

Assessing and Post-Processing Black Box Large Language Models for Knowledge Editing

Anonymous ACL submission

Abstract

The task of Knowledge Editing (KE) is aimed at efficiently and precisely adjusting the behavior of large language models (LLMs) to update specific knowledge while minimizing any adverse effects on other knowledge. Current research predominantly concentrates on editing white-box LLMs, neglecting a significant scenario: editing black-box LLMs, where access is limited to interfaces and only textual output is provided. In this paper, we initially officially introduce KE on black-box LLMs, followed by presenting a thorough evaluation framework aimed at addressing the shortcomings of current evaluations, which are inadequate for black-box LLMs editing and lack comprehensiveness. To address privacy leaks of editing data and style over-editing in existing approaches, we propose a new postEdit framework, ensuring privacy through downstream processing and maintaining textual style consistency via fine-grained editing. Experiments and analysis conducted on two benchmarks show that postEdit surpasses all baselines and exhibits robust generalization, notably enhancing style retention by an average of +20.82%.¹

1 Introduction

As large language models (LLMs) are widely applied to knowledge-intensive tasks and the world’s state evolves, the requirements of updating LLMs to rectify obsolete information or incorporate new knowledge to maintain their relevance is constantly emerging (Zhao et al., 2023; Liu et al., 2023a; Bian et al., 2023; Wang et al., 2023a). Frequent retraining is impractical due to intensive computational overload and time consumption. To address this issue, the concept of knowledge editing (KE) has been proposed, aiming to efficiently and precisely modify the behavior of LLMs to update specific knowledge without negatively influencing other

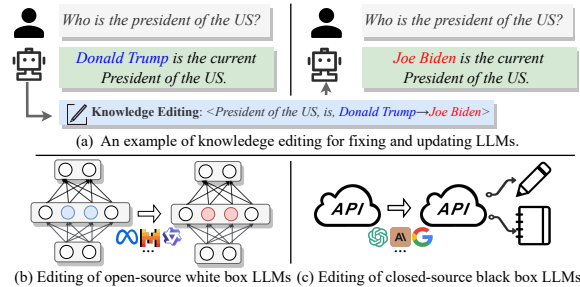


Figure 1: Illustration of Knowledge Editing and comparison of two editing scenarios, where black-box LLMs editing constrains LLMs to only obtain textual output.

knowledge (Yao et al., 2023; Wang et al., 2023b; Zhang et al., 2024), as illustrated in Fig 1(a).

A prevalent approach to KE involves manipulating the internals of LLMs through gradients or causal analysis (De Cao et al., 2021; Mitchell et al., 2021; Meng et al., 2022a,b; Huang et al., 2023), as depicted in Fig 1(b). While these methods have shown promise, they require LLMs to be locally deployed and parameter-transparent, termed white-box LLMs editing. In more typical scenarios, LLMs are provided via APIs by upstream manufacturers (e.g., OpenAI, Google) for downstream services, with inaccessible internal workings and text-only output. We refer to KE on such LLMs as **black-box LLMs editing**, as shown in Fig 1(c). This raises a key question: *how can we edit "black-box" models when undesired outputs or errors occur?* Furthermore, existing KE evaluation protocols rely on changes in the model’s logits before and after editing, and are unattainable for black-box LLMs, prompting another question: *how can we comprehensively evaluate black-box KE methods?*

There are some studies based on external memory that can be applied to black-box LLM editing scenarios. SERAC (Mitchell et al., 2022) utilizes a surrogate model to generate edited responses when queries are classified within the editing scope (INS), while relying on the base LLM for queries out of the editing scope (OOS). IKE (Zheng et al.,

¹We will release our code after blind review.

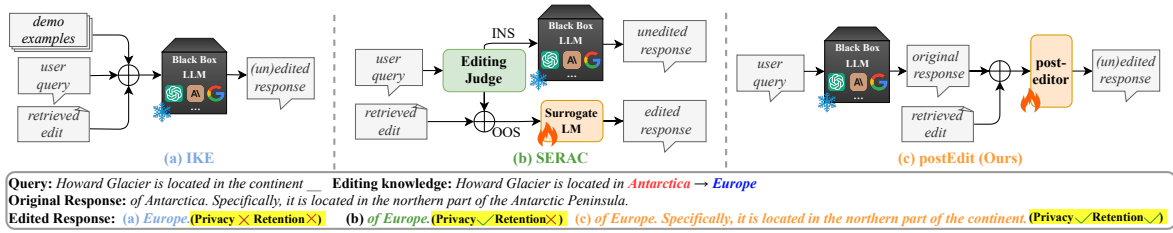


Figure 2: Comparison of different KE frameworks for black-box LLM editing. IKE operates on LLM input, and SERAC performs editing using a surrogate model parallel to LLM, while our postEdit edits after the output of LLM and achieves both privacy protection and style retention.

2023) facilitates in-context learning (Dong et al., 2022) of LLM itself by demonstrating exemplars to learn the ability to discern the need of editing and how to edit. However, as depicted in Fig 2(a)(b), these methods encounter two crucial drawbacks: (1) **Privacy leakage of editing data.** IKE inputs recall data from the demonstration library and edit memory to LLMs, inevitably disclosing downstream private editing data to upstream LLM providers. (2) **Style over-editing.**² One of the core objectives of KE is to ensure localized editing, whereby KE methods should only edit the knowledge of LLMs while keeping the original output style unchanged. Specifically, the different scales or types between the surrogate model and base LLM result in stylistic differences for SERAC, while LLM’s sensitivity to prompts and demonstrations (Chen et al., 2023) leads to style over-editing in IKE. Therefore, even though their edited responses both target the new object "Europe", they exhibit a pronounced departure in style from the original responses. An ideal black-box editing method should preserve downstream data privacy while achieving commendable editing performance and style retention.

In this paper, we firstly revisit the existing evaluation of KE and formulate an improved general evaluation framework for black-box LLM editing. In addition to the traditional lexical evaluation of knowledge editing, our framework incorporates the assessment of style retention for the first time and conducts a comprehensive evaluation from both textual and semantic perspectives. (see Section 3). To solve the problems of existing methods mentioned above, we propose a novel post-editing approach termed **postEdit**, applied after the output of LLMs, as illustrated in Fig 2(c). Diverging from previous approaches, on the one hand, the post-processing mechanism allows postEdit to be deployed as a post-plugin at the downstream end, safeguarding

²In this paper, the style extensively covers the expressive forms, conciseness, length, information, etc., of the text.

the privacy of editing data. On the other hand, an expert model called post-editor, guided by editing knowledge, makes fine-grained modifications to original responses generated by LLM, thereby effectively preserving the original style. As the role of post-editor is to discern and precisely edit the original response rather than storing new knowledge, we integrate edit memory and a retriever into postEdit, like IKE and SERAC, for efficient knowledge injection. We leave the detailed exposition in Section 4. Finally, we conduct comprehensive experiments and analysis to demonstrate that postEdit achieves outstanding performance in both editing and style retention, exhibiting robust generalization across various aspects, including LLMs, data, and scales in Section 5 and 6.

Our contributions are three-fold: (1) We officially introduce knowledge editing on black-box LLMs and propose a comprehensive KE evaluation framework, incorporating the assessment of style retention for the first time. (2) We propose a novel postEdit method to post-edit the output of LLMs through an expert model in a plug-in manner. Our postEdit can both maintain the privacy of downstream editing data and achieve commendable editing performance and style retention. (3) Experiments and analysis on two benchmarks demonstrate that our postEdit outperforms all baselines in both editing and style retention (Retention Score +20.82% ↑), showing robust generalization.

2 Related Work

2.1 Knowledge Editing

White-box LLMs Editing The initial KE methods involve updating parameters using constrained fine-tuning (Sinitsin et al., 2020; Zhu et al., 2020). Recent studies mostly center around hyper-network and attribution. Hyper-network-based approaches (De Cao et al., 2021; Mitchell et al., 2021) train a hyper-network to capture gradient changes for specific edits, while attribute-based methods (Dai

et al., 2022; Meng et al., 2022a,b; Wu et al., 2023; Li et al., 2024) locate neuron activation in networks for targeted parameter updates. However, these approaches exclusively focus on editing in white-box LLM scenarios, overlooking concerns on black-box LLMs editing.

Memory-based Editing In addition to injecting edits as parameters into LLM, memory-based KE methods store edits in explicit memory and utilize retrieval-augmented methods to adjust the model’s final predictions based on relevant edits. Unlike conventional Retrieval-Augmented Generation (RAG) methods (Gao et al., 2024) focus on enhancing document retrieval, KE methods concentrate on modifying knowledge for INS queries and maintain output consistency for OOS queries. Therefore, SERAC (Mitchell et al., 2022) introduces an INS/OOS judge model, while IKE (Zheng et al., 2023) uses demonstrations with INS and OOS examples to determine whether to edit or maintain knowledge. Although applicable to black-box editing scenarios, these methods face challenges related to privacy and style over-editing.

2.2 Post-processing Methods

Some post-processing methods have been applied to other tasks. Cao et al. (2020) fine-tune a BART model to improve factual consistency in abstractive summarization by using summaries with errors as input and original or gold summaries as training targets. Thorne and Vlachos (2021) fine-tune a T5 model to correct factual errors by recovering masked statements based on retrieved evidence. RARR (Gao et al., 2023) employs PaLM with few-shot demonstrations for error correction and attribution report generation. Different from these studies, postEdit applies post-processing to the knowledge editing task, fine-tuning a post-editor to simultaneously determine query relevance within the editing scope and make fine-grained modifications.

