

Lifelong Knowledge Editing for LLMs with Retrieval-Augmented Continuous Prompt Learning

Qizhou Chen^{1,2*}, Taolin Zhang^{2*}, Xiaofeng He^{1,3}, Dongyang Li^{1,2},
Chengyu Wang^{4†}, Longtao Huang², Hui Xue²

¹ East China Normal University, Shanghai, China ² Alibaba Group, Hangzhou, China

³NPPA Key Laboratory of Publishing Integration Development, ECNUP, Shanghai, China

⁴ Alibaba Cloud Computing, Hangzhou, China

chen_qizhou@outlook.com, zhangt10519@gmail.com, chengyu.wcy@alibaba-inc.com

Abstract

Model editing aims to correct outdated or erroneous knowledge in large language models (LLMs) without the need for costly retraining. Lifelong model editing is the most challenging task that caters to the continuous editing requirements of LLMs. Prior works primarily focus on single or batch editing; nevertheless, these methods fall short in lifelong editing scenarios due to catastrophic knowledge forgetting and the degradation of model performance. Although retrieval-based methods alleviate these issues, they are impeded by slow and cumbersome processes of integrating the retrieved knowledge into the model. In this work, we introduce RECIPE, a RetriEval-augmented Continuous Prompt LEarning method, to boost editing efficacy and inference efficiency in lifelong learning. RECIPE first converts knowledge statements into short and informative continuous prompts, prefixed to the LLM’s input query embedding, to efficiently refine the response grounded on the knowledge. It further integrates the Knowledge Sentinel (KS) that acts as an intermediary to calculate a dynamic threshold, determining whether the retrieval repository contains relevant knowledge. Our retriever and prompt encoder are jointly trained to achieve editing properties, i.e., reliability, generality, and locality. In our experiments, RECIPE is assessed extensively across multiple LLMs and editing datasets, where it achieves superior editing performance. RECIPE also demonstrates its capability to maintain the overall performance of LLMs alongside showcasing fast editing and inference speed.¹

1 Introduction

Large language models (LLMs) (Touvron et al., 2023; Roumeliotis and Tselikas, 2023; Zeng et al.,

2023) have become key techniques in NLP. However, once trained, the knowledge encapsulated within LLMs becomes static (Petroni et al., 2019). This can lead to outputs that are outdated or even erroneous as time progresses (Yao et al., 2023). In response, model editing techniques have been developed (Meng et al., 2022, 2023; Hartvigsen et al., 2022; Tan et al., 2023; Hu et al., 2024; Jiang et al., 2024), aimed at efficiently updating and correcting LLMs without the necessity of retraining with large-scale parameters. This concept is economically advantageous as it reduces computational costs and enhances the accuracy of outputs produced by LLMs (Cao et al., 2021; Mitchell et al., 2022; Meng et al., 2023; Mishra et al., 2024).

Previous efforts in model editing have primarily focused on single and batch edits. Notable examples include ROME (Meng et al., 2022), MEND (Mitchell et al., 2022), and MEMIT (Meng et al., 2023), which achieve edits by applying offsets to part of the model’s parameters. However, in the real world, LLMs frequently require continuous knowledge updates to stay abreast of emerging knowledge. Thus, the concept of lifelong editing has been introduced (Hartvigsen et al., 2022). As shown in the upper part of Figure 1, with continuous editing, the accumulating offsets on parameters can result in model performance degradation or even failure (Hartvigsen et al., 2022; Huang et al., 2023; Han et al., 2023; Hu et al., 2024). Some techniques (Dong et al., 2022; Huang et al., 2023) address the challenges by integrating extra model parameters. Nevertheless, as shown in the middle of Figure 1, the increase in additional parameters leads to diminished model performance and reduced inference efficiency.

Retrieval-based methods separate knowledge from the model. As shown in the lower part of Figure 1, they temporarily integrate retrieved knowledge into the model during each inference, mitigating knowledge forgetting and performance degra-

* Q. Chen and T. Zhang contributed equally to this work.

† Corresponding author.

¹Source codes is available at <https://github.com/qizhou000/RECIPE>.

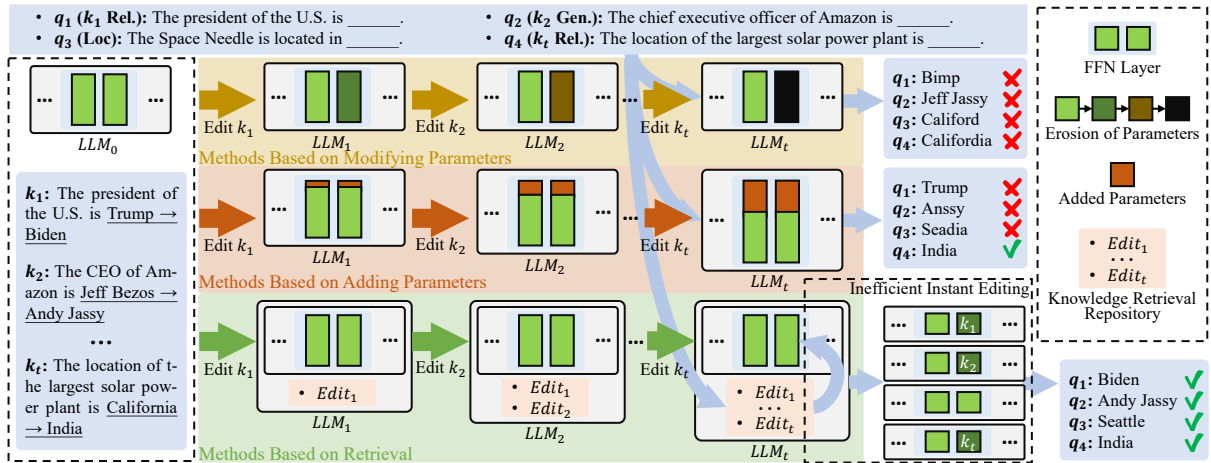


Figure 1: Comparison among three types of methods in lifelong editing scenarios. Modifying parameters and adding extra parameters result in the degradation of LLM performance as editing progresses. In contrast, retrieval-based editors store knowledge in a repository and apply knowledge editing on the fly, which maintains the LLM unchanged and relieves it from accumulating parameter offsets or adding extra parameters. (Best viewed in color)

dition. However, these methods face certain limitations: retrieval at each token prediction (Hartvigsen et al., 2022), caching large update matrices (Han et al., 2023; Yu et al., 2024), and overly long editing prefixes (Jiang et al., 2024). These shortcomings affect the model’s inference efficiency.

In this paper, we introduce RECIPE, a novel *RetriEval-augmented Continuous Prompt Learning* framework to enhance editing efficacy and inference efficiency for LLMs in lifelong learning scenarios. Two key techniques of RECIPE are introduced as follows:

Knowledgeable Continuous Prompt Learning:

In retrieval-based methods, LTE (Jiang et al., 2024) effectively avoids the shortcomings of earlier approaches. However, their overly long editing prefixes can reduce model inference speed, and full-parameter fine-tuning also increases the risk of overfitting. Therefore, we aim to explore the shortest possible prefixes to enable the LLM to follow editing instructions. In RECIPE, we achieve this by constructing a continuous short prompt encoder that transforms each piece of editing knowledge (expressed as text) into a knowledgeable continuous prompt. This approach is grounded in prior research, as exemplified by P-tuning (Li and Liang, 2021; Liu et al., 2022), which demonstrated that continuous prompts enable LLMs to perform downstream tasks more effectively. Here, we conceptualize each knowledge edit as a distinct mini-task. To ensure editing efficacy, our prompt encoder is trained to align with three key editing properties including reliability, generality, and locality (Yao

et al., 2023).

