

Editing the Mind of Giants: An In-Depth Exploration of Pitfalls of Knowledge Editing in Large Language Models

Anonymous ACL submission

Abstract

Knowledge editing is a rising technique for efficiently updating factual knowledge in large language models (LLMs) with minimal alteration of parameters. However, recent studies have identified concerning side effects, such as knowledge distortion and the deterioration of general abilities, that have emerged after editing. This paper conducts a comprehensive study of these side effects, providing a unified view of the challenges associated with knowledge editing in LLMs. We discuss related work and summarize potential research directions to overcome these limitations. Our experiments highlight the limitations of current knowledge editing methods, emphasizing the need for deeper understanding of inner knowledge structures of LLMs and improved knowledge editing methods.

1 Introduction

Recent advancements in large language models (LLMs) have significantly improved NLP applications, enabling LLMs to understand and generate language at a human-like level. However, the mechanisms of knowledge storage in LLMs remain unclear, raising concerns about the reliability of their output, particularly in applications like chatbots. To address these issues, researchers have explored various methods. Traditional methods like fine-tuning, continual learning, and retraining are computationally expensive and may degrade LLM performance. *Knowledge editing* has emerged as a promising alternative, offering efficient adjustments with minimal computational costs and fewer alterations (Cao et al., 2021; Dai et al., 2022; Meng et al., 2022, 2023; Dong et al., 2022; Mitchell et al., 2022a,b; Hartvigsen et al., 2023; Huang et al., 2023; Yu et al., 2024; Zheng et al., 2023; Li et al., 2023; Tan et al., 2024; Gupta et al., 2024b; Wang et al., 2024). This method allows precise LLMs refinement, enhancing their practical and reliable use in real-world

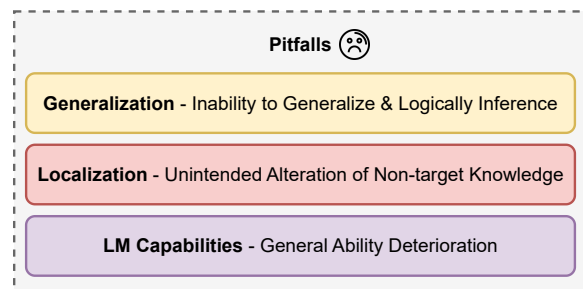


Figure 1: An overview of pitfalls in current knowledge editing methods. The subsequent sections dive into three key challenges: generalization issues (Section 3.1), locality issues (Section 3.2), and deterioration of general LLM abilities (Section 3.3).

applications.

Knowledge editing can be divided into two main categories: parameter-modifying and parameter-preserving. Both aim to refine LLM knowledge efficiently while avoiding the drawbacks of previous tuning methods (Yao et al., 2023). Parameter-modifying methods, including meta-learning (Cao et al., 2021; Mitchell et al., 2022a; Tan et al., 2024) and locate-and-edit techniques (Dai et al., 2022; Meng et al., 2022, 2023; Li et al., 2023; Gupta et al., 2024b), strive to update model parameters effectively. By contrast, parameter-preserving methods introduce external components, like knowledge bases (Mitchell et al., 2022b; Zhong et al., 2023) or extra model parameters (Dong et al., 2022; Huang et al., 2023; Hartvigsen et al., 2023; Yu et al., 2024) to maintain the integrity of pre-trained LLMs while updating their knowledge.

Despite the success of knowledge editing, challenges remain. Knowledge editing can have unintended side effects, potentially damaging the general abilities and intrinsic structures of LLMs. Previous research has mainly focused on performance

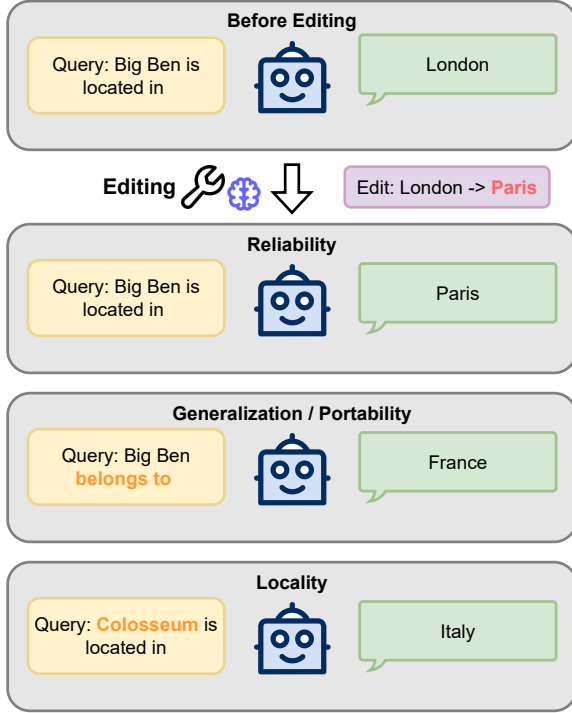


Figure 2: Illustration of properties that knowledge editing methods should satisfy: reliability, generalizability/portability, and locality.

improvements and innovations within knowledge editing methods, with limited attention to potential drawbacks. Consequently, this survey aims to provide a holistic view of current issues in the knowledge editing paradigm and encourage further investigations into the pitfalls and intrinsic knowledge structures of LLMs. A brief overview of the discussed pitfalls is shown in Figure 1.

This paper is organized as follows: Section 2 introduces the definition and methods of knowledge editing. Section 3 discusses current challenges and corresponding benchmarks. In Section 4, we present experimental results evaluating different editing methods. Finally, Section 5 explores related studies and future research directions. We summarize our contributions as follows:

1. We are the first to provide a comprehensive analysis of the side effects associated with existing knowledge editing techniques.
2. We systematically organize previous research and conduct experiments to benchmark the side effects of knowledge editing, providing a unified perspective on this issue.
3. We discuss related studies and potential directions to address existing challenges, encouraging further exploration in this field.

2 Overview of Knowledge Editing

2.1 Problem Definition

Knowledge editing for LLMs entails modifying the output of LLMs in response to specific edit queries, with the aim of minimizing alterations to their original behavior (Yao et al., 2023; Mazzia et al., 2023; Zhang et al., 2024a). In this section, we follow the notation from Mazzia et al. (2023).

We denote the input and output space as \mathbb{X} and \mathbb{Y} , respectively. The function space $\mathbb{F} : \mathbb{X} \rightarrow \mathbb{Y}$ is estimated by the base model f_{θ_0} parameterized by $\theta_0 \in \Theta$. Finally, let $Z_e = \{(x_e, y_e) \mid f_{\theta_0}(x_e) \neq y_e\}$ be the set of edit queries we would like to apply to the base model. The goal of knowledge editing is to efficiently derive the edited model f_{θ_e} from the base model that satisfies the following:

$$f_{\theta_e}(x_e) = y_e, \forall (x_e, y_e) \in Z_e \quad (1)$$

The ideal edited model f_{θ_e} should satisfy three properties: **reliability**, **generalization**, and **locality**. An illustration is shown in Figure 2.