3 Evaluation Framework

3.1 Problem Formulation

A knowledge entry is typically shown as a triple (subject, relationship, object). Following Wang et al. (2023b), an edit can be defined as $e = (t, t^*) = (s, r, o \rightarrow o^*)$, denoting the update of an old knowledge triple t to the new one t^* . As multiple input-output pairs can be associated with the same tuple, the input set associated with edit e is denoted as $\mathcal{X}_e = I(s, r)$, referred to as in-

scope (INS) input space, the target output set associated with o^* is denoted as $\mathcal{Y}_e^* = O^*(s, r, o^*)$, and the corresponding original output set is denoted as $\mathcal{Y}_e = O(s, r, o)$. For a base LLM $f_{base} : \mathcal{X} \rightarrow \mathcal{Y}$, given an edit e , the goal of KE is to modify the original output $y_o \in \mathcal{Y}_e$ to $y_e \in \mathcal{Y}_e^*$ for input $x \in \mathcal{X}_e$, while keeping the output unaffected for out-of-scope (OOS) queries, i.e., $y_e = y_o$ if $x \notin \mathcal{X}_e$.

Furthermore, we define KE on black-box LLMs as the editing on a certain class of LLMs, where we have no access to anything other than textual outputs of LLMs. It should be noted that this scenario only restricts the base LLM to be edited, with no limitations imposed on auxiliary models or tools.

3.2 Evaluation Protocol

3.2.1 Existing Logit-based Evaluation

Previous studies (Meng et al., 2022a; Mitchell et al., 2022; Zheng et al., 2023) primarily assess KE based on three metrics: **Efficacy**, **Generalization**, and **Specificity**, by calculating the change in logits of the model before and after editing.³ On the one hand, the inaccessibility of logits for black-box LLMs renders these metrics ineffective. On the other hand, KE should only modify spans in the response involving the edit, while keeping the rest and style unchanged to minimize negative impacts of editing. However, this aspect has been fully overlooked, leading to incomplete evaluation.

3.2.2 Improved Multi-perspective Evaluation

For black-box LLMs editing, the evaluation of KE focuses on what changes and what remains in the edited output y_e compared to original output y_o . Therefore, we formulate the evaluation framework from both the aspects of editing and retention.

Editing The Editing metric is designed to evaluate the editing for INS input and non-editing for OOS input. When $x \in \mathcal{X}_e$, the expected output space of f_{base} transitions from \mathcal{Y}_e to \mathcal{Y}_e^* . From the perspective of textual editing (**TE**), \mathcal{Y}_e^* discards the old target o and incorporates the new target o^* . From the perspective of semantic editing (**SE**), the joint text composed of \mathcal{X}_e and \mathcal{Y}_e^* implies the new knowledge t^* and contradicts the old knowledge t . When $x \notin \mathcal{X}_e$, the situation is reversed. We formalize TE as follows:

$$\text{TE} = \begin{cases} \frac{1}{2} \{ \text{ctn}(y_e, o^*) + (1 - \text{ctn}(y_e, o)) \} & x \in \mathcal{X}_e \\ \frac{1}{2} \{ \text{ctn}(y_e, o) + (1 - \text{ctn}(y_e, o^*)) \} & x \notin \mathcal{X}_e \end{cases} \quad (1)$$

³We provide details of these metrics in Appendix A.1.

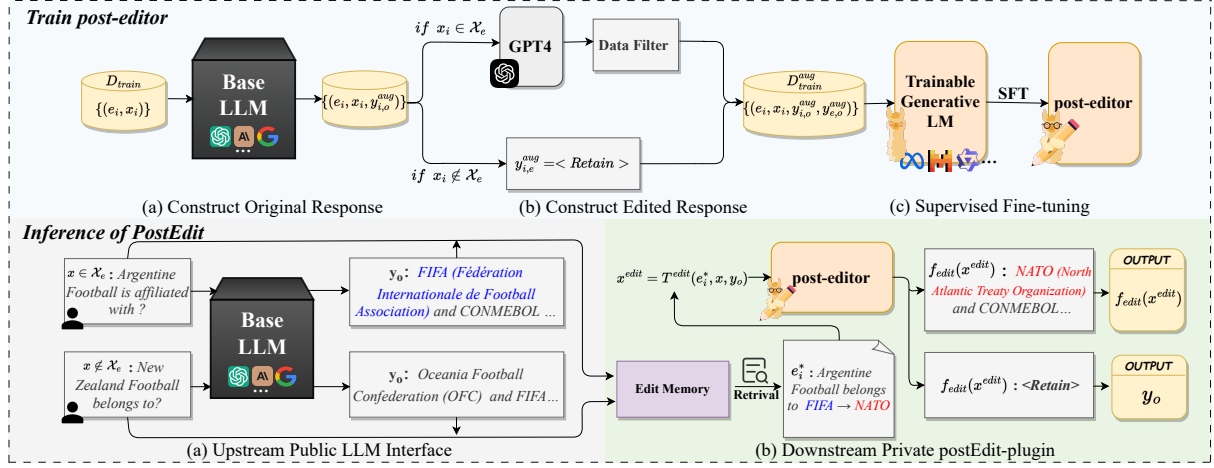


Figure 3: The overall architecture of postEdit. The post-editor is trained to learn: (1) distinguish between INS and OOS queries; (2) edit the output of INS queries while preserving style. Pseudo-code is provided in Appendix B.1.

where $\text{ctn}(a, b) = 1$ if **a** contains **b**, otherwise 0. Similarly, SE is formalized as follows:

$$\text{SE} = \begin{cases} \frac{1}{2} \{ \text{ent}([x, y_e], t^*) + (1 - \text{ent}([x, y_e], t)) \} & x \in \mathcal{X}_e \\ \frac{1}{2} \{ \text{ent}([x, y_e], t_o) + (1 - \text{ent}([x, y_e], t^*)) \} & x \notin \mathcal{X}_e \end{cases} \quad (2)$$

where $\text{ent}(a, b) = 1$ if **a** entails **b**, otherwise 0 by using the Natural Language Inference (NLI) model, $[x, y_e]$ denotes the concatenation of input-output pair, and t_o indicates the knowledge tuple associated with OOS input-output pair $[x, y_o]$.

Retention To assess the extent to which the edited output preserves the original style, we introduce Retention as an adversarial metric for Editing. We separately evaluate textual retention (**TR**) and semantic retention (**SR**) using ROUGE scores (Lin, 2004) and the SBERT model (Reimers and Gurevych, 2019), formalized as follows:

$$\text{TR} = \begin{cases} \text{ROUGE}(M(y_e, o^*), M(y_o, o)) & x \in \mathcal{X}_e \\ \text{ROUGE}(y_e, y_o) & x \notin \mathcal{X}_e \end{cases} \quad (3)$$

$$\text{SR} = \begin{cases} \text{sim}(M(y_e, o^*), M(y_o, o)) & x \in \mathcal{X}_e \\ \text{sim}(y_e, y_o) & x \notin \mathcal{X}_e \end{cases} \quad (4)$$

where $M(a, b)$ denotes masking the span relevant to **b** in **a**. For $x \in \mathcal{X}_e$, we employ a masking operation to extract text unrelated to editing.

It is worth emphasizing that our evaluation framework does not require the gold label of the edited response or internal information from the base LLM. This enables its applicability to a wide range of scenarios beyond black-box LLM editing.

Due to space limitations, we further elaborate on and discuss the proposed evaluation framework in Appendix A.2 and provide pseudo-code in Appendix A.3. Subsequently, we conduct extensive

experiments on the consistency between these metrics and human evaluation in Appendix A.4, where excellent Pearson Consistency scores validate the rationality of the proposed metrics.

4 Methodology

4.1 Overall Architecture

To solve the problems of privacy leakage of editing data and style over-editing, as illustrated in Fig 3, postEdit is deployed downstream and post-processes the output of base LLM, comprising three components: an edit-memory $M_e = \{e_i\}$ for storing editing knowledge, a retriever f_{retr} for recalling an edit, and a trained generative model named post-editor f_{edit} for executing the edit⁴. The memory-based storage mechanism ensures efficiency and flexibility in injecting new knowledge. During the inference phase, the retriever first recalls the edit with the highest similarity to user input from M_e . Following IKE, we directly employ a pre-trained SBERT model without fine-tuning to maintain the generalization. Finally, the post-editor performs the editing guided by recalled edit.

4.2 Train post-editor

Original Response Augment The training dataset of KE typically consists of editing knowledge, along with queries covering both INS and OOS input, denoted as $D_{train} = \{(e_i, x_i)\}$. Previous studies (Mitchell et al., 2022; Zheng et al., 2023) usually directly use the new object o_i^* in e_i as the target output for editing, resulting in stylistic differences between the editor and base LLM.

⁴In the main experiment, we fine-tune LLaMa 2-7B (Touvron et al., 2023) as the post-editor and conduct an analysis of performance at various scales in Section 6.5.

To address this gap, we first construct the original response $y_{i,o}^{aug} = f_{base}(x_i)$ via base LLM for each sample.

Edited Response Augmentation In order to construct the training output targets for post-editor, we utilize both GPT-4 and rules to further augment the training dataset. For INS inputs, the objective is to modify the original response. Thus, given edit e_i , input x_i , and original output $y_{i,o}^{aug}$ are aggregated using an editing template T^{aug5} and fed into GPT-4 to obtain the edited output $y_{i,e}^{aug}$. For OOS inputs, the goal is to maintain the original response without modification. Therefore, we introduce a special token $\langle Retain \rangle$ as the target output, denoting no need for editing. We formulate this process as:

$$y_{i,e}^{aug} = \begin{cases} f_{gpt4}(T^{aug}(e_i, x_i, y_{i,o}^{aug})) & x_i \in \mathcal{X}_e \\ \langle Retain \rangle & x_i \notin \mathcal{X}_e \end{cases} \quad (5)$$

Recent studies (Zhou et al., 2023; Lu et al., 2023; Liu et al., 2023b) have proven that the quality of training data is often more crucial than quantity. To further enhance the quality of augmented data and alleviate training burden, we evaluate and filter the edited responses obtained through GPT-4 augment. Based on the joint evaluation using the Editing metrics TE and SE, formalized as $\mathbf{1}_{\{TE=1 \& SE=1\}} y_{i,e}^{aug}$, augmented samples with poor quality are discarded. Ultimately, we obtain the augmented training set $D_{train}^{aug} = \{(e_i, x_i, y_{i,o}^{aug}, y_{i,e}^{aug})\}$.