Dynamic Prompt Retrieval with Knowledge Sentinel:

To avoid unnecessary costs from introducing additional retrievers (Han et al., 2023; Jiang et al., 2024), we map knowledge statements and queries into the same representational space through the prompt encoder to compute retrieval similarity. To determine whether the retrieval repository contains knowledge related to an input query, manually setting a fixed similarity threshold is a common practice (Han et al., 2023). However, this method does not take into account that different queries often require distinct thresholds due to semantic variations. Therefore, we introduce the Knowledge Sentinel (KS), a trainable embedding representation, as an intermediary to dynamically compute the threshold for each query. Employing a specifically designed contrastive learning mechanism, the KS module is jointly trained with the prompt encoder to align retrieval with model editing.

The contributions of our paper are summarized as follows:

- We propose a novel RetriEval-augmented Continuous Prompt Learning framework, RECIPE, to facilitate lifelong model editing. As far as we know, we are the first to adapt prompt learning for model editing, and explore the shortest possible editing prefixes.
- Within RECIPE, the prompt learning encoder effectively shortens the editing prefixes, enhancing the inference efficiency of the post-edit LLM. The KS improves editing retrieval,

thus enhancing overall editing performance.

- We conduct comparative experiments of life-long editing across multiple backbones and editing datasets. The results demonstrate the superiority of RECIPE.

2 Related Works

2.1 Model Editing

We categorize model editing methods into three types: modifying parameters, adding extra parameters, and retrieval-based methods.

Methods modifying model parameters can be further divided into Locate-then-Edit (L&E) and meta-learning-based methods. For L&E, ROME (Meng et al., 2022) identifies the LLMs’ edit-sensitive layers through causal tracing and proposes rank-one model editing to modify parameters. MEMIT (Meng et al., 2023) and WILKE (Hu et al., 2024) respectively use multi-layer allocation and dynamic localization to alleviate the single matrix update burden of ROME. In meta-learning-based methods, KnowledgeEditor (Cao et al., 2021) and MEND (Mitchell et al., 2022) respectively transform editing knowledge and the gradient decomposition of LLM to the offsets of the weights to be edited. MALMEN (Tan et al., 2023) and DAFNet (Zhang et al., 2024) enhance MEND’s parameter fusion in multiple editing by utilizing normal equations and constructing interactions between edits, respectively. Although these methods show success in single or batch editing scenarios, in a lifelong editing situation, as the number of edits increases, the accumulating mismatches of parameter offsets can lead to model degradation or failure (Hu et al., 2024).

Methods adding extra parameters, such as CaLiNet (Dong et al., 2022) and T-Patcher (Huang et al., 2023), achieve model editing by introducing additional neurons to the LLM for each piece of editing knowledge, thereby avoiding modifications to the original model parameters. However, in the lifelong editing scenario, the continuous addition of neurons can progressively dominate the LLM’s inference process. This can lead to a reduction in inference speed and model capability.

Retrieval-based editors effectively circumvent the issue of accumulated parameter offsets and the potentially unbounded addition of neurons. GRACE (Hartvigsen et al., 2022) performs editing retrieval for each token prediction by calculating

the linear distance between representations, which impacts the efficiency of long sequence predictions. RASE (Han et al., 2023) develops an editing retrieval model to improve sequential editing. MELO (Yu et al., 2024) introduces a batch editing version of GRACE using LoRA (Hu et al., 2022). Both methods’ instant editing schemes require updating the corresponding FFN matrices separately for each sample in the input batch, which increases memory load as the batch size grows. LTE (Jiang et al., 2024) fine-tunes the LLM to respond to knowledge when prefixed with editing information and retrieves relevant content using the off-the-shelf backbone (Reimers and Gurevych, 2019). However, lengthy editing prefixes still impact inference efficiency. Additionally, fine-tuning the model itself is likely to lead to overfitting, thereby affecting its original performance. To address the above issues, I propose RECIPE to explore more efficient retrieval and instant editing, as well as smaller interventions to the model to avoid overfitting on the editing dataset.

2.2 Prompt Tuning

Prompt tuning is a typical parameter-efficient learning method that only requires updating a relatively small number of parameters. There are two types of prompt tuning methods: discrete and continuous. The discrete methods (Gao et al., 2021; Levy et al., 2023; Wang et al., 2023b; Duan et al., 2023) guide the model to generate relevant outputs for specific tasks by designing fixed, text-based prompts. Continuous methods (Li and Liang, 2021; Liu et al., 2022, 2021a,b; Mu et al., 2023; Xu et al., 2023; Zhang et al., 2023), more relevant to RECIPE, utilize trainable word embedding vectors as prompts. Building on the foundations of these works, our approach is duly justified, encoding individual pieces of knowledge as continuous prompts.

3 Background

In this section, we first formally present the model editing task and its lifelong version. Then, we introduce the evaluation properties in model editing.

An LLM $f_{llm} \in \mathcal{F}$ can be regarded as a function $f_{llm} : \mathcal{Q} \mapsto \mathcal{A}$ that maps an input query q to its predicted answer $a = f_{llm}(q)$. Given an edit example pair (q_e, a_e) that $f_{llm}(q_e) \neq a_e$, a model editor $\mathbf{ME} : \mathcal{F} \times \mathcal{Q} \times \mathcal{A} \mapsto \mathcal{F}$ outputs a post-edit model f'_{llm} such that:

$$f'_{llm} = \mathbf{ME}(f_{llm}, q_e, a_e) \quad (1)$$

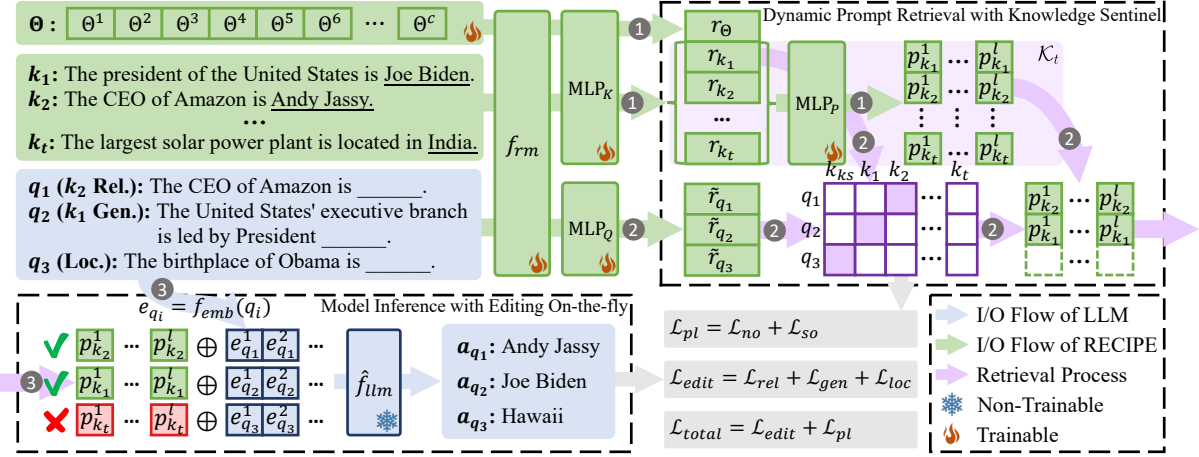


Figure 2: Illustration of the RECIPE framework. Process 1 constructs and updates the knowledge retrieval repository \mathcal{K}_t . During the inference stage, Process 2 retrieves query-related prompts from \mathcal{K}_t . Process 3 utilizes the retrieved continuous prompts to correct the LLM’s response. For lifelong editing, the repository can be continuously updated (e.g., from \mathcal{K}_{t-1} to \mathcal{K}_t) with each new insertion of knowledge and prompts.

Given an initial model f_{llm}^0 , ME will iteratively implement editing as the demands of editing continue to emerge in a lifelong editing scenario:

$$f_{llm}^t = \text{ME}(f_{llm}^{t-1}, q_{e_t}, a_{e_t}), t = 1, 2, 3, \dots \quad (2)$$

At any timestep t in the lifelong editing process, a good ME should make the edited model f_{llm}^t meet the following three criteria (Yao et al., 2023):

Reliability requires f_{llm}^t to correctly remember all the previously edit samples themselves:

$$\mathbb{E}_{(q_e, a_e) \sim \{(q_{e\tau}, a_{e\tau})\}_{\tau=1}^t} \mathbb{I} \{ f_{llm}^t(q_e) = a_e \} \quad (3)$$

where the \mathbb{I} is the indicator function.