Reliability Given an edit query (x_e, y_e) , the edited model f_{θ_e} should output the target answer y_e when given the target input x_e , i.e. $f_{\theta_e}(x_e) = y_e$. The reliability of a editing method is measured by calculating the average edit success rate:

$$\mathbb{E}_{(x'_e, y'_e) \sim Z_e} \mathbb{1}\{f_{\theta_e}(x'_e) = y'_e\} \quad (2)$$

Generalization The edited model should generalize the edited knowledge to relevant instances. The generalization metric is commonly formulated as the average success rate on the neighboring set:

$$\mathbb{E}_{(x'_e, y'_e) \sim N(x_e, y_e)} \mathbb{1}\{f_{\theta_e}(x'_e) = y'_e\}, \quad (3)$$

where $N(x_e, y_e)$ is the set of neighboring instances of an edit query (x_e, y_e) . Earlier works evaluate this metric by rephrasing the input prompts (Mitchell et al., 2022a; Meng et al., 2022; Huang et al., 2023).

Locality The editing process should not affect instances unrelated to the edit queries. The locality set of an edit query (x_e, y_e) can be defined as $L(x_e) = \{(x_{loc}, y_{loc}) \in \mathbb{X} \times \mathbb{Y} \text{ s.t. } x_{loc} \notin N(x_e, y_e) \wedge f_{\theta_0}(x_{loc}) = y_{loc}\}$. The locality, also known as specificity, of an editing method is measured by calculating the level of invariance of model output before and after the edits, which can be calculated as follows:

$$\mathbb{E}_{(x_{loc}, y_{loc}) \sim L(x_e)} \mathbb{1}\{f_{\theta_e}(x_{loc}) = y_{loc}\} \quad (4)$$

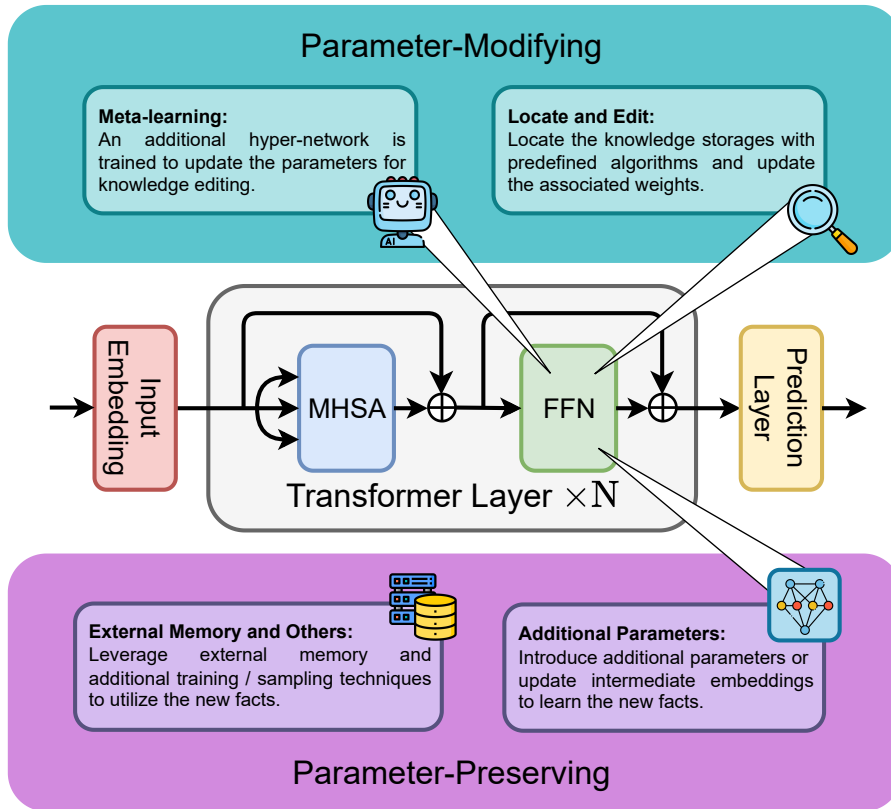


Figure 3: Illustration of the two categories of model editing methods in transformer-based large language models, which includes parameter-modifying (meta-learning and locate-and-edit) and parameter-preserving (additional parameters, external memory, in-context learning, and decoding) methods. MHSA and FFN stand for multi-head self-attention and feed-forward network, respectively.

2.2 Current Methods

Current knowledge editing methods are categorized into parameter-modifying (Section 2.2.1) and parameter-preserving (Section 2.2.2) editing methods, each containing several strategies. An overview and illustration of current methods are included in Table 1 and Figure 3, respectively.

2.2.1 Parameter-Modifying

Meta-learning Meta-learning methods train a hyper-network to predict network parameter updates. For instance, KnowledgeEditor (Cao et al., 2021) trains a deep network to predict weight updates. MEND (Mitchell et al., 2022a) decomposes the gradient matrix into two rank-one matrices and utilized a hyper-network to update these matrices, thereby accelerating the editing process. Built upon MEND, MALMEN (Tan et al., 2024) refines the process by formulating the aggregation of parameter shifts into a least-squares problem, further improving the scalability of meta-learning methods.

Locate and Edit Locate-and-edit methods identify specific knowledge locations in LLMs for con-

sequent editing. KN (Dai et al., 2022) utilizes the proposed knowledge attribution method to pinpoint neurons expressing relational facts, allowing efficient updates or erasures without fine-tuning. ROME (Meng et al., 2022) proposes causal tracing method to identify neuron activations linked to specific knowledge. The authors demonstrate the significance of middle-layer feed-forward networks (FFNs) in factual predictions when processing the subject’s last token. Built upon the hypothesis that the FFN modules in a transformer layer can be viewed as key-value memories (Geva et al., 2021), ROME injects new knowledge into the key-value memories by deriving the closed form solution from the least-squares problem. MEMIT (Meng et al., 2023) scales up ROME by editing a set of MLPs from consecutive middle-layers via solving a normal equation. PMET (Li et al., 2023) proposes to update multi-head self-attention (MHSA) modules in addition to FFNs. EMMET (Gupta et al., 2024b) on the other hand, integrates the objectives of ROME and MEMIT into a unified preservation-memorization objective, facilitating batch-editing

Category	Strategy	Method
Parameter-modifying	Meta-learning	Knowledge Editor (Cao et al., 2021) MEND (Mitchell et al., 2022a) MALMEN (Tan et al., 2024)
	Locating and editing	Knowledge Neuron (Dai et al., 2022) ROME (Meng et al., 2022) MEMIT (Meng et al., 2023) PMET (Li et al., 2023) EMMET (Gupta et al., 2024b)
Parameter-preserving	Additional parameters	CaliNET (Dong et al., 2022) T-Patcher [†] (Huang et al., 2023) GRACE [†] (Hartvigsen et al., 2023) MELO [†] (Yu et al., 2024)
	External memory	SERAC [†] (Mitchell et al., 2022b) MeLLO [†] (Zhong et al., 2023)
	In-context learning	IKE [†] (Zhong et al., 2023)
	Decoding	DeepEdit [†] (Wang et al., 2024)

Table 1: Overview of knowledge editing methods. The methods are categorized into two major families, parameter-modifying and parameter-preserving methods, each containing several strategies. Methods marked with [†] have the ability to process sequential edits.

capabilities for both methodologies.