Supervised Fine-tuning (SFT) After data augment and filtering, the post-editor is trained in a supervised fine-tuning manner, where the query, edit, and original response are aggregated as input using an editing template T^{edit} (distinct from T^{aug}), with $y_{i,e}^{aug}$ as the output target. After tokenizing $y_{i,e}^{aug}$ as $\{y_{i,e_1}^{aug}, y_{i,e_2}^{aug}, \dots, y_{i,e_T}^{aug}\}$, the loss function of SFT can be formalized as follows:

$$\mathcal{L}_{sft} = - \sum_{i=1}^{|D_{train}^{aug}|} \sum_{t=0}^{T-1} \log P(y_{i,e_{t+1}}^{aug} | x_i^{edit}, y_{i,e_{\leq t}}) \quad (6)$$

where $x_i^{edit} = T^{edit}(e_i, x_i, y_{i,o}^{aug})$.

4.3 Inference of PostEdit

For a user query $x \in D_{test}$, the original response $y_o = f_{base}(x)$ is obtained through the upstream LLM interface. On the downstream side, the retriever recalls the most similar edit e_{i^*} to x from

the edit memory:

$$i^* = \operatorname{argmax}_{0 \leq i < |M_e|} \operatorname{sim}(x, e_i) \quad (7)$$

Next, we obtain the input $x^{edit} = T^{edit}(e_{i^*}, x, y_o)$ by populating the editing template T^{edit} and transmit it to the post-editor to yield the output $f_{edit}(x^{edit})$. Finally, by discerning whether $f(x^{edit})$ contains the special token $\langle Retain \rangle$, we determine the ultimate output:

$$y_e = \begin{cases} f_{edit}(x^{edit}) & f_{edit}(x^{edit}) \neq \langle Retain \rangle \\ y_o & f_{edit}(x^{edit}) = \langle Retain \rangle \end{cases} \quad (8)$$

5 Experiments

5.1 Experiment Setting

Datasets We conduct experiments on two widely-used KE datasets, CounterFact (Meng et al., 2022a) and zsRE (Levy et al., 2017), where edits in the training and test sets don't overlap. Each entry comprises an edit and three types of queries: **Simple** queries to validate the success of knowledge injection, **Rephrase** queries to assess the generalization of the edit, and **out-of-scope (OOS)** queries to verify the local effect of the edit. Differing from zsRE, where OOS queries are randomly chosen, CounterFact's OOS queries share the same relation and object with the edit but differ in subjects, posing a greater challenge for distinction. We provide details and processing procedures in Appendix C.1.

Baselines We employ ChatGPT (gpt-3.5-turbo) as the base LLM and extensively compare postEdit with methods applicable to black-box LLM editing, including PROMPT (Zheng et al., 2023), IKE (Zheng et al., 2023), SERAC (Mitchell et al., 2022), and SERAC(ChatGPT). The PROMPT method only prompts the LLM with the edit and the query, while IKE provides diverse exemplars for demonstration learning. SERAC employs a fine-tuned surrogate model⁶ to respond to queries within the editing scope, and SERAC(ChatGPT) is a variant where the surrogate model is changed to ChatGPT. Detailed introduction of baselines are shown in Appendix C.2 and more baselines from other tasks are compared in Appendix D.1.

Test Procedure The default test procedure of KE involves editing a single piece of knowledge, assessing it, and then rolling back the system to original state before moving on to the next edit. This

⁵All templates mentioned are shown in Appendix B.2.

⁶For a fair comparison, the surrogate model uses the same pre-trained model and training data as the post-editor.

Method	Textual Editing (TE)				Semantic Editing (SE)				Textual Retention (TR)				Semantic Retention (SR)			
	Simple	Rephrase	OOS	AVG (HM)	Simple	Rephrase	OOS	AVG (HM)	Simple	Rephrase	OOS	AVG (HM)	Simple	Rephrase	OOS	AVG (HM)
PROMPT	85.17	86.73	63.8	78.57 _(76.62)	83.1	84.57	61.97	76.54 _(74.65)	21.42	21.54	18.11	20.36 _(20.19)	53.14	54.86	51.37	53.13 _(53.05)
IKE	94.2	85.8	85.4	88.47 _(88.29)	93.2	84.5	85.3	87.67 _(87.5)	24.14	18.98	22.81	21.97 _(21.75)	53.45	48.94	57.69	53.36 _(53.12)
SERAC	95.4	87.4	96.1	92.97 _(92.79)	94.6	87.3	96.2	92.7 _(92.53)	35.66	37.62	96.01	56.43 _(46.13)	65.51	64.64	97.04	75.73 _(73.1)
SERAC (ChatGPT)	95.23	85.8	98.6	93.2 _(92.87)	95.3	86	98.6	93.31 _(92.98)	23.43	26.71	96.41	48.85 _(33.08)	55.04	56.88	97.91	69.95 _(65.26)
postEdit (ours)	96.8	94.7	99.4	96.97 _(96.93)	92.5	92.1	99.4	94.67 _(94.55)	88.65	89.66	99.64	92.65 _(92.39)	93.9	94.02	99.82	95.91 _(95.84)

Table 1: Performance comparison on CounterFact. AVG is the direct average, while HM is the harmonic mean. We bold the best and underline the second-best results. Results are averaged over three random runs.

Method	Textual Editing (TE)				Semantic Editing (SE)				Textual Retention (TR)				Semantic Retention (SR)			
	Simple	Rephrase	OOS	AVG (HM)	Simple	Rephrase	OOS	AVG (HM)	Simple	Rephrase	OOS	AVG (HM)	Simple	Rephrase	OOS	AVG (HM)
PROMPT	88.83	86.87	58.37	78.02 _(74.53)	86.5	84.97	60.27	77.24 _(74.29)	47.76	45.35	34.93	42.68 _(41.51)	73.4	74.62	61.29	69.77 ₍₆₉₎
IKE	98.1	97.6	78	91.23 _(90.2)	97.7	94.7	83.1	91.83 _(91.38)	19.72	16.36	27.83	21.3 _(20.3)	42.26	38.67	58.53	46.49 _(45.04)
SERAC	98.7	95.1	100	97.93 _(97.89)	97.6	93.3	100	96.97 _(96.89)	68.02	66.06	100	78.03 _(75.3)	86.84	85.91	100	90.92 _(90.48)
SERAC (ChatGPT)	94.7	87.5	100	94.07 _(93.77)	96.17	88.53	100	94.9 _(94.61)	52.22	52.01	100	68.08 _(61.75)	75.2	77.56	100	84.25 _(82.69)
postEdit (ours)	<u>98.4</u>	98.6	100	99 _(98.99)	96.2	95.4	100	97.2 _(97.16)	95.76	96.13	100	97.3 _(97.26)	97.69	97.89	100	98.53 _(98.52)

Table 2: Performance comparison on zsRE.

setting keeps the edit memory size at 1, turning the retriever into an "oracle" to encourage methods to prioritize editing and locality capabilities. We compare methods under various memory sizes in Section 6.4 and discuss the efficiency of methods in Appendix E.

5.2 Main Results

Table 1 and Table 2 show the main results of postEdit and comparable baselines on two benchmark KE datasets. In general, our postEdit method consistently outperforms all baselines with a large margin, both in terms of Editing and Retention scores. Next, we analyze the results from three aspects:

(1) **Comparison of different methods.** We can see that postEdit achieves nearly all optimal Editing scores, along with a significant surpassing of baselines in Retention scores. On CounterFact, postEdit outperforms the suboptimal baselines by 3.77% (TE), 1.36% (SE), 36.22% (TR), and 20.18% (SR) in average scores. On zsRE, postEdit surpasses the suboptimal baselines by 1.07% (TE), 0.23% (SE), 19.27% (TR), and 7.61% (SR). This shows that postEdit can accurately locate and modifies spans in the text related to editing, while maintaining other content, thereby achieving high performance in both Editing and Retention.

(2) **Comparison of different query types.** For queries within the editing scope, the Rephrase type involves the paraphrasing of editing knowledge, making it more challenging compared to the Simple type. Concerning CounterFact, discernible decrements in Rephrase performance are observed for IKE and SERAC in contrast to the Simple type (e.g., TE score, IKE: 94.2→85.8, SERAC: 95.5→87.4), whereas postEdit performance remains stable (96.8→94.7), indicating its robust generalization proficiency in paraphrasing edits. For OOS queries, while SERAC and postEdit excel

on the zsRE dataset, postEdit surpasses SERAC on more challenging CounterFact, showcasing its precise differentiation of queries requiring editing without additional editing judge module.

(3) **Comparison of different metrics.** Comparing the Editing and Retention of baselines reveals a serious issue of style over-editing. For example, the Editing performance of IKE surpasses that of PROMPT, while the Retention lags behind PROMPT, indicating a negative impact of demonstration on IKE’s style retention. Despite achieving commendable Editing scores, SERAC and SERAC (ChatGPT) still fall short in terms of Retention. This highlights that effective editing does not guarantee good retention, emphasizing the need for a comprehensive evaluation of knowledge editing.

6 Analysis

6.1 Generalization of PostEdit

In Section 4.1, we fine-tune the post-editor to acquire the ability of discriminating and executing edits. Therefore, it is imperative to validate the generalization of post-editor’s abilities. For postEdit and baselines, we initially utilize ChatGPT as the base LLM and CounterFact as the training set or demonstration library. Subsequently, we conduct testing under different base LLMs and datasets **without re-training**, as illustrated in Fig 4.