Generality requires f_{llm}^t to correctly answer queries belonging to relevant neighbors of previously edited samples:

$$\mathbb{E}_{(q_e, a_e) \sim \{(q_{e\tau}, a_{e\tau})\}_{\tau=1}^t} \mathbb{E}_{(q_g, a_g) \sim N(q_e, a_e)} \mathbb{I}_g(q_g, a_g) \quad (4)$$

s.t. $\mathbb{I}_g(q_g, a_g) = \mathbb{I} \{ f_{llm}^t(q_g) = a_g \}$

where $N(q_e, a_e)$ is the relevant neighbors of edit sample (q_e, a_e) .

Locality requires f_{llm}^t to maintain consistency with the initial model f_{llm}^0 on queries unrelated to previously edited samples:

$$\mathbb{E}_{(q_e, a_e) \sim \{(q_{e\tau}, a_{e\tau})\}_{\tau=1}^t} \mathbb{E}_{(q_l, a_l) \sim O(q_e, a_e)} \mathbb{I}_l(q_l, a_l) \quad (5)$$

s.t. $\mathbb{I}_l(q_l, a_l) = \mathbb{I} \{ f_{llm}^t(q_l) = f_{llm}^0(q_l) \}$

where $O(q_e, a_e)$ is the irrelevant samples set w.r.t. the edit sample (q_e, a_e) . Note that the locality metric implicitly includes the preservation of the general performance of f_{llm}^t relative to f_{llm}^0 .

4 The Proposed Approach

In this section, we formally introduce the RECIPE framework, with the overall architecture in Figure 2. First, RECIPE maintains a knowledge retrieval repository, which stores representations of editing knowledge mapped to their knowledgeable continuous prompts described in Sec. 4.1. Next, we introduce a dynamic retrieval technique with the KS to facilitate knowledge retrieval to filter out irrelevant knowledge in Sec. 4.2. To ensure the LLMs adhere to the edited knowledge related to the query efficiently, RECIPE prefixes the retrieved continuous prompt to the word embeddings of the LLM’s input query, as detailed in Sec. 4.3. Finally, we describe the joint training procedure of the RECIPE framework in Sec. 4.4. The algorithms for RECIPE are detailed in Appendix A.

4.1 Construction and Update of Knowledge Retrieval Repository

The knowledge retrieval repository is initialized as empty, i.e., $\mathcal{K}_0 = \{\}$, and is updated from \mathcal{K}_{t-1} to \mathcal{K}_t by adding a new key-value pair corresponding to new editing knowledge, k_t , at each timestep t in our lifelong editing setting.

Specifically, at timestep t , given a new knowledge statement k_t , the knowledge representation $r_{k_t} \in \mathcal{R}^{d_r}$ is achieved through an encoder f_{rm} (e.g., RoBERTa (Liu et al., 2019)) stacked with a multilayer perceptron (MLP) MLP_K :

$$r_{k_t} = \text{MLP}_K(f_{rm}(k_t)) \quad (6)$$

where f_{rm} concatenates the maximum, minimum,

and average pooling of its output token representations (including the [CLS] token) into a vector to maximally retain the semantic information of the input. Then, the continuous prompt $p_{k_t} \in \mathbb{R}^{l \times d_{llm}}$ is generated through another MLP, i.e., MLP_P :

$$p_{k_t} = f_{resp}(\text{MLP}_P(r_{k_t})) \quad (7)$$

where l and d_{llm} are the length of the continuous prompt and the dimension of the LLM’s word embedding, respectively. In other words, l is the number of Continuous Prompt Tokens (CPTs) leveraged for LLM inference. f_{resp} is the reshape operation that maps the vector into a matrix with shape $l \times d_{llm}$. Finally, the knowledge retrieval repository is updated from \mathcal{K}_{t-1} to \mathcal{K}_t : $\mathcal{K}_t = \mathcal{K}_{t-1} \cup \{(r_{k_t}, p_{k_t})\}$ where (r_{k_t}, p_{k_t}) is the key-value pair for knowledge retrieval.

4.2 Dynamic Prompt Retrieval with Knowledge Sentinel

The existence of a query-related prompt in the repository is usually determined by using a manually set similarity threshold (Han et al., 2023). However, using a fixed threshold does not account for the fact that the sensitivity to similarity with related knowledge varies among different queries due to semantic differences. The Knowledge Sentinel (KS) serves as an intermediary leveraged to dynamically compute similarity thresholds for various queries. To be specific, KS $\Theta \in \mathcal{R}$ is a trainable word embedding of f_{rm} with token length c . It is transformed into the knowledge representation space as: $r_\Theta = \text{MLP}_K(f_{rm}(\Theta))$. Given a query q and the knowledge retrieval repository $\mathcal{K}_t = \{(r_{k_\tau}, p_{k_\tau})\}_{\tau=1}^t$, the prompt retrieval process is as follows:

$$\tilde{r}_q = \text{MLP}_Q(f_{rm}(q)) \quad (8)$$

$$\text{KS}(q) = \begin{cases} p_{k_j} & \tilde{r}_q^T \cdot r_{k_j} > \tilde{r}_q^T \cdot r_\Theta \\ \emptyset & \text{otherwise} \end{cases} \quad (9)$$

where $j = \text{argmax}_{\tau=1, \dots, t} \tilde{r}_q^T \cdot r_{k_\tau}$, which can be efficiently searched via modern vector databases or search engines (e.g., Chen et al. (2021)). MLP_Q is the MLP that maps the query representation to the knowledge representation space. If the retrieved continuous prompt is not sufficiently similar to the query compared to KS, an empty set is returned. Hence, the inference of LLMs is not modified.

4.3 Model Inference with Editing On-the-fly

Previous retrieval-based methods suffer from cumbersome editing processes and post-retrieval knowledge integration (Hartvigsen et al., 2022; Jiang et al., 2024). To address this challenge, we prefix the retrieved continuous prompt to the word embedding of the input query to efficiently correct the response of the LLM.

Specifically, we consider the LLM to be edited as $f_{llm} : \mathcal{Q} \mapsto \mathcal{A}$, where \hat{f}_{llm} is f_{llm} with the embedding layer f_{emb} removed. Given an input query q , and the retrieved continuous prompt $p_{k_\tau} = \text{KS}(q)$, the inference process is reformulated as: $a_q = \hat{f}_{llm}(p_{k_\tau} \oplus f_{emb}(q))$ where \oplus denotes the concatenation of the retrieved continuous prompt matrix and the word embedding matrix of q .

The feasibility of our approach is supported by previous work such as P-Tuning (Li and Liang, 2021; Liu et al., 2022), which demonstrates the efficacy of training continuous prompt embeddings to enhance the performance of LLMs on downstream tasks. In RECIPE, we treat the editing of each knowledge statement as a mini-task. Instead of fine-tuning a specific prompt encoder for each mini-task, we accomplish the objectives of these mini-tasks by training RECIPE modules that generate continuous prompts, ensuring the LLM adheres to the corresponding knowledge.

4.4 Model Training

The losses are formulated to ensure adherence to the editing of generated continuous prompts and effective retrieval of query-related knowledge for the LLM. Given a batch of training data consisting of b editing sample pairs $\{(q_{e_i}, a_{e_i})\}_{i=1}^b$ and their corresponding sampled generality and locality pairs $\{(q_{g_i}, a_{g_i})\}_{i=1}^b, \{(q_{l_i}, a_{l_i})\}_{i=1}^b$, the losses are formulated as follows.