2.2.2 Parameter-Preserving

Additional Parameters Some methods utilize additional parameters, such as adding new neurons or employing parameter-efficient techniques. CaliNET (Dong et al., 2022) extends the FFN modules with calibration memory slots to adjust the predicted token distribution. T-Patcher (Huang et al., 2023) adds neurons in the FFN’s last layer to rectify classification errors and incorrectly generated tokens, activating only in response to associated mistakes. GRACE (Hartvigsen et al., 2023) wraps a selected layer with an Adaptor that includes a codebook and deferral mechanism, learning to decode desired outputs while caching embeddings of error inputs. The GRACE layer stores the edits and could be updated continuously over long deployments. MELO (Yu et al., 2024) utilizes Dy-LoRA (Valipour et al., 2023) modules to learn edits, indexing them in an inner vector database to dynamically activate corresponding LoRA blocks during inference.

External Memory Other methods utilize external memories for editing. SERAC (Mitchell et al., 2022b) leverages a scope classifier to determine whether an user-supplied edit example stored in its memory is related to the inputs. If no example exists, the inputs are passed to the base model; otherwise, a counterfactual model generates modified answers using the inputs and the related example.

MeLLO (Zhong et al., 2023) decomposes a multi-hop question into subquestions iteratively. The model then checks if the tentative answer generated by the base model contradicts the most relevant facts retrieved from the edited fact memory and adjusts the outputs accordingly.

In-Context Learning and Decoding Certain strategies require no additional parameters. IKE (Zheng et al., 2023) edits factual knowledge via in-context learning with demonstrations to guide the language model. DeepEdit (Wang et al., 2024) employs decoding constraints, including filtering step candidates, depth-first search to store valid candidates in a stack, and a greedy search to output the optimal path for multi-hop reasoning.

3 Challenges of Knowledge Editing

While knowledge editing methods have been extensively researched, comprehensive studies on related challenges are lacking. In this section, we discuss the pitfalls of knowledge editing from three perspectives: inability to logically infer and robustly generalize (Section 3.1), unintended alteration of non-target knowledge (Section 3.2), and deterioration of general LLM abilities (Section 3.3).

3.1 Inability to Logically Inference and Robustly Generalize

When a fact is updated, it is crucial not only to revise the specific piece of knowledge but also to evaluate the impact on the related reasoning chain. Recently the term **portability** has been proposed in (Yao et al., 2023) to evaluate whether an edited fact can be logically inferred within the knowledge chain, and to further assess the robustness of generalization. In their study, they introduce three metrics to evaluate portability: Subject Replace (checking if synonyms of the subject are edited), Reversed Relation (checking if the reversed relation of the target is edited), and One Hop (assessing if modified knowledge is usable for further derivation). Similarly, RippleEdits benchmark as well as corresponding Logical Generalization and Compositionality metrics are proposed to examine whether edited knowledge can be inferred in composite relations of facts (Cohen et al., 2023). Additionally, ReCoE benchmark is proposed to assess the propagation of updates in interconnected facts using various reasoning schemes in complex question-answering datasets (Hua et al., 2024). Furthermore, MQuAKE benchmark is introduced to

Challenge	Benchmark	Metric
Portability and Generalization	RippleEdits (Cohen et al., 2023)	Logical Generalization, Compositionality I, Compositionality II
	ConflictEdit (Li et al., 2024)	Conflict Score, Conflict Magnitude, Success Score
	MQuAKE (Zhong et al., 2023)	Edit-wise Success Rate, Instance-wise Accuracy, Multi-hop Accuracy
	ReCoE (Hua et al., 2024)	QA Accuracy
	ZsRE + CounterFact [†] (Yao et al., 2023)	Subject-Replace, Reverse-Relation, One-Hop
Locality	RippleEdits (Cohen et al., 2023)	Subject Aliasing, Preservation, Relation Specificity
	RoundEdit (Li et al., 2024)	Success Score, Distortion (↓), Ignore Rate (↓), Failure Rate (↓), Tied Fact Damage (↓)
	ZsRE + CounterFact [†] (Yao et al., 2023)	Other-Attribution, Distract-Neighbor, Other-Task
	CounterFact (Meng et al., 2022)	Locality, Neighborhood Score, Neighborhood Magnitude
	CounterFact+ (Hoelscher-Obermaier et al., 2023)	Neighborhood KL Divergence

Table 2: Performance benchmarks and evaluation metrics addressing generalization/portability and locality issues in knowledge editing methods. Unless specifically indicated by a downward arrow, higher values signify better performance in those evaluation metrics. CounterFact benchmark is proposed by (Meng et al., 2022), and CounterFact with [†] mark is modified by (Yao et al., 2023) to further examine the proposed metrics.

evaluate more complex reasoning and inference ability on multi-hop questions (Zhong et al., 2023).

When editing multiple logically related facts simultaneously, models may suffer from confusion due to conflicts. ConflictEdit benchmark is proposed to examine different editing methods on conflicted edit facts (Li et al., 2024). The different benchmarks and corresponding metrics and are arranged systematically in Table 2.

3.2 Unintended Alteration of Non-Target Knowledge

Locality is conventionally assessed using a locality dataset to evaluate the impact of edits on unrelated facts by measuring the Neighborhood Score and Neighborhood Magnitude (NS & NM; Meng et al., 2022, 2023). However, current evaluation methods do not adequately capture the post-edit effects on content beyond the locality dataset, which means the edited model could still contain unintended alterations. For example, while the location of the Louvre is successfully modified from Paris to London, the edited model might also output London in an unrelated context or increase the probability of words semantically related to London (e.g., Big Ben) when mentioning the Louvre. Some modified benchmark (CounterFact+) and corresponding metric (Neighborhood KL Divergence) (Hoelscher-Obermaier et al., 2023) is then designed to disclose these previously implicit pitfalls. Another

study (Yao et al., 2023) extends this exploration to three facets of locality: Other Relations (evaluating the retention of other attributes of the updated subject), Distract Neighborhood (assessing whether model will be swayed by edited cases when they are concatenated before unrelated inputs), and Other Tasks (examining the influence of edits on the performance of other tasks).