We can see that whether generalizing from CounterFact to zsRE or from ChatGPT to PaLM2⁷ and LLaMA2-70B-chat⁸, postEdit consistently demonstrates optimal performance in Editing and Retention. The robust generalization of post-editor highlights its plug-and-play applicability across diverse scenarios, requiring no retraining when faced with a new set of editing requests or when replacing

⁷<https://ai.google/discover/palm2>

⁸<https://huggingface.co/meta-LLaMA>

ID	Edit	Query	Original Response	Edited Response		
				IKE	SERAC	postEdit
1	The nationality of Marcel Maupi was what? <i>French</i> → <i>Italian</i>	What was Marcel Maupi's nationality?	Marcel Maupi was a <i>French</i> national.	Italian.	Marcel Maupi was Italian.	Marcel Maupi was an <i>Italian</i> national.
2	Football Association of Malaysia is affiliated with the <i>FIFA</i> → <i>Avengers</i>	Football Association of Malaysia is affiliated with _	- Asian Football Confederation (AFC) - <i>FIFA (Fédération Internationale de Football Association)</i>	Avengers.	Avengers.	- Asian Football Confederation (AFC) - <i>Avengers</i>
3	What is the programming language for CLPython? <i>Lisp</i> → <i>JavaScript</i>	What is the programming language for CLPython?	CLPython is an implementation of Python language in <i>Common Lisp</i> , which means that the programming language used for CLPython is <i>Common Lisp</i> .	JavaScript.	The programming language for CLPython is JavaScript.	CLPython is an implementation of Python language in <i>JavaScript</i> , which means that the programming language used for CLPython is <i>JavaScript</i> .
4	Howard Glacier is located in <i>Antarctica</i> → <i>Europe</i>	Howard Glacier is located in the continent _	of <i>Antarctica</i> . Specifically, it is located in the <i>northern part of the Antarctic Peninsula</i> .	Europe.	of Europe.	of <i>Europe</i> . Specifically, it is located in the <i>northern part of the continent</i> .
5	In what city or state did the formation of I Am Kloot occur? <i>Manchester</i> → <i>Bilbao</i>	In which city or state was the formation of I Am Kloot?	I Am Kloot was formed in <i>Manchester, England</i> .	Bilbao.	I Am Kloot was formed in Bilbao.	I Am Kloot was formed in <i>Bilbao, Spain</i> .

Table 3: Editing cases sampled from CounterFact and zsRE under different methods.

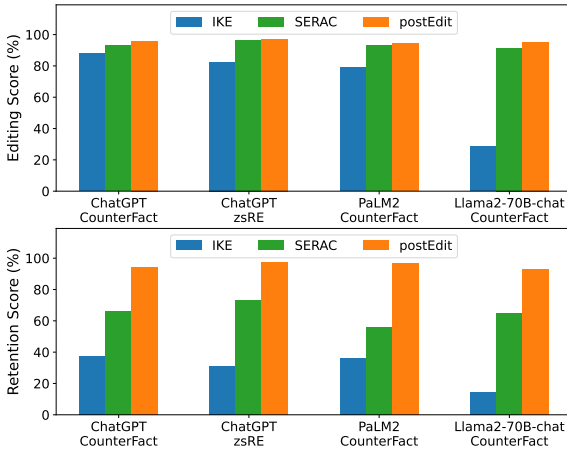


Figure 4: Performance under different base LLMs and datasets, where Editing Score is the average of TE and SE, and Retention Score is the average of TR and SR.

the base LLM. In contrast, both IKE and SERAC exhibit performance fluctuations, particularly evident in a significant decline when IKE is applied to LLaMA2-70B-chat. Further analysis reveals that conflicts between editing data and the intrinsic knowledge of LLaMA2-70B-chat lead to frequent refusals to generate responses based on edits. However, postEdit successfully mitigated the impact of knowledge conflicts through post-processing.

6.2 Case Study

To visually demonstrate the editing and style retention of postEdit and baselines, we conduct the case study in Table 3. In Case 1, postEdit accurately identifies and modifies "*French*" to "*Italian*" while maintaining the rest of the text unchanged to keep the style to the greatest extent. In contrast, IKE only responds with "*Italian*" and SERAC replies with "*Marcel Maupi was Italian*" without referencing the original response, revealing serious style over-editing. In Cases 2 and 3, postEdit respectively replaces "*FIFA (Fédération Internationale de Football Association)*" with "*Avengers*" and modifies "*Common Lisp*" to "*JavaScript*". This demonstrates

Method	Semantic Editing (SE)				Semantic Retention (SR)			
	Simple	Rephrase	OOS	AVG	Simple	Rephrase	OOS	AVG
postEdit	92.5	92.1	99.4	94.67	93.9	94.02	99.82	95.91
<i>Module Ablation</i>								
-w/o data filter	90.6	90.6	99.4	93.53	94.19	93.76	99.82	95.92
post-editor→ChatGPT	89.73	87.8	70.77	82.54	89.39	88.78	83.27	86.26
GPT4→ChatGPT	93.2	91.8	99.4	94.80	90.04	89.54	99.81	93.13
SBERT Judgement	92.2	85.2	96.3	91.23	94.47	92.49	98.97	95.31
<i>Training Data Ablation</i>								
-w/o Simple	91.8	91.2	99.5	94.17	93.96	94.21	99.89	96.02
-w/o Rephrase	92	12.9	99.8	68.23	94.37	71.67	99.95	88.66
-w/o OOS	92.2	91.5	4.7	62.8	94.47	94.12	75.01	87.86

Table 4: Ablation Study on CounterFact.

that postEdit can locate and edit spans semantically related to editing knowledge, going beyond a rudimentary replacement of old objects with new ones. Furthermore, it is evident that postEdit can handle spans logically associated with the editing. In Case 4, the location changes from "*Antarctica*" to "*Europe*", and the span in the original response, describing the location as "*the northern part of the Antarctic Peninsula*", is correspondingly adjusted to "*the northern part of the continent*". Similarly, in Case 5, as "*Manchester*" is changed to "*Bilbao*", the country is also edited from "*England*" to "*Spain*".

6.3 Ablation Study

To understand the roles of each component and training data type in postEdit, we conduct ablation study in Table 4.

Module Ablation In our postEdit framework, we utilize GPT-4 to generate edited responses and subsequently perform data filtering. After removing data filtering, the SE score for INS queries exhibits a decline (Simple -1.9 and Rephrase -1.5), indicating that data filtering effectively enhances the quality of training data. Replacing the post-editor with ChatGPT results in a noticeable decline in performance across different types. This suggests that LLMs like ChatGPT are not proficient performing such editing tasks, highlighting the need for fine-tuning the post-editor. Substituting GPT-4 with ChatGPT for edited response augmentation results in a slight SE score increase (avg +0.13) but a significant SR score decrease (avg -2.78). This

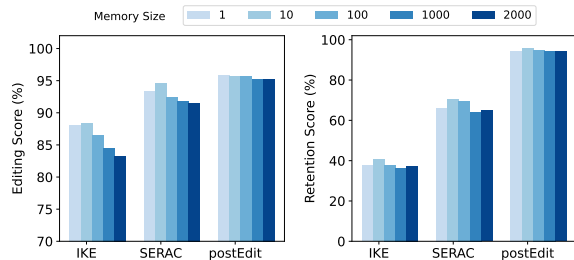


Figure 5: Performance of methods under different Edit Memory size on CounterFact.

indicates that ChatGPT lacks the fine-grained granularity in editing compared to GPT-4, thereby resulting in a coarser-grained post-editor. Finally, we introduce the editing judging module, the same as SERAC, through comparing the SBERT semantic similarity with a threshold. The observed decrease in Rephrase and OOS scores demonstrates the superior discriminative capability of the post-editor. **Training Data Ablation** We further conduct data ablation by removing each type of data from the training set. We observe that removing Simple data has no notable impact, while the removal of Rephrase data leads to a significant drop (-79.2) in the SE metric. This indicates that Rephrase data plays a crucial role in improving the post-editor’s ability for editing knowledge injection and generalization, while relying solely on Simple data doesn’t suffice for achieving the post-editor’s generalization. After removing OOS data, although there is a noticeable decline in OOS metrics, the metrics for Simple and Rephrase do not show a discernible improvement. This indicates that post-editor doesn’t excessively compromise its ability to perform edits when learning to discriminate editing.

6.4 Effect of Memory Size

In real-world scenarios, as the world evolves, edited knowledge should be continuously infused and preserved, i.e., the size of Edit Memory will continue to expand⁹. For the edit retrieved from Edit Memory, IKE utilizes the base LLM itself, SERAC applies a similarity threshold, and postEdit employs the post-editor to determine whether the query is within the scope of editing. We evaluate the performance of these methods under varying memory sizes in Fig 5. With the same retriever, postEdit exhibits the highest robustness among methods in both Editing and Retention scores, substantiating the superiority of the postEdit mechanism in discerning the necessity of editing.

⁹In some studies, this corresponds to Batch Editing and Sequence Editing.

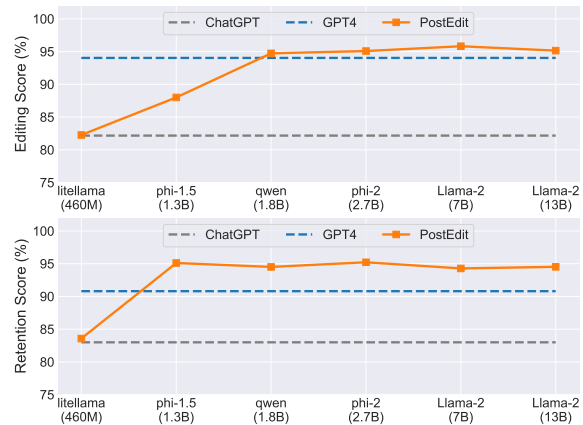


Figure 6: Performance curves of the post-editor at different scales on CounterFact.