Editing: The editing loss aims to ensure that the generated continuous prompt guides the LLM to follow the properties of reliability, generality, and locality (Yao et al., 2023). Based on the pairs (q_{e_i}, a_{e_i}) , the sample-wise losses corresponding to these three properties are defined as follows:

$$\mathcal{L}_{rel}^{(i)} = -\log \hat{f}_{llm}(a_{e_i} | p_{k_i} \oplus f_{emb}(q_{e_i})) \quad (10)$$

$$\mathcal{L}_{gen}^{(i)} = -\log \hat{f}_{llm}(a_{g_i} | p_{k_i} \oplus f_{emb}(q_{g_i})) \quad (11)$$

$$\mathcal{L}_{loc}^{(i)} = \text{KL} \left(f_{llm}(q_{l_i}) || \hat{f}_{llm}(p_{k_i} \oplus f_{emb}(q_{l_i})) \right) \quad (12)$$

where p_{k_i} is the continuous prompt transformed

through Eq. 6 and Eq. 7 using knowledge k_i that is the concatenation of q_{e_i} and a_{e_i} . The KL denotes the Kullback-Leibler divergence, which is chosen to better fit the LLM’s original prediction distribution on the locality data. The batch-wise loss function for model editing is derived as follows:

$$\mathcal{L}_{edit} = \frac{1}{b} \sum_{i=1}^b \left(\mathcal{L}_{rel}^{(i)} + \mathcal{L}_{gen}^{(i)} + \mathcal{L}_{loc}^{(i)} \right). \quad (13)$$

Prompt Learning: The training losses for prompt learning are based on contrastive learning (van den Oord et al., 2018; He et al., 2020) and are aligned with the properties of reliability, generality, and locality (Yao et al., 2023). For a batch of samples, the loss functions for learning continuous prompts are formulated as follows:

$$\mathcal{L}_{no}^{(i)} = \delta(\tilde{r}_{q_{e_i}}, r_{k_i}, R) + \delta(\tilde{r}_{q_{g_i}}, r_{k_i}, R), \quad (14)$$

$$\begin{aligned} \mathcal{L}_{so}^{(i)} = & \delta(\tilde{r}_{q_{l_i}}, r_{\Theta}, R) + \delta(\tilde{r}_{q_{e_i}}, r_{\Theta}, R \setminus k_i) \\ & + \delta(\tilde{r}_{q_{g_i}}, r_{\Theta}, R \setminus k_i), \end{aligned} \quad (15)$$

$$\mathcal{L}_{pl} = \frac{1}{b} \sum_{i=1}^b (\mathcal{L}_{no}^{(i)} + \mathcal{L}_{so}^{(i)}), \quad (16)$$

where $R = \{r_{k_i}\}_{i=1}^b \cup \{r_{\Theta}\}$ and $R \setminus k_i = R \setminus \{r_{k_i}\}$. r_{k_i} is the representation of the editing knowledge k_i transformed through Eq. 6. The query representations $\tilde{r}_{q_{e_i}}, \tilde{r}_{q_{g_i}}, \tilde{r}_{q_{l_i}}$ for $q_{e_i}, q_{g_i}, q_{l_i}$ are attained via Eq. 8, respectively. δ is the InfoNCE loss (van den Oord et al., 2018), formulated as:

$$\delta(q, k_+, \{k_i\}_{i=1}^n) = -\log \frac{\exp(q \cdot k_+ / \tau)}{\sum_{i=1}^n \exp(q \cdot k_i / \tau)}, \quad (17)$$

where τ is the temperature, typically set to 1 by default. In our work, the *neighbor-oriented* loss $\mathcal{L}_{no}^{(i)}$ encourages higher similarity between the editing knowledge and the corresponding reliability or generality queries. The *sentinel-oriented* loss $\mathcal{L}_{so}^{(i)}$ ensures that input queries yield the highest similarity with the KS in cases where the retrieval repository lacks relevant knowledge.

Thus, the total training loss is: $\mathcal{L}_{total} = \mathcal{L}_{edit} + \mathcal{L}_{pl}$. During training, the parameters of the LLM llm are kept frozen. The trainable modules include only f_{rm} , MLP_K , MLP_Q , MLP_P , and Θ , which renders our approach highly lightweight.

5 Experiments

In this section, we present the experimental results on LLMs including LLAMA-2 (7B), GPT-J

(6B), and GPT-XL (1.5B). We compare RECIPE against strong baselines including Fine-Tune (FT) (Yao et al., 2023), MEND (Mitchell et al., 2022), ROME (Meng et al., 2022), MEMIT (Meng et al., 2023), MALMEN (Tan et al., 2023), WILKE (Hu et al., 2024), T-Patcher (TP) (Huang et al., 2023), GRACE (Hartvigsen et al., 2022), R-ROME (Han et al., 2023), and LTE (Jiang et al., 2024). For the evaluation of lifelong editing, we use datasets including ZSRE (Mitchell et al., 2022), CF (Meng et al., 2022), and RIPE (Cohen et al., 2023). For the evaluation of general performance after editing, we apply datasets including CSQA (Talmor et al., 2019), ANLI (Nie et al., 2020), MMLU (Hendrycks et al., 2021), and SQuAD-2 (Rajpurkar et al., 2018). For more detailed description of datasets, baselines, and model settings, please refer to Appendix B and Appendix C.

5.1 The Performance of RECIPE

Editing Performance: Table 1 and Appendix D present the overall editing performance across various numbers of edits to simulate a lifelong editing scenario. From the single-edit perspective, our method exhibits optimal performance in most testing scenarios. In the lifelong editing scenarios, we have the following observations: (1) Methods that modify the parameters of LLMs show outstanding editing performance in a single edit. Yet, they exhibit a significant decline in editing performance as the number of edits increases. This trend aligns with the toxic accumulation issue highlighted by Hu et al. (2024). (2) Method introducing additional parameters maintains a degree of reliability and generality in the lifelong editing process. However, the cumulative addition of extra parameters compromises the original inference process, evidenced by the pronounced deterioration in locality observed in ZSRE. (3) Retrieval-based approaches demonstrate robustness against the increasing number of edits. Among them, our method achieves the best results, affirming the strengths of retrieval as well as validating the efficacy of our strategy.

From the perspective of metrics, for reliability and generality, the co-occurrence of correct retrieval and the model’s faithful adherence to editing instructions is a prerequisite for editing success. Therefore, compared to other retrieval-based methods, we consider the additional introduction of KS in RECIPE, which enhances retrieval performance and thereby improves editing efficacy, to be a crucial factor for RECIPE’s superior perfor-

Editor	CSQA	MMLU	ANLI	SQUAD-2	Average
N/A	38.91	41.54	34.04	36.43	37.73
FT	19.27	29.93	33.33	0.59	20.78
MEND	20.31	24.68	33.07	0.04	19.52
ROME	19.97	23.03	33.47	0.41	19.22
MEMIT	19.68	23.23	33.39	0.01	19.08
TP	19.62	22.84	33.37	0.96	19.20
GRACE	38.60	41.20	33.93	36.28	37.50
R-ROME	38.50	41.12	33.90	36.31	37.46
MALMEN	20.85	24.83	33.03	0.27	19.75
LTE	19.45	23.21	33.41	25.25	25.33
WILKE	19.87	23.37	33.37	0.07	19.17
RECIPE	38.76	41.40	34.13	36.50	37.70

Table 2: Performance of LLAMA-2 after 1,000 edits. “N/A” denotes performance without any edits. Bold font highlights the optimal post-editing performance.

Type	Editor	Edit Time	Infer. Time	Total Time
MP	FT	1.7205	0.0589	1.7794
	MEND	0.0987	0.0590	0.1577
	ROME	17.1639	0.0586	17.2225
	MEMIT	33.6631	0.0591	33.7222
	MALMEN	2.3418	0.0589	2.4007
	WILKE	38.7165	0.0587	38.7752
AP	TP	5.9061	0.0615	5.9676
RB	GRACE	12.5343	0.0936	12.6279
	R-ROME	17.3135	0.0606	17.3741
	LTE	0.0076	0.0634	0.0710
	RECIPE	0.0078	0.0598	0.0676

Table 3: Average time (s) taken for a single edit and model inference after 10,000 edits.

served that non-retrieval-based methods lead to a significant reduction in general capabilities. This can be attributed to the accumulation of pattern mismatches caused by external interventions of editing. Among retrieval-based methods, LTE also exhibits performance degradation. In contrast, our RECIPE does not involve direct intervention on LLM parameters but instead relies on concatenating a short prompt to guide the LLM’s adherence to knowledge. It demonstrates the best preservation of general performance, suggesting that it inflicts minimal harm on the model.