Unintended edits to unrelated facts may occur because a single edit can implicitly change the predictive distribution among objects associated with the same (*subject - relation*) pair. After multiple consecutive edits, these alterations can accumulate and distort the stored knowledge. To evaluate this condition, the concept of Knowledge Distortion has been introduced by Li et al. (2024), which estimates the Jensen–Shannon divergence of the object set distribution before and after editing. This can be further extended to metrics such as the Ignore Rate, measuring how objects other than the target in the object set are neglected after editing, and the Failure Rate, which measures the proportion of instances where over half of the objects in the set are overlooked.

3.3 Deterioration of General LLM Abilities

Current evaluation metrics are primarily limited to scenarios where editing is performed only once or infrequently, prompting some studies to extend evaluations to the outcomes after consecutive edits.

A study by Gupta et al. (2024a) discovers that post-edit models exhibit susceptibility to both gradual forgetting and catastrophic forgetting in sequential editing scenarios. Notably, their findings indicate that the extent of knowledge forgetting is more pronounced in meta-learning-based methods compared to locate-and-edit methods. Additionally, models with parameters modified successively show a decline in performance across various downstream NLP tasks (Gu et al., 2024). Furthermore, perplexity is found to increase after consecutive edits across all parameter-modified methods and different LLMs, and is proposed as another metric to indicate model collapse (Yang et al., 2024). These findings further corroborate that model editing aimed at modifying parameters adversely affects the general capabilities of the original LLMs.

4 Experiments

The experiments are done to evaluate robust generalization and locality (Section 4.1.1 as well as deterioration of general LLM abilities (Section 4.1.2 across different editing methods.

4.1 Experimental Setup

4.1.1 Robust Generalization and Locality

We use GPT-J (Wang and Komatsuzaki, 2021) as the baseline model for editing and implement six distinct editing methodologies to assess robust generalization and locality: MEND (meta-learning), ROME and MEMIT (locate-and-edit), SERAC (external memory), and IKE (prompting).

Given the overlap in benchmarks for robust generalization and locality, we select a subset for our experiments. Robust generalization is evaluated in single edit (modifying a single fact) and multiple edit (altering multiple logically interconnected facts) settings. Single edit metrics include Subject-Replace, Reverse-Replace, and One-Hop reasoning (Yao et al., 2023). Multiple edit metrics include multi-hop editing accuracy (Zhong et al., 2023), and Conflict Score and Conflict Magnitude for Reverse Conflict and Composite Conflict respectively (Li et al., 2024). For locality, single edit metrics include Other-Attribution, Distract-Neighbor, and Other-Task (Yao et al., 2023), while multiple edit metrics encompass Success Rate, Distortion, Ignore Rate, and Failure Rate (Li et al., 2024).

4.1.2 Deterioration of General LLM Abilities

Following the settings of (Gu et al., 2024), we assess deterioration of general LLM abilities post-

editing using six methodologies: ROME, MEMIT, SERAC, MEND, KN, and GRACE. We evaluate general abilities across four NLP downstream tasks: open-domain question answering, sentiment analysis, reasoning, and summarization. These tasks are assessed after 10 to 40 edits on the Zero-Shot Relation Extraction (ZsRE) dataset (Levy et al., 2017), comparing the results against pre-editing benchmarks. More details on the selected downstream tasks are in Appendix B.

4.2 Experimental Results and Discussion

In general, current editing methodologies show sub-optimal performance in both robust generalization and locality. Regarding robust generalization (Table 3), IKE, which leverages prompt demonstrations, excels in single edit but declines with multiple edits. This suggests that prompt demonstrations may become confused when editing multiple logically related facts. Conversely, fine-tuning and meta-learning-based methods are less susceptible to confusion after editing multiple related facts.

Regarding locality (Table 4), IKE maintains stable performance across metrics in single edit settings. Parameter-modifying methods excel in Other Attribution but decline in other metrics, except MEMIT, which remains stable across all metrics. In multiple edit scenarios, all methods except SERAC show similar performance. In the multiple edit scenario, all methods except SERAC exhibit relatively similar performance. SERAC displays low edit success rate and distortion rate, suggesting its scope classifier does not adopt most edits in this scenario. This may be attributed to its weakness in recovering edited facts, which is crucial in this metric setting.

In terms of general LLM abilities (Figure 4), the number of edits affects methods differently. Meta-learning methods like MEND degrade significantly after 10-20 edits. Locate-and-edit methods such as ROME and KN degrade after 10 edits, while MEMIT remains stable after 40 edits. This disparity can be attributed to MEMIT’s strategy of adjusting parameters across multiple layers, as opposed to ROME’s single-layer edits and KN’s approach of modifying a few neurons. This distribution of parameter modifications across layers may help mitigate deterioration.

GRACE, which stores edited facts with additional parameters, shows no performance change in downstream tasks after edits. One possible explanation is that the edits are conducted on the

Methods	Single Edit			Multiple Edit				
	One-Hop			Multiple-Hop	Reverse Conflict		Composite Conflict	
	SR	RR	OH		CS	CM	CS	CM
FT	72.96	8.05	1.34	1.6	80.28	71.11	75.45	64.28
MEND	42.45	0.00	11.34	9.2	88.89	60.50	84.85	43.45
ROME	37.42	46.42	50.91	7.6	65.92	-0.65	71.70	37.04
MEMIT	27.73	47.67	52.74	8.1	51.40	-1.60	57.15	-1.50
SERAC	17.79	1.30	5.53	7.9 [†]	50.89 [†]	-0.02 [†]	50.84 [†]	-0.02 [†]
IKE	88.77	92.96	55.38	8.3 [†]	58.20 [†]	-1.00 [†]	50.52 [†]	-0.99 [†]

Table 3: Experimental results for portability and generalization. SR: Subject-Replace, RR: Reverse-Relation, OH: One-Hop Accuracy, MH: Multi-hop Accuracy, CS: Conflict score, CM: Conflict magnitude. Higher values indicate better performance for all metrics in this table. Results marked with † are obtained in our own experiments, and other results are taken from previous studies.

Methods	Single Edit			Multiple Edit			
	OA	DN	OT	Succ.	D (↓)	IR (↓)	FR (↓)
FT	12.88	9.48	49.56	100.0	16.12	97.48	97.32
MEND	73.50	32.96	48.86	99.12	14.35	87.64	86.56
ROME	78.94	50.35	52.12	99.80	13.95	78.98	77.60
MEMIT	86.78	60.47	74.62	99.72	13.50	72.03	70.44
SERAC	99.50	39.18	74.84	50.14 [†]	3.78 [†]	99.62 [†]	99.64 [†]
IKE	84.13	66.04	75.33	100.0 [†]	13.43 [†]	73.53 [†]	73.00 [†]

Table 4: Experimental results for locality. OA: Other-Attribution, DN: Distract-Neighbor, OT: Other-Task, Succ.: Success rate, D: Distortion, IR: Ignore rate, FR: Failure rate. Unless specifically indicated by a downward arrow, higher values signify better performance in those evaluation metrics. Results marked with † are obtained in our own experiments, and other results are taken from previous studies.

