pages 1437–1447.	Tom Kwiatkowski, Jennimaria Palomaki, Olivia Red-	766
713	Neil Houlsby, Andrei Giurgiu, Stanislaw Jastrzebski,	field, Michael Collins, Ankur Parikh, Chris Alberti,	767
714	Bruna Morrone, Quentin De Laroussilhe, Andrea	Danielle Epstein, Illia Polosukhin, Jacob Devlin, Ken-	768
715	Gesmundo, Mona Attariyan, and Sylvain Gelly. 2019.	ton Lee, et al. 2019. Natural questions: a benchmark	769
716	Parameter-efficient transfer learning for nlp. In <i>In-</i>	for question answering research. <i>Transactions of the</i>	770
717	<i>ternational conference on machine learning</i> , pages	<i>Association for Computational Linguistics</i> , 7:453–	771
718	2790–2799. PMLR.	466.	772
719	Edward J Hu, Yelong Shen, Phillip Wallis, Zeyuan	Dmitry Lepikhin, HyounJoong Lee, Yuanzhong Xu,	773
720	Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang,	Dehao Chen, Orhan Firat, Yanping Huang, Maxim	774
721	and Weizhu Chen. 2021. Lora: Low-rank adap-	Krikun, Noam Shazeer, and Zhifeng Chen. 2020.	775
722	tation of large language models. <i>arXiv preprint</i>	Gshard: Scaling giant models with conditional com-	776
723	<i>arXiv:2106.09685</i> .	putation and automatic sharding. <i>arXiv preprint</i>	777
		<i>arXiv:2006.16668</i> .	778
724	Zeyu Huang, Yikang Shen, Xiaofeng Zhang, Jie Zhou,	Patrick Lewis, Ethan Perez, Aleksandra Piktus, Fabio	779
725	Wenge Rong, and Zhang Xiong. Transformer-	Petroni, Vladimir Karpukhin, Naman Goyal, Hein-	780
726	patcher: One mistake worth one neuron. In <i>The</i>	rich Küttler, Mike Lewis, Wen-tau Yih, Tim Rock-	781
727	<i>Eleventh International Conference on Learning Rep-</i>	täschel, et al. 2020. Retrieval-augmented generation	782
728	<i>resentations</i> .	for knowledge-intensive nlp tasks. <i>Advances in Neu-</i>	783
		<i>ral Information Processing Systems</i> , 33:9459–9474.	784
729	Gautier Izacard and Edouard Grave. 2020. Leverag-	Xiang Lisa Li and Percy Liang. 2021. Prefix-tuning:	785
730	ing passage retrieval with generative models for	Optimizing continuous prompts for generation. <i>arXiv</i>	786
731	open domain question answering. <i>arXiv preprint</i>	<i>preprint arXiv:2101.00190</i> .	787
732	<i>arXiv:2007.01282</i> .		
733	Robert A Jacobs, Michael I Jordan, Steven J Nowlan,	Xiaopeng Li, Shasha Li, Shezheng Song, Jing Yang, Jun	788
734	and Geoffrey E Hinton. 1991. Adaptive mixtures of	Ma, and Jie Yu. 2024. Pmet: Precise model editing	789
735	local experts. <i>Neural computation</i> , 3(1):79–87.	in a transformer. In <i>Proceedings of the AAAI Con-</i>	790
		<i>ference on Artificial Intelligence</i> , volume 38, pages	791
736	Albert Q Jiang, Alexandre Sablayrolles, Arthur Men-	18564–18572.	792
737	sch, Chris Bamford, Devendra Singh Chaplot, Diego	Xi Victoria Lin, Xilun Chen, Mingda Chen, Weijia Shi,	793
738	de las Casas, Florian Bressand, Gianna Lengyel, Guil-	Maria Lomeli, Rich James, Pedro Rodriguez, Jacob	794
739	laume Lample, Lucile Saulnier, et al. 2023a. Mistral	Kahn, Gergely Szilvasy, Mike Lewis, et al. 2023.	795
740	7b. <i>arXiv preprint arXiv:2310.06825</i> .	Ra-dit: Retrieval-augmented dual instruction tuning.	796
		<i>arXiv preprint arXiv:2310.01352</i> .	797
741	Albert Q Jiang, Alexandre Sablayrolles, Antoine	Hongyin Luo, Yung-Sung Chuang, Yuan Gong, Tian-	798
742	Roux, Arthur Mensch, Blanche Savary, Chris Bam-	hua Zhang, Yoon Kim, Xixin Wu, Danny Fox, He-	799
743	ford, Devendra Singh Chaplot, Diego de las Casas,	len Meng, and James Glass. 2023. Sail: Search-	800
744	Emma Bou Hanna, Florian Bressand, et al. 2024.	augmented instruction learning. <i>arXiv preprint</i>	801
745	Mixtral of experts. <i>arXiv preprint arXiv:2401.04088</i> .	<i>arXiv:2305.15225</i> .	802
746	Zhengbao Jiang, Frank F Xu, Luyu Gao, Zhiqing	Kevin Meng, David Bau, Alex Andonian, and Yonatan	803
747	Sun, Qian Liu, Jane Dwivedi-Yu, Yiming Yang,	Belinkov. 2022a. Locating and editing factual as-	804
748	Jamie Callan, and Graham Neubig. 2023b. Ac-	sociations in gpt. <i>Advances in Neural Information</i>	805
749	active retrieval augmented generation. <i>arXiv preprint</i>	<i>Processing Systems</i> , 35:17359–17372.	806
750	<i>arXiv:2305.06983</i> .		
751	Mandar Joshi, Eunsol Choi, Daniel S Weld, and Luke	Kevin Meng, Arnab Sen Sharma, Alex Andonian,	807
752	Zettlemoyer. 2017. Triviaqa: A large scale distantly	Yonatan Belinkov, and David Bau. 2022b. Mass-	808
753	supervised challenge dataset for reading comprehen-	editing memory in a transformer. <i>arXiv preprint</i>	809
754	sion. <i>arXiv preprint arXiv:1705.03551</i> .	<i>arXiv:2210.07229</i> .	810