5.2 Efficiency Comparison

To underscore the efficiency of RECIPE, we conduct a comparative analysis on editing and inference time after 10,000 edits, as delineated in Table 3. Among methods leveraging edit-specific training such as MEND, MALMEN, LTE, and RECIPE, a notable reduction in editing time is observed when compared to techniques necessitating multiple iterations of back-propagation during editing. For inference speed, methods that mod-

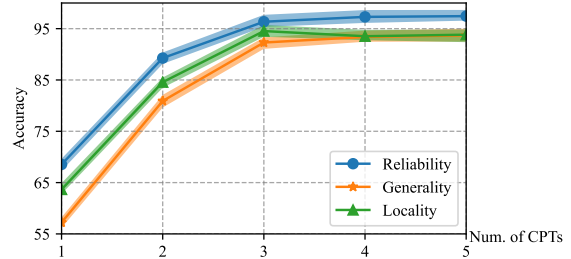


Figure 3: Impact of the number of CPTs on editing performance of RECIPE.

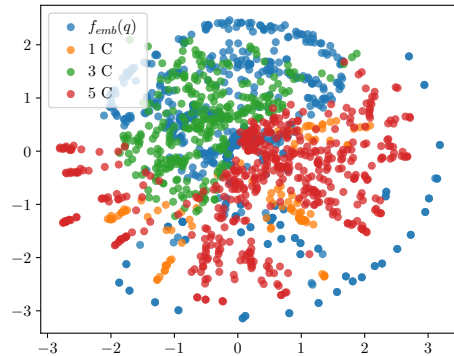


Figure 4: Visualization of word embeddings with varying numbers of CPTs.

ify model parameters maintain consistent speeds as they do not alter the original inference pipeline. T-Patcher slows down the inference speed due to the accumulating neurons. Among retrieval-based methods, GRACE reduces the parallelism in model inference due to its unique dictionary pairing mechanism. R-ROME and LTE need to calculate editing matrices on the fly and concatenate long editing instructions, respectively. In contrast, RECIPE effectively preserves the LLM’s original inference speed by concatenating short continuous prompts for editing. The shortest total time also highlights RECIPE’s efficiency advantage.

5.2.1 Number of Continuous Prompt Tokens

To assess whether an increase in Continuous Prompt Tokens (CPTs) can enhance the editing performance of RECIPE (Kaplan et al., 2020), Figure 3 illustrates the average impact of varying CPTs on editing efficacy across the editing benchmarks after 1,000 edits. The results show a noticeable performance dip with a single CPT, particularly in generality, indicating that fewer tokens limit representational capacity and lead to learning overly-simple patterns. Optimal editing performance is observed with three CPTs. Beyond this, while reliability and generality improve modestly, locality slightly decreases. This suggests that more CPTs

expand representational capabilities but also introduce additional LLMs’ interference.

Regarding the peak editing performance with three CPTs, we suggest that this is because the information carried by edit facts can be succinctly represented as relational triples (Head Entity, Relation, Tail Entity), and these triples can be represented as three word-level token embeddings. Thus, we further visualize LLAMA-2’s word embeddings of subjects and objects of 100 edit facts in CF, along with the corresponding representations of 1, 3, and 5 CPTs, reduced to two dimensions using t-SNE (Van der Maaten and Hinton, 2008). From Figure 4, the representations with three CPTs are closer to word embeddings than the others, indicating that the granularity of information carried by three CPTs is more akin to that of word embeddings of LLMs.

5.3 Ablation Study

We conduct an ablation study using LLAMA-2 on ZSRE (Mitchell et al., 2022), CF (Meng et al., 2022), and RIPE (Cohen et al., 2023). Average results are detailed in Table 4. Without CPTs, we resort to using word embeddings of knowledge statements as retrieval prompts from the knowledge repository. Excluding KS involved applying a conventional contrastive learning loss to align reliability and generality sample representations closer to editing knowledge while distancing those of locality samples. Upon completion of training, we employ an absolute similarity threshold decision strategy (Han et al., 2023) for filtering irrelevant knowledge. Despite its high locality, the omission of CPTs significantly impairs RECIPE’s reliability and generality. It can be observed that the results are nearly identical to those obtained without using an editor at all. This underscores that merely using raw concatenated knowledge prefixes fails to make LLMs comply with editing directives. Conversely, CPTs aid LLM adherence to specified edits. Additionally, discarding KS leads to a deterioration in editing efficacy, particularly impacting generality and locality. The reason is that an absolute similarity threshold fails to adequately address the diverse thresholds required by distinct queries.

6 Conclusion

We propose RECIPE, an effective and efficient LLM editor that includes two essential modules. Continuous prompt learning prefixes transformed

Settings	100 Edits			1000 Edits		
	Rel.	Gen.	Loc.	Rel.	Gen.	Loc.
N/A	27.30	26.07	100.00	27.30	26.07	100.00
RECIPE	97.29	93.74	97.38	96.05	92.34	95.36
- CPT	27.42	26.18	99.98	27.38	26.15	99.97
- KS	95.55	89.10	92.45	94.01	86.63	88.55
- BOTH	27.41	26.17	99.96	27.35	26.12	99.94

Table 4: Ablation study of RECIPE.

knowledge to input query to achieve efficient post-retrieval editing. Dynamic prompt retrieval with KS retrieves and determines whether the repository contains relevant knowledge without fixed similarity thresholds. In lifelong editing, RECIPE demonstrates exceptional editing performance and efficiency while simultaneously preserving LLM functionality without degradation.

Limitations

Due to the limitation in machine resources, we have not experimented on larger knowledge encoders apart from RoBERTa (Liu et al., 2019) and larger LLMs. We speculate that either a larger encoder or a larger LLM may yield better editing performance. Additionally, the current editing experiments are primarily implemented on QA-based datasets. We will expand our RECIPE framework to other types of editing tasks and larger models in the future.

Acknowledgments

This work is supported by the National Key Research and Development Program of China under Grant No. 2022ZD0120302.

References

- Nicola De Cao, Wilker Aziz, and Ivan Titov. 2021. [Editing factual knowledge in language models](#). In *EMNLP*, pages 6491–6506.
- Qi Chen, Bing Zhao, Haidong Wang, Mingqin Li, Chuanjie Liu, Zengzhong Li, Mao Yang, and Jingdong Wang. 2021. [SPANN: highly-efficient billion-scale approximate nearest neighborhood search](#). In *NeurIPS*, pages 5199–5212.
- Roi Cohen, Eden Biran, Ori Yoran, Amir Globerson, and Mor Geva. 2023. [Evaluating the ripple effects of knowledge editing in language models](#). *CoRR*, abs/2307.12976.
- Qingxiu Dong, Damai Dai, Yifan Song, Jingjing Xu, Zhifang Sui, and Lei Li. 2022. [Calibrating factual knowledge in pretrained language models](#). In *EMNLP*, pages 5937–5947.