ZsRE dataset, which is distinct from the requirements of downstream tasks, leading to the stored facts not being retrieved during inference. Similarly, SERAC, utilizing external memory for edited facts, preserves general NLP abilities post-editing. This preservation stems from SERAC being trained once before editing begins, solely performing inference during editing, thereby preventing changes in the model’s output, even after multiple edits.

Overall, parameter-modifying methods degrade downstream task performance by altering pre-trained LLM parameters. In contrast, parameter-preserving methods maintain the original parameters, resulting in stable downstream task performance even after multiple edits.

5 Future Prospects

5.1 Leveraging Information Retrieval and External Memory

Research shows that using external knowledge bases, rather than relying solely on internal knowledge, benefits LLMs by guiding content generation based on predefined facts. External knowledge sources, such as text corpora, structured tables,

or key-value databases, can be utilized either to finetune LLMs for improved information retrieval or to employ prompting techniques for querying these sources. These approaches separate factual knowledge from inference process, thus preserves the original model parameters and minimizes post-editing damage. Moreover, they ensure that generated content aligns with predefined knowledge bases, thereby enhancing accountability and accuracy.

5.2 Improving Understandings of LLMs’ Internal Knowledge Structures

While identifying where factual knowledge is stored in LLMs has been extensively explored (Meng et al., 2022, 2023; Dai et al., 2022; Hernandez et al., 2024; Geva et al., 2021), the correlation between knowledge location and editing success remains low (Hase et al., 2023). Additionally, despite evidence suggesting a strong connection between factual knowledge and the feed-forward network layers (Meng et al., 2022; Geva et al., 2021, 2022), recent findings indicate that updates to multi-head self-attention layers also improve outcomes (Li et al., 2023). This suggests that

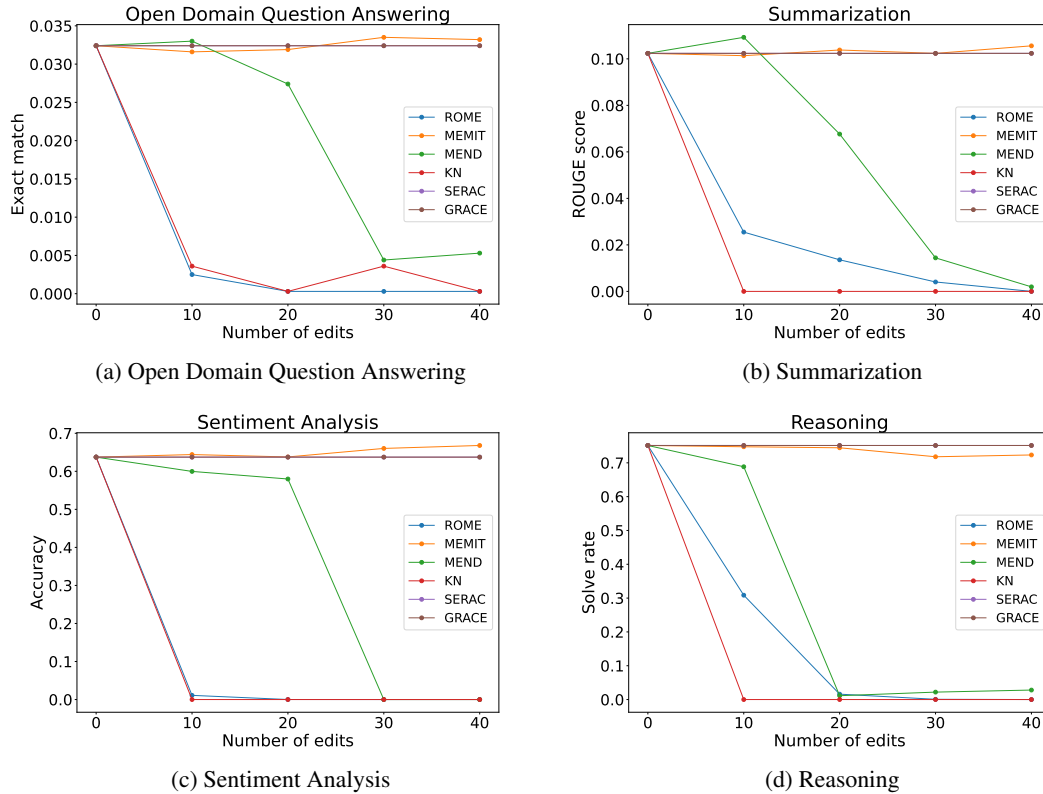


Figure 4: The experimental results for the deterioration of general abilities were obtained by editing GPT-J with various editing algorithms, including ROME, MEMIT, MEND, KN, SERAC, and GRACE, each applied 10 to 40 times. The edited models were subsequently evaluated on four downstream tasks, including open-domain question answering, sentiment analysis, summarization, and reasoning. The results for SERAC and GRACE are overlapping.

465 locating fact storage alone doesn't fully explain
 466 knowledge structures in LLMs. Further research
 467 is needed to understand how knowledge locations
 468 interact with model predictions in order to enhance
 469 LLM interpretability and controllability.

470 Preserving LLMs' general capabilities is also
 471 crucial for model editing, as discussed in Sec-
 472 tion 3.3. Recent breakthroughs in identifying re-
 473 gions within models that correlate with general lin-
 474 guistic abilities have opened up a direction for fu-
 475 ture research in model editing (Zhang et al., 2024b).
 476 By making targeted modifications, we can poten-
 477 tially prevent the deterioration of general abilities
 478 and improve the specificity and effectiveness of
 479 model editing methods.

480 5.3 Improving Robustness of Knowledge 481 Editing

482 Even after achieving fair scores on the existing met-
 483 rics, models may revert to pre-edit versions or pro-
 484 vide ambiguous answers if the altered knowledge
 485 is conflicted with inherited concepts. Experiments
 486 show that more popular knowledge is easier for

487 modified models to revert to (Ma et al., 2024), in-
 488 dicated the lack of robustness in current editing
 489 strategies. A deeper understanding of how LLMs
 490 store and process interconnected knowledge enti-
 491 ties is crucial for more robust editing and warrants
 492 future research.