6.5 Effect of Post-editor Scale

To investigate the effect of post-editor scale on performance, we compare evaluation scores across models ranging from 460M to 13B in size. As illustrated in Fig 6, it is evident that with the increase in post-editor scale, editing scores gradually improve (significant from 460M to 1.8B, followed by slower gains beyond 1.8B), while retention score remains stable after reaching 1.3B. This suggests that editing ability is more influenced by the model scale, and a larger post-editor can enhance editing performance while maintaining the retention. We also compare the effectiveness of post-editor with zero-shot ChatGPT and GPT-4. Similar to the findings in Section 6.3, LLMs like ChatGPT are not proficient in executing the editing task. Therefore, on CounterFact, the performance of the 460M post-editor is comparable to ChatGPT, and the 1.8B post-editor surpasses GPT-4. This indicates that the postEdit framework does not rely on a large-scale post-editor, and small-sized editors can achieve satisfactory performance and high efficiency.

7 Conclusion

In this paper, we firstly introduce a comprehensive evaluation framework for knowledge editing under black-box LLMs, incorporating multiple perspectives and considering the style retention. Next, we propose a novel postEdit framework to address existing issues in privacy leakage of editing data and style over-editing in current methods by post-processing the output of LLMs. Finally, experiments on two benchmarks and thorough analysis demonstrate that postEdit outperforms all baselines and achieves strong generalization.

591 Limitations

592 This paper primarily investigates the assessment
593 and methodology of knowledge editing in black-
594 box LLM scenarios. The proposed evaluation
595 framework can comprehensively assess edited re-
596 sponses from multiple perspectives, and the postE-
597 dit method effectively addresses issues related to
598 privacy concerns of editing data and style over-
599 editing. However, our work also has several limita-
600 tions: (1) Although our proposed evaluation frame-
601 work and postEdit method mainly focus on knowl-
602 edge editing in black-box LLM scenarios, they can
603 be equally applied to editing in white-box LLM sce-
604 narios. Due to constraints in length and the focus
605 of the paper, we haven't thoroughly explored this
606 in the paper. (2) Although the postEdit framework
607 does not require retraining when injecting editing
608 knowledge, it still necessitates an initial fine-tuning
609 phase to enable the post-editor to learn the ability
610 to discern whether a query is within the editing
611 scope and how to perform the editing, resulting in a
612 certain computational load. (3) Our study primarily
613 investigates the application of knowledge editing
614 in knowledge question answering tasks, similar to
615 previous research. We believe that our framework
616 can be extended to other scenarios, such as fact-
617 checking and sentiment editing. We leave these
618 explorations for future research.

619 Ethic Consideration

620 In this paper, we propose a knowledge editing ap-
621 proach that can be flexibly applied downstream to
622 post-process the outputs of LLMs, effectively safe-
623 guarding the privacy of downstream private editing
624 data and maintaining consistency in the style of the
625 LLM. While the purpose of knowledge editing is
626 to rectify errors or outdated knowledge in LLMs,
627 malicious knowledge editing may lead to the gen-
628 eration of harmful or inappropriate outputs by the
629 model. Therefore, ensuring secure and responsi-
630 ble practices in knowledge editing is of paramount
631 importance. The application of these techniques
632 should be guided by ethical considerations, with
633 safeguard measures in place to prevent misuse and
634 mitigate the potential for harmful outcomes. Ad-
635 ditionally, due to the difficulty in obtaining contin-
636 uously up-to-date knowledge, some KE datasets
637 such as CounterFact use counterfactual knowledge
638 to validate the effectiveness of methods. Further-
639 more, the base LLM, such as ChatGPT used in this
640 work, merely serves as a demonstration of research

on knowledge editing in black-box model scenar- 641
ios. We emphasize that these datasets and LLMs 642
are solely for academic exploration and do not in- 643
volve actual applications in real-world scenarios, 644
nor do they include content modification or attacks 645
on commercially used LLMs. 646

References 647

- Ning Bian, Xianpei Han, Le Sun, Hongyu Lin, Yaojie 648
Lu, and Ben He. 2023. Chatgpt is a knowledgeable 649
but inexperienced solver: An investigation of com- 650
monsense problem in large language models. *arXiv* 651
preprint arXiv:2303.16421. 652
- Meng Cao, Yue Dong, Jiapeng Wu, and Jackie Chi Kit 653
Cheung. 2020. [Factual error correction for abstrac-](#) 654
[tive summarization models](#). In *Proceedings of the* 655
2020 Conference on Empirical Methods in Natural 656
Language Processing (EMNLP), pages 6251–6258, 657
Online. Association for Computational Linguistics. 658
- Dhivya Chandrasekaran and Vijay Mago. 2021. [Evolu-](#) 659
[tion of semantic similarity—a survey](#). *ACM Comput-* 660
ing Surveys, 54(2):1–37. 661
- Yanda Chen, Chen Zhao, Zhou Yu, Kathleen McKe- 662
own, and He He. 2023. [On the relation between](#) 663
[sensitivity and accuracy in in-context learning](#). In 664
Findings of the Association for Computational Lin- 665
guistics: EMNLP 2023, pages 155–167, Singapore. 666
Association for Computational Linguistics. 667
- Damai Dai, Li Dong, Yaru Hao, Zhifang Sui, Baobao 668
Chang, and Furu Wei. 2022. [Knowledge neurons in](#) 669
[pretrained transformers](#). In *Proceedings of the 60th* 670
Annual Meeting of the Association for Computational 671
Linguistics (Volume 1: Long Papers), pages 8493– 672
8502, Dublin, Ireland. Association for Computational 673
Linguistics. 674
- Nicola De Cao, Wilker Aziz, and Ivan Titov. 2021. [Edit-](#) 675
[ing factual knowledge in language models](#). In *Pro-* 676
ceedings of the 2021 Conference on Empirical Meth- 677
ods in Natural Language Processing, pages 6491– 678
6506, Online and Punta Cana, Dominican Republic. 679
Association for Computational Linguistics. 680
- Qingxiu Dong, Lei Li, Damai Dai, Ce Zheng, Zhiy- 681
ong Wu, Baobao Chang, Xu Sun, Jingjing Xu, and 682
Zhifang Sui. 2022. A survey for in-context learning. 683
arXiv preprint arXiv:2301.00234. 684
- Luyu Gao, Zhuyun Dai, Panupong Pasupat, Anthony 685
Chen, Arun Tejasvi Chaganty, Yicheng Fan, Vincent 686
Zhao, Ni Lao, Hongrae Lee, Da-Cheng Juan, and 687
Kelvin Guu. 2023. [RARR: Researching and revising](#) 688
[what language models say, using language models](#). 689
In *Proceedings of the 61st Annual Meeting of the* 690
Association for Computational Linguistics (Volume 1: 691
Long Papers), pages 16477–16508, Toronto, Canada. 692
Association for Computational Linguistics. 693

694	Yunfan Gao, Yun Xiong, Xinyu Gao, Kangxiang Jia, Jinliu Pan, Yuxi Bi, Yi Dai, Jiawei Sun, Meng Wang, and Haofen Wang. 2024. Retrieval-augmented generation for large language models: A survey.	748
695		749
696		750
697		751
698	Edward J. Hu, Yelong Shen, Phillip Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang, and Weizhu Chen. 2021. Lora: Low-rank adaptation of large language models.	752
699		753
700		754
701		
702	Zeyu Huang, Yikang Shen, Xiaofeng Zhang, Jie Zhou, Wenge Rong, and Zhang Xiong. 2023. Transformer-patcher: One mistake worth one neuron. <i>arXiv preprint arXiv:2301.09785.</i>	755
703		756
704		757
705		758
706	Diederik P. Kingma and Jimmy Ba. 2017. Adam: A method for stochastic optimization.	759
707		
708	Omer Levy, Minjoon Seo, Eunsol Choi, and Luke Zettlemoyer. 2017. Zero-shot relation extraction via reading comprehension. In <i>Proceedings of the 21st Conference on Computational Natural Language Learning (CoNLL 2017)</i> , pages 333–342, Vancouver, Canada. Association for Computational Linguistics.	760
709		761
710		762
711		763
712		764
713		765
714	Xiaopeng Li, Shasha Li, Shezheng Song, Jing Yang, Jun Ma, and Jie Yu. 2024. Pmet: Precise model editing in a transformer. <i>Proceedings of the AAAI Conference on Artificial Intelligence</i> , 38(17):18564–18572.	766
715		767
716		768
717		769
718	Chin-Yew Lin. 2004. Rouge: A package for automatic evaluation of summaries. In <i>Text summarization branches out</i> , pages 74–81.	770
719		771
720		772
721	Jiacheng Liu, Wenya Wang, Dianzhuo Wang, Noah A Smith, Yejin Choi, and Hannaneh Hajishirzi. 2023a. Vera: A general-purpose plausibility estimation model for commonsense statements. <i>arXiv preprint arXiv:2305.03695.</i>	773
722		774
723		775
724		776
725		777
726	Wei Liu, Weihao Zeng, Keqing He, Yong Jiang, and Junxian He. 2023b. What makes good data for alignment? a comprehensive study of automatic data selection in instruction tuning. <i>arXiv preprint arXiv:2312.15685.</i>	778
727		779
728		780
729		781
730		782
731	Keming Lu, Hongyi Yuan, Zheng Yuan, Runji Lin, Junyang Lin, Chuanqi Tan, Chang Zhou, and Jingren Zhou. 2023. # instag: Instruction tagging for analyzing supervised fine-tuning of large language models. <i>arXiv e-prints</i> , pages arXiv–2308.	783
732		784
733		785
734		786
735		787
736	Kevin Meng, David Bau, Alex Andonian, and Yonatan Belinkov. 2022a. Locating and editing factual associations in gpt. <i>Advances in Neural Information Processing Systems</i> , 35:17359–17372.	788
737		789
738		790
739		791
740	Kevin Meng, Arnab Sen Sharma, Alex Andonian, Yonatan Belinkov, and David Bau. 2022b. Mass-editing memory in a transformer. <i>arXiv preprint arXiv:2210.07229.</i>	792
741		793
742		794
743		795
744	Sewon Min, Kalpesh Krishna, Xinxu Lyu, Mike Lewis, Wen-tau Yih, Pang Koh, Mohit Iyyer, Luke Zettlemoyer, and Hannaneh Hajishirzi. 2023. FActScore: Fine-grained atomic evaluation of factual precision	796
745		797
746		798
747		799
		800
		801
		802
	in long form text generation. In <i>Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing</i> , pages 12076–12100, Singapore. Association for Computational Linguistics.	
	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>	
	Eric Mitchell, Charles Lin, Antoine Bosselut, Christopher D Manning, and Chelsea Finn. 2022. Memory-based model editing at scale. In <i>International Conference on Machine Learning</i> , pages 15817–15831. PMLR.	
	Kishore Papineni, Salim Roukos, Todd Ward, and Wei-Jing Zhu. 2002. Bleu: a method for automatic evaluation of machine translation. In <i>Proceedings of the 40th Annual Meeting of the Association for Computational Linguistics</i> , pages 311–318, Philadelphia, Pennsylvania, USA. Association for Computational Linguistics.	
	Nils Reimers and Iryna Gurevych. 2019. Sentence-BERT: Sentence embeddings using Siamese BERT-networks. In <i>Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)</i> , pages 3982–3992, Hong Kong, China. Association for Computational Linguistics.	
	Anton Sinitsin, Vsevolod Plokhotnyuk, Dmitriy Pyrkov, Sergei Popov, and Artem Babenko. 2020. Editable neural networks. <i>arXiv preprint arXiv:2004.00345.</i>	
	James Thorne and Andreas Vlachos. 2021. Evidence-based factual error correction. In <i>Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume 1: Long Papers)</i> , pages 3298–3309, Online. Association for Computational Linguistics.	
	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>	
	Cunxiang Wang, Sirui Cheng, Zhikun Xu, Bowen Ding, Yidong Wang, and Yue Zhang. 2023a. Evaluating open question answering evaluation. <i>arXiv preprint arXiv:2305.12421.</i>	
	Song Wang, Yaochen Zhu, Haochen Liu, Zaiyi Zheng, Chen Chen, et al. 2023b. Knowledge editing for large language models: A survey. <i>arXiv preprint arXiv:2310.16218.</i>	
	Xinwei Wu, Junzhuo Li, Minghui Xu, Weilong Dong, Shuangzhi Wu, Chao Bian, and Deyi Xiong. 2023. DEPN: Detecting and editing privacy neurons in pre-trained language models. In <i>Proceedings of the 2023</i>	