- Haonan Duan, Adam Dziedzić, Nicolas Papernot, and Franziska Boenisch. 2023. [Flocks of stochastic parrots: Differentially private prompt learning for large language models](#). In *NeurIPS*.
- Tianyu Gao, Adam Fisch, and Danqi Chen. 2021. [Making pre-trained language models better few-shot learners](#). In *ACL*, pages 3816–3830.
- Xiaoqi Han, Ru Li, Hongye Tan, Yuanlong Wang, Qinghua Chai, and Jeff Z. Pan. 2023. [Improving sequential model editing with fact retrieval](#). In *EMNLP*, pages 11209–11224.
- Thomas Hartvigsen, Swami Sankaranarayanan, Hamid Palangi, Yoon Kim, and Marzyeh Ghassemi. 2022. [Aging with GRACE: lifelong model editing with discrete key-value adaptors](#). *CoRR*, abs/2211.11031.
- Kaiming He, Haoqi Fan, Yuxin Wu, Saining Xie, and Ross B. Girshick. 2020. [Momentum contrast for unsupervised visual representation learning](#). In *CVPR*, pages 9726–9735.
- Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. 2016. [Deep residual learning for image recognition](#). In *CVPR*, pages 770–778.
- Dan Hendrycks, Collin Burns, Steven Basart, Andy Zou, Mantas Mazeika, Dawn Song, and Jacob Steinhardt. 2021. [Measuring massive multitask language understanding](#). In *ICLR*.
- Chenhui Hu, Pengfei Cao, Yubo Chen, Kang Liu, and Jun Zhao. 2024. [Wilke: Wise-layer knowledge editor for lifelong knowledge editing](#). *CoRR*, abs/2402.10987.
- Edward J. Hu, Yelong Shen, Phillip Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang, and Weizhu Chen. 2022. [Lora: Low-rank adaptation of large language models](#). In *ICLR*.
- Zeyu Huang, Yikang Shen, Xiaofeng Zhang, Jie Zhou, Wenge Rong, and Zhang Xiong. 2023. [Transformer-patcher: One mistake worth one neuron](#). In *ICLR*.
- Yuxin Jiang, Yufei Wang, Chuhan Wu, Wanjun Zhong, Xingshan Zeng, Jiahui Gao, Liangyou Li, Xin Jiang, Lifeng Shang, Ruiming Tang, Qun Liu, and Wei Wang. 2024. [Learning to edit: Aligning llms with knowledge editing](#). *CoRR*, abs/2402.11905.
- Jared Kaplan, Sam McCandlish, Tom Henighan, Tom B. Brown, Benjamin Chess, Rewon Child, Scott Gray, Alec Radford, Jeffrey Wu, and Dario Amodei. 2020. [Scaling laws for neural language models](#). *CoRR*, abs/2001.08361.
- Itay Levy, Ben Bogin, and Jonathan Berant. 2023. [Diverse demonstrations improve in-context compositional generalization](#). In *ACL*, pages 1401–1422.
- Omer Levy, Minjoon Seo, Eunsol Choi, and Luke Zettlemoyer. 2017. [Zero-shot relation extraction via reading comprehension](#). In *CoNLL*, pages 333–342.
- Mike Lewis, Yinhan Liu, Naman Goyal, Marjan Ghazvininejad, Abdelrahman Mohamed, Omer Levy, Veselin Stoyanov, and Luke Zettlemoyer. 2020. [BART: denoising sequence-to-sequence pre-training for natural language generation, translation, and comprehension](#). In *ACL*, pages 7871–7880.
- Xiang Lisa Li and Percy Liang. 2021. [Prefix-tuning: Optimizing continuous prompts for generation](#). In *ACL*, pages 4582–4597.
- Xiao Liu, Kaixuan Ji, Yicheng Fu, Zhengxiao Du, Zhilin Yang, and Jie Tang. 2021a. [P-tuning v2: Prompt tuning can be comparable to fine-tuning universally across scales and tasks](#). *CoRR*, abs/2110.07602.
- Xiao Liu, Kaixuan Ji, Yicheng Fu, Weng Tam, Zhengxiao Du, Zhilin Yang, and Jie Tang. 2022. [P-tuning: Prompt tuning can be comparable to fine-tuning across scales and tasks](#). In *ACL*, pages 61–68.
- Xiao Liu, Yanan Zheng, Zhengxiao Du, Ming Ding, Yujie Qian, Zhilin Yang, and Jie Tang. 2021b. [GPT understands, too](#). *CoRR*, abs/2103.10385.
- Yinhan Liu, Myle Ott, Naman Goyal, Jingfei Du, Mandar Joshi, Danqi Chen, Omer Levy, Mike Lewis, Luke Zettlemoyer, and Veselin Stoyanov. 2019. [Roberta: A robustly optimized BERT pretraining approach](#). *CoRR*, abs/1907.11692.
- Kevin Meng, David Bau, Alex Andonian, and Yonatan Belinkov. 2022. [Locating and editing factual associations in GPT](#). In *NeurIPS*.
- Kevin Meng, Arnab Sen Sharma, Alex J. Andonian, Yonatan Belinkov, and David Bau. 2023. [Mass-editing memory in a transformer](#). In *ICLR*.
- Abhika Mishra, Akari Asai, Vidhisha Balachandran, Yizhong Wang, Graham Neubig, Yulia Tsvetkov, and Hannaneh Hajishirzi. 2024. [Fine-grained hallucination detection and editing for language models](#). *CoRR*, abs/2401.06855.
- Eric Mitchell, Charles Lin, Antoine Bosselut, Chelsea Finn, and Christopher D. Manning. 2022. [Fast model editing at scale](#). In *ICLR*.
- Jesse Mu, Xiang Li, and Noah D. Goodman. 2023. [Learning to compress prompts with gist tokens](#). In *NeurIPS*.
- Yixin Nie, Adina Williams, Emily Dinan, Mohit Bansal, Jason Weston, and Douwe Kiela. 2020. [Adversarial NLI: A new benchmark for natural language understanding](#). In *ACL*, pages 4885–4901.
- Fabio Petroni, Tim Rocktäschel, Sebastian Riedel, Patrick S. H. Lewis, Anton Bakhtin, Yuxiang Wu, and Alexander H. Miller. 2019. [Language models as knowledge bases?](#) In *EMNLP*, pages 2463–2473.
- Pranav Rajpurkar, Robin Jia, and Percy Liang. 2018. [Know what you don’t know: Unanswerable questions for squad](#). In *ACL*, pages 784–789.