811	Eric Mitchell, Charles Lin, Antoine Bosselut, Chelsea Finn, and Christopher D Manning. 2021. Fast model editing at scale. <i>arXiv preprint arXiv:2110.11309</i> .	865
812		866
813		867
814	Xiaonan Nie, Xupeng Miao, Shijie Cao, Lingxiao Ma, Qibin Liu, Jilong Xue, Youshan Miao, Yi Liu, Zhi Yang, and Bin Cui. 2021. Evomoe: An evolutionary mixture-of-experts training framework via dense-to-sparse gate. <i>arXiv preprint arXiv:2112.14397</i> .	868
815		869
816		870
817		871
818		872
819	Oded Ovadia, Menachem Brief, Moshik Misha'eli, and Oren Elisha. 2023. Fine-tuning or retrieval? comparing knowledge injection in llms. <i>arXiv preprint arXiv:2312.05934</i> .	873
820		874
821		875
822		876
823	Bowen Pan, Yikang Shen, Haokun Liu, Mayank Mishra, Gaoyuan Zhang, Aude Oliva, Colin Raffel, and Rameswar Panda. 2024. Dense training, sparse inference: Rethinking training of mixture-of-experts language models. <i>arXiv preprint arXiv:2404.05567</i> .	877
824		878
825		879
826		880
827		881
828	P Rajpurkar. 2016. Squad: 100,000+ questions for machine comprehension of text. <i>arXiv preprint arXiv:1606.05250</i> .	882
829		883
830		884
831	Carl Rasmussen and Zoubin Ghahramani. 2001. Infinite mixtures of gaussian process experts. <i>Advances in neural information processing systems</i> , 14.	885
832		886
833		887
834	Noam Shazeer, Azalia Mirhoseini, Krzysztof Maziarczyk, Andy Davis, Quoc Le, Geoffrey Hinton, and Jeff Dean. 2017. Outrageously large neural networks: The sparsely-gated mixture-of-experts layer. <i>arXiv preprint arXiv:1701.06538</i> .	888
835		889
836		890
837		891
838		892
839	Freda Shi, Xinyun Chen, Kanishka Misra, Nathan Scales, David Dohan, Ed H Chi, Nathanael Schärli, and Denny Zhou. 2023. Large language models can be easily distracted by irrelevant context. In <i>International Conference on Machine Learning</i> , pages 31210–31227. PMLR.	893
840		894
841		895
842		896
843		897
844		898
845	Kihyuk Sohn. 2016. Improved deep metric learning with multi-class n-pair loss objective. In <i>Advances in Neural Information Processing Systems</i> , volume 29. Curran Associates, Inc.	899
846		900
847		901
848		902
849	Chenmian Tan, Ge Zhang, and Jie Fu. 2023. Massive editing for large language models via meta learning. <i>arXiv preprint arXiv:2311.04661</i> .	903
850		904
851		905
852	Qwen Team. 2024. Qwen2.5: A party of foundation models.	906
853		907
854	Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Nikolay Bashlykov, Soumya Batra, Prajjwal Bhargava, Shruti Bhosale, et al. 2023. Llama 2: Open foundation and fine-tuned chat models. <i>arXiv preprint arXiv:2307.09288</i> .	908
855		909
856		910
857		911
858		912
859		913