803 *Conference on Empirical Methods in Natural Lan-*
804 *guage Processing*, pages 2875–2886, Singapore. As-
805 sociation for Computational Linguistics.

806 Yunzhi Yao, Peng Wang, Bozhong Tian, Siyuan Cheng,
807 Zhoubo Li, Shumin Deng, Huajun Chen, and Ningyu
808 Zhang. 2023. [Editing large language models: Prob-](#)
809 [lems, methods, and opportunities](#). In *Proceedings*
810 *of the 2023 Conference on Empirical Methods in*
811 *Natural Language Processing*, pages 10222–10240,
812 Singapore. Association for Computational Linguis-
813 tics.

814 Yuheng Zha, Yichi Yang, Ruichen Li, and Zhiting Hu.
815 2023. [AlignScore: Evaluating factual consistency](#)
816 [with a unified alignment function](#). In *Proceedings*
817 *of the 61st Annual Meeting of the Association for*
818 *Computational Linguistics (Volume 1: Long Papers)*,
819 pages 11328–11348, Toronto, Canada. Association
820 for Computational Linguistics.

821 Ningyu Zhang, Yunzhi Yao, Bozhong Tian, Peng
822 Wang, Shumin Deng, Mengru Wang, Zekun Xi,
823 Shengyu Mao, Jintian Zhang, Yuansheng Ni, Siyuan
824 Cheng, Ziwen Xu, Xin Xu, Jia-Chen Gu, Yong Jiang,
825 Pengjun Xie, Fei Huang, Lei Liang, Zhiqiang Zhang,
826 Xiaowei Zhu, Jun Zhou, and Huajun Chen. 2024. [A](#)
827 [comprehensive study of knowledge editing for large](#)
828 [language models](#).

829 Yue Zhang, Yafu Li, Leyang Cui, Deng Cai, Lemao Liu,
830 Tingchen Fu, Xinting Huang, Enbo Zhao, Yu Zhang,
831 Yulong Chen, Longyue Wang, Anh Tuan Luu, Wei
832 Bi, Freda Shi, and Shuming Shi. 2023. [Siren’s song](#)
833 [in the ai ocean: A survey on hallucination in large](#)
834 [language models](#).

835 Wayne Xin Zhao, Kun Zhou, Junyi Li, Tianyi Tang,
836 Xiaolei Wang, Yupeng Hou, Yingqian Min, Beichen
837 Zhang, Junjie Zhang, Zican Dong, et al. 2023. A
838 survey of large language models. *arXiv preprint*
839 *arXiv:2303.18223*.

840 Ce Zheng, Lei Li, Qingxiu Dong, Yuxuan Fan, Zhiyong
841 Wu, Jingjing Xu, and Baobao Chang. 2023. [Can](#)
842 [we edit factual knowledge by in-context learning?](#)
843 In *Proceedings of the 2023 Conference on Empiri-*
844 *cal Methods in Natural Language Processing*, pages
845 4862–4876, Singapore. Association for Computa-
846 tional Linguistics.

847 Zexuan Zhong, Zhengxuan Wu, Christopher Manning,
848 Christopher Potts, and Danqi Chen. 2023. [MQuAKE:](#)
849 [Assessing knowledge editing in language models via](#)
850 [multi-hop questions](#). In *Proceedings of the 2023*
851 *Conference on Empirical Methods in Natural Lan-*
852 *guage Processing*, pages 15686–15702, Singapore.
853 Association for Computational Linguistics.

854 Chunting Zhou, Pengfei Liu, Puxin Xu, Srini Iyer, Jiao
855 Sun, Yuning Mao, Xuezhe Ma, Avia Efrat, Ping Yu,
856 Lili Yu, et al. 2023. Lima: Less is more for alignment.
857 *arXiv preprint arXiv:2305.11206*.

858 Chen Zhu, Ankit Singh Rawat, Manzil Zaheer, Srinadh
859 Bhojanapalli, Daliang Li, Felix Yu, and Sanjiv Kumar.
2020. [Modifying memories in transformer models](#).
arXiv preprint arXiv:2012.00363.

A Details of Evaluation

A.1 Details of Existing Metrics

There are three metrics based on logits mainly used to evaluate the performance of knowledge editing in previous work, namely Efficacy, Generalization, and Specificity.

- **Efficacy** measures the accuracy of knowledge editing using **ES** (Efficacy Score) and **EM** (Efficacy Magnitude). For Simple type queries, the meaning of ES is $E[I[P(o^*) > P(o)]]$, and EM is obtained by $E[P(o^*) - P(o)]$.
- **Generalization** measures the accuracy of knowledge editing on Rephrase queries by using **RS** (Rephrase Score) and **RM** (Rephrase Magnitude). For Rephrase type queries, RS and RM are actually calculated to derive ES and EM under the condition of rephrasing queries.
- **Specificity** uses **NS** (Neighborhood Score) and **NM** (Neighborhood Magnitude) to measure the ability of knowledge editing to preserve unrelated knowledge. When dealing with OOS queries beyond the editing scope, no editing should take place, and the original facts should be preserved. Therefore, NS is obtained by $E[I[P(o) > P(o^*)]]$, and NM is obtained by $E[P(o) - P(o^*)]$.

A.2 Elaboration and Discussion of Evaluation Framework

While some knowledge-related fields, including Hallucination (Zhang et al., 2023) and Retrieval-Augmented Generation (RAG) (Gao et al., 2024), involve metrics related to fact-checking or validation, such as FactScore (Min et al., 2023) and AlignScore (Zha et al., 2023), it is important to emphasize that Knowledge Editing assessment involves a generated text and two conflicting knowledge references: the pre-editing old knowledge and the post-editing new knowledge, which fundamentally distinguishes the evaluation from metrics in these fields. For INS, the goal is to thoroughly replace old knowledge and introduce new knowledge, whereas for OOS, it is the opposite. This distinction renders the motivation and formulation of the proposed metrics (TE, SE) markedly different from those in other fields, although they may also utilize NLI or Contain function as the basic component.

Additionally, one of the core demands of KE is to maintain locality. Previous works focused solely on whether edited knowledge preserves the pre-

Human Score	Auto Metric	Pearson Correlation
Editing	TE	0.7644
	SE	0.7784
	Editing	0.8074
Retention	TR	0.9195
	SR	0.8868
Overall	Retention	0.9255
	Editing	0.5356
	Retention	0.7612
	Overall	0.839

Table 5: The Pearson correlation coefficient between auto metrics and manual scores. For the auto metrics, Editing is the average of TE and SE; Retention is the average of TR and SR; Overall is the average of Editing and Retention.

vious state for OOS queries, neglecting whether information in other segments of the output remains consistent or is disrupted, which we term as Style Retention/Over-editing. To measure the extent of style retention in edited output compared to the original output, we introduce TR and SR metrics. The design of TR and SR is inspired by the widespread use of N-gram/semantic overlap in the NLP community to measure consistency between generated text and reference text (Papineni et al., 2002; Lin, 2004; Chandrasekaran and Mago, 2021). For INS, we calculate the consistency of the remaining text before and after masking new entities, while for OOS, it is calculated directly.

The rationality of these metrics is validated in Appendix A.4.

A.3 Pseudo-code of Evaluation Framework

We summarize the pseudo-code of our proposed evaluation framework in Algorithm 1.