- Pranav Rajpurkar, Jian Zhang, Konstantin Lopyrev, and Percy Liang. 2016. [Squad: 100, 000+ questions for machine comprehension of text](#). In *EMNLP*, pages 2383–2392.
- Nils Reimers and Iryna Gurevych. 2019. [Sentence-bert: Sentence embeddings using siamese bert-networks](#). In *EMNLP*, pages 3980–3990.
- Konstantinos I. Roumeliotis and Nikolaos D. Tselikas. 2023. [Chatgpt and open-ai models: A preliminary review](#). *Future Internet*, 15(6):192.
- Alon Talmor, Jonathan Herzig, Nicholas Lourie, and Jonathan Berant. 2019. [Commonsenseqa: A question answering challenge targeting commonsense knowledge](#). In *NAACL*, pages 4149–4158.
- Chenmien Tan, Ge Zhang, and Jie Fu. 2023. [Massive editing for large language models via meta learning](#). *CoRR*, abs/2311.04661.
- Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier Martinet, Marie-Anne Lachaux, Timothée Lacroix, Baptiste Rozière, Naman Goyal, Eric Hambro, Faisal Azhar, Aurélien Rodriguez, Armand Joulin, Edouard Grave, and Guillaume Lample. 2023. [Llama: Open and efficient foundation language models](#). *CoRR*, abs/2302.13971.
- Aäron van den Oord, Yazhe Li, and Oriol Vinyals. 2018. [Representation learning with contrastive predictive coding](#). *CoRR*, abs/1807.03748.
- Laurens Van der Maaten and Geoffrey Hinton. 2008. [Visualizing data using t-sne](#). *Journal of machine learning research*, 9(11).
- Peng Wang, Ningyu Zhang, Xin Xie, Yunzhi Yao, Bozhong Tian, Mengru Wang, Zekun Xi, Siyuan Cheng, Kangwei Liu, Guozhou Zheng, and Huajun Chen. 2023a. [Easyedit: An easy-to-use knowledge editing framework for large language models](#). *CoRR*, abs/2308.07269.
- Yizhong Wang, Yeganeh Kordi, Swaroop Mishra, Alisa Liu, Noah A. Smith, Daniel Khashabi, and Hannaneh Hajishirzi. 2023b. [Self-instruct: Aligning language models with self-generated instructions](#). In *ACL*, pages 13484–13508.
- Ziyun Xu, Chengyu Wang, Minghui Qiu, Fuli Luo, Runxin Xu, Songfang Huang, and Jun Huang. 2023. [Making pre-trained language models end-to-end few-shot learners with contrastive prompt tuning](#). In *WSDM*, pages 438–446.
- Yunzhi Yao, Peng Wang, Bozhong Tian, Siyuan Cheng, Zhoubo Li, Shumin Deng, Huajun Chen, and Ningyu Zhang. 2023. [Editing large language models: Problems, methods, and opportunities](#). *CoRR*, abs/2305.13172.
- Lang Yu, Qin Chen, Jie Zhou, and Liang He. 2024. [MELO: enhancing model editing with neuron-indexed dynamic lora](#). In *AAAI*, pages 19449–19457.
- Aohan Zeng, Xiao Liu, Zhengxiao Du, Zihan Wang, Hanyu Lai, Ming Ding, Zhuoyi Yang, Yifan Xu, Wendi Zheng, Xiao Xia, Weng Lam Tam, Zixuan Ma, Yufei Xue, Jidong Zhai, Wenguang Chen, Zhiyuan Liu, Peng Zhang, Yuxiao Dong, and Jie Tang. 2023. [GLM-130B: an open bilingual pre-trained model](#). In *ICLR*.
- Taolin Zhang, Qizhou Chen, Dongyang Li, Chengyu Wang, Xiaofeng He, Longtao Huang, Hui Xue’, and Jun Huang. 2024. [DAFNet: Dynamic auxiliary fusion for sequential model editing in large language models](#). In *ACL*, pages 1588–1602.
- Zhenru Zhang, Chuanqi Tan, Haiyang Xu, Chengyu Wang, Jun Huang, and Songfang Huang. 2023. [Towards adaptive prefix tuning for parameter-efficient language model fine-tuning](#). In *ACL*, pages 1239–1248.
- 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?](#) *CoRR*, abs/2305.12740.
- Kaijie Zhu, Qinlin Zhao, Hao Chen, Jindong Wang, and Xing Xie. 2023. [Promptbench: A unified library for evaluation of large language models](#). *CoRR*, abs/2312.07910.

Algorithm 1 Training of RECIPE

```

1: Input: LLM to be edited  $f_{llm}$ ; initialized collection of RECIPE parameters  $\mathcal{M}$ ; training set  $\mathcal{D} = \left\{ \left( (q_{e_i}, a_{e_i}), \{q_{g_i}^j, a_{g_i}^j\}_{j=1}^{N_{g_i}}, \{q_{l_i}^j, a_{l_i}^j\}_{j=1}^{N_{l_i}} \right) \right\}_{i=1}^N$ ; maximum iteration number  $I_{max}$ ; batch size of training samples  $b$ ; learning rate  $\eta$ .
2: Output: trained RECIPE parameters  $\mathcal{M}$ .
3: while  $iter < I_{max}$  do
4:    $\{(q_{e_i}, a_{e_i}), (q_{g_i}, a_{g_i}), (q_{l_i}, a_{l_i})\}_{i=1}^b \leftarrow$  Sample  $b$  training samples from  $\mathcal{D}$ 
5:   for  $i \leftarrow 1$  to  $b$  do
6:     # Get editing knowledge
7:      $k_i = q_{e_i} + a_{e_i}$  # String concatenation
8:     # Get knowledge representation
9:      $r_{k_i} \leftarrow$  Transform  $k_i$  using Eq. 6
10:    # Get continuous prompts
11:     $p_{k_i} \leftarrow$  Transform  $r_{k_i}$  using Eq. 7
12:    # Get query representations
13:     $\tilde{r}_{q_{e_i}} \leftarrow$  Transform  $q_{e_i}$  using Eq. 8
14:     $\tilde{r}_{q_{g_i}} \leftarrow$  Transform  $q_{g_i}$  using Eq. 8
15:     $\tilde{r}_{q_{l_i}} \leftarrow$  Transform  $q_{l_i}$  using Eq. 8
16:  end for
17:  # Get knowledge representation of KS
18:   $r_{\Theta} \leftarrow$  Transform  $\Theta$  using Eq. 6
19:  # Compute loss and update parameters
20:   $\mathcal{L}_{edit}, \mathcal{L}_{pl} \leftarrow$  Compute losses using Eq.13 and Eq.16
21:   $\mathcal{L}_{total} = \mathcal{L}_{edit} + \mathcal{L}_{pl}$ 
22:   $\mathcal{M} \leftarrow$  Adam  $(\nabla_{\mathcal{M}} \mathcal{L}_{total}, \eta)$ 
23: end while
24: return  $\mathcal{M}$ 

```

A Algorithms of RECIPE

The training and editing algorithms for RECIPE are detailed in Alg. 1 and Alg. 2, respectively.

Algorithm 2 Editing of RECIPE in a Lifelong Scenario

1: **Input:** Knowledge retrieval repository $\mathcal{K}_{t-1} = \{(r_{k_\tau}, p_{k_\tau})\}_{\tau=1}^{t-1}$; editing knowledge (q_{e_t}, a_{e_t}) .
2: **Output:** updated knowledge retrieval repository \mathcal{K}_t .
3: # Get editing knowledge
4: $k_t = q_{e_t} + a_{e_t}$ # String concatenation
5: # Get knowledge representation
6: $r_{k_t} \leftarrow$ Transform k_t using Eq. 6
7: # Get continuous prompts
8: $p_{k_t} \leftarrow$ Transform r_{k_t} using Eq. 7
9: # Update knowledge retrieval repository
10: $\mathcal{K}_t = \mathcal{K}_{t-1} \cup \{(r_{k_t}, p_{k_t})\}$
11: **return** \mathcal{K}_t

Algorithm 3 Inference of LLM Equipped with RECIPE

1: **Input:** LLM to be edited f_{llm} , including the embedding layer f_{emb} and the transformer module \hat{f}_{llm} ; knowledge retrieval repository $\mathcal{K}_t = \{(r_{k_\tau}, p_{k_\tau})\}_{\tau=1}^t$; knowledge representation of KS r_Θ ; input query q .
2: **Output:** LLM’s output with RECIPE intervened a_q .
3: # Get query representation
4: $\tilde{r}_q = \text{MLP}_Q(f_{rm}(q))$
5: # Get the index of knowledge with the largest similarity
6: $j = \arg \max_{\tau=1, \dots, t} \tilde{r}_q^T \cdot r_{k_\tau}$
7: # Filter irrelevant knowledge and get output
8: **if** $\tilde{r}_q^T \cdot r_{k_j} > \tilde{r}_q^T \cdot r_\Theta$ **then**
9: $a_q = \hat{f}_{llm}(p_{k_j} \oplus f_{emb}(q))$
10: **else**
11: $a_q = f_{llm}(q)$
12: **end if**
13: **return** a_q

The inference process of the LLM equipped with RECIPE is described in Alg. 3.

B Datasets and Baselines

B.1 Model Editing Datasets

We employ three public model editing datasets, including ZSRE (Mitchell et al., 2022), CounterFact (CF) (Meng et al., 2022), and Ripple Effect (RIPE) (Cohen et al., 2023) as our experimental datasets. For methods that require edit training, including MEND (Mitchell et al., 2022), MALMEN (Tan et al., 2023), LTE (Jiang et al., 2024), and our RECIPE, we utilize the above training sets to learn their parameters.

ZSRE (Levy et al., 2017) is generated through question-answering with BART (Lewis et al., 2020) and manual filtering, including 162,555 training and 19,009 testing samples. Each sample comprises an editing sample and its rephrased and irrelevant counterparts, matching the reliability, generality, and locality editing properties.

CF (Meng et al., 2022) is characterized by the editing of false facts and includes 10,000 training

and 10,000 testing samples. These false facts are more likely to conflict with the original knowledge within LLMs, making the editing process more challenging and thus providing a robust evaluation of the editors’ ability to enforce edits.