493 6 Conclusion

494 Although model editing techniques appear promis-
 495 ing for cost-effectively updating knowledge, they
 496 still have significant pitfalls. Current editing meth-
 497 ods often struggle with making logical inferences
 498 based on the edited facts, introducing unintended al-
 499 terations of non-target knowledge and deterioration
 500 in model performance, particularly with parameter-
 501 modified methods. By harnessing information re-
 502 trieval techniques and delving into how models
 503 store and process knowledge, deviations in model
 504 abilities can be mitigated, and the controllability
 505 of edited facts can be enhanced, ultimately leading
 506 to greater robustness. We hope our work illumi-
 507 nates potential directions for future improvements
 508 in knowledge editing.

7 Limitations

The field of knowledge editing is advancing at an impressive pace, with numerous innovations in editing methodologies and evaluation metrics being proposed. Despite our efforts to collect and organize previous work, some contributions may not be included in this paper. However, we will continue to monitor the latest developments in this field and update our GitHub repository with recent related works.

References

Yonatan Bisk, Rowan Zellers, Ronan Le bras, Jianfeng Gao, and Yejin Choi. 2020. [Piqa: Reasoning about physical commonsense in natural language](#). *Proceedings of the AAAI Conference on Artificial Intelligence*, 34(05):7432–7439.

Nicola De Cao, Wilker Aziz, and Ivan Titov. 2021. [Editing factual knowledge in language models](#). In *Conference on Empirical Methods in Natural Language Processing*.

Karl Cobbe, Vineet Kosaraju, Mohammad Bavarian, Mark Chen, Heewoo Jun, Lukasz Kaiser, Matthias Plappert, Jerry Tworek, Jacob Hilton, Reiichiro Nakano, Christopher Hesse, and John Schulman. 2021. [Training verifiers to solve math word problems](#). *ArXiv*, abs/2110.14168.

Roi Cohen, Eden Biran, Ori Yoran, Amir Globerson, and Mor Geva. 2023. [Evaluating the ripple effects of knowledge editing in language models](#). *Transactions of the Association for Computational Linguistics*, 12:283–298.

Damai Dai, Li Dong, Yaru Hao, Zhifang Sui, Baobao Chang, and Furu Wei. 2022. [Knowledge neurons in pretrained transformers](#). In *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 8493–8502, Dublin, Ireland. Association for Computational Linguistics.

Qingxiu Dong, Damai Dai, Yifan Song, Jingjing Xu, Zhifang Sui, and Lei Li. 2022. [Calibrating factual knowledge in pretrained language models](#). In *Findings of the Association for Computational Linguistics: EMNLP 2022*, pages 5937–5947, Abu Dhabi, United Arab Emirates. Association for Computational Linguistics.

Mor Geva, Avi Caciularu, Kevin Wang, and Yoav Goldberg. 2022. [Transformer feed-forward layers build predictions by promoting concepts in the vocabulary space](#). In *Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing*, pages 30–45, Abu Dhabi, United Arab Emirates. Association for Computational Linguistics.

Mor Geva, Roei Schuster, Jonathan Berant, and Omer Levy. 2021. [Transformer feed-forward layers are key-value memories](#). In *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing*, pages 5484–5495, Online and Punta Cana, Dominican Republic. Association for Computational Linguistics.

Bogdan Gliwa, Iwona Mochol, Maciej Biesek, and Aleksander Wawer. 2019. [SAMSum corpus: A human-annotated dialogue dataset for abstractive summarization](#). In *Proceedings of the 2nd Workshop on New Frontiers in Summarization*, pages 70–79, Hong Kong, China. Association for Computational Linguistics.

Jia-Chen Gu, Hao-Xiang Xu, Jun-Yu Ma, Pan Lu, Zhen-Hua Ling, Kai-Wei Chang, and Nanyun Peng. 2024. [Model editing can hurt general abilities of large language models](#).

Akshat Gupta, Anurag Rao, and Gopala Anumanchipalli. 2024a. [Model editing at scale leads to gradual and catastrophic forgetting](#).

Akshat Gupta, Dev Sajnani, and Gopala Anumanchipalli. 2024b. [A unified framework for model editing](#).

Thomas Hartvigsen, Swami Sankaranarayanan, Hamid Palangi, Yoon Kim, and Marzyeh Ghassemi. 2023. [Aging with grace: Lifelong model editing with discrete key-value adapters](#). In *Advances in Neural Information Processing Systems*.

Peter Hase, Mohit Bansal, Been Kim, and Asma Ghandeharioun. 2023. [Does localization inform editing? surprising differences in causality-based localization vs. knowledge editing in language models](#). In *Thirty-seventh Conference on Neural Information Processing Systems*.

Evan Hernandez, Arnab Sen Sharma, Tal Haklay, Kevin Meng, Martin Wattenberg, Jacob Andreas, Yonatan Belinkov, and David Bau. 2024. [Linearity of relation decoding in transformer language models](#). In *The Twelfth International Conference on Learning Representations*.

Jason Hoelscher-Obermaier, Julia Persson, Esben Kran, Ioannis Konstas, and Fazl Barez. 2023. [Detecting edit failures in large language models: An improved specificity benchmark](#). In *Findings of the Association for Computational Linguistics: ACL 2023*, pages 11548–11559, Toronto, Canada. Association for Computational Linguistics.

Wenyue Hua, Jiang Guo, Mingwen Dong, Henghui Zhu, Patrick Ng, and Zhiguo Wang. 2024. [Propagation and pitfalls: Reasoning-based assessment of knowledge editing through counterfactual tasks](#).

Zeyu Huang, Yikang Shen, Xiaofeng Zhang, Jie Zhou, Wenge Rong, and Zhang Xiong. 2023. [Transformer-patcher: One mistake worth one neuron](#). In *The Eleventh International Conference on Learning Representations*.