A.4 Consistency with Human Evaluation

In Section 3.2.2, we proposed a comprehensive evaluation framework, incorporating editing metrics (TE, SE) and retention metrics (TR, SR) to evaluate the quality of output text after knowledge editing. Prior to employing these metrics for evaluation, it was imperative to ensure their validity and necessity. To address this, we sample 300 data points from the test set (comprising Simple, Rephrase, and OOS examples in a 1:1:1 ratio) and enlist human evaluators to independently score them from the perspectives of editing, retention, and overall assessment.

The rules for human scorers scoring the effective-

ness of knowledge editing are as follows: in terms of editing, for INS queries, scoring is as follows: 0 points if there is no editing at all; 0.5 points if there are partial edits, and the sentence still retains old knowledge or exhibits logical inconsistencies; 1 point for perfect knowledge editing with no issues. For OOS queries, the scoring rules are reversed. In the retention aspect, after disregarding content related to the edited knowledge in the sentence, for responses within the editing scope: 0 points for very poor consistency between new and old responses; 0.5 points for ordinary consistency; 1 point for excellent consistency. In the overall aspect, human scorers are required to consider the overall impact of knowledge editing and assign scores within the range of 0, 1, 2, 3, 4 to the edited outputs. Then, we conduct Pearson correlation analyses between these human scores and our automated metrics.

As shown in Table 5, both textual metrics (TE, TR) and semantic metrics (SE, SR) demonstrate commendable consistency scores with human ratings, affirming the effectiveness of the proposed metrics. Moreover, Whether for editing or retention, the consistency score of the joint assessment of textual and semantic dimensions surpasses that of any individual metric. This underscores the necessity of incorporating both textual and semantic metrics in the evaluation process. Finally, the Pearson correlation coefficient between auto editing and human overall score is a mere 0.5356. However, a combined evaluation of editing and retention metrics yield a significantly higher consistency score of 0.839 with human judgments. This suggests that effective alignment with human preferences cannot rely solely on editing scores but requires a comprehensive assessment integrating both editing and retention metrics.

B Details of Method

B.1 Pseudo-code of PostEdit

We summarize the pseudo-code for training post-editor and inference of postEdit in Algorithm 2 and Algorithm 3, respectively.

B.2 Details of Prompts

We demonstrate the two prompt templates T^{aug} and T^{edit} used in the postEdit method as follows:

Prompt Template T^{aug}

For the following query and original response, you need to follow in order: Firstly, locate all spans related to the **old fact: {s} {r} {o}** in original reply; Secondly, modify these spans according to **new fact: {s} {r} {o*}**. Thirdly, output the edited response based on the modified spans (Do not output other content).
 ### The query:
 {x}
 ### Original response:
 {y_o}
 ### Edited response:

Prompt Template T^{edit}

Instruction:
 You will assume the role of an editor. For the following query and original response, if the new fact impacts the query or original response, incorporate the new fact into the original response. If not, simply output the following word: retain.
 ### New fact:
 The answer of {s} {r} has been updated from {o} to {o*}.
 ### The query:
 {x}
 ### Original response:
 {y_o}
 ### Edited response:

C Details of Experiments Setup

C.1 Details of Datasets

In this work, we mainly used two datasets: zsRE and CounterFact.

- **zsRE** (Levy et al., 2017) is one of the most popular question answering (QA) datasets which use question rephrasing as the equivalence neighborhood. These queries of Rephrase type are generated by back-translation. In zsRE, the relationship between entities is associated with a set of crowd-sourced generated questions. Additionally, zsRE associates questions with randomly generated sentences to add out-of-editing scope examples.
- **CounterFact** (Meng et al., 2022a) is a more chal-

Dataset	Data Type	Train Number	Test Number	Length of Original Response (mean/max)
CounterFact	ALL	30000	1500	51.34/436
	Simple	10000	500	50.40/436
	Rephrase	10000	500	53.03/374
	OOS	10000	500	50.59/367
zsRE	ALL	30000	1500	22.39/406
	Simple	10000	500	14.84/119
	Rephrase	10000	500	18.38/257
	OOS	10000	500	33.96/406

Table 6: Statistical information on the sampled datasets.

lenging dataset than zsRE, the expected output of which is contradictory to the fact. It is built to distinguish superficial alterations in the word selections and significant, generalized modifications in its foundational factual knowledge. In CounterFact, the edited answer to the question can sometimes be counterfactual to real world, which makes it harder for the model to predict desired answer and avoid the effects of pre-trained LLMs knowing these desired facts before editing.

Following the previous work (Zheng et al., 2023), for CounterFact, we designate data with edit id numbers ranging from 0 to 2000 as the test set for knowledge edit, while the remaining data constitute the training set. As we adopt ChatGPT as our base LLM in main experiments, in order to control the dataset size, we randomly sampled 30,000 examples (10,000 each for Simple, Rephrase, and OOS) from the original training set. These samples constitute our training set. Additionally, we randomly selected 1,500 examples (500 each for Simple, Rephrase, and OOS) from the original test set to create our query test set. The original response for INS test queries are ensured to hit the old knowledge object before editing, and the OOS are ensured to have no wrong knowledge before editing. We present the statistical information of the datasets after sampling in Table 6, and show a training sample and test sample from zsRE respectively as follows:

Sample From zsRE Training Set

```
{
  "edit_id": 15000,
  "edit": "Denis Dyack » Denys de La Tour || Who is the designer of Too Human?",
  "query": "Who is the designer from Too Human?",
  "query_type": "rephrase",
  "original_response_by_gpt3.5": "The
```

```
designer of Too Human is Denis Dyack.",
  "edited_response_by_gpt4": "The designer of Too Human is Denys de La Tour."
}
```

Sample From zsRE Test Set

```
{
  "edit_id": 70,
  "edit": "Serpens » Andromeda || Which constellation is NGC 6604 in?",
  "query": "Which constellation does NGC 6604 belong to?",
  "query_type": "rephrase",
  "original_response": "NGC 6604 belongs to the constellation of Serpens."
}
```

C.2 Details of Baselines

- **IKE** (Zheng et al., 2023) is a method of knowledge editing that does not involve modifying the parameters of LLMs. It defines three types of demonstration formatting templates including copy, update, and retain. These templates serve distinct functions and act as guiding principles for the language model, enabling it to edit knowledge through in-context learning, allowing IKE to maintain both efficiency and excellent generalization and specificity. This opens up the possibility of employing IKE for the task of knowledge editing even in scenarios involving black-box models.
- **PROMPT** (Zheng et al., 2023) is similar to IKE, as a method of knowledge editing through in-context learning. However, unlike IKE, PROMPT doesn't require constructing three types of demonstrations but directly provides new knowledge to the LLM for knowledge editing.
- **SERAC** (Mitchell et al., 2022) is a memory-

based method of knowledge editing. This method stores edits in explicit memory and learns to reason about these edits as needed to adjust the predictions of the base LLM without modifying parameters. SERAC uses an explicit cache of user-provided edit descriptors, alongside a scope classifier and surrogate model. When presented with a query, SERAC uses the scope classifier to determine if the query falls within the editing scope. If it does, the output is predicted via the surrogate model; otherwise, it defers to the base LLM for the output.

- **SERAC (ChatGPT)** In SERAC, the surrogate model is obtained by fine-tuning a smaller language model compared to the base LLM. We utilize ChatGPT as the surrogate model to derive a SERAC variant that requires no additional training.

C.3 Details of Implementation

As described in Section 3.2.2, our evaluation framework employs a NLI model for computing SE, ROUGE scores for computing TR, and a SBERT model for computing SR. In details, SE utilizes albert-xxlarge-v2-snli_mnli_fever_anli_R1_R2_R3-nli¹⁰ as the NLI model; ROUGE score is implemented through the rouge library¹¹, using the F1 score of ROUGE-1; SR uses all-MiniLM-L6-v2¹² as the SBERT model.

For training of post-editor, we employ ChatGPT (gpt-3.5-turbo-0301) for original response augment and GPT-4 (gpt-4-0613) for edited response augment¹³, with the default temperature coefficient ($t = 0.1$). In order to enhance training efficiency and reduce the number of updated parameters, we adopt the LoRA strategy (Hu et al., 2021) to finetune LLaMA 2-7B. Specifically, the rank of LoRA is set to 8, with *lora_alpha* at 16 and *lora_dropout* at 0.05. The LoRA update matrix is applied to the self-attention and FFN layers, with *target_modules* as ["q_proj", "k_proj", "v_proj", "o_proj", "gate_proj", "down_proj", "up_proj"]. We train 5 epochs to optimize post-editor, employing a batch size of 128 and a learning rate of 5e-2. We also use the warmup and cosine annealing strategy, with a warmup ratio

¹⁰https://huggingface.co/ynie/albert-xxlarge-v2-snli_mnli_fever_anli_R1_R2_R3-nli

¹¹<https://pypi.org/project/rouge>

¹²<https://huggingface.co/sentence-transformers/all-MiniLM-L6-v2>

¹³<https://platform.openai.com/docs/models>

of 0.1 and the Adam optimizer (Kingma and Ba, 2017).

For retriever of postEdit, consistent with all baselines, we use all-MiniLM-L6-v2 to encode queries and edit knowledge, while employing dot product as the similarity function. For base LLM, we use ChatGPT (gpt-3.5-turbo-0301) in main experiments, with a temperature coefficient of 0.1. During inference of post-editor, we set the temperature coefficient of 0.1 and use beam search to decode the output, where *num_beams* is set to 4. To further improve the inference speed, we apply 8-bit quantization when loading post-editor.

In terms of baselines, for SERAC, we fine-tune the surrogate model using the same LLAMA2-7B as post-editor and the similarity discrimination threshold is set at 0.7, determined through hyperparameter search on the training set (ranging from 0.1 to 0.9 with a step size of 0.1). To better maintain consistency between baselines and postEdit implementations, we adopt training output targets consistent with postEdit for the surrogate model of SERAC, i.e., GPT-4 augmented edited response, rather than new objects of editing knowledge, aiming to achieve higher stylistic retention. For IKE, we set the number of demonstration examples to 32. The rest of the hyperparameter settings for the baselines follow the default configurations in their original papers. All experiments use a single Nvidia A100 GPU (80 GB of memory).