860	Mojtaba Valipour, Mehdi Rezagholizadeh, Ivan Kobzyev, and Ali Ghodsi. 2022. Dylora: Parameter efficient tuning of pre-trained models using dynamic search-free low-rank adaptation. <i>arXiv preprint arXiv:2210.07558</i> .	914
861		915
862		916
863		917
864		918
		919
	Lean Wang, Lei Li, Damai Dai, Deli Chen, Hao Zhou, Fandong Meng, Jie Zhou, and Xu Sun. 2023. Label words are anchors: An information flow perspective for understanding in-context learning. In <i>Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing</i> , pages 9840–9855, Singapore. Association for Computational Linguistics.	
	Renzhi Wang and Piji Li. 2024a. LEMoE: Advanced mixture of experts adaptor for lifelong model editing of large language models. In <i>Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing</i> , pages 2551–2575, Miami, Florida, USA. Association for Computational Linguistics.	
	Renzhi Wang and Piji Li. 2024b. Memoe: Enhancing model editing with mixture of experts adaptors. <i>arXiv preprint arXiv:2405.19086</i> .	
	Xun Wu, Shaohan Huang, and Furu Wei. 2024. Mixture of lora experts. <i>arXiv preprint arXiv:2404.13628</i> .	
	Shitao Xiao, Zheng Liu, Peitian Zhang, Niklas Muenighoff, Defu Lian, and Jian-Yun Nie. 2024. C-pack: Packed resources for general chinese embeddings. In <i>Proceedings of the 47th international ACM SIGIR conference on research and development in information retrieval</i> , pages 641–649.	
	Shi-Qi Yan, Jia-Chen Gu, Yun Zhu, and Zhen-Hua Ling. 2024. Corrective retrieval augmented generation. <i>arXiv preprint arXiv:2401.15884</i> .	
	Aiyuan Yang, Bin Xiao, Bingning Wang, Borong Zhang, Ce Bian, Chao Yin, Chenxu Lv, Da Pan, Dian Wang, Dong Yan, et al. 2023. Baichuan 2: Open large-scale language models. <i>arXiv preprint arXiv:2309.10305</i> .	
	Lang Yu, Qin Chen, Jie Zhou, and Liang He. 2024. Melo: Enhancing model editing with neuron-indexed dynamic lora. In <i>Proceedings of the AAAI Conference on Artificial Intelligence</i> , volume 38, pages 19449–19457.	
	Wenhao Yu, Chenguang Zhu, Zaitang Li, Zhiting Hu, Qingyun Wang, Heng Ji, and Meng Jiang. 2022. A survey of knowledge-enhanced text generation. <i>ACM Computing Surveys</i> , 54(11s):1–38.	
	Wayne Xin Zhao, Kun Zhou, Junyi Li, Tianyi Tang, Xiaolei Wang, Yupeng Hou, Yingqian Min, Beichen Zhang, Junjie Zhang, Zican Dong, et al. 2023. A survey of large language models. <i>arXiv preprint arXiv:2303.18223</i> .	
	Ce Zheng, Lei Li, Qingxiu Dong, Yuxuan Fan, Zhiyong Wu, Jingjing Xu, and Baobao Chang. 2023. Can we edit factual knowledge by in-context learning? <i>arXiv preprint arXiv:2305.12740</i> .	
	Zexuan Zhong, Zhengxuan Wu, Christopher D Manning, Christopher Potts, and Danqi Chen. 2023. Mquake: Assessing knowledge editing in language models via multi-hop questions. <i>arXiv preprint arXiv:2305.14795</i> .	