618	Tom Kwiatkowski, Jennimaria Palomaki, Olivia Redfield, Michael Collins, Ankur Parikh, Chris Alberti, Danielle Epstein, Illia Polosukhin, Jacob Devlin, Kenton Lee, Kristina Toutanova, Llion Jones, Matthew Kelcey, Ming-Wei Chang, Andrew M. Dai, Jakob Uszkoreit, Quoc Le, and Slav Petrov. 2019. Natural questions: A benchmark for question answering research . <i>Transactions of the Association for Computational Linguistics</i> , 7:452–466.	1631–1642, Seattle, Washington, USA. Association for Computational Linguistics.	674 675
627	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.	Chenmian Tan, Ge Zhang, and Jie Fu. 2024. Massive editing for large language models via meta learning . In <i>International Conference on Learning Representations</i> .	676 677 678 679
633	Xiaopeng Li, Shasha Li, Shezheng Song, Jing Yang, Jun Ma, and Jie Yu. 2023. Pmet: Precise model editing in a transformer . In <i>AAAI Conference on Artificial Intelligence</i> .	Mojtaba Valipour, Mehdi Rezagholizadeh, Ivan Kobyzev, and Ali Ghodsi. 2023. DyLoRA: Parameter-efficient tuning of pre-trained models using dynamic search-free low-rank adaptation . In <i>Proceedings of the 17th Conference of the European Chapter of the Association for Computational Linguistics</i> , pages 3274–3287, Dubrovnik, Croatia. Association for Computational Linguistics.	680 681 682 683 684 685 686 687
637	Zhoubo Li, Ningyu Zhang, Yunzhi Yao, Mengru Wang, Xi Chen, and Huajun Chen. 2024. Unveiling the pitfalls of knowledge editing for large language models . In <i>The Twelfth International Conference on Learning Representations</i> .	Ben Wang and Aran Komatsuzaki. 2021. GPT-J-6B: A 6 Billion Parameter Autoregressive Language Model . https://github.com/kingoflolz/mesh-transformer-jax .	688 689 690 691
642	Xinbei Ma, Tianjie Ju, Jiyang Qiu, Zhuosheng Zhang, hai zhao, lifeng Liu, and Yulong Wang. 2024. Is it possible to edit large language models robustly? In <i>ICLR 2024 Workshop on Large Language Model (LLM) Agents</i> .	Yiwei Wang, Muhao Chen, Nanyun Peng, and Kai wei Chang. 2024. Deepedit: Knowledge editing as decoding with constraints . <i>ArXiv</i> , abs/2401.10471.	692 693 694
647	Vittorio Mazzia, Alessandro Pedrani, Andrea Caciolai, Kay Rottmann, and Davide Bernardi. 2023. A survey on knowledge editing of neural networks . <i>ArXiv</i> , abs/2310.19704.	Wanli Yang, Fei Sun, Xinyu Ma, Xun Liu, Dawei Yin, and Xueqi Cheng. 2024. The butterfly effect of model editing: Few edits can trigger large language models collapse .	695 696 697 698
651	Kevin Meng, David Bau, Alex Andonian, and Yonatan Belinkov. 2022. Locating and editing factual associations in gpt . In <i>Neural Information Processing Systems</i> .	Yunzhi Yao, Peng Wang, Bo Tian, Siyuan Cheng, Zhoubo Li, Shumin Deng, Huajun Chen, and Ningyu Zhang. 2023. Editing large language models: Problems, methods, and opportunities . In <i>Conference on Empirical Methods in Natural Language Processing</i> .	699 700 701 702 703
655	Kevin Meng, Arnab Sen Sharma, Alex Andonian, Yonatan Belinkov, and David Bau. 2023. Mass editing memory in a transformer . <i>The Eleventh International Conference on Learning Representations (ICLR)</i> .	Lang Yu, Qin Chen, Jie Zhou, and Liang He. 2024. Melo: Enhancing model editing with neuron-indexed dynamic lora . <i>Proceedings of the AAAI Conference on Artificial Intelligence</i> , 38(17):19449–19457.	704 705 706 707
660	Eric Mitchell, Charles Lin, Antoine Bosselut, Chelsea Finn, and Christopher D Manning. 2022a. Fast model editing at scale . In <i>International Conference on Learning Representations</i> .	Ningyu Zhang, Yunzhi Yao, Bozhong Tian, Peng Wang, Shumin Deng, Mengru Wang, Zekun Xi, Shengyu Mao, Jintian Zhang, Yuansheng Ni, Siyuan Cheng, Ziwen Xu, Xin Xu, Jia-Chen Gu, Yong Jiang, Pengjun Xie, Fei Huang, Lei Liang, Zhiqiang Zhang, Xiaowei Zhu, Jun Zhou, and Huajun Chen. 2024a. A comprehensive study of knowledge editing for large language models .	708 709 710 711 712 713 714 715
664	Eric Mitchell, Charles Lin, Antoine Bosselut, Chelsea Finn, and Christopher D. Manning. 2022b. Memory-based model editing at scale . In <i>International Conference on Machine Learning</i> .	Zhihao Zhang, Jun Zhao, Qi Zhang, Tao Gui, and Xuanjing Huang. 2024b. Unveiling linguistic regions in large language models . Association for Computational Linguistics.	716 717 718 719
668	Richard Socher, Alex Perelygin, Jean Wu, Jason Chuang, Christopher D. Manning, Andrew Ng, and Christopher Potts. 2013. Recursive deep models for semantic compositionality over a sentiment treebank . In <i>Proceedings of the 2013 Conference on Empirical Methods in Natural Language Processing</i> , pages	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? In <i>The 2023 Conference on Empirical Methods in Natural Language Processing</i> .	720 721 722 723 724
673		Zexuan Zhong, Zhengxuan Wu, Christopher D Manning, Christopher Potts, and Danqi Chen. 2023. MQuAKE: Assessing knowledge editing in language models via multi-hop questions . In <i>The 2023 Conference on Empirical Methods in Natural Language Processing</i> .	725 726 727 728 729

A Detailed Explanation of Evaluation Metrics and Examples

A.1 Portability / Generalization

Single Edit In the single edit scenario, we modify only one fact in the logical chain with each edit.

- **One-Hop:** This setting focuses on evaluating the impact of a single edit on direct, one-hop reasoning tasks.

For one-hop evaluations, we adopt the methods proposed by (Yao et al., 2023). These include:

- **Subject Replace:** This metric tests the model’s generalization ability by replacing the subject in the question with an alias or synonym, assessing if the edited attribute is generalized to other descriptions of the same subject.
- **Reversed Relation:** This metric evaluates the model’s capability to handle reversed relations by filtering for suitable relations (e.g., one-to-one relation) and asking the reverse question to check if the target entity is also updated.
- **One-Hop Test:** This metric assesses the edited language model’s performance on downstream tasks that require one-hop reasoning.

Multiple Edits In the multiple edits scenario, we evaluate the model’s performance after applying several logically related edits. This part consists of:

- **Multi-Hop editing:** Evaluate whether the model can infer edited knowledge in multi-hop questions.
- **Conflict editing:** Assess how the model handles multiple conflicting edits.

In the multi-hop setting, we assess the model’s performance on multi-hop questions using the evaluation methods proposed by (Zhong et al., 2023), which include:

- **Edit-wise Success Rate (EW):** This metric measures how many facts can be successfully recalled from the edited language model.

$$EW = \mathbb{1}\{f^*(s) = o^*\} \quad (5)$$

where f^* is the model after editing, s refers to the edited subject, and o refers to target object.