D More Experiments

D.1 Comparison with more Baselines

In Section 5, we compared methods that have the same scenario as postEdit. In this section, we transfer some methods from other task scenarios as baselines to further enrich the experiments:

- MeLLO (Zhong et al., 2023) is a method specifically designed for multi-hop reasoning scenarios in knowledge editing, storing edited facts externally and iteratively prompts LLMs to generate answers consistent with the edited facts.
- RARR (Gao et al., 2023) aims to reduce hallucinations in LLM outputs by scrutinizing and revising. It initially uses search engines for evidence and attribution, then corrects unsupported content while preserving the original output, achieved through few-shot demonstrations. We replace the search engine with edit memory.
- In addition to PROMPT and IKE, similar to the

Method	Textual Editing (TE)				Semantic Editing (SE)				Textual Retention (TR)				Semantic Retention (SR)			
	Simple	Rephrase	OOS	AVG (HM)	Simple	Rephrase	OOS	AVG (HM)	Simple	Rephrase	OOS	AVG (HM)	Simple	Rephrase	OOS	AVG (HM)
MeLLO	42.42	32.87	37.07	37.55 ^(37.05)	43.61	35.11	44.3	41.11 ^(40.55)	16.42	11.22	15.59	14.47 ^(14.01)	38.5	31.61	41.58	37.32 ^(36.74)
RARR	53.9	49.47	85.67	63.17 ^(59.48)	55.9	50.96	86.48	64.6 ^(61.13)	54.18	54.9	63.19	57.44 ^(57.15)	62	62.98	71.13	65.39 ^(65.12)
RAG-8shot	99.7	99.79	9.35	69.32 ^(23.62)	98.9	95.64	11.79	68.54 ^(28.47)	26.2	23.98	4.57	18.21 ^(10.04)	55.32	53.5	25.01	44.54 ^(39.09)
postEdit (ours)	96.8	94.7	99.4	96.97 ^(96.93)	92.5	92.1	99.4	94.67 ^(94.55)	88.65	89.66	99.64	92.65 ^(92.39)	93.9	94.02	99.82	95.91 ^(95.84)

Table 7: Performance comparison on CounterFact.

conventional RAG approach, we utilize few-shot <query, edit, edited output> prompts to enhance the base LLM’s utilization of editing knowledge, where all demonstration samples belong to the INS type, referred to as RAG-8shot.¹⁴

The results are shown in Table 7. Overall, postEdit still outperforms all baselines. We can further observe that: Firstly, since MeLLO and RARR are not designed specifically for general knowledge editing scenarios, they perform poorly on CounterFact. Secondly, leveraging the impressive in-context learning capabilities of ChatGPT, RAG-8shot achieves near-perfect INS Editing scores, but faces significant challenges on OOS Editing due to the lack of OOS demonstrations. This emphasizes the need for a INS/OOS judgment mechanism on top of RAG. Lastly, post-processing methods (postEdit, RARR) achieve higher Retention scores compared to pre-processing methods (MeLLO, RAG-8shot), highlighting the advantage of post-processing for style retention.

D.2 Does Post-editor just Remember the Patterns of Training Data for Testing?

In the experiment setup of KE, the edits in the training set and the test set are completely non-overlapping. Therefore, the post-editor can not rely on edits seen during training for testing. However, another risk of overfitting to the training data occurs when post-editor directly memorize patterns of INS and OOS data rather than making judgments based on recalled edits. To address this, we test the performance of postEdit when the edit memory is empty. As shown in Table 8, when edit memory is empty, post-editor tends to classify queries as OOS type, leading to nearly 100% OOS editing scores and nearly 0% INS (Simple and Rephrase) editing scores. This demonstrates that post-editor relies on edit knowledge guidance for INS/OOS judgment and revisions, rather than memorizing patterns from the training data.

¹⁴Since in the standard KE experimental setup, the size of edit memory is set to 1, serving as an "oracle" retrieval setting to encourage methods to focus more on editing and locality capabilities. Therefore, we don’t compare with some RAG methods that focus on improving retrieval recall.

Types	CounterFact		zsRE	
	TE	SE	TE	SE
Simple	0.0	0.0	0.0	0.67
Rephrase	0.0	0.0	0.0	0.33
OOS	100.0	98.59	100.0	100.0
AVG	33.33	32.86	33.33	33.67

Table 8: Test results for CounterFact and zsRE when Edit Memory is empty. We simulate this scenario by replacing the recalled edit with an empty string "".

E Discussion on Efficiency

Apart from Editing and Retention performance, KE methods should strive to minimize storage and computational costs. For memory-based black-box LLM editing, in addition to Edit Memory and the retriever, storage overhead also encompasses the demonstration library for IKE, the judge model and surrogate model for SERAC, and the post-editor for postEdit. Furthermore, although memory-based methods do not incur computational overhead for editing, they do introduce inference expenses. Specifically, for IKE, the inference cost increases from $f_{base}(x)$ to $f_{retr}(x, M_e) + f_{base}(demos, e, x)$; for SERAC, the additional cost is $f_{retr}(x, M_e) + f_{judge}(x, e_{retr})$; and for postEdit, it is $f_{retr}(x, M_e) + f_{edit}(e, x, y_o)$. To further reduce post-editing overhead, one approach is to improve the reasoning efficiency of the post-editor. As highlighted in Section 6.5, a small-scale post-editor can also achieve commendable performance. Another potential option is to employ white-box parameter-editing methods to directly integrate editing knowledge into the post-editor. The post-editor can then use its knowledge to modify the original response of base LLM, exchanging editing costs for memory storage and retrieval expenses.

Algorithm 1: Pseudo-code of Evaluation Framework in a Python-like style.

```
# x: the input of LLM (All text is processed in lowercase, the same below.)
# x_label: "INS" if x in editing scope else "OOS"
# y_o, y_e: the original and edited output of LLM
# o_old, o_new: the object of old knowledge  $t$  and new knowledge  $t^*$  for editing
# k_old, k_new: text format of  $t$  and  $t^*$ 
# k_self: text format of LLM's self-knowledge  $t_o$  and is equivalent to  $[x, y_o]$ 
# func_entail(a,b): return True if a entails b else False by using a NLI model
# func_rouge(a,b): return the ROUGE score of a and b
# func_sim(a,b): return the similarity of a and b using a SBERT model

def TE(y_e, x_label, o_old, o_new):
    ctn_old=1 if o_old in y_e else 0
    ctn_new=1 if o_new in y_e else 0
    if x_label=="INS":
        TE_score=0.5*ctn_new + 0.5*(1-ctn_old)
    else:
        TE_score=0.5*ctn_old + 0.5*(1-ctn_new)
    return TE_score

def SE(x_label, x, y_e, k_old, k_new, k_self, func_entail):
    ent_new=1 if func_entail(x+" "+y_e,k_new) else 0
    if x_label=="INS":
        ent_old=1 if func_entail(x+" "+y_e,k_old) else 0
        SE_score=0.5 * ent_new + 0.5 * (1-ent_old)
    else:
        ent_old=1 if func_entail(x+" "+y_e,k_self) else 0
        SE_score=0.5*ent_old + 0.5*(1-ent_new)
    return SE_score

def TR(x_label, y_o, y_e, o_old, o_new, func_rouge):
    if x_label=="INS":
        TR_score=func_rouge(y_o.replace(o_old,"mask"), y_e.replace(o_new,"mask"))
    else:
        TR_score=func_rouge(y_o,y_e)
    return TR_score

def SR(x_label, y_o, y_e, o_old, o_new, func_sim):
    if x_label=="INS":
        SR_score=func_sim(y_o.replace(o_old,"mask"), y_e.replace(o_new,"mask"))
    else:
        SR_score=func_sim(y_o,y_e)
    return SR_score
```

Algorithm 2: Train post-editor

Data: training dataset $D_{train} = \{(e_i, x_i)\}$

Require: base LLM f_{base} , GPT-4 f_{gpt4} , trainable generative model f_{edit} , training epoch \mathbf{E} , batch size \mathbf{B}

for i in $1, \dots, |D_{train}|$ **do**

$y_{i,o}^{aug} = f_{base}(x_i)$

▷Original Response Augment

if $x_i \in \mathcal{X}_e$ **then**

$y_{i,e}^{aug} = f_{gpt4}(T^{aug}(e_i, x_i, y_{i,o}^{aug}))$

▷Edited Response Augment

if $TE(y_{i,e}^{aug}) \neq 1$ or $SE(y_{i,e}^{aug}) \neq 1$ **then**

delete $(e_i, x_i, y_{i,o}^{aug}, y_{i,e}^{aug})$

end

else

$y_{i,e}^{aug} = \langle Retain \rangle$

end

end

$D_{train}^{aug} = \{(e_i, x_i, y_{i,o}^{aug}, y_{i,e}^{aug})\}$

for epoch in $1, \dots, \mathbf{E}$ **do**

for iter=0, 1, 2, \dots **do**

 sample a mini-batch \mathbf{B} from D_{train}^{aug}

▷Supervised Fine-tuning

 compute \mathcal{L}_{sft} by equation 6 and optimize f_{edit}

end

end

Output: trained post-editor f_{edit}

Algorithm 3: Inference of PostEdit

Input: use query x

Require: Edit Memory M_e , base LLM f_{base} , post-editor f_{edit} , SBERT retriever f_{retr}

get original response: $y_o = f_{base}(x)$

retrieve the most similar edit index: $i^* = \operatorname{argmax}_{0 \leq i < |M_e|} \operatorname{sim}(x, e_i)$

get post-editor's output: $f_{edit}(x_{edit}) = f_{edit}(T^{edit}(e_{i^*}, x, y_o))$

if $f_{edit}(x_{edit}) \neq \langle Retain \rangle$ **then**

$y_e = f_{edit}(x_{edit})$

else

$y_e = y_o$

end

Output: final response y_e