A RAG Template

In this section, we present the input formats for *without-RAG*, *Direct-RAG*, and *Prompt-RAG*.

Input Format of Without-RAG
Question: { <i>question</i> }
Answer:

(a) Input Format of Without-RAG.

Input Format of Direct-RAG
Knowledge: { <i>Top-1 paragraph</i> }
Question: { <i>question</i> }
Answer:

(b) Input Format of Direct-RAG.

Input Format of Prompt-RAG
Knowledge: { <i>Top-1 paragraph</i> }
Base above knowledge, answer the following question with a very short phrase, such as "1998", "May 16th, 1931", or "James Bond", to meet the criteria of exact match datasets.
Question: { <i>question</i> }
Answer:

(c) Input Format of Prompt-RAG

Table 4: The input format of the LLM.

B Synthetic Q-A pairs

The process consists of the following steps:

1. We leverage the `spacy`⁴ library to perform Named Entity Recognition (NER) on the passage p , obtaining a set of entities:

$$A = \text{spacy}(p) \quad (15)$$

2. For each entity $a \in A$, we use the trained question generator G_θ to generate a corresponding question.

$$q_{\text{edited}} = G_\theta(q|\mathcal{T}(p, a)) \quad (16)$$

⁴<https://spacy.io/>

Table 6 presents an example of synthetic QA generation using G_θ , where the passage is the retrieved paragraph, followed by several generated QA pairs. As observed, the quality of the generated pairs is not perfect, with irrelevant or incorrect questions, such as "Where does the story *Don Quixote* take place?" with the answer "Western," which is factually incorrect, or "Sancho Panza was the steed of which fictional character?" with the answer "Don Quixote's", which is a misleading formulation. These inaccuracies and mismatches between questions and answers may significantly impact the overall system performance. This highlights a key factor affecting the results, making it evident that generating high-quality and factually accurate edited knowledge is crucial. Therefore, we are currently exploring methods to improve the accuracy and relevance of these generated QA pairs, with a focus on more reliable techniques for editing and validating the knowledge. The need to generate high-quality edited knowledge is urgent to ensure that the knowledge injection process leads to better, more precise performance in downstream tasks.

C Re-ranker Results

This section presents the performance of applying a re-ranker to refine the retrieval results of DPR (Karpukhin et al., 2020), comparing the accuracy before and after re-ranking. In a question-answering system, the accuracy of the retrieval source plays a crucial role in the overall result. Additionally, the scoring produced by the re-ranker is important for merging between experts. The improvements were evaluated using Top-K accuracy, which measures the presence of a golden passage, ensuring a stable document source for responding answer.

Top-K	Type	NQ	TQA
@1	Before	44.60	56.53
	After	64.67	76.33
@2	Before	55.73	65.27
	After	71.73	79.60
@4	Before	64.47	72.07
	After	77.87	82.80
@8	Before	72.93	76.73
	After	81.47	84.53

Table 5: Comparison of DPR retrieval accuracy results before and after applying the re-ranker we trained.

Passage: "forced to deceive him at certain points. The novel is considered a satire of orthodoxy, veracity and even nationalism. In exploring the individualism of his characters, Cervantes helped move beyond the narrow literary conventions of the chivalric romance literature that he spoofed, which consists of straightforward retelling of a series of acts that redound to the knightly virtues of the hero. The character of Don Quixote became so well known in its time that the word 'quixotic' was quickly adopted by many languages. Characters such as Sancho Panza and Don Quixote's steed, Rocinante, are emblems of Western literary culture."

"question": "Who wrote the novel don quixote and the adventures of sancho panza?",
"answer": "Cervantes"

"question": "Sancho Panza is a character in which novel?"
"answer": "Don Quixote"

"question": "What was the name of Don Quixote's faithful companion?"
"answer": "Sancho Panza"

"question": "Sancho Panza was the steed of which fictional character?"
"answer": "Don Quixote's"

"question": "What was the name of Don Quixote's horse?"
"answer": "Rocinante"

"question": "Where does the story don quixote take place?"
"answer": "Western"

Table 6: An example of synthetic QA pairs.