- **Instance-wise Accuracy (IW):** This metric tests how many multi-hop instances the model can recall all the individual single-hop facts. This metric is crucial for multi-hop performance, as the model must encode each fact to answer the multi-hop question.

$$IW = \mathbb{1}\left\{\bigwedge_{(s,r,o^*) \in C^*} [f^*(s) = o^*]\right\} \quad (6)$$

where $C^* = \langle (s_1, r_1, o_1), \dots, (s_n, r_n, o_n) \rangle$ is the chain of facts of a multi-hop question. In this chain, the object of the i^{th} fact is the subject of the next fact. (i.e., $o_i = s_{i+1}$)

- **Multi-hop Accuracy (MH):** This metric assesses the accuracy of the original and edited language models on multi-hop questions. In the MQuAKE dataset (Zhong et al., 2023), there are three generated multi-hop questions for each instance. If any of the three questions is correctly answered by the model, we consider it accurate.

$$MH = \mathbb{1}\left\{\bigvee_{q \in Q} f^*(q) = a^*\right\} \quad (7)$$

where Q is a set of similar multi-hop questions with the same answer a^* .

As for Conflict editing, we use the setting and evaluation methods from (Li et al., 2024). The settings consist of:

- **Reverse Conflict:** This setting introduces conflicts by editing facts with reverse relations. For example:
edit 1: $(s_1, r_1, o_1 \rightarrow o_2)$
Hamlet was written by Shakespeare \rightarrow *Agatha Christie*.
edit 2: $(o_2, r_2, s_1 \rightarrow s_2)$
The notable work of Agatha Christie is Hamlet \rightarrow *Odyssey*
the updated knowledge then could be represented as:

$$\begin{cases} k_o = (s_1, r_1, o_2) \\ k_n = (s_2, r_1, o_2) \end{cases} \quad 810$$

where k_o refers to old knowledge, and k_n refers to new knowledge.

- **Composite Conflict:** This explores more complex situations where the edits are associated with a fact that is not influenced by the

editing (**tied fact**). For example:

edit 1: $(s_1, r_1, o_1 \rightarrow o_2)$

Hamlet was written in English \rightarrow *French*

edit 2: $(s_2, r_2, o_2 \rightarrow o_3)$

Shakespeare wrote in French \rightarrow *German*

tied fact: (s_1, r, s_2)

The notable work of Shakespeare is Hamlet

where $r \wedge r_1 \rightarrow r_2$ is a logical rule. The updated knowledge then could be represented as:

$$\begin{cases} k_f = (s_1, r, s_2) \\ k_0 = (s_1, r_1, o_2) \\ k_n = (s_1, r_1, o_3) \end{cases}$$

where k_f refers to a tied fact.

The evaluation methods include:

- **Conflict Score (CS):** Measures how well a knowledge editing method handles knowledge conflicts by calculating the ratio that the new fact is more probable than the old fact after knowledge editing.

$$CS = \mathbb{1}\{p_{f'_\theta}(k_n) > p_{f'_\theta}(k_o)\} \quad (8)$$

- **Conflict Magnitude (CM):** Estimates the decrease in probability of the old fact after editing.

$$CM = \frac{p_{f_{\theta^m}}(k_o) - p_{f_{\theta'}}(k_o)}{p_{f_{\theta^m}}(k_o)} \quad (9)$$

θ^m is the intermediate model parameters after *edit 1*.

A.2 Locality

Single Edit In the single edit scenario for locality, we adopt the methods proposed by (Yao et al., 2023), including:

- **Other Attribution (OA):** The modified **ZsRE** and **CounterFact** datasets are applied to test whether the non-target attributes of the edited subjects remained the same. For example, if we reset *Lionel Messi* as a basketball player, his nationality should stay the same.
- **Distract Neighbor (DN):** Previous studies indicate that if edit cases are concatenated with unrelated context, the model tends to output content related to the edit cases. For example, if the original prompt is "Windows 11 is a product of ___", an edit case is added in front and be "Windows 11 is a product of

Google. Office 365, developed by ___". It testifies whether the model prediction would be "distracted" by the edit case.

- **Other Task (OT)** The edited model is tested on the multiple-choice QA task **Physical Interaction QA** (PIQA, Bisk et al. (2020)) and the performance is evaluated by accuracy.

Multiple Edits We also test the model's locality in the multiple edits scenario by adopting the methods and evaluations from (Li et al., 2024). The settings consist of:

- **Round Edit:** This edits the knowledge triplet back-and-forth, for example:

edit 1: $(s, r, o_1 \rightarrow o^*)$

edit 2: $(s, r, o^* \rightarrow o_1)$

where o^* is an intermediate object.

The evaluation metrics include:

- **Distortion (D) (Li et al., 2024):**

$$D = JS(p_{f_\theta}(\text{Obj} | (s, r)), p_{f_{\theta'}}(\text{Obj} | (s, r))) \quad (10)$$

estimates the JS divergence of the objects distribution before and after edit.

- **Ignore Rate (IR) (Li et al., 2024):**

$$IR = \frac{1}{|\text{Obj}| - 1} \sum_{o \in \text{Obj} \setminus \{o_1\}} \mathbb{1}\{p_{f_\theta}(o | (s, r)) > p_{f_{\theta'}}(o | (s, r))\} \quad (11)$$

measures the extent to which objects in Obj set (excluding the target object o_1) are disregarded or overlooked after the process of knowledge editing.

- **Failure Rate (FR) (Li et al., 2024):**

$$FR = \mathbb{1}\{IR > 0.5\} \quad (12)$$

calculates the rate when Ignore Rate > 0.5

- **Tied Fact Damage (TDF) (Li et al., 2024):**

$$TDF = \frac{p_{f_{\theta^m}}(k_f) - p_{f_{\theta'}}(k_f)}{p_{f_{\theta^m}}(k_f)} \quad (13)$$

k_f denotes the tied facts and θ^m is the intermediate model parameters after *edit 1*.

Other Locality Metrics

- **Neighborhood KL Divergence (Hoelscher-Obermaier et al., 2023):**

$$\text{NKL} \stackrel{\text{def}}{=} \sum_{w \in W} \log \left(\frac{P(w)}{P^*(w)} \right) \quad (14)$$

- **Neighborhood Score (NS) (Meng et al., 2022):** collect a set of "neighborhood" subjects and evaluate the success fraction for $\mathbb{P}[o^c] > \mathbb{P}[o^*]$, while the o^c denotes the correct facts and o^* denotes the false facts.
- **Neighborhood Magnitude (NM) (Meng et al., 2022):** the differences of $\mathbb{P}[o^c]$ and $\mathbb{P}[o^*]$ for the "neighborhood" subjects.

B Detailed Explanation of experiments for deterioration of general LLM abilities

We follow the settings of (Gu et al., 2024) for this part of experiments. Different evaluation metrics were applied for each downstream task: Exact Match for open-domain question answering on the Natural Question dataset (Kwiatkowski et al., 2019), accuracy for sentiment analysis on the SST2 dataset (Socher et al., 2013), solve rate for reasoning on the GSM8K dataset (Cobbe et al., 2021), and ROUGE score for summarization on the SAM-Sum dataset (Gliwa et al., 2019).