RIPE (Cohen et al., 2023) differentiates the generality and locality properties into fine-grained types, comprising 3,000 training and 1,388 testing samples. The generality of each sample includes logical generalization, combination I, combination II, and subject aliasing, while the locality data cover forgetfulness and relation specificity.

B.2 General Datasets of LLMs

To evaluate the damage of editors to the general performance of LLMs, we select four prevalent benchmarks to assess LLMs’ general capabilities. They are CSQA (Talmor et al., 2019) to evaluate commonsense knowledge, ANLI (Nie et al., 2020) for reasoning abilities, MMLU (Hendrycks et al., 2021) to gauge exam capabilities, and SQuAD-2 (Rajpurkar et al., 2018) for comprehension skills. PromptBench (Zhu et al., 2023) is utilized as the evaluation framework for this experiment.

CSQA (CommonSense Question Answering) (Talmor et al., 2019) is designed to evaluate LLMs’ commonsense knowledge through multiple-choice questions. It includes 12,102 samples, split into 9,741 for training, 1,221 for validation, and 1,140 for testing.

ANLI (Adversarial Natural Language Inference) (Nie et al., 2020) evaluates LLMs’ natural language reasoning capacity by determining whether the relationship between a premise and a hypothesis is one of entailment, contradiction, or neutrality. The difficulty of the tasks increases across three rounds. It includes a total of 169,265 samples, with 162,865 for training, 3,200 for validation, and 3,200 for testing.

MMLU (Massive Multitask Language Understanding) (Hendrycks et al., 2021) tests LLMs’ mastery of specialized domain knowledge through multiple-choice questions covering 57 different academic fields and disciplines, such as history, literature, law, and biology. The dataset comprises a total of 6,783 questions distributed across testing, validation, and development sets, containing 5,871, 627, and 285 samples, respectively.

SQuAD-2 (Stanford Question Answering Dataset version 2) (Rajpurkar et al., 2018) assesses the reading comprehension abilities of LLMs by posing questions based on paragraphs taken from

over 500 Wikipedia articles. Compared to its first version (Rajpurkar et al., 2016), its challenge lies in the inclusion of questions that do not have answers derivable from the text. The dataset contains a total of 142,192 questions, with 130,319 in the training set and 11,873 in the validation set. We report the performance on its validation set with default hyper-parameter settings.

Among them, we use selection accuracy as the evaluation metric for CSQA, MMLU, and ANLI, and the reciprocal of the model’s perplexity (PPL) on the predicted sequence as the evaluation metric for SQuAD-2. All metrics are scaled to a range of 0-100.

B.3 Baselines

In addition to fine-tuning (FT) as the basic baseline, we compare our RECIPE approach with various strong editing baselines. **MEND** (Mitchell et al., 2022) trains an MLP to transform the low-rank decomposition of the gradients of the model to be edited with respect to the editing samples. **ROME** (Meng et al., 2022) first uses causal mediation analysis to locate the layer that has the greatest impact on the editing sample. **MEMIT** (Meng et al., 2023) expands the editing scope to multiple layers based on ROME, thereby improving editing performance and supporting batch editing. **T-Patcher** (Huang et al., 2023) (TP) attaches and trains additional neurons in the FFN of the last layer of the model to be edited. **MALMEN** (Tan et al., 2023) formulates the parameter shift aggregation as a least square problem, subsequently updating the LM parameters using the normal equation. **WILKE** (Hu et al., 2024) selects the editing layer based on the pattern matching degree of editing knowledge across different layers.

We also leverage competitive retrieval-based editing methods to validate the effectiveness further. **GRACE** (Hartvigsen et al., 2022) proposes retrieval adapters for continuous editing, which maintains a dictionary-like structure to construct new mappings for potential representations that need to be modified. **RASE** (Han et al., 2023) leverages factual information to enhance editing generalization and guide the identification of edits by retrieving related facts from the fact-patch memory. In our baseline settings, we use the ROME (Meng et al., 2022) model as the specific basic editor for RASE to perform the editing task, named **R-ROME**. **LTE** (Jiang et al., 2024) elicits the capabilities of LLMs to follow knowledge editing in-

structions, thereby empowering them to effectively leverage updated knowledge to answer queries.

C Model Settings and Training Details

RECIPE. (1) **Hyper-parameter Settings:** For our proposed RECIPE, we use the same hyper-parameter settings across different backbones, including LLAMA-2², GPT-J³, and GPT2-XL⁴. The number of continuous prompt tokens and KS tokens are set as $l = 3$ and $c = 10$, respectively. MLP_K , MLP_Q , and MLP_P are each composed of two linear layers, with an intermediate dimension set to 4096 and are connected in a residual manner (He et al., 2016). The dimensions of the knowledge and query representations are also set to 4096. The total numbers of RECIPE’s training parameters for GPT2-XL, GPT-J, and LLAMA-2 are 220M, 250M, and 250M, respectively. (2) **Training Details:** We set the learning rate ($\eta = 1e - 5$), the batch size to 8, and the maximum number of iterations to 150,000. A checkpoint is saved every 5000 iterations, and ultimately, the one with the smallest loss is selected for evaluation. The training process requires approximately 3 days on an NVIDIA A800 GPU. These experiments are presented on average with 5 random runs, using different random seeds but the same hyper-parameters. **Baseline Models.** For R-ROME (Han et al., 2023) and LTE (Jiang et al., 2024), we implement the settings mentioned in their respective papers and trained them on the same datasets as ours. For the other baselines, we follow the same settings as described in EasyEdit (Wang et al., 2023a) for training and evaluation.

D Results with Different Backbones

Lifelong editing experiments on GPT-J (6B) and GPT2-XL (1.5B) are presented in Table 5 and Table 6. The results also show a similar conclusion with general results, demonstrating the efficacy of our method. Notably, comparing Tables 1, 5, and 6, ROME and MEMIT exhibit a significant performance decline on LLAMA-2 compared to the other two backbones. MEMIT is an improved version of ROME, and both methods aim to edit the weights of the selected FFN layers. We speculate that this performance discrepancy may stem from structural

²<https://huggingface.co/meta-llama/llama-2-7b-hf>

³<https://huggingface.co/EleutherAI/gpt-j-6b>

⁴<https://huggingface.co/openai-community/gpt2-xl>

differences in the FFN layers between LLAMA-2, GPT2-XL, and GPT-J. The original ROME and MEMIT papers assume a common two-layer structure for the LLM’s FFN layers, which is exactly the configuration of GPT2-XL and GPT-J. However, the FFN of LLAMA-2 consists of three linear layers, which may account for the ineffectiveness of ROME and MEMIT in adapting to the LLAMA-2 architecture.

E Summary of Innovations

We summarize the innovations of this work as follows.

First, to the best of our knowledge, RECIPE is the first method to employ prompt learning for solving model editing. Indeed, previous works have used string-constructed prefixes to accomplish model editing, such as IKE (Zheng et al., 2023) and LTE (Jiang et al., 2024). However, we are the first to explore how to achieve model editing with the shortest editing prefixes through prompt learning. Second, a direct application of prompt learning might involve training a continuous prompt as a fixed prefix, which can append any editing string behind it to form a complete instruction for the LLM to follow. This approach treats the entire editing task as a prompt learning task. However, the concatenation of the continuous prompt with the editing string results in a lengthy prefix. Moreover, it can be reasonably hypothesized that a smaller amount of trainable parameters might make it challenging to converge to satisfactory editing performance. In contrast, our RECIPE innovatively treats the editing of each piece of knowledge as a mini prompt learning task, and considers the training of the encoder that generates continuous editing prompts as a meta-task. This strategy significantly reduces the length of editing prefixes while enhancing the representational capabilities of the editing training, thus ensuring editing performance. Third, RECIPE completely decouples knowledge editing from model parameters through prompt learning, thereby avoiding model degradation due to external intervention in model parameters or overfitting on the editing dataset